



7450 Ethernet Service Switch
7750 Service Router
Virtualized Service Router
Release 23.10.R1

Multiservice ISA and ESA Guide

3HE 19229 AAAC TQZZA 01
Edition: 01
October 2023

Nokia is committed to diversity and inclusion. We are continuously reviewing our customer documentation and consulting with standards bodies to ensure that terminology is inclusive and aligned with the industry. Our future customer documentation will be updated accordingly.

This document includes Nokia proprietary and confidential information, which may not be distributed or disclosed to any third parties without the prior written consent of Nokia.

This document is intended for use by Nokia's customers ("You"/"Your") in connection with a product purchased or licensed from any company within Nokia Group of Companies. Use this document as agreed. You agree to notify Nokia of any errors you may find in this document; however, should you elect to use this document for any purpose(s) for which it is not intended, You understand and warrant that any determinations You may make or actions You may take will be based upon Your independent judgment and analysis of the content of this document.

Nokia reserves the right to make changes to this document without notice. At all times, the controlling version is the one available on Nokia's site.

No part of this document may be modified.

NO WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY OF AVAILABILITY, ACCURACY, RELIABILITY, TITLE, NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE, IS MADE IN RELATION TO THE CONTENT OF THIS DOCUMENT. IN NO EVENT WILL NOKIA BE LIABLE FOR ANY DAMAGES, INCLUDING BUT NOT LIMITED TO SPECIAL, DIRECT, INDIRECT, INCIDENTAL OR CONSEQUENTIAL OR ANY LOSSES, SUCH AS BUT NOT LIMITED TO LOSS OF PROFIT, REVENUE, BUSINESS INTERRUPTION, BUSINESS OPPORTUNITY OR DATA THAT MAY ARISE FROM THE USE OF THIS DOCUMENT OR THE INFORMATION IN IT, EVEN IN THE CASE OF ERRORS IN OR OMISSIONS FROM THIS DOCUMENT OR ITS CONTENT.

Copyright and trademark: Nokia is a registered trademark of Nokia Corporation. Other product names mentioned in this document may be trademarks of their respective owners.

© 2023 Nokia.

Table of contents

1	Getting started.....	19
1.1	About this guide.....	19
1.2	ISA and ESA configuration process.....	19
1.3	Conventions.....	21
1.3.1	Precautionary and information messages.....	21
1.3.2	Options or substeps in procedures and sequential workflows.....	22
2	ISA and ESA hardware.....	23
2.1	In this section.....	23
2.2	MS-ISA2 overview.....	23
2.3	MS-ISM overview.....	23
2.4	ESA overview.....	24
2.5	Application Assurance hardware features.....	27
2.5.1	AA system support.....	27
2.5.2	Host IOM support for AA on ISAs.....	27
2.5.3	Host IOM support for AA on ESAs.....	28
2.6	Configuring an ESA with CLI.....	28
2.6.1	Provisioning an ESA and ESA-VM.....	28
3	Application Assurance.....	30
3.1	AA overview.....	30
3.1.1	AA: inline policy enforcement.....	30
3.1.2	AA integration in subscriber edge gateways.....	31
3.1.3	Fixed residential broadband services.....	33
3.1.3.1	DS-Lite.....	33
3.1.3.2	6to4 /6RD.....	35
3.1.4	Wireless LAN gateway broadband services.....	36
3.1.5	Application-aware business VPN services.....	37
3.1.6	Mobile Backhaul.....	38
3.1.7	Stateful firewall service.....	40
3.2	AA system architecture.....	40
3.2.1	AA ISA resource configuration.....	40
3.2.1.1	AA ISA groups.....	40
3.2.1.2	Redundancy.....	42

3.2.1.3	ISA load balancing.....	43
3.2.1.4	Asymmetry removal.....	44
3.2.1.5	ISA overload detection.....	51
3.2.2	AA packet processing.....	52
3.2.2.1	Divert of traffic and subscribers.....	53
3.2.2.2	Traffic identification.....	71
3.2.2.3	Statistics and accounting.....	79
3.2.2.4	AQP.....	107
3.2.2.5	AA TCP Optimization.....	113
3.2.2.6	AA Dynamic Experience Management.....	117
3.2.2.7	AA HTTP redirect.....	121
3.2.2.8	AA HTTPS policy redirect.....	124
3.2.2.9	ICAP - large scale category-based URL filtering.....	124
3.2.2.10	Web-service URL classification.....	125
3.2.2.11	Local URL-list filtering.....	132
3.2.2.12	HTTP header enrichment.....	134
3.2.2.13	ACR HTTP enrichment.....	139
3.2.2.14	HTTP in browser notification.....	140
3.2.2.15	AA firewall.....	142
3.2.3	Service monitoring and debugging.....	161
3.2.4	CPU utilization.....	162
3.2.5	CLI batch: begin, commit and abort commands.....	162
3.3	Configuring AA with CLI.....	163
3.3.1	Provisioning AA ISA MDA.....	163
3.3.2	Configuring an AA ISA group.....	164
3.3.2.1	Configuring watermark parameters.....	166
3.3.3	Configuring a group policy.....	166
3.3.3.1	Beginning and committing a policy configuration.....	166
3.3.3.2	Aborting a policy configuration.....	167
3.3.3.3	Configuring an IP prefix list.....	167
3.3.3.4	Configuring AA session filters.....	168
3.3.3.5	Configuring an application group.....	170
3.3.3.6	Configuring an application.....	170
3.3.3.7	Configuring an application filter.....	170
3.3.3.8	Configuring an application profile.....	171
3.3.3.9	Configuring suppressible app-profile with SRRP.....	172

3.3.3.10	Configuring application service options.....	173
3.3.3.11	Configuring a policer.....	173
3.3.3.12	Configuring an application QoS policy.....	174
3.3.3.13	Configuring a charging filter.....	175
3.3.3.14	Configuring an application and DNS IP cache for URL content charging strengthening.....	176
3.3.3.15	Configuring an HTTP error redirect.....	178
3.3.3.16	Configuring HTTP header enrichment.....	179
3.3.3.17	Configuring an HTTP redirect policy.....	180
3.3.3.18	Configuring an HTTPS policy redirect.....	181
3.3.3.19	Configuring a captive redirect HTTP redirect policy.....	182
3.3.3.20	Configuring ICAP URL filtering.....	184
3.3.3.21	Configuring web-service URL classification.....	187
3.3.3.22	Configuring local URL-list filtering.....	190
3.3.3.23	Configuring HTTP notification.....	191
3.3.3.24	Configuring tethering detection.....	192
3.3.4	Configuring AA volume accounting and statistics.....	193
3.3.4.1	Configuring cflowd collector.....	194
3.3.4.2	Configuring AA comprehensive, RTP performance, TCP performance, and volume reporting.....	196
4	IP tunnels.....	200
4.1	IP tunnels overview.....	200
4.1.1	Tunnel ISAs.....	202
4.1.1.1	Public tunnel SAPs.....	203
4.1.1.2	Private tunnel SAPs.....	204
4.1.1.3	IP interface configuration.....	204
4.1.1.4	GRE and IP-IP tunnel configuration.....	205
4.1.1.5	IP fragmentation and reassembly for IP tunnels.....	207
4.1.1.6	TCP MSS adjustment.....	208
4.1.1.7	MTU propagation.....	208
4.1.2	Operational conditions.....	210
4.1.2.1	Dynamic configuration change support for IPsec gateway.....	211
4.1.3	QoS interactions.....	213
4.1.4	OAM interactions.....	214
4.1.5	Redundancy.....	214
4.1.6	Statistics collection.....	215

4.1.7	Security.....	215
4.1.7.1	GRE tunnel multicast support.....	215
4.1.7.2	IPv6 over IPv4 GRE tunnel.....	216
4.1.8	IKEv2.....	217
4.1.8.1	IKEv2 traffic selector and TS-list.....	217
4.1.8.2	IKEv2 fragmentation.....	219
4.1.9	SHA2 support.....	219
4.1.10	IPsec client lockout.....	220
4.1.11	IPsec tunnel CHILD_SA rekey.....	220
4.1.12	Multiple IKE/ESP transform support.....	221
4.2	Using certificates for IPsec tunnel authentication.....	222
4.2.1	IKEv2 digital signature authentication.....	222
4.3	Trust-anchor profile.....	223
4.4	Cert-profile.....	223
4.5	Video wholesale example.....	224
4.6	1:1 Multi-chassis IPsec redundancy.....	225
4.6.1	Architecture.....	226
4.6.2	MIMP.....	226
4.6.2.1	MIMP protocol states.....	227
4.6.2.2	Election logic.....	227
4.6.2.3	Protection status.....	228
4.6.2.4	Other details.....	228
4.6.3	Routing.....	229
4.6.3.1	Routing in public service.....	229
4.6.3.2	Routing in private services.....	229
4.6.3.3	Other details about shunting.....	230
4.6.4	MC-IPsec aware VRRP.....	230
4.6.5	Synchronization.....	231
4.6.5.1	Automatic CHILD_SA rekey.....	231
4.6.5.2	Encryption of Synced States.....	231
4.6.6	MC-IPsec Responder only.....	233
4.7	N:M MC-IPsec redundancy.....	233
4.7.1	Redundancy domain.....	234
4.7.2	Redundancy role.....	234
4.7.3	ISA tunnel-member pool.....	234
4.7.4	Redundancy state and protection status.....	235

4.7.5	Election.....	235
4.7.6	Revertive.....	236
4.7.7	DA versus DS.....	237
4.7.8	State synchronization.....	237
4.7.9	Routing.....	238
4.7.10	VRRP.....	238
4.7.11	Shunting.....	238
4.7.12	Responder only.....	239
4.7.13	Coexisting tunnel groups.....	239
4.7.14	Provisioning requirements.....	239
4.7.15	Configuring N:M.....	239
4.8	IPsec deployment requirements.....	244
4.9	IKEv2 remote-access tunnel.....	245
4.9.1	IKEv2 remote access tunnel – RADIUS-based PSK/certificate authentication.....	246
4.9.1.1	IKEv2 remote-access tunnel – EAP authentication.....	249
4.9.2	IKEv2 remote-access tunnel – authentication without RADIUS.....	251
4.9.3	IKEv2 remote-access tunnel – address assignment.....	252
4.9.3.1	DHCPv4 address assignment.....	253
4.9.3.2	DHCPv6 address assignment.....	253
4.9.3.3	DHCPv4/v6 usage notes.....	254
4.9.4	IPv6 IPsec support.....	255
4.9.4.1	IPv6 as payload.....	256
4.9.4.2	IPv6 as payload: static LAN-to-LAN tunnel.....	256
4.9.4.3	IPv6 as payload: dynamic LAN-to-LAN tunnel.....	256
4.9.4.4	IPv6 as payload: remote-access tunnel.....	256
4.9.4.5	IPv6 as encapsulation.....	257
4.10	MLDv2 over IPsec.....	257
4.10.1	MLDv2 over IPsec – traffic selector.....	257
4.10.2	MLDv2 over IPsec – configuration.....	258
4.11	Secured interface.....	258
4.12	IPsec client database.....	259
4.13	IPsec transport mode protected IP tunnel.....	262
4.13.1	Packet format.....	262
4.13.2	IKEv2 and SA.....	263
4.13.3	Traffic selector.....	263
4.13.4	Fragmentation and reassembly.....	263

4.13.5	QoS marking.....	263
4.13.6	Configuring IPsec transport mode protected tunnels.....	264
4.14	Configuring IPsec with CLI.....	265
4.14.1	Provisioning a tunnel ISA.....	265
4.14.2	Configuring a tunnel group.....	265
4.14.3	Configuring router interfaces for IPsec.....	265
4.14.4	Configuring IPsec parameters.....	266
4.14.5	Configuring IPsec in services.....	266
4.14.6	Configuring X.509v3 certificate parameters.....	267
4.14.7	Configuring MC-IPsec.....	269
4.14.7.1	Configuring MIMP.....	269
4.14.7.2	Configuring multi-chassis synchronization.....	269
4.14.7.3	Configuring routing for MC-IPsec.....	270
4.14.7.4	Configuring MCS encryption.....	271
4.14.8	Configuring and using CMPv2.....	272
4.14.9	Configuring OCSP.....	273
4.14.10	Configuring IKEv2 remote — access tunnel.....	274
4.14.11	Configuring IKEv2 remote — access tunnel with local address assignment.....	276
4.14.12	Configuring secured interfaces.....	277
5	L2TPV3 tunnels.....	279
5.1	L2TPv3 overview.....	279
5.2	Control plane.....	280
5.3	Public SAP.....	280
5.4	Private SAP.....	281
6	Video services.....	282
6.1	Video services.....	282
6.1.1	Video groups.....	282
6.1.2	Video SAP.....	282
6.1.3	Video interface.....	283
6.1.3.1	Video interface properties.....	283
6.1.4	Multicast information policies.....	284
6.1.5	Perfect stream protection.....	285
6.1.6	Perfect stream selection.....	285
6.1.6.1	Stream identification.....	285

6.1.6.2	Initial sequence identification.....	286
6.1.6.3	Packet selection.....	286
6.1.6.4	Clock recovery.....	287
6.1.6.5	Playout.....	287
6.1.6.6	Loss of transport.....	288
6.1.6.7	Perfect stream in relation to FCC/RET.....	288
6.1.6.8	Perfect stream in relation to VQM.....	288
6.1.7	Video quality monitoring.....	288
6.1.7.1	VoIP/video/teleconferencing performance measurements.....	296
6.1.7.2	Mean Opinion Score (MOS) performance measurements solution architecture..	296
6.2	Retransmission and Fast Channel Change.....	297
6.2.1	RET and FCC overview.....	297
6.2.1.1	Retransmission.....	297
6.2.1.2	FCC.....	298
6.2.1.3	RET and FCC server concurrency.....	302
6.2.2	Separate timers for FCC and RET.....	303
6.2.3	Peak bandwidth and sessions per ISA.....	303
6.2.4	Video support on ESA.....	304
6.3	IP Video for Live.....	304
6.3.1	Delayed IGMP join for FCC.....	305
6.4	Configuring video service components with CLI.....	305
6.4.1	Video services overview.....	305
6.4.1.1	Configuring an ISA module.....	306
6.4.1.2	Configuring a video group.....	307
6.4.1.3	Configuring a video SAP and video interface in a service.....	308
6.4.1.4	Basic multicast information policy configuration.....	309
6.4.2	Sample configurations.....	310
6.5	Configuring RET/FCC video components with CLI.....	321
6.5.1	Configuring RET/FCC video features in the CLI.....	321
6.5.1.1	Configuring the RET server.....	321
6.5.1.2	Configuring the FCC server.....	324
6.5.1.3	Logging and accounting collection for video statistics.....	328
7	Network Address Translation.....	330
7.1	Terminology.....	330
7.2	Network Address Translation (NAT) overview.....	331

7.2.1	Principles of NAT.....	332
7.2.2	Traffic load balancing.....	332
7.2.3	Application compatibility.....	335
7.3	Large Scale NAT.....	335
7.3.1	Port range blocks.....	336
7.3.1.1	Reserved ports and priority sessions.....	336
7.3.1.2	Preventing port block starvation.....	336
7.3.2	Pools with flexible port allocations.....	339
7.3.2.1	Free port limit.....	341
7.3.2.2	Restrictions.....	341
7.3.3	Association between NAT subscribers and IP addresses in a NAT pool.....	342
7.3.4	Timeouts.....	343
7.4	L2-Aware NAT.....	343
7.4.1	Port block extensions.....	345
7.4.1.1	Managing port block space.....	346
7.4.2	L2-Aware NAT bypass.....	348
7.4.2.1	Full ESM host bypass.....	348
7.4.2.2	Selective L2-Aware NAT bypass.....	349
7.4.3	L2-Aware NAT destination-based multiple NAT policies.....	353
7.4.3.1	Logging.....	354
7.4.3.2	Static port forwards.....	357
7.4.3.3	L2-Aware ping.....	357
7.4.3.4	UPnP.....	358
7.4.3.5	L2-Aware NAT and multicast.....	359
7.5	NAT pool addresses and ICMP Echo Request/Reply (ping).....	360
7.6	Traffic steering to NAT.....	362
7.6.1	Routing approach in LSN44.....	362
7.6.1.1	Static NAT routes.....	362
7.6.1.2	Dynamic routing to LSN44.....	364
7.6.1.3	Combination of static and dynamic routes.....	366
7.6.1.4	Scale and logging notes.....	366
7.6.2	NAT steering through IP filters.....	366
7.7	L2-Aware support for residential gateway types.....	369
7.8	One-to-one (1:1) NAT.....	371
7.8.1	Static 1:1 NAT.....	371
7.8.1.1	Protocol agnostic behavior.....	372

7.8.1.2	Modification of parameters in static 1:1 NAT.....	373
7.8.1.3	Load distribution over ISAs in static 1:1 NAT.....	373
7.8.1.4	NAT-policy selection.....	374
7.8.1.5	Mapping timeout.....	375
7.8.1.6	Logging.....	375
7.8.1.7	Restrictions.....	375
7.8.2	ICMP.....	375
7.9	Deterministic NAT.....	375
7.9.1	Overview.....	375
7.9.2	Supported deterministic NAT types.....	376
7.9.3	Number of subscribers per outside IP and per pool.....	376
7.9.4	Referencing a pool.....	376
7.9.5	Outside pool configuration.....	377
7.9.6	Mapping rules and the map command in deterministic LSN44.....	380
7.9.7	Hashing considerations in deterministic LSN44.....	383
7.9.7.1	Distribution of outside IP addresses across MS-ISAs in an MS-ISA NA group...	384
7.9.8	Sharing of deterministic NAT pools.....	385
7.9.9	Simultaneous support of dynamic and deterministic NAT.....	385
7.9.10	Selecting traffic for NAT.....	385
7.9.11	Inverse mappings.....	385
7.9.11.1	MIB approach.....	385
7.9.11.2	Off-line approach to obtain deterministic mappings.....	386
7.9.12	Logging.....	387
7.9.13	Deterministic DS-Lite.....	387
7.9.13.1	Hashing considerations in DS-Lite.....	389
7.9.13.2	Order of configuration steps in deterministic DS-Lite.....	390
7.10	Destination Based NAT (DNAT).....	392
7.10.1	Combination of SNAPT and DNAT.....	392
7.10.2	Forwarding model in DNAT.....	393
7.10.3	DNAT traffic selection via NAT classifier.....	394
7.10.4	Configuring DNAT.....	394
7.10.4.1	DNAT traffic selection and destination IP address configuration.....	394
7.10.4.2	Micro-netting original source (inside) IP space in DNAT-only case.....	395
7.11	LSN – multiple NAT policies per inside routing context.....	396
7.11.1	Restrictions.....	396
7.11.2	Multiple NAT policies per inside routing context.....	396

7.11.3	Routing approach for NAT diversion.....	398
7.11.4	Filter-based approach.....	399
7.11.5	Multiple NAT policies and deterministic NAT.....	399
7.11.5.1	Combination of deterministic LSN44, non-deterministic LSN44, and MNP.....	399
7.11.6	Multiple NAT policies with DS-Lite and NAT64.....	401
7.11.7	Default NAT policy.....	402
7.11.8	Scaling considerations.....	402
7.11.9	Multiple NAT policies and SPF configuration considerations.....	402
7.11.9.1	Multiple NAT policies and forwarding considerations.....	403
7.12	NAT policy selection in non-deterministic NAT.....	404
7.13	Default DMZ Host.....	406
7.14	NAT and CoA.....	408
7.14.1	CoA and NAT policies.....	408
7.14.2	CoA and DNAT.....	409
7.14.3	Modifying an active NAT prefix list or NAT classifier via CLI.....	412
7.15	Watermarks.....	414
7.16	Port forwards.....	414
7.16.1	Static port forwards.....	415
7.16.2	Port Control Protocol (PCP).....	416
7.16.3	PORT_SET option.....	417
7.16.3.1	Terminology.....	418
7.16.3.2	Enabling the PORT_SET option.....	418
7.16.3.3	Port allocation scheme.....	419
7.16.3.4	Limits and quotas.....	419
7.16.3.5	Port overlaps.....	420
7.16.3.6	Port allocation example.....	420
7.16.3.7	Operational considerations.....	421
7.16.4	Universal plug and play Internet Gateway Device service.....	422
7.16.4.1	Configuring UPnP IGD service.....	423
7.17	NAT Point-to-Point Tunneling Protocol (PPTP) ALG.....	423
7.17.1	PPTP protocol.....	423
7.17.1.1	Supported control messages.....	424
7.17.1.2	GRE tunnel.....	424
7.17.2	PPTP ALG operation.....	425
7.17.3	Multiple sessions initiated from the same PPTP client node.....	427
7.17.4	Selection of call IDs in NAT.....	427

7.18	Modifying active NAT prefix list or NAT classifier via CLI.....	427
7.19	NAT logging.....	430
7.19.1	Syslog/SNMP/local-file logging.....	430
7.19.1.1	Filtering LSN events to system memory.....	430
7.19.1.2	NAT logging to a local file.....	436
7.19.2	SNMP trap logging.....	437
7.19.3	NAT syslog.....	438
7.19.4	LSN RADIUS logging.....	439
7.19.4.1	Periodic RADIUS logging.....	446
7.19.4.2	RADIUS buffer management on ISA or ESA-VM.....	448
7.19.5	Summarization logs and bulk operations.....	450
7.19.5.1	Summarization logs and RADIUS logging.....	451
7.19.6	Integrated L2-Aware NAT RADIUS logging and BNG accounting.....	451
7.19.6.1	Enabling RADIUS logging for L2-Aware NAT subscribers.....	456
7.19.6.2	Timestamp interpretation.....	456
7.19.6.3	High logging rates.....	458
7.19.6.4	Intra-chassis redundancy.....	458
7.19.7	LSN and L2-Aware NAT flow logging.....	459
7.19.7.1	IPFIX flow logging.....	459
7.19.7.2	Template formats.....	460
7.19.7.3	Template format 1 and format 2.....	462
7.19.7.4	Configuration example.....	464
7.19.7.5	Syslog flow logging.....	466
7.20	DS-Lite and NAT64 fragmentation.....	470
7.20.1	Overview.....	470
7.20.2	IPv6 fragmentation in DS-Lite.....	470
7.20.3	NAT64.....	471
7.21	DS-Lite reassembly.....	472
7.21.1	Interpreting fragmentation statistics.....	472
7.21.2	Support for small first fragments.....	475
7.21.2.1	Upstream reassembly with small first IPv6 fragments less than 1280 bytes....	475
7.21.2.2	Downstream fragmentation with small first IPv6 fragment.....	475
7.22	Histogram.....	477
7.23	NAT redundancy.....	481
7.23.1	NAT stateless dual-homing.....	482
7.23.1.1	Configuration considerations.....	484

7.23.1.2	Troubleshooting commands.....	485
7.23.2	Active-active ISA redundancy model.....	488
7.23.2.1	Startup conditions.....	490
7.23.2.2	Recovery.....	490
7.23.2.3	Adding additional ISAs in the ISA group.....	490
7.23.3	L2-Aware bypass.....	490
7.23.3.1	Sharing IP addresses in L2-Aware NAT.....	492
7.23.3.2	Recovery.....	493
7.23.3.3	Default bypass during reboot or MS-ISA provisioning.....	493
7.23.3.4	Logging.....	493
7.23.4	Stateful inter-chassis NAT redundancy.....	493
7.23.4.1	Health status and failure events.....	496
7.23.4.2	Route advertisements.....	498
7.23.4.3	Flow synchronization.....	498
7.23.4.4	Rapid consecutive switchovers.....	500
7.23.4.5	ISA-to-ISA communication.....	500
7.23.4.6	Preemption.....	500
7.23.4.7	Message delivery prioritization.....	501
7.23.4.8	Subscriber-aware NAT.....	501
7.23.4.9	Matching configuration on redundant pair of nodes.....	501
7.23.4.10	Online configuration changes.....	502
7.23.4.11	Scenario with monitoring ports.....	503
7.23.4.12	Configuring stateful inter-chassis NAT redundancy.....	507
7.24	ISA feature interactions.....	509
7.24.1	MS-ISA use with service mirrors.....	509
7.24.2	Network Address Translation.....	509
7.24.3	Subscriber aware Large Scale NAT44.....	509
7.25	Mapping of Address and Port using Translation (MAP-T).....	517
7.25.1	MAP-T rules.....	519
7.25.2	A+P mapping algorithm.....	520
7.25.3	Routing considerations.....	521
7.25.4	Forwarding considerations in the BR.....	522
7.25.4.1	IPv6 addresses.....	523
7.25.4.2	1:1 translations and IPv4 prefix translations.....	523
7.25.4.3	Hub-and-spoke topology.....	524
7.25.4.4	Rule prefix overlap.....	524

7.25.5	BMR rules implementation example.....	524
7.25.6	ICMP.....	526
7.25.7	Fragmentation.....	526
7.25.7.1	Fragmentation in the downstream direction.....	526
7.25.7.2	Fragmentation in the upstream direction.....	527
7.25.7.3	Fragmentation statistics.....	528
7.25.8	Maximum Segment Size (MSS) adjust.....	529
7.25.9	Statistics collection.....	530
7.25.10	Logging.....	530
7.25.11	Licensing.....	531
7.25.12	Configuration.....	531
7.25.12.1	Modifying MAP-T parameters when the MAP-T domain is active.....	532
7.25.13	Inter-chassis redundancy.....	532
7.26	Configuring NAT.....	533
7.26.1	ISA redundancy.....	533
7.26.2	NAT Layer 2-Aware configurations.....	534
7.26.3	Large scale NAT configuration.....	536
7.26.4	NAT configuration examples.....	538
7.27	Configuring VSR-NAT.....	541
7.27.1	VSR-NAT licensing.....	541
7.27.2	Statistics collection For LSN bindings.....	542
7.27.3	Statistics collection for LSN bandwidth.....	542
7.27.4	VSR-NAT show command examples.....	543
7.28	VSR scaling profiles on BB-ISA.....	546
7.28.1	Scaling profiles on the VSR.....	546
7.28.2	Scale profile modification.....	546
7.29	NAT scaling profiles on ESA.....	547
7.29.1	Scaling profiles for NAT on ESA.....	547
7.29.2	Scale profile modification.....	547
7.30	Expanding a NAT group.....	547
8	Residential firewall.....	549
8.1	Residential firewall overview.....	549
8.1.1	Supported protocols and extension headers.....	549
8.1.1.1	Unknown protocols.....	549
8.1.1.2	TCP and UDP.....	550

8.1.1.3	ICMPv6.....	550
8.1.2	Application Layer Gateway.....	550
8.1.3	Additional filtering control.....	550
8.1.4	TCP MSS adjustment.....	550
8.1.5	Static port forwards and DMZ.....	551
8.2	Residential firewall provisioning.....	551
8.2.1	Domains and addressing.....	552
9	TCP MSS adjustment.....	553
9.1	Overview.....	553
9.2	TCP MSS adjustment for ESM hosts.....	553
9.3	TCP MSS adjustment for NAT services.....	554
10	L2TP network server.....	555
10.1	Subscriber agg-rate-limit on LNS.....	555
10.2	LNS reassembly.....	556
10.2.1	Overview.....	556
10.2.2	Reassembly function.....	557
10.2.3	Load sharing between the ISAs.....	558
10.2.4	Inter-chassis ISA redundancy.....	558
10.3	MLPPPoE, MLPPP(oE)oA with LFI on LNS.....	558
10.3.1	Terminology.....	559
10.3.2	LNS MLPPPoX.....	559
10.3.3	MLPPP encapsulation.....	559
10.3.4	MLPPPoX negotiation.....	559
10.3.5	Enabling MLPPPoX.....	560
10.3.6	Link Fragmentation and Interleaving (LFI).....	561
10.3.6.1	MLPPPoX fragmentation, MRRU and MRU considerations.....	561
10.3.7	LFI functionality implemented in LNS.....	562
10.3.7.1	Last mile QoS awareness in the LNS.....	563
10.3.7.2	BB-ISA processing.....	564
10.3.7.3	LNS-LAC link.....	565
10.3.7.4	AN-RG link.....	565
10.3.7.5	Home link.....	565
10.3.7.6	Optimum fragment size calculation by LNS.....	565
10.3.8	Upstream traffic considerations.....	567

10.3.9	Multiple links MLPPPoX with no interleaving.....	567
10.3.10	MLPPPoX session support.....	568
10.3.11	Session load balancing across multiple BB-ISAs.....	568
10.3.12	BB-ISA hashing considerations.....	569
10.3.13	Last mile rate and encapsulation parameters.....	569
10.3.14	Link failure detection.....	571
10.3.15	CoA support.....	571
10.3.16	Accounting.....	572
10.3.17	Filters and mirroring.....	572
10.3.18	PTA considerations.....	572
10.3.19	QoS considerations.....	572
10.3.19.1	Dual-pass.....	572
10.3.19.2	Traffic prioritization in LFI.....	573
10.3.19.3	Shaping based on the last mile wire rates.....	574
10.3.19.4	Downstream bandwidth management on egress port.....	575
10.3.20	Sub/sla-profile considerations.....	575
10.3.21	Example of MLPPPoX session setup flow.....	575
10.3.22	Other considerations.....	576
10.4	LNS support on ESA.....	577
10.5	Configuration notes.....	577
11	Standards and protocol support.....	579
11.1	Access Node Control Protocol (ANCP).....	579
11.2	Bidirectional Forwarding Detection (BFD).....	579
11.3	Border Gateway Protocol (BGP).....	579
11.4	Bridging and management.....	581
11.5	Broadband Network Gateway (BNG) Control and User Plane Separation (CUPS).....	582
11.6	Certificate management.....	582
11.7	Circuit emulation.....	583
11.8	Ethernet.....	583
11.9	Ethernet VPN (EVPN).....	583
11.10	gRPC Remote Procedure Calls (gRPC).....	584
11.11	Intermediate System to Intermediate System (IS-IS).....	584
11.12	Internet Protocol (IP) Fast Reroute (FRR).....	585
11.13	Internet Protocol (IP) general.....	585
11.14	Internet Protocol (IP) multicast.....	587

11.15	Internet Protocol (IP) version 4.....	588
11.16	Internet Protocol (IP) version 6.....	589
11.17	Internet Protocol Security (IPsec).....	590
11.18	Label Distribution Protocol (LDP).....	591
11.19	Layer Two Tunneling Protocol (L2TP) Network Server (LNS).....	592
11.20	Multiprotocol Label Switching (MPLS).....	592
11.21	Multiprotocol Label Switching - Transport Profile (MPLS-TP).....	593
11.22	Network Address Translation (NAT).....	593
11.23	Network Configuration Protocol (NETCONF).....	594
11.24	Open Shortest Path First (OSPF).....	594
11.25	OpenFlow.....	595
11.26	Path Computation Element Protocol (PCEP).....	595
11.27	Point-to-Point Protocol (PPP).....	595
11.28	Policy management and credit control.....	596
11.29	Pseudowire (PW).....	596
11.30	Quality of Service (QoS).....	597
11.31	Remote Authentication Dial In User Service (RADIUS).....	597
11.32	Resource Reservation Protocol - Traffic Engineering (RSVP-TE).....	597
11.33	Routing Information Protocol (RIP).....	598
11.34	Segment Routing (SR).....	598
11.35	Simple Network Management Protocol (SNMP).....	599
11.36	Timing.....	602
11.37	Two-Way Active Measurement Protocol (TWAMP).....	602
11.38	Virtual Private LAN Service (VPLS).....	602
11.39	Voice and video.....	603
11.40	Wireless Local Area Network (WLAN) gateway.....	603
11.41	Yet Another Next Generation (YANG).....	603
11.42	Yet Another Next Generation (YANG) OpenConfig Models.....	603

1 Getting started

1.1 About this guide

This guide describes details pertaining to Integrated Services Adapters (ISAs) and Extended Services Appliances (ESAs) and the services they provide. ISA may refer to ISA2 or an ESA-VM unless otherwise specified.

This guide is organized into functional chapters and provides concepts and descriptions of the implementation flow, as well as Command Line Interface (CLI) syntax and command usage.



Note: Unless otherwise indicated, this guide uses classic CLI command syntax and configuration examples.

The topics and commands described in this document apply to the:

- 7450 ESS
- 7750 SR
- Virtualized Service Router

For a list of unsupported features by platform and chassis, and for services supported by ISAs and ESAs, see the *SR OS R23.x.Rx Software Release Notes*, part number 3HE 19269 000 x TQZZA.

Command outputs shown in this guide are examples only; actual displays may differ depending on supported functionality and user configuration.



Note: The SR OS CLI trees and command descriptions can be found in the following guides:

- *7450 ESS, 7750 SR, 7950 XRS, and VSR Classic CLI Command Reference Guide*
- *7450 ESS, 7750 SR, 7950 XRS, and VSR Clear, Monitor, Show, and Tools Command Reference Guide* (for both MD-CLI and Classic CLI)
- *7450 ESS, 7750 SR, 7950 XRS, and VSR MD-CLI Command Reference Guide*



Note: This guide generically covers Release 23.x.Rx content and may contain some content that will be released in later maintenance loads. See the *SR OS R23.x.Rx Software Release Notes*, part number 3HE 19269 000 x TQZZA for information about features supported in each load of the Release 23.x.Rx software.

1.2 ISA and ESA configuration process

The ESA is specialized hardware that hosts ESA Virtual Machines (ESA-VMs). Each ESA-VM is configured as an integrated service type. ESA extends the proven Integrated Services Adapter (ISA) system implementation architecture and related control processing module (CPM) functions on the 7750 SR systems to include ESA-VM-based virtual ISA (v-ISA) functionality.

[Table 1: ISA and ESA configurations and supported features](#) lists the ISA and ESA configurations with descriptions and their supported functionalities.

Table 1: ISA and ESA configurations and supported features

Configuration	Description	Supported features
mda type: isa2-aa esa vm-type: aa	Application Assurance	<ul style="list-style-type: none"> Per-flow stateful deep packet inspection on OSI Layers 3 to 7 Release-independent protocol signatures, applications, application groups, and charging groups Flow attribute classification using deterministic and heuristic machine-learning algorithms Per-application and per-attribute policy enforcement and charging Layer 7 stateful firewall to block unsolicited traffic, with full application-level gateway (ALG) URL filtering and web classification based filtering Access network congestion detection and control using Dynamic Experience Management Application and protocol based volume and performance reporting
mda type: isa2-bb esa vm-type: bb	Broadband	<ul style="list-style-type: none"> CGN: LSN44, DS-Lite and NAT64 L2-Aware NAT for tight integration between BNG subscribers and NAT44 LNS WLAN-GW vRGW Generic re-assembly and TCP MSS-adjust
mda type: isa2-tunnel esa vm-type: tunnel	IP tunnels	<ul style="list-style-type: none"> IPsec tunnel: Secure network traffic on IP level for site-to-site, remote-access, mobile backhaul GRE/IP-in-IP tunnel: Overlay IP interface with transport as GRE/IP-in-IP tunnel L2TPv3 tunnel: Pseudowire for VPLS and Routed VPLS
mda type: isa2-video	Video	<ul style="list-style-type: none"> Fast Channel Change (FCC) Video Packet Retransmission (RET) Video Quality Monitoring (VQM) Perfect Stream Multicast (S,G) NAT

[Table 2: Configuration details](#) is a summary of the ISA and ESA virtual machines (ESA-VMs) guide structure by task. Specific configuration details for a software area, CLI syntax and command usage to configure parameters for each function are contained within each section.

Table 2: Configuration details

Area	Task	Section
Application Assurance	Configure Application Assurance entities	Configuring AA with CLI
IP tunnels	Determine IPsec deployment requirements	IPsec deployment requirements
	Configure IPsec	Configuring IPsec with CLI
L2TPV3 tunnels	Configure the L2TPV3 control plane	Control plane
	Configure public SAP	Public SAP
	Configure private SAP	Private SAP
Video services	Configure video services components	Configuring video service components with CLI
	Configure REF/FCC video components	Configuring RET/FCC video components with CLI
Network Address Translation	Configure NAT on SR	Configuring NAT
	Configure NAT on VSR	Configuring VSR-NAT
Residential firewall	Configure the residential firewall	Residential firewall provisioning
TCP MSS adjustment	Configure TCP MSS adjustments for BB	TCP MSS adjustment
	Configure TCP MSS adjustments for tunnel-ISA	TCP MSS adjustment
	Configure AA TCP MSS adjustments	AQP
L2TP network server	Configure subscriber aggregate rate limit on LNS	Subscriber agg-rate-limit on LNS
	Configure LNS reassembly	LNS reassembly
	Configure MLPPPoE and MLPPP(oE)oA with LFI on LNS	MLPPPoE, MLPPP(oE)oA with LFI on LNS

1.3 Conventions

This section describes the general conventions used in this guide.

1.3.1 Precautionary and information messages

The following information symbols are used in the documentation.



DANGER: Danger warns that the described activity or situation may result in serious personal injury or death. An electric shock hazard could exist. Before you begin work on this equipment, be aware of hazards involving electrical circuitry, be familiar with networking environments, and implement accident prevention procedures.



WARNING: Warning indicates that the described activity or situation may, or will, cause equipment damage, serious performance problems, or loss of data.



Caution: Caution indicates that the described activity or situation may reduce your component or system performance.



Note: Note provides additional operational information.



Tip: Tip provides suggestions for use or best practices.

1.3.2 Options or substeps in procedures and sequential workflows

Options in a procedure or a sequential workflow are indicated by a bulleted list. In the following example, at step 1, the user must perform the described action. At step 2, the user must perform one of the listed options to complete the step.

Example: Options in a procedure

1. User must perform this step.
2. This step offers three options. User must perform one option to complete this step.
 - This is one option.
 - This is another option.
 - This is yet another option.

Substeps in a procedure or a sequential workflow are indicated by letters. In the following example, at step 1, the user must perform the described action. At step 2, the user must perform two substeps (a. and b.) to complete the step.

Example: Substeps in a procedure

1. User must perform this step.
2. User must perform all substeps to complete this action.
 - a. This is one substep.
 - b. This is another substep.

2 ISA and ESA hardware

2.1 In this section

This section provides an overview of Nokia's implementation of the ISA and ESA hardware.



Note: Cards must be configured using the commands described in the *7450 ESS, 7750 SR, 7950 XRS, and VSR Interface Configuration Guide*.



Note:

The following conditions apply to the ISA and ESA hardware:

- ISAs (ISA2s alone or on MS-ISM cards) and ESA-VMs cannot be intermixed within the same ISA group. This limitation applies to all ISA group types.
- All ISA group types allow ESA-VMs to be hosted on ESAs of any hardware version.

2.2 MS-ISA2 overview

The MS-ISA2 (or ISA2-MS in CLI) is a second generation Integrated Services Adapter for multiservice processing, as a resource module within the router system providing packet buffering and packet processing.

The MS-ISA2 fits in an MDA or ISA slot on an IOM4-e and has no external ports, so all communication passes through the Input/Output Module (IOM), making use of the network processor complex on the host IOM for queuing and filtering functions like other MDAs and ISAs.

The actual ingress and egress throughput varies depending on the buffering and processing demands of a specific application, but the MS-ISA2 hardware can support 40 Gb/s of throughput processing. The processed rate (up to 40 Gb/s) is the sum of the upstream and downstream rates (for example, 10 Gb/s up and 30 Gb/s down, or 20 Gb/s up and 20 Gb/s down).

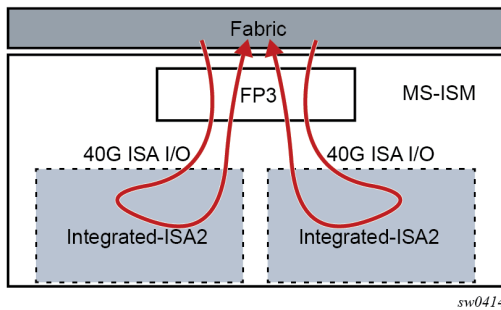
2.3 MS-ISM overview

The Multiservice Integrated Services Module (MS-ISM) card contains two ISA2 processing modules providing increased packet processing throughput and scale compared to the MS-ISA platform. Each ISA2 processing module supports a 40G datapath for packet processing; as with ISA1 the actual throughput varies by function. The processed rate (up to 40 Gb/s) is the sum of the upstream and downstream rates (for example, 10 Gb/s up and 30 Gb/s down, or 20 Gb/s up and 20 Gb/s down).

The IOM base card is an imm-2pac-fp3 with two embedded positions for ISA2s. Hot swap or field replacement of the ISA2s within an MS-ISM assembly is not supported. IMM cards offering 10x10GE media plus one ISA2, or 1x100GE media plus one ISA2.

The following shows the ISA2 processing modules in the MG-ISM card.

Figure 1: MS-ISM with ISA2s



The MS-ISA2 remains as a common base hardware assembly to be used as a generic CPU processing platform for multiple applications. The functions supported on the MS-ISA2 and MS-ISM include the following software based capabilities:

- Application Assurance (AA)
- Tunnel (IPsec, GRE)
- Broadband (NAT, LNS)
- Video (FCC, RET)

2.4 ESA overview

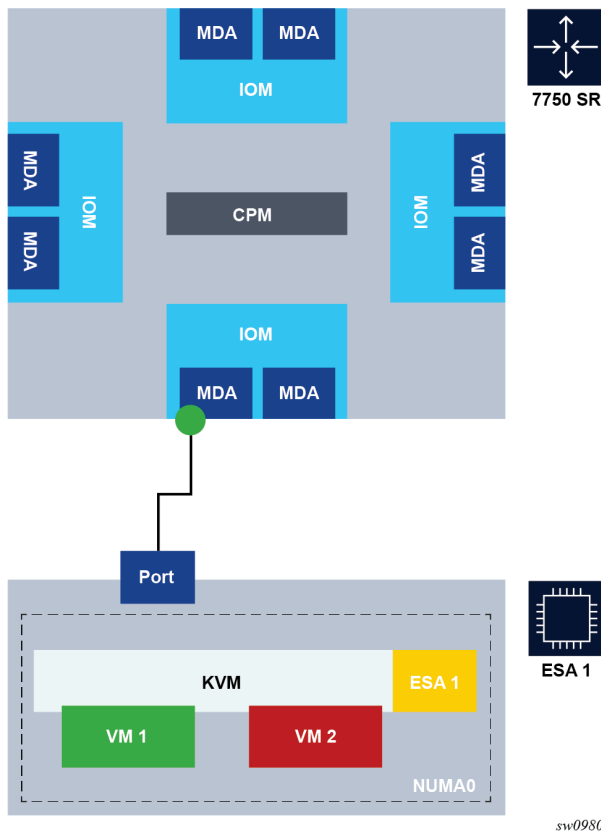
An Extended Services Appliance (ESA) is a server that attaches to a host 7750 SR over standard SR system interface ports, and which has one to four Virtual Machine (VM) instances to perform multiservice processing. The ESA provides packet buffering and processing and is logically part of the router system. The ESA 100G-2 includes one 20-core Intel fourth generation (Sapphire Rapids) 4416+ processor and 128 Gbytes of memory. The ESA 400G (Revision BA) includes two 32-core Intel fourth generation (Sapphire Rapids) 6438N processors and 512 Gbytes of memory.

The ESA processing rate is the sum of the upstream and downstream rates (for example, 80 Gb/s up and 20 Gb/s down, or 50 Gb/s up and 50 Gb/s down).

The ESA 100G-2 hardware can support up to 100 Gb/s of throughput processing, and the ESA 400G up to 400 Gb/s of processing. However, the maximum ESA ingress and egress throughput varies depending on the buffering and processing demands of a specific application.

The following shows an ESA connected to a 7750 SR.

Figure 2: ESA connection to 7750 SR



A direct local fiber connection must be used to connect an ESA port to a 7750 SR port. As with other MDAs and ISAs, all communication passes through the 7750 SR Input/Output Module (IOM), making use of the network processor complex on the host IOM for queuing and filtering functions.

The ESA 100G-2 includes one Mellanox Connect X6 2-port 100 Gb/s NIC with QSFP28 optics connectors. Each NIC has a maximum 200 Gb/s throughput per NIC, but the CPU capacity of the ESA 100G-2 limits the number of useful links to one 100 Gb/s port. Either of the two ESA NIC ports can be used to connect to the 7750 SR port.

The ESA 400G includes two Mellanox Connect X6 2-port 100 Gb/s NIC with QSFP28 optics connectors. Each NIC has a maximum 200 Gb/s throughput per NIC, and any of the four ESA NIC ports can be used to connect to the 7750 SR port.

The following SR to ESA port speeds are supported:

- 100GE (using QSFP28 optics in both the SR and ESA)
- 40GE (using QSFP+ optics in both the SR and ESA)
- 25GE (using a QSFP28 - SFP28/SFP+ Adapter and SFP28 optics in both SR and ESA)
- 10GE (using a QSFP28 - SFP28/SFP+ Adapter and SFP+ optics in both SR and ESA)

ESA 400G performance may be enhanced by configuring up to four ESA VMs for a single ESA across two CPUs. The two ESA NICs each connect to only one NUMA cell (CPU socket). For each ESA VM, reserve at least one port for SR interconnect. The most common ESA 400G deployment scenarios are as follows:

- one port and one ESA VM – one port per NIC and one ESA VM per CPU socket
- two ports and two ESA VMs – one port per NIC and one ESA VM per CPU socket to ensure maximum port and ESA VM performance
- four ports and four ESA VMs – two ports per NIC and two ESA VMs per CPU socket for maximum performance and density



Note: When using four ports and four ESA VMs, because each CPU socket is shared by two VMs, the throughput for each VM is slightly less than when one VM is used.

Ports for an ESA may be from the same or from different IOMs, XMAAs, or MDAs. Any combination of supported port speeds may be used on an ESA. If at least one host-port between the SR and the ESA is up, the ESA instance stays up.

An ESA-VM must be associated with one specific 7750 SR port. One physical 7750 SR port can be used by multiple VMs within an ESA. ESA-VMs may be configured as different types or the same type.

As each ESA-VM may only be associated with one 7750 SR port, LAG cannot be used between ports to an ESA. ESA-to-SR link resilience is handled by provisioning more VM instances than the processing requires (using the ISA group N+1 redundancy model). Functional sparing capacity is also handled by provisioning more VM instances than required.

Each ESA is managed by one 7750 SR. The ESA software (`hypervisors.tim` file, located on the active CPM from the 7750 SR host) can only be instantiated by a 7750 SR and cannot be instantiated in any other virtualized environment. Creation, configuration, deletion, resource allocation, and upgrade of a ESA-VM are all controlled by the 7750 SR CPM.

SR system LLDP must be enabled for ESA use, as LLDP is used to verify connectivity between the configured SR ESA host-ports and the matching configured ESA port for an ESA-VM. To set up an ESA in a 7750 SR system, complete the following actions in any order:

- Install the ESA hardware in a rack, then apply power to the ESA hardware.
- Connect the ESA hardware to a compatible 7750 SR chassis, IOM, or MDA using the appropriate optics.
- From the 7750 SR, configure ESA host and ESA-VM ports; see [Configuring an ESA with CLI](#).

See the 7750 SR ESA Chassis Installation Guide for more information about the first two items in the preceding list.



Note: After the ESA host-port is assigned, the port defaults are automatically modified. The new port defaults cannot be changed by the operator until the port is unassigned as an ESA host-port.

The ESA hardware is then booted by the 7750 SR CPM and available resources are discovered by the 7750 SR. ESA-VMs are configured as a type and size (number of cores and amount of memory). ESA-VM types include services that also run on ISAs, thereby providing a virtualized ISA function as an ESA-VM within the SR system and as part of an ISA group. An ISA group can only contain physical ISAs or ESA-VMs. Traffic for an ESA-VM enters the 7750 SR and is forwarded to the ESA-VM in a manner identical to that of a traditional ISA.

Multiple ESAs may be configured per IOM and per system as needed for scale.

ESA 100G-2 and 400G provide CLI, SNMP, and YANG support for the following hardware monitoring states:

- ESA health – unknown, OK, degraded, or critical
- PSU health – unknown, OK, degraded, or critical

- Fan redundancy – unknown, redundant, non-redundant, or failed-redundant
- Fan health – unknown, OK, degraded, or critical
- Power supply mismatch – true or false
- Power supply redundancy – unknown, redundant, non-redundant, or failed-redundant
- Temperature health – unknown, OK, degraded, or critical

ESA hardware monitoring events and states are integrated into the SR OS system facility alarms.

2.5 Application Assurance hardware features

2.5.1 AA system support

The Application Assurance Integrated Services Adapter (AA ISA) is a resource adapter, which means that there are no external interface ports on the AA ISA itself. Similarly, ESAs only do processing functions for traffic on the ESA interconnect ports to the SR system. Traffic on the SR system is forwarded to ISAs or ESA from any other IOMs on a system in which the AA ISA or ESA is installed, with a divert mechanism used to switch traffic internally to the AA ISA or ESA-VM.

See the SR OS R23.x.Rx Software Release Notes for information about the ESA platform support.

The following table describes Application Assurance support on the 7750 SR and 7450 ESS.

Table 3: AA system support

System	AA on MS-ISM	AA on MS-ISA2
7750 SR-12	Yes	Yes
7750 SR-12e	Yes	Yes
7750 SR-7	Yes	Yes
7750 SR-1e	No	Yes
7750 SR-2e	No	Yes
7750 SR-3e	No	Yes
7450 ESS-12	Yes	Yes
7450 ESS-7	Yes	Yes

2.5.2 Host IOM support for AA on ISAs

The AA MS-ISA2 is supported on IOM4-e, IOM4-e-B, IOM4-e-HS, and on 7750 SR-1e, 7750 SR-2e, and 7750 SR-3e (IOM-e). The MS-ISM versions contain one or two ISA2s embedded on a IMM card.

Each IOM can support a maximum of two AA ISA2 modules. To maximize AA ISA redundancy, deployment of AA ISAs on separate host IOMs is recommended as it provides IOM resilience. Traffic from any supported IOM (for example, IOM4-e, a fixed port IOM (IMM)) can be diverted to an AA ISA host IOM.

The MS-ISA2 is field replaceable and supports hot insertion and removal. An SR system can support up to 15 active ISA2s for AA, each providing up to 40 Gb/s processing and 600 Gb/s total per system.

AA ISA software upgrades are part of the ISSU functionality. Upgrades to AA ISA software, for example to activate new protocol signatures, do not impact the second MDA slot for the IOM carrying the AA ISA, nor do upgrades impact the router itself (for example a new AA ISA software image can be downloaded without a need to upgrade other software images).

2.5.3 Host IOM support for AA on ESAs

ESA port connectivity is supported on most FP3-based IOMs and all FP4-based (or later) cards. For a list of supported platforms or cards, contact your Nokia representative.

An SR system can support up to 15 active and one standby ESA-VMs for AA.

AA ESA-VM software upgrades are part of the ISSU functionality. Upgrades to AA software, for example to activate new protocol signatures, do not impact other ESA-VMs on the same ESA or on other traffic on the same IOM, nor do upgrades impact the router itself (for example, a new AA software image can be downloaded to an ESA-VM without a need to upgrade other software images).

The ESA version must match the build release version of the host IOM.

2.6 Configuring an ESA with CLI

This section provides information to configure an ESA using the command line interface from a 7750 SR. It is assumed that the user is familiar with the basic concepts of configuring policies.

2.6.1 Provisioning an ESA and ESA-VM

Use the following syntax to provision an ESA.



Note: Each ESA host-port and ESA-VM port must each be associated with a dedicated 7750 SR 100G port.

```
config>esa esa-id
  vm vm-id
  vm-type {aa | bb | tunnel | video}
```

The following example shows an ESA containing both a VM-type AA and a VM-type BB.

```
configure
  esa 1 create
    description "Esa for AA-BB"
    host-port 7/1/c6/1
    vm 1 create
      description "Application-Assurance ISA"
      vm-type aa
      host-port 7/1/c6/1
```

```

cores 12
memory 20
no shutdown
exit
vm 2 create
description "Broadband ISA"
vm-type bb
host-port 7/1/c6/1
cores 9
memory 40
no shutdown
exit
    
```

The following output displays an ESA and ESA-VM for the preceding configuration example.

```

show esa
=====
Extended Services Appliance Summary
=====
ESA  Description                               Admin  Oper
      State                               State
-----
1                                           up     up
=====

show esa detail
=====
ESA 1
=====
Description           : Esa for AA-BB
Admin State           : up
Operational State     : up
Oper flags            : none
IOM Host Port         : 7/1/c6/1
Hardware Data
  System manufacturer : Nokia Solutions and Networks
  System product name  : ESA-100G
  System serial number : QTFCT99040103
  Software Version     : TiMOS-H-19.10.S24 hypervisor/esa Copyright (c)
                        : 2000-2019 Nokia. All rights reserved. All use
                        : subject to applicable license agreements. Built
                        : on Wed Oct 23 20:35:01 PDT 2019 by builder in /
                        : builds/c/19105/S24/panos/hypervisors
Time of last boot     : 2019/10/24 14:49:58 UTC
Cores available       : 23
Cores allocated       : 23
Cores remaining       : 0
Memory available      : 192 GB
Memory allocated      : 60 GB
Memory remaining      : 132 GB
Performance enabled   : yes
Export restricted     : no
=====
    
```

3 Application Assurance

3.1 AA overview

Network operators are transforming broadband network infrastructures to accommodate unified architecture for IPTV, fixed and mobile voice services, business services, and High Speed Internet (HSI), all with a consistent, integrated awareness and policy capability for the applications using these services.

As bandwidth demand grows and application usage shifts, the network must provide consistent application performance that satisfies the end customer requirements for deterministic, managed quality of experience (QoE), according to the business objectives for each service and application. AA is the enabling network technology for this evolution in the service router operating system.

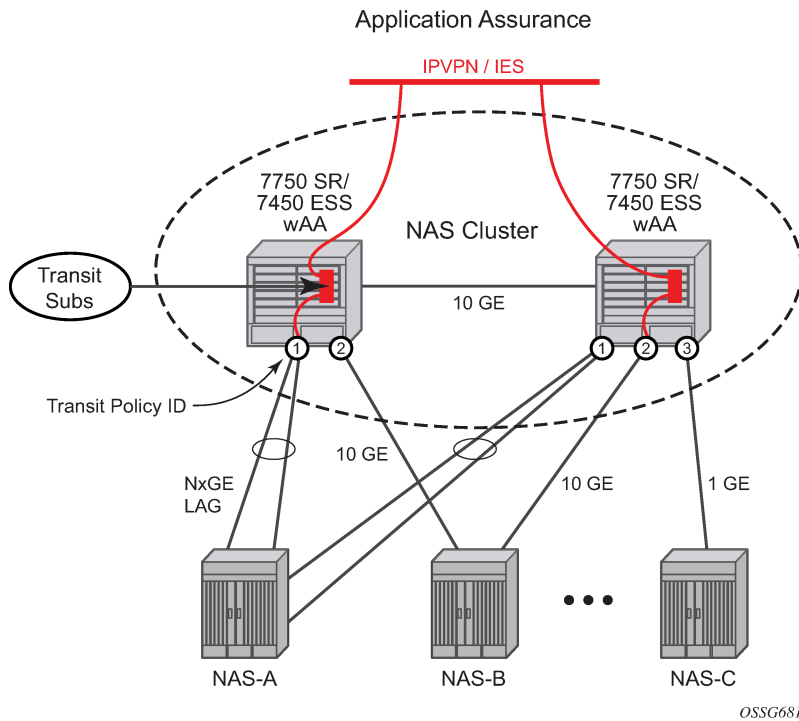
AA, coupled with subscriber or VPN access policy control points enables any broadband network to provide application-based services. For service providers, this unlocks:

- the opportunity for new revenue sources
- content control varieties of service
- control over network costs incurred by various uses of HSI
- complementary security aspects to the existing network security
- improved quality of service (QoS) sophistication and granularity of the network
- the ability to understand and apply policy control on the transactions traversing the network

3.1.1 AA: inline policy enforcement

The following shows AA ISA inline identification, classification, and control.

Figure 3: AA ISA inline identification, classification and control



The integrated solution approach for AA recognizes that a per-AA subscriber and per-service capable QoS infrastructure is a pre-condition for delivering application-aware QoS capabilities. Enabling per-application QoS in the context of individual subscriber's VPN access points maximizes the ability to monetize the application service, because a direct correlation can be made between customers paying for the service and the performance improvements obtained from it. By using an integrated solution there is no additional cost related to router port consumption, interconnect overhead or resilience to implement in-line application-aware policy enforcement.

3.1.2 AA integration in subscriber edge gateways

Multiple deployment models are supported for integrating AA in the various subscriber edge and VPN PE network topologies (Figure 4: AA deployment topologies). In all cases, AA can be added by in-service upgrade to the installed base of equipment instead of needing to deploy and integrate a whole new set of equipment and vendors into the network for Layer 4-7 awareness.

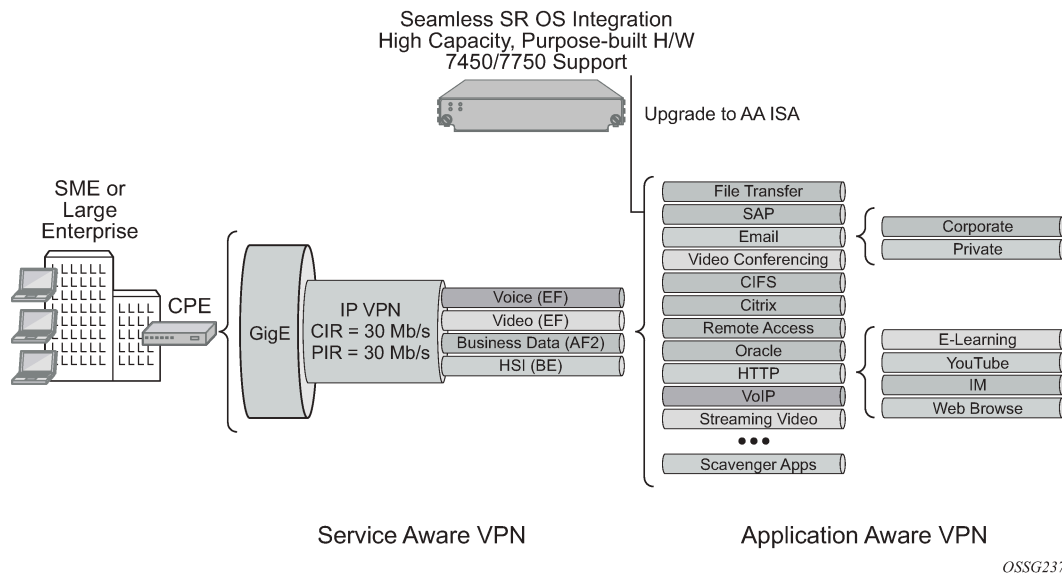
Integrating Layer 4-7 application policy with the 7750 SR or 7450 ESS subscriber edge policy context is the primary solution to address both residential broadband edge and Layer 2/Layer 3 application aware business VPN. Placement of Layer 4-7 analysis at the distributed subscriber edge policy point simplifies AA deployments in the following ways:

- For residential markets, CO-based deployment allows deployment-driven scaling of resources to the amount of bandwidth needed and the amount of subscribers requiring application-aware functionality.
- For AA business VPNs, a network deployment allows large scale application functionality at a VPN provider edge access point, vastly reducing complexity, cost, and time-to-market required to offer application-aware VPN services.

- Traffic asymmetry is avoided. Any subscriber traffic usually passes through one CO subscriber edge element so there is no need for flow paths to be recombined for stateful analysis.
- PE integration provides a single point of policy enforcement.
- SeGW integration provides firewall protection for NMS, MME and SGW.

The following shows AA deployment topologies.

Figure 4: AA deployment topologies



There are residential topologies where it is not possible or practical to distribute ISAs into the same network elements that run ESM, including for legacy edge BRASs that still need AA policy (reporting and control) for the same Internet services, and which needs to be aligned and consistent with the ESM AA policy. This is supported using transit AA subscribers, typically in the first routed element behind the legacy edge.

AA enables per AA subscriber (a residential subscriber, or a Layer 2 or Layer 3 SAP or spoke SDP), per application policy for all or a subset of AA subscriber's applications. This provides the ability to:

- implement Layer 4-7 identification of applications using a multitude of techniques from a simple port-based/IP address based identification to behavioral techniques used to identify, for example, encrypted or evasive applications
- when identified, apply a QoS policy on either an aggregate or a per-AA subscriber, per-application basis
- provide reports on the identification made, the traffic volume and performance of the applications, and policies implemented

An integrated AA module allows the SR and ESS product families to provide application-aware functions that previously required standalone devices (either in residential or business environment) at a fraction of the cost and operational complexity that additional devices in a network required.

A key benefit of integrating AA in the existing IP/MPLS network infrastructure (as opposed to an in-line appliance) is the ability to select traffic for treatment on a granular, reliable basis. Only traffic that requires AA treatment is simply and transparently diverted to the ISA. Other traffic from within the same service or interface follows the normal forwarding path across the fabric. In the case of ISA failure, ISA redundancy is

supported and in the case where no backup ISAs are available, the AA traffic reverts to the normal fabric matrix forwarding, also known as "fail to fabric".

[Table 4: Traffic diversion to the ISA](#) lists ISA traffic diversion information.

Table 4: Traffic diversion to the ISA

Deployment case	System divert ID	AA subscriber type	App-profile on:
Residential Edge (BNG)	ESM Sub-ID	ESM	ESM sub (All IPs, not per-host)
vRGW Bridged Residential Gateway (BRG) subscriber	ESM Sub-ID	ESM	ESM sub (All IPs, not per-host)
vRGW BRG session	ESM-MAC	ESM-MAC	ESM-MAC (by device, for any hosts assigned to a device)
Wireless LAN GW	ESM or DSM	ESM or DSM	ESM or DSM
Business Edge	L2/L3 SAP	SAP	SAP (Aggregate)
Residential Transit	Parent L3 SAP or spoke SDP	Transit AA	Transit Sub
Spoke Attached Edge	Spoke SDP	Spoke SDP	Spoke SDP (Aggregate)
SeGW	Parent SAP or spoke SDP or L2/L3 SAP	Transit AA SAP	Transit AA SAP

3.1.3 Fixed residential broadband services

Fixed residential HSI services as a single edge Broadband Network Gateway (BNG), virtualized Residential Gateway (vRGW), or as part of the Triple Play Service Delivery Architecture (TPSDA) are a primary focus of AA performance, subscriber and traffic scale.

To the service provider, application-based service management offers:

- application aware usage metering packages (quotas, 0-rating and so on)
- new revenue opportunities to increase ARPU (average revenue per user) (for gaming, peer-to-peer, Internet VoIP and streaming video, and so on)
- fairness (aligns usage of HSI network resources with revenue on a per-subscriber basis)
- operational visibility into the application usage, trends, and pressure points in the network

To the C/ASP, service offerings can be differentiated by improving the customer's on-line access experience. The subscriber can benefit from this by gaining a better application experience, while paying only for the value (applications) that they need and want.

3.1.3.1 DS-Lite

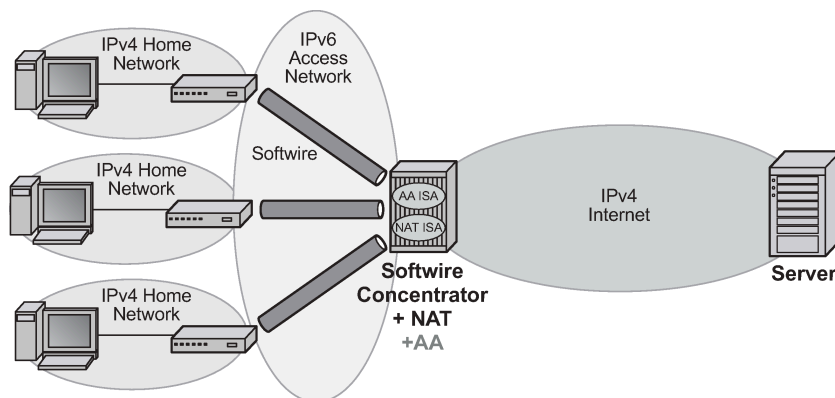
Dual-Stack Lite (DS-Lite) is an IPv6 transition technique that allows tunneling of IPv4 traffic across an IPv6-only network. Dual-stack IPv6 transition strategies allow service providers to offer IPv4 and IPv6 services and save on OPEX by allowing the use of a single IPv6 access network instead of running concurrent IPv6 and IPv4 access networks. DS-Lite has two components: the client in the customer network (the Basic Bridging BroadBand element (B4)) and an Address Family Transition Router (AFTR) deployed in the service provider network.

DS-Lite leverages a network address and port translation (NAPT) function in the service provider AFTR element to translate traffic tunneled from the private addresses in the home network into public addresses maintained by the service provider. On the 7750 SR, this is facilitated through the Carrier Grade NAT function.

When a customer's device sends an IPv4 packet to an external destination, DS-Lite encapsulates the IPv4 packet in an IPv6 packet for transport into the provider network. These IPv4-in-IPv6 tunnels are called softwires. Tunneling IPv4 over IPv6 is simpler than translation and eliminates performance and redundancy concerns.

The following figure shows the DS-lite deployment

Figure 5: DS-Lite deployment



al_0182

The IPv6 source address of the tunnel represents a unique subscriber. Only one tunnel per customer (although more is possible), but the IPv6 addresses cannot overlap between different customers. When encapsulated traffic reaches the softwire concentrator, the device treats the source-IP of the tunnel to represent a unique subscriber. The softwire concentrator performs IPv4 network address and port translation on the embedded packet by re-using Large Scale NAT and L2-Aware NAT concepts.

Advanced services are offered through AA multi service ISA to the DS-Lite connected customers. Subscribers' traffic (ESMs or transit-ip) are diverted to AA ISA for Layer 3 to Layer 7 identification and classifications, reporting and control based on the IPv4 packets (transported within the IPv6 DS-Lite tunnel). This AA classification, reporting and control of subscribers' traffic take effect before any NAT44 functions. In specific, AA sites on the subscriber side of NAT44.

The absence of a control protocol for the IP-in-IP tunnels simplifies the operational/management model, because any received IPv6 packet to the AA ISA can be identified as a DS-Lite tunneled packet if:

- Protocol 4 is indicated in the IPv6 header.
- The embedded IP packet is IPv4 (inside).

Fragmented IPv4 packets are supported only if tunneled through non-fragmented IPv6 packets.

Fragmentation at the IPv6 layer is not supported by AA ISA (when used to tunnel fragmented or non-fragmented IPv4 packets). These packets are cut-through with sub-default policy applied with a possibility of re-ordering.

If DSCP AQP action is applied to DS-Lite packet, both IPv4 and IPv6 headers are modified. AQP mirroring action is applied at the IPv6 layer. All collected statistics include the tunnel over-head bytes (also known as IPv6 header size).

3.1.3.2 6to4 /6RD

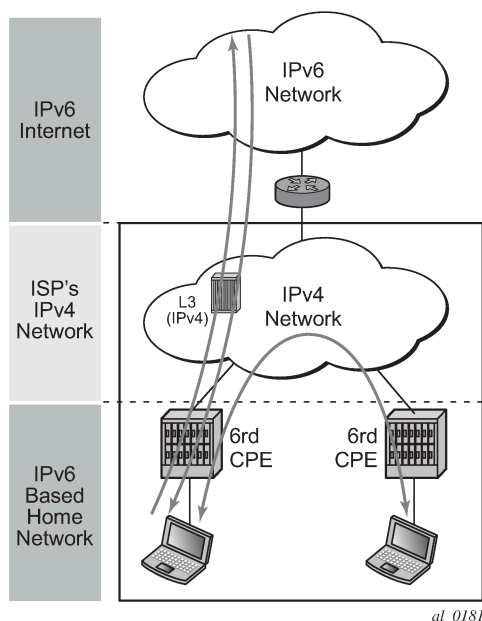
6RD/6to4 tunneling mechanism allows IPv6 sites to communicate over an IPv4 network without the need to configure explicit tunnels, as well as and for them to communicate with native IPv6 domains via relay routers. Effectively, 6RD/6to4 treats the wide area IPv4 network as a unicast point-to-point link layer. Both ends of the 6RD/6to4 tunnel are dual-stack routers. Because 6RD/6to4 does not build explicit tunnels, it scales better and is easier to manage after setup.

6to4 encapsulates an IPv6 packet in the payload portion of an IPv4 packet with protocol type 41. The IPv4 destination address for the encapsulating IPv4 packet header is derived from the IPv6 destination address of the inner packet (which is in the format of 6to4 address) by extracting the 32 bits immediately following the IPv6 destination address's 2002:: prefix. The IPv4 source address in the encapsulating packet header is the IPv4 address of the outgoing interface (not system IP address).

6RD is very similar to 6to4; the only difference is that the fixed 2002 used in 6to4 prefix is replaced by a configurable prefix.

The following figure shows an important deployment of 6RD/6to4 deployment is in access network.

Figure 6: 6to4 in access network deployment



To provide IPv6 services to subscribers, 6RD is deployed in these access networks to overcome the limitations of IPv4 only access network gear (for example, DSLAMs) with no dual stack support.

From an AA ISA point of view, deployment of 6RD in the access network is similar to that of the general deployment case between IPv6 islands with the added simplification that each 6RD tunnel carries traffic of a single subscriber.

When AA ISA sees an IPv4 packet with protocol type 41 and a payload that includes an IPv6 header, it detects that this is a 6rd/6to4 tunneled packet.

AA ISA detects, classifies, reports, and applies policies to 6rd/6to4 packet for ESM, SAP, spoke-SDP, and transit-IP (**ip-policy**) AA subscriber types.

Fragmented IPv6 are supported only if tunneled through non-fragmented IPv4 packets.

Fragmentation at the IPv4 layer is not supported by AA ISA (when used to tunnel fragmented or non-fragmented IPv6 packets). These packets are cut-through with sub-default policy applied with a possibility of re-ordering.

If the packet has IPv4 options then AA ISA does not look into the IPv6 header. The packet is classified as IPv4 "unknown TCP/UDP". Furthermore, TCP/UDP checksum error detection is only applied for lpipe and routed services.

If the DSCP AQP action is applied to 6RD6to4 packets, both IPv4 and IPv6 headers are modified. AQP mirroring action is applied at the IPv4 layer. All collected statistics include the tunnel over-head bytes, aka. IPv4 header size.

3.1.4 Wireless LAN gateway broadband services

AA enables a variety of use cases important for Wireless LAN Gateway deployments in residential, public WiFi or VPN wireless LAN services. These include:

- **HTTP redirect for subscriber authentication with HTTP allowlist**

This redirects all non-authenticated subscriber HTTP traffic to an authentication portal and blocks the rest of Internet access, but allows user access to specific allowed websites, download Apps and software needed to authenticate.

- **HTTP/HTTPS redirect by policy**

This is URL or application blocking or usage threshold notification. Redirects some or all subscriber HTTP/HTTPS traffic to an portal landing site based on static or dynamic policy. This can be done while not interrupting selected HTTP/HTTPS based services such as streaming video.

- **inline HTTP browser notification**

This provides messaging in the form of web banners, overlays, or http-redirection. This can always be enabled, One-time per sub at authentication (greeting message "Welcome to our WiFi Service"), one time per COA, or periodically.

- **ICAP for large scale URL filtering**

An ICAP client in AA interacts with offline ICAP URL filtering services for parental control or large denylists. This reduces cost as only URLs for specific flows are sent to server, instead of full inline traffic.

- **analytics**

This provides the operator insight into the following: Application and App-group volume usage by time of day/day of week, top subs, devices, and so on.

- **traffic control for fair use policy**

This prevents some users of the hotspot from consuming a disproportionate amount of resources by limiting to volume of such use across all subscribers as a traffic management tool, or on a per-subscriber basis.

- **stateful firewall**

This prevents unsolicited sessions from attacking devices.

- **web-service URL classification**

AA communicates with a web-service which offers URL categorization and provides parental control services. For the web-service model, AA receives the category of the URL and makes local policy decisions.

3.1.5 Application-aware business VPN services

AA for business services can be deployed at the Layer 2 or Layer 3 network provider edge (PE) policy enforcement point for the service or at Layer 2 aggregation policy enforcement point complimentary to the existing Layer 3 IP VPN PE. In a business environment, an AA subscriber represents a VPN access point. A typical business service can have a much larger average bandwidth rate than the residential service and is likely to have a smaller AA subscriber count than a residential deployment.

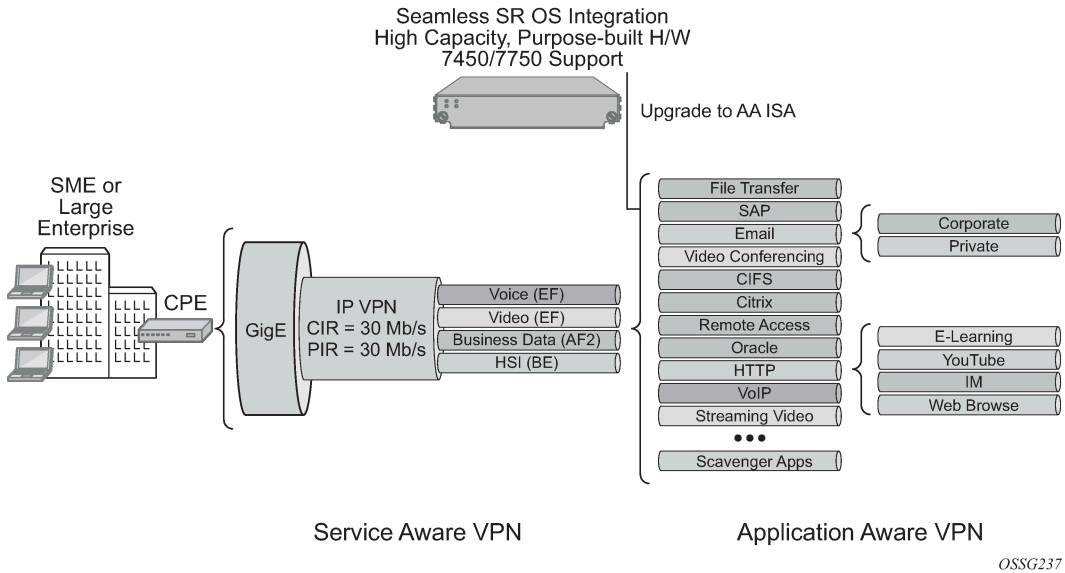
Multiple ISA2s can be deployed per PE, each incrementally processing up to 40 Gb/s. The in-network scalability is a key capability that allows a carrier to be able to grow the service bandwidth without AA throughput affecting the network architecture (more edge nodes, application-aware devices).

Application-aware Layer 2 and Layer 3 VPNs implemented using AA ISA equipped 7750 SR and 7450 ESS together with rich network management (NSP NFM-P, 5750 RAM, end customer application service portals) give operators a highly scalable, flexible, and cost effective integrated solution for application-based services to end customers. These services may include the following:

- rich application reporting with VPN, access site visibility
- right-sizing access pipes into a VPN service to improve/ensure application performance
- application-level QoS (policing, session admission, remarking, and so on) to ensure application-level performance, end-customer QoE objectives are met.
- value-added services such as application verification, new application detection, application mirroring
- performance reporting for real time (RTP) and non-real time (TCP) based applications
- Dual Stack IPv4 – IPv6 support
- GTP, 6RD tunneling support
- control unauthorized or recreational applications by site, by time of day.

The following figure shows AA BVS services integrated into the provider edge.

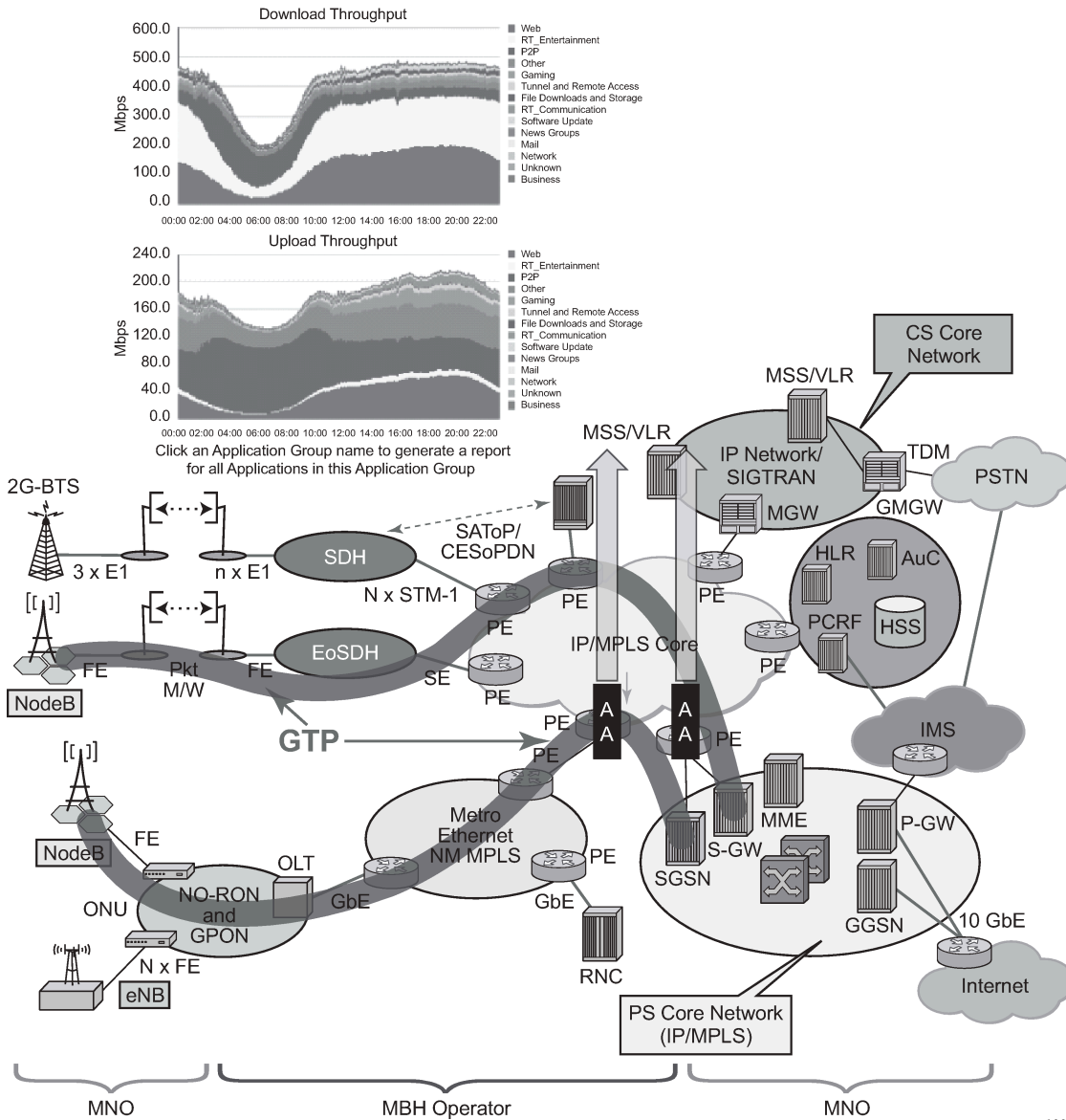
Figure 7: AA BVS services integrated into the provider edge



3.1.6 Mobile Backhaul

This section discusses Mobile Backhaul (MBH). The following figure shows a GTP-MBH AA deployment.

Figure 8: GTP-MBH AA deployment



In addition to SeGW FireWall deployments that require AA to support handling of GTP encapsulated traffic (S1-U interface), there are a number of deployments that require AA to support detection such as, classification and control of traffic encapsulated within GTP tunnels. These deployments are very similar in nature to AA support for other tunneling mechanisms such as 6RD, 6to4, DS-Lite. and so on For GTP tunnels, two main deployment use cases are identified: WiFi offload and mobile backhaul.

In Mobile Backhaul (MBH) deployment, operators provide business-GW network services called Mobile Data Roaming traffic service (that is, GPRS roaming exchange/service) to Mobile Network operators (MNOs) using MPLS network. MNOs, in turn, use MBH networks to create GTP tunnels across the MBH network between their mobile access network (for example, eNBs/SGSN/SGW) and PGSN/PGW network devices.

MNOs look into their MBH network providers to provide more analytical reporting of the applications running over the GTP-U tunnels.

AA-ISA is used to report on diverted business SAPs, regardless of how the traffic is encapsulated (GTP-U and 6RD, for example). From AA-ISA point of view, the diverted business SAP represents the subscriber. The subscriber is the MNO itself. No transit AA subscriber support is required in this deployment.

In this situation, multiple GTP-U tunnels are carried per SAP. AA reports on the actual content of these tunnels and not on the GTP-U tunnel themselves. For example, AA reports on the applications per SAP and not applications per GTP-u tunnel.

While this use case does not require any form of AA control functions, all AA actions/control functions can be used except for actions that require packet modifications (such as HTTP enrichments, HTTP redirect, remarking, DSCP Remark, HTTP notification).

3.1.7 Stateful firewall service

AA supports stateful firewall, which may be used for Gi firewall, GRX (GTP Roaming) firewall, or SeGW firewall deployed within a 7750 SR Security gateway in ultra-broadband access networks (3G/4G/Femto) and provides the operator with back-end core network security protection. For more information, see [AA firewall](#).

3.2 AA system architecture

3.2.1 AA ISA resource configuration

AA ISAs are flexible embedded, packet processing resource cards that require configuration such that services may be associated with the resources. This includes assigning ISAs to groups, optionally defining group partitions, and setting the redundancy model. Load balancing is affected by how ISAs are grouped.

3.2.1.1 AA ISA groups

An AA ISA group allows operators to group multiple AA ISAs into one of several logical groups for consistent management of AA resources and policies across multiple AA ISA cards configured for that group. The following operations can be performed at the group level:

- Define one or multiple AA ISA groups to allow AA resource partitioning/reservation for different types of AA service.
- Define the AA subscriber scale mode for the group. Various scales for residential, VPN, and lightweight-Internet (used for WLAN-GW distributed subscriber management) modes are supported.
- Assign physical AA ISAs to a group.
- Select forwarding classes to be diverted for inspection by the AA subscribers belonging to the group and select the AA policy to be applied to the group.
- Configure redundancy and bypass mode features to protect against equipment failure.
- Configure QoS on IOMs which host AA ISAs for traffic toward AA ISAs and from AA ISAs.
- Configure ISA capacity planning using low and high thresholds.

- Enable partitions of a group.
- Configure the ISA traffic overload behavior for the group to either back pressure to the host IOM (resulting in possible network QoS-based discards) or to cut-through packets through the ISA without full AA processing. Cut-through is typically enabled for AA VPN groups but not for residential groups.

Residential services is an example where all AA services may be configured as part of a single group encompassing all AA ISAs, for operator-defined AA service. This provides management of common applications and reporting for all subscribers and services, with common or per customer AQP (using ASOs characteristics to divide AA group's AQP into per app-profile QoS policies).

Multiple groups can be further used to create separate services based on different sets of common applications, different traffic divert needs (such as for capacity planning) or different redundancy models. Multiple groups may be used:

- when there is a mix of residential and business customers
- among different business VPN verticals
- for business services with a common template base but different levels of redundancy, different FC divert, or scaling over what is supported per single group
- when system level status statistics have AA ISA group/partition scope of visibility

3.2.1.1.1 AA group partitions

VPN-specific AA services are enabled using operator defined partitions of an AA Group into AA policy partitions, typically with one partition for each VPN-specific AA service. The partition allows VPN specific custom protocols/application/application group definition, VPN specific policy definition and VPN specific reporting (some VPNs with volume-only reports, while others with volume and performance reports). Each partition's policy can be again divided into multiple application QoS policies using ASOs.

The use of ISA groups and partitions also improves scaling of policies, as needed with VPN-specific AA policies.

If partitions are not defined, all of the AA group acts as a single partition. When partitions are configured, application identification, policy and statistics configuration applies only to the specific partition and not any other partitions configured under the same AA group.

The definition of application profiles (and related ASO characteristics/values) are within the context of a specific partition (however, application profiles names must have node-wide uniqueness).

The definition of applications, application groups and AQP are also specific to a partition. This allows:

- the definition of unique applications and app-groups per partition
- the definition of AQP policy per partition
- the definition of common applications and app-groups per partition with per partition processing and accounting

Partitions also enable accounting/reporting customization for every AA subscriber associated with a partition, for example:

- the ability to define different types of reporting/accounting policies for different partitions in a single AA group, such as uniquely defining which application, protocols, app groups are being reported on for every AA subscriber that uses a given partition.
- AA group level protocol statistics with partition visibility (for example, protocol counts reported for each partition of the group separately)

The system provides independent editing and committing of each partition config (separate begin/commit/abort commands).

Policer templates allow group-wide policing, and can be referenced by partition policies.

3.2.1.1.2 Bypass modes

If no active AA ISA is available (for example, because of an operational failure, misconfiguration) the default behavior is to forward traffic as if no AA was configured, the system does not send traffic to the AA ISA (equivalent to fail to closed). Alarms are raised to flag this state externally. There is an optional "fail to open" feature where AA ISA service traffic is dropped if no active AA ISA is present (such as no AA ISA is present and operationally up).

3.2.1.2 Redundancy

AA ISA group redundancy is supported, to protect against card failure and to minimize service interruption during maintenance or protocol signature upgrades.

3.2.1.2.1 No AA ISA group redundancy

AA can be configured with no ISA redundancy within the AA group. All AA ISAs are configured as primary with no backup (up to the limit of active AA ISAs per node). There is a no fault state indicating that a spare AA ISA is missing. If a primary is configured but not active, there is a "**no aa-isa**" fault.

3.2.1.2.2 Failure to fabric

If no ISA redundancy is deployed or insufficient ISAs are available for needed sparing, the system implements "failure to fabric". When the ISA status shows the ISA is not available and there is no redundant ISA available, the ingress IOMs simply do not divert the packets that would have been sent to that ISA, but instead these proceed to the next hop directly across the fabric. When the ISA becomes available, the divert eligible packets resume divert through the ISA. This behavior is completely internal to the system, without affecting the forwarding or routing configuration and behavior of the node or the network.

3.2.1.2.3 N+1 AA ISA card warm redundancy

The system supports N+1 AA ISA equipment warm redundancy (N primary and 1 backup). If a backup is configured and there is no ISA available (a primary and backup failed), there is a "**no aa-isa**" fault. The backup AA ISA is pre-configured with `isa-aa.tim` and the group policies. Data path traffic is only sent to active AA ISAs, so the backup has no flow state. If a backup ISA is unavailable, there is a "backup missing" fault.

An AA subscriber is created and assigned to a primary AA ISA when an application profile is assigned to a subscriber, SAP, or spoke SDP. By default, AA subscribers are balanced across all configured primary AA ISAs.

Upon failure of a primary AA ISA, all of its AA subscribers and their traffic are operationally moved to the newly active backup AA ISA but the current flow states are lost (warm redundancy). The new AA ISA identifies any session-based active flows at a time of switchover as an existing protocol, while the other

flows are re-identified. The existing protocol-based application filters can be defined to ensure service hot redundancy for a subset of applications. After the backup AA ISA has taken control, it waits for operator control to revert activity to the failed primary AA ISA module.

The user can disable a primary AA ISA for maintenance by triggering a controlled AA ISA activity switch to do the AA ISA software field upgrade (a shutdown of an active AA ISA is recommended to trigger an activity switch).

The activity switch experiences the following AA service impact:

- All flow states for the primary ISA are lost, but existing flows can be handled with special AQP rules for the existing flows by the newly active backup AA ISA until sessions end.
- All statistics gathered on the active AA ISA since the last interval information that was sent to the CPM is lost.

3.2.1.3 ISA load balancing

Capacity-cost based load balancing allows a cost to be assigned to diverted AA subscribers (by the app-profile). Load Balancing uses the total allocated costs on a per-ISA basis to assign the subscriber to the lowest sum cost ISA resource. Each ISA supports a threshold as the summed cost value that notifies the operator if or when capacity planning has been exceeded.

The load balancing decision is made based on the AA capacity cost of an AA subscriber. The capacity cost is configured against the app-profile. When assigning a new diverted AA subscriber to an ISA, the ISA with the lowest summed cost (that also has sufficient resources) is chosen. Examples of different load-balancing approaches that may be implemented using this flexible model include:

- **AA subscriber count balancing**
Configure the capacity cost for each app-profile to the same number (for example, 1).
- **AA subscriber stats resource balancing**
Configure the capacity cost to the number of stats collected for AA subscribers using the app-profile. This may be used if different partitions have significantly different stats requirements.
- **bandwidth balancing**
Configure the capacity cost to the total bandwidth in both directions (in kb/s) expected for those AA subscribers. This may be used if different AA subscribers have highly varying bandwidth needs.

Load balancing operates across ISAs within an AA group, and does not balance across groups. The system ensures that app-profiles assigned to AA subscribers (ESM subscribers, SAPs and spoke SDPs) that are within a single VPLS, Epipe, IES, and VPRN services are all part of the same AA group (partitions within an AA group are not checked or relevant).

Users can replace the app-profile assigned to an AA subscriber with another app-profile (from the same group/partition) that has a different capacity cost.

Regardless of the preferred choice of ISA, the system takes into account the following:

- Resource counts have per-ISA limits. If exceeded on the ISA of choice, that ISA cannot be used and the next best is chosen.
- Divert IOM service queuing resources may limit load-balancing. If queuing resources are exhausted, the system attempts to assign the AA subscriber to the ISA where the first AA subscriber within that service (VPLS or Epipe) or service type (IES or VPRN) was allocated.

For prefix transit AA subscriber deployments using the **remote-site** command, traffic for the remote transit subs are processed a second time. The ISA used by the parent AA subscriber is used by all transits within the parent. In remote-site cases there may be a need to increase capacity cost of parent because the transits stay on same ISA as the parent.

Prefix transit AA subscribers are all diverted to the same Group and partition as the parent SAP.

3.2.1.4 Asymmetry removal

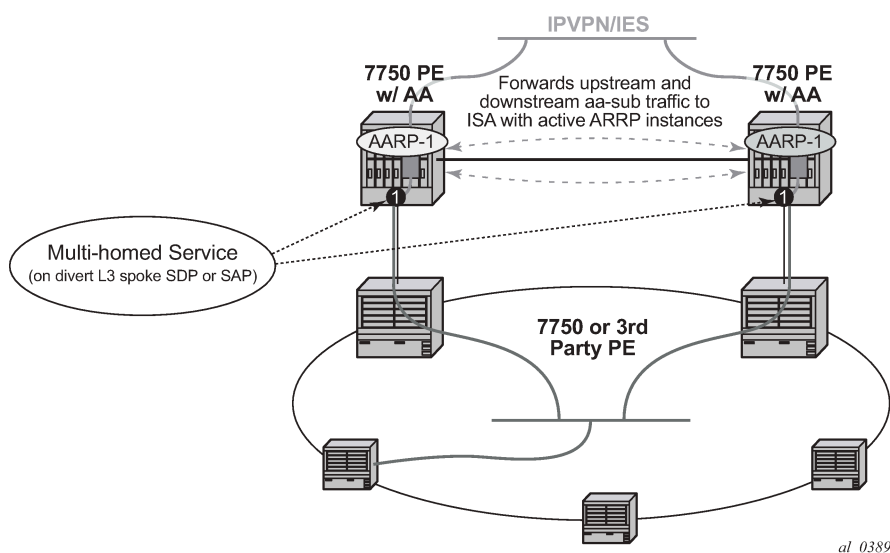
Asymmetry removal is only supported on 7750 SR routed services. Asymmetry removal is a means of eliminating traffic asymmetry between a set of multi-homed SAP or spoke SDP endpoints. This can be across endpoints within a single node or across a pair of inter-chassis link connected routers. Asymmetry means that the two directions of traffic for a flow (to-sub and from-sub) take different paths through the network. Asymmetry removal ensures that all packets for each flow, and all flows for each AA subscriber are diverted to an AA ISA.

Traffic asymmetry is created when there are multi-homed links for a service, and the links are simultaneously carrying traffic. In this scenario packets for flows use any reachable paths, therefore creating dynamic and changing asymmetry. Single node or multi-chassis asymmetry removal is used for any case where traffic for an AA subscriber may be forwarded over diverse paths on active AA divert links in a multi-homed topology. This includes support for SAP or spoke AA subscribers as well as business and residential transit AA subscribers within the diverted service.

Asymmetry removal must be implemented in the first routed hop on the network side of the subscriber management point, such that there is a deterministic and fixed SAP and spoke SDP association between the downstream subscriber management the parent divert service.

Asymmetry removal allows support for the SAP or spoke SDPs to the downstream element to be multi-homed on active links to redundant PE AA nodes as shown in [Figure 9: Transit sub SAP and spoke SDP multi-homing with asymmetry](#).

Figure 9: Transit sub SAP and spoke SDP multi-homing with asymmetry

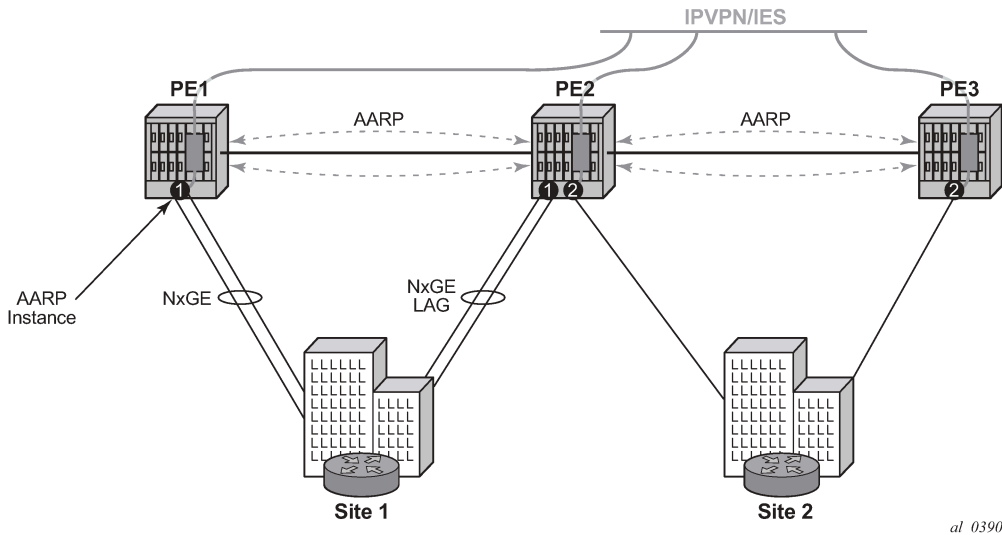


AA for transit-ip subscribers is commonly deployed behind the point of the subscriber policy edge after aggregation. This includes cases where AA needed is behind:

- any node running ESM but where there is no need or space to deploy distributed AA ISAs
- legacy BRAS that do not support integrated application policy

Asymmetry removal also allows a VPN site (Figure 10: VPN site multi-homing with asymmetry) to be connected with multi-homed, dual-active links while offering AA services with the ISA.

Figure 10: VPN site multi-homing with asymmetry



Asymmetry removal is supported for Layer 3 AA divert services:

- IES SAP and spoke SDP
- VPRN SAP and spoke SDP

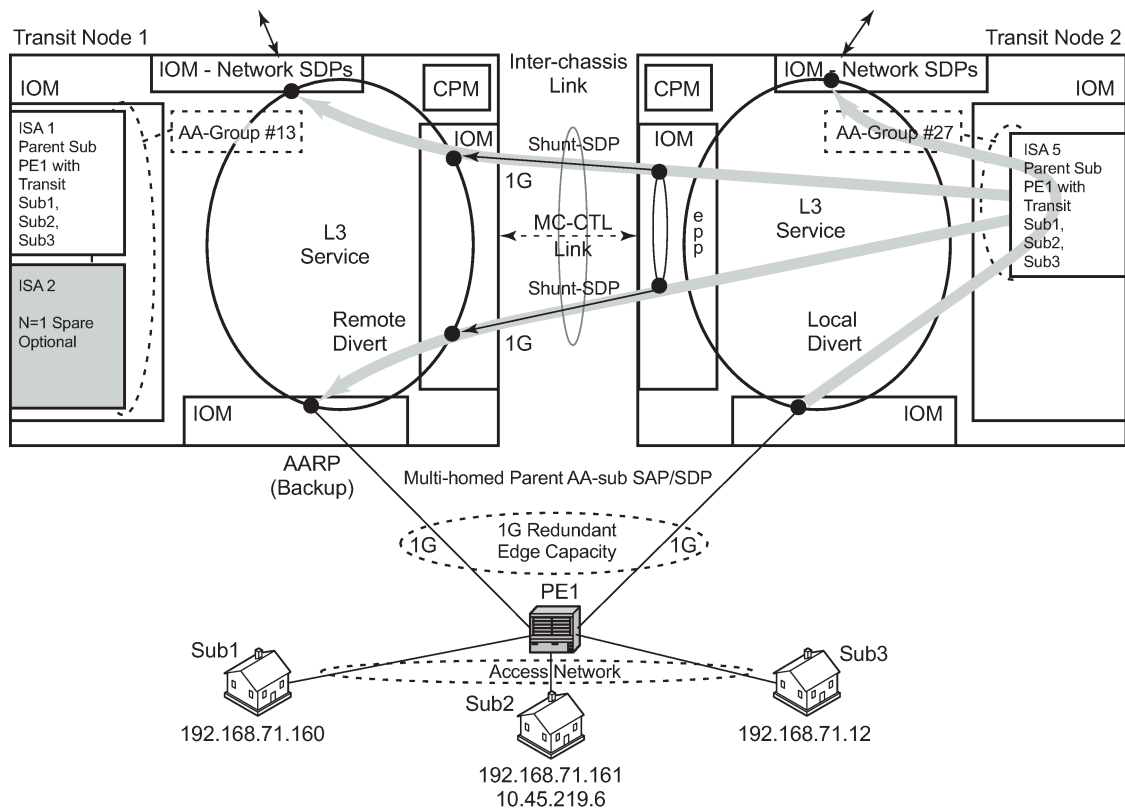
When asymmetry exists between multi-chassis redundant systems, Ipipe spoke SDPs are used to interconnect these services between peer nodes over an Inter-Chassis Link (ICL).

Asymmetry removal supports multiple endpoints of a service with a common AARP instance, with a primary endpoint assigned the app-profile (and transit policy for transit subs). The remaining endpoints are defined as secondary endpoints of the AARP instance. All SAP or spoke endpoints within a specific AARP instance are load balanced to the same ISA in that node. Multi-endpoint AARP instances allow single-node asymmetry removal and multi-chassis asymmetry removal with multiple active links per node.

3.2.1.4.1 Asymmetry removal overview

Figure 11: Multi-chassis asymmetry removal functional overview shows an overview of the multi-chassis asymmetry removal functionality.

Figure 11: Multi-chassis asymmetry removal functional overview



For a multi-homed parent AA subscriber, the parent SAP/SDP that is Active/Inactive per chassis is agreed by the inter-chassis AA Redundancy Protocol (AARP). For single node multi-homed endpoints, the AARP state is determined within a single node, as described later in the AARP operational states section.

- Divert AA subscribers are cost-based load balanced across ISAs in each chassis/AA group (node-local decision).
- Divert AA subscriber multi-homed pairing is supported by AA Redundancy Protocol (AARP) over inter-chassis link.
 - The same AARP ID is assigned to the divert SAP in both nodes.
 - AARP state in one node is master when all AARP conditions are met.
 - Packets arriving on node with the master AARP ID divert locally to ISA.
 - From sub packets on a node with backup AARP ID remote diverted over the subscriber side shunt, appearing to the ISA as if it was a local packet from the AA subscriber and returned to the network side interface spoke SDP shunt after ISA processing. To-sub packets on node with backup AARP ID remote divert over the network side shunt, appearing to the ISA as if it was a local network side divert packet for the AA subscriber, then returned to the subscriber side interface spoke SDP shunt after ISA processing.
 - All packets are returned to the original node for system egress (sent back over the inter-chassis shunts).
- If ISA N+1 sparing is available in a node, ISA sparing activates before AARP activity switch.

- Supports asymmetry for business SAPs and spoke SDPs, with or without transit AA subs.
- The AARP master-selection-mode is in minimize-switches mode by default, which is non-revertive and does not factor endpoint status. This can be configured per AARP instance using the master-selection-mode. The priority-rebalance configuration rebalances based on priority after the master failure condition is repaired. The inter-chassis-efficiency mode does priority based rebalance and includes the EP status for cases where an AARP activity switch is preferred to sustained ICL traffic load (when peer nodes are geographically remote).

3.2.1.4.2 Failure modes

Failure modes include the following:

- **AARP infrastructure failure including shunts**

For AARP to remove asymmetry, the AARP link must be synchronized between peers and all components of the shunts (Ipipe shunts and interface shunts) must be up and operational. If any of those components has failed, each AARP ID operates as standalone and diverts locally. Asymmetry is not removed.

- **failure of one of the interfaces to the dual homed site**

Routing moves all traffic to the remaining link/node if this is the master AARP peer node no action is required. For any traffic the backup node, inter-chassis shunting is used. There is no change to the AARP master/backup state. Traffic is still processed by the same ISA as before the failure.

- **network reachability fails to master AARP node**

AARP node loses reachability on the network side. This does not trigger an AARP activity switch, the shunt is used to move traffic from the backup node to the master node for the duration of the reachability issue. Routing should take care of traffic reconvergence. However, if the peer AARP is also not reachable, both nodes go on standalone mode and there is no asymmetry removal.

- **master AA ISA failure**

AARP activity flips for all the master AARP instances linked to this local ISA if there is no local spare available. Any traffic arriving on the node with the failed ISA uses the shunt to reach the master ISA.

3.2.1.4.3 AARP peered node/instance configuration

The multi-homed diverted AA subscriber in each peer node must be configured with the following parameters set in each node of the peer pair as displayed in the following table.

Table 5: Parameter values for peer nodes

Parameter	Value
Service ID	Node specific
Interface	Node specific
SAP or spoke/SDP ID	Node specific
AA-group ID	Node specific

Parameter	Value
App-profile name	Content must be the same in both peers as to not affect behavior. Nokia recommends using same name and content.
Transit policy ID	Same in both (only applies if transits are used)
AARP ID	Same in both
shunt-sdp <i>sdp-id:vc-id</i>	Node specific but must properly cross-connect the local AA subscriber service with the peer lpipe/service shunt interface to operate properly for asymmetry removal for remote divert traffic. Peer AARPs stay in standalone mode until cross-connect is configured properly.
Master-selection mode	Same in both
Other ISA-AA group configuration	Same in both, including fail-to, divert FC, and so on
IOM traffic classification into a FC	Same in both (can affect AA divert because this is conditioned by the FC). This includes sub side, network side and shunt IOMs.

AARP operation has the following required dependencies:

- For multi-chassis, shunt links are configured and operationally Up.
- For multi-chassis, peer communications established.
- Dual-homed SAP or spoke configured.
- app-profile configured against SAP or spoke with divert (making the subscriber an AA subscriber). This endpoint is called the primary endpoint if more than one endpoint is configured for an AARP instance.
- All endpoints within an AARP instance must be of the same type (SAP or spoke).
- All endpoints with an AARP instance must be within the same service.

3.2.1.4.4 MC-CTL

A multi-chassis control link (MC-CTL) is automatically established between peer AARP instances to exchange configuration and status information. Information exchanged includes configured service, protecting SAP, spoke, redundant-interface name, shunt-sdp, app-profile, priority and operational states.

AARP requires configuration of the peer IPv4 system address to establish a session between the two node's system IPv4 addresses.

3.2.1.4.5 Multi-chassis datapath shunts

When traffic needs to be remotely diverted it flows over shunts that are provisioned as *sdp-id:vc-id* between the dual-homed AA subscriber local service and a remote vc-switching lpipe.

- **subscriber to network direction**

The traffic is either handled locally (diverted to a local ISA when the AARP state is Master) or at the peer 7750 SR (redirect over the shunt Ipipe when the local AARP state is Backup or Remote). When traffic arrives on the subscriber side spoke SDP of the shunt-Ipipe, the system uses the AARP ID of the Ipipe to associate with an app-profile, therefore triggering Ipipe divert. It is diverted to the same ISA used to service the dual-homed SAP or spoke SDP. The ISA then treats this traffic the same as if it was received locally on the dual-homed SAP or spoke SDP context. After ISA processing, the traffic returns on the network side of the Ipipe to the peer. When the traffic returns to the original 7750 SR, the shunt Ipipe terminates into the routed service and it makes a new routing decision.

- **network to subscriber direction**

The traffic is either handled locally (diverted to a local ISA when the AARP state is Master) or at the peer 7750 SR (remote divert over the shunt Ipipe when the local AARP state is Backup or Remote). When traffic arrives on the shunt Ipipe from the peer with an AARP ID and associated app-profile, it is diverted through AA on the way to the subscriber-side spoke SDP. After AA processing, the traffic returns on the subscriber side of the Ipipe to the peer. When the traffic returns to the original 7750 SR, the shunt Ipipe terminates into the routed service and it makes a new routing decision to go out the dual-homed SAP or spoke SDP.

- **AARP operational states**

In single node operation, there are two operational states, Master or Standalone. A single node AARP instance is Master when all these conditions are met, otherwise AARP is in the standalone state with no asymmetry removal occurring:

- Dual-homed (primary) and dual-homed-secondary endpoints are configured.
- Divert capability is up.
- App profile is diverting.
- AA subscriber is configured.

With multi-chassis operation there are 4 operational states for an AARP instance: master, backup, remote and standalone:

- **master**

In multi-chassis operation, an AARP instance can only become operationally master when the inter-chassis link datapath is operational and the control path is or was up, the received peer node status indicating the peer's AARP instance and similar dependencies is or was up, and the AARP priority is higher than the peer. When the priority is equal then higher system interface IP address is used as a tiebreak.

The master state is immediately switched to remote for AARP related failures that result in the instance being not ready. ICL datapath shunt SDP failures cause the peer AARP go standalone. A shunt/Ipipe SDP failure is determined by the failure detection protocol used (BFD on routes, keep-alive on SDPs, LDP/RSVP, and so on).

When a SAP or spoke SDP with an AARP instance is shut down nothing changes for AARP, as packets can still use the AARP interface. When the SAP or spoke SDP is deleted, AARP is disassociated from the spoke SDP or SAP before deleting. The AARP instance still exists after deleting the SAP or spoke but without an association to an AA subscriber, the AARP state goes standalone.

- **backup**

Backup is the AARP state when all required conditions of the AARP instance are met except the master/backup priority evaluation.

- **remote**

When an AARP instance is operating with remote divert set for the protecting SAP or spoke AA subscriber. The peer AARP instance is the Master, there is no backup as the local system is not ready. This state is entered as a result of a failure in a local resource on the AARP instance, which triggers the divert traffic to the remote peer, such as a ISA failure without ISA backup). AA subscriber traffic is diverted over shunts to the peer.

- **standalone**

AARP is not operational between the multi-chassis pair, with AA operating with local AA divert to the ISAs within that node. There is no master or backup. This is the starting initial state for the AARP instance, or as a result of a failure in a dependent ICL resource (MC-CTL communication link or shut down).

An AARP instance activity switch is when one node moves from master to remote or backup mode, with the peer node becoming master. This can occur on a per-instance basis using the re-evaluate tool, or for all instances on an ISA that fails. On an AARP activity switch, AA divert changes from local to remote or remote to local, such that any packet is not seen by both nodes, resulting in no missed packet counts or double counts against the AA subscriber.

AARP activity is non-revertive to maximize the ID accuracy of flows. When an AARP instance toggles activity, packets are diverted to the newly active divert ISA and are processed as new flows, which for mid-session flows, often results in "unknown" traffic ID until those flows terminate. When the condition that triggered the AARP activity switch is resolved and the instance remains in backup state to not cause an additional application ID impacting event. This is consistent with AA N+1 ISA activity switches.

Because AA ISA availability is a criteria for AARP switches, any ISA failure or shutdown moves all AARP instance activity to ISAs in the master peer nodes, such as during software upgrades of ISAs. Depending on the nature of the failure or sequence of an upgrade procedure, all AA traffic may be processed by ISAs in one of the peers with no traffic being processed by ISAs on the other node.

If rebalancing the ISA load between the peer nodes is necessary, use the following command option to rerun AARP activity evaluation on a per-ISA basis to determine the master and backup AARP instances based on configured priority.

```
tools perform application-assurance aarp force-evaluate
```

The following table shows the interaction and dependencies between AARP states between a local node and its peer.

Table 6: Interaction and dependencies between AARP states

Local AARP operation state	Peer AARP operational state	Description
Master	Backup	<p>Inter-Chassis Link (ICL) Communication established between AARP peers.</p> <p>AARP dependent resources are up (to-sub/from-sub shunt, aarp control link, dual-homed SAP or spoke SDP).</p> <p>AARP instances have negotiated initial state assignment using configured priority/system IP address.</p> <p>AA service is available for the dual-homed SAP or spoke subscriber.</p>

Local AARP operation state	Peer AARP operational state	Description
		All to-sub/from-sub traffic specific to the dual-homed SAP or spoke SDP is serviced on the local node. Peer node is available to takeover in the event of a AA service failure on the local node. Asymmetry is removed for the dual-homed SAP or spoke subscribers, serviced by AA on the local (master) node.
Master	Remote	Same as Master/Backup except: AA service is available on the local node. AA service is unavailable on the peer node.
Standalone	Standalone	Initial state of the AARP instances upon creation or a result of a failure in any of the AARP dependent resources. All to-sub/from-sub traffic for the dual-homed SAP or spoke is serviced on each node independently. aarp instance operational state is outOfService on both sides. Asymmetry is not removed for the dual-homed SAP or spoke subscribers (traffic ID is not optimal).

3.2.1.5 ISA overload detection

Capacity cost resource counting does not have a hard per-ISA limit, because the cost values are decoupled from actual ISA resources. However, the value of the total summed cost per-ISA can be reported, and a threshold value can be set which raises an event when exceeded.

ISA capacity overload detection and events are supported within the system resource monitoring / logging capabilities if the traffic and resource load crosses the following high and low load thresholds on a per-ISA basis:

- ISA capacity cost
- flow table consumption (number of allocated flows)
- flow setup rate
- traffic volume
- host IOM egress weighted average shared buffer pool use (within the egress QoS configuration for each group). These thresholds are also used for overload cut-through processing.

While an app-profile is assigned to AA subscribers, the capacity-cost for that app-profile can be modified. The system makes updates in terms of the load balancing summary, but this does not trigger a re-balance.

In the absence of user configuration, the app-profile default capacity cost is 1. The range for capacity cost is 1 to 65535 (for example, for bandwidth based balancing the value 100 could represent 100 kb/s). 0 is an invalid value.

If the re-balancing of AA subscribers is required (for instance after the addition of new ISAs), there is a **tools** command to re-balance AA subscribers between ISAs within a group. Re-balance affects which AA subscribers divert to which ISAs based on capacity cost. Transit subs cannot be rebalanced independent

of the parent (they move with the parent divert), and DSM subs cannot be load balanced as all subs on an ISA-AA are from the associated ISA-BB pair. The system attempts to move AA subscribers from the fullest ISA to the least full ISA based on the load balancing mode. If the load becomes balanced or an AA subscriber move fails because of ISA resources or divert IOM service queuing resources, the load balancing terminates.

Alternatively, load balancing can be manually accomplished by the AA subscriber being removed and re-added. This triggers a load balancing decision based on capacity-cost. For ESM, SAP, and spoke SDP subscriber types, this can be accomplished by removing and re-applying the AA subscriber's app-profile. In the case of ESM AA subscribers, shutting down and re-enabling either sub-sla-mgmt or the hosts has the same effect. Dynamic ESM AA subscribers re-balance naturally over time as subscribers come and go from the network.

For transit AA subscriber deployments, the parent divert SAPs are load-balanced based on AA capacity cost from the app-profile configured against the SAP or SDP. The parent capacity cost should be configured to represent the maximum expected cost when all transit subs are present.

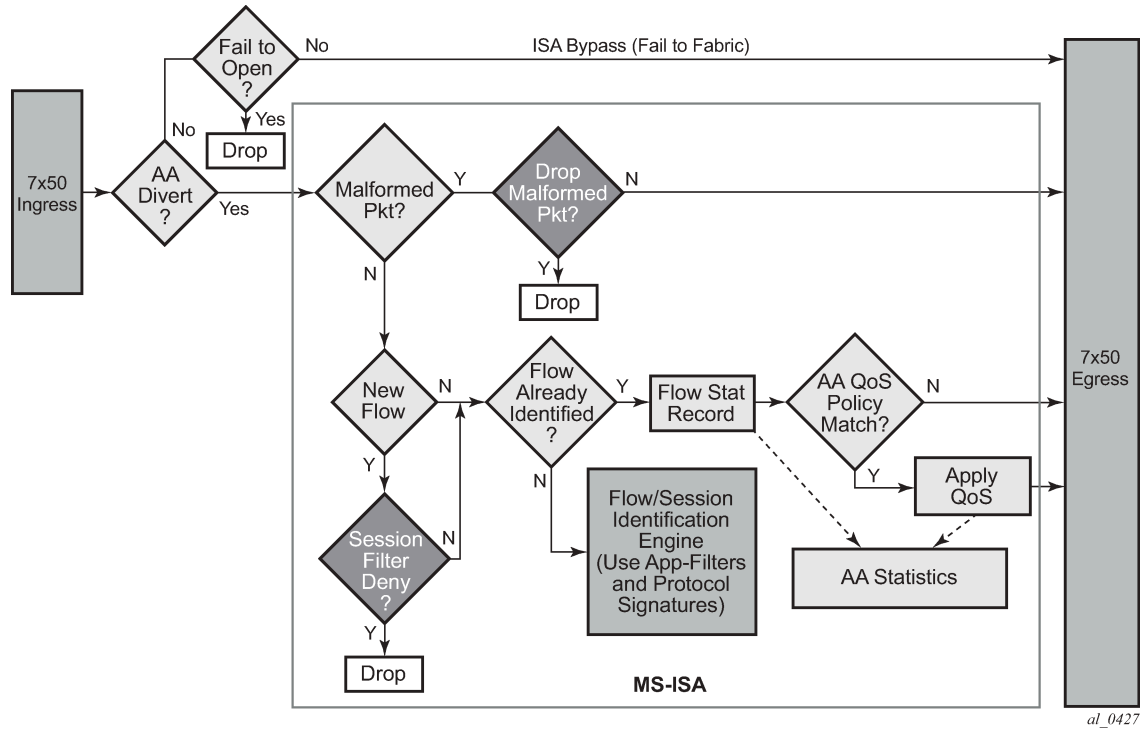
All traffic not matching a configured transit subscribers is dealt with as a member of the parent SAP and according to its app-profile.

3.2.2 AA packet processing

There are four key elements of AA packet processing ([Figure 12: Application Assurance high level functional components](#)):

- divert (selection of traffic to be diverted to the AA ISA)
- identification of the traffic on a per flow (session) basis
- reporting of the traffic volume and performance
- policy treatment of the identified traffic

Figure 12: Application Assurance high level functional components

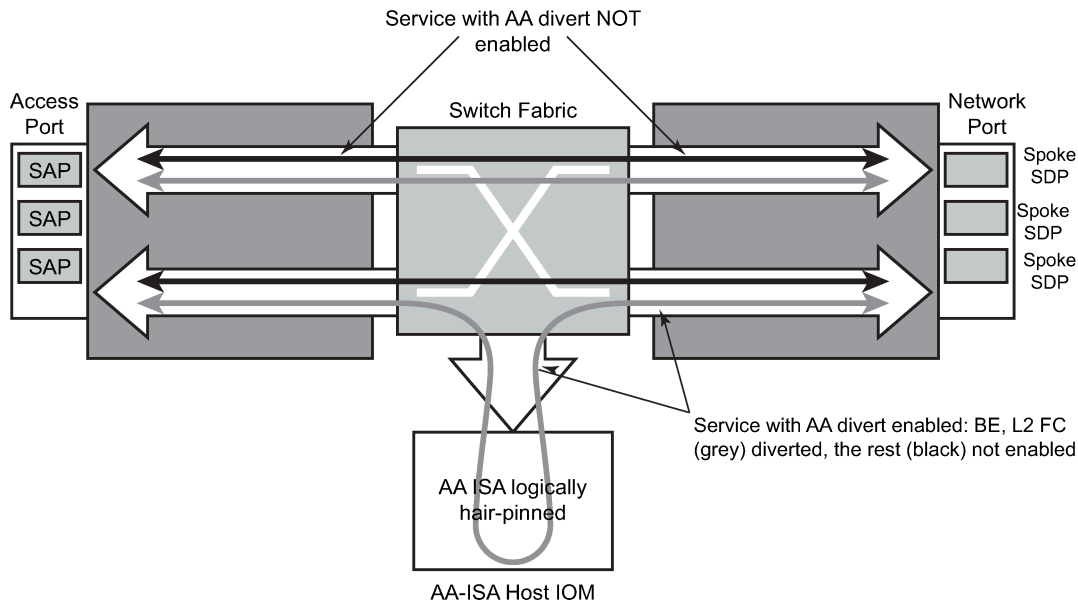


3.2.2.1 Divert of traffic and subscribers

Any traffic can be diverted for application-aware processing. AA is enabled through the assignment of an application profile as part of either an enhanced subscriber management or static configuration. This process enables the AA functionality for all traffic of interest to and from a subscriber, SAP or spoke SDP. Which traffic is deemed of interest, is configured through an AA ISA group-specific configuration of forwarding classes (FCs) to be diverted to AA and enabled on a per subscriber, SAP, spoke SDP using application profiles.

Figure 13: Application Assurance ingress datapath shows the general mechanism for filtering traffic of interest and diverting this traffic to the appropriate AA ISA module residing on an IOM (referred as the host IOM). This traffic management divert method applies to both bridged and routed configurations.

Figure 13: Application Assurance ingress datapath



Fig_47

For a SAP, subscribers with application profiles enabling AA, the traffic is diverted to the active AA ISA using ingress QoS policy filters, identifying forwarding and sub-forwarding classes that could be diverted to the AA. Only single point (SAP, ESM, or DSM subscriber, spoke SDP) configuration is required to achieve divert for both traffic originated by and destined for an AA subscriber. Diversion (divert) to the AA ISA is conditional based on the AA ISA status (enabled, failed, bypassed, and so on).

Unless the AA subscriber's application profile is configured as "divert" using Application Profiles and the FC is selected to be diverted as well, the normal ingress forwarding occurs. Traffic that is filtered for divert to AA ISAs is placed in the appropriate location for that system's AA ISA destination.

Users can leverage the extensive QoS capabilities of the router when deciding what IP traffic is diverted to the AA system for inspection. Through AA ISA group-wide configuration, at least one or more QoS forwarding classes with the "divert" option can be identified. The forwarding classes can be used for any AA subscriber traffic the service provider wants to inspect with AA.

3.2.2.1.1 Services and AA subscribers

The 7750 SR and 7450 ESS AA ISA provides, for Layer 3 to Layer 7, packet processing used by the AA feature set. AA is applied to IPv4 and IPv6 traffic on a per-AA subscriber basis, where an AA subscriber is one of the following:

- ESM subscriber
- ESM-MAC subscriber (bridged residential/vRGW device)
- distributed sub management (DSM) subscriber
- SAP or spoke
- transit

Non-IPv4 and IPv6 traffic is not diverted to AA and is forwarded as if AA was not configured; however, AA divert is supported for IP over PPPoE on Layer 2 (Epipe or VPLS) SAPs. An AA subscriber can be contained in the following services:

- IES
- VPLS
- VLL (Epipe and Ipipe)
- VPRN

AA is supported with:

- bridged CO
- routed CO
- multi-homed COs
- Layer 2/Layer 3 VPN service access points and spoke SDPs

The AA ISA feature set uses existing 7750 SR or 7450 ESS QoS capabilities and further enhances them to provide application-aware traffic reporting and management on per individual AA subscriber, AA subscriber-type or group. A few examples of per-application capabilities within the above AA subscriber contexts include:

- per AA subscriber, application traffic monitoring and reporting
- per application bandwidth shaping/policing/prioritization
- throttling of flow establishment rate
- limiting the number of active flows per application (such as BitTorrent, video or teleconference sessions, and so on)
- application-level classification to provide higher or lower (including drop) level traffic management in the system (for example, IOM QoS) and network

The following restrictions are noted — AA is not supported for tunneled transit traffic (GRE or L2TP tunnels using PPP or DHCP based policy) destined for a remote BRAS.

3.2.2.1.2 Spoke SDPs

AA on spoke SDP services allows AA divert of the spoke SDP, logically representing a remote service point, typically used where the remote node does not support AA. A SAP or spoke SDP can be assigned an app-profile, and when this app-profile is enabled for **divert**, all packets to and from that SAP or spoke SDP are diverted to an AA ISA (for forwarding classes that are configured as divert eligible).

[Table 7: Spoke SDP divert](#) shows spoke SDP divert capabilities.

Table 7: Spoke SDP divert

Access node service (spoke SDP type)	Connected to service				
	Epipe	VPLS	IES	VPRN	Ipipe
Epipe (Ethernet spoke)	Y	Y	Y	Y	Y
Ipipe (IP spoke)	N/A	N/A	Y	Y	Y

Access node service (spoke SDP type)	Connected to service				
	VPLS (Ethernet spoke)	N/A	Y	Y	Y

3.2.2.1.3 Transit AA subscribers

A transit AA subscriber is an ISA local AA subscriber contained within a parent AA subscriber. There are two types of transit AA subs:

transit IP AA subscribers defined by transit IP policy as one or more /32 IP addresses per sub

transit prefix AA subscribers defined by transit prefix policy as one or more prefix IP addresses, used in business VPNs

A transit AA subscriber incorporates the following attributes:

- name
- IP address (one or more hosts)
- app-profile (the divert/no divert and capacity cost setting of the app-profile does not affect transit AA subscribers because divert occurs only against the parent SAP)

When a SAP or spoke-SDP diverted to AA is configured with transit subs, that SAP or Spoke-SDP is referred to as the parent AA subscriber. Transit AA subscribers are supported on the parent SAPs or spoke SDPs that support AA divert.

Table 8: [Transit AA subscribers support](#) lists the Transit AA subscriber support.

Table 8: Transit AA subscribers support

Transit subscriber type	Epipe	VPLS	IES	VPRN	lpipe
Transit IP	Y	N/A	Y	Y	N/A
Transit prefix	Y	Y	Y	Y	Y

The transit AA subscribers within a parent AA subscriber can be displayed using the **show application-assurance group transit-ip-policy** or **transit-prefix-policy** command.

For transit IP AA subscribers, all packets are accounted for when they are in the ISA records. Therefore, transit IP AA subscriber counts do not count against the parent SAP in reporting. For transit prefix AA subscriber deployments using the **remote-site** command, traffic for the remote transit subs are processed and counted for both the local parent and the remote transit subscriber.

3.2.2.1.3.1 Transit AA subscriber app-profile

The app-profile assigned to the aa-sub-id affects both stats and control of the policy. App-profiles are assigned to the transit AA subscribers either explicitly when the transit-aa-sub is created, or by default (when not specified) according to a default app-profile configured in a transit-ip-policy or transit-prefix-policy. This allows transit AA subscribers to be treated with a different default app-profile than the app-

profile (default or specified) set against the parent aa sub. The number of AA subscriber stats used per ISA is proportional to the number of AA subscribers including transit subscribers subs are added.

ASO policy override is supported for static transit subs.

3.2.2.1.3.2 Transit IP policy and transit prefix policy

A transit policy is associated with the parent (divert) SAP/SDP to define how transit AA subscribers are created within that parent. The transit policy must be defined in the **configure application-assurance group partition transit-prefix-policy** or **configure application-assurance group partition transit-ip-policy** context before it can be assigned to a parent. Transit IP subs can be created using the methods described in [Table 9: Transit IP subscriber types and creation methods](#)

Table 9: Transit IP subscriber types and creation methods

Transit IP subscriber type	Creation method
Static	CLI/SNMP configuration of a transit AA subscriber is done within the transit-ip-policy
Dynamic	DHCP authentication
Dynamic	RADIUS accounting to Policy and Charging Rules Function (PCRF) or AAA
Dynamic	seen-IP transit auto-create

Transit prefix subs are created by static CLI/SNMP configuration of a transit AA subscriber within the transit-prefix-policy. The transit prefix policy follows IP filter conventions for first match and ordering of entries. While for residential /32 transits if there is an IP address conflict between any static prefix transit subs, the latter configuration is blocked, for business transit subs multiple overlapping address entries are allowed to enable longest match within subnets. IP addresses for a VPN site as an AA subscriber are configured with the transit prefix policy. There are two options:

- **aa-sub-ip** is used when the site is on the same side of the system as the parent SAP.
- **network-ip** is used when the site is on the same opposite of the system as the parent SAP.

A transit prefix subscriber may only have either **aa-sub-ip** entries or **network-ip** entries but not both.

The IP addresses defined in the **transit-ip-policy** for a transit sub are full /32 IP addresses. The IP addresses defined in the transit-prefix-policy for a transit sub are any length from /0 to /32.

Multiple IP addresses (from any prefix/pool) can be assigned to a single transit AA sub. IP addresses must be unique within a transit policy, but can be re-used in separate policies (because they have parent specific context).

The transit policy contains the default app-profile for the transit sub if a transit policy is created but app-profile is not specified. An app-profile can be later explicitly assigned to the transit sub after the sub is created (using RADIUS COA, DHCP or static).

For dynamic transit IP subs, a sub-ident-policy (also used by ESM to associate sub ID policies to a SAP) can now also be associated with the AA subscriber parent by defining the sub-ident policy in the transit IP policy. This determines how sub identifying strings are derived from DHCP option 82 fields. The policy also contains app-profile-map which maps the strings to the defined app-profiles. Transit subs do not use the sla-profile or sub-profile aspects of the sub-ident-map.

In the case of multi-homed transit subs, the transit-ip-policy must be the same on both nodes of the multi-homed parent link to ensure consistency of sub context and policy.

There is no configurable limit to the number of hosts per sub (this is similar to lease-populate which limits the number of dynamic hosts per SAP) and no limit to the number of transit subs per transit IP policy (parent). This is a function of the PE doing subscriber management.

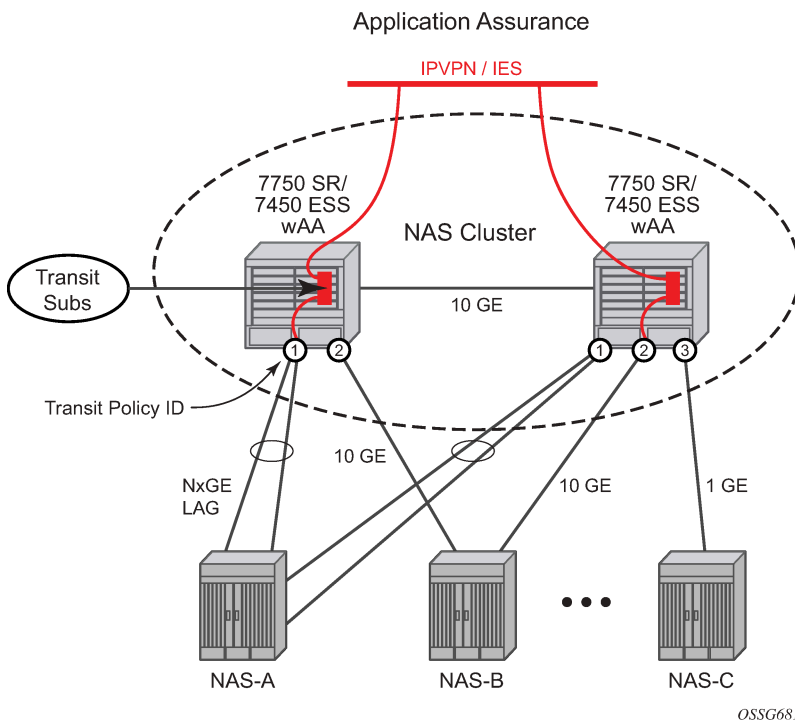
If transit sub resource limits are exceeded (hosts per sub, or subs per ISA) the transit sub creation is blocked (for both static and dynamic models).

There is a per-ISA group/partition show list of AA subscribers in a transit-ip-policy which includes a parent field for transit subs (static versus dynamic identified).

Persistent AA statistics is supported dynamic transit AA subs, ensuring that accounting usage information is not lost when the sub disconnects before reporting interval end.

3.2.2.1.3.3 Static-remote-aa-sub command

Figure 14: Static-remote-aa-sub usage topology



This command enables unique ISA treatment of transit prefix subscribers configured on the opposite (remote) side of the system from the parent SAP or spoke SDP. Provisioning a transit sub as remote-aa-sub within a transit prefix policy enables the ISA to treat any network IP-based transit subs in the following ways:

- treats packets for the parent AA subscriber independent of whether transits are also configured (stats and policers for parent work as usual)
- subsequently, treats the same packet as a transit-sub packet when matching to a configured transit sub (stats, policers)

- allows natural direction of the packet for both the parent AA subscriber and the transit-AA subscriber, as shown in [Figure 14: Static-remote-aa-sub usage topology](#), where a packet from a remote client to a local server is seen as to-sub for the parent, and from-sub for the transit sub that is logically at the far end site
- corrects directionality of packet ID for all AA subscribers and allows for correct operation of app-filter flow-setup-direction

3.2.2.1.3.4 Static transit AA subscriber provisionings

Static (through CLI/SNMP) provisioning of transit AA subscribers is supported. A profile policy override to set policy characteristics by ASO (as opposed to within an app-profile) is supported only for statically configured transit AA subs.

If there is an IP address conflict between a static and dynamic transit sub, the static takes precedence (per ESM). If the static is configured first, the dynamic transit sub is rejected. If the dynamic is created first, a warning is provided before removing the dynamic transit sub and notifying the sub-manager by CoA failure.

3.2.2.1.3.5 DHCP transit IP AA subscribers at DHCP relay node

DHCP-based transit sub creation provides a sub ID and lease time for IP addresses correlated to the ESM/subscriber context in the PE.

The 7750 DHCP relay agent creates dynamic DHCP AA subscribers when the DHCP ACK is received from the DHCP server, including the sub name, IP address and app-profile from DHCP Option 67 (if present) when the DHCP ACK messages passes through AA node to the downstream subscriber-edge node. If there is no app-profile assigned when the transit AA subscriber is created, a default transit AA subscriber app-profile is used (configured in the transit-ip-policy assigned against the divert parent AA subscriber).

This is compatible with the ESM router edge as well as third-party BRAS and CMTS.

Dynamic AA subscriber stats records are persistent across modem reset/session releases. The end of accounting records are created when transit subs are released.

Multiple IPs per transit AA subscriber are determined by seeing a common the DHCP Option 82 cct ID.

3.2.2.1.3.6 RADIUS transit AA subscribers

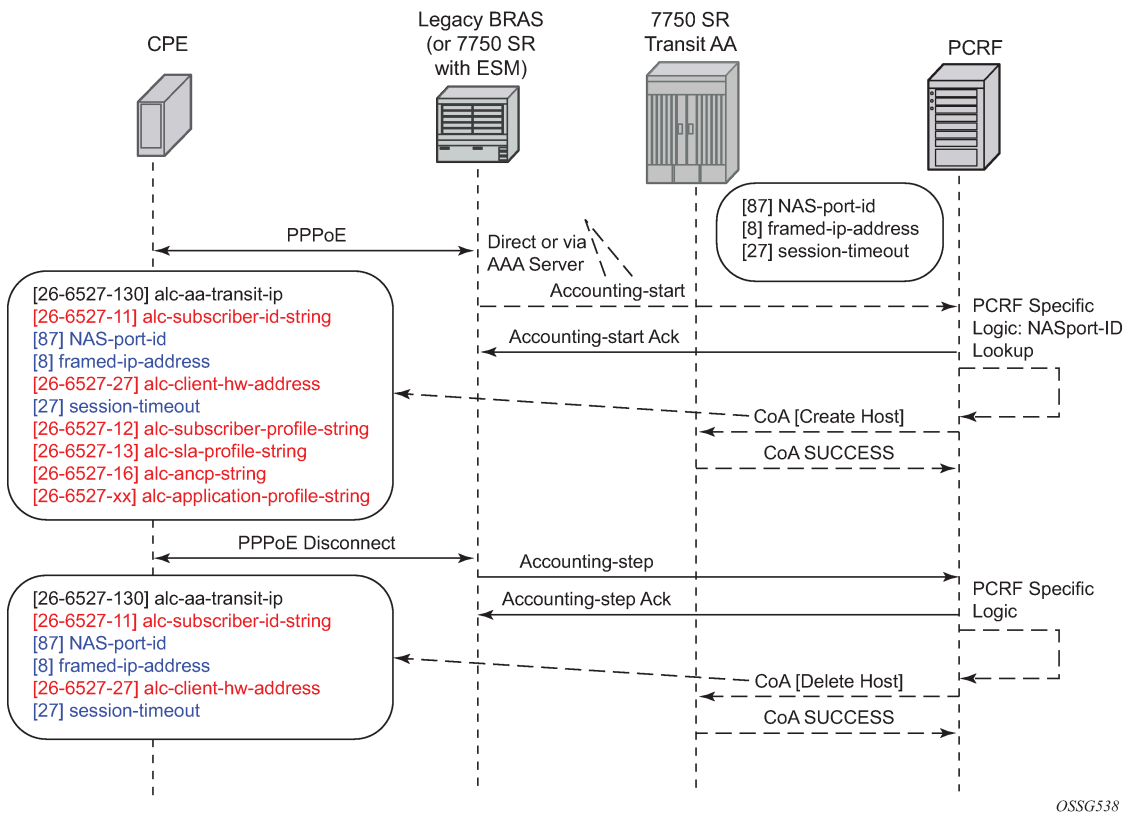
Transit subs can be dynamically provisioned by RADIUS accounting start messages forwarded by the RADIUS AAA server to a RADIUS sub-manager function at the OSS layer. This RADIUS sub manager manages dynamic transit AA subscribers on the appropriate ISA and transit-ip-policy based on the RADIUS accounting information. The interface for the sub manager to configure transit AA subscribers is RADIUS CoA messages, which are acknowledged with a CoA success message to the sub manager.

If a dynamic transit sub cannot be created as requested by a CoA because of resource constraints or conflicts, the node replies to the sub manager with a CoA fail message so that retries do not continue. This message should contain information as to the cause of the rejection. Multiple IPs per sub are allowed when common sub-ID names are seen, but with differing IP hosts.

When a RADIUS update or CoA message is seen, it could contain a modified IP address or app-profile for an existing transit sub which is accepted without affecting transit AA subscriber statistics. These transit AA subscribers are removed by the sub manager when a RADIUS accounting stop message is received.

Figure 15: RADIUS CoA example shows a RADIUS CoA Example.

Figure 15: RADIUS CoA example



The attributes in RADIUS CoA that identify the downstream transit AA subscribers are:

- downstream BRAS/ CMTS: NAS-port-ID
- IP address (framed-ip-address)
- subscriber ID (per RADIUS accounting sub-id-string)

3.2.2.1.3.7 Seen-IP RADIUS notification

Seen-IP transit subscriber notification provides RADIUS Accounting Start notification of the IP addresses and location of active subscribers within a parent AA service. This allows a PCRF to dynamically manage RADIUS AA subscriber policy (create, modify, delete) without requiring static network topology mapping of a subscriber edge gateway to the parent transit service.

When **detect-seen-IP** is enabled within a transit policy, the ISA detects IP flows on an AA parent subscriber that do not map to an existing transit AA subscriber. It then uses a simple RADIUS Accounting Start notification from the transit AA node to the PCRF to initiate subscriber creation, providing information about the location of the transit subscriber traffic. This provides notice for subscriber authentication changes, such as new subscriber sessions or new host IP addresses added to an existing AA subscriber, while being independent of the network topology for how the BNG is homed into the transit AA nodes.

The RADIUS Accounting Start message is sent to the RADIUS Server referenced by the specified **seen-ip-radius-acct-policy**. This RADIUS message contains the following information about the flow:

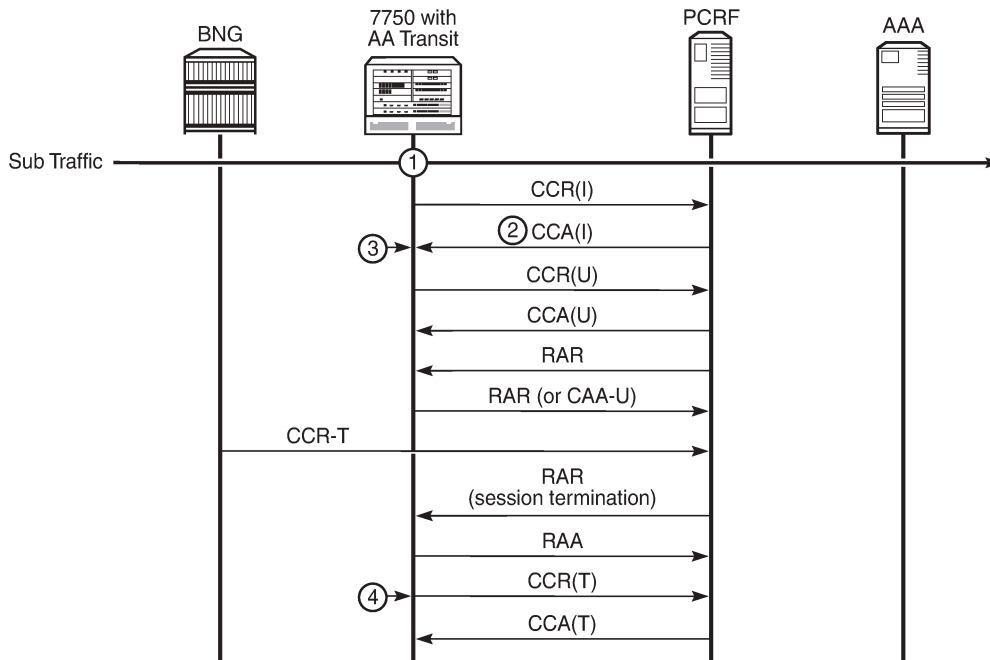
- subscriber-side IP address
- parent SAP or spoke-SDP ID (NAS port ID)
- IP address of node making the request
- peer SAP or spoke-SDP ID (NAS port ID), if configured
- peer IP address of SR making the request, if configured
- AARP ID, if configured

3.2.2.1.3.8 Diameter transit AA subscriber

For Diameter transit AA subscribers, AA auto-detects new IP addresses and notifies the PCRF of new subscribers via a Gx CCR-I message. The PCRF locates the subscriber’s AA policy and returns the information via CCA-I message to AA.

Figure 16: 3GPP pull model shows a 3GPP pull model, whereby AA initiates the Gx session. Table 10: 3GPP pull model description describes the figure. The PCRF can, at any time after the session is established, push new policies using a RAR message. New policies can include new usage-monitoring or AA ASO values.

Figure 16: 3GPP pull model



25621

Table 10: 3GPP pull model description

Legend	Description
1	AA detects a new IP address, and sends a CCR-I containing the subscriber-side IP address

Legend	Description
2	CCA(I) contains subscriber AA policy information
3	AA applies an AA appProfile, ASOs, and any AA usage monitoring
4	AA reports usage counters for all specified or enabled usage monitoring keys, and removes the sub

The CCR-I message from the 7750 SR node to PCRF contains the following:

- session ID
- subscriber-side IP address
- IP-CAN-Type AVP (if enabled) with the value "tbc"
- subscription ID AVP, with the following characteristics (if *avp-subscription-id* is configured as *subscriber-id*):
 - Type is END_USER_E164 (private by default)
 - ID is an auto-generated unique ID

The CCA-I message from PCRF to the 7750 SR node contains the following:

- session ID
- Charging-Rule-Install containing the following information:

```

Charging-Rule-Definition ::= <AVP Header: 1003>
    {Charging-Rule-Name}
    [TDF-Application-Identifier]
    [Monitoring-Key]
    [AA-Functions {
        AA profile
        AA-App-Service-Options {
            AA-App-Service-Options-Name
            AA-App-Service-Options-Value
        }
        .....
    }
    [AVP]
    
```

- Charging-Rule-Name
 - Usage monitoring starts with "AA-UM:"
 - AppProfile and ASOs start with "AA-Functions:"
- AA-Functions (AVPs to set AA profile and ASO values)
- TDF-application-identifier; this field specifies a predefined AA charging group, application group, or application name for which usage monitoring functionality is required

Termination of the Gx session is only done after AA receives an RAR-T message from PCRF with the session-release-cause AVP meeting the configured threshold. After replying to an RAR message with an RAA message, AA sends a CCR-T message with reports of usage counters, if any are enabled.

[Table 11: Used AVPs](#) lists the AVPs used for Diameter transit AA subs.

Table 11: Used AVPs

AVP	Category	Details	User configurable
Session-Id	M	Globally unique Generated for each session as: <peer identity>;<high 32 bits>;<low 32 bits>; [<optional value>;]<subscriber ip>	N
Auth-Application-Id	M	Set as Gx (16777238)	N
CC-Request-Type	M	Set to INITIAL_REQUEST (1) when initiating a new session Set to TERMINATION_REQUEST (3) when ending a session Set to UPDATE_REQUEST (2) in all other situations	N
CC-Request-Number	M	Generated internally according to Gx specifications Request numbering starts at 0	N
Subscription-Id	M	Configurable using Subscription-Id-Type and Subscription-Id-Data	—
Subscription-Id-Type	M	Configurable under subscriber-mgmt>diameter-application-policy>gx>avp-subscriber-id origin	Y
Subscription-Id-Data	M	Configurable under subscriber-mgmt>diameter-application-policy>gx>avp-subscriber-id origin [type type]	Y
Framed-IP-Address	M	Set to the subscriber's IP address as seen by AA-ISA in the from-sub direction of the data traffic	N

When the Subscription-Id-Type is "Subname", then Subscription-Id-Data is auto-generated by AA to be unique node-wide, using the transit IP policy, SAP, and sub IP address.

Unlike AA ESM Diameter-controlled subscribers, transit Gx AA subscribers are not required to support ADC rules over Gx.

Transit Gx AA subscribers use PCC rules as per ESM Diameter AA subscriber implementation, and therefore uses AA-Function AVP.

For transit Gx AA subs, similarly to ESM Gx-controlled subs, the PCRF can set the subID in a CCR-I by sending a PCC rule with the name of the charging rule prefixed with "sub-id:". The AVP appears as follows:

```
charging-rule-install (1001) VM----- [44]
  vendor-id TGPP
  data [32] (Grouped)
```

```
charging-rule-name (1005) VM----- [30]
  vendor-id TGPP
  data [18] (UTF8String) : Sub-Id:subID-name
```

In addition to using the AA Function AVP, AA supports the configuration of the application profile and ASOs by the PCRF via a CCR-I, CCR-U, or RAR that has PCC rules with the name of the charging rule prefixed with "AA-Functions:App:" or "AA-Functions:ASO", such as:

```
charging_rule_install[0].charging_rule_name[0] = AA-Functions:App:<name>
charging_rule_install[0].charging_rule_name[1] = AA-Functions:Aso:<char>:<val>
...
charging_rule_install[0].charging_rule_name[n] = AA-Functions:Aso:<char>:<val>
```

AA allows for the definition of up to 32 ASOs. If the number of ASOs is larger than what can fit within a single charging-rule-install AVP, multiple charging-rule-install AVPs can be used in the CCR-I message.

As the Gx protocol is already supported by the 7750 SR/VSR system, there are no new configurations required. All existing configurations introduced to support ESM Gx control on a BNG can be re-used in AA transit deployment.

3.2.2.1.3.9 Transit AA subscriber auto-create

For dynamic transit AA subscribers, AA can automatically detect new IP addresses and create a local subscriber context with no interaction with RADIUS or Diameter policy control. When transit-auto-create is enabled within a transit policy, the ISA detects IP flows on an AA parent subscriber that do not map to an existing transit AA subscriber. When auto-create is enabled, AA subscriber contexts are auto-created under the parent diverted SAPs and spokes using the transit-ip-policy name and subscriber IP address as the AA-sub name. The default app-profile configured against the transit-ip-policy is applied to these subscribers.

By default, auto-created subscribers are persistent and never removed (without removal of the AA group). ISAs may periodically be shut down and then no shutdown to clear the aa-subs to avoid running out of sub context, or, inactive subscribers may be automatically removed after a timeout period by using the transit-auto-create context command for inactivity-monitoring. With this, periodically AA removes any inactive auto-created subscriber where an inactive sub is defined as having no active flows in the last period.

3.2.2.1.3.10 Transit AA subscriber persistence

Transit AA subscribers can be persistent within a single node, because in some cases, there is not a dual-node BNG subscriber redundancy configuration. This allows a single node that has dynamically created transit subs to retain the subscriber state, context, and statistics across a node or ISA reboot.

If dynamic transit AA subscribers are released, renewed, or otherwise changed during an outage or reboot of a transit AA node, the sub manager notifies the transit node of these changes.

Prefix transit subs are not affected by persistence because they can only be statically configured.

3.2.2.1.3.11 Transit diameter AA subscriber geo-redundancy

If there is no Multi Chassis Synchronization (MCS) between the two peer nodes, the two geo-redundant nodes are configured as two distinct realm nodes from the PCRF point of view. Each node acts independently of the other node. After a switch-over, AA on the newly active node detects new traffic flows

with new IP addresses. AA notifies PCRF with a CCR-I message to retrieve subscriber policies. After PCRF confirms the same IP address used on a different Gx session ID, it deletes the old session for that IP address.

Similar behavior takes place if an MCS is configured with session IDs and OSI states journaled across, as per ESM implementation. The two geo-redundant nodes are configured with the same host realm, and appear to PCRF as one node. After a switch-over, AA-ISA on the newly active node detects new traffic flows with new IP addresses. AA-ISA notifies PCRF with a CCR-I message to retrieve subscriber policies. The Gx session ID used is unique and different from the session ID used on the previously active node. After PCRF confirms the same IP address used on a different Gx session ID, it deletes the old session for that IP address.

3.2.2.1.3.12 Policers for transit AA subscribers

AA subscriber per-subscriber policers can provide per-SAP policing for the parent SAP, with transit AA subscribers each supporting distinct per-subscriber policers within the parent (packets are only processed once against one AA subscriber, the parent or the transit subscriber). Packets matching transit AA subscribers and policers are not included in a parent policer.

There is no policer hierarchy unless system wide policers are referred to by both the parent AA subscriber and transit AA subscriber. When the remote-site configuration is not used, system policers can be used to police all traffic for a site containing transits, subject to constraints on system policer scale.

When the remote-AA subscriber config is used, the parent owns all packets for stats and policing, so any transit subscriber configuration within the parent does not affect the stats or policers. AA policers are supported on a transit subscriber basis, across all (multiple) IP prefixes per sub.

3.2.2.1.3.13 ISA host IOM for transit subs

The AA divert IOM is not impacted by transit AA subscribers in the divert parent. The ISA host IOM egress datapath functions to convert the parent SAP into transit AA subscribers that are then handled by the ISA consistent with all other AA subscriber features. The ISA itself treats all AA subscribers equally regardless of whether the AA subscriber is from ESM, from DSM, from an SAP, or from a transit subscriber in a parent SAP or spoke.

Prefix transit subs can only be created on IOM4 cards supporting ISAs (and IOM-e).

3.2.2.1.4 AA subscriber application services

This section describes AA subscriber application services.

3.2.2.1.4.1 Application profile

Application profiles enable AA service for an ESM or DSM subscriber, Service Access Point or spoke SDP (AA subscriber). Each application profile is unique in the system and defines the AA service that the AA subscriber receives. An ESM subscriber can be assigned to an application profile which affects every host of the particular subscriber. For SAP or spoke SDP AA subscribers, an application profile can be assigned, which affects all traffic originated/destined over that SAP or spoke SDP. By default, ESM and DSM subscribers, SAPs or spoke SDPs are not assigned an application profile.

The following are main properties of application profiles:

- One or more application profiles can be configured in the system.
- Application profiles specify whether AA subscriber's traffic is to be diverted to AA.
- Application profiles are defined by an operator that can reference the configured application service options (ASO) characteristics. See [Application Service Options \(ASOs\)](#).
- Application profiles must only be assigned after AA resources (AA ISA cards) are configured.
- Application profiles can be assigned a capacity cost used for subscriber load balancing among ISAs within the AA group (see [ISA load balancing](#)).
- Application profiles can be assigned to a scope from a subscriber or session, which controls whether the application profile is applied to the entire subscriber or to a device.

ESM and DSM policy includes an application profile string. The string points to an application profile pre-provisioned within the router and is derived by:

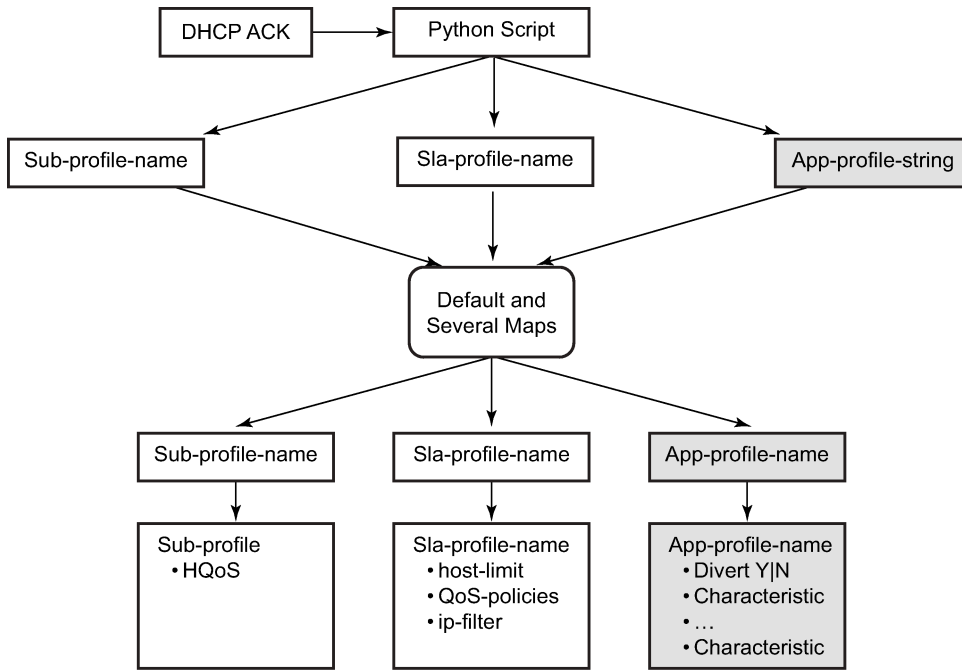
- parsing the DHCP Option 82 sub-option 1 circuit ID payload, vendor specific sub-option 9, or customer-defined option different from option 82, during authentication and the DHCPDISCOVER, as well as re-authentication and the subscriber's DHCPREQUEST
- RADIUS using a new VSA: [26-6527-45] Alc-App-Prof-Str
- DIAMETER using "AA-profile-name" AVP under ADC rule
- inheritance from defaults in the **sap>sub-sla-mgmt** context, to allow default application profile assignment if no application profile was provided
- static configuration

Mid-session (PPP/DHCP) changes to the application profile string allows:

- modification of the application profile a subscriber is mapped to and pushes the change into the network as opposed to waiting for the subscriber to re-authenticate to the network
- change to the subscribers application profile inline, without a need for the subscriber to re-authenticate to RADIUS or perform any DHCP message exchange (renew or discover) to modify their IP information

[Figure 17: Determining the subscriber profile, SLA profile and application profile of a host](#) shows the process for determining the subscriber profile, SLA profile, and application profile of a host.

Figure 17: Determining the subscriber profile, SLA profile and application profile of a host



OSSG170

3.2.2.1.4.2 Application profile map

AA adds new map (**app-profile-map**) application profile command to associate an *app-profile-string* from dynamic subscriber management to a specific application profile using its *app-profile-name* that has been pre-provisioned. The application profile map is configured in the **config>subscr-mgmt>sub-ident-pol** context.

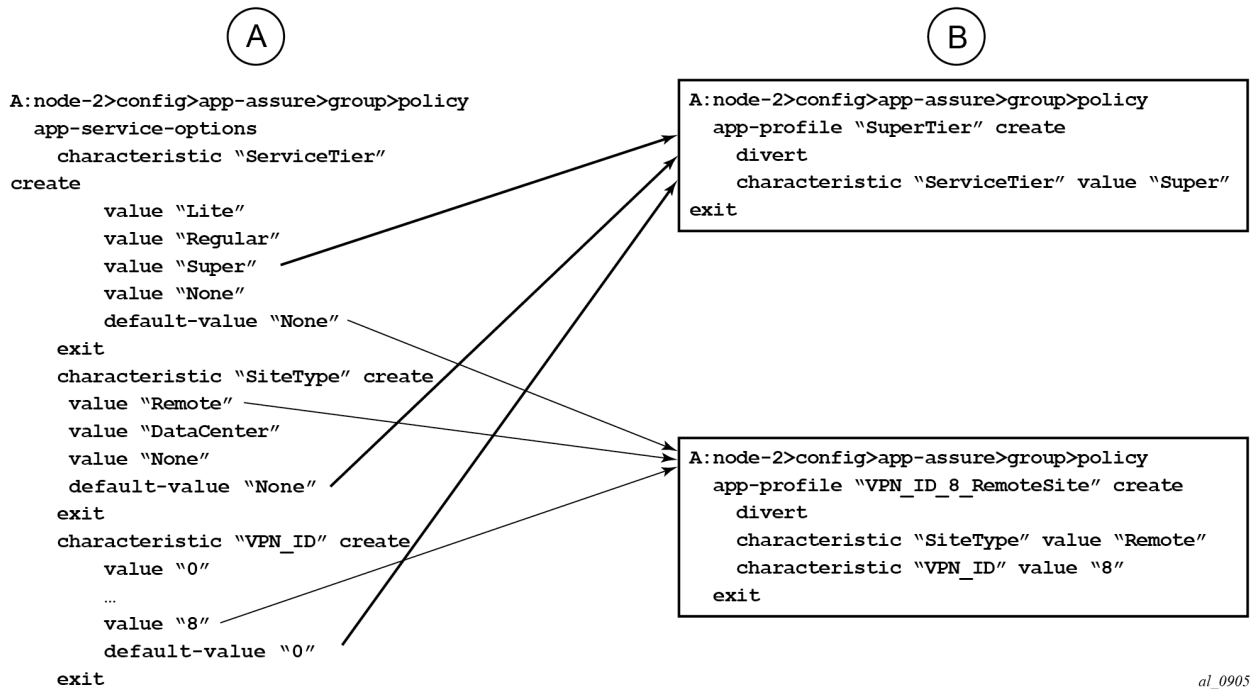
The pre-defined subscriber identification policy has to be assigned to a SAP, which determines the sub-id, sub-, sla-, and app-profiles.

3.2.2.1.4.3 Application Service Options (ASOs)

ASOs are used to define service provider or customer visible network control (policy) that is common between sets of AA subscribers (for example, upstream/downstream bandwidth for a tier of AA service). ASO definition decouples every AA subscriber from needing subscriber-specific entries in the AQP for standard network services.

As an example, an operator can define an ASO called ServiceTier to define various HSI services (Super, Lite, and so on) (Figure 18: Configuration example A). The operator can then reference these defined ASOs when creating the App Profiles that are assigned to AA subscribers (Figure 18: Configuration example B).

Figure 18: Configuration example



Then, the defined ASOs are used in the AQP definition to determine the needed treatment or policy (Figure 19: AQP definition example).

Figure 19: AQP definition example

```

app-qos-policy
  entry 50 create
    description "Limit downstream b/w for Super subscribers"
    match
      traffic-direction network-to-subscriber
      characteristic "ServiceTier" eq "Super"
    exit
    action
      bandwidth-policer "SuperDown"
    exit
    no shutdown
  exit
  entry 110 create
    match
      application-group eq "Tunneling"
      characteristic "SiteType" eq "Remote"
    exit
    action
      remark fc af
    exit
    no shutdown
  exit
  
```

al_0906

Alternatively, if ASOs were not used in the previous example, then the operator would have to define a unique AQP entry for every subscriber. Each of these AQPs has its "match" criteria setup to point to the subscriber ID, while the action for all of these unique AQPs is the same for the same service (for Tier 1 service, the policer bandwidth is the same for all Tier 1 AA subscribers) (Figure 20: Single ASO example).

Figure 20: Single ASO example

```

7750SR>config>aa>group>policy>aqp>
  entry 100 create
  match
    aa-sub eq " sub_1"
  exit
  action
    bandwidth-policer "superDown"
  exit
  no shutdown
exit

  entry 101 create
  match
    aa-sub eq " sub_2"
  exit
  action
    bandwidth-policer "superDown"
  exit
  no shutdown
exit

  entry 102 create
  match
    aa-sub eq " sub_3"
  exit
  action
    bandwidth-policer "superDown"
  exit
  no shutdown

```

al_0907

The single ASO example shows how the use of a single ASO can prevent the user from provisioning an AQP entry every time a subscriber is created.

Other example uses of ASO entries include:

- entry per application group that is to be managed, such as VoIP, P2P, HTTP
- several entries where specific applications within an application group can individually be managed as service parameters, for example, HTTP content from a specific content provider, or streaming video from network television or games
- HSI tiers (for example, Gold, Silver, and Bronze for specifying bandwidth levels)
- VPN customer ID

Application characteristics are defined as specific to the services offered within the operator network. The operator defines ASO characteristics and assigns to each ASO one or more values to define service offering to the customers.

The following are the main elements of an ASO:

- A unique name is applied to each characteristic.

- The name is unique to the group partitionpolicy, but the expectation is that characteristics are consistent network wide.
- Operator-defined values (variables) are defined for each characteristic and are unique to each characteristic. A default value must be specified from the set of the values configured.

The following lists how ASO characteristics are used:

- Application service options are used as input to application profiles.
- AQP rule sets also use the ASO characteristics to influence how specific traffic is inspected and policies applied.
- Multiple ASO characteristic values are allowed in a single rule.

Syntax checking is performed when defining application profiles and AQPs that include application characteristics. This ensures:

- The characteristic is correctly identified.
- When specifying a characteristic in an **app-profile** and **app-qos-policy**, the value must be specified. The default value applies if a characteristic is not specified within an **app-profile**.

3.2.2.1.4.4 ASO overrides

This feature enables individual attributes/values to be set against an AA subscriber complementary to using app-profiles. The AA subscriber types supported that can have ASO overrides by CLI/SNMP are provisioned business AA SAPs and spoke SDPs, and statically-provisioned transit AA subs. ESM and transit-aa AA subscribers can have ASO overrides applied by RADIUS override VSAs.

Application profile assignment is still used to obtain the following information:

- the application-assurance group (and partition) to which the AA subscriber is being assigned to
- whether the traffic should be diverted
- capacity-cost (for load balancing to a multi-isa group)

The information configured in the app-profile is also used, but ASO characteristics and values (these are from the policy defined in the group and partition) can be overridden.

The overrides are specific to a single AA subscriber. An ASO override does not affect any other AA subscriber or the app-profile config itself.

Typically the ASO characteristics in the app-profile would not be specified, therefore leaving all characteristics at their default values. This is not mandatory though and the app-profile could specify any ASO characteristic and non-default value.

The AQP has entries that can refer to ASO characteristics (attributes) and values in their match criteria. In the absence of any individual attribute/value override, an AA subscriber continues to work as before - using the ASO characteristics/values defined inside the app-profile assigned to them. With overrides, the AA subscriber attributes used in AQP lookups are the combination of the following:

- the characteristics/values from the app-profile
- any specific characteristics and values overridden for that AA subscriber

Show command output displays the combined set of attributes that apply to the AA subscriber.

The **override** commands can only be used if there is already an app-profile assigned to the AA subscriber, otherwise, the overrides are rejected.

The app-profile attribute override is assigned to a specific AA subscriber (SAP, spoke SDP) within the AA Group:partition with where the subscriber exists. While subscriber names are unique, the Group:partition policy context where apps, app-profiles and ASO characteristics are defined is relevant to the override context. Override for ESM subscribers can be triggered via DIAMETER or RADIUS.

3.2.2.1.4.5 AA subscriber scale mode and ASOs

When an AA policy is administered using a per-AA subscriber policy attribute assignment (ASO override) as opposed to using a service profile-based model, the number of attributes and values of ASOs that are needed in an AA Group can be much larger than the ASO scale needed for app-profile based deployments.

In conjunction with the App-profile ASO override, set the AA-group to a mode optimized for the needed ASO, subscriber, and flow scale requirements.

3.2.2.2 Traffic identification

AA traffic identification classifies per-flow traffic into applications and application groups, and assigns flow attributes to the traffic.

Application identification means there is sufficient flow information to provide the network operator with a view of the underlying nature and value of the content based on the end-user application related to each flow. The application ID does not include:

- anti-virus signatures per IPS/UTM
- content inspection (e-mail, text, picture, or video images). The payload data content of flows is not examined as part of the application identification.

AA can identify and measure non-encrypted IP traffic flows using any available information from Layer 2 to Layer 7, and encrypted IP traffic flows using heuristic and pattern-match techniques. Nokia maintains a global application database (App-DB) to keep application definitions (app filters and applications) up-to-date in customer deployments.

AA attempts to positively identify the protocols and applications for flows based on a pattern signature observation of the setup and initial packets in a flow. The system correlates control and data flows belonging to the same application. In parallel, statistical and behavioral techniques are also used to identify the application. Until identified, the flow does not have a known application and is treated according to the default policies (AQP policies defined using all or any ASO characteristics, subscriber ID and traffic direction as match criteria) for traffic for that AA subscriber, app-profile and direction (packets are forwarded unless an action is configured otherwise). If the identification beyond OSI Layer 2 is not successful, the flow is flagged as an unknown protocol type, (for example `unknown_tcp` or `unknown_udp`). The unknown traffic is handled as part of all application statistics and policy, including generation of stats on the volume of unknown traffic.

AA allows operators to optionally define port-based applications for trusted TCP or UDP ports. Operators must explicitly identify a TCP/UDP port or ports in an application filter used for trusted port application definition and specify if protocol signature-based application identification is to be performed on a flow. Two options are available:

- If no protocol signature processing is required (expected to be used only when (A) AQP policy must be performed from the first packet seen, (B) the protocol signature processing requires more than one packet to positively identify a protocol or application, and (C) no other application traffic runs over a

TCP/UDP port), the first packet seen by AA ISA for a flow on that TCP/UDP port allows application identification. The traffic for a flow is identified as "trusted tcp/trusted_udp" protocols.

- If protocol signature verification of an application is required, that is, it is expected to be used only when:
 - AQP policy must be performed from the first packet seen.
 - The protocol signature processing requires more than one packet to positively identify a protocol or application.
 - Other application traffic may run over a TCP/UDP port, for example TCP port 80

The first packet seen identifies the application but protocol signature-based analysis continues. After the identification completes, the application is re-evaluated against the remaining application filters allowing detection and policy control of unexpected applications on a trusted port.

Flow attributes provide optional flow classification metadata extracted from traffic by AA stateful flow processing, which complements and is orthogonal to the AA application and app-group classification.

When a flow attribute is enabled, every flow is assigned a confidence level for each enabled attribute. Confidence levels are used to condition control actions and for analytics. AA cflowd record exports include the 0-100% confidence level for each attribute, and are null for attributes that are not enabled. This allows analytics of traffic by flow attribute, including confidence. As with AA protocol signatures, flow attribute signatures are provided as a set in the AA .tim file.

At AA system startup or after an AA ISA activity switch, traffic diverted to the ISA may contain flows that are already in progress. Typically, application identification relies on patterns in the initial set of packets exchanged. If these initial packets were not seen by the ISA, in progress TCP flows are marked with the existing protocol signature and have a policy applied according to an application based on the existing protocol until they end or the identification of an in-progress flow is possible. Statistics are generated.

From the first packet of a flow, a default per-AA subscriber AQP policy is applied to every packet. After an application is identified, subsequent packets for a flow apply AA subscriber and application-specific AQP. The AA-generated statistics for the flow with AA subscriber and application context are collected based on the final determination of the flow's application. A subset of the applications may be monitored on an ongoing basis to further refine the identification of applications carried with the traffic flow and to identify applications using an external application wrapper to evade detection.

3.2.2.2.1 AA identification components

[Figure 21: Policy structure](#) shows the relationship between the AA system components used to identify applications and configure AA-related capabilities. Each ID-related component is defined as follows:

- protocol signatures
- application filters
- applications
- application groups
- charging filters
- charging groups
- flow attributes

Figure 21: Policy structure

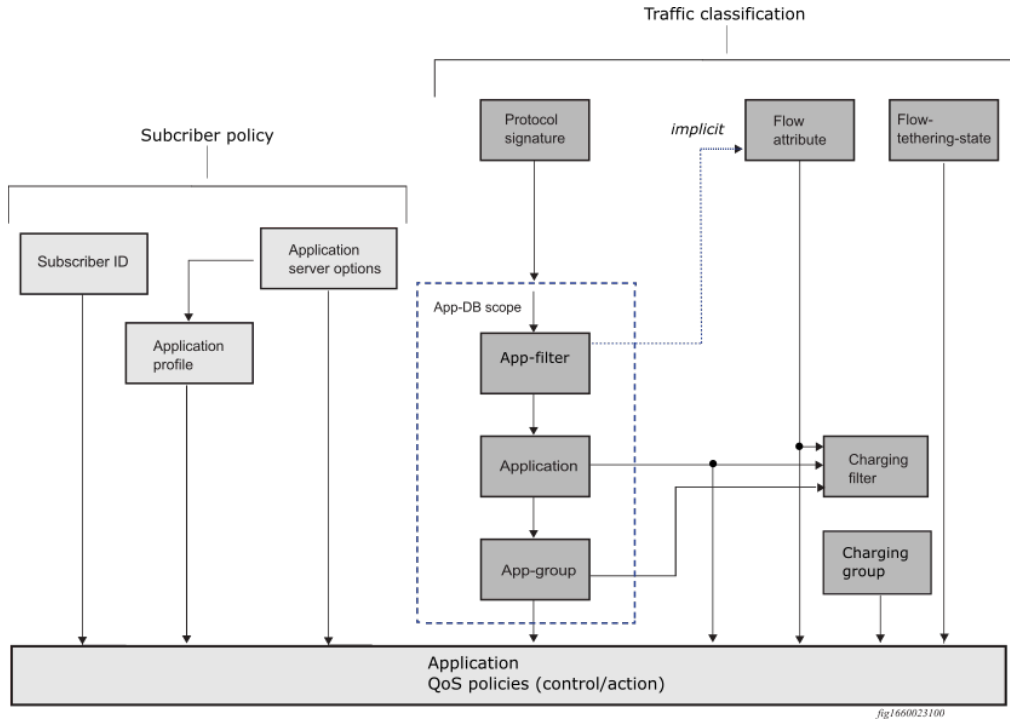


Table 12: AA classification elements provides an overview of how those various components are used in AA to recognize types of flows and sessions.

Table 12: AA classification elements

Term	Definition	Examples
Protocol signature	The Nokia proprietary component of AA flow identification provided as part of AA S/W load to identify protocols used by clients. Where a protocol is defined as an agreed on format for transmitting data between two devices.	Tftp, iMap, msn-msgr, RTP, emule, http_video, bittorrent, SIP. The Nokia protocol signatures do not rely on IP port numbers to identify a TCP/UDP port based protocols and applications to avoid false-positives but allow users to define application filters if a port-based identification is deemed adequate (see the example below).
Application filter	User configurable, optional component of AA flow identification that uses any combination of protocol signatures, server IP address and port, flow set-up direction, configurable expressions (for example, an HTTP string match) to identify user's traffic.	The http_video + IP address of a partner's video server or the http_video + an HTTP string to identify a partner's video content as TCP or UDP + TCP/UDP port number to identify a TCP- or UDP-based protocol or application.

Term	Definition	Examples
Application	User configurable, optional component of AA flow identification that allows defining any specific forms of traffic to and from end user clients by combining application filter entries.	Google Talk, POP3, YouTube, iTunes, Shoutcast.
Application charging group	User configurable, optional component of AA flow identification that allows grouping of similar end use applications using user-defined names and groups.	IM, mail, multimedia, P2P, tunneling, Web, other.
Clients	End user programs that generate user traffic for applications and protocols, and that are used in a process of AA flow identification verification.	The list of clients is constantly evolving as new clients or versions are introduced in the marketplace. The following example illustrates clients that may be used to generate Application traffic matching Bit Torrent application defined using Bit Torrent and DHT protocol signatures: Limewire, BitTorrent, Azureus, Ktorrent, Transmission, Utorrent.
Flow attributes	Nokia-provided algorithms based on machine learning, which assign flow attribute confidence levels to all traffic flows, whether encrypted or not. The attribute results may be used for analytics or control purposes.	Video, audio, download, upload, abr_service, encrypted, ESNI, real_time_communication
Charging filter	User configurable, optional component of AA flow identification that uses any combination of flow attribute, application, or app-group specified charging group, and flow tethering detection to determine the charging group.	Charging video traffic for a specified application differently from other traffic of the same application that has different flow attributes.

3.2.2.2.2 Protocol signatures

The set of signatures used to identify protocols is generated by Nokia and included with the AA software load. The signature set includes:

- The protocols that can be identified with this load, using a combination of pattern and behavioral techniques. The protocols are used in generating statistics by protocol, and are used as input in combination with other information to identify applications.
- Pattern signatures are the set of pattern-match signatures used in analysis.
- Behavior signatures are the set of diagnostic techniques used in analysis.

Dynamic upgrades of the signatures in the system are implemented by invoking an **admin application-assurance upgrade** command and then performing AA ISA activity switches.

The protocol signatures are included in `aa-isa.tim` software load which is not tightly coupled with software releases allowing for protocol signature updates without upgrading and impacting of routing/forwarding engines as part of an ISSU upgrade that updates only the AA ISA software. See upgrade procedures described in the *SR OS R23.x.Rx Software Release Notes* for more information.

Because protocol signatures are intended to be the most basic block of Application Identification, other AA components like Application Filters are provided to further customize Protocol Signatures allowing operators to customize their applications and to reduce a need for a new Protocol Signature load when a new Application may need to be identified. This architecture gives operators more flexibility in responding to ever changing needs in application identifications.

Signature upgrade without a router upgrade is allowed within a major router release independently of system ISSU limits. An AA ISA signature upgrade is supported before the first ISSU router release (for example, operators can upgrade signatures for pre-ISSU minor releases).

In addition, any router release from ISSU introduction release can run any newer `aa-isa.tim` image within the same major release by performing an `aa-isa.tim` single step upgrade. For example, Release 8.4 may be upgraded in a single step to run Release 8.14 of `isa-aa.tim`.

Each protocol, except internal protocols used for special-case processing statistic gathering (cut-through, for example), can be referenced in the definition of one or multiple applications (through the App-Filter definition). Assignment of a supported protocol to an app-filter or application is not mandatory. Protocols not assigned to an application are automatically mapped by the system to the default **Unknown** application.

3.2.2.2.3 Custom protocols

Custom protocols are supported using configurable strings (up to 16 hex octets) for pattern-matched application identification in the payload of TCP or UDP based applications (mutually exclusive to other string matches in an app-filter).

The match is specified for the client-to-server, server-to-client, or any direction for TCP based applications, and in the any direction for UDP based applications.

There is a configurable description and custom protocol ID for a protocol, with configurable shutdown. When disabled, traffic is identified as if the protocol was not configured.

Custom protocols and ALU-provided protocols are functionally equivalent. Custom protocols are used in application definition without limitations (all app-filter entries except strings are supported). Collection of custom protocol statistics on a partition/ISA group/special study sub level is supported.

3.2.2.2.4 Protocol shutdown

The protocol **shutdown** feature provides the ability for signature upgrades without automatically affecting policy behavior, especially if some or even all new signatures are not required for a service. All new signatures are disabled on upgrade by default to ensure no policy/service impact because of the signature update.

All protocols introduced at the R1 stage of a specific release are designated as "Parent" signatures for a specific release and cannot be disabled.

Within a major release, all protocols introduced post-R1 of a major release as part of any `isa-aa.tim` ISSU upgrade are by default **shutdown**. They must be enabled on a per-protocol basis (system-wide) to take effect.

When shutdown, post R1-introduced protocols do not change AA behavior (app-id, policy, statistics are as before the protocol introduction), for example, traffic maps to the parent protocol on which the new signature is based. In cases where there is more than one parent protocol, all traffic is mapped to a single, most-likely, parent protocol. For example if 80% of a new protocol has traffic mapping to `unknown_tcp`, and 20% mapping to another protocols, `unknown_tcp` would be used as parent.

Enabling/disabling of a new protocol takes affect for new flows only. The current status (enabled/shutdown) of a signature and the parent protocol is visible to an operator as part of retrieving protocol information through CLI/SNMP.

3.2.2.2.5 Supported protocol signatures

Protocol signatures are release independent and can be upgraded independently from the router's software and without impacting router's operations as part of an ISSU upgrade. A separate document describes signatures supported for each signature software load (`isa-aa.tim`). New signature loads are distributed as part of the SR/ESS maintenance cycle. Traffic identified by new signatures are mapped to an **Unknown** application until the AA policy configuration changes to make use of the newly introduced protocol signatures.

3.2.2.2.6 Application groups

Application groups are defined as a container for multiple applications. The only application group created by default is **Unknown**. Any applications not assigned to a group are automatically assigned to the default **Unknown** group. Application groups are expected to be defined when a common policy on a set of applications is expected, yet per each application visibility in accounting is required. The application group name is a key match criteria within application QoS policy rules.

3.2.2.2.7 Charging groups

Charging groups allow usage accounting by application or application groups in a manner that does not affect app to application group mapping.

Charging groups are user-configurable groupings of applications or application groups that are charged in a similar manner.

For example, AA application group statistics for "Streaming Video" includes all streaming applications, independent of whether any specific application is 0-rated for charging. AA charging groups are used for charging-related statistics.

As with application groups, charging groups are defined under an AA policy context for an AA group or partition. After being defined, individual applications and application groups can be directly associated with the chosen charging group. The charging-group name is a key match criteria within application QoS policy rules.

A default charging group can be specified for the AA policy to associate a charging group to any applications or application groups that are not explicitly assigned to a charging group.

Charging groups are also assigned an export-id number for accounting export purposes.

If no export ID is assigned, that charging group cannot be added to the AA subscriber stats RADIUS export type. After a charging-group index is referenced, it cannot be deleted without removing the reference.

The priority to determine the charging group for a flow is as follows:

1. charging group assigned by a charging filter (if configured)
2. default tethered charging group (if configured)
3. charging group configured at the application level
4. charging group configured at the application group level
5. default charging group

3.2.2.2.8 Applications

The application context defines and assigns a description to the application names supported by the application filter entries, and assigns applications to application groups.

- Application name is a key match criteria within application QoS policy rules, which are applied to a subscribers IP traffic.
- Each application can be associated with one of the application groups provided by AA.

The AA system provides no pre-defined applications other than **Unknown**. Applications must be explicitly configured. Any protocols not assigned to an application are automatically assigned to the default **Unknown** application. Nokia provides sets of known-good application/app-group configurations upon request. Contact the technical support staff for further information.

The applications are used by AA to identify the type of IP traffic within the subscriber traffic.

The network operator can:

- Define unique applications.
- Associate applications with an application group. The application group must already be configured.

3.2.2.2.9 Application filters

Application filters (app-filter) are provided as an indirection between protocols and applications to allow the addition of variable parameters (port number, IP addresses, and so on) into an application definition. An application filter is a numbered rule entry that defines the use of protocol signatures and other criteria to define an application. Multiple rules can be used to define what constitutes an application but each rule maps to only one application definition.

The system concept of application filters is similar to IP filters. Match of a flow to multiple rules is possible and is resolved by picking the rule with the lowest entry number that matches. A flow is only ever assigned to one application.

The following criteria can be assigned to an application filter rule entry:

- unique entry ID number
- application name
- flow setup direction
- server IP address (or server IP filter list)
- HTTP port (or HTTP port list) used by HTTP proxies
- server port (or server port list)
- protocol signature
- IP protocol number

- string matches against Layer 5+ protocol header fields (for example, a string expression against HTTP header fields)

The application must be pre-configured before using it in an app-filter. After defined, the new application names can be referenced.

3.2.2.2.10 HTTP

- **HTTP protocol**

The Hyper-Text Transfer Protocol (HTTP) has become the most significant protocol used on the Internet and has expanded its role beyond web browsing with a large number of applications using HTTP for a variety of functions on both desktop and mobile devices.

AA provides the tools required by residential, mobile and business VPN service providers to accurately classify any web-based applications regardless of where the content is stored and how it is delivered. This is done by using either the default protocol signatures delivered with the AA ISA software or by defining string based signatures from the HTTP header information fields included in the HTTP request messages to further refine the detection.

- **HTTP session persistency**

HTTP can use both non persistent connections and persistent connections. Non-persistent connection uses one TCP connection per HTTP request while persistent connection can reuse the same TCP connection for multiple HTTP request to the same server.

Nowadays most applications are using HTTP/1.1 and persistent connection but HTTP/1.0 and non-persistent connections remains on older software and mobile devices.

HTTP flows are classified in a particular application using the first HTTP request of the flow only by default. Optionally, the MS-ISA offers the flexibility to classify each HTTP request within the same flow independently using **http-match-all-request** feature.

- **HTTP proxy support**

AA also supports traffic classification of HTTP between a subscriber and a web proxy. This feature is enabled by default, the ISA monitors and detects HTTP proxy flows automatically, each request within the same persistent connection to the proxy server is classified independently.

3.2.2.2.11 AA IP prefix lists

AA ISA allows the match section of session filters, AQPs entries and application filters to include matching against a configured IP filter list or lists. Each IP filter list (aka IP pools) can have up to 64 IP address entries with a configurable mask for each entry.

3.2.2.2.12 Shallow flow inspection

When application awareness is not required, but requires other AA value added functions (for example, TCP-O, DEM), significant performance can be gained by not performing pattern or behavior-based identification. A shallow inspection configuration option can disable AA Layer 7 classification to increase throughput performance for deployments that can operate using only AA Layer 3 and Layer 4 shallow flow inspection. This configuration disables all signature-based flow inspection. This configuration can be used with TCP optimization, Dynamic Experience Management, Layer 3 and Layer 4 application filter classification, and Dynamic Experience Management.

3.2.2.13 Flow attributes

Each flow attribute can be enabled for use in the CLI:

```
config>isa>application-assurance-group <x>
  [no] flow-attribute <flow-attribute>
```

The classification techniques used by the flow-attribute algorithms include:

- **behavioral machine learning**

This is statistical analysis of flow features to train classifiers, independent of the transport protocol. For example, flows for the Skype protocol have an encrypted attribute, but do not use TLS or QUIC. The confidence level depends on both the algorithm and traffic.

- **protocol based**

Stateful packet payload inspection is used against any number of protocols to satisfy the attribute conditions. For example, the TLS protocol has the encrypted attribute. The confidence level is 0 or 100.

- **app-filter based**

The AA app-DB classification of traffic into applications may be used to explicitly set flow attributes that always apply to specific app-filter entries. For example, video bearer flows for an application that match video component app-filters can be assigned that attribute. The confidence level is 0 or 100.

AQP traffic control policies may include flow attributes as a match condition to only affect traffic matching or not matching configurable flow attribute confidence level:

```
config>app-assure>group>policy>aqp>entry>match> flow-attribute <flow-attribute-name> confidence
  {lt | gte | eq } <confidence>
```

3.2.2.3 Statistics and accounting

AA statistics provide the operator with information to understand application usage within a network node.

AA XML record accounting aggregates the flow information into per application group, per application, per protocol reports on volume usage during the last accounting interval. This information is then sent to a statistics collector element for network wide correlation and aggregation into customized graphical usage reports. AA uses and benefits from the rich 7750 SR or 7450 ESS accounting infrastructure and the functionality it provides to control accounting policy details.

The following types of accounting volume records are generated and can be collected:

- per ISA group and partition record for each configured application group
- per ISA group and partition record for each configured application
- per ISA group and partition record for each configured protocol
- per each AA subscriber record with operator-configurable field content using custom AA records for operator-selected subset of protocols, applications and application groups
- per AA subscriber per each configured application record (special study mode)
- per AA subscriber per each supported protocol record (special study mode)
- per ISA AA-performance record, containing information about the traffic and resources of each ISA

- per AA partition stats record for counts of traffic by Layer 3 protocol used to transport L4 protocols. This includes TCP, UDP and NonTcpUdp carried by IPv4, IPv6, DS_Lite, 6to4/6RD, GTP, and Teredo protocols

AA supports RADIUS accounting export of per AA subscriber charging group statistics.

Each AA group:partition can be configured for AA subscribers stats export by referencing both an accounting policy (for XML statistics) or a RADIUS accounting policy. To determine how to export various counters for subscriber AA statistics, an export-using keyword is used when enabling AA subscriber level stats export to specify the export method to be used for each, whether accounting policy or RADIUS accounting policy or diameter-based usage monitoring.

Per AA flow statistics are provided as described in [Cflowd AA records](#).

See the *7450 ESS*, *7750 SR*, *7950 XRS*, and *VSR System Management Guide* for information about general accounting functionality.

3.2.2.3.1 Per-AA subscriber special study

The system can be configured to generate statistical records for each application and protocol that the system identifies for specific AA subscribers. These capabilities are disabled by default but can be enabled for a subset of AA subscribers to allow detailed monitoring of those AA subscriber's traffic.

Per-AA subscriber per-application and per-AA subscriber per-protocol records are enabled by assigning individual AA subscribers to special study service lists. The system and ISA group limit the number of AA subscribers in this mode to constrain the volume of stats generated. When an AA subscriber is in a special study mode, one record for every application or one record for every protocol that are configured in the system are generated for that subscriber. For example, if 500 applications are configured and 200 protocols are identified, 700 records per AA subscriber are generated, if the AA subscriber is listed in both the per-aa-sub-application and per-aa-sub protocol lists.

3.2.2.3.2 System aspects

AA uses the existing redundant accounting and logging capability of the 7750 SR and 7450 ESS for sending application and subscriber usage information, in-band or out-of-band. AA statistics are stored using compressed XML format with other system and subscriber statistics in compact flash modules on the redundant SF/CPMs. A large volume of statistics can be expected under scaled scenarios when per-AA subscriber statistics/accounting is enabled.

AA accounting and statistics can be deployed as part of other system functionality as long as the system's function is compatible with AA accounting or as long as the system-level statistics can become application-aware because of, for example, AA ISA-based classification. An example of this feature interaction includes volume and time-based accounting where AA-based classification into IOM queues with volume and time accounting enabled can, for instance, provide different quota/credit management for off-net and on-net traffic or white/grey applications.

3.2.2.3.3 AA XML volume statistics and accounting

AA is configured to collect and report on the following statistics when at least one AA ISA is active. The default AA statistics interval is 15 minutes.

Statistics to be exported from the node are aggregated into accounting records, which must be enabled to be sent. By default, no records are sent until enabled. Each record template type is enabled individually

to control volume of statistics to the wanted level of interest. Only non-zero records are written to the accounting files for all AA subscriber based statistics to reduce the volume of data.

The operator can further select a subset of the fields to be included in per-AA subscriber records and whether to send records if no traffic was present for a specific protocol or application, for example, sending only changed records.

Each record generated contains the record fields as described in [Table 13: AA statistics fields generated per record \(accounting file\)](#). The header row represents the record type.

Table 13: AA statistics fields generated per record (accounting file)

Record fields	Description	Group/ partition app group	Group/ partition application	Group/ partition protocol	AA subscriber custom	AA subscriber special study per app	AA subscriber special study protocol	XML name
Application Group	Name	X						data name
Application	Name		X			X		data name
Protocol	Name			X			X	data name
Aggregation Type ID	ID (can be protocol, application, charging group or application group record)				X			agg- type- name
# Active Subscribers	# of subscribers who had a flow of this category during this interval	X	X	X				nsub
# allowed flows from-sub	# of new flows that were identified and allowed	X	X	X	X	X	X	sfa
# allowed flows to- sub	As above in opposite direction	X	X	X	X	X	X	nfa
# denied flows from-sub	the # of new flows that were identified and denied	X	X	X	X	X	X	sfd
# denied flows to- sub	As above in opposite direction	X	X	X	X	X	X	nfd
# Active flows from-sub	# of flows that were either: closed, opened and closed, opened, or	X	X	X	X	X	X	saf

Record fields	Description	Group/ partition app group	Group/ partition application	Group/ partition protocol	AA subscriber custom	AA subscriber special study per app	AA subscriber special study protocol	XML name
	continued during this interval							
# active flows to-sub	As above, in opposite direction	X	X	X	X	X	X	naf
Total packets from-sub		X	X	X	X	X	X	spa
Total packets to-sub		X	X	X	X	X	X	npa
Total bytes from-sub		X	X	X	X	X	X	sba
Total bytes to-sub		X	X	X	X	X	X	nba
Total discard packets from-sub		X	X	X	X	X	X	spd
Total short flows	Number of flows with duration <= 30 seconds that completed up to the end of this interval	X	X	X	X	X	X	sdf
Total medium flows	Number of flows with duration <= 180 seconds that completed up to the end of this interval	X	X	X	X	X	X	mdf
Total long flows	Number of flows with duration > 180 seconds that completed up to the end of this interval	X	X	X	X	X	X	ldf
Total discard packets to-sub		X	X	X	X	X	X	npd
Total discard bytes from-sub		X	X	X	X	X	X	sbd
Total discard bytes to-sub		X	X	X	X	X	X	nbd
Total flows completed	# of to- and from-subscriber flows that have been completed	X	X	X	X	X	X	tfc

Record fields	Description	Group/ partition app group	Group/ partition application	Group/ partition protocol	AA subscriber custom	AA subscriber special study per app	AA subscriber special study protocol	XML name
	up to the reported interval.							
Total flow duration	Duration, in seconds, of all flows that have been completed up to the reported interval.	X	X	X	X	X	X	tfd
From AA Sub: Maximum throughput byte count	Maximum of all total byte counts recorded for throughput intervals within this accounting interval for traffic originated by AA subscriber for an application/app-group. AA ISA discarded traffic is not included.				X			sbm
From AA Sub: Packet count corresponding to the max. throughput byte count interval.	Packet count for the throughput interval with the maximum byte count value for traffic originated by AA subscriber for the application/app-group. AA ISA discarded traffic is not included.				X			spm
To AA Sub: Max throughput time slot index	UTC time that corresponds to the end of the 5-minute throughput interval where the max throughput byte count was detected.				X			smt
From AA Sub: Forwarding-class	Observed forwarding-class bits.	X	X	X	X	X	X	sfc
To AA Sub: Forwarding-class	Observed forwarding-class bits.	X	X	X	X	X	X	nfc
To AA Sub:	Maximum of all total byte counts recorded				X			nbm

Record fields	Description	Group/ partition app group	Group/ partition application	Group/ partition protocol	AA subscriber custom	AA subscriber special study per app	AA subscriber special study protocol	XML name
Maximum throughput byte count	for throughput intervals within this accounting interval for traffic originated from Network toward AA subscriber for an application/ app-group. AA ISA discarded traffic is not included.							
To AA Sub: Packet count corresponding to the max. Throughput byte count interval.	Packet count for the throughput interval with the maximum byte count value for traffic originated from network toward AA subscriber for an application / app-group. AA ISA discarded traffic is not included.				X			npm
From AA Sub: Max throughput time slot index	UTC time that corresponds to the end of the 5-minute throughput interval where the max throughput byte count was detected.				X			nmt
From AA Sub: Forwarding-class	Observed forwarding-class bits.	X	X	X	X	X	X	X
From AA Sub: Maximum throughput byte count	Maximum of all total byte counts recorded for throughput intervals within this accounting interval for all traffic originated by AA subscriber. AA ISA discarded traffic is not included.				X			sbm
From AA Sub: Packet count corresponding	Packet count for the throughput interval with the maximum byte count value for				X			spm

Record fields	Description	Group/ partition app group	Group/ partition application	Group/ partition protocol	AA subscriber custom	AA subscriber special study per app	AA subscriber special study protocol	XML name
to the max. Throughput byte count interval.	traffic originated by AA subscriber. AA ISA discarded traffic is not included.							
From AA Sub: Max throughput time slot index	UTC time that corresponds to the end of the 5-minute throughput interval where the max throughput byte count was detected.				X			smt
To AA Sub: Maximum throughput byte count	Maximum of all total byte counts recorded for throughput intervals within this accounting interval for traffic originated from network toward AA subscriber. AA ISA discarded traffic is not included.				X			nbm
To AA Sub: Packet count corresponding to the max. Throughput Byte Count interval.	Packet count for the throughput interval with the maximum byte count value for traffic originated from network toward AA subscriber. AA ISA discarded traffic is not included.				X			npm
To AA Sub: Max throughput time slot index	UTC time that corresponds to the end of the 5-minute throughput interval where the max throughput byte count was detected.				X			nmt
Forwarding Class		X						fc
App-Profile	AA subscriber application profile name				X			app- profile

Record fields	Description	Group/ partition app group	Group/ partition application	Group/ partition protocol	AA subscriber custom	AA subscriber special study per app	AA subscriber special study protocol	XML name
App-Service-Options	List of the app-service-options characteristics and values per AA subscriber				X			app-service-option

The records are generated per ISA group and partition, with an ISA group identified by the group ID (XML field name "aaGroup"), partition identified by the partition ID (XML field name "aaPart name" and per AA subscriber (if applicable) with the AA subscriber identified by the ESM, DSM, or transit subscriber name, SAP ID (XML field name "subscriber name", "sap name" or "spoke SDP ID" respectively).

The date, time, and system ID for the records are visible as part of the existing accounting log capability, therefore it does not need to be contained inside the AA records themselves.

The Forwarding Class is included in AA XML records as generally a VPN interconnection SLA is a combination of Bandwidth connection at the site level and Forwarding Class to transport the traffic over the MPLS network, by mapping the end-customer DSCP or 802.1P traffic value into a FC.

AA accounting stats of the application/application-group volume usage per forwarding class shows the exact volume of each application at the per FC level and better ties the AA reports to the VPN services and SLA.

This can also identify key applications using a non-optimal FC over a VPN or ite and allow the option for AA to remark these into a higher traffic class, with reporting per FC to show resulting use.

3.2.2.3.4 AA partition traffic type statistics

AA-ISA provides, at the AA partition level, traffic volume visibility of the Layer 3 protocols used to transport the different Layer 4 protocols. These include a traffic volume break down of TCP, UDP and Non-TCP-UDP carried by IPv4 and IPv6, DS_Lite, 6to4/6RD, GTP, and Teredo protocols.

Traffic-type statistics are broken down by "family" and "protocol":

- Family includes IPv4, IPv6, DS-Lite, 6RD/6to4, Teredo, and GTP (IP v4/v6 in v4/v6).
- Protocol includes TCP, UDP, and Other.

Therefore, AA-ISA traffic type record provides a collection of 15 sets of traffic volume (bytes) statistics figures, as shown in [Table 14: Families and protocols](#)

Table 14: Families and protocols

Family	Protocols
IPv4	TCP, UDP, Other
IPv6	TCP, UDP, Other
DS-Lite	TCP, UDP, Other (IPv4 tunneled inside IPv6)

Family	Protocols
6to4/6RD	TCP, UDP, Other (IPv6 tunneled inside IPv4)
Teredo	TCP, UDP, Other (IPv6 tunneled inside IPv4 and UDP)
v4inv4GTP	TCP, UDP, Other (IPv4 tunneled inside IPv4 GTP)
v4inv6GTP	TCP, UDP, Other (IPv4 tunneled inside IPv6 GTP)
v6inv4GTP	TCP, UDP, Other (IPv6 tunneled inside IPv4 GTP)
v6inv6GTP	TCP, UDP, Other (IPv6 tunneled inside IPv6 GTP)

These statistics are always counted. There is no configuration required to enable or disable tracking. However, the operator has the option to enable or disable export of these statistics via XML.

[Table 15: AA-partition traffic type statistics](#) lists the statistic record fields per AA partition.

Table 15: AA-partition traffic type statistics

Record name	Type	Description
nsa	cumulative	sessions admitted (to-sub)
ssa	cumulative	sessions admitted (from-sub)
nca	cumulative	chunks admitted (to-sub)
sca	cumulative	chunks admitted (from-sub)
sba	cumulative	octets admitted (from-sub)
spa	cumulative	packets admitted (from-sub)
sbd	cumulative	octets denied (from-sub)
spd	cumulative	packets denied (from-sub)
nba	cumulative	octets admitted (to-sub)
npa	cumulative	packets admitted (to-sub)
nbd	cumulative	octets denied (to-sub)
npd	cumulative	packets denied (to-sub)
sfa	cumulative	flows admitted (from-sub)
sfd	cumulative	flows denied (from-sub)
saf	intervalized	active flows (from-sub)
nfa	cumulative	flows admitted (to-sub)

Record name	Type	Description
nfd	cumulative	flows denied (to-sub)
naf	intervalized	active flows (to-sub)
tfc	cumulative	total terminated flows
tfd	cumulative	total terminated flow duration
sdf	cumulative	short duration flows
mdf	cumulative	medium duration flows
ldf	cumulative	long duration flows
sfc	cumulative	forwarding-class bitmap (from-sub)
nfc	cumulative	forwarding-class bitmap (to-sub)
tet	cumulative	number of subscribers tethered
nnt	cumulative	number of subscribers not tethered

3.2.2.3.5 Configurable AA subscriber statistics collection

Existing average volume statistics collected over an accounting interval are extended to provide the maximum volume (bytes or packets) recorded for a throughput measurement period (5 minutes) within an accounting interval. These additional statistics improve accuracy for the access-pipe right-sizing service.

Maximum throughput statistics can be enabled for the selected applications or application groups enabled for custom per AA statistics. In addition, the operator can enable (disabled by default) per AA subscriber "Max-throughput" statistics for total (/aggregate) subscriber traffic, independent of defined applications/ application-groups.

Maximum throughput statistics records are allocated from the 2048K records available for use for per subscriber records.

Maximum throughput statistics are not provided for the protocols enabled for custom per AA statistics.

3.2.2.3.6 AA-performance record for ISA load

The AA-performance statistics record provides visibility of ISA loading related statistics to allow operational monitoring and planning of ISA overload:

- provides end of reporting interval snapshot of current values of the parameters listed in below into a per AA ISA Planning record. "Current" is the value of a counter at the end of the reporting interval, for rate based values this is the ~10sec short term current rate used in CLI statistics.
- provides time-based averages during record interval of the above values: Average(I)
- provides peak values of the above values in the reporting interval: Peak(I)

The NSP Analytics provide further analysis and thresholding triggers based on these ISA statistics, suitable for long-range planning trends such as average number of subs or peak numbers of flows.

The node per-ISA planning record values are cleared on accounting read (per all accounting records). Not reading the records means that the average and peak values are the values for the last reporting interval. The time last read is indicated in the record.

The following table lists AA performance planning record fields.

Table 16: AA performance planning record fields

Parameter	Current	Average	Peak
ISA ID			
active subs (with flows)	# subs	# subs	# subs
downloaded subs	# subs	# subs	# subs
ISA AA sub stats resource allocation	# stats records		
ISA capacity cost	sum of cost of active AA subs		
ISA Transit Subs	# subs		
diverted traffic	(packets, octets)		
entered ISA	(packets, octets)		
policy discards in ISA	(packets, octets)		
congestion discards in ISA	(packets, octets)		
error discards in ISA	(packets, octets)		
policy bypass errors	(packets, octets)		
returned traffic	(packets, octets)		
Volume cflowd			
Records reported	# records		
Reports dropped	# records		
Packets sent	# packets		
Comprehensive cflowd			
Records reported	# records		
Reports dropped	# records		
Packets sent	# packets		
TCP performance cflowd			
Flows not allocated	#flows		
Records reported	# records		

Parameter	Current	Average	Peak
Reports dropped	# records		
Packets sent	# packets		
RTP performance cflowd			
Flows not allocated	#flows		
Records reported	# records		
Reports dropped	# records		
Packets sent	# packets		
Number of synchronization sources that had to be aborted	#SSRC aborted		
URL-filter			
url-filter http-requests sent	# http-requests		
url-filter - http-request errors	# http-requests		
url-filter - http-requests dropped	# http-requests		
url-filter - http-requests permitted	# http-requests		
url-filter - http-requests redirected	# http-requests		
url-filter - http-requests blocked	# http-requests		
url-filter - http default actions	# http-requests		
url-filter - subscriber count	# subs		
url-list permits	url local list #http requests allowed		
url-list redirects	url local list #http requests redirected		
url-list drops	url local list #http requests dropped		
ICAP			
icap requests	# messages		
icap request errors	# messages		
icap permits	# messages		
icap redirects	# messages		

Parameter	Current	Average	Peak
icap drops	# messages		
icap late responses	# messages		
icap average rtt	seconds		
icap tcp connections	# icap sessions		
Web Service			
webServRequests (wsreq)	# successful requests		
httpRequestErrors (wsersp)	# unsuccessful requests		
webServResponses (wsrsp)	# responses		
webServLateResponses (wslrsp)	# late responses		
webServAvgRtt (wsrtt)	Average SRTT		
webServCacheHits (wsch)	# requests served by cache		

The following table lists AA performance records.

Table 17: AA performance records

Record name	Type	Description	MIB object (if applicable)
tmo	Cumulative	Octets to MDA	tmnxBsxGrpStatusOctsToMda
tmp	Cumulative	Packets to MDA	tmnxBsxGrpStatusPktsToMda
fmo	Cumulative	Octets from MDA	tmnxBsxGrpStatusOctsFromMda
fmp	Cumulative	Packets from MDA	tmnxBsxGrpStatusPktsFromMda
dco	Cumulative	Octets discarded because of congestion in MDA	tmnxBsxGrpStatusOctsDisCongMda
dcp	Cumulative	Packets discarded because of congestion in MDA	tmnxBsxGrpStatusPktsDisCongMda
dpo	Cumulative	Octets discarded because of policy in MDA	tmnxBsxGrpStatusOctsDiscPolicy
dpp	Cumulative	Packets discarded because of policy in MDA	tmnxBsxGrpStatusPktsDiscPolicy
deo	Cumulative	Octets discarded because of error	tmnxBsxGrpStatusOctsDiscErrors
dep	Cumulative	Packets discarded because of error	tmnxBsxGrpStatusPktsDiscEnors
pbo	Cumulative	Octets policy bypass	tmnxBsxGrpStatusOctsPolicyByps

Record name	Type	Description	MIB object (if applicable)
pbp	Cumulative	Packets policy bypass	tmnxBsxGrpStatusPktsPolicyBybys
nfl	Cumulative	Number of flows	tmnxBsxGrpStatusFlows
caf	Intervalized	Current active flows	tmnxBsxGrpStatusFlowsCurrent
aaf	Intervalized	Average active flows	tmnxBsxGrpStatusFlowsAverage
paf	Intervalized	Peak active flows	tmnxBsxGrpStatusFlowsPeak
cfr	Intervalized	Current flow setup rate	tmnxBsxGrpStatusFlowSetupRate
afr	Intervalized	Average flow setup rate	tmnxBsxGrpStatusFlowSetupRateAvg
pfr	Intervalized	Peak flow setup rate	tmnxBsxGrpStatusFlowSetupRatePk
ctr	Intervalized	Current traffic rate	tmnxBsxGrpStatusTrafficRate
atr	Intervalized	Average traffic rate	tmnxBsxGrpStatusTrafficRateAvg
ptr	Intervalized	Peak traffic rate	tmnxBsxGrpStatusTrafficRatePeak
cpr	Intervalized	Current packet rate	tmnxBsxCflowdStatusPktRateCurr
apr	Intervalized	Average packet rate	tmnxBsxGrpStatusPacketRateAvg
ppr	Intervalized	Peak packet rate	tmnxBsxGrpStatusPacketRatePeak
cas	Intervalized	Current active subscribers (with flows)	tmnxBsxGrpStatusSubsCurrent
aas	Intervalized	Average active subscribers (with flows)	tmnxBsxGrpStatusSubsAverage
pas	Intervalized	Peak active subscribers (with flows)	tmnxBsxGrpStatusSubsPeak
cds	Intervalized	Current diverted subscribers	tmnxBsxGrpStatusSubsDiverted
ads	Intervalized	Average diverted subscribers	tmnxBsxGrpStatusSubsDivertedAvg
pds	Intervalized	Peak diverted subscribers	tmnxBsxGrpStatusSubsDivertedPk
rfl	Intervalized	Flows in use	tmnxBsxGrpStatusFlowResInUse
rcc	Cumulative	ISA capacity cost	tmnxBsxGrpMdaCapacityCost
rss	Cumulative	Subscriber statistics count	tmnxBsxGrpMdaStatsResourceCount
rti	Cumulative	Transit IP address count	tmnxBsxGrpMdaTransitipAddr
rtp4	Cumulative	Transit prefix v4 address count	tmnxBsxGrpMdaTransPrefV4Entr
rtp6	Cumulative	Transit prefix v6 address count	tmnxBsxGrpMdaTransPrefV6Entr

Record name	Type	Description	MIB object (if applicable)
rtp6r	Cumulative	Transit prefix v6 remote address count	tmnxBsxGrpMdaTransPrefV6RemEntr
srs	Cumulative	Seen IP, requests sent	tmnxBsxGrpStatusHCSeenIpReqSenp
srd	Cumulative	Seen IP, requests dropped	tmnxBsxGrpStatusHCSeenIpReqDrop
tsc	Cumulative	Total subscribers created	tmnxBsxGrpStatusHCSubsCreated
tsd	Cumulative	Total subscribers deleted	tmnxBsxGrpStatusHCSubsDeleted
tsm	Cumulative	Total subscribers modified	tmnxBsxGrpStatusHCSubsModified
vrr	Cumulative	Volume cflowd, records reported	tmnxBsxCflowdStatusRecReported
vrđ	Cumulative	Volume cflowd, records dropped	tmnxBsxCflowdStatusRecDropped
vps	Cumulative	Volume cflowd, packets sent	tmnxBsxCflowdStatusPktsSent
crr	Cumulative	Comprehensive cflowd, records reported	tmnxBsxCflowdStatusRecReported
crđ	Cumulative	Comprehensive cflowd, records dropped	tmnxBsxCflowdStatusRecDropped
cps	Cumulative	Comprehensive cflowd, packets sent	tmnxBsxCflowdStatusPktsSent
trr	Cumulative	TCP performance cflowd, records reported	tmnxBsxCflowdStatusRecReported
trđ	Cumulative	TCP performance cflowd, records dropped	tmnxBsxCflowdStatusRecDropped
tps	Cumulative	TCP performance cflowd, packets sent	tmnxBsxCflowdStatusPktsSent
tfn	Cumulative	TCP performance cflowd, flows but no cflowd resources available	tmnxBsxCflowdStatusFlowsNoRes
rrr	Cumulative	RTP performance cflowd, records reported	tmnxBsxCflowdStatusRecReported
rrđ	Cumulative	RTP performance cflowd, records dropped	tmnxBsxCflowdStatusRecDropped
rps	Cumulative	RTP performance cflowd, packets sent	tmnxBsxCflowdStatusPktsSent
rfn	Cumulative	RTP performance cflowd, flows but no cflowd resources available	tmnxBsxCflowdStatusFlowsNoRes

Record name	Type	Description	MIB object (if applicable)
rsr	Cumulative	RTP performance cflowd, number of synchronization sources that had to be aborted	tmnxBsxCflowdStatusHCUSupSSRCSt
res	Cumulative	srfLOW collector, records sent The new data name is the collector address and port inserted into the XML record.	tmnxBsxCflowdCollStatRecSent
hrs	Cumulative	URL filter, HTTP requests sent	tmnxBsxUrIFlTrStatsHttpRequests
hre	Cumulative	URL filter, HTTP request errors	tmnxBsxUrIFlTrStatsHttpReqErrors
hri	Cumulative	URL filter, HTTP requests dropped	n/a
hrp	Cumulative	URL filter, HTTP requests permitted	tmnxBsxUrIFlTrStatsHttpRespAllow
hrr	Cumulative	URL filter, HTTP requests redirected	tmnxBsxUrIFlTrStatsHttpRespRedir
hrb	Cumulative	URL filter, HTTP requests blocked	tmnxBsxUrIFlTrStatsHttpRespBlock
hda	Cumulative	URL filter, HTTP default actions	tmnxBsxUrIFlTrStatsHttpRespDef
irs	Cumulative	ICAP, icap requests	tmnxBsxIcapServerStatsRequests
ire	Cumulative	ICAP, icap request errors	tmnxBsxIcapServerStatsReqErrors
irp	Cumulative	ICAP, icap permits	tmnxBsxIcapServerStatsReapAllow
irr	Cumulative	ICAP, icap redirects	tmnxBsxIcapServerStatsRespRedir
ird	Cumulative	ICAP, icap drops	tmnxBsxIcapServerStatsRespBlock
ilr	Cumulative	ICAP, icap late responses	tmnxBsxUrIFlTrStatsIcapLateResp
irt	Cumulative	ICAP, icap average rtt	tmnxBsxIcapServerStatsRoundTrip
itc	Cumulative	ICAP, icap TCP connections	tmnxBsxIcapServerStatsConnEst
ifs	Cumulative	URL filter, subscriber count	n/a
rtp4r	Cumulative	Transit prefix, v4 remote address count	tmnxBsxGrpMdaTransPrefV4RemEntr
lrp	Cumulative	URL list permits	tmnxBsxUrIFlTrStatsHttpRespAllow
lrr	Cumulative	URL list redirects	tmnxBsxUrIFlTrStatsHttpRespRedir
lrd	Cumulative	ULR list drops	tmnxBsxUrIFlTrStatsHttpRespBlock
fra	Intervalized	Flow resources average	tmnxBsxGrpStatusFlowResAvg

Record name	Type	Description	MIB object (if applicable)
frp	Intervalized	Flow resources peak	tmnxBsxGrpStatusFlowResPeak
frs	Intervalized	Flow resources alarm state	tmnxBsxGrpStatusFlowResState
fre	Intervalized	Flow resources alarm count	tmnxBsxGrpStatusFlowResRsdCount
frtm	Intervalized	Flow resources alarm time	tmnxBsxGrpStatusFlowResRaisdTime
feo	Cumulative	Flow exhaust octets	tmnxBsxGrpStatusFlwResCtThruOcts
fep	Cumulative	Flow exhaust packets	tmnxBsxGrpStatusFlwResCtThruPkts
fss	Intervalized	Flow setup rate alarm state	tmnxBsxGrpStatusFlowSetupState
fse	Intervalized	Flow setup rate alarm count	tmnxBsxGrpStatusFlowSetupRsdCnt
fstm	Intervalized	Flow setup rate alarm time	tmnxBsxGrpStatusFlowSetupRsdTime
brs	Intervalized	Bitrate alarm state	tmnxBsxGrpStatusBitRateState
bre	Intervalized	Bitrate alarm count	tmnxBsxGrpStatusBitRateRsdCount
brtm	Intervalized	Bitrate alarm time	tmnxBsxGrpStatusBitRateRsdTime
prs	Intervalized	Packet rate alarm state	tmnxBsxGrpStatusPktRateState
pre	Intervalized	Packet rate alarm count	tmnxBsxGrpStatusPktRateRsdCount
prtm	Intervalized	Packet rate alarm time	tmnxBsxGrpStatusPktRateRaisedTime
ocs	Intervalized	Overload alarm state	tmnxBsxGrpStatusWaSBfFmSubState tmnxBsxGrpStatusWaSBfToSubState
oce	Intervalized	Overload alarm count	tmnxBsxGrpStatusWaSBfFmSubRsdCnt tmnxBsxGrpStatusWaSBfToSubRsdCnt
octm	Intervalized	Overload alarm time	tmnxBsxGrpStatusWaSBfFmSubRsdTm tmnxBsxGrpStatusWaSBfToSubRsdTm
oco	Cumulative	Overload cut-through octets	tmnxBsxGrpStatusOvrlDctThruOcts
ocp	Cumulative	Overload cut-through packets	tmnxBsxGrpStatusOvrlDctThruPkls
mcpua	Intervalized	Management CPU average	tmnxBsxGrpStatusMgmlCpuAvg
mcpup	Intervalized	Management CPU peak	tmnxBsxGrpStatusMgmtCpuPeak
dcpua	Intervalized	DP CPU average	tmnxBsxGrpStatusDatapathCpuAvg
dcpup	Intervalized	DP CPU peak	tmnxBsxGrpStatusDatapathCpuPeak
dcpus	Intervalized	DP CPU alarm state	tmnxBsxGrpStatusDatapathCpuState

Record name	Type	Description	MIB object (if applicable)
dcpue	Intervalized	DP CPU alarm count	tmnxBsxGrpStatusDatapathCpuRsdCt
dcpum	Intervalized	DP CPU alarm time	tmnxBsxGrpStatusDatapathCpuRSDTm

3.2.2.3.7 AA partition traffic type statistics

AA ISA provides, at the AA partition level, traffic volume visibility of the Layer 3 protocols used to transport the different Layer 4 protocols. These include a traffic volume break down of TCP, UDP and Non-TCP-UDP carried by IPv4, IPv6, DS_Lite, 6to4/6RD and Teredo protocols.

Traffic-type statistics are broken down by family and protocol:

- Family includes IPv4, IPv6, DS-Lite, 6RD/6to4, and Teredo.
- Protocol includes TCP, UDP, and Other.

Therefore, AA ISA traffic type record provides a collection of 15 sets of traffic volume (Bytes) statistics figures as shown in [Table 18: Families and protocols](#)

Table 18: Families and protocols

Family	Protocols
IPv4	TCP, UDP, Other
IPv6	TCP, UDP, Other
DS-Lite	TCP, UDP, Other (IPv4 tunneled inside IPv6)
6to4/6RD	TCP, UDP, Other (IPv6 tunneled inside IPv4)
Teredo	TCP, UDP, Other (IPv6 tunneled inside IPv4 and UDP)
v4inv4GTP	TCP, UDP, Other (IPv4 tunneled inside IPv4 GTP)
v4inv6GTP	TCP, UDP, Other (IPv4 tunneled inside IPv6 GTP)
v6inv4GTP	TCP, UDP, Other (IPv6 tunneled inside IPv4 GTP)
v6inv6GTP	TCP, UDP, Other (IPv6 tunneled inside IPv6 GTP)

These statistics are always counted. There is no configuration required to enable/disable tracking. However, the operator has the option to enable/disable export of these statistics via XML.

[Table 19: Per AA partition stats record fields](#) lists the record fields.

Table 19: Per AA partition stats record fields

Record name	Type	Description	MIB object (if applicable)
sba	cumulative	octets admitted (from-sub)	tmnxBsxTrafStatOctsAdmFmSb

Record name	Type	Description	MIB object (if applicable)
spa	cumulative	packets admitted (from-sub)	tmnxBsxTrafStatPktsAdmFmSb
sbd	cumulative	octets denied (from-sub)	tmnxBsxTrafStatOctsDnyFmSb
spd	cumulative	packets denied (from-sub)	tmnxBsxTrafStatPktsDnyFmSb
nba	cumulative	octets admitted (to-sub)	tmnxBsxTrafStatOctsAdmToSb
npa	cumulative	packets admitted (to-sub)	tmnxBsxTrafStatPktsAdmToSb
nbd	cumulative	octets denied (to-sub)	tmnxBsxTrafStatOctsDnyToSb
npd	cumulative	packets denied (to-sub)	tmnxBsxTrafStatPktsDnyToSb
sfa	cumulative	flows admitted (from-sub)	tmnxBsxTrafStatFlwsAdmFmSb
sfd	cumulative	flows denied (from-sub)	tmnxBsxTrafStatFlwsDnyFmSb
saf	intervalized	active flows (from-sub)	tmnxBsxTrafStatActFlwsFmSb
nfa	intervalized	active flows (to-sub)	tmnxBsxTrafStatActFlwsToSb
nfd	cumulative	flows denied (to-sub)	tmnxBsxTrafStatFlwsDnyToSb
naf	intervalized	active flows (from-sub)	tmnxBsxTrafStatActFlwsFmSb
afc	cumulative	total terminated flows	tmnxBsxTrafStatTermFlws
afd	cumulative	total terminated flow duration	tmnxBsxTrafStatTermFlwDur
sdf	cumulative	short duration flows	tmnxBsxTrafStatShrtDurFlws
mdf	cumulative	medium duration flows	tmnxBsxTrafStatMedDurFlws
ldf	cumulative	long duration flows	tmnxBsxTrafStatLngDurFlws
sfc	cumulative	forwarding-class bitmap (from-sub)	N/A
nfc	cumulative	forwarding-class bitmap (to-sub)	N/A

3.2.2.3.8 AA partition admit–deny statistics

At the partition level, AA provides counters that capture events associated with various application QoS policy (AQP) actions related to packet or flow drops and admit actions. These statistics are exported via XML using configured accounting policies.

When enabled at the partition level, AA reports the statistics listed below:

- **AQP drop actions**

Drop and admit counters for “to” and “from” subscriber directions are provided for the following AQP commands:

- **error-drop**

- **overload-drop**
- **tcp-validate**
- **fragment-drop all**
- **fragment-drop out-of-order**
- **gtp-sanity-drop**
- **flow policers**

Drop and admit event counters for both “to” and “from” subscriber directions for flow count and flow rate policers, operating at the system or subscriber level.
- **hit counters**

Counters for “to” and “from” subscriber directions are provided for:

 - GTP filters for each hit on entry of a GTP filter as well as drops related to the GTP maximum size and default action. The GTP maximum size and SCTP PPID range action hit counts only report drop statistics and not permit statistics.
 - SCTP filters for each hit on entry of an SCTP filter as well as hits on PPID range and default actions
 - session filters for each hit on entry within a session filter and default action

[Table 15: AA-partition traffic type statistics](#) lists the record names used for AA admit-deny statistics.

3.2.2.3.8.1 Admit–deny threshold crossing alerts

AA supports Threshold Crossing Alerts (TCAs) that can be configured against any of the statistics counters listed in [AA partition admit–deny statistics](#). A high watermark and a low watermark can be configured for each counter. After the counter value reaches the configured high watermark within any 60 second interval, an event (Trap is set) is raised. The event is cleared if the counter goes below the low watermark threshold in any subsequent 60 seconds interval.

3.2.2.3.9 RADIUS accounting AA records

AA RADIUS accounting provides per- level, AA subscriber charging group statistics as well as application-group (AG) and application statistics as part of the RADIUS accounting infrastructure. The primary use of this is to enhance RADIUS accounting with AA information useful for usage-based billing plans, providing flexibility to charge and rate application content using IP subnets, HTTP URLs, SIP URIs, and other AA-identified applications.

The system can export AA accounting statistics using accounting policy records exported with RADIUS accounting. For non-DSM subscribers, AA RADIUS accounting is AA subscriber ID-based, where the AA subscriber context IPv4 and IPv6 host addresses for the sub are not reflected in RADIUS accounting. For DSM subscribers, the AA counters are included in the BB RADIUS session which is based on the BB sub and reflects the BB host context.

AA RADIUS accounting is implemented using ALU Vendor Specific Attributes (VSAs). This provides all charging group counters for a subscriber to be exported with a common accounting session ID. The following statistics are included in each record. Accounting values are for forwarded packets:

- input octets (from-sub)
- input packets (from-sub)

- output octets (to-sub)
- output packets (to-sub)

AA RADIUS accounting is supported for ESM, DSM, transit, and SAP or spoke SDP AA-subtypes. RADIUS accounting is used to export AA charging group, app-group, and application values according to the RADIUS accounting policy interval. Charging group statistics are exported in RADIUS accounting independent of application groups (either or both can be enabled).

For DSM subscribers, RADIUS accounting records can be configured to be exported under the Broadband ISA (BB) configuration. In this case, the AA charging group, application group, application and sub aggregate (total AA traffic) counters are passed to the BB ISA for export to the BB RADIUS accounting sessions.

3.2.2.3.10 AA Gx based usage monitoring

Using 3GPP (third generation Partnership Project) diameter (Gx) functionality, AA ISA upon receiving requests from Policy and Charging Rules Function (PCRF), can monitor application usage at the subscriber's level and report back to PCRF whenever the usage exceeds the thresholds set by the PCRF.

Usage-monitoring can be used by operators to report to PCRF when:

- AA ISA detects the start of a subscriber application (by setting usage threshold to be very low).
- A pre-set usage volume per subscriber application is exceeded.

AA can monitor subscriber's traffic for any defined:

- application
- application group
- charging group

AA ISA Gx-based usage monitoring is restricted to AA ESM and transit AA subscribers' type therefore it is only supported on 7750 SR.

The AA ISA Gx usage monitoring feature builds on 3GPP Release 11 defined Application Detection and Control (ADC) Gx attributes. In addition, AA ISA is compliant with 3GPP Release 12, whereby the ADC rule functionality is integrated in the PCC rules.

AA ISA reports accumulated usage when:

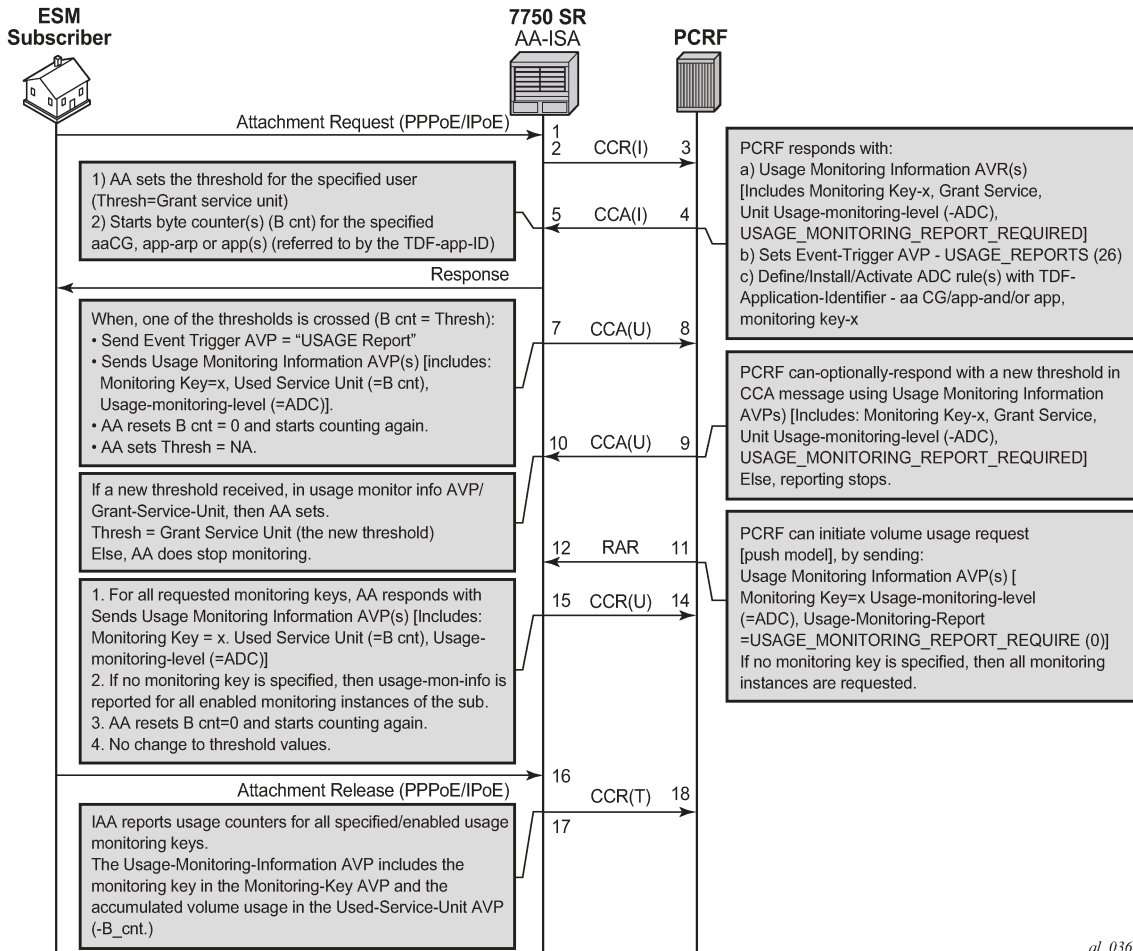
- A usage threshold is reached.
- The PCRF explicitly disables usage monitoring.
- The PCRF requests for a report.
- When the ADC or PCC rule associated with the monitoring instance is removed or deactivated.
- When a session is terminated.

An AA defined application, application group or charging group is automatically allowed to be referenced by an ADC rule for the purpose of usage monitoring only if:

- It is already selected for either XML or RADIUS per subscriber accounting.
- It is explicitly enabled by the operator for per sub statistics collection.
- Usage monitoring is enabled for the specific AA group:partition.

Figure 22: Usage monitoring illustrates the different messaging pt call flows involved in application level usage monitoring. For details about the supported AVPs used in these messages, see section [Supported AVPs](#).

Figure 22: Usage monitoring



al_0360

AA ISA (the PCEF) supports **Usage-Thresholds** AVPs that refer to the thresholds (in byte) at which point an event needs to be sent back to the PCRF ([Figure 22: Usage monitoring](#)).

No time based thresholds are supported.

AA supports **grant-service-unit** AVP using the following possible values (AVP):

- CC-Input-Octets AVP (code 412): From Subscriber total byte count threshold
- CC-Output-Octet AVP (code 414): To subscriber total byte count threshold
- CC-Total-octets AVP (code 421): Threshold of aggregate traffic (Input and Output byte counters)

As shown in [Figure 22: Usage monitoring](#) (T=7), AA sends a CCR message with a USAGE_REPORT Event-Trigger AVP to the PCRF when the usage counter reaches the configured usage monitoring threshold for a subscriber (and an application group). AA counters are reset (to zero) when the monitoring threshold is reached (and an event is sent back to PCRF). The counters, however, do not stop counting

newly arriving traffic. AA counters only include “admitted” packets. Any packets that got discarded by AA because of –say- policing actions- are not counted for usage-monitoring purposes.

The TDF-Application-Identifier AVP–within the ADC or PCC rule- refers to AA Charging group, AA application group or to an AA application.

TDF-Application-Identifiers (such as charging-groups) have to be manually entered at the PCRF to match AA charging groups configured on the 7750 SR.

If the TDF-Application-Identifiers refers to a name that is used for both a charging group and an application (or application group), AA monitors the charging group. In other words, AA charging group has higher precedence than AA application group.

3.2.2.3.11 Supported AVPs

- **ADC Rule AVP**

The ADC Rule install appears in the CCA and RAR messages from PCRF toward AA ISA.

- For installing a new ADC rule or modifying an ADC rule already installed, ADC-Rule-Definition AVP shall be used.
- For activating a specific predefined ADC rule, ADC-Rule-Name AVP shall be used as a reference for that ADC rule.

```
ADC-Rule-Definition ::= < AVP Header: 1094 >
  { ADC-Rule-Name }
  [ TDF-Application-Identifier ]
  ; AA charging group /application group / application name
  [ Flow-Status ]*
  [ QoS-Information ]*
  [ Monitoring-Key ]
  [ Redirect-Information ] ::= < AVP Header: 1085 >*
    [ Redirect-Support ] ; *
    [ Redirect-Address-Type ];*
    [ Redirect-Server-Address ];*
  [ Mute-Notification ]*

*[ AVP ]
```

The AVPs marked by an asterisk in the above example are not supported by AA ISA.

The TDF-application-Identifier field specifies a predefined AA charging group, application group or application name for which usage monitoring functionality is required (for a subscriber).

The Monitoring-Key AVP (AVP code 1066), refers to a predefined (by PCRF) USAGE Monitoring AVP.

The value of the monitoring key is random. However, it should be noted that a monitoring key instance can only be used in a single ADC rule (for example, single app/app-grp/chg-grp). While the standards allow for a monitoring instance to be referenced by one or more ADC rules, AA ISA implementation restricts this to one ADC rule. Hence, if a monitoring key is referenced in one ADC rule, it cannot be referenced by another.

- **PCC Rule AVP**

The PCC rule install appears in the CCA and RAR messages from PCRF toward AA-ISA.

- For installing a new PCC rule or modifying an PCC rule already installed, the ADC-Rule-Definition AVP shall be used.

- For activating a specific predefined ADC rule, ADC-Rule-Name AVP shall be used as a reference for that ADC rule.

```
Charging-Rule-Definition ::= < AVP Header: 1003 >
{ Charging-Rule-Name }
  [ TDF-Application-Identifier ]
  [ Monitoring-Key ]
  .....
  *[ AVP ]
```

Charging-Rule-Name is the name of the charging rule that contains a rule related to usage monitoring of a TDF_application_id has to start with: "AA-UM:". For example, AA-UM: Peer to peer traffic for APN x".

TDF-application-Identifier specifies a predefined AA charging group, application group or application name for which usage monitoring functionality is required (for a subscriber).

The Monitoring-Key AVP (AVP code 1066) refers to a predefined (by PCRF) USAGE Monitoring AVP.

The value of the monitoring key is random. However, it should be noted that a monitoring key instance can only be used in a single PCC rule (for example, single app/app-grp/chg-grp). While the standards allow for a monitoring instance to be referenced by one or more PCC rules, AA ISA implementation restricts this to one PCC rule. Hence, if a monitoring key is referenced in one PCC rule, it cannot be referenced by another.

- **Usage-Monitoring-Information AVP**

The Usage-Monitoring-Information AVP (AVP code 1067) is of type Grouped, and it contains the usage monitoring control information.

The Monitoring-Key AVP identifies the usage monitoring control instance.

```
Usage-Monitoring-Information ::= < AVP Header: 1067 >
  [ Monitoring-Key ]
  [ Granted-Service-Unit ]
  [ Used-Service-Unit ]
  [ Usage-Monitoring-Level ]
  [ Usage-Monitoring-Report ]
  [ Usage-Monitoring-Support ]
  *[ AVP ]
```

- **Monitoring-Key-AVP**

The Monitoring-Key AVP (AVP code 1066) is of type OctetString and is used for usage monitoring control purposes as an identifier to a usage monitoring control instance.

- **Granted-Service-Unit AVP**

The Granted-Service-Unit AVP shall be used by the PCRF to provide the threshold level to the PCEF.

The CC-Total-Octets AVP shall be used for providing threshold level for the total volume, or the CC-Input-Octets or CC-Output-Octets AVPs shall be used for providing threshold level for the uplink volume or the downlink volume.

```
Granted-Service-Unit ::= < AVP Header: 431 >
  [ Tariff-Time-Change ]*
  [ CC-Time ]*
  [ CC-Money ]*
  [ CC-Total-Octets ]
  [ CC-Input-Octets ]
```

```
[ CC-Output-Octets ]
[ CC-Service-Specific-Units ]*
*[ AVP ]*
```

The AVPs marked by an asterisk in the above example are not supported by AA ISA.

- **Used-Service-Unit AVP**

This AVP is used by AA_ISA (the PCEF) to provide the measured usage to the PCRF. Reporting is done, as requested by the PCRF, in CC-Total-Octets, CC-Input-Octets or CC-Output-Octets AVPs of Used-Service-Unit AVP.

The Used-Service-Unit AVP contains the amount of used units measured from the point when the service became active or, if interim interrogations are used during the session, from the point when the previous measurement ended.

```
Used-Service-Unit ::= < AVP Header: 446 >
[ Tariff-Change-Usage ]*
[ CC-Time ]*
[ CC-Money ]*
[ CC-Total-Octets ]
[ CC-Input-Octets ]
[ CC-Output-Octets ]
[ CC-Service-Specific-Units ]*
*[ AVP ]*
```

The AVPs marked by an asterisk in the above example are not supported by AA ISA.

The CC-Total-Octets AVP (AVP Code 421) is of type Unsigned64 and contains the total number of requested, granted, or used octets regardless of the direction (sent or received).

The CC-Input-Octets AVP (AVP Code 412) is of type Unsigned64 and contains the number of requested, granted, or used octets that can be/have been received from the end user.

The CC-Output-Octets AVP (AVP Code 414) is of type Unsigned64 and contains the number of requested, granted, or used octets that can be/have been sent to the end user.

- **Usage-Monitoring-Level AVP**

The Usage-Monitoring-Level AVP (AVP code 1068) is of type Enumerated and is used by the PCRF to indicate the level to which the usage monitoring instance applies.

If Usage-Monitoring-Level AVP is not provided, its absence shall indicate the value PCC_RULE_LEVEL (1).

The following values are defined (by the standard):

- SESSION_LEVEL (0) is not applicable for AA-ISA.
- If PCC_RULE_LEVEL (1) is provided within an RAR or CCA command by the PCRF, it indicates that the usage monitoring instance applies to one or more PCC rules. This is used in 3GPP Release 12 by the AA Usage Monitoring feature.
- If ADC_RULE_LEVEL (2) is provided within an RAR or CCA command by the PCRF, it indicates that the usage monitoring instance applies to one or more ADC rules. This is used in 3GPP Release 11 by the AA Usage Monitoring feature.

- **Usage-Monitoring-Report AVP**

The Usage-Monitoring AVP (AVP code 1069) is of type Enumerated and is used by the PCRF to indicate that accumulated usage is to be reported by AA ISA (the PCEF) regardless of whether a usage threshold is reached for a specific usage monitoring key (within a Usage-Monitoring-Information AVP).

If USAGE_MONITORING_REPORT_REQUIRED (0) is provided within an RAR or CCA command by the PCRF, it indicates that accumulated usage shall be reported by the PCEF.

If no monitoring keys are set, AA ISA reports all enabled monitoring instances for the subscriber.

- **Usage-Monitoring-Support AVP**

The Usage-Monitoring-Support AVP (AVP code 1070) is of type Enumerated and is used by the PCRF to indicate whether usage monitoring shall be disabled for a specific Monitoring Key.

USAGE_MONITORING_DISABLED (0) indicates that usage monitoring is disabled for a monitoring key.

- **Event-Trigger AVP (all access types)**

The Event-Trigger AVP (AVP code 1006) is of type Enumerated. When sent from the PCRF to the PCEF (AA ISA) the Event-Trigger AVP indicates an event that can cause a re-request of ADC rules. When sent from the PCEF to the PCRF the Event-Trigger AVP indicates that the corresponding event has occurred at the gateway.

USAGE_REPORT (26) is used in CCA and RAR commands by the PCRF when requesting usage monitoring at the PCEF (AA ISA). The PCRF also provides in the CCA or RAR command the Usage-Monitoring-Information AVPs including the Monitoring-Key AVP and the Granted-Service-Unit AVP.

When used in a CCR command, this value indicates that AA ISA (the PCEF) generated the request to report the accumulated usage for one or more monitoring keys. AA ISA provides the accumulated usage volume using the Usage-Monitoring-Information AVPs including the Monitoring-Key AVP and the Used-Service-Unit AVP.

The usage_report event must be set by the PCRF, otherwise AA ISA does not report usage-monitoring when a threshold is crossed.

- **Usage-Monitoring disabled**

When enabled, the PCRF may explicitly disable usage monitoring as a result of receiving a CCR from AA ISA which is not related to reporting usage, but related to other external triggers (such as subscriber profile update), or a PCRF internal trigger.

When the PCRF disables usage monitoring, AA ISA reports the accumulated usage which has occurred while usage monitoring was enabled since the last report.

To disable usage monitoring for a monitoring key, the PCRF sends the Usage-Monitoring-Information AVP including only the applicable monitoring key within the Monitoring-Key AVP and the Usage-Monitoring-Support AVP set to USAGE_MONITORING_DISABLED.

When the PCRF disables usage monitoring in a RAR or CCA command, AA ISA sends a new CCR command with CC-Request Type AVP set to the value UPDATE_REQUEST and the Event-Trigger AVP set to USAGE_REPORT to report accumulated usage for the disabled usage monitoring keys.

- **termination session**

At AA ISA subscriber's session termination, AA ISA sends the accumulated usage information for all monitoring keys for which usage monitoring is enabled in the CCR command with the CC-Request-Type AVP set to the value TERMINATION_REQUEST.

3.2.2.3.12 Cflowd AA records

AA ISA allows cflowd records to be exported to an external cflowd collector. The cflowd collector parameters (such as IP address and port number) are configured per application assurance group. Operators can choose to export cflowd records directly inband on a configurable VLAN from AA or via

the CPM, similar to the way system cflows are exported. By exporting directly inband, a higher rate of cflowd records can be exported compared to export via CPM, as inband export by-passes CPM and therefore avoids the CPM bottleneck that could potentially lead to cflowd packets discards. All cflowd records collected are exported to the configured collectors. AA ISA supports cflowd Version 10/ IPFIX.

A cflowd record is only exported to the collector after the flow is closed/terminated.

For each of the supported cflowd templates described in this section, the operator can customize which fields to include in the exported templates. For example, an operator can include the following fields: flow attributes, HTTP/SNI hostname, AA protocol/application/application-group, and DEM-related fields.

- **volume statistics**

AA ISA allows an operator to collect per flow volume statistics to be exported for any group partition. The packet sampling rate is configurable per AA- ISA-group level. For example, a packet sample rate of 10 means that one of every 10 packets is selected for volume statistics collection. If a flow has at least a single packet sampled for cflowd volume statistics, its per-flow cflowd volume record is exported to the configured collector upon flow closure.

The volume cflowd record includes flow statistics, flow related L3 to L7 information such as IP 5tuple, and a large set of fields that the operator can selectively choose from to include in the exported volume cflowd records; for example, flow duration, application ID, application group ID, device information, flow attributes, HTTP/SNI hostname, and so on.

- **comprehensive statistics**

AA ISA allows an operator to collect per flow comprehensive statistics to be exported through cflowd v10/IPFIX.

Unlike AA volume cflowd, which is packet sampled and enabled at the AA-partition level, (covering all traffic within a partition, which prohibits the use of high sampling rates), AA comprehensive flow uses a flow (instead of packet) sampling cflowd mechanism, and allows operators to target applications and application groups for sampling. Therefore, providing finer control at the application/application group level, instead of at the partition level (which is the case for volume cflowd).

The operator can decide to collect comprehensive statistics for sampled flows within an enabled group-partition application/application group. The operator can specify parameters to include in the exported comprehensive cflowd record, such as applications/application groups, host fields (applicable to HTTP traffic only), subscriber device type (when available), flow duration, device information, flow attributes, HTTP/SNI hostname, and so on, as well as other general statistics such as a flow's byte or packet counts.

The flow sampling rate is configurable on a per-ISA group level. For example, a flow sample rate of 10 means that every 10th flow is selected for comprehensive statistics collection. Anytime a flow is sampled (selected for comprehensive statistics collection), its mate flow in the reverse direction is also selected. The two flows are exported in a single cflowd record.

Per-flow comprehensive can be enabled (or disabled), using one of two configurable sampling rates, per application/app-group per partition per AA ISA-group.

Applications or application groups selected for comprehensive statistics gathering can use one of these two sampling rates. For example, important applications are assigned high sampling rates, while other applications are subjected to a lower flow sampling rate.

- **TCP application performance**

AA ISA allows an operator to collect per flow TCP performance statistics to be exported through cflowd v10/IPFIX.

The operator can decide to collect TCP performance for sampled flows within a TCP enabled group-partition-application/application-group. The flow sampling rate is configurable on per ISA-group level. For example a flow sample rate of 10 means that every 10th TCP flow is selected for TCP performance statistics collection. Anytime a flow is sampled (selected for TCP performance statistics collection) its mate flow in reverse direction is also selected. This allows collectors to correlate the results from the two flows and provide additional statistics (such as round-trip delay). Per-flow cflowd TCP performance records are exported to the configured collectors upon flow closure.

Two configurable TCP flow sampling rates are available per AA ISA group. Applications or application groups selected for TCP performance monitoring can use one of these two sampling rates. For example, important applications are assigned high sampling rates, while other TCP applications are subjected to TCP performance monitoring using a lower flow sampling rate.

Per-flow TCP performance can be enabled (or disabled), using one of two configurable sampling rates, per application/app-group per partition per AA ISA-group.

- **audio/video (A/V) application performance**

AA ISA integrates a third party audio/video performance measurement software stack to perform VoIP and video conferencing MOS-related measurements for RTP based A/V applications.

A passive monitoring technology estimates transmission quality of voice and video over packet technologies by considering the effects of packet loss, jitter and delay in addition to the impairments caused by encoding/decoding technology. A rich set of diagnostic data is provided that can be used to help network managers identify a variety of problems that could impact the quality of voice and video streams or service level agreements (SLAs).

This feature provides:

- call quality analysis using optimized ITU-T G.107, such as listening and conversational quality MOS and R-factor scores (MOS-LQ, MOS-CQ R-LQ and R-CQ)
- measurements of perceptual effects of burst packet loss and recency using ETSI TS 101 29-5 Annex E Extensions
- reporting of RTCP XR (RFC 3611, *RTP Control Protocol Extended Reports (RTCP XR)*) VoIP metrics payloads

After a flow terminates, AA ISA formats the flow MOS parameters into a cflowd record and forwards the record to a configured IPFIX /10 cflowd collector. The collector then summarizes these records using route of interest information (source/destinations). In addition, RAM provides the user with statistics (minimum, maximum, and average values) for the different performance parameters that are summarized.

Two configurable RTP flow sampling rates are available per AA ISA group. Applications or Application groups selected for RTP performance monitoring can use one of these two sampling rates. For example, important applications (such as Cisco's Telepresence video conferencing or operator's VoIP service) are assigned high sampling rates, while other RTP applications are subjected to RTP performance monitoring using a lower flow sampling rate.

Like TCP performance, per flow audio/video performance can be enabled (or disabled), using one of two configurable sampling rates, per application/app-group per partition per AA ISA-group.

The operator can decide to collect RTP A/V performance for sampled RTP flows within an RTP A/V enabled group-partition-application/application-group. The two available flow sampling rates is configurable on per ISA group level. For example a flow sample rate of 10 means that every 10th RTP flow is selected for RTP performance statistics collection. Anytime a flow is sampled (selected for RTP A/V performance statistics collection) its mate flow reverse direction is also selected. When

RTP dynamic payload types (RTP "PT") are used, only flows that use SIP to signal RTP codec can be selected for RTP performance measurement. Flows that use static RTP payload types can be selected for performance measurement regardless of the signaling channel used to setup the call.

3.2.2.3.13 vRGW nested router detection

Bridged residential gateway operators may require the vRGW to detect if and when in-home devices are used to route access to network services, thereby acting as a nested router in the home's LAN to hide multiple end devices behind the router MAC address. Data traffic coming from nested router devices is typically much higher than what an individual device generates or consumes. Nested routers also may violate terms of service for a BRG managed home on the operator's network.

For a vRGW, each device in the home that is diverted to AA becomes an **esm-mac** AA subscriber type. When AA tethering detection is enabled, an **esm-mac** AA subscriber that has traffic behavior representing multi-device traffic patterns is detected by the AA process and a "tethering state" is placed against the AA subscriber, thereby identifying a potential nested router.

The operator can install policies to handle nested router (tethering state) devices as appropriate, including but not limited to: applying different charging, blocking, rate limiting or redirecting the traffic from the device to a web portal. Per-subscriber tethering state is an indication of devices that are operating as a nested router and can be included in the AA subscriber cflowd record export.

3.2.2.4 AQP

An AQP is an ordered set of entries defining application-aware policy (actions) for IP flows diverted to a specific AA ISA group. The IP flow match criteria are based on application identification (application or application group name) but are expected to use additional match criteria such as ASO characteristic value, IP header information or AA subscriber ID, for example.

When ASO characteristic values are used in application profiles, the characteristics values can be further used to subdivide an AQP into policy subsets applicable only to a subset of AA subscribers with a specific value of an ASO characteristic in their profile. This allows to, for example, subdivide AQP into policies applicable to a specific service option (MOS iVideo Service), specific subscriber class (Broadband service tier, VPN, Customer X), or a combination of both.

A system without AQP defined has statistics generated but does not impact the traffic that is flowing through the system. However, it is recommended that an AQP policy is configured with at least default bandwidth and flow policing entries to ensure a fair access to AA ISA bandwidth/flow resources for all AA subscribers serviced by a specific AA ISA.

AQP rules consist of match and action criteria:

- Match refers to application identification determined by application and application group configuration using protocol signatures and user-configurable application filters that allow customers to create a wide range of identifiable applications. To further enhance system-wide per subscriber/service management user configurable application groups are provided.
- An AQP consists of a numbered and ordered set of entries each defining match criteria including AND, NOT and wild card conditions followed by a set of actions.

AQP Entry <#> = <Match Criteria> AND <Match Criteria> <action> <action>

- OR match conditions are supported in AQP through defining multiple entries. Multiple match criteria of a single AQP entry form an implicit AND function. An AQP can be defined for both recognized and unrecognized traffic. IP traffic flows that are in the process of being identified have a default policy

applied (AQP entries that do not include application identification or IP header information). Flows that do not match any signatures are identified as unknown-tcp or unknown-udp and can have specific policies applied (as with any other protocol).

- Actions define AA actions to be applied to traffic, a set of actions to apply to the flows like bandwidth policing, packet discards, QoS remarking and flow count or/and rate limiting.

3.2.2.4.1 AQP match criteria

Match criteria consists of any combination of the following parameters:

- the source/destination IP address and port/port-list, or IP-prefix list
- application name
- application group name
- charging group name
- one or more ASO characteristic and value pairs
- direction of traffic (subscriber-to-network, network-to-subscriber, or both, or spoke SDP)
- DSCP name
- AA subscriber (ESM, DSM, or transit subscriber, SAP or spoke SDP)
- **ip-protocol-num** field, which when used in AQP matches allows more precise control of match criteria; for example, to specify port or IP address matches specifically for either TCP or for UDP.
- subscriber tethering state

AQP entries with match criteria that exclusively use any combination of ASO characteristic and values, direction of traffic, and AA subscriber define default policies. All other AQP entries define application aware policies. Both default and application aware policies. Until a flow's application is identified only default policies can be applied.

3.2.2.4.2 AQP actions

An AQP action consists of the following action types. Multiple actions are supported for each rule entry (unlike ip-filters):

- dual or single-bucket bandwidth rate limit policer
- drop (discard)
- error drop
- flow count limit policer
- flow setup rate limit policer
- fragment drop
- HTTP enrichment
- HTTP error redirect
- HTTP notification
- HTTP redirect
- HTTPS redirect

- source mirror for an existing mirror service
- remark QoS (one or a combination of discard priority, forwarding class name, DSCP). When applied, ingress marked FC and discard priority is overwritten by AA ISA and the new values are used during egress processing (for example, egress queueing or egress policy DSCP remarking). For MPLS class-based forwarding, ingress-marked FC is still used to select an egress tunnel.
- none (monitor and report only)
- session filter
- URL-Filter (ICAP or web-service based URL filtering)
- GTP filter
- SCTP filter
- TCP MSS adjust

The value entered should be the MSS value needed for IPv4 packets. IPv6 packets are automatically adjusted to 20 bytes less reflecting the larger IP header.

- TCP validate

Any flow diverted to an ISA group is evaluated against all entries of an AQP defined for that group at flow creation (default policy entries), application identification completion (all entries), and an AA policy change (all flows against all entries as a background task). Any one flow can match multiple entries, in which case multiple actions are selected based on the AQP entry's order (lowest number entry, highest priority) up to a limit of:

- 1 drop action
- Any combination of (applied only if no drop action is selected):
 - up to 1 mirror action
 - up to 1 FC, 1 priority and 1 DSCP remark action
 - up to 4 BW policers
 - up to 12 flow policers

AQP entries the IP flow matched, that would cause the above per-IP-flow limits to be exceeded are ignored (no actions from that rule are selected).

Examples of some policy entries may be:

- Limit the subscriber to 20 concurrent Peer To Peer (P2P) flows max.
- Rate limit upstream total P2P application group to 400 kb/s.
- Remark the voice application group to EF.

3.2.2.4.3 AA policers

The rate limit (policer) policy actions provide the flow control mechanisms that enable rate limiting by application or AA subscribers.

There are six types of policers:

- Flow rate policer monitors a flow setup rate.
- Flow count limits control the number of concurrent active flows.
- Single-rate bandwidth policers monitor bandwidth using a single rate and burst size parameters.

- Dual-rate bandwidth rate policers monitor bandwidth using CIR/PIR and CBS/MBS. These can only be used at the per-subscriber granularity.
- Time of day overrides the default policer values at the specified time of day.
- Congestion override policers apply when the subscriber is in a congestion state.

After a policer is referred to by an AQP action for one traffic direction, the same policer cannot be referred to in the other direction. This also implies that AQP rules with policer actions must specify a traffic direction other than the "both" direction.

[Table 20: Policer's hardware rate steps for AA ISA](#) illustrates a policer's hardware rate steps for AA ISA.

Table 20: Policer's hardware rate steps for AA ISA

Hardware rate steps	Rate range (rate step x 0 to rate step x 127 and max)
0.5 Gbytes/s	0 to 64 Gbytes/s
100 Mb/s	0 to 12.7Gbytes/s
50 Mb/s	0 to 6.4 Gbytes/s
10 Mb/s	0 to 1.3 Gbytes/s
5 Mb/s	0 to 635 Mb/s
1 Mb/s	0 to 127 Mb/s
500 kb/s	0 to 64 Mb/s
100 kb/s	0 to 12.7 Mb/s
50 kb/s	0 to 6.4 Mb/s
10 kb/s	0 to 1.2Mb/s
8 kb/s	0 to 1 Mb/s
1 kb/s	0 to 127 kb/s

Policers are unidirectional and are named with these attributes:

- policer name
- policer type (single or dual bucket bandwidth, flow rate limit, flow count limit)
- granularity (select per-subscriber, system-wide, or ANL)
- parameters for flow setup rate (flows per second rate)
- parameters for flow count (maximum number of flows)
- rate parameters for single-rate bandwidth policer (PIR)
- parameters for two-rate bandwidth policer (CIR, PIR)
- PIR and CIR adaptation rules (min, max, closest)
- burst size (CBS and MBS)
- conformant action (allow) (mark as in-profile)

- non-conformant action (discard, or mark with options being in profile and out of profile)

Policers allow temporary over subscription of rates to enable new sessions to be added to traffic that may already be running at peak rate. Existing flows are impacted with discards to allow TCP backoff of existing flows, while preventing full capacity from blocking new flows.

Policers can be based on an AQP rule configuration to allow per-app-group, per-AA subscriber total, per AA profile policy per application, and per system per app-group enforcement.

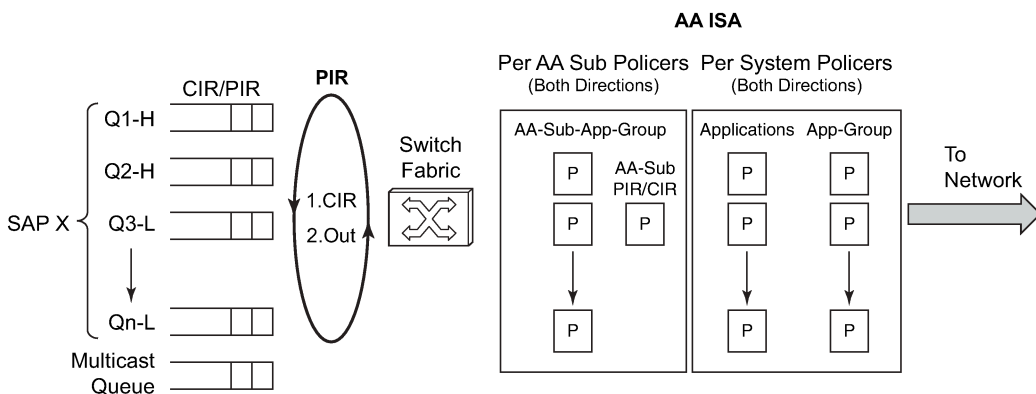
Policers are applied with two levels of hierarchy (granularity):

- **per individual AA subscriber**
 - per-AA subscriber per app group/application or protocol rate
 - per-AA subscriber per application rate limit for a small selection of applications
 - per-AA subscriber PIR/CIR. This allows the AA ISA to emulate IOM ingress policers in from-sub direction
- **per system (AA ISA or a group of AA subscribers)**
 - total protocol/application rate
 - total app group rate
- **Per ANL**
 - per-ANL per application group/application or protocol rate

Flows may be subject to multiple policers in each direction (from-subscriber-to-network or from network-to-subscriber).

In [Figure 23: From-AA subscriber application-aware bandwidth policing](#), AA policers are applied after ingress SAP policers. Configuration of the SAP ingress policers can be set to disable ingress policing or to set PIR/CIR values such that AA ISA ingress PIR/CIR are invoked first. This enables application aware discard decisions, ingress policing at SAP ingress is application blind. However, this is a design/implementation guideline that is not enforced by the node.

Figure 23: From-AA subscriber application-aware bandwidth policing

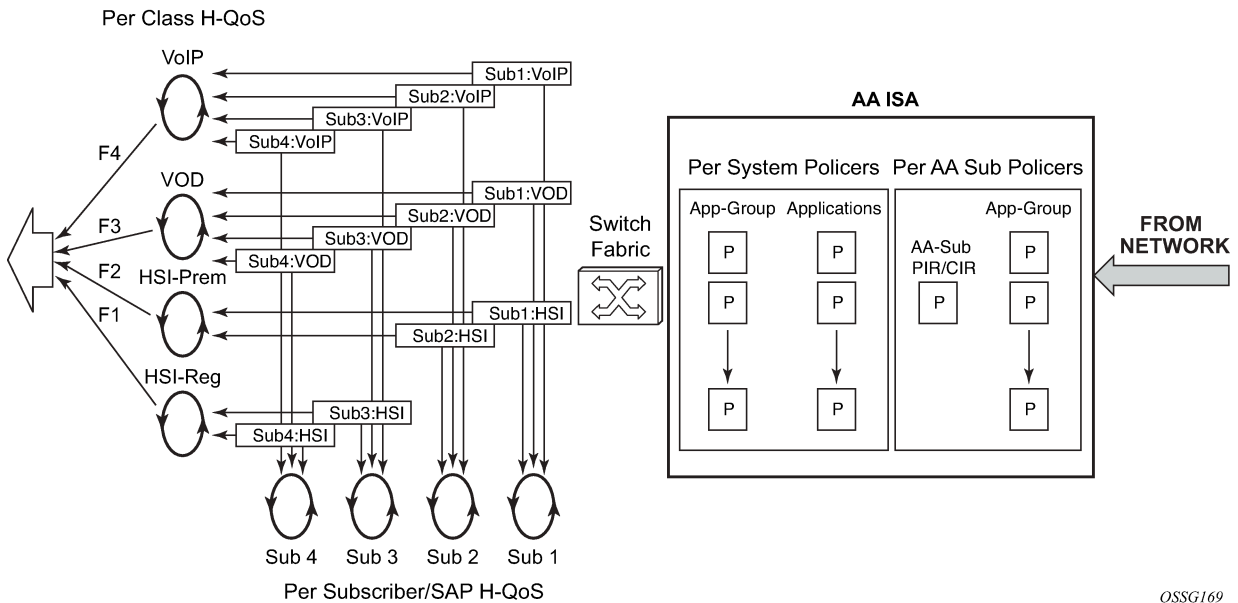


OSSG166

In the to-AA subscriber direction ([Figure 24: To-AA subscriber application-aware bandwidth policing](#)), traffic hits the AA ISA policers before the SAP egress queuing and scheduling. This allows application aware flow, AA subscriber and node traffic policies to be implemented before the Internet traffic is mixed with the

other services at node egress. AA ISA policers may remark out-of-profile traffic which allows preferential discard at an IOM egress congestion point only upon congestion.

Figure 24: To-AA subscriber application-aware bandwidth policing



OSSG169

3.2.2.4.3.1 Time of day policing adjustments

Time-of-day changes to AA policing rates are configured using time-of-day overrides in the policers. Up to eight overrides can be configured per policer, each using either a daily or weekly time-range. The adjusted policing limits are applied immediately to all flows.

3.2.2.4.3.2 Congestion override policing

As part of Dynamic Experience Management (DEM), congested Access Network Locations (ANLs) or Non-location Based DEM (NLB-DEM) can, if configured, trigger a policing override of the per-subscriber bandwidth policers.

When a subscriber is declared to be in a congestion state, the per-subscriber congestion policer rates are triggered and override any existing per-subscriber policer rates, including time-of-day policer rates. These per-sub congestion policer rates are applied for the duration of time that the subscriber is in a congestion state.

After the subscriber's state is changed to uncongested, the per-subscriber congestion policer rates are no longer applied to the subscriber's traffic. The adjusted policing limits are applied immediately to all subscriber flows.

The per-sub congestion override policers are only applicable to bandwidth policers, both single and dual leaky buckets. They are not applicable to per-subscriber flow count or flow rate policers.

To configure the per-subscriber bandwidth policer override rates, use the following commands:

- **config>app-assure>group>policer>congestion-override**

- **config>app-assure>group>policer>congestion-override>cbs**
- **config>app-assure>group>policer>congestion-override>mbs**
- **config>app-assure>group>policer>congestion-override>pir**
- **config>app-assure>group>policer>congestion-override>cib**

3.2.2.4.4 Charging filters

Charging filters are an optional mechanism to allow the charging group to be assigned based on additional criteria, including flow attributes and tethering detection. Charging filters enable the system to differentiate charging between different traffic component types even within a common application.

For example, video traffic from an application may be charged differently from other traffic types of the same application.

The match criteria are followed by an action that specifies the charging group for these match criteria. If a user configures more than one match condition, the conditions will be ANDed. Charging filters use a first-match list of entries (like app-filter entries), evaluated in numerical order, that allows flexible priority in setting the charging criteria.



Note:

The system does not support adc-start-stop notifications for charging groups assigned to charging filters.

3.2.2.5 AA TCP Optimization

AA TCP Optimization (TCPO) allows operators to leverage their existing access network deployments to provide higher throughput and improve users' experience. AA TCPO overcomes the inherited inefficiencies of TCP when faced with larger round-trip delays, high bandwidth, and non-congestion related packet loss, which are common in wireless access networks.

In particular, TCP connections have several factors that affect the effective throughput and users' experience:

- Round Trip Time (RTT) delay and available bandwidth
- TCP window size
- TCP reaction to error in both the slow start and congestion avoidance that result in random TCP packet loss not because of congestion (for example, because of interference on the access side)
- buffer bloating in the mobile access network



Note: AA TCPO is offered on the ESA platform running in lightweight subscale mode.

3.2.2.5.1 AA TCPO receiver window size

For TCP session throughputs to reach the full rate of the available bandwidth, the Receiver Window size (RWIN) must be equal to or larger than the Bandwidth Delay Product (BDP).

Calculate BDP as follows: $BDP \text{ (bytes)} = \text{total available bandwidth} \times \text{round trip time}$

AA TCPO dynamically sets RWIN toward the network side large enough to fit or exceed the maximum available bandwidth times maximum anticipated delay, therefore achieving max throughput for a specific bandwidth speed. AA TCPO performs automatic adjustment of RWIN based on the various dynamic factors such as available buffers and delays.

3.2.2.5.2 AA TCPO Round Trip Time

The basic problem limiting TCP performance arises from how the TCP congestion avoidance algorithm interacts with networks having large BDP. Congestion avoidance increases the sender's TCP window by only a single packet for each successful round-trip acknowledgment.

When the TCP window is small, increasing it by a single packet is reasonable, but if a window is very large (for hundreds of packets), then each additional round-trip acknowledgment adds a small increase to the sender's TCP window. In this case, it takes an extraordinarily large number of round trips to rebuild the TCP window in response to a single packet loss, and this leads to slow TCP behavior.

Data shows that when the bandwidth is above a specific level (4M), it has no effect on the webpage loading. On the other hand, reduced RTT directly benefits page loading.

During the TCP congestion avoidance phase, lower RTT decreases the BDP threshold, which means higher throughput for the same bandwidth and RWIN values. The throughput is then limited by the segment with the highest RTT.

AA TCPO reduces RTT, as retransmissions because of loss in the access network are handled by the AA instead of the content server. AA allocates enough buffers in the downlink direction according to the BDP.

The main benefit of reducing RTT is a faster recovery from a packet loss event.

3.2.2.5.3 Effect of packet loss in mobile networks

Traditional TCP stacks consider a packet loss event as a sign of congestion. However, in mobile and wireless networks, AA TCPO implements TCP Illinois, TCP BBR, or TCP Westwood (TCPW) toward the subscriber to address packet loss because of non-congestion events. These TCP stacks use bandwidth and delay estimates to enhance the window control (cwin) and back-off process which ensures both faster recovery and more effective congestion avoidance in RAN deployments.

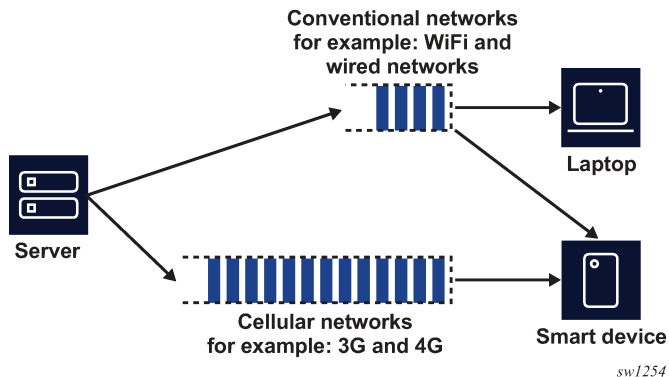
3.2.2.5.4 AA TCPO congestion window size

The initial congestion window size (cwnd) is a key component of TCP slow start, also known as the exponential growth phase. Using standard TCP, a period of three RTT is required for an average HTTP web load during slow start, which can happen after a packet loss. The majority of web transactions are on the small side, under 16 kbytes. The transactions are very short and complete before a slow start gets a chance to ramp up. To make full use of the significant bandwidth offered in mobile networks, AA TCPO sets the initial size of the cwnd, by default, high enough to allow small web transactions to take place within one RTT period. AA TCPO also offers the operator the ability to configure the initial cwnd size as part of the AA TCPO configuration policy.

3.2.2.5.5 AA TCPO buffer bloating in mobile networks

Mobile data networks suffer from buffer bloating. Buffer-bloating phenomena works against the TCP congestion control mechanism and results in poor performance. See [Figure 25: Over-buffering in mobile networks](#).

Figure 25: Over-buffering in mobile networks



The two issues that contribute to the bad effect buffer bloating has on TCP performance are:

- Standard TCP congestion control mechanism is loss based.
- Wireless networks employ over-buffering to compensate for burst traffic and channel changes (soft handover). These large buffers tend to conceal packet loss from the sender.

The preceding two factors result in a TCP sender continuing to increase its sending rate even if it has already exceeded the bottleneck link bandwidth capacity. The extra packets are all absorbed by the buffers and the rest are dropped, adding several seconds to the RTT. Networks that suffer from buffer bloating report fluctuation in RTT over time and large bursts of discards or drops.

There are two possible solutions to solve the problem with buffer bloating, while still maintaining high throughput close to BDP:

- Use a dynamic RWIN mechanism at the handsets.
- Deploy delay-based or bandwidth-based TCP sender stacks. This deployment overcomes buffer-bloating issues because these TCP congestion mechanisms use RTT or available bandwidth, which is based on RTT measurements, as in TCPW, to adjust the cwnd size.

3.2.2.5.6 AA TCPO configuration

AA TCPO is an AA session filter configurable action that refers to a configurable TCPO policy.

Using this session-filter action, operators can optionally target some IP servers for optimization, or conversely exclude specific sites (servers IP addresses) from optimization.

Using the TCPO policy, the operator can select the TCP stack to be used toward the access network, as well as TCP congestion algorithm parameters, such as cwnd and the initial slow start threshold.

Operators can enable or disable Delayed ACK (DACK) as part of the TCPO configuration policy. However, DACK timeout is not configurable. It is set at 200 ms.

Operators have the option to enable TCPO for only those TCP flows that have a network side delay above a configurable threshold. This provides an option, for example, to disable TCPO for content that is colocated with the TCP optimizer.

In addition, operators can use the ADP action to abandon TCP optimization, which disables the TC optimizer for the flows that match the configured AQP match conditions, such as detected application and application groups. This provides a mechanism for the operator to target, or exclude, specific applications or application groups from undergoing TCPO.

TCP Selective ACK (SACK) is supported by AA TCPO if both the server and client negotiate SACK successfully.

3.2.2.5.7 AA TCPO operation

AA TCPO implements the TCP new-reno stack toward the core network and TCPW, TCP BBR, or TCP Illinois configurable toward the access network.

AA TCPO does not modify or change the advertised MSS setting.

For interactions between AA TCPO and the rest of the AA feature set (such as policers), AA TCPO is logically implemented upstream from AA toward the core network. For example, if policing is enabled for a subscriber in the "from subscriber" direction, the subscriber packets are subjected to policing before any AA TCPO related action is applied to these packets. Similarly, "from subscriber" statistics maintain counts of flows before hitting the TCPO. On the other hand, in the downstream direction, the "to subscriber" statistics reflect post AA TCPO behavior, while "to subscriber" policing is applied to packets generated by AA TCPO toward the subscriber. The same concept applies to other AA features such as TCP performance measurements and HTTP enrichments.

3.2.2.5.8 AA TCPO restrictions

AA TCPO does not perform any optimization for flows that have fragmented packets, IPv6 extension headers or unsupported or unknown TCP options, such as TCP encrypt. If such packets are received in the middle of the session, AA TCPO ceases optimization for the rest of the session duration.

Nokia recommends enabling the error-drop AQP action to avoid having error packets bypassing AA TCPO and reaching the end systems. Similarly, Nokia recommends that the operator configures the drop out-of-order fragments AQP action.

Operators are recommended to configure session-filter entries to drop unsupported IPv6 extension headers to keep the state machine of AA TCPO synchronized with the TCP stacks at the endpoints.

The unsupported extension headers are:

- 139, Experimental use, Host Identity Protocol [RFC 5201]
- 140, Level 3 Multihoming Shim Protocol for IPv6, [RFC 5533]
- 253, Assigning Experimental and Testing Numbers Considered Useful [RFC 3692], Experimental Values in IPv4, IPv6, ICMPv4, ICMPv6, UDP, and TCP Headers [RFC 4727]
- 254, Assigning Experimental and Testing Numbers Considered Useful [RFC 3692], Experimental Values in IPv4, IPv6, ICMPv4, ICMPv6, UDP, and TCP Headers [RFC 4727]

If these packets arriving in the middle of a flow are not dropped, they arrive at the far end TCP stack without going through AA TCPO. Problems arise when packets of the same session are sent without having EH, arrive at TCPO, but AA TCPO has missed the earlier packet exchanges for that session.

3.2.2.6 AA Dynamic Experience Management

Dynamic Experience Management (DEM) is a feature of AA software that monitors user plane traffic to build a network-wide view of congestion at the subscriber or ANL levels. This enables real-time dynamic actions, such as rate limiting and application blocking. It provides the optimum user experience within the actual, overall network capabilities.

In situations of high network load and congestion, the subscriber quality of experience (QoE) degrades because of restricted resources across the network (such as in the radio transport layers). In this context, background traffic and real-time traffic cannot be differentiated efficiently and dynamically. The ability to differentiate background traffic from real-time traffic is important for delay-sensitive applications, such as video.

The following implementations of the DEM gateway are available, depending on the network access deployment:

- Wireless LAN DEM (WLAN-DEM) for WIFI wireless LAN gateway deployment
- Non-Location Based DEM (NLB-DEM) for any access network

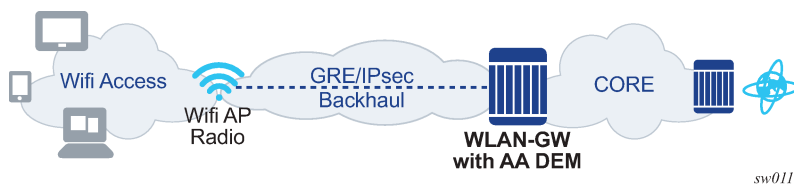
3.2.2.6.1 WLAN-DEM

Dynamic Experience Management (DEM) is a Wireless LAN Gateway (WLAN GW) capability that monitors user plane traffic to build a network-wide view of congestion on the subscriber, application, and access point radio levels. DEM enables making real-time decisions and dynamic actions, such as rate limiting or blocking of specific applications. It provides a managed, optimal user experience within the actual, overall network capabilities.

In situations of high network load and congestion, application Quality of Experience (QoE) degrades because of restricted resources across the network (for example, in radio or transport). In this context, operators cannot differentiate background traffic from real-time traffic efficiently and dynamically. This differentiation is especially important for delay-sensitive applications such as video.

In Radio Access Networks (RAN), the network congestion points are typically located in the access point WiFi radio. See [Figure 26: WiFi network congestion at AP radio](#).

Figure 26: WiFi network congestion at AP radio



Increased penetration of WiFi-enabled devices (for example, mobile handsets, tablets, laptops, and TVs) and widespread use of streaming video results in frequent data plane congestion events in WiFi networks. This congestion results in service degradation for WiFi subscribers attached to congested access points and creates challenges in implementing fair usage policies to manage network congestion in the access network.

DEM provides the capability of managing WiFi access congestion points at the WLGW to provide some level of QoS guarantees to different applications, which otherwise poses challenges as the loading of the

different access points at any point in time is different, both in quantity (Bandwidth) and application types (for example, video, web, or mail).

3.2.2.6.2 Intelligent network congestion control

DEM technology is an implementation of intelligent congestion control. If congestion is predicted or detected, the DEM gateway automatically scales back the less delay-sensitive traffic and gives priority to more delay-sensitive applications. Applications are managed to their respective resource needs to provide the best QoE. Over The Top (OTT) applications and users are managed to their respective resource needs and configured preferences.

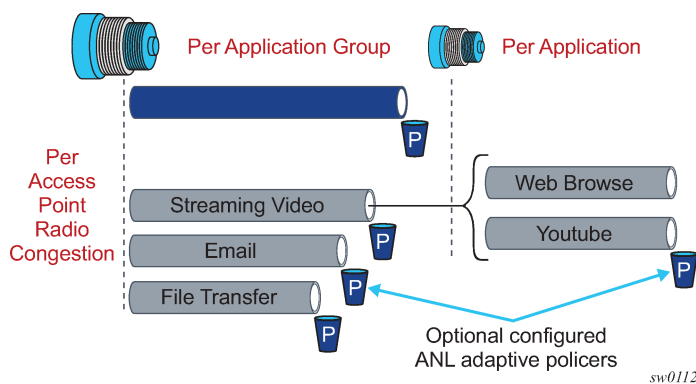
A DEM-GW builds on AA Layer 3 to Layer 7 DPI capabilities to detect applications per AA subscriber as well as per congestion point. It allows the DEM-GW AA to take intelligent actions when congestion occurs in the access network.

3.2.2.6.3 Multi-point congestion enforcement

The DEM technology allows the DEM-GW to detect congestion within the access network.

If congestion is detected at any point, DEM-GW can employ policies per application, per application group, or per subscriber to limit the impact of low-priority traffic on QoE-sensitive applications. See [Figure 27: DEM-GW multi-point congestion control](#).

Figure 27: DEM-GW multi-point congestion control



A DEM-GW is integrated directly into the WLGW using AA. The DEM-GW models the congestion points, called ANLs, that it learns from the WLGW subscriber attributes, and manages them accordingly to achieve the configured QoE/QoS target.

The DEM-GW achieves congestion control by:

- running DPI to classify flows into applications, including encrypted traffic
- dynamically learning access network congestion points and estimating their maximum capacity:
 - through real-time detection, sniffing, measurements and profiling
 - continuous monitoring of UEs locations and associating them to the right access point radio congestion points
- QoE enforcement (efficient access point radio congestion detection, localization and management provided via configurable adaptive policers)

The DEM-GW actively runs intelligent congestion control. It relies on location information relayed by WLGW sub management for Access Point MAC and VLAN.

For AP congestion detection, the DEM-GW runs an algorithm-based on measurements of Round Trip Time (RTT) to determine congestion state.

The DEM-GW uses location-awareness of all UEs to apply traffic management at specific impacted access sites, while unrestricting users during times of non-congestion. This ensures different applications within an AP radio get fair share of available resources, while controlling low-value traffic during times of congestion.

The inherited subscriber or application awareness at the DEM-GW (SSG/PGW/GGSN), when integrated with AA application detection and control, results in entitlement-based enforcements of specific applications for specified users or UEs, allowing the operators to provide differentiated services.

The end-to-end DEM solution can involve PCRF for opt-in policy control and off-line reporting platforms to facilitate some additional value-add use-cases.

3.2.2.6.4 NLB-DEM

Non-Location Based DEM (NLB-DEM) operates in any access network, within the scope of a subscriber. As such, no location information is required.

NLB DEM-GW runs a congestion detection algorithm at the subscriber level using its RTT mechanism, independent of the user location information or the location of congestion within the access network (for example, WiFi, fixed wireless, DSL, Cable, and so on). Per-subscriber bandwidth policers, if configured, can be triggered if subscriber traffic congestion is detected.

Unlike WLAN-DEM, NLB-DEM does not offer reporting or actions at the ANL level. However, NLB-DEM still offers per-subscriber policing congestion override and does not require any policy interface to pass location information. This makes NLB-DEM flexible and light weight, allowing NLB-DEM to be deployed in any type of access network.

3.2.2.6.5 Access-Network-Location policers

DEM-GW employs adaptive bandwidth policer variants of AA single leaky bucket bandwidth policers, called Access-Network-Location policers. These policers are used exclusively with DEM-GW congestion points (WLGW AP radio). They are similar to existing single bucket policers, but differ in the following aspects:

- The policer rate is configured using a ratio (/%) instead of absolute rates.
- The ratios are applied against the total estimated measured capacity of the congestion point to derive the actual policer's rate. For example, for measured capacity at congestion time of 1.5Mbps, or a configured policer rate of 30%, the actual policer rate applied: $1.5 * 30\% = 0.5$ Mbps.
- Adaptive policers are applied only in the downstream traffic direction.
- Adaptive policers run only while the associated ANL is in congestion state. No action is taken when there is no congestion.

These policers are invoked using existing AQP mechanisms that match configured parameters such as apps or app-groups and execute the configured actions.

Adaptive-policers are used to throttle traffic going through access point radios during congestion state. Multiple adapt-policers can be configured per congestion point-type (type =MAC+VLAN). For example:

- adapt-policer 1,rate=20(%), backhaul links — called from AQP entry with "email" app-group match condition

- adapt-policer 2,rate=10(%), backhaul links — called from AQP entry with OTT video app-group match condition
- adapt-policer 3,rate=0(%), backhaul links — called from AQP entry with p2p app-group match condition (this effectively drops p2p traffic during congestion)

3.2.2.6.6 DEM-GW per-subscriber congestion override policers

Although ANL adaptive policers apply to all traffic going through the ANL to maintain a positive customer experience and ensure priority traffic is not starved during congestion, they do not differentiate between identical traffic classes belonging to different subscribers.

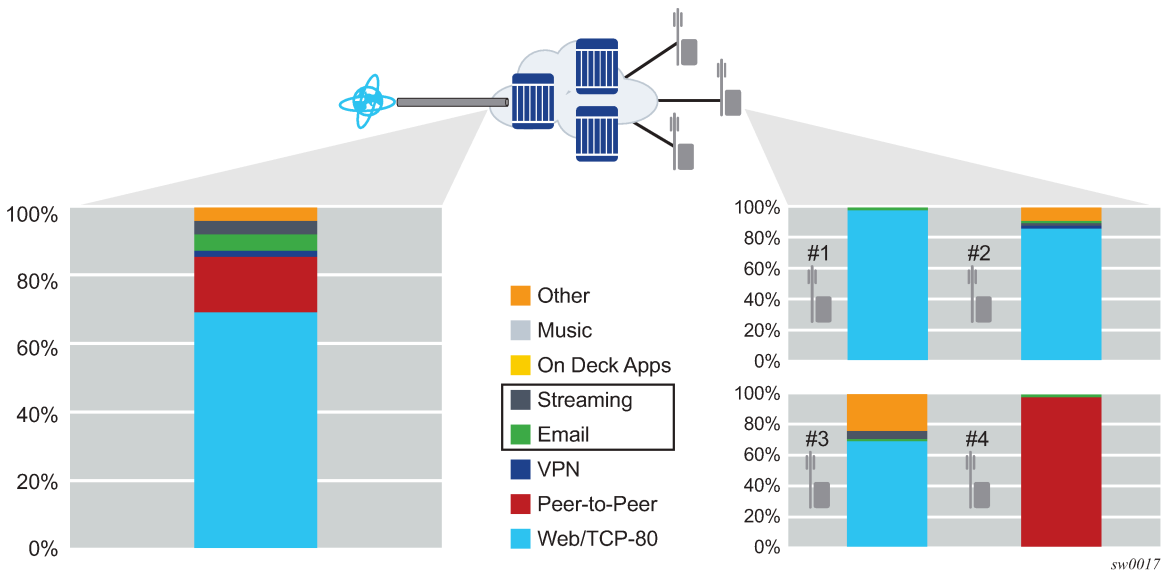
The per-subscriber policers are enabled automatically when the **congestion-override** command is enabled and the subscriber is in a congestion state. Typically, the policing rate is set so that it mostly affects heavy users. The congestion override policers can be used for all DEM use cases, such as NLB-DEM and ANL-based DEM. The operator can configure a second-stage congestion override policer. Second stage policers are applied when congestion persists even after applying the override congestion policers. Typically, the operators set a stricter bandwidth control in second stage policers to relieve congestion conditions.

Very similar to Time-of-Day (ToD) policers, per-subscriber congestion policers can be applied to all the traffic of a subscriber, specific applications, or application groups as configured in the matching section of AQPs.

3.2.2.6.7 Location-based analytics

Location-based analytics provides the operator with an accurate view of the subscriber's location (ANL) and application usage for a specified location in WiFi networks for the purpose of data-mining. See [Figure 28: Access point radio per application reporting](#).

Figure 28: Access point radio per application reporting



To provide an accurate reporting of the subscriber location via analytics tools such as the Network Services Platform, AA exports location information and congestion status in both volume and comprehensive cflowd reports. The off-line cflowd collector then allows per ANL (Access Point and AP radio) per application or application groups statistics.

3.2.2.7 AA HTTP redirect

- **AA HTTP policy redirect**

With AA ISA HTTP policy based redirect feature, when HTTP flows are blocked, the user is directed to a web portal that displays relevant messages to indicate why the HTTP traffic is blocked, such as: time of day policy to block youtube.com, top-up request, and so on.

Without HTTP policy redirect, when HTTP flows are blocked, the subscriber application retries and before it times-out.

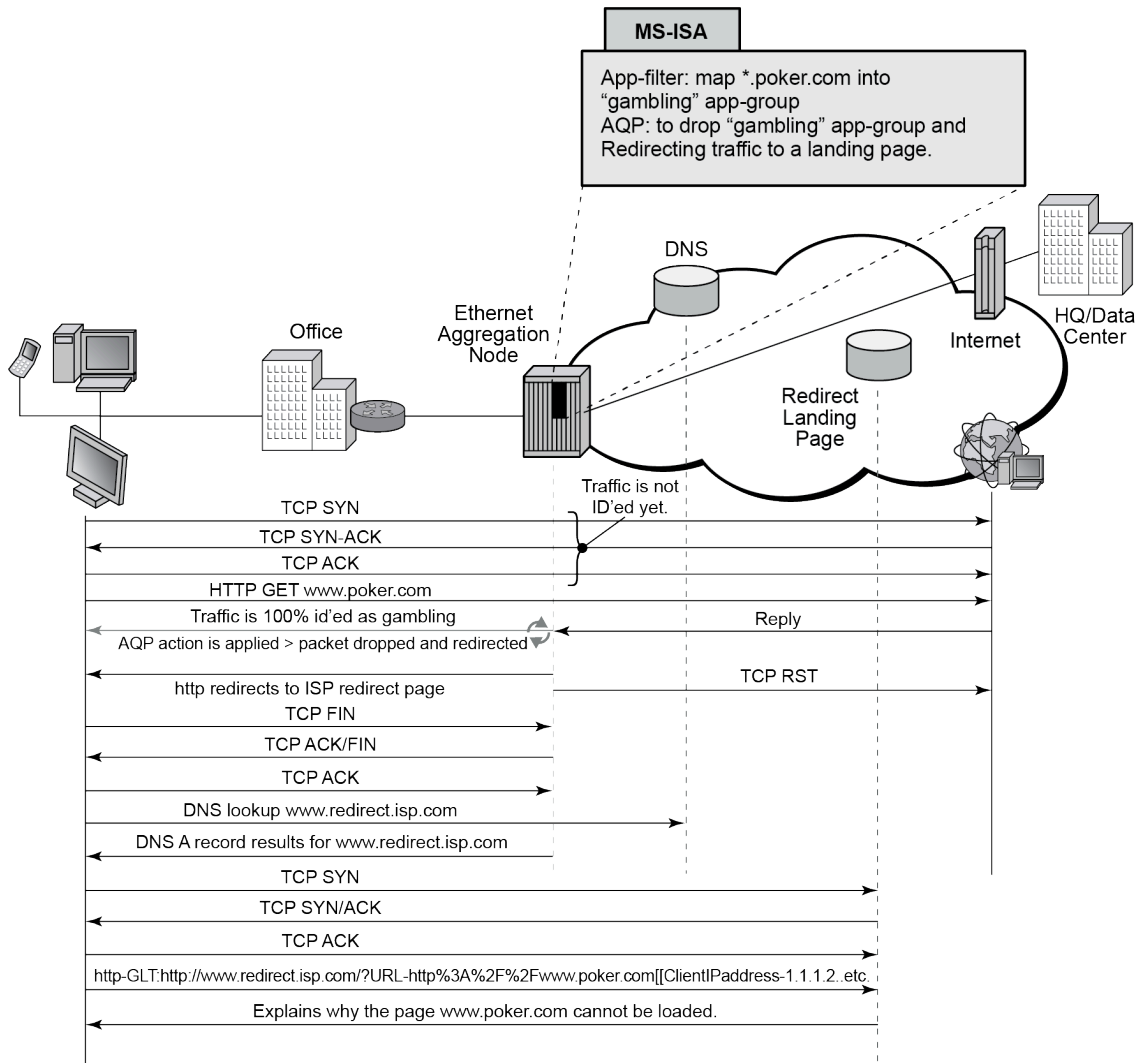
AA ISA provides full customer control to configure an AQP action that redirects traffic that matches the AQP match criteria. Hence, the HTTP redirect service can be applied at any level (application, application group, specific subscribers, specific source IP addresses) or any other AQP match criteria.

To illustrate, say the operator configures www.poker.com as a "gambling" app-group.

The operator configures AA_ISA to drop and redirect all HTTP traffic classified under "gambling" app-group to www.redirect.isp.com. When a client/subscriber initiates an HTTP GET for www.poker.com. Traffic to poker.com is dropped at the AA ISA. AA ISA issues a redirect to the client. [in the redirect, information about the user are encoded in the PATH message, such as www.poker.com, sub-ID, sub-type, reason for redirect (=AQP drop action) AA application name]. Client, unaware of the drop, responds to the redirect.

Redirect landing page describes to the user why the page at www.poker.com is not accessible. See [Figure 29: HTTP redirect because of URL block](#).

Figure 29: HTTP redirect because of URL block



al_0186

AA ISA allows the operator to configure HTTP redirect policies. An HTTP redirect policy contains, most importantly, the HTTP host to be used for the redirect. Within the AQP actions, such polices can be linked (like policers). Redirect takes place only if the AQP configured matching criteria is met and the HTTP flow is dropped (because of other AQP actions, such as "drop", flow-count/rate policers). The redirect only applies to HTTP traffic. Non-HTTP flows (even if the conditions above are met) are not redirected (no redirect for RTSP traffic).

The HTTP redirect policy includes an option for TCP-client-reset. This is used to improve the end-user experience when TCP traffic that cannot be HTTP redirected is blocked. Resetting the client TCP session avoids the client waiting for tcp session timeout. The ISA initiates a TCP reset toward the client if the AA policy results in an http-redirect with packet drop but the HTTP redirect cannot be delivered. Scenarios for this include blocked HTTPs (TLS) sessions, blocking of non-HTTP traffic, and blocking of existing flows after a policy re-evaluate of an existing subscriber. A mid-session policy change to

redirect and block traffic for a sub causes a TCP reset of existing non-http tcp sessions when the next packet for those sessions arrives. For example, when the packet is dropped.

- **AA HTTP 404 redirect**

HTTP status code-based redirect feature provides error resolution and search technology that enhances the Internet experience for end customers while generating new revenue stream for the ISP.

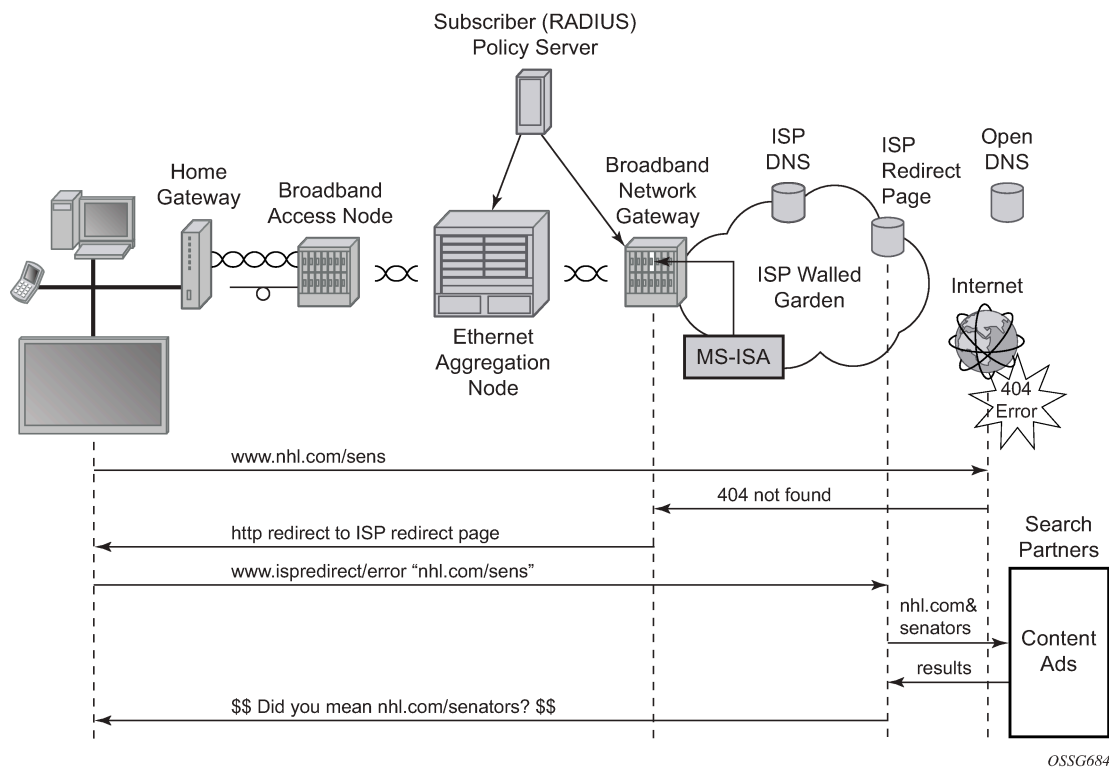
Nokia’s AA ISA HTTP status code-based redirect feature, along with its partners Barefruit or Xercole, replaces unhelpful DNS and HTTP error messages with relevant alternatives, giving the user a search solution instead of no direction. Customers benefit from an improved surfing experience as they are served relevant results that can help them find what they were looking for. The ISP, on the other hand, receives a share of the search revenue.

Every time an end-user clicks on a broken link (Page Not Found), an error page displays. Frequently, a search provider produces results, through a browser plug-in, for that user. This generates revenue for the search provider if the user clicks on a paid link.

With AA ISA HTTP status code-based redirect feature, the user sees high-quality, relevant search results. In addition, instead of the search provider receiving all of the revenue, the ISP is paid every time a user clicks on a sponsored link.

AA ISA provides full customer control to configure an AQP action that redirects traffic that matches the AQP match criteria (Figure 30: AQP actions). Hence, the HTTP redirect service can be applied at any level (application, application group, specific subscribers, specific source IP addresses) or any other AQP match criteria.

Figure 30: AQP actions



OSSG684

HTTP headers are intercepted by AA ISA on the return path from the requested web site. If the HTTP status code is a non-custom 404, then the response is replaced with JavaScript that redirects the client to the Contextual Analysis Servers (Barefruit server). This redirect contains details of the original URI that gave rise to the 404 error.

The operator can configure AA ISA HTTP 404 redirect to use either Barefruit or Xerocole partner contextual analysis servers. A redirect policy can be defined at the AA group level (similar to policers), and then referenced as many times as needed in AQP actions. The system allows a maximum of one HTTP 404 redirect policy per AA group.

3.2.2.8 AA HTTPS policy redirect

The majority of traffic is HTTPS. Many users are not aware of the differences between HTTP and HTTPS and the operator should be able to inform users why access to a web page is not allowed. If access to an HTTPS page is blocked, the user cannot know that this was the result of a policy decision.

AA supports policy-based redirection of HTTPS traffic to a landing page that displays relevant messages to indicate why the traffic is blocked. Similar to HTTP Redirect, the operator can configure an AQP to redirect traffic matching the AQP criteria to an informative web portal.

When a user attempts to access an HTTPS site, an SSL tunnel (between the user's browser and the web server) is established. AA analyses the traffic, and if the configured HTTPS redirect AQP matches, AA returns a Nokia certificate. After the certificate is accepted, the user is redirected to the informative portal.

For configuration details, see [Configuring an HTTPS Policy Redirect](#).

3.2.2.9 ICAP - large scale category-based URL filtering

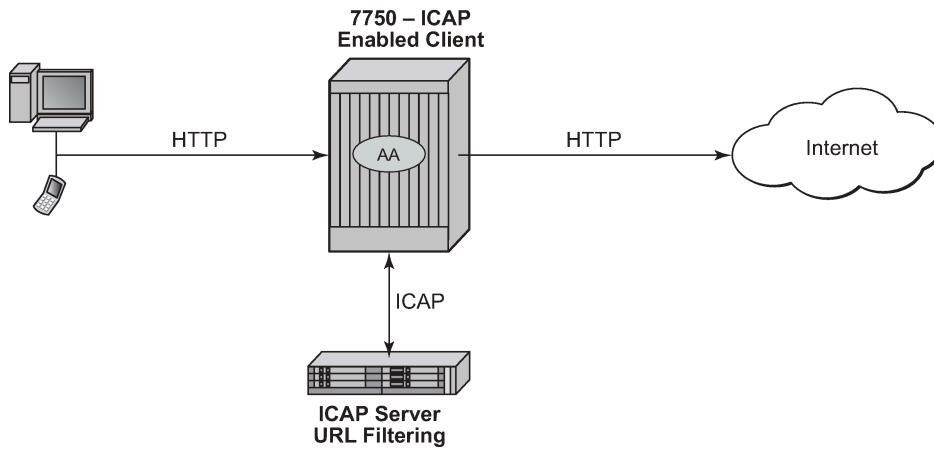
ICAP and the use of the AA-interface is only supported on the 7750 SR. Large scale URL filtering is a common content filtering requirement from broadband, mobile, and business VPN operators that allows them to solve various use cases such as:

- category based URL filtering; this is typically offered as an opt-in service by broadband or mobile operators to protect the subscribers from accessing selected category of URLs, such as, gambling, drugs, pornography, racism and so on
- managed URL filtering service for Business VPN to prevent employee from accessing specific content

AA provides both a cost efficient and best of breed content filtering solution to solve these use cases by enabling off-line dedicated web filtering servers through the Internet Content Adaptation Protocol (ICAP). Using application assurance the operator does not need to deploy costly inline filtering appliances or a limited client software solution requiring maintenance and updates for a growing number of computing devices and operating systems (for example, laptop, smartphone, smartTV, tablets).

A high level packet flow diagram of the solution is shown in [Figure 31: ICAP high level flow diagram](#). The AA ISA is the ICAP client and performs inline Layer 7 packet processing functions while the ICAP application server is used for URL filtering off-line, therefore the application server does not need to be inserted in the data flow:

Figure 31: ICAP high level flow diagram



al_0185

The 7750 SR uses the AA capabilities to extract the URL from the subscriber's HTTP/HTTPs request and send an ICAP rating request to the ICAP server along with the subscriber-id information. The ICAP server can then return an accept or redirect response based on various criteria such as subscriber profile, URL categories, allowlist, denylist, time of the day.

The ICAP response received by the 7750 SR ICAP client is used to either accept, redirect, or block the flow.

To handle the instance where an Internet server's reply arrives before the ICAP server's response, AA blocks traffic from the Internet server until the response from the ICAP server is received. This ensures that the appropriate action is always applied to the Internet traffic.

- Each HTTP request within a TCP flows are sent to the ICAP server for rating.
- HTTPs (SSL/TLS) ICAP URL-Filtering is limited to the domain name information.
- HTTPs Redirection can only be performed if the Client Hello message contains an SNI, to match the filter and proceed with the redirect action.

3.2.2.10 Web-service URL classification

AA provides a URL filtering mechanism, as an alternative to ICAP, that is used to provide content filtering. Similar to the ICAP-based solution, web-service URL classification can be used to restrict opt-in users from accessing specific (configurable) URL categories. But in contrast to the ICAP-based solution, the decision is made by the AA. There is no involvement by an external element. Nokia has partnered with a leading URL categorization service provider to help provide the categorization; however, Nokia provides the end-to-end service. The operator does not need to integrate with any third party server; all communication between the AA and the URL categorization database is transparent to the operator.

The URL categorization database contains tens of millions of URLs and is constantly being updated. Along with the categorization of new URLs, the database is also updated with malicious URLs. Operators are therefore able to offer protection against phishing and spam, and other similar sites.

Apart from enabling the operator to restrict content to opt-in users, web-service URL classification provides access to the Internet Watch Foundation (IWF) List. The IWF List contains websites which host the physical or sexual abuse of children. Using the Web-service URL classification, operators can restrict access to those sites. The IWF list is updated in real-time, so the list is always up-to-date.

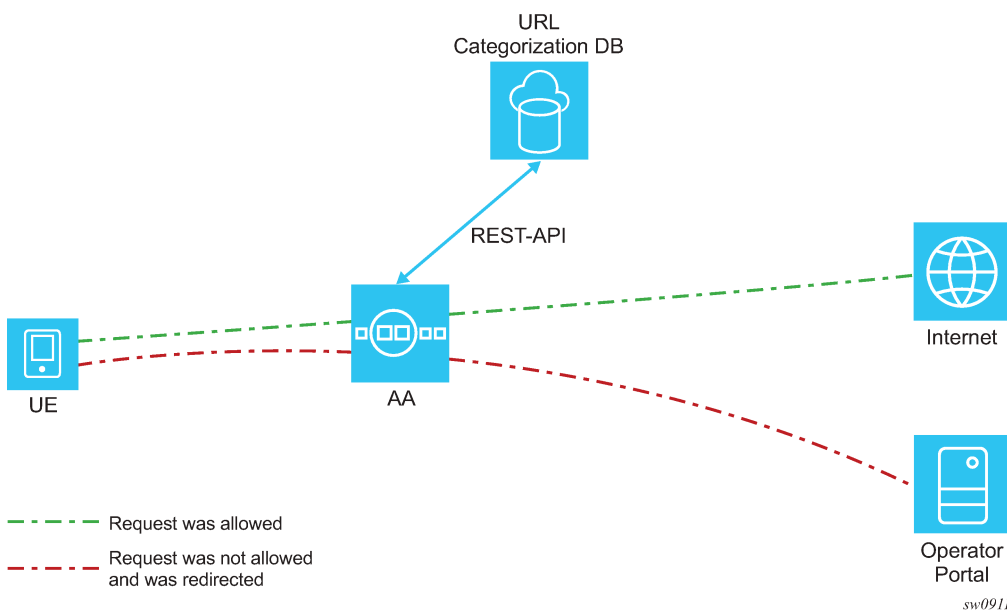
Using the web-service URL classification, operators can offer the following:

- parental control** typically offered as an opt-in service; restricts users from accessing URLs categorized as, for example, drugs, gambling, violence, and weapons
- security and threat protection** blocks URLs categorized as compromised or phishing/fraud
- IWF** blocks sites belonging to the IWF list (mandated by law in several countries)

The URL categorization database is hosted on the Cloud. The DNS service ensures that the AA always connects to the fastest server, ensuring minimum latency. In addition, AA implements a cache containing URLs and their categories. The vast majority of categorization requests are served by the cache and do not affect the user experience.

Figure 32: AA URL categorization displays a high level diagram of URL categorization in AA.

Figure 32: AA URL categorization



For every HTTP, HTTPS, HTTP2c, or QUIC request an opt-in user makes, AA first checks if the category of the requested URL is available in the cache. If available, AA checks if that category is allowed for the subscriber and acts accordingly.

If the request is not present in the cache, AA makes a Rest-API call to the URL categorization database and asks for the URL’s category. After the AA receives the response, the AA decides whether the request should be allowed or redirected.

The AA removes user-identifiable data before sending the URL to the categorization database. As an example, if the user requests “http://www.service.com/request.php?client=123”, the AA sends the following URL for categorization: “http://www.service.com/request.php”.

When the response from the URL categorization database arrives after the web server’s reply, the AA is always a restrictive filter; that is, AA always blocks the traffic in a way similar to the ICAP-based solution. This ensures that operators are never in breach of contract and a late URL categorization response does not result in a website displayed to the user.

The operator can define up to eight profiles containing categories. Using ASOs, a profile is mapped to a user and defines which URL categories are not allowed for that user.

The operator can manually set the category of a hostname to one of the supported categories. This is used in cases where the Categorization database has categorized a hostname as "Unknown" and operator either knows the category or a hostname has been misclassified. In either case, the operator must contact Nokia and request reclassification of the hostname.

[Table 21: Example URL category profile](#) shows three example profiles with different levels of restriction. The high-level profile example is the most restrictive profile. The medium-level profile is less restrictive and contains moderate category blocking. The low-level profile contains only a few URL categories to block.

Table 21: Example URL category profile

High	Medium	Low
IWF List	IWF List	IWF List
Phishing/Fraud	Phishing/Fraud	Phishing/Fraud
Spyware and Malicious Sites	Spyware and Malicious Sites	Spyware and Malicious Sites
Illegal Drugs	Illegal Drugs	—
Violence	Violence	—
Weapons	Weapons	—
Nudity	—	—
Alcohol	—	—
Criminal Skills/Hacking	—	—
Hate Speech	—	—

Profiles can be modified at any time. Profiles can be dynamically mapped to users using PCRF/AAA.

Web-service URL Filtering is supported both for HTTP (using the entire URL) and HTTPS (using the hostname only) traffic.

A detailed configuration example is available in the [Configuring web-service URL classification](#) section.



Note: An additional license is needed to use the feature.

3.2.2.10.1 URL filtering categories for web-service URL classification

[Table 22: URL filtering categories](#) lists all the URL filtering categories supported by the AA for web-service URL classification. For information about web-service URL classification, see [Web-service URL classification](#). For information about configuring web-service URL classification, see [Configuring web-service URL classification](#).

Table 22: URL filtering categories

ID	Name	Description
1	Compromised	Web pages that have been compromised by someone other than the site owner, which appear to be legitimate, but house malicious code
2	Criminal Skills/Hacking	Web pages depicting activities that violate human rights including murder, sabotage, bomb building and so on Information about illegal manipulation of electronic devices, encryption, misuse, and fraud. For example, Warez and other illegal software distribution
3	Hate Speech	Web pages that promote extreme right/left wing groups, sexism, racism, religious hate and other discrimination
4	Illegal Drugs	Web pages that promote the use or information of common illegal drugs and the misuse of prescription drugs and compounds
5	Phishing/Fraud	Manipulated web pages and emails used for fraudulent purposes, also known as phishing
6	Spyware and Malicious Sites	Web sites or software that installs on a user's computer that have the intent to collect information or make system changes without the user's consent
7	Nudity	Web pages that display full or partial nudity with no sexual references or intent
8	Mature	Web sites that are not appropriate for children, includes sites with content about alternative lifestyles, profanity and so on
9	Pornography/Sex	Web pages containing explicit sexual content unsuitable for persons under the age of 18
10	Violence	Web pages that promote questionable activities such as violence and militancy
11	Weapons	Web pages that include guns and weapons when not used in a violent manner

ID	Name	Description
12	Anonymizer	Web pages that promote proxies and anonymizers for surfing websites with the intent of circumventing filters
13	Computers and Technology	Web sites with information about computers, software, hardware, peripheral, and computers services
14	Download Sites	Web sites containing Shareware, Freeware, and other software. Also, P2P sites and software
15	Translator	Web pages that translate languages
16	Alcohol	Web pages that promote, advocate, or sell alcohol including beer, wine, and hard liquor
17	Health	Web pages supporting personal health and medical services including pages with information about equipment, procedures, and so on; not including drugs
18	Pharmacy	Web pages which include prescribed medications and information about approved drugs and their medical use
19	Tobacco	Web pages promoting the use of tobacco related products (for example, cigarettes, cigars, and pipes)
20	Gambling	Web pages which promote gambling, lotteries, casinos, and betting agencies involving chance
21	Games	Web pages consisting of computer games, game producers, and online gaming
22	Cars/Transportation	Web pages about vehicles including the selling, promoting, or discussion of vehicles
23	Dating and Relationships	Web pages that promote relationships such as dating and marriage
24	Home/Leisure	Web sites with information about home improvement and decorating, family, gardening, hobbies and so on
25	Personal Webpages	Web sites about or hosted by personal individuals. Also includes communication through blogs and guestbook servers, and

ID	Name	Description
		information about personal hobbies and activities
26	Restaurants	Web sites about food, dining, and catering services including sites that provide reviews, advertisement, or other promotion
27	Sports and Recreation	Web sites about sports teams, fan clubs and news. Sites supporting recreation activities including zoos, public recreation centers, pools, and amusement parks
28	Travel	Web pages which provide travel and tourism information, online booking, and travel services such as airlines, car rentals, and hotels
29	Government	Web sites for government organizations, departments, or agencies, including police, fire, and hospitals
30	Military	Web pages sponsored by the armed forces and government-controlled agencies
31	Non-profits	Web pages supporting clubs, communities, unions, and non-profit organizations
32	Politics and Law	Web sites that promote political parties and interest groups. Sites containing information about elections and legislation, and sites that offer legal information and advice
33	Religion	Web sites containing religious information, including information about sects, cults, occultism and religious fundamentalism
34	Education	Web sites for educational institutions and schools and for educational and reference materials, including dictionaries and encyclopedias
35	Art	Web sites about theater, museums, exhibits, photography, and digital graphic resources
36	Entertainment and Videos	Web sites about videos, TV, and motion picture, including celebrity sites and entertainment news
37	Humor	Web pages which include comics, jokes, and other humorous content

ID	Name	Description
38	Music	Web pages that include Internet radio and streaming media, musicians, bands, MP3, and media downloads
39	News	Web pages with general news information such as newspapers and magazines
40	Finance	Web sites for bank and insurance companies and other financial institutions, and for active trading of certificates and stocks
41	Internet Watch Foundation List	Web pages that show the physical or sexual abuse of children, including URLs reported by the Internet Watch Foundation (IWF); examples are child pornography, pedophilia, and child abuse
42	Shopping	Online shopping websites, catalogs, and online ordering. Also includes, auction sites, advertising, and classified ads. Excludes shopping for products and services exclusively covered by another category such as health
43	Chat/IM	Communication through chat or IM services as well as sites with information about IM communication or chat rooms
44	Community Sites	Newsgroup sites and postings including forums and bulletin boards
45	Social Networking	Social networking web pages and online communities built around communities of people where users connect to other users
46	Web-based E-mail	Web pages that enable users to send and receive e-mails through a web-accessible e-mail account
47	Portal Sites	General web pages with customized personal portals, including white and yellow pages
48	Search Engines	Web pages that enable searching of web, newsgroups, pictures, directories, and other online content
49	Online Ads	Web pages supporting advertising graphics, banners, and pop-up ad content

ID	Name	Description
50	Business/Services	General business web pages
51	Job Search	Web pages supporting job searches, agency searches, career planning, and human resources
52	Real Estate	Web pages with information about renting, purchasing, selling, or financing real estate including homes, apartments, office space, and so on.
53	Spam	Products and web pages promoted through spam techniques
54	Miscellaneous	Web pages that do not clearly fall into any other category
55	Uncategorized	Web pages for which there was no categorization provided
56	Marijuana	Web pages about marijuana or smoking marijuana. Includes web pages on legalization, medicinal use, facts and info, and pictures endorsing the drug. Does not include government sponsored web pages such as the Drug Enforcement Agency.
57	Provocative Attire	Refers to photos and videos where the person who is the subject of the photo or video is wearing sexually provocative clothing such as lingerie. Examples are bikini, bustier, negligee, and so on.

3.2.2.11 Local URL-list filtering

Service providers may need to apply network-wide URL filtering policies and either allow or prevent some content for all subscribers. AA supports both **allow-list** and **deny-list** local URL-list filtering.

Local URL-list filtering is performed on both HTTP and HTTPS traffic.

The system supports both unencrypted and OpenSSL 3DES encrypted file formats to protect the contents of the list.

Operators can specify the size of the URL list to be filtered. The size can be set to either standard or extended. If the specified URL list is configured as extended, support is provided for filtering on a larger number of URLs.

The hostnames of a local list may contain wildcards.

File format

The following characters are considered invalid and result in a failure to load the URL list:

- non-printable ASCII characters other than \n and \r
- space characters in the URL

When specifying a URL, do not include schema such as https:// or ftp://. The schema http:// is allowed but is not necessary.

The following is an example of URL list file content:

```
# Comment line for domain1 URL not using http:// schema
www.domain1.com/URI1
# Comment line for domain2 URL using http:// schema
http://domain2.com/URI2
```

3.2.2.11.1 Deny-lists

Operators want to prevent subscribers from accessing illegal content in the following situations:

- court-ordered URL takedown
- child pornography related content
- government-mandated URL takedown list

Operators can use AA to comply with these regulatory requirements, typically driven by government or court order to control the access to specific URLs hosting illegal content. In the context of child protection the operator may be required or incited to provide this filtering.

Local URL-list filtering is applied network-wide to all subscribers. This solution provides a cost-efficient method by storing the list of URLs to be filtered on the system compact flash. Therefore, using the AA-ISA ICAP functionality along with an external server is not necessary.

The ISA-AA local url-list filtering policy provides URL control capability using a list of URLs contained in a file stored on one of the system's compact flash cards. The router uses the AA capabilities to extract the URL from the subscriber's HTTP request and compares it to the list of URLs contained in this local file. If a match is found, the subscriber flow is redirected to a preconfigured web server landing page typically describing why the access to this resource was denied.

3.2.2.11.2 Allow-lists

Operators may have a list of hostnames for which they do not want to perform any web-service URL classification (or ICAP-based URL filtering). Sites may include the operator's portals or portals which the operator trusts as safe.

Similar to the deny-lists, the globally allowed sites are included in a file and in a local filtering url-filter. If a subscriber's HTTP request is included in the **allow-list**, then access to the site is allowed. The system does not check any additional URL filters (if configured).

3.2.2.11.3 URL-list update

The system supports a flexible mechanism to upgrade a local URL list automatically using either CRON or the NSP NFM-P to comply with the regulatory requirements for list upgrade frequency.

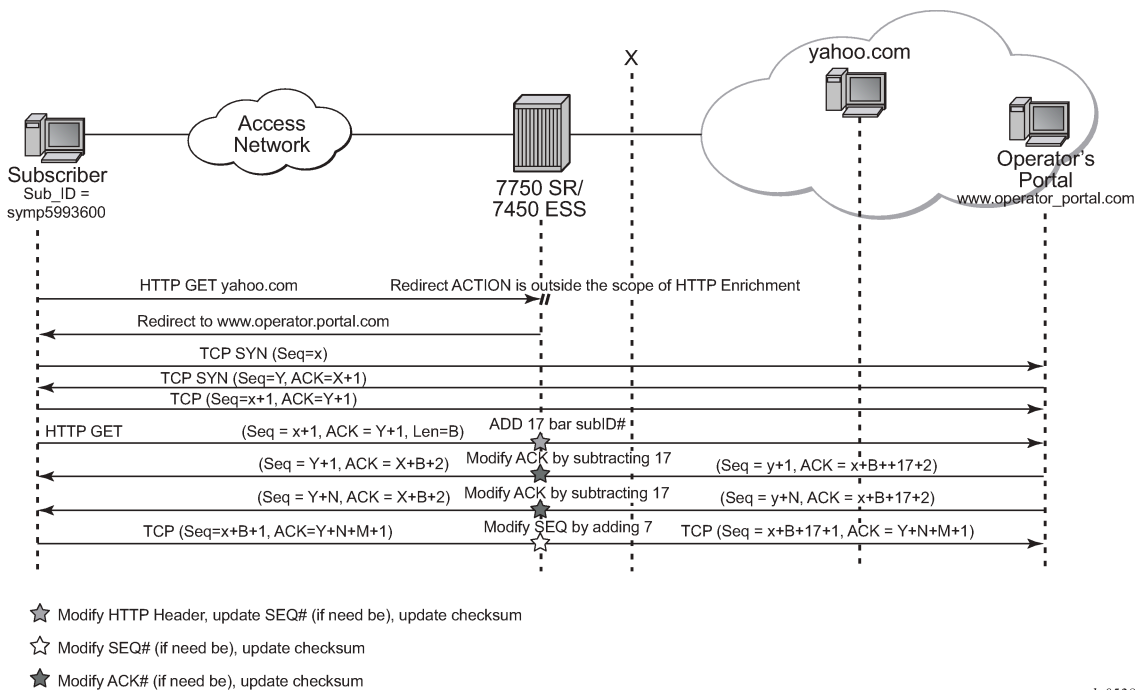
3.2.2.11.4 HTTP/HTTPS

Each HTTP request within a TCP flow is filtered by the AA ISA. For HTTPS traffic, the system extracts the domain name information contained in the SSL/TLS server name.

3.2.2.12 HTTP header enrichment

AA ISA supports modifications of the HTTP header for traffic going to specific user configured sites (URLs/IPs) to add network-based information, such as subscriber ID, to the HTTP header. These sites use this information to authenticate the user or present the user with user-specific information and services.

Figure 33: HTTP enrichment



In [Figure 33: HTTP enrichment](#), the operator configures the AA ISA to insert the subscriber ID into the HTTP header for all the HTTP traffic destined for the operator portal (designated by server IP or HTTP hostname). Traffic going to other destinations, such as yahoo.com, does not get enriched. To support this, AA introduces a new AQP action called **HTTP_enrich** that allows the operator to enrich traffic that satisfies the AQP-matching conditions.

The operator can configure multiple HTTP enrichment policies that are applied to traffic going to different destinations. For example, HTTP traffic destined for xyz.com, gets the user's IP inserted into the header, while traffic going to billing.xyz.com gets enriched with the subscriber ID and the user's IP address.

The AA ISA supports the insertion of one or more fields listed in [Table 23: HTTP header enrichment fields and formats](#) into the HTTP header.



Note: If a field that is supported in Fixed Wireless Access (FWA) mode only is used in another deployment mode, AA either enriches the header with the default value or the enrichment is not performed.

Table 23: HTTP header enrichment fields and formats

Field	Format	Supported in all deployments	Supported in FWA mode only
apn	Complete APN string	—	✓
apn-ni	APN Network Identifier (APN-NI) or the complete APN if AA cannot decode the APN properly	—	✓
billing-type	4-byte value IE: 0001	—	✓
dynamic-acr	Dynamic Anonymous Customer Reference (ACR) ¹	—	✓
imei-hyphenated ²	Subscriber's IMEI with format AABBBBBB-CCCCC-EE	—	✓
imei-sv	Subscriber's IMEI with format AABBBBBBCCCCCEE	—	✓
imsi ²	Subscriber's IMSI	—	✓
msisdn	Subscriber's MSISDN	—	✓
msisdn-ts	Subscriber's MSISDN appended with the UNIX timestamp	—	✓
msisdn-without-cc	Subscriber's MSISDN without the country code	—	✓
pgw_ggsn_address	PGW IP Address in IPv4 or IPv6 format, if the user is not in 5G UPF IP Address, if the user is in 5G	—	✓
plmin-id	MCC/MNC values as defined in 3GPP 24.301	—	✓

¹ See [ACR HTTP enrichment](#) for more information.

² AA can insert two copies of this parameter in the same HTTP header (with a different header name).

Field	Format	Supported in all deployments	Supported in FWA mode only
rat-type	4-byte value as defined in 3GPP 29.212	—	✓
static-acr	Static ACR ¹	—	✓
static-string ²	As configured by the operator	✓	—
subscriber-id	Subscriber's ID	✓	—
subscriber-ip	Subscriber's IP address in IPv4 or IPv6 format	✓	—
timestamp	8-byte UNIX timestamp when enrichment took place	—	✓
user-location	See Table 24: User-location encoding and enrichment examples	—	✓
user-location-3gpp	ULI encoded as defined in 3GPP TS2.061	—	✓
user-location-raw ²	ULI in raw format: <uli-type1>[+<uli-type2>]= <ULI data in hex> Example: x-locinfo: TAI+ECGI=1300622c46130062014adf16	—	✓

[Table 24: User-location encoding and enrichment examples](#) lists user-location encoding and enrichment examples.

Table 24: User-location encoding and enrichment examples

Identity-type format	Enrichment example ³
CGI mcc-mnc-lac-ci	x-user-loc: CellId=310-053-01a1-100a
SAI mcc-mnc-lac-sac	x-user-loc: ServiceId=310-054-01a2-200b
RAI mcc-mnc-lac-rac	x-user-loc: RoutingId=310-066-01a3-300c

³ In the enrichment examples, the header field name is configured as **x-user-loc**.

Identity-type format	Enrichment example ³
TAI mcc-mnc-tac	x-user-loc: TrackingId=310-220-01a4
ECGI mcc-mnc-eci	x-user-loc: EutranCellId=623-01-1234567
LAI mcc-mnc-lac	Not supported
Not-available	x-user-loc: User Location unavailable

**Note:**

- The ULI type can change and when it does, the mobile gateway reports it as a change to the AA-sub. AA stores and uses the latest ULI for enrichment. If the ULI type changes between the original packet transmission and the TCP retransmission, the latest ULI type is used on retransmission (that is, the ULI for retransmission is not cached).
- The mcc-mnc values of the ULI are decimal digits (BCD), while the trailing portion of the ULI are hexadecimal values.
- There can be more than one ULI identity type specified for an AA-sub, in which case the ECGI format is used. The only combination ULI supported by the mobile gateway is TAI + ECGI.
- LAI is not supported, therefore this value is not used for enrichment.
- The ULI is an optional parameter and may not be reported for an AA-sub. In this case, when enrichment is requested, the header is enriched with "User Location Unavailable".

The text preceding an inserted field is fully configurable. For example, sub-ID = 1243534666 or x-sub-ID = 1243534666.

AA supports enrichment of all HTTP methods, such as GET, POST, and so on. AA enriches HTTP traffic without having to terminate the TCP session (for example, it does not act as a proxy). In this way, the AA enrichment function does not intervene with other TCP acceleration functions or appliances that could be deployed by the operator.

For configured enriched fields, operators can optionally configure AA ISA to perform MD5 hashing, RC4 encryption, AES encryption, and anti-spoofing. When hashing is configured, the operator can optionally configure a string which is appended to the parameter before being hashed.

To perform encryption, the operator can configure an encryption key, which is used to encrypt the header values using the RC4 or AES algorithm, or install a certificate and configure the header enrichment to use that certificate. In the case of certificate-based header enrichment, the system uses the key contained in the HTTP header and performs encryption using the RSA algorithm. The system supports a maximum of 20 certificate profiles per group. The resulting string can be optionally encoded in base64 before it is inserted into the header.

Anti-spoofing, if enabled, ensures that only the fields enriched by AA are valid. Anti-spoofing is applicable only to the subscriber ID field. When configured, anti-spoofing analyzes existing headers. You must also apply the following configuration.

³ In the enrichment examples, the header field name is configured as **x-user-loc**.

Example: MD-CLI

```
[ex:/configure application-assurance group 1 http-enrich "example"]
A:admin@node-2# info
  admin-state enable
  field "msisdn" {
    anti-spoof true
    name "x-msisdn"
  }
```

Example: classic CLI

```
A:node-2>config>app-assure>group>http-enrich$ info
-----
      field "msisdn"
          name "x-msisdn"
          anti-spoof
      exit
      no shutdown
-----
```

If the packet contains a header named "x-msisdn", AA appends an "X" character to it, and the name of the header becomes "x-msisdnX". AA subsequently enriches the MSISDN header, which ensures that the portals recognize only the header inserted by AA.



Note: AA anti-spoofs up to six headers in the same request, but you can configure anti-spoofing for more than six headers. If the HTTP header contains more than six headers that are anti-spoofed, the connection resets. The maximum of six applies to both different headers and headers of the same type (for example, an MSISDN header inserted seven times).

AA statistics reflect post header enrichment packet sizes.

3.2.2.12.1 HTTP enrichment exceptions

AA HTTP enrichment functionality has the following exceptions:

- To handle the case of TCP retransmission, AA ISA implements an enrichment window of size = 5. If a retransmission of a packet occurs outside the last five enriched packets, no enrichment takes place.
- Corrupted packets; AA ISA-cut-through and out-of-order fragments are not enriched.
- Out-of-sequence packets are not enriched. For example, if AA –ISA receives out-of-sequence HTTP requests: REQ2,REQ1,REQ3; only REQ2 and REQ3 can be enriched.
- No enrichment takes place if, by enriching, the resulting packet size exceeds the configured MTU size. AA ISA does not perform fragmentation. Verify the maximum HTTP enriched packet size configured using the following command, and ensure that the MTU configured on all interfaces is greater:

– **MD-CLI**

```
configure isa application-assurance-group http-enrich-max-packet-size
```

– **classic CLI**

```
configure isa application-assurance-group http-enrich-max-pkt
```

If the MTU is exceeded, the packet is forwarded but not enriched.

- The length of an encrypted header is directly analogous to the length of the encryption key. If a 2048-bit key is used, the encrypted header becomes 512 bytes long. Operators must be cautious when defining the key length and selecting which fields are encrypted and enriched to ensure that the configured MTU size is not exceeded.
- AA ISA does not support header enrichment for WAP1.x, RTSP or SIP headers.
- AA ISA does not support header enrichment for L2 services.
- AA TCP performance measurements cannot coexist with HTTP enrichment. Enriched flows are ineligible for TCP performance sampling. If a flow is selected for TCP performance measurements and is later enriched, then TCP performance measurements cease to continue.
- Enrichment can be applied as an action to any AQP entry, subject to the following conditions:
 - The matching conditions for the AQPs cannot include a specific HTTP protocol (such as, protocol eq HTTP_video). In other words, applications that require a specific HTTP protocol type (such as video or Flash) are not considered for enrichment.
 - Within the same AQP entry, the enrichment action cannot coexist with any other AQP action (such as mark or police).
 - The AQP hit counter is not updated based on executing an HTTP enrichment action of an AQP.
- If it cannot extract the APN Network Identifier (APN-NI), AA performs enrichment using the entire APN string.

3.2.2.13 ACR HTTP enrichment

Operators can provide an ID to portals using a unique identifier, without exposing the user's secure information (for example, the MSISDN). For this purpose, AA supports header enrichment with Anonymous Customer Record (ACR) of two types: static ACR and dynamic ACR.

A static ACR is always the same for a user. The content provider is not able to retrieve the user's MSISDN, but can track the number of times the same user has connected. To make the encryption result deterministic for the static ACR, the encryption must have padding with null ASCII characters. This ensures the same sequence of bytes is produced every time.

A dynamic ACR is created during session establishment and remains the same while the session is active. A new dynamic ACR is generated when the user reconnects. With a dynamic ACR, the operator cannot track the user's MSISDN or the number of times the user accessed the portal.

[Table 25: ACR formats](#) describes how an ACR is constructed.



Note: The exceptions mentioned in [HTTP enrichment exceptions](#), also apply in ACR HTTP enrichment.

Table 25: ACR formats

Code	Description	Format	Length	Notes
CC	Country Code	NUM	3 digits	Example: 234 (UK)
NC	Network Code	NUM	3 digits	Example: 015
T	ACR Type	CHAR	4 characters	"STAT" or "DYNM"

Code	Description	Format	Length	Notes
RC	Date and time of transactions	"RSV" CCYYMMD DT hh:mm:ssZ	23 characters	The string "RSV" followed by the date time (ISO8601 date), followed by the character "Z" with no spaces Example: RSV2009-07-09T15:51:15Z
Encrypted	Encryption of the MSISDN and timestamp	—	344 characters	STAT = encrypted MSISDN DYNM = encrypted [ISO8601 date + MSISDN] The timestamp defines the creation time of the ACR The format of MSISDN is without '+' and leading '0'

3.2.2.14 HTTP in browser notification

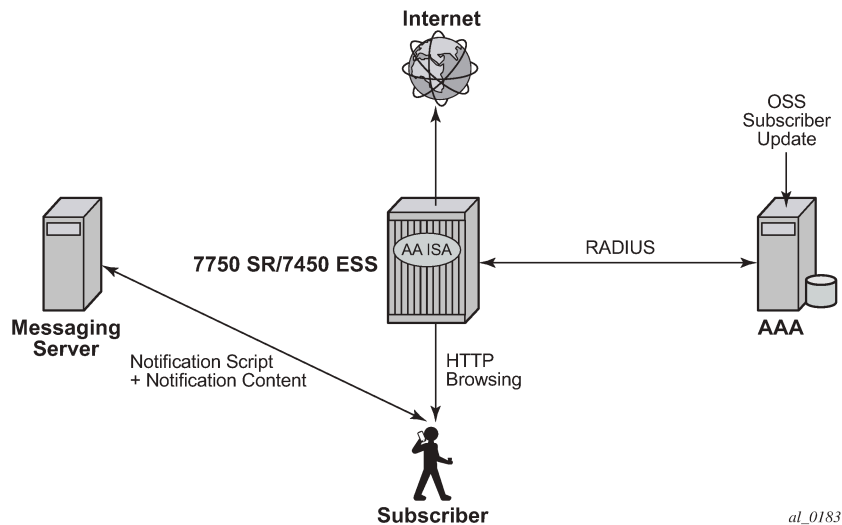
The AA ISA HTTP notification policy-based feature enables the operator to send in browser notification messages to their subscribers. The notification format can either be an overlay, a web banner, or a splash page, which makes HTTP notification less disruptive than standard HTTP redirection for the subscriber; both the original content and the notification message can be displayed at the same time while browsing.

There is a wide range of notification use cases in Broadband and WiFi networks to use this functionality such as fair use policy threshold warning, marketing and monetization messages, late bill payment notice, copyright infringement notice and operational outages.

The solution is based on two primary components, the AA ISA responsible for specific packet manipulation and a messaging server. The messaging server controls the message format and its content while the AA ISA modifies selected HTTP flows so that the subscriber transparently downloads a script located on the messaging server. This script is then executed by the web browser to display the notification message. The AA ISA only select specific HTTP request flows meeting the criteria of a web browser session compatible with in browser notification messages.

A high level view of the typical network elements involved in HTTP in browser notifications are described in [Figure 34: HTTP in browser notification - high level](#).

Figure 34: HTTP in browser notification - high level



The AA ISA provides full subscriber control to configure an AQP action enabling HTTP notification policy based on specific app-profile attributes (ASO characteristics), application, or application group. The operator can dynamically modify the subscriber policy from the policy manager to enable or disable HTTP notification during the lifetime of the subscriber.

- **notification interval**

The notification can be configured to be displayed either once during the lifetime of the subscriber or at configured minimum interval (in minutes). When an interval in minutes is selected, the subscriber continues to receive notifications messages while browsing.

- **success verification**

The system identifies successful and failed notifications. In the event the notification is not successful, the system automatically retries notifying the subscriber at the next flow that meets the criteria of a web browser session.

- **HTTP notification example**

To illustrate how HTTP notification works, the steps below describe a typical usage quota notification example with a subscriber reaching its monthly quota:

- AAA identifies that a particular subscriber is now over its quota.
- A RADIUS CoA message is sent from the AAA to the 7750 SR to modify the subscriber app-profile to enable HTTP notification.
- The AA ISA modifies the subscriber profile and enable HTTP notification for this subscriber.
- The notification message is displayed in the subscriber web browser while browsing (in the form of an overlay or web banner). The content of the notification includes a link to the operator web portal to acknowledge the reception of the overage notification.
- Until the subscriber clicks on the acknowledgment link, the AA ISA continues to execute the same policy so that notification messages are displayed in the subscriber web browser at the configured interval.

- After the subscriber has clicked on the link provided in the notification message, the provider OSS system updates the AAA which then sends a new CoA message to the 7750 SR to modify the subscriber app-profile.
- The AA ISA modifies the subscriber app-profile and disables HTTP notifications for this subscriber.

- **HTTP notification customization through RADIUS VSA**

The operator can customize the notification message per subscriber through the use a new radius VSA returned either at the subscriber creation time or within a CoA. This new VSA is a string appended automatically at the end of the script-url request made by the subscriber toward the messaging server, and it is not interpreted by the AA ISA. When received by the messaging server, it can be used to return specific content to the subscriber.

As an example, the HTTP Notification can be customized using the RADIUS VSA to display location based information, and the messaging server can use this data to display content based on the wanted location:

- Alc-AA-Sub-Http-Uri-Param RADIUS VSA: location=SohoStation
- Configured Script-URL: http://10.1.1.1/notification.js
- Subscriber HTTP request to the messaging server:
http://10.1.1.1/notification.js?subId=<aa-subscriber-id>&VSA=&location=SohoStation

3.2.2.15 AA firewall

The AA firewall (FW) feature extends AA ISA application level analysis to provide an in-line integrated stateful service that protects subscribers from malicious security attacks. Using the AA stateful packet filtering feature combined with AA Layer 7 classifications and control empowers operators with advanced, next generation firewall functionalities that integrated are within. AA stateful firewall and application firewall run on the AA ISA. In a stateful inspection, the AA FW does not only inspect packets at Layers 3 — 7, but also monitors and keeps track of the connection's state. If the operator configures a "deny" action within a session filter then the matching packets (matching both the AQP and associated session filter match conditions) are dropped and no flow session state/context is created.

AA FW can be used in all deployments of AA ISA (on the related diverted AA subscriber context):

- BNG (ESM)
- WLAN Gateway (ESM or DSM)
- Transit-subscriber (SAP or spoke-sdp)
- Business AA (SAP or spoke-sdp)
- Gi firewall (SAP or spoke-sdp)
- SeGW firewall (SAP)
- GRX FW firewall (SAP)

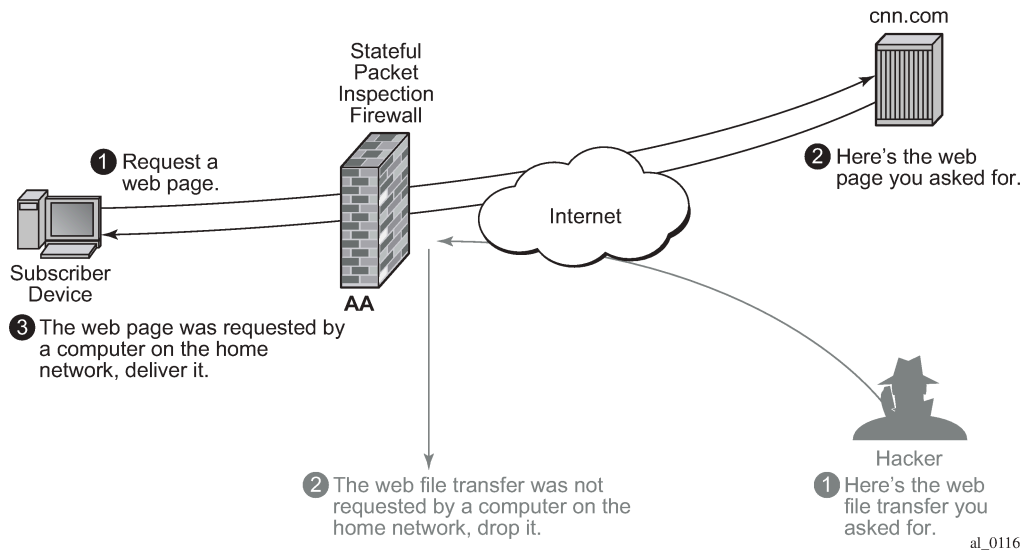
AA FW enabled solution provides:

- stateful/stateless packet filtering and inspection with Application-Level Gateway (ALG) support
- security gateway ([SeGW firewall protection](#) for S1, MME (SCTP), S1-U (GTP-U) and OAM traffic protection)
- [GTP backhaul roaming firewall protection](#)

Stateful flow processing and inspection uses IP Layers 3/4 header information to build a state of the flow within AA ISA. Layer 7 inspection is used to provide ALG support. Stateful flow/session processing takes note of the originator of the session and can therefore allow traffic to be initiated from the subscriber while denying, if configured, traffic originating from the network. Packets received from the network are inspected against the session filter and only those that are part of a subscriber-initiated-session are allowed.

Figure 35: Stateful firewall shows stateful firewall processing.

Figure 35: Stateful firewall



Stateless packet filtering does not take note of session initiator and therefore, it discards or allows packets independent of the any previous packets. Stateless packet filtering can be performed in the system using IOM ACLs.

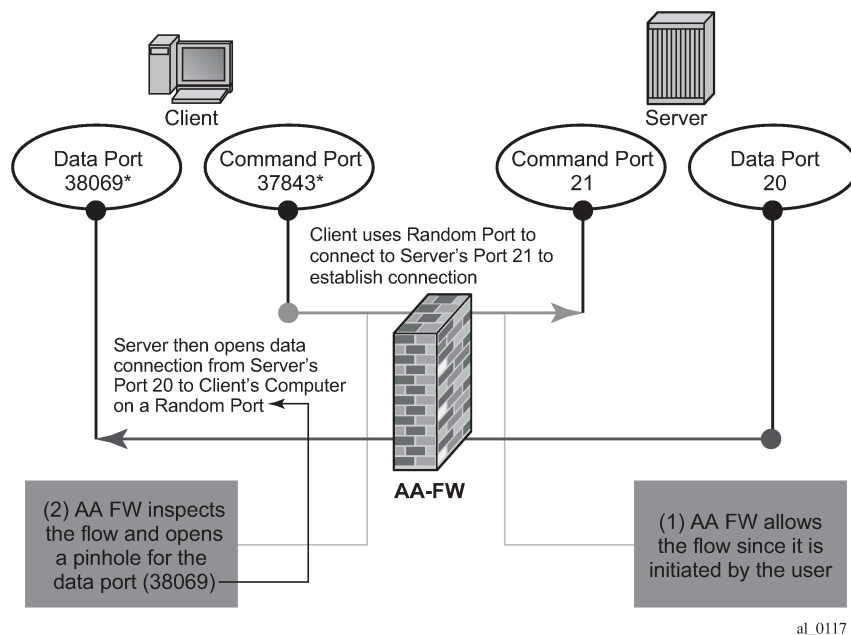
AA FW inspection of packets at Layer 7 offers Application Layer Gateway functionality for the following applications:

- rtsp
- sip
- h323 (IPv4 only)
- googletalkvoice
- ftp
- tftp
- pptp
- citrix
- sybase
- msexchange
- skinny
- ares
- bittorrent

- dns
- irc
- mailru
- qvod
- R commands
- sc2
- socks
- vudu
- winmx
- xunlei

Figure 36: Application layer gateway support shows application layer gateway support.

Figure 36: Application layer gateway support



These applications make use of control channels and flows that spun other flows. AA FW inspects the payload of these control flows so that it can open a pinhole for the associated required flows.

3.2.2.15.1 DoS protection

Denial of Service (DoS) attacks work by consuming network and system resources, making them unavailable for legitimate network applications. Network flooding attacks, malformed packets, and port scans are examples of such DoS attacks.

The aim of AA FW DoS protection is to protect subscribers and prevent any abuse of network resources.

Using AA FW stateful session filters, operators can protect their subscribers from any port scan scheme by configuring the session filters to disallow any traffic that is initiated from the network.

Furthermore, AA ISA provides configurable flow policers. These policers, when configured, prevent all sorts of flooding attacks (for example, ICMP PING flooding, UDP flooding, SYN Flood Attack). These policers provide protection at multiple levels; per system per application/application groups and per subscriber per applications/applications groups. AA ISA flow policers has two flavors; flow setup rate policers and flow count policers. Flow setup rate policers limit the number of new flows, while flow count policers limit the total number of active flows.

To protect hosts and network resources, AA_FW validates/checks the following parameters, if any fails, it declares the packet to be invalid (/Errored):

- IP layer validation:
 - IP version is not 4 nor 6
 - checksum error (IPv4)
 - header length check
 - packet length check
 - TTL/Hop limit (not equal to zero) check
 - fragment offset check (teardrop and ping of death protection)
- class D/E (>=224.0.0.0)
- broadcast 255.255.255.255 (multicast source address)
- 127.x.x.x (invalid source address)
- invalid subnet (subnet, 0) [unless /31 point-to-point interface]
- invalid subnet multicast (subnet, -1) unless /31 point-to-point interface
- IPv4 destination address checks:
 - broadcast 255.255.255.255, 0.x.x.x, 127.x.x.x
- IPv6_source address checks:
 - multicast source address (FFxx:xxxx:.....:xxxx)
- IPv6_destination address checks:
 - invalid destination address (=:)
- TCP/UDP validation:
 - header checksum
 - source or destination ports (not equal to zero) check
(only dest port is checked for UDP)
 - invalid TCP flags
 - TCP FIN Only (only the FIN flag set)
 - TCP No Flags (no flags are set)
 - TCP FIN RST (both FIN and RST are set)
 - TCP SYN URG (both SYN and URG are set)
 - TCP SYN RST (both SYN and RST are set)
 - TCP SYN FIN (both SYN and FIN are set)
 - validates that the first packet of a TCP flow does not contain RST or FIN flags

The above complements ESM enhanced security features, such as IP (or mac) anti-spoofing protection (for example, protecting against "LAND attack") and network protocols DoS protections. The combination provides a world class carrier grade FW function.

3.2.2.15.2 TCP validation

Operators can configure AA AQP actions to monitor TCP packet exchanges and ensure that they follow TCP handshake procedures. AA drops packets that do not conform to these procedures. AA FW checks for corrupted TCP packets and invalid TCP flag settings for the different TCP session states.

For example, if the SACK Permitted or MSS option is detected, but the calculated option length is incorrect, AA flags the packet as malformed and drops it. TCP sessions that start without a SYN and packets received after a FIN are discarded as well.

Furthermore, if strict **tcp-validation** is configured, AA checks and drops TCP packets with invalid sequence or acknowledgment numbers.

Drops because of TCP validation policies are recorded under permit-deny statistics. Therefore, TCAs can be configured against these statistics. Optionally, the operator can also capture TCP validation drop activity by enabling event logging.

3.2.2.15.3 Policy partitioned AA FW

AA FW can provide up to 128 virtual/partitioned FWs, each with its own FW policies. This is achieved through the use of AA-Partitions. Different VPNs can have different FW policies/rules.

3.2.2.15.4 Configuring AA FW

AA ISA AQPs are enhanced with few new AQP actions that provide session filtering functionality. As is the case of AQPs, these have partition level scope. This allows different FW policies to be implemented by utilizing AA partitions concepts within the same AA ISA group. Hence, multiple virtual AA FW instances can be realized. There is no need for multiple physical instances of FWs to implement different FW policies.

The AA FW stateful session filter consists of multiple entries (similar to ACLs) with a **match** and **action** per entry. Actions are **deny** or **permit**. A **deny** action results in packets discarded without creating a session / flow context. **match** conditions include IP 5 tuples info. An overall default action is also configurable, in case of a no match to any session filter entry.

AQPs with session filter actions, need to have, as a matching condition, traffic direction, ASOs or subscriber name. It cannot have any references to applications or application groups.

AA FW offers, in addition to session-filter actions, a variety of AQP actions to that are aimed to allow or deny: errored/malformed packets, fragmented packets or first packet out-of-order fragments and overload traffic.

Statistics are incremented when packets are dropped by a session filter. These are accounted against:

- protocol = denied by default policy
- application= unknown
- application group = unknown

A session-filter hit-count counter is maintained by AA ISA and can be viewed via CLI. There is no current support for export of session-filter entry hit counters via XML to SAM.

3.2.2.15.5 AA FW logging

AA ISA can be configured, per AQP or per session filter, to log events related to how the packets are processed (either allowed or denied). AA supports event logging locally on the node or remotely via syslog. AA ISA FW logs contain the following information:

- group partition
- timestamp
- 5-tuple
- direction
- subscriber info (if available)
- log source/type (session-filter or AQP)
- action (allow/drop)
- session-filter specific
 - session-filter name
 - session-filter entry
- AQP specific
 - drop reason
 - fragment offset (if applicable)
 - fragment ID (if applicable)
 - TCP validation policy (if applicable)
- If an out of order fragment triggers the log, then whatever 5-tuple information is available is included.

For AQPs, only **drop** events are captured in the log. The logs do not capture drops because of flow policers.

The operator can configure up to one event log per partition. For offline logging via syslog, the operator needs to configure the IP address of the syslog server and the VLAN ID to be used to connect to the server.

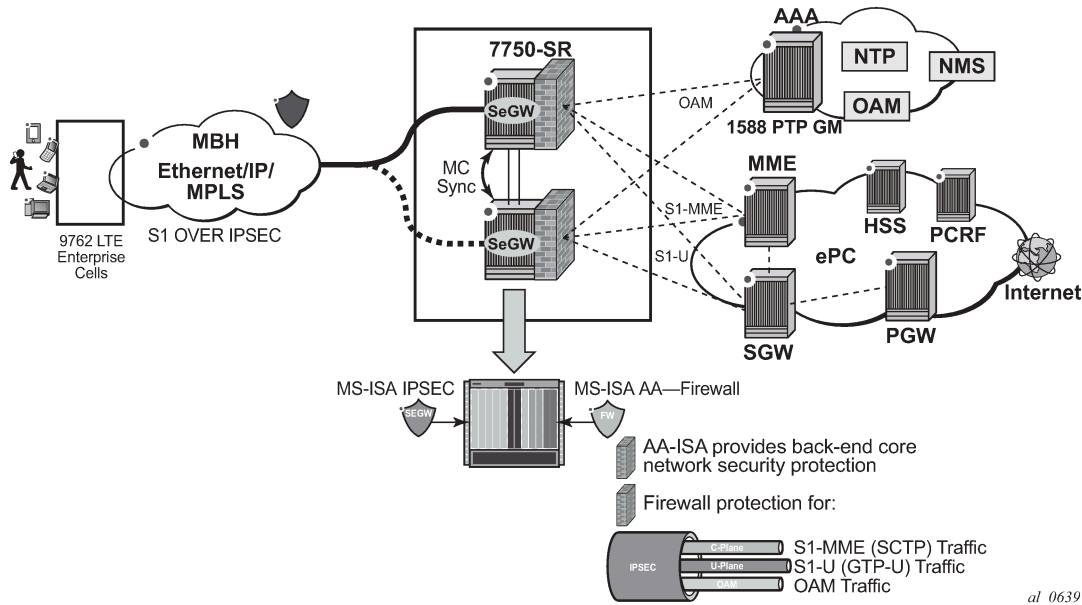
3.2.2.15.6 SeGW firewall protection

The 7750 SR SeGW with AA firewall (AA FW) deployed in 3G/4G/Femto access networks provides the operator with back-end core network security protection. AA FW provides protection for:

- S1-MME (SCTP) traffic
- S1- U (GTP-U) traffic
- OAM traffic

[Figure 37: SeGW firewall deployment](#) shows an example of an SeGW firewall deployment.

Figure 37: SeGW firewall deployment



SAPs on the private side of tunnel ISA are diverted to AA for firewall protection. If per eNB/ Femto Access Point (FAP) control is needed, then AA auto-configures or instantiates subscribers based on the “seen-ip” transit-AA subscriber model (no RADIUS interaction is required). This auto-creates a subscriber per eNB/ FAP. Alternatively, AA applies firewall rules to the diverted SAP (for all eNBs/FAPs) at the aggregate level (for all eNBs/FAPs).

One AA ISA is supported per tunnel ISA group. Therefore, all private side SAPs that are diverted to AA for firewalling service go to the same AA ISA module with no need to load balance the traffic into different AA ISAs. If the capacity of one AA ISA is not sufficient, then the IPsec tunnel group is split into two (or more) groups. Each group is served by an AA ISA.

3.2.2.15.6.1 OAM traffic protection

The aim of AA Firewall protection is to protect and prevent any abuse of OAM network resources (such as NMS).

Network flooding attacks, malformed packets and port scans are examples of such attacks that can be carried out using a compromised eNB/Femto Access Points (FAP).

- **ports scan attacks**

Using AA FW stateful session filters, operators can allow traffic only on certain IP addresses and port numbers.

For example, operator can configure AA to only allow traffic that is initiated by NMS toward the FAPs. Therefore, a compromised FAP cannot initiate an attack on NMS infrastructure.

Operator can limit the type of traffic allowed based on Layer 3 — Layer 7 classification. Operator can allow only HTTP with a specific URL/domain, DNS, PTP, FTP (independent of the port number used) and block all other traffic.

- **flood attacks**

The operator can limit the type of traffic allowed based on Layer 3 — Layer 7 classification. The operator can allow only HTTP with a specific URL/domain, DNS, PTP, FTP. The AA ISA provides configurable flow policers that can act on FW permitted sessions. These policers, when configured prevent all sort of flooding attacks, such as ICMP PING flooding, UDP flooding, SYN Flood Attack, and so on, of the port number used) and block all other traffic.

- These policers provide protection at multiple levels; per system per application/application groups and per FAP (or per NMS) per applications or applications groups.
- There are three types of AA ISA policers:
 - flow setup rate policers to limit the number of new flows
 - flow count policers to limit the total number of active flows
 - bandwidth policers to limit the total OAM bandwidth allowed by a FAP toward NMS

- **malformed packets attacks**

To protect Hosts and network resources, AA FW performs validation on IP packets, at the IP layer and TCP/UDP layer, to ensure that the packets are valid. Invalid packets are discarded (a configurable option). This provides protection against well-known attacks such as LAND attack. See [Stateful firewall service](#) for a complete description. AA allows the operator to optionally drop fragmented or out-of-order fragmented IP packets.

In addition, for OAM traffic, all AA functionalities including Layer 7 analytics and control as well as Application Layer Gateway (ALG) are supported.

For more details on OAM Traffic protection, see the [AA firewall](#) and the [DoS protection](#) sections.

3.2.2.15.6.2 S1-MME (SCTP) firewall

Network flooding attacks, malformed packets and port scans are examples of DoS attacks that can be carried out using a compromised eNB/FAP. AA FW provides inspection of SCTP (the protocol used to communicate to MME). Such inspection includes checking for SCTP protocol ID, source or destination ports, PPID, SCTP chunk checking and malformed SCTP packet (such as checksum validation).

SCTP chunk checking includes checking for:

- invalid length values (frames with invalid length value are dropped regardless of the chunk type)
- data chunks with length value that is too small to accommodate PPID (such frames are dropped as invalid/badly formed)
- data chunks with length value that is too large for chunk (such frames are dropped as invalid/badly-formed)

For S1-MME traffic, the operator can configure various AA actions:

- Drop packets with invalid checksum, src/dest IP or port numbers (malformed Packet protection) by appropriately configuring session filters with the following command.

```
configure application-assurance group policy app-qos-policy entry action error-drop
```

- Use the **sctp-filter** command for PPID filtering.
- Rate limit the amount of S1-MME traffic (flooding protection) in terms of bandwidth (b/s), using AA bandwidth policers.
- Limit the number of concurrent SCTP flows (flooding protection) using AA flow count policers.

- Limit the SCTP flow setup rate (flows per second) to protect against DoS flooding using AA flow rate policers.
- Drop fragmented packets or drop out of order fragmented packets using the following command.

```
configure application-assurance group policy app-qos-policy entry action fragment-drop {all
| out-of-order}
```

The actions above can be applied per eNB/FAP IP address or per MME (to control aggregate traffic per MME).

3.2.2.15.6.3 SCTP PPID filtering

AA allows the operator to configure PPID filters that contain a list of PPIDs to allow or deny using the following commands.

```
configure application-assurance group sctp-filter ppid default-action
configure application-assurance group sctp-filter ppid default-action entry action
```

The filter can then be used within an AQP action.

AA identifies data chunks within SCTP payloads (for example, as first, nth or last chunk) and filters according to the configure PPID filter. If any chunk PPID matches a PPID on the configured blocked PPID list, the whole SCTP packet is dropped.

SCTP packets without data chunks are not impacted or accounted for by an SCTP Filter.

For IP fragmentation, and in the case where the operator did not configure AA ISA to drop "all fragmented traffic", only the first IP fragment is inspected and subject to the PPID filtering. Any action applied to the first fragment is also applied to the remaining fragments. Out-of-order fragments appearing before the first fragment receive the default action (for example, drop action of "out-of-order-Frag").

3.2.2.15.6.4 S1-U GTP traffic protection

The 7750 SR SeGW with AA FW provides protection of SGW/SGSN infrastructure against an attack from a compromised eNB/FAP. AA FW offers the following:

- protection against malformed GTP packets attack. For GTP-v1 traffic carried over UDP port number port 2152, AA performs various packet sanity checks, such as:
 - UDP destination port is 2152
 - version (GTP-U should always be version 1)
 - protocol type bit should be 1
 - invalid/missing mandatory header fields
 - invalid optional/spare header fields
 - invalid/missing header extensions
 - invalid length

For S1-U interface, only GTP-v1 is supported. No support for GTP-v2 (as there is no signaling on S1-U interface).

Details of the various GTP sanity checks that are performed for different GTP-U message types are shown in [Table 26: GTP-U message types](#).

Table 26: GTP-U message types

Payload size	Encapsulated data checks	IE checks	Header extension checks	Optional HEADER check	GTP mandatory header checks					
					If E, S or PN = 1	Length	TEID	Spare	PT	
>0	Payload Size is assumed to be the size of the remainder of the packet, unless the packet is fragmented No checking of the encapsulated data	No checks	Valid types = Service Class Indicator and PDCP PDU Number Extension size= 4*# of extensions	Optional Size = 8 IF E= 0, ExtSize = 0	Optional Size + Extension Size + Payload Size	<>0	0	1	1	G-PDU (Encapsulated Data Delivery) – Message Type 255
	No payload after the IEs	Only private extensions are allowed.	No external header allowed.	No option headers allowed.	IE Size	0	0	1		Echo Request – Message Type 1
	No payload after the IEs	Recovery ID is present Private extensions allowed.	No external header allowed.	No option headers allowed.	IE Size	0	0	1	1	Echo Response – Message Type 2
	No payload after the IEs	Extension Header Type List IE is	No external	No option	IE Size	0	0	0	1	Supported Extension Headers

Payload size	Encapsulated data checks	IE checks	Header extension checks	Optional HEADER check	GTP mandatory header checks						
					If E, S or PN = 1	Length	TEID	Spare	PT		Version
		present Private extensions allowed No checking on the extension header value	header allowed.	headers allowed.							Notification - Message Type 31
	No payload after the IEs	TEID IE and GTP-U Peer Address IE are present IE type and length are verified Private extensions allowed	Only the UDP Port Extension Header is valid	Optional Size = 8	Optional Size + Extension Size +IE Size	<>0	0	1	1		Error Indication – Message Type 26
	No payload after the IEs	Only Private extensions are allowed	no valid external header allowed.	Optional Size = 8 IF E = 0, ExtSize = 0	IE Size	<>0	0	1	1		End Marker – Message Type 254

To enable GTP packet sanity checks, the operator must configure:

```
configure application-assurance group gtp
```

When the **gtp** command is issued for a partition, AA treats traffic with UDP destination port number 2152 as GTP. It applies the different GTP level firewall functions as configured by the operator. However, it does not look beyond the GTP header for further inner L3-L7 packet classifications and actions. For example, Ipfix record for GTP traffic contains the 5 tuples of the GTP-u tunnel (eNB, SGW IPs and port numbers, and so on, no TEID).

- protection against unsupported GTP messages
- AA allows the operator to configure a GTP filter to indicate which GTP message types are to be allowed or denied as well as the maximum allowed GTP message length:

```
configure application-assurance group gtp gtp-filter
```

```
config>app-assure>group <aa-group-id>[:<partition>]>gtp
  gtp-filter <gtp-filter-name> [create]
  max-payload-length <bytes> // [0..65535]
  message-type
  default-action {permit|deny}
  entry <entry-id> value <gtp-message-value> action {permit|deny}
```

- There are approximately 67 valid message names to enter in the above GTP filter. Both names and numbers are accepted as input (for user convenience), but the CLI info always shows the name: echo-request, echo-response, error-indication, g-pdu, end-marker and supported-extension-headers-notification.
- After a GTP filter is configured, it can then be included as an AQP action:

```
config>app-assure>group <aa-group-id>[:<partition>]> policy
  app-qos-policy
  entry <entry-id> [create]
  action
  gtp-filter <gtp-filter-name>
```

- extensive GTP header sanity checks (included in [Table 26: GTP-U message types](#)) that are based on different GTP message types are only performed when these GTP messages are permitted by the GTP filter. If no GTP filter is configured, then no extensive GTP-U header checks are performed. In other words, if the operator wants to allow all GTP-U packets and perform all GTP header sanity checks, then the operator needs to configure a GTP filter with default action of **permit** and no values, such as:

```
config>app-assure>group 1:100> gtp
  gtp-filter "allow-all" create
  message-type
  default-action permit
```

- protection against flooding attacks; AA can be configured to drop all fragments and/or out of order fragments, using AQP action: **fragment-drop {all | out-of-order}**.
- In the case that the IP **fragment-drop** command is not set, then the following conditions apply to the way AA inspects GTP traffic:
 - Permit/deny decisions are entirely based on the first fragment. The first fragment contains the entire GTP header in almost all of the cases.
 - Max packet length check is not done across fragments. Only the first fragment length is checked. In other words, AA ISA may allow a packet that is larger than the max packet allowed if it is fragmented, with the first fragment smaller than the configured maximum packet size allowed.
 - First fragmented packet is discarded (and logged), as well as subsequent fragments:
 - If the first packet is too small to contain the mandatory header (12 bytes, ending with the TEID).
 - If the mandatory header indicates there should be an optional header, and the fragment is too small to contain the optional header (mandatory + optional = 16 bytes).

- GTP-in-GTP protection; GTP-in-GTP is a spoofing method that uses GTP-in-GTP encapsulation. After receiving the GTP packet in the upstream, the Serving GPRS Support Node (SGSN) encodes the packet again and forwards the packet to the Gateway GPRS Support Node (GGSN), through the relative PDP context. The embedded GTP packet may get decoded by the GGSN and allow an attacker to spoof GTP packets.

AA provides a mechanism to detect and drop GTP-in-GTP GTP-U packets:

```
*A:Dut-C>config>app-assure>group>gtp#
+---gtp-filter <gtp-filter-name> [create]
| +---gtp-in-gtp {permit|deny}
```

By default, GTP-in-GTP checking is disabled.

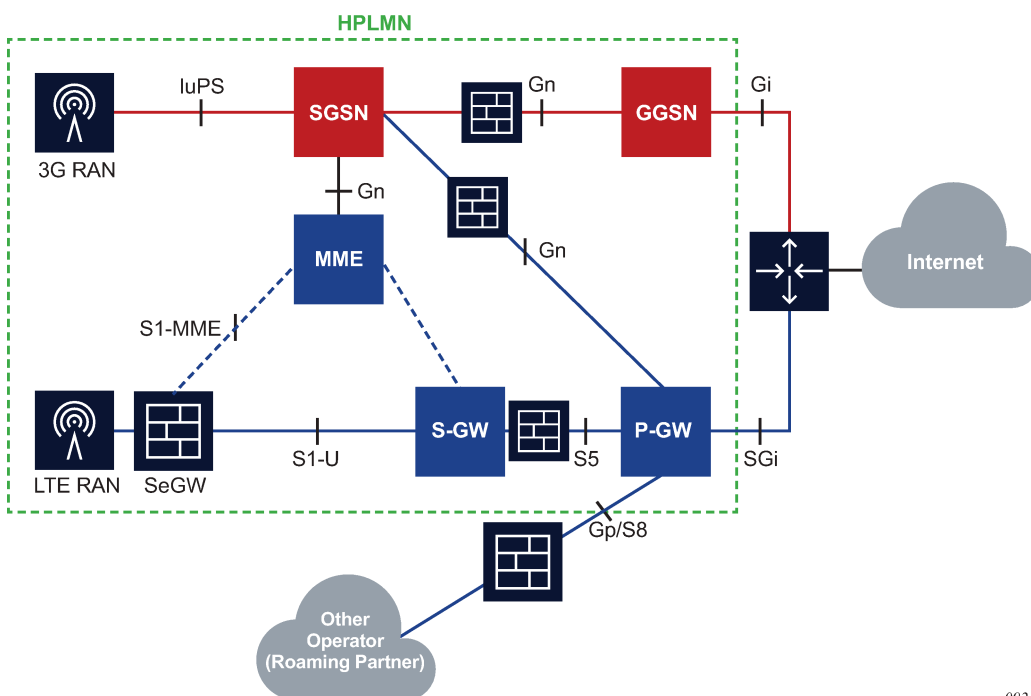
3.2.2.15.7 GTP backhaul roaming firewall protection

Wireless network operators rely on the GPRS Tunneling Protocol (GTP) for the delivery of mobile data services across the access network. GTP is not designed to be secure and is exposing the mobile access network to risk from both its own subscribers and its partners' networks attacks.

While roaming is essential to mobile operators, it comes with its own additional unique security risks when providing connectivity to roaming partners' networks and the end customers.

Figure 38: S5/S8/Gn/Gp AA firewall deployment shows the S5/S8/Gn/Gp AA Firewall deployment.

Figure 38: S5/S8/Gn/Gp AA firewall deployment



sw0924

AA is deployed as a GTP firewall on S8/Gp (or S5/Gn) interfaces either as part of the 7750 SR router in the form of the AA-ISA hardware module or as a separate Virtual SR (VSR) appliance. The AA firewall provides network security features, such as:

- GTP protocol validation (checks for anomaly attacks that involve malformed, corrupt, or spoofed traffic):
 - header length checks
 - IE length validation
 - invalid reserved field validation
 - reserved information elements validation
 - missing mandatory information elements validation
 - out of state message/information elements validation (GTPv1-c only)
 - sequence number validation
 - TEID validation (blocks GTP tunnel creations that have not been signaled correctly)
- GGSN/PGW and SGSN/SGW redirection protection
- GTP-in-GTP check
- handover control to prevent session hijacking
- source address—UE – anti-spoofing protection
- protection against unauthorized Public Land Mobile Network (PLMN) access and unauthorized APN access by filtering based on APN, International Mobile Station Identity (IMSI) prefix, GTP src IP address prefix list
- protection against unsupported GTP message types by filtering messages based on message type and message length
- protection against flooding attacks:
 - GTP traffic bandwidth policing (limits the GTP bandwidth from a roaming partner's SGSN or SGW)
 - GTP tunnel limiting (limits the number of concurrent GTP tunnels and the setup rate of these tunnels from a roaming partner's SGSN or SGW)
- protection against IP fragmentation-based attacks by using drop rules for IP fragmentation of GTP messages

AA FW supports GTPv1 and GTPv2.

3.2.2.15.7.1 UE IP address anti-spoofing

Source address spoofing (also known as Overbilling Attacks) is initiated by a malicious UE that hijacks (spoofs) an IP address of another UE and invokes a download from a malicious server on the Internet. After the download begins, the malicious UE exits the session. The UE under attack, which is receiving the download traffic, gets charged for traffic it did not solicit.

AA FW associates the GTP-C message's End user address IE with the GTP-U packets to make sure the packets carried in the upstream have the correct source IP address (inner IP within the GTP-U tunnel). When the UE address is negotiated within the PDP Context creation handshake, all the packets originating from the UE that contain a different source address are detected by AA FW and dropped.

To enable UE IP anti-spoofing protection, the operator must enable **validate-source-ip-addr**:

```
*A:Dut-C>config>app-assure>group>
+---gtp-filter <gtp-filter-name> [create]
|   +---validate-source-ip-addr
```

By default, **validate-source-ip-addr** is disabled.

3.2.2.15.7.2 GTP TEID validation

Compromised GSNs can send storms of GTP traffic with invalid GTP Tunnel Identifications (TEIDs) to cause a DoS attack. By inspecting GTP-C messages, AA FW supports stateful correlation of upstream and downstream GTP flows (DstIP + TEID) of the same PDN session.

AA drops packets with TEIDs that have not been negotiated correctly.

By default, TEID validation is disabled. The operator can enable AA to drop GTP traffic with invalid TEID using the following command sequence.

```
*A:Dut-C>config>app-assure>group>
+---gtp-filter <gtp-filter-name> [create]
|   +---validate-gtp-tunnels
```

3.2.2.15.7.3 GTP-C out-of-state message-type protection

GTP is a stateful protocol. Consequently, some message types can only be sent in specific states. For example, PDP context update messages are not allowed for PDP contexts that do not exist or have been closed.

AA performs stateful GTP protocol validation and allows only packets that are allowed for any state or a specific deployment.

[Table 27: Invalid message types in GTP FW roaming deployments](#) lists the message types that are invalid in GTP FW roaming deployments. When AA FW GTP-C inspection is enabled, the packets with the message types listed in [Table 27: Invalid message types in GTP FW roaming deployments](#) are dropped and the associated event logs include a "wrong interface" indication.



Note: The packets are dropped regardless of the configuration in the **message-type** or **message-type-gtpv2** filter.

Table 27: Invalid message types in GTP FW roaming deployments

GTP version	GTP-U port	GTP-C port
GTPv1	no invalid message types	GTPU PDU GTPV1_END_MARKER GTPV1_MSG_ERR_IND GTPV1-ALL-MBMS message-types GTPV1-ALL-Location management message-types
GTPv2	not applicable	GTP_PKT_ERROR_INDICATION GTP_PKT_DNLK_DATA_FAIL_INDICATION GTP_PKT_STOP_PAGING_INDICATION GTP_PKT_CRE_INDR_TNL_REQ GTP_PKT_CRE_INDR_TNL_RSP

GTP version	GTP-U port	GTP-C port
		GTP_PKT_DEL_INDR_TNL_REQ
		GTP_PKT_DEL_INDR_TNL_RSP
		GTP_PKT_RELEASE_BEARERS_REQ
		GTP_PKT_RELEASE_BEARERS_RSP
		GTP_PKT_DNLK_DATA
		GTP_PKT_DNLK_DATA_ACK
		GTP_PKT_MOD_ACCESS_BEARERS_REQ
		GTP_PKT_MOD_ACCESS_BEARERS_RSP
		GTP_PKT_REMOTE_UE_RPRT_NOTF
		GTP_PKT_REMOTE_UE_RPRT_ACK

AA does not perform GTP-C inspection by default. To enable GTP-C inspection, use the following command:

```
*A:Dut-C>config>app-assure>group>
+---gtpc-inspection
```

3.2.2.15.7.4 GTP anomaly prevention (sequence number checks)

Protocol anomaly attacks involve malformed or corrupt packets that typically fall outside of the protocol specifications. Packets are denied by AA FW if they fail the sanity check. Examples of GTP sanity checks are: invalid GTP header length, invalid Information Element (IE) length, invalid reserved fields, invalid sequence number, missing mandatory IEs, out-of-state message type.

In addition to the GTP-C inspection and GTP-U protocol validation described in [UE IP address anti-spoofing](#), [GTP TEID validation](#), and [GTP-C out-of-state message-type protection](#), AA FW performs sequence number validation, whereby AA FW ensures that there are no out-of-sequence GTP packets. By default, sequence number validation is disabled. To enable sequence number validation, use the following CLI command:

```
*A:Dut-C>config>app-assure>group>
+---gtpc-inspection
+---gtp-filter <gtp-filter-name> [create]
|   | +---validate-sequence-number
```

GTP Packets with wrong sequence numbers are dropped when **validate-sequence-number** is enabled.

3.2.2.15.7.5 GTP message type filtering

In addition to performing stateful GTP validation, in which packets with invalid message types (that is, message types that are not applicable to the roaming interfaces) are denied, AA FW allows the operator to further restrict allowed message types by configuring entries for GTP message type filters to deny (or permit) the message types listed in [Table 28: Allowed message types that can be denied](#).

Table 28: Allowed message types that can be denied

GTP version	GTP-U port	GTP-C port
GTPv1	GTPV1_MSG_ECHO_REQ GTPV1_MSG_ECHO_RESP GTPV1_SUPP_EXT_HDR_NOTIF GTPV1_MSG_ERR_IND GTPV1_END_MARKER GTPU_PDU	GTPV1_MSG_ECHO_REQ GTPV1_MSG_ECHO_RESP GTPV1_SUPP_EXT_HDR_NOTIF GTPV1_MSG_VER_NOT_SUPP_IND GTPV1_MSG_PDP_CREATE_REQ GTPV1_MSG_PDP_CREATE_RESP GTPV1_MSG_PDP_UPD_REQ GTPV1_MSG_PDP_UPD_RESP GTPV1_MSG_PDP_DEL_REQ GTPV1_MSG_PDP_DEL_RESP GTPV1_MSG_NET_INIT_REQ GTPV1_MSG_NET_INIT_RESP GTPV1_MSG_MSINFO_REQ GTPV1_MSG_MSINFO_RESP
GTPv2	N/A	GTP_PKT_ECHO_REQ GTP_PKT_ECHO_RSP GTP_PKT_VERSION_NOT_SUPPORTED GTP_PKT_CRE_SES_REQ GTP_PKT_CRE_SES_RSP GTP_PKT_MOD_BEARER_REQ GTP_PKT_MOD_BEARER_RSP GTP_PKT_DEL_SES_REQ GTP_PKT_DEL_SES_RSP GTP_PKT_CHG_NOT_REQ GTP_PKT_CHG_NOT_RSP GTP_PKT_MOD_BEARER_CMD GTP_PKT_MOD_BEARER_FAIL_INDICATION GTP_PKT_DEL_BEARER_CMD GTP_PKT_DEL_BEARER_FAIL_INDICATION GTP_PKT_BEARER_RESOURCE_CMD GTP_PKT_BEARER_RESOURCE_FAIL_INDICATION GTP_PKT_SUSPEND_NOTIFICATION

GTP version	GTP-U port	GTP-C port
		GTP_PKT_SUSPEND_ACK GTP_PKT_RESUME_NOTIFICATION GTP_PKT_RESUME_ACK GTP_PKT_CRE_BEARER_REQ GTP_PKT_CRE_BEARER_RSP GTP_PKT_UPD_BEARER_REQ GTP_PKT_UPD_BEARER_RSP GTP_PKT_DEL_BEARER_REQ GTP_PKT_DEL_BEARER_RSP GTP_PKT_TRACE_SESSION_ACTIVATION GTP_PKT_TRACE_SESSION_DEACTIVATION GTP_PKT_UPDATE_PDN_CONNECTION_SET_REQ GTP_PKT_UPDATE_PDN_CONNECTION_SET_RSP GTP_PKT_DELETE_PDN_CONNECTION_SET_REQ GTP_PKT_DELETE_PDN_CONNECTION_SET_RSP

By default, GTP message filtering allows all of the GTP messages.

To configure GTPv2 message filtering, use the following command:

```
*A:Dut-C>config>app-assure>group>
+---gtp-filter <gtp-filter-name> [create]
  +---gtpc-inspection
  | +---message-type-v2
  | |
  | | +---default-action {permit|deny}
  | |
  | | +---entry <entry-id> value <gtpv2-message-value> action {permit|deny}
  | |
  | | no entry <entry-id>
```

To configure GTPv1 message filtering, use the following command:

```
*A:Dut-C>config>app-assure>group>
| +---message-type
| |
| | +---default-action {permit|deny}
| |
| | +---entry <entry-id: 1..255> value <gtpv1-message-value> action {permit|deny}
| |
| | no entry <entry-id>
```



Note: If the operator configures a message type invalid for the roaming interface and the user configures the message type to be denied, the message type is dropped and counted under that filter entry (and not tagged dropped because of "wrong-interface" in the event log). However, configuring the message-type filter to "permit" a message type that is invalid for the roaming

interface does not take effect, because the packet with the specified message type is dropped by the GTP-C protocol inspection process.

3.2.2.15.7.6 Unauthorized APN attacks and APN filtering

Access Point Name (APN) filtering checks GTP-C messages to determine if a roaming subscriber is allowed to access a specified external network (also known as APN).

create-session-request and "create pdp request" GTP message types contain an APN IE in the header of a GTP packet. The APN IE consists of an external network-ID (for example, Nokia.com). Optionally, the APN IE can include a unique ID that identifies the operators' PLMN.

APN filtering prevents malicious UEs from initiating a "Create PDP/session Request" flood attack toward the PGW/GGSN for invalid or disallowed APNs.

The AA GTP filter can be configured to perform APN filtering to restrict roaming subscribers' access to specific external networks.

An APN filter, an IMSI prefix, and an SGSN Address pool can be used together to filter GTP packets as shown in the following example.

An APN filter entry can also be optionally combined with an SGSN or SGW IP address (or IP address prefix list) to further restrict allowed APN access by configured SGSN or SGW nodes.

```
*A:Dut-C>config>app-assure>group>
+---gtp-filter <gtp-filter-name> [create]
|   +---imsi-apn-filter //NEW and all its children attributes
|   |   +---default-action {permit|deny} //default permit
|   |   |   +---entry <entry-id: 1031.2030> create
|   |   |   |   + apn <string 0|1..32 characters>
|   |   |   |   |---no apn
|   |   |   |   + src-gsn ip-prefix-list <ip-prefix-list>
|   |   |   |   + src-gsn <ip address prefix>
|   |   |   |   |---no src-gsn
|   |   |   |   + action {permit|deny} //default permit
|   |   |   +---no entry <entry-id>
```

An APN filter and an IMSI prefix filter (see [Unauthorized PLMN access—IMSI prefix filtering](#)) can be used together to filter GTP packets.

By default, AA FW allows all of the APNs.

3.2.2.15.7.7 Unauthorized PLMN access—IMSI prefix filtering

The PLMN of a subscriber's home network is identified by combining the Mobile Country Code (MCC) and the Mobile Network Code (MNC). MCC-MNC is also known as the IMSI prefix. IMSI prefix acts as a PLMN identifier.

GTP IMSI prefixes filters can be configured to deny GTP incoming traffic from invalid roaming partners. Conversely, the filters can be configured to allow only incoming traffic from those network operators that have signed roaming agreements. The GTP packets with IMSI prefixes that do not match the configured prefixes are dropped.

An IMSI filter entry can also be optionally combined with an SGSN or SGW IP address (or IP address prefix list) to further restrict allowed IMSI prefix traffic to specific SGSN or SGW nodes.

```
*A:Dut-C>config>app-assure>group>
```



```

+---gtp-filter <gtp-filter-name> [create]
|   +---imsi-apn-filter //NEW and all its children attributes
|   |   +---default-action {permit|deny} //default permit
|   |   +---entry <entry-id: 1031..2030> create
|   |   |   + mcc-mnc-prefix <string of 1 to 6 decimal digits>
|   |   |   |---no mcc-mnc-prefix
|   |   |   + src-gsn ip-prefix-list <ip-prefix-list>
|   |   |   + src-gsn <ip address prefix>
|   |   |   |---no src-gsn
|   |   |   + action {permit|deny} //default permit
|   |   +---no entry <entry-id>

```

An IMSI filter, an APN, and an SGSN Address can be used together to filter GTP packets.

By default, AA FW does not perform IMSI prefix filtering.

3.2.2.15.7.8 Unauthorized network access

An attacker using an unauthorized GSN can cause a denial of service attack using spoofed PDP Context Delete messages (DoS attack) or using a spoofed Update PDP Context Request to hijack existing sessions. Such attacks can also spoof a Create PDP Context Request to gain unlawful Internet access. Session hijacking can come from the SGW/SGSN or the PGW/GGSN. An unauthorized GSN can hijack GTP tunnels or cause a denial of service by intercepting another GSN and redirecting traffic to it.

Operators can use AA-FW to configure pools of trusted GSN IP addresses (using an AA IP-Prefix-list) to stop spoofed requests from untrusted GSNs.

AA IP-Prefix-lists can be configured to model GSN groups as follows:

```

*A:Dut-C>config>app-assure>group#
  ip-prefix-list ip-prefix-list-name [create]
    prefix ip-prefix/ip-prefix-length [name prefix-name]

```

The configured AA IP-Prefix-lists are then referenced in session-filters, such that only sessions that match the lists are "permitted".

```

*A:Dut-C>config>app-assure>group# session-filter
  default-action deny
  entry          + Configure an entry in the session filter
  match
  src-ip        // Configure IP addresses that correspond to authorized SGW/SSGN
  action
  permit

```

3.2.2.15.7.9 Typical S8/Gp AA FW network deployment model

AA GTP FW is typically deployed as an L3 and VPRN service. SAPs and spokes are diverted to AA for GTP FW. L2 and VPLS connectivity is supported by AA. AA transit subscribers (identified by SGW and SSGN IPs) are auto-created under the parent diverted SAPs and spokes. Operators need to enable inactivity monitoring to remove inactive subscribers. This can be done using the following command:

```

config>app-assure>group>transit-ip-policy>transit-auto-create
[no] inactivity-mon

```

3.2.3 Service monitoring and debugging

Operators can use AA-specific tools in addition to system tools that allow them to monitor, adjust, debug AA services. The following are examples of some of the available functions:

- Display and monitor AA ISA group status and statistics (AA ISA status and capacity planning/monitoring).
- Clear AA ISA group statistics (clears all system and per-AA subscriber statistics).
- Enable special study mode for real-time monitoring of AA subscriber traffic (ESM or transit subscriber, SAP or spoke SDP).
- Display per AQP entry statistics for number of hits (flow matching the entry) and conflicts (actions ignored because of per-flow-action-limit exceeded).
- Mirror (all or any subset of traffic seen by an AA ISA group).
- Display all the per-ISA statistics from the aa-performance record, for examining resource loading of each ISA.
- Display the top active AA subscribers per ISA by bytes, packets or flows, for traffic in each direction.

3.2.4 CPU utilization

The ISA show status command displays per ISA CPU utilization by main tasks, to provide insight into what aspects of load may be loading the ISA. These are split into 2 main areas:

- **management CPU**

This includes all tasks related to communication between the CPM and the ISA, with the following usage percentage reported:

- system (various infrastructure and overhead work)
- management (managing AA policy, AA subscriber and trap configurations, and handling tools requests)
- statistics (collecting and reporting statistics and Cflowd reporting)
- idle

- **datapath CPUs**

This includes all tasks related to datapath packet and flow processing on the ISA, with the following usage percentage reported:

- system (various infrastructure and overhead work)
- packet processing (receiving, associating with flows, applying application QoS policy, and transmitting)
- application ID (using protocol signatures and other techniques to identify application/app-group and determine the application QoS policy)

3.2.5 CLI batch: begin, commit and abort commands

The AA uses CLI batch capability in policy definition. To start editing a policy, a begin command must be executed. To finish editing either abort (discard all changes) or commit (accept all changes) needs to be executed. CLI batch state is preserved on an HA activity switch.

To enter the mode to create or edit policies, the **begin** keyword must be entered at the prompt. Other editing commands include:

- The **commit** command saves changes made to policies during a session. The newly committed policy takes effect immediately for all new flows. Existing flows transition onto the new policy shortly after the commit.
- The **abort** command discards changes that have been made to policies during a session.

To allow flexible order for policy editing, the **policy>commit** function cross references policy components to verify, among others:

- whether all ASO characteristics have a default value and are defined in the app-profile
- checks whether limits are adhered

3.3 Configuring AA with CLI

This section provides information to configure Application Assurance entities using the command line interface. It is assumed that the user is familiar with basic configuration of policies.

3.3.1 Provisioning AA ISA MDA

The following illustrates syntax to provision AA ISA2 and configure ingress IOM QoS parameters (the egress IOM QoS is configured in the **config>isa>application-assurance-grp>qos** context).

- classic CLI
configure card slot-number mda mda-slot mda-type isa2-aa
- MD-CLI
configure card slot-number mda mda-slot mda-type isa2-aa

The following output displays an AA ISA configuration example.

Classic CLI:

```
*A:cpm-a>config>app-assure# show mda 1/1
=====
MDA 1/1
=====
Slot  Mda  Provisioned      Equipped      Admin      Operational
      Mda  Mda-type        Mda-type        State      State
-----
1     1     isa2-aa         isa2-ms         up         up
=====
*A:cpm-a>config>app-assure#

*A:cpm-a>config>card# info
-----
card-type iom4-e-b
mda 1
mda-type isa2-aa
exit
-----
*A:cpm-a>config>card#
```

MD-CLI:

```
(gl)[/]
A:admin@Dut-E# show mda 9/1
=====
MDA 9/1
=====
Slot  Mda   Provisioned Type           Admin   Operational
      Mda   Equipped Type (if different) State   State
-----
 9     1     isa2-aa                    up      up
      me-isa2-ms
=====

(gl)[/]
A:admin@Dut-E# configure card 9
(gl)[/configure card 9]
A:admin@Dut-E# info
  card-type iom4-e-b
  mda 1 {
    mda-type isa2-aa
  }
  fp 1 {
  }
```

3.3.2 Configuring an AA ISA group

About this task

This section describes how to enable AA on the router.

Procedure

- Step 1.** Create an AA ISA group.
- Step 2.** Assign active and optional backup AA ISAs to an AA ISA group.
- Step 3.** Select the forwarding classes to divert.
- Step 4.** Enable the group.
- Step 5.** Perform the following optional steps:
 - a. Enable group policy partitioning.
 - b. Configure capacity cost threshold values.
 - c. Configure the number of transit prefix IP policies.
 - d. Configure IOM egress queues to the MS-ISA.
 - e. Enable overload cut through and configure the high and low watermarks values.
 - f. Configure performance statistics accounting.

Example: AA ISA group configuration

The following example illustrates AA ISA group configuration with:

- primary AA ISA and warm redundancy provided by the backup AA ISA
- “fail-to-wire” behavior configured on group failure

- BE forwarding class selected for divert
- default IOM QoS for logical ISA egress ports. The ISA ingress QoS is configured as part of ISA provisioning (**config>card>mda>network>ingress>qos**).

The following commands illustrate the AA ISA group configuration context:

```

config>isa>application-assurance-group isa-aa-group-id [create] [aa-sub-scale sub-scale]
backup mda-id
description description
divert-fc fc-name
no fail-to-open
http-enrich-max-pkt size-in-octets
isa-capacity-cost-high-threshold threshold
isa-capacity-cost-low-threshold threshold
isa-overload-cut-through
partitions
minimum-isa-generation min-isa-generation
overload-sub-quarantine
[no] shutdown
partitions
primary mda-id
qos
  egress
    from-subscriber
      pool [name]
        resv-cbs percent-or-default
        slope-policy slope-policy-name
        port-scheduler-policy port-scheduler-policy-name
        queue-policy network-queue-policy-name
        wa-shared-high-wmark percent
        wa-shared-low-wmark percent
    to-subscriber
      pool [name]
        resv-cbs percent-or-default
        slope-policy slope-policy-name
        port-scheduler-policy port-scheduler-policy-name
        queue-policy network-queue-policy-name
        wa-shared-high-wmark percent
        wa-shared-low-wmark percent
  shared-resources
    tcp-adv-func size
  [no] shutdown
  statistics
    performance
      accounting-policy acct-policy-id
      collect-stats
  transit-prefix-ipv4-entries entries
  transit-prefix-ipv4-remote-entries entries
  transit-prefix-ipv6-entries entries
  transit-prefix-ipv6-remote-entries entries

```

The following output displays an AA ISA group configuration example:

```

A:ALU-A>config>isa>aa-grp# info detail
-----
no description
primary 1/2
backup 2/2
no fail-to-open
isa-capacity-cost-high-threshold 4294967295
isa-capacity-cost-low-threshold 0
no partitions

```

```

divert-fc be
qos
  egress
    from-subscriber
      pool
        slope-policy "default"
        resv-cbs default
      exit
      queue-policy "default"
      no port-scheduler-policy
    exit
  to-subscriber
    pool
      slope-policy "default"
      resv-cbs default
    exit
    queue-policy "default"
    no port-scheduler-policy
  exit
exit
exit
no shutdown
-----
A:ALU-A>config>isa>aa-grp#

```

3.3.2.1 Configuring watermark parameters

Use the following CLI syntax to configure thresholds for logs and traps when under high consumption of the flow table. The flow table has a limited size and these thresholds can be established to alert the user that the table is approaching capacity. These flow table watermarks represent the number of flow contexts allocated on the ISA, which is slightly higher than the actual number of existing flows at the point when the watermark is reached.

The low threshold is used while the high threshold is used as an alarm.

```

config>application-assurance
  flow-table-high-wmark high-watermark
  flow-table-low-wmark low-watermark

```

3.3.3 Configuring a group policy

3.3.3.1 Beginning and committing a policy configuration

To enter the mode to create or edit Application Assurance policies, you must enter the **begin** keyword at the **config>app-assure>group>policy** prompt. The **commit** command saves changes made to policies during a session. Changes do not take effect in the system until they have performed the commit function. The **abort** command discards changes that have been made to policies during a session.

The following error message displays when creating or modifying a policy without entering **begin** first:

```

A:ALA-B>config>app-assure>group>policy#
MINOR: AA #1005 Invalid Set - Cannot proceed with changes when in non-edit mode

```

There are no default policy options. All parameters must be explicitly configured.

Use the following CLI syntax to begin a policy configuration:

```
config>app-assure# group group-id
policy
  begin
```

Use the following CLI syntax to commit a policy configuration:

```
config>app-assure# group group-id
policy
  commit
```

3.3.3.2 Aborting a policy configuration

Use the following CLI syntax to abort a policy configuration:

```
config>app-assure# group group-id
policy
  abort
```

3.3.3.3 Configuring an IP prefix list

An operator can use IP lists to define a list of IP addresses (along with any masks). This list can be later referenced in AQPs, application filters or session-filters.

Use the following CLI syntax to configure an application filter entry:

```
config>aa>group>policy>app-assurance>group <aa-group-id>[:<partition>]
  ip-prefix-list <prefix-list-name> [create]
  no ip-prefix-list <prefix-list-name>
  description <description>
  no description
  prefix <address/mask> [name <prefix-name>]
  no prefix <address/mask>
```

The following example displays an IP prefix list configuration:

```
*A:Dut-A>config>app-assure>group# ip-prefix-list AllowedLAN1Hosts create
*A:Dut-A>config>app-assure>group>pfx>$ description "allowed hosts"
*A:Dut-A>config>app-assure>group>pfx>$ prefix 10.10.8.2/32
*A:Dut-A>config>app-assure>group>pfx>$ prefix 10.10.8.180/32
*A:Dut-A>config>app-assure>group>pfx>$ prefix 10.10.8.231/32
*A:Dut-A>config>app-assure>group>pfx>$ exit
*A:Dut-A>config>app-assure>group#

*A:Dut-A>config>app-assure>group# ip-prefix-list "AllowedLan1Hosts"
*A:Dut-A>config>app-assure>group>pfx># info
-----
      description "allowed hosts"
      prefix 10.10.8.2/32
      prefix 10.10.8.180/32
      prefix 10.10.8.231/32
```

```
-----
*A:Dut-A>config>app-assure>group>px>#
```

3.3.3.4 Configuring AA session filters

Session filters can be configured to allow stateful firewall use-cases.

Use the following CLI syntax to configure an AA session filter:

```
*A:Dut-A>config>app-assure>group# session-filter <session-filter-name> [create]
  default-action {permit | deny} [event-log <event-log-name>]
  description <description-string>
  entry <entry-id> [create]
    action {permit | deny} [event-log <event-log-name>]
    match
      dst-ip <ip-address>
      dst-ip ip-prefix-list <ip-prefix-list-name>
      no dst-ip
      dst-port {eq | gt | lt} <port-num>
      dst-port range <start-port-num> <end-port-num>
      dst-port port-list <port-list-name>
      no dst-port
      ip-protocol-num <ip-protocol-number>
      no ip-protocol-num
      src-ip <ip-address>
      no src-ip
      src-ip ip-prefix-list <ip-prefix-list>
      src-port {eq | gt | lt} <port-num>
      src-port range <start-port-num> <end-port-num>
      src-port port-list <port-list-name>
      no src-port
```

```
*A:Dut-A>config>app-assure>group# session-filter " denyUnsolicitedwMgntCntrl" create
  description "S-FW opted-in sub – allow ISP access"
  default-action deny event-log "FW_log"
  entry 10 create
    description "allow ICMP access from ISP LAN#1"
    match
      ip-protocol-num icmp
      src-ip 10.10.8.0/24
    exit
    action permit
  exit
  entry 30 create
    description "allow all TCP (e.g. FTP/telnet)access from ISP LAN#2"
    match
      ip-protocol-num tcp
      src-ip 192.168.0.0/24
    exit
    action permit
  entry 40 create
    description "allow TCP on port 80 /HTTP access from a IP List on ISP LAN#1"
    match
      ip-protocol-num tcp
      src-ip ip-prefix-list AllowedLAN1Hosts
      dst-port eq 80
    exit
    action permit
  exit
```



```

*A:Dut-A>config>app-assure>group>sess-fltr$ info
-----
description "S-FW opted-in sub . allow ISP access"
default-action deny event-log "FW_Log"
entry 10 create
  description "allow ICMP access from ISP LAN#1"
  match
    ip-protocol-num icmp
    src-ip 10.10.8.0/24
  exit
  action permit
exit
entry 20 create
  description "allow ICMP access from ISP LAN#2"
  action deny
exit
entry 30 create
  description "allow all TCP (e.g. FTP/telnet)access from ISP LAN#2"
  match
    ip-protocol-num tcp
    src-ip 192.168.0.0/24
  exit
  action permit
exit
entry 40 create
  description "allow TCP on port 80 /HTTP access from a IP List on
ISP LAN#1"
  match
    ip-protocol-num tcp
    src-ip ip-prefix-list "AllowedLan1Hosts"
    dst-port eq 80
  exit
  action permit
exit
-----
*A:Dut-A>config>app-assure>group>sess-fltr$

*A:Dut-A>config>app-assure>group>policy>aqp>
entry 110 create
  description "FW for managed opted-in subs"
  match
    traffic-direction network-to-subscriber
  exit
  action
    session-filter " denyUnsolicitedwMgntCntrl "
    fragment-drop all event-log "FW_log"
    error-drop event-log "FW_log"
    overload-drop

  exit
exit

*A:Dut-A>config>app-assure>group>policy>aqp>entry# info
-----
description "FW for managed opted-in subs."
match
  traffic-direction network-to-subscriber
exit
action

```

```

        session-filter "denyUnsolicitedwMgmtCntrl"
        fragment-drop all event-log "FW_log"
        error-drop event-log "FW_log"
        overload-drop

        exit
        no shutdown
-----
*A:Dut-A>config>app-assure>group>policy>aqp>entry#

```

3.3.3.5 Configuring an application group

An operator can configure an application group to group multiple application into a single application assurance entity by referencing those applications in the group created.

Use the following CLI syntax to configure an application group:

```

config>app-assure>group>policy# app-group application-group-name [create]
description description

```

The following example displays an application group configuration:

```

*A:ALA-48>config>app-assure>group>policy# app-group "Peer to Peer" create
*A:ALA-48>config>app-assure>group>policy>app-grp# info
-----
description "Peer to Peer file sharing applications"
-----
*A:ALA-48>config>app-assure>group>policy>app-grp#

```

3.3.3.6 Configuring an application

An operator can configure an application to group multiple protocols, clients or network applications into a single Application Assurance application by referencing it later in the created application filters as display in other sections of this guide.

Use the following CLI syntax to configure an application:

```

config>app-assure>group>policy# application application-name [create]
app-group app-group-name
description description

```

The following example displays an application configuration:

```

*A:ALA-48>config>app-assure>group>policy# application "SQL" create
*A:ALA-48>config>app-assure>group>policy>app# info
-----
description "SQL protocols"
app-group "Business Critical Applications"
-----
*A:ALA-48>config>app-assure>group>policy>app#

```

3.3.3.7 Configuring an application filter

An operator can use an application filter to define applications based on ALU protocol signatures and a set of configurable parameters like IP flow setup direction, IP protocol number, server IP address and server TCP/UDP port. An application filter references an application configured as previously shown.

Use the following CLI syntax to configure an application filter entry:

```
config>app-assure>group>policy# app-filter
  entry entry-id [create]
    application application-name
    description description-string
    expression expr-index expr-type {eq | neq} expr-string
    flow-setup-direction {subscriber-to-network | network-to-subscriber | both}
    http-match-all-requests
    ip-protocol-num {eq | neq} protocol-id
    network-address {eq | neq} ip-address
    network-address {eq | neq} ip-prefix-list ip-prefix-list-name
    protocol {eq | neq} protocol-signature-name
    server-address {eq | neq} ip-address
    server-address {eq | neq} dns-ip-cache dns-ip-cache-name
    server-address {eq | neq} ip-prefix-list ip-prefix-list-name
    server-port {eq | neq | gt | lt} server-port-number
    server-port {eq | neq} range start-port-num end-port-num
    server-port {eq} {port-num | range start-port-num end-port-num} first-packet-trusted |
    first-packet-validate}
    no shutdown
```

The following example displays an application filter configuration:

```
*A:ALA-48>config>app-assure>group>policy>app-filter# entry 30 create
*A:ALA-48>config>app-assure>group>policy>app-filter>entry# info
-----
      description "DNS traffic to local server on expected port #53"
      protocol eq "dns"
      flow-setup-direction subscriber-to-network
      ip-protocol-num eq *
      server-address eq 192.0.2.0/32
      server-port eq 53
      application "DNS_Local"
      no shutdown
-----
*A:ALA-48>config>app-assure>group>policy>app-filter>entry#
```

3.3.3.8 Configuring an application profile

Use the following CLI syntax to configure an application profile:

```
config>app-assure>group>policy# app-profile app-profile-name [create]
  characteristic characteristic-name value value-name
  description description-string
  [no] aa-sub-suppressible
  divert
```

The following example displays an application profile configuration:

```
*A:ALA-48>config>app-assure>group>policy# app-profile "Super" create
*A:ALA-48>config>app-assure>group>policy>app-prof# info
```

```

-----
description "Super User Application Profile"
divert
characteristic "Server" value "Prioritize"
characteristic "ServiceBw" value "SuperUser"
characteristic "Teleworker" value "Yes"
characteristic "VideoBoost" value "Priority"
-----
*A:ALA-48>config>app-assure>group>policy>app-prof#

```

3.3.3.9 Configuring suppressible app-profile with SRRP

For information about SRRP, see the *7450 ESS, 7750 SR, and VSR Triple Play Service Delivery Architecture Guide*.

In the context of an ESM SRRP deployment, the operator can define at the app-profile level if the subscriber is diverted to the ISA-AA card per SRRP group interface. This can be useful to reduce the total number of ISA cards required in the event of a switch-over from a primary to backup SRRP node when AA is used as a value-add service for selected subscribers.

To configure the network for suppressible app-profiles in the context of SRRP the operator needs to:

- Enable the capability to suppress AA subscribers on a specific SRRP group interface, typically by selecting backup SRRP group interfaces.
- ESM subscribers with a valid app-profile are diverted to AA by default, to suppress selected group of subscribers using AA for optional value-add services. The operator then specifies which app-profile is suppressed and therefore not diverted to AA.

Use the following CLI syntax to enable the capability to suppress ESM subscribers from a backup SRRP group interface:

```

config>service>vprn>sub-if>grp-if# suppress-aa-sub [create]
characteristic characteristic-name value value-name
description description-string
[no] aa-sub-suppressible
divert

```

The following example displays an application profile configuration used for premium subscribers, this type of subscriber is always be diverted to Application Assurance, it is also the default configuration:

```

7750>config>app-assure>group>policy# info
-----
app-profile "Premium" create
characteristic "Parental-Control" eq "Yes"
divert
exit
-----

```

The following example displays an application profile configuration for best effort / value-add subscribers not diverted to Application Assurance on the SRRP group interface configured with "suppress-aa-sub":

```

7750>config>app-assure>group>policy# info
-----
app-profile "1-default" create
divert
aa-sub-suppressible
exit
-----

```

3.3.3.10 Configuring application service options

Use the following CLI syntax to configure application service options:

```
config>app-assure>group>policy# app-service-options
  characteristic characteristic-name [create]
  default-value value-name
  value value-name
```

The following example displays an application service options configuration:

```
*A:ALA-48>config>app-assure>group>policy>aso# info
-----
      characteristic "Server" create
        value "Block"
        value "Permit"
        value "Prioritize"
        default-value "Block"
      exit
      characteristic "ServiceBw" create
        value "Lite_128k"
        value "Power_5M"
        value "Reg_1M"
        value "SuperUser"
        default-value "Reg_1M"
      exit
      characteristic "Teleworker" create
        value "No"
        value "Yes"
        default-value "No"
      exit
      characteristic "VideoBoost" create
        value "No"
        value "Priority"
        default-value "No"
      exit
-----
*A:ALA-48>config>app-assure>group>policy>aso#
```

3.3.3.11 Configuring a policer

Use the following CLI syntax to configure a policer:

```
config>app-assure>group>policy# policer policer-name type type granularity granularity create
  action {priority-mark | permit-deny}
  adaptation-rule pir adaptation-rule
  description description-string
  mbs maximum burst size
  rate pir-rate
  tod-override tod-override-id [create]
```

The following example displays an Application Assurance policer configuration.

```
*A:ALA-48>config>app-assure>group# policer "RegDown_Policer" type dual-bucket-
```

```
bandwidth granularity subscriber create

*A:ALA-48>config>app-assure>group>policer# info
-----
description "Control the downstream aggregate bandwidth for Regular
1Mbps subscribers"
rate 1000 cir 500
mbs 100
cbs 50
-----
*A:ALA-48>config>app-assure>group>policer#
```

3.3.3.12 Configuring an application QoS policy

Use the following CLI syntax to configure an application QoS policy:

```
config>app-assure>group>policy# app-qos-policy
entry entry-id [create]
action
  bandwidth-policer policer-name
  drop
  error-drop [event-log event-log-name]
  flow-count-limit policer-name
  flow-rate-limit policer-name
  fragment-drop {all | out-of-order} [event-log event-log-name]
  http-error-redirect redirect-name
  mirror-source [all-inclusive] mirror-service-id
  overload-drop [event-log event-log-name]
  remark
    dscp in-profile dscp-name out-profile dscp-name
    fc fc-name
    priority priority-level
  url-filter url-filter-name characteristic characteristic-name
description description-string
match
  aa-sub sap {eq | neq} sap-id
  aa-sub esm {eq | neq} sub-ident-string
  aa-sub spoke-sdp {eq | neq} sdp-id:vc-id
  app-group {eq | neq} application-group-name
  application {eq | neq} application-name
  characteristic characteristic-name {eq} value-name
  dscp {eq | neq} dscp-name
  dst-ip {eq | neq} ip-address[/mask]
  dst-ip {eq | neq} ip-prefix-list ip-prefix-list-name
  dst-port {eq | neq} port-num
  dst-port {eq | neq} range start-port-num end-port-num
  src-ip {eq | neq} ip-address[/mask]
  src-ip {eq | neq} ip-prefix-list ip-prefix-list-name
  src-port {eq | neq} port-num
  src-port {eq | neq} range start-port-num end-port-num
  traffic-direction {subscriber-to-network | network-to-subscriber | both}
no shutdown
```

The following example displays an application QoS policy configuration:

```
*A:ALA-48>config>app-assure>group>policy>aqp# entry 20 create
-----
description "Limit downstream bandwidth to Reg_1M subscribers"
match
```

```

        traffic-direction network-to-subscriber
        characteristic "ServiceBw" eq "Reg_1M"
    exit
    action
        bandwidth-policer "RegDown_Policer"
    exit
    no shutdown
-----
*A:ALA-48>config>app-assure>group>policy>aqp#

```

The following example displays an AQP entry configuration to mirror all positively identified only P2P traffic (AppGroup P2P) for a subset of subscribers with ASO characteristic **aa-sub-mirror** enabled:

```

A:ALA-48>config>app-assure>group>policy>aqp#
-----
entry 100 create
match
    app-group eq P2P
    characteristic aa-sub-mirror eq enabled
exit
action
    # mirror to an existing mirror service id
    mirror-source 100
exit
no shutdown
exit
-----
A:ALA-48>config>app-assure>group>policy>aqp#

```

The following example displays an AQP entry to mirror all P2P traffic (all positively identified P2P traffic and any unidentified traffic that may or may not be P2P - AppGroup P2P) for a subset of subscribers with ASO characteristic **aa-sub-mirror** enabled (the order is significant).

```

A:ALA-48>config>app-assure>group>policy>aqp>entry#
-----
entry 100 create
match
    app-group eq P2P
    characteristic aa-sub-mirror value enabled
exit
action
    mirror-source all-inclusive 100
exit
no shutdown
exit
-----
A:ALA-48>config>app-assure>group>policy>aqp#

```

3.3.3.13 Configuring a charging filter

Use the following CLI syntax to configure a charging filter:

```

config>app-assure>group>policy#
charging-filter entry <num>
match
    application eq|neq <app-name>
    app-group eq|neq <app-group-name>
    tethered-flow
    flow-attribute <name> confidence {eq|lt|gte} <value>

```

```
charging-group <charging-group-name>
[no] shutdown
[no] description
```

For example, the following charging filters match the Whatsapp video and audio traffic, and assign them to different charging groups.

Example: Charging filter configuration (classic CLI)

```
*A:node-2>config>app-assure>group>policy>chrg-fltr# info
entry 1 create
  description "Charging-filter entry for WhatsApp Video"
  match
    application eq "Whats App"
    flow-attribute "video" confidence gte 100
  exit
  charging-group "cgVideo"
  no shutdown
exit
entry 4 create
  description "Charging-filter entry for WhatsApp Audio Only"
  match
    application eq "Whats App"
    flow-attribute "audio" confidence gte 80
    flow-attribute "video" confidence eq 0
  exit
  charging-group "cgAudio"
  no shutdown
exit
```

3.3.3.14 Configuring an application and DNS IP cache for URL content charging strengthening

In the context of URL content charging, also known as zero rating, the DNS IP cache (**dns-ip-cache** command) feature ensures that only legitimate traffic is classified in an application and charging-group. Subscribers' DNS responses matching a list of domain names used for content charging populate the DNS IP cache. The system can then be configured to create app-filters matching HTTP or HTTPS expressions as well as the IP cache ensuring that traffic is properly classified. If the operator uses proxies in their network, they may also configure a maximum of 8 IP addresses which match the IP addresses of the proxies used. Traffic whose destination IP address matches one of the configured proxies is assumed to be legitimate traffic.

To configure the system for URL content charging strengthening with a dns-ip-cache the operator needs to:

- Create an application of interest and its related app-filter's URL expressions. This application is typically mapped into a charging-group.
- Create a **dns-ip-cache**. Configure parameters so the IP cache is populated by the domain names from the application mapped to the zero rating charging group and specify which DNS server IP addresses the IP cache listens from.
- Configure a AQP to enable the dns-ip-cache.
- Optionally, configure static IP addresses matching the IP addresses of the trusted proxies.

Use the following CLI syntax to create a dns-ip-cache:

```
config>app-assure>group#
```



```

dns-ip-cache dns-ip-cache-name [create]
  dns-match
    description <description-string>
    no description
    domain <domain-name> expression <expression>
    no domain <domain-name>
    server-address <server-address> [name <server-name>]
    no server-address <server-address>
  ip-cache
    size <cache-size>
    high-watermark <percent>
    low-watermark <percent>
    [no] static-address <static-ip-address>
  [no] shutdown

```

The following example displays a configuration for a **dns-ip-cache** configured to snoop DNS responses for two different domains ***.domain1.com** and ***.domain2.com** which are zero rated or charged specifically by the operator. The configuration only uses DNS responses from the DNS server addresses configured within the **dns-match** to populate the **ip-cache**:

```

7750>config>app-assure>group# info
-----
dns-ip-cache "dns-ip-cache1" create
  description "DNS IP Cache #1"
  dns-match
    domain "Sponsor#1-Domain#1" expression "*.domain1.com$"
    domain "Sponsor#1-Domain#2" expression "*.domain2.com$"
    server-address 10.8.4.4 name "CompanyName"
    server-address 10.8.8.8 name "CompanyName"
    server-address 192.168.100.11 name "OperatorX-DNS1"
    server-address 192.168.100.12 name "OperatorX-DNS2"
  exit
  ip-cache
    size 1000
    high-wmark 90
    low-wmark 80
  exit
  no shutdown
exit
-----

```

The domains configured in the **dns-ip-cache** must match the **app-filter** expressions for the applications zero rated or charged specifically by the operator. The following example displays the charging-group **Zero Rated** and application **Sponsor Content #1** configuration:

```

7750>config>app-assure>group>policy# info
-----
charging-group "Zero Rated" create
  description "Zero Rated Content"
  export-id 10
exit
app-group "Web" create
exit
application "Sponsor Content #1" create
  description "Application#1 - Content Zero Rated"
  app-group "Web"
  charging-group "Zero Rated"
exit
app-filter
  entry 100 create
    expression 1 http-host eq "*.sponsor1-domain1.com$"

```

```

        server-address eq dns-ip-cache "dns-ip-cache1"
        application "Sponsor Content #1"
        no shutdown
    exit
    entry 110 create
        expression 1 http-host eq "/*.domain2.com$"
        server-address eq dns-ip-cache "dns-ip-cache1"
        application "Sponsor Content #1"
        no shutdown
    exit
exit
-----

```

The following example displays the AQP entry to enable the **dns-ip-cache** to snoop DNS responses; this can be optionally based on ASO characteristics:

```

A:7750>config>app-assure>group>policy>aqp# entry 100 create
    match
        characteristic "dns-ip-cache" eq "yes"
    exit
    action
        action dns-ip-cache "dns-ip-cache1"
    exit
    no shutdown

```

3.3.3.15 Configuring an HTTP error redirect

Use the following CLI syntax to configure an HTTP error redirect policy:

```

config>app-assure>group>
    http-error-redirect redirect-name create
        no http-error-redirect redirect_name
        description description-string
        no description
        error-code error-code [custom-msg-size custom-msg-size]
        no error-code error-code
        http-host http-host // eg. www.demo.barefruit.com
        no http-host
        participant-id participant-id // 32-char string used by template 1
        no participant-id
        no] shutdown
        template template-id // {1, 2} one for Barefruit, 2= Xerocole
        no template

```

The following example displays an Application Assurance HTTP redirect configuration:

```

*A:ALA-48>config>app-assure>group# http-error-redirect "redirect-404" create
    description "redirect policy of 404 to Barefruit servers"
    error-code 404
    http-host
        att.barefruit.com
    participant-id att-ISP
    template 1

*A:ALA-48>config>app-assure>group> http-error-redirect# redirect-404 info
-----
        description "redirect policy of 404 to Barefruit servers"
        template 1

```

```

http-host "att.barefruit.com"
participant-id "att-ISP"

error-code 404

*A:ALA-48>config>app-assure>group>http-error-redirect#

```

3.3.3.16 Configuring HTTP header enrichment

Use the following CLI syntax to configure an AA HTTP header enrichment policy:

```

config>app-assure>group> http-enrich <http_enrich_name> [create]
[no] description <description-string>
[no] shutdown
[no] field <field_name> name <header_name>
// Where "Field name" can be:
// subscriber-ip: Header name for subscriber IP
// subscriber-id: Header name for the subscriber ID
// static-string: Header name for inserted string
// static-string-2: Header name for inserted string
//imsi: Header name for subscriber IMSI
//imsi-2: Header name for subscriber IMSI
//msisdn: Header name for subscriber MSISDN
//imei-sv: Header name for subscriber IMEI-SV
//rat-type: Radio Access technology
//pgw_ggsn-address: Header name for Packet Gateway (GGSN or UPF) IP address serving the
UE
//apn: Header name for APN used by the UE
//user-location: Header name for UE location(ULI)
//billing-type: Header name for charging type
//plmn-id: Public land mobile network ID of the SGSN/MME
//timestamp: Header name for the timestamp (inserted in unix time format. For example:
1531204313)
//apn-ni: Header name for APN-NI (Network Identifier) used by the UE
//msisdn-without-cc: Header name for subscriber MSISDN without country code
//imei-hyphenated: Subscriber's IMEI with format AABBBBBB-CCCCC-EE
//imei-hyphenated-2: Subscriber's IMEI with format AABBBBBB-CCCCC-EE
//user-location-raw: Header for ULI in raw format <uli-type1>[+<uli-type2>]=<ULI data
in hex>
//user-location-raw-2: Header for ULI in raw format <uli-type1>[+<uli-type2>]=<ULI data
in hex>
//user-location-3gpp: ULI encoded as defined in 3GPP TS2.061
//static-acr: static ACR
//dynamic-acr: dynamic ACR
[no] http-enrich <http_enrich_name>

```

The following example displays an AA HTTP header enrichment configuration:

```

*A:BNG>config>app-assure>group# http-enrich enrich_example create
*A:BNG>config>app-assure>group>http-enrich$ description "enrich HTTP headers with
subscriber IP and subscriber ID"
*A:BNG>config>app-assure>group>http-enrich$ field "static-string" name "x-string"
*A:BNG>config>app-assure>group>http-enrich$ field "static-string" static-string
"orange"
*A:BNG>config>app-assure>group>http-enrich$ field "subscriber-id" name "x-subID"
*A:BNG>config>app-assure>group>http-enrich$ field "subscriber-id" anti-spoof
*A:BNG>config>app-assure>group>http-enrich$ field "subscriber-ip" name x-subIP
*A:BNG>config>app-assure>group>http-enrich$ field "subscriber-ip" encode type md5 key
"secret10"
-----

```

```
*A:BNG>config>app-assure>group>http-enrich$ info
-----
      field "static-string"
        name "x-string"
        static-string "orange"
      exit
      field "subscriber-id"
        name "x-subID"
        anti-spoof
      exit
      field "subscriber-ip"
        name "x-subIP"
        encode type md5 key
"bF0sZZDNT8DbZoVJHD1vrYr5mJaEggEqWbSvPhgIcPW6hym0sc080." hash2
      exit
-----
*A:BNG>config>app-assure>group>http-enrich$
```

In addition, the following **show** routine displays various HTTP enrichment-related statistics:

```
show application-assurance group 1 http-enrich "MyTemplate"
=====
Application Assurance Group 1 HTTP Enrichment "MyTemplate"
=====
Description   : (Not Specified)
Admin Status  : Up
AQP Referenced: No

-----
Name           Field           Enabled
                Features
-----
subscriber-ip  MyField         M
subscriber-id  Sub-ID          AC
-----
                                     A=anti-spoof,C=encode-cert,M=encode-md5,R=encode-rc4

-----
Group          Enriched          Not Enriched
-----
1              0                 0
-----
Total          0                 0
-----
=====
```

3.3.3.17 Configuring an HTTP redirect policy

Use the following CLI syntax to configure an HTTP redirect policy:

```
config>app-assure>group# http-redirect redirect-name [create]
  description <description-string>
  no description
  [no] redirect-https
  [no] template <template-id>
  redirect-url URL // redirect URL e.g. www.isp.com/redirect.html
  no redirect-url
  [no] shutdown
  [no] tcp-client-reset
  no http-redirect <redirect-name>
```

The following example displays an AA HTTP redirect configuration:

```
*A:ALA-48>config>app-assure>group# http-redirect "redirect1" create
  description "redirect policy for blocked http content traffic without url
  parameters"
  template 3
  redirect-url http://www.isp.com/redirect.html
  no shutdown
```

The following example displays an Application Assurance **http-redirect** configuration using macro substitution to append url parameters within the redirect url:

```
*A:ALA-48>config>app-assure>group# http-redirect "redirect2" create
  description "redirect policy for blocked http content traffic with url parameters"
  template 5
  redirect-url "http://www.isp.com/
  redirect.html?requestedurl=$URL&subscriberid=$SUB&subscriberip=$IP&routerid=$RTRI
  D&vsa=$URLPRM"
  no shutdown
```

The following example displays AQP entry to block all http gaming traffic (AppGroup BlockedContent) and perform redirect:

```
A:ALA-48>config>app-assure>group>policy>aqp>entry#
-----
  entry 100 create
    match
      app-group eq BlockedContent
    exit
    action
      drop
      http-redirect redirectgaming flow-type dropped-flows
    exit
  no shutdown
  exit
-----
A:ALA-48>config>app-assure>group>policy>aqp#
```

3.3.3.18 Configuring an HTTPS policy redirect

Use the following CLI syntax to configure an HTTPS redirect policy:

```
config>app-assure>group# http-redirect redirect-name [create]
  description <description-string>
  no description
  [no] redirect-https
  redirect-url <redirect-url>
  no redirect-url
  [no] shutdown
  [no] tcp-client-reset
  template <template-id>
  no template
  no http-redirect <redirect-name>
```

The following example displays the configuration of an HTTPS redirect object to redirect traffic to the operator's portal:

```
A:ALA-1>config>app-assure>group# http-redirect "https-redirect-portal" create
description "simple https redirect policy"
redirect-url https://www.portal.com/info.html
redirect-https
no shutdown
```

After creating the HTTP Redirect object, an AQP is used to configure the traffic to be redirected:

```
A:ALA-1>config>app-assure>group>policy>aqp>entry#
-----
entry 100 create
match
  application eq BlockedHTTPS_page
exit
action
  drop
  http-redirect https-redirect-portal flow-type dropped-flows
exit
no shutdown
exit
-----
```

3.3.3.19 Configuring a captive redirect HTTP redirect policy

The traditional HTTP redirect policy is used to redirect flows on the HTTP response packet, meaning the TCP three-way handshake and the original HTTP request are allowed by the 7750 SR to the Internet before the subscriber is redirected. The captive redirect HTTP redirect policy is used to redirect flows without sending any traffic to the Internet unless it matches a configurable allowlist by terminating TCP sessions in the ISA-AA cards, in which case HTTP flows are redirected to a predefined redirect URL while non-HTTP TCP flows are TCP reset.

3.3.3.19.1 Captive redirect and HTTPS flows redirection

The captive redirect HTTP redirect policy can be optionally configured to redirect HTTPS sessions in addition to HTTP to a pre-defined redirect landing page, typically the captive-portal URL in the context of a WiFi network. This capability is particularly useful when the router is used to provide a captive-portal type of access, as it allows the operator to improve the user experience by redirecting the subscriber's web browser sessions to the needed captive-portal landing page when the user first connects to the network using HTTPS instead of HTTP.

Before the introduction of this feature, users opening their web browsers to an HTTPS URL when first connecting to a new Wi-Fi network and expecting to be redirected to a captive portal were instead presented with an error page automatically generated by the web browser because the session was dropped or reset by the network, therefore ultimately preventing the user from connecting. Most non-technical users connecting to a captive-portal network may not know the difference between HTTP and HTTPS when it comes to login/redirection, and a number of subscribers may not connect or may get frustrated trying multiple different links before a successful Wi-Fi authentication.

When the system is configured for captive-redirect **redirect-https**, it terminates transport layer security (TLS) TCP sessions in the ISA-AA cards and return a self-signed certificate back to the user. Upon the

user acceptance of the security warning generated by the web browser, the web session then automatically redirects to the configured captive-portal landing page.

Captive redirect policy supports redirection for HTTP, HTTPS, HTTP2, SPDY, and TCP Fast Open connections.

A **session-filter** is used to define the criteria for permitting or redirecting flows using the captive redirect HTTP redirect policy. Typically the operator needs to allow UDP on port 53 for DNS and they can optionally allow other content based on IP address, port number, IP prefix list, or DNS IP cache therefore allowing specific on-net or off-net applications through the captive redirect policy.

To configure the system for captive redirect HTTP redirect the operator needs to:

- Create an http-redirect policy. If the ISA group aa-sub-scale mode is configured for residential or VPN, then configure the http-redirect policy for captive-redirect and associate the appropriate VLAN ID AA Interface (an aa-interface routable within the subscriber's service must be created for each ISA-AA card in the system). If the ISA group aa-sub-scale mode is configured for DSM, then there is no need to associate the http-redirect policy to a VLAN ID and no need to create an AA Interface.
- Create a session filter policy to allow at the minimum UDP on port 53. Additional traffic can be allowlisted based on a statically defined IP prefix list or a dynamic DNS IP cache policy. The redirect landing page should be configured using IP prefixes.
- The last action of the session filter should be set to http-redirect the remaining flows using a predefined captive redirect HTTP redirect policy.

Use the following CLI syntax to create a captive redirect HTTP redirect policy:

```
config>app-assure>group# http-redirect <redirect-name> [create]
    captive-redirect
        vlan-id <service-port-vlan-id>
        no vlan-id
    description <description-string>
    no description
    [no] redirect-https
    redirect-url <redirect-url>
    no redirect-url
    [no] shutdown
    [no] tcp-client-reset
    template <template-id>
    no template
no http-redirect <redirect-name>
```

The following example displays a typical configuration for a session filter user in the context of captive redirect:

```
A:7750# configure application-assurance group 1:1 create
A:7750>config>app-assure>group# info
-----
    session-filter "wifi-unauthenticated" create
        default-action deny
        entry 5 create
            match
                ip-protocol-num udp
                dst-port eq 53
            exit
            action permit
        exit
        entry 10 create
            match
                dst-ip dns-ip-cache "whitelist"
```

```

        exit
        action permit
    exit
    entry 15 create
        description "Allow traffic to the redirect landing page server"
        match
            ip-protocol-num tcp
            dst-port eq 80
            dst-ip 172.16.70.100/32
        exit
        action permit
    exit
    entry 20 create
        match
            ip-protocol-num tcp
        exit
        action http-redirect "redirect-portal"
    exit
exit
-----

```

The following example displays a typical configuration for the AA interface used by the captive redirect HTTPS redirect policy for ESM Subscribers (DSM does not require the configuration of the AA Interface):

```

A:7750# configure service ies 1 customer 1 create
A:7750>config>service>ies# info
-----
    aa-interface "aa-if-captive-redirect-isa_1-2" create
        description "AA Interface for ISA-AA card 1/2"
        address 172.16.3.1/31
        sap 1/2/aa-svc:20 create
            no shutdown
        exit
    no shutdown
exit
-----

```

The following example displays a typical configuration for the HTTPS redirect policy for ESM Subscribers (DSM does not require the configuration of the VLAN ID):

```

A:7750# configure application-assurance group 1
A:7750>config>app-assure>group>http-redir# info
-----
    template 5
        tcp-client-reset
        redirect-https
        redirect-url "http://172.16.70.100/Redirect/redirect-portal.html?RequestedURL=$URL"
        captive-redirect
            vlan-id 20
    exit
    no shutdown
-----

```

3.3.3.20 Configuring ICAP URL filtering

To configure the system for ICAP URL filtering, the operator needs to:

- Create an AA interface and assign an IP address to each AA ISA within an IES or VPRN service. This routed interface is then used by the system to establish TCP communication with the ICAP server.
- Create an HTTP redirect policy (used by the URL filter to redirect HTTP or HTTPS traffic).
- Create a URL filter, configure the ICAP server IP address, redirect policy, and default action.
- Verify that the AA interface and URL filter are operationally up.

Use the following CLI syntax to configure the AA interface for each AA ISA:

```
config>service>vprn# aa-interface <aa-if-name> [create]
config>service>vprn>aa-if# aa-interface interface <ip-int-name> [create]
description <description-string>
no description
address <ipv4_subnet/31>
no address
  sap <card/mda/aa-svc:vlan> [create]
  description <description-string>
  no description
  egress
    [no] filter
    [no] qos
  exit
  ingress
    [no] qos
  exit
  [no] shutdown
exit
```

The following example displays an AA interface created for the ISA card 1/2 using IP address 172.16.2.1/31:

```
A:7750>config>service>ies# info
-----
aa-interface "aa-if1" create
  address 172.16.2.1/31
  sap 1/2/aa-svc:10 create
  egress
    filter ip 10
  exit
  no shutdown
exit
no shutdown
exit
```

In the example above, 172.16.2.1 is used by the IOM side of the interface while the ISA itself is automatically assigned 172.16.2.0 based on the /31 subnet.

Consider the following recommendations:

- More than one AA interface can be configured for each AA ISA card, however, the operator must use the same service VLAN across all these interfaces for a URL filter object.
- Configure an egress IP filter under the SAP toward the ISA AA interface to only allow selected IP addresses or subnets (for example, ICAP servers or network management).

Use the following CLI syntax to configure the URL filter:

```
config>app-assure>group#
url-filter <url-filter-name> [create]
  default-action {allow | block-all | block-http-redirect <redirect-name>}
no default-action
```

```

http-redirect <http-redirect-name>
no http-redirect
http-request-filtering {all | first}
no apply-function-specific-behavior
icap
  custom-x-header <x-header-name>
  [no] custom-x-header
  vlan-id <service-port-vlan-id>
  no vlan-id
  default-action {allow | block-all | block-http-redirect <redirect-name>}
  no default-action
  http-redirect <http-redirect-name>no http-redirect
  server <ip-address[:port]> [create]
    description <description-string>
    no description
    [no] shutdown
  no server <ip-address[:port]>
no url-filter <url-filter-name>

```

If the **apply-function-specific-behavior** command is enabled, the operator may configure an ICAP-specific **default-action** and **redirect**. This is done by configuring **default-action** and **http-redirect** in **config>app-assure>group>url-filter>icap** context.

The following example displays a URL filter configuration:

```

*A:7750>config>app-assure>group# url-filter "filter1" create
  apply-function-specific-behavior
  icap
    vlan-id 10
    default-action block-http-redirect "http-redirect-portal"
    server 172.16.1.101 create
      no shutdown
    exit
  exit
no shutdown

```

Enabling the **apply-function-specific-behavior** command allows the operator to configure a **default-action** and **HTTP redirect** which are specific to ICAP URL filtering only. Alternatively, the operator configures the same **default-action** and **http-redirect** for all **url-filter** functions by disabling the **apply-function-specific-behavior** and configuring a **default-action** and **http-redirect** in the **config>app-assure>group>url-filter** context.

The following example displays an AQP entry to enable ICAP URL filtering for opted-in subscribers based on ASO characteristics:

```

A:7750>config>app-assure>group>policy>aqp# entry 100 create
  match
    characteristic "url-filter" eq "yes"
  exit
  action
    url-filter "filter1"
  exit
no shutdown

```

Optionally, the operator can add a custom x-header to the ICAP request in order for the ICAP server to filter traffic based on this new x-header value instead of filtering based on subscriber names. This is done by defining a new ASO characteristic for the different ICAP filtering service packages used in the network and referring the characteristic name in the URL filter AQP action.

The following example displays a URL filter configuration with the **custom-x-header** field added to the ICAP request:

```
A:7750>config>app-assure>group# url-filter "filter1" create
  apply-function-specific-behavior
  icap
  default-action block-http-redirect "http-redirect-portal"
  http-redirect "http-redirect-portal"
  custom-x-header "Filtering-Policy"
  vlan-id 10
  server 172.16.1.101 create
    no shutdown
  exit
exit
no shutdown
```

The following example displays the ASO characteristic used to define the type of filtering policy available:

```
A:7750>config>app-assure>group>policy>aso# info
-----
      characteristic "url-filter-policy" create
        value "filtering-policy-1" #less than 10 years old
        value "filtering-policy-2" # less than 16 years old
        value "example1"
        value "none"
        value "example2"
        default-value "none"
      exit
-----
```

The following example displays the AQP action required to add the appropriate ASO value to the ICAP **custom-x-header** field:

```
A:7750>config>app-assure>group>policy>aqp# entry 100 create
  match
    characteristic "url-filter" eq "yes"
  exit
  action
    url-filter "filter1" characteristic "url-filter-policy"
  exit
no shutdown
```

3.3.3.21 Configuring web-service URL classification

About this task

The following example describes how to configure web-service URL classification.

Procedure

Step 1. Configure the shared resources to be used for the cache.

In the following example configuration, 100% of shared resources are allocated to web-service URL classification:

```
config>isa>aa-grp>shr-res-pool# info detail
-----
      web-service-url-filter 100
```

Step 2. Configure the HTTP redirect policy.

Traffic that belongs to one of the blocked categories (as defined in the profile activated for the user), is redirected.

The operator can append the category name and category ID to the redirect URL (that is, the category that resulted in the redirect action).

The following example displays an HTTP redirect policy configuration:

```
config>app-assure>group 1
-----
http-redirect "redirect-ws-filter" create
description "Redirect for Web Service URL Filtering"
template 5
tcp-client-reset
redirect-https
redirect-url "http://10.154.90.1/
redirect.html?catname=$CATNAME&catid=$CATID"
no shutdown
exit
-----
```

Step 3. Configure the URL filter and category profiles.

The operator defines a new URL filter with the following attributes:

- classifier = "web-service-1"
defines the URL categorization database to use
- category-set-id = 1
defines the URL categories to use
- fqdn = nokia-api.webtitancloud.com
defines the URL categorization database endpoint that the AA connects to
- dns-server
defines the DNS server to resolve the FQDN
- profile
defines up to eight profile names
- category
defines the category to be blocked in the configured profile

In the following URL filter configuration example, the three profiles defined in [Table 21: Example URL category profile](#) of the [Web-service URL classification](#) section are configured. Two **classification-overrides** entries are also configured:

```
configure application-assurance group 1 url-filter "ws-filter"
config>app-assure>group>url-filter# info
-----
apply-function-specific-behavior
web-service
classifier web-service-1 category-set-id 1
vlan-id 100
default-action block-http-redirect "redirect-ws-filter"
http-redirect "redirect-ws-filter"
```

```

fqdn nokia-api.webtitancloud.com
dns-server 8.8.8.8
classification-overrides
entry 1 expression www.site1.abc category "Phishing/Fraud"
entry 2 expression www.site2.def category "Illegal Drugs"
profile low create
    category "Internet Watch Foundation List" block
    category "Spyware And Malicious Sites" block
    category "Phishing/Fraud" block
profile medium create
    category "Internet Watch Foundation List" block
    category "Spyware And Malicious Sites" block
    category "Phishing/Fraud" block
    category "Illegal Drugs" block
    category "Violence" block
    category "Weapons" block
profile high create
    category "Internet Watch Foundation List" block
    category "Spyware And Malicious Sites" block
    category "Phishing/Fraud" block
    category "Illegal Drugs" block
    category "Violence" block
    category "Weapons" block
    category "Nudity" block
    category "Alcohol" block
    category "Criminal Skills/Hacking" block
    category "Hate Speech" block
default-category-profile "low"
no shutdown
-----

```

Enabling **apply-function-specific-behavior** allows the operator to configure a **default-action** and **http-redirect** which are specific to web-service URL Classification only. Alternatively, the operator may configure the same **default-action** and **http-redirect** for all **url-filter** functions by disabling the **apply-function-specific behavior** (which is the default) and configuring a **default-action** and **http-redirect** in the **config>app-assure>group>url-filter**.

Step 4. Configure the ASO.

The ASO is used to dynamically allocate a profile to a user.

The following output displays the configuration of an example ASO:

```

config>app-assure>group>policy>aso# info
-----
characteristic "url-filter-policy" create
value "high"
value "medium"
value "low"
default-value "low"
exit
-----

```

Step 5. Configure the AQP.

The AQP is used to execute the URL filter policy. The URL filter uses the ASO value that is active for the user to select the category profile.

In the AQP defined in the following example, there is no match condition. Therefore, the web-service URL classification is applied to all subscribers:

```

config>app-assure>group>policy>aqp# entry 100 create
action

```

```
url-filter "ws-filter" characteristic "url-filter-policy"
exit
no shutdown
```

3.3.3.22 Configuring local URL-list filtering

To configure the system for local URL-list filtering, the operator needs to:

- Create a URL-list policy referencing a valid file located on the compact flash.
- Create a url-filter policy for local-filtering by referencing this URL-list.
- Create an AQP to apply this url-filter policy.

Use the following CLI syntax to create a URL-list:

```
config>app-assure>group# url-list <url-list-name> [create]
description <description-string>
no description
decrypt-key <key | hash-key | hash2-key> [hash | hash2 | custom]
no decrypt-key
file <file-url>
no file
[no] shutdown
size <url-list-size>
[no] expression-match
[standard | extended] - Default : standard
```

Wildcards are supported on hostname entries. To enable wildcard support, the **url-list** must have **expression-match** (in **config>app-assure-group>url-list**) enabled (the default is disabled). An entry may contain the following:

- head anchors character set [^ *]
- tail anchors character set [\$ *]
- mid expression character set [d \. **]
- hex escape characters [x00-\xFF]

The same capabilities with those described in *Application Filters* section are provided.

When **expression-match** is set to enabled, the list should contain hostnames only, with wildcards.

The decryption key is optional. If the decryption key is not specified, the system assumes that the file is not encrypted. To encrypt a file in Linux using the supported encryption format, use the following command:

```
Linux# openssl des3 -nosalt -in <input-file.txt> -out <output.enc>
```

The following example displays a URL list configuration:

```
A:7750>config>app-assure>group# url-list url-list1 create
-----
description "Local List for URL Filtering"
decrypt-key ".i84/PluS0lMGoQkae7mAV20j10n726Z" hash2
file "cf3:\url-list1.enc"
no shutdown
-----
```

Use the following CLI syntax to create a url-filter policy for local-filtering:

```
config>app-assure>group# url-filter <url-filter-name> [create]
url-filter <url-filter-name> [create]
description <description-string>
no description
default-action {allow | block-all | block-http-redirect <redirect-name>}
no default-action
[no] http-redirect <redirect-name>
http-request-filtering {all | first}
[no] apply-function-specific-behavior
local-filtering
    deny-list <url-list-name>
        default-action {allow | block-all | block-http-redirect <redirect-name>}
        no default-action
        [no] http-redirect <redirect-name>
    [no] allow-list <url-list-name>
[no] shutdown
```

The following example displays a deny-list URL filter configured for local-filtering:

```
A:7750>config>app-assure>group# url-filter "url-denylist1" create
A:7750>config>app-assure>group>url-filter# info
-----
apply-function-specific-behavior
local-filtering
    default-action allow
    http-redirect "http-redirect-portal"
    deny-list "url-list1"
exit
no shutdown
-----
```

The default action should always be configured to "allow" when the url-filter is configured for local-filtering. The default-action in this context represents the action the system takes in case the local-list file is not accessible; this scenario may happen if the source file was corrupted or if the compact flash card was not accessible.

The following example displays the AQP entry to enable ICAP url-filtering for opted-in subscribers based on ASO characteristics:

```
A:7750>config>app-assure>group>policy>aqp# entry 100 create
    match
        characteristic "child-protection" eq "yes"
    exit
    action
        url-filter "url-blacklist1"
    exit
no shutdown
```

3.3.3.23 Configuring HTTP notification

Use the following CLI syntax to configure an HTTP Notification policy:

```
config>app-assure>group#
    http-notification <http-notification-name> [create]
    description <description-string>
    no description
```

```

script-url <script-url-name>
no script-url
interval {one-time | <minimum-interval>}
template <template-id>
no template
[no] shutdown
no http-notification <http-notification-name>

```

The following example displays an HTTP notification policy configured with a minimum interval of 5 minutes:

```

A:7750>config>app-assure>group# http-notification "in-browser-notification" create
A:7750>config>app-assure>group>http-notif# info
-----
description "In Browser Notification Example"
template 1
script-url "http://10.1.1.1/In-Browser-Notification/script.js"
interval 5
no shutdown
-----

```

The operator then needs to enable the http-match-all-req feature for any HTTP request sent the messaging server domain which is used to monitor HTTP notification success/failures. This is done by creating a new application and enabling http-match-all-req within the app-filter.

```

A:7750>config>app-assure>group>policy# application "IBN Messaging Server" create
A:7750>config>app-assure>group>policy>app$ app-group "Web"

A:7750>config>app-assure>group>policy# app-filter entry 100 create
A:7750>config>app-assure>group>policy>app-filter>entry$ info
-----
expression 1 http-host eq "^10.1.1.1$"
http-match-all-req
application "IBN Messaging Server"
no shutdown
-----

```

The following example displays the AQP entry required to match this policy based on an ASO characteristic:

```

A:7750>config>app-assure>group>policy>aqp# info
-----
entry 200 create
match
characteristic "in-browser-notification" eq "yes"
exit
action
http-notification "in-browser-notification"
exit
no shutdown
exit
-----

```

3.3.3.24 Configuring tethering detection

To configure tethering detection, enable the tethering detection option and specify the tethering state of the subscriber for which the specified application QoS policy match entry is applied.

Use the following CLI syntax to configure tethering detection:

```
configure application-assurance group group-number tethering-detection no shutdown
```

Use the following CLI syntax to specify the tethering state of the subscriber:

```
configure application-assurance group group-number policy
  begin
    app-qos-policy entry match aa-sub-tethering-state {detected | not-detected}
  commit
exit
```

3.3.4 Configuring AA volume accounting and statistics

A network operator can configure AA volume statistic collection and accounting on both AA ISA system and subscriber levels.

The following commands illustrate the configuration of statistics collection and accounting policy on an AA group/partition aggregate level (without subscriber context).

CLI syntax:

```
config>app-assure>group>statistics>app-group
  accounting-policy act-policy-id
  collect-stats
```

CLI syntax:

```
config>app-assure>group>statistics>application
  accounting-policy act-policy-id
  collect-stats
```

CLI syntax:

```
config>app-assure>group>statistics>protocol
  accounting-policy act-policy-id
  collect-stats
```

The following commands illustrate the configuration of statistics collection and accounting policy for each AA subscriber in the system:

```
config>app-assure>group>statistics>aa-sub
  accounting-policy acct-policy-id
  aggregate-stats
  app-group app-group-name export-using export-method [export-method]
  application application-name export-using export-method [export-method]
  charging-group charging-group-name export-using export-method [export-method]
  collect-stats
  exclude-tcp-retrans
  max-throughput-stats
  protocol protocol-name export-using export-method
  radius-accounting-policy rad-acct-plcy-name
```

The following commands illustrate configuration of special study mode for a subset of AA subscribers (configured) to collect all protocol or application statistics with an AA subscriber context:

```
config>app-assure>group>statistics# aa-sub-study {application | protocol}
  accounting-policy acct-policy-id
  collect-stats
```

For details on accounting policy configuration (including among others AA record type selection and customized AA subscriber record configuration), see the *7450 ESS, 7750 SR, 7950 XRS, and VSR System Management Guide*.

The following output illustrates per AA-subscriber statistics configuration that elects statistic collection for a small subset of all application groups, applications, protocols:

```
*A:ALU-40>config>app-assure>group>statistics>aa-sub# info
-----
      accounting-policy 4
      collect-stats
      app-group "File Transfer"
      app-group "Infrastructure"
      app-group "Instant Messaging"
      app-group "Local Content"
      app-group "Mail"
      app-group "MultiMedia"
      app-group "Business_Critical"
      app-group "Peer to Peer"
      app-group "Premium Partner"
      app-group "Remote Connectivity"
      app-group "Tunneling"
      app-group "Unknown"
      app-group "VoIP"
      app-group "Web"
      app-group "Intranet"
      application "BitTorrent"
      application "eLearning"
      application "GRE"
      application "H323"
      application "TLS"
      application "HTTP"
      application "HTTPS"
      application "HTTPS_Server"
      application "HTTP_Audio"
      application "HTTP_Video"
      application "eMail_Business"
      application "eMail_Other"
      application "Oracle"
      application "Skype"
      application "SAP"
      application "SIP"
      application "SMTP"
      application "SQL_Alltypes"
      application "TFTP"
      protocol "bittorrent"
      protocol "dns"
      protocol "sap"
      protocol "skype"
-----
*A:ALU-40>config>app-assure>group>statistics>aa-sub#
```

3.3.4.1 Configuring cflowd collector

The following output displays an AA cflowd collector configuration example:

```
*A:ALA-48# configure application-assurance group 1 cflowd collector 192.168.131.149:55000
create
  *A:ALA-48>config>app-assure>group>cflowd>collector$description
  "cflowd_collector_NewYork"
  *A:ALA-48>config>app-assure>group>cflowd>collector# no shutdown
  *A:ALA-48>config>app-assure>group>cflowd>collector# exit
```

The following output displays the configuration of the newly created example AA cflowd collector:

```
*A:ALA-48>config>app-assure>group>cflowd# info
-----
collector 192.168.131.149:55000 create
  description "cflowd_collector_NewYork"
  no shutdown
-----
*A:ALA-48>config>app-assure>group>cflowd#
```

Use the following CLI syntax to configure an AA cflowd collector:

```
config>application-assurance>group isa-aa-group-id cflowd
  collector <ip-address[:port]> [create]
  no collector <ip-address[:port]>
  [no] description - Configure description for this cflowd collector
  [no] shutdown - Administratively enable/disable this cflowd collector
  direct-export
  collector <collector-id> [create]
  no collector <collector-id>
  [no] address + Configure cflowd direct export collector remote address
  [no] description - Configure description for this cflowd direct export
```

The following output displays an AA direct export cflowd collector configuration example:

```
*A:Dut-C>config>app-assure>group>cflowd# direct-export vlan-id 300
*A:Dut-C>config>app-assure>group>cflowd# direct-export collector 1 create
*A:Dut-C>config>app-assure>group>cflowd>dir-exp-coll$ description "direct export to collector
in Toronto/CA"
*A:Dut-C>config>app-assure>group>cflowd>dir-exp-coll$ address 10.10.1.1:55000
*A:Dut-C>config>app-assure>group>cflowd>dir-exp-coll>addr$ no shutdown
*A:Dut-C>config>app-assure>group>cflowd>dir-exp-coll>addr$ exit
*A:Dut-C>config>app-assure>group>cflowd>dir-exp-coll$ exit
*A:Dut-C>config>app-assure>group>cflowd# no shutdown
*A:Dut-C>config>app-assure>group>cflowd# info
-----
no shutdown
direct-export
  vlan-id 300
  collector 1 create
  address 10.10.1.1:55000
  no shutdown
  exit
  description "direct export to collector in Toronto/CA"
  exit
exit
-----
*A:Dut-C>config>app-assure>group>cflowd#
```



Note: The VLAN-ID must match the one configured under the AA ISA interface. For an example of how to configure an AA interface, see [Configuring ICAP URL filtering](#).

3.3.4.2 Configuring AA comprehensive, RTP performance, TCP performance, and volume reporting

The following commands show the configuration of AA comprehensive reporting:

```
config>application-assurance>group isa-aa-group-id cflowd comprehensive
[no] flow-rate - Configure cflowd comprehensive flow sample rate
[no] flow-rate2 - Configure secondary cflowd comprehensive flow sample rate
[no] shutdown - Administratively enable/disable comprehensive sampling
        template + Configure comprehensive template|
        [no] dynamic-fields + Configure the list of dynamic fields
                [no] field <field-name> - Configure a dynamic field
                [no] shutdown - Administratively enable/disable comprehensive template
        [no] field-selection <field-selection>      Configure the field selection :
legacy|dynamic
```

The following commands show the configuration of AA RTP performance reporting:

```
config>application-assurance>group isa-aa-group-id cflowd rtp-performance
[no] flow-rate - Configure cflowd RTP performance flow sample rate
[no] flow-rate2 - Configure secondary cflowd RTP performance flow sample rate
[no] shutdown - Administratively enable/disable RTP performance sampling

        audio-template + Configure rtp performance audio fields template
                [no] dynamic-fields + Configure the list of dynamic fields
                        [no] field <field-name> - Configure a dynamic field
                        [no] shutdown - Administratively enable/disable audio template
                [no] field-selection <field-selection>      Configure the field selection :
legacy|dynamic

        video-template + Configure rtp performance video fields template
                [no] dynamic-fields + Configure the list of dynamic fields
                        [no] field <field-name> - Configure a dynamic field
                        [no] shutdown - Administratively enable/disable video template
                [no] field-selection <field-selection>      Configure the field selection :
legacy|dynamic

        voice-template + Configure rtp performance voice template
                [no] dynamic-fields + Configure the list of dynamic fields
                        [no] field <field-name> - Configure a dynamic field
                        [no] shutdown - Administratively enable/disable voice template
                [no] field-selection <field-selection>      Configure the field selection :
legacy|dynamic
```

The following commands show the configuration of AA TCP performance reporting:

```
config>application-assurance>group isa-aa-group-id cflowd tcp-performance
[no] flow-rate - Configure cflowd TCP performance flow sample rate
[no] flow-rate2 - Configure secondary cflowd TCP performance flow sample rate
[no] shutdown - Administratively enable/disable TCP Performance sampling
        template + Configure tcp performance fields template
                [no] dynamic-fields + Configure the list of dynamic fields
                        [no] field <field-name> - Configure a dynamic field
                        [no] shutdown - Administratively enable/disable tcp-performance template
                [no] field-selection <field-selection>      Configure the field selection : legacy|
dynamic
```

The following commands show the configuration of AA volume reporting:

```
config>application-assurance>group isa-aa-group-id cflowd volume
  [no] rate - Configure cflowd volume sampling rate
  [no] shutdown - Administratively enable/disable volume sampling
  template + Configure volume template
    [no] dynamic-fields + Configure the list of dynamic fields
      [no] field <field-name> - Configure a dynamic field
      [no] shutdown - Administratively enable/disable volume template
    [no] field-selection <field-selection> Configure the field selection : legacy|
dynamic
```

The following commands show the configuration of AA comprehensive, RTP performance, TCP performance, and volume reporting:

```
config>application-assurance group isa-aa-group-id[:partition [create]]
  *A:Dut-C>config>app-assure>group>cflowd#
    comprehensive + Configure cflowd comprehensive export
      [no] app-group app-group-name [flow-rate | flow-rate2]
      [no] application application-name [flow-rate | flowrate2]
      [no] shutdown - Administratively enable/disable comprehensive sampling

    rtp-performance + Configure cflowd RTP performance export
      [no] app-group app-group-name [flow-rate | flow-rate2]
      [no] application application-name [flow-rate | flowrate2]
      [no] shutdown - Administratively enable/disable RTP performance sampling

    tcp-performance + Configure cflowd TCP performance export
      [no] app-group app-group-name [flow-rate | flow-rate2]
      [no] application application-name [flow-rate | flowrate2]
      [no] shutdown - Administratively enable/disable TCP performance sampling

    volume + Configure cflowd volume export
      [no] shutdown - Administratively enable/disable volume sampling
```

The following example shows a configuration that includes the following:

- Enables per-flow volume statistics for group 1, partition 1 and configures sampling rate to 1/1000.
- Enables per-flow TCP performance statistics for web_traffic application within group 1, partition 1 and configures TCP sampling rate to 1/500.
- Enables per-flow TCP performance statistics for citrix_traffic application within group 1, partition 1 using TCP sampling rate2 to 1/100.
- Enables per-flow RTP A/V performance statistics for voip_traffic application within group 1, partition 1 and configures RTP sampling rate to 1/10.
- Enables per-flow comprehensive statistics for web_traffic.

Example:

```
*A:ALA-4# configure application-assurance group 1 cflowd
*A:ALA-4>config>app-assure>group>cflowd# volume rate 1000
*A:ALA-4>config>app-assure>group>cflowd# tcp-performance flow-rate 500
*A:ALA-4>config>app-assure>group>cflowd# tcp-performance flow-rate2 100
*A:ALA-4>config>app-assure>group>cflowd# comprehensive flow-rate 5
*A:ALA-4>config>app-assure>group>cflowd# rtp-performance flow-rate 10
*A:ALA-4>config>app-assure>group>cflowd# comprehensive flow-rate 5
*A:ALA-4>config>app-assure>group>cflowd# no shutdown
*A:Dut-C>config>app-assure>group>cflowd# info
```

```

-----
        shutdown
        direct-export
            vlan-id 300
            collector 1 create
                address 10.10.1.1:55000
                shutdown
            exit
        description "direct export to collector in Toronto/CA"
        exit
    exit
volume
    rate 1000
exit
rtp-performance
    flow-rate 10
    flow-rate2 5
exit
tcp-performance
    flow-rate 500
    flow-rate2 100
exit
-----
*A:Dut-C>config>app-assure>group>cflowd#
-----
*A:ALA-48>config>app-assure>group>cflowd#
*A:ALA-4# configure application-assurance group 1:1 cflowd
*A:ALA-4>config>app-assure>group>cflowd#
*A:ALA-4>config>app-assure>group>cflowd# volume no shutdown
*A:ALA-4>config>app-assure>group>cflowd# tcp-performance application "web_traffic"
*A:ALA-4>config>app-assure>group>cflowd# tcp-performance application "citrix" flow-rate2
*A:ALA-4>config>app-assure>group>cflowd# tcp-performance no shutdown
*A:ALA-4>config>app-assure>group>cflowd# rtp-performance application "voip_traffic"
*A:ALA-4>config>app-assure>group>cflowd# rtp-performance no shutdown
*A:ALA-4>config>app-assure>group>cflowd# info
-----
volume
    no shutdown
exit
rtp-performance
    no shutdown
    application "voip_traffic"
exit
tcp-performance
    no shutdown
    application "web_traffic"
    application "citrix" flow-rate2
exit
-----
*A:ALA-4>config>app-assure>group>cflowd#

```

Because no template selection is made in the above example for any of the configured cflowd templates (volume, RTP/TCP performance), the legacy template is used.

Alternatively, if the operator wants to have specific fields within, for example, volume templates, then the dynamic option must be selected under:

```

config>application-assurance>group isa-aa-group-id cflowd volume
    template + Configure volume template
    [no] dynamic-fields + Configure the list of dynamic fields

```

```
[no] field-selection <field-selection> Configure the field selection : legacy|
dynamic
```

The following example displays a volume template configured to use dynamic field selection:

```
*A:Dut-C>config>app-assure>group>cflowd>volume# template
*A:Dut-C>config>app-assure>group>cflowd>volume>template# field-selection dynamic
*A:Dut-C>config>app-assure>group>cflowd>volume>template# dynamic-fields
*A:Dut-C>config>app-assure>group>cflowd>volume>template>dynamic-fields# field "hostName"
*A:Dut-C>config>app-assure>group>cflowd>volume>template>dynamic-fields# field "aaApp"
*A:Dut-C>config>app-assure>group>cflowd>volume>template>dynamic-fields# field "aaAppGrp"
*A:Dut-C>config>app-assure>group>cflowd>volume>template>dynamic-fields# exit
*A:Dut-C>config>app-assure>group>cflowd>volume>template# info
-----
                field-selection dynamic
                dynamic-fields
                shutdown
                field "aaApp"
                field "aaAppGrp"
                field "hostName"
                exit
-----
*A:Dut-C>config>app-assure>group>cflowd>volume>template#
```

4 IP tunnels

4.1 IP tunnels overview

This section discusses IP Security (IPsec), GRE tunneling, and IP-IP tunneling features supported by the MS-ISA2/ESA-VM (ISA is used to refer to any of these hardware). In these applications, the ISA functions as a resource module for the system, providing encapsulation and (for IPsec) encryption functions. The IPsec encryption functions provided by the ISA are applicable for many applications including mobile backhaul, encrypted SDPs, video wholesale, site-to-site encrypted tunnel, and remote access VPN concentration.

Figure 39: 7750 SR IPsec implementation architecture shows an example of an IPsec deployment, and the way this would be supported inside a 7750 SR. GRE and IP-IP tunnel deployments are very similar. IP tunnels have two flavors GRE/IP-IP, in all but a few area the information for IP Tunnels applies to both types.

Figure 39: 7750 SR IPsec implementation architecture

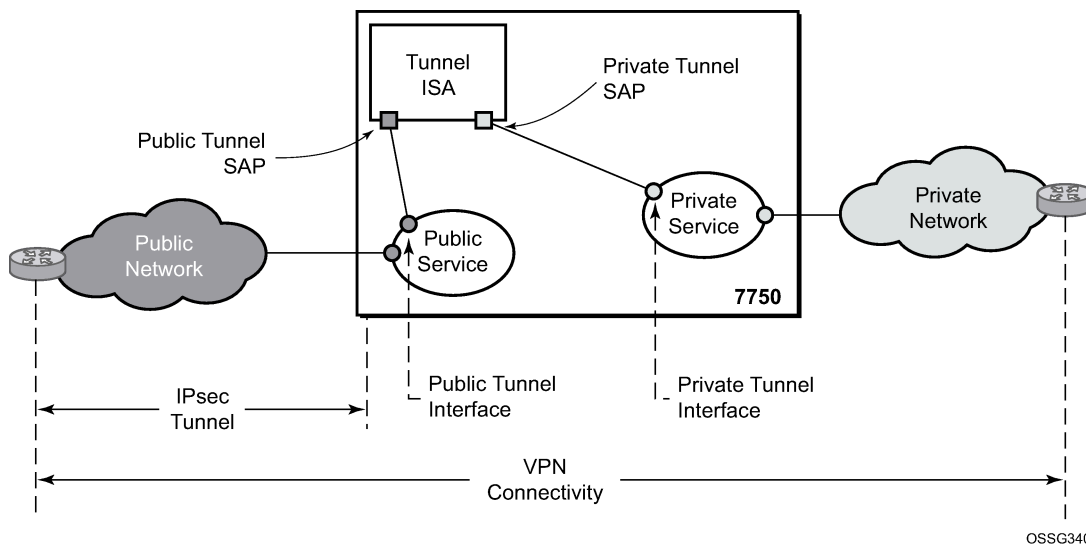


Figure 39: 7750 SR IPsec implementation architecture, the public network is typically an "insecure network" (for example, the public Internet) over which packets belonging to the private network in the diagram cannot be transmitted natively. Inside the 7750 SR, a public service instance (IES or VPRN) connects to the public network and a private service instance (typically a VPRN) connects to the private network.

The public and private services are typically two different services, and the ISA is the only "bridge" between the two. Traffic from the public network may need to be authenticated and encrypted inside an IPsec tunnel to reach the private network. In this way, the authenticity/confidentiality/integrity of accessing the private network can be enforced. If authentication and confidentiality are not important then access to the private network may alternatively be provided through GRE or IP-IP tunnels.

The ISA provides a variety of encryption features required to establish bidirectional IPsec tunnels including:

control plane

- manual keying
- dynamic keying: IKEv1/v2
- IKEv1 mode: main and aggressive
- authentication: Pre-Shared-Key /xauth with RADIUS support/X.509v3 Certificate/EAP
- Perfect Forward Secrecy (PFS)
- DPD
- NAT-Traversal
- security policy
- DH-Group: 1/2/5/14/15/19/20/21

data plane

- ESP (with authentication) tunnel mode
- authentication algorithm: MD5/SHA1/SHA256/SHA384/SHA512/AES-XCBC
- encryption algorithm: DES/3DES/AES128/AES192/AES256/AES-GCM128/AES-GCM192/AES-GCM256/AES-GMAC128/AES-GMAC192/AES-GMAC256
- anti-replay protection
- N:M IPsec ISA card redundancy

SR OS uses a configured authentication algorithm for the Pseudorandom Function (PRF). IPsec features are supported on the 7750 SR, the 7450 ESS, and VSR.

There are two types of tunnel interfaces and SAPs. See [Table 29: Tunnel interfaces and SAPs](#) for more information.

Table 29: Tunnel interfaces and SAPs

Tunnel interface/SAP	Association/configuration
Public tunnel interface	configured in the public service; outgoing tunnel packets have a source IP address in this subnet
Public tunnel SAP	associated with the public tunnel interface; a logical access point to the ISA card in the public service
Private tunnel interface	configured in the private service; can be used to define the subnet for remote access IPsec clients
Private tunnel SAP	associated with the private tunnel interface, a logical access point to the ISA card in the private service

Traffic flows to and through the ISA card as follows:

• **upstream direction**

The encapsulated (and possibly encrypted) traffic is forwarded to a public tunnel interface if its destination address matches the local or gateway address of an IPsec tunnel or the source address of

a GRE or IP-IP tunnel. Inside the ISA card, encrypted traffic is decrypted, the tunnel header is removed, the payload IP packet is delivered to the private service, and from there, the traffic is forwarded again based on the destination address of the payload IP packet.

- **downstream direction**

Unencapsulated/clear traffic belonging to the private service is forwarded into the tunnel by matching a route with the IPsec/GRE/IP-IP tunnel as next-hop. The route can be configured statically, learned by running OSPF on the private tunnel interface (GRE tunnels only), learned by running BGP over the tunnel (IPsec and GRE tunnels only), or learned dynamically during IKE negotiation (IPsec only). After clear traffic is forwarded to the ISA card, it is encrypted if required, encapsulated per the tunnel type, delivered to the public service, and from there, the traffic is forwarded again based on the destination address of the tunnel header.

4.1.1 Tunnel ISAs

A tunnel-group is a collection of MS-ISA2s (mda-type **isa2-tunnel**) or ESA-VM (vm-type tunnel) configured to handle the termination of one or more IPsec, GRE or IP-IP tunnels. Two example tunnel-group configurations are shown below:

```
config isa
  tunnel-group 1 create
    primary 1/1
    backup 2/1
    no shutdown
  exit

config isa
  tunnel-group 2 create
    multi-active
    mda 3/1
    mda 3/2
    no shutdown

config isa
  tunnel-group 3 create
    multi-active
    esa-vm 3/1
    esa-vm 4/1
    no shutdown
```

A GRE, IP-IP, or IPsec tunnel belongs to only one tunnel group. There are two types of tunnel groups:

- **single-active tunnel-group**

A single-active tunnel-group can have one tunnel-ISA designated as primary and optionally one other tunnel-ISA designated as backup. If the primary ISA fails the affected failed tunnels are re-established on the backup (which is effectively a cold standby) if it is not already in use as a backup for another tunnel-group.

- **multi-active tunnel-group**

A multi-active tunnel-group can have multiple tunnel-ISAs designated as primary. Only one ISA is supported on VSR.

A multi-active tunnel is the recommended tunnel-group type. Certain features like MC-IPsec are only supported with a multi-active tunnel-group.

The ESA-VM is only supported in a multi-active tunnel-group.

Note that the ESA-VM and ISA/ISA2 cannot coexist in the same tunnel-group.

The **show isa tunnel-group** command allows the operator to view information about all configured tunnel groups. This command displays the following information for each tunnel-group: group ID, primary tunnel-ISAs, backup tunnel-ISAs, active tunnel-ISAs, admin state and oper state.

There are three thresholds that are used to monitor memory usage in a tunnel ISA:

- **max-threshold**

When the memory usage of an ISA exceeds this threshold, any new IKE states are rejected.

- **high-watermark**

When the memory usage of an ISA exceed this threshold, a trap is generated.

- **low-watermark**

When the memory usage of an ISA fall below this threshold, a clear trap is generated.

These three thresholds are fixed, not configurable.

A tunnel-group has an **isa-scale-mode**, which defines the maximum number of all tunnels (all types combined) which can be established on each ISA of the tunnel group. This is currently fixed at 32,000 tunnels per ISA. This value is different on VSR and vSIM, see the corresponding User Guides for details.

4.1.1.1 Public tunnel SAPs

A VPRN or IES service (the delivery service) must have at least one IP interface associated with a public tunnel SAP to receive and process the following types of packets associated with GRE, IP-IP and IPsec tunnels:

- GRE (IP protocol 47)
- IP-IP (IP protocol 4)
- IPsec ESP (IP protocol 50)
- IKE (UDP)

The public tunnel SAP type has the format **tunnel-tunnel-group.public:index**, as shown in the following CLI example.

```
*A:Dut-C>config>service# info
-----
customer 1 create
  description "Default customer"
exit
ies 1 customer 1 create
  interface "public" create
    address 192.168.12.1/24
    tos-marking-state untrusted
    sap tunnel-1.public:200 create
  exit
exit
no shutdown
exit
vprn 2 customer 1 create
  route-distinguisher 10.1.1.1:65007
  interface "greTunnel" tunnel create
    address 10.0.0.1/24
```

```

        dhcp
        no shutdown
    exit
    sap tunnel-1.private:210 create
        ip-tunnel "toCel" create
            dest-ip 10.0.0.2
            gre-header
            source 192.168.12.100
            remote-ip 10.251.12.2
            backup-remote-ip 10.251.12.22
            delivery-service 1
            no shutdown
        exit
    exit
    exit
    no shutdown
exit
-----
*A:Dut-C>config>service#

```

4.1.1.2 Private tunnel SAPs

The private service must have an IP interface to a GRE, IP-IP, or IPsec tunnel to forward IP packets into the tunnel, causing them to be encapsulated (and possibly encrypted) per the tunnel configuration and to receive IP packets from the tunnel after the encapsulation has been removed (and decryption). That IP interface is associated with a private tunnel SAP.

The private tunnel SAP has the format **tunnel-tunnel-group.private:index**, as shown in the following CLI example where a GRE tunnel is configured under the SAP.

```

*A:Dut-A# show ip tunnel
=====
IP Tunnels
=====
TunnelName                SapId                SvcId    Admn
Local Address             OperRemoteAddress   DlvrySvcId Oper
-----
tun-1-gre-tunnel         tunnel-1.private:1   201      Up
192.168.1.2              192.168.3.2         1201     Up
-----
IP Tunnels: 1
=====

```

4.1.1.3 IP interface configuration

In the configuration example of the previous section the IP address 10.0.0.1 is the address of the GRE tunnel endpoint from the perspective of payload IP packets. This address belongs to the address space of the VPRN 1 service and is not exposed to the public IP network carrying the GRE encapsulated packets. An IP interface associated with a private tunnel SAP does not support unnumbered operation.

It is possible to configure the IP MTU (M) of a private tunnel SAP interface. This sets the maximum payload IP packet size (including IP header) that can be sent into the tunnel – for example, it applies to the packet size before the tunnel encapsulation is added. When a payload IPv4 packet that needs to be forwarded into the tunnel is larger than M bytes the payload packet is IP fragmented (before tunnel encapsulation)

if the DF bit is clear, otherwise the packet is discarded. When a payload IPv6 packet that needs to be forwarded into the tunnel is larger than M bytes the packet is discarded if its size is less than 1280 bytes otherwise it is forwarded and encapsulated intact.

4.1.1.4 GRE and IP-IP tunnel configuration

To bind an IP/GRE or IP-IP tunnel to a private tunnel SAP, the **ip-tunnel** command should be added under the SAP. To configure the tunnel as an IP/GRE tunnel, the **gre-header** command must be present in the configuration of the **ip-tunnel**. To configure the tunnel as an IP-IP tunnel, the **ip-tunnel** configuration should have the **no gre-header** command. When configuring a GRE or IP-IP tunnel, the **dest-ip** command specifies an IPv4 or IPv6 address (private) of the remote tunnel endpoint. A tunnel can have up to 16 dest-ip addresses. If any of the dest-ip addresses are not contained by a subnet of the local private endpoint, then the tunnel does not come up. In the CLI sub-tree under **ip-tunnel**, there are commands to configure the following:

- **source address of the GRE or IP-IP tunnel**

This is the source IPv4 address of GRE or IP-IP encapsulated packets sent by the delivery service. It must be an address in the subnet of the associated public tunnel SAP interface.

- **remote IP address**

If this address is reachable in the delivery service (there is a route), then this is the destination IPv4 address of GRE or IP-IP encapsulated packets sent by the delivery service.

- **backup remote IP address**

If the remote IP address of the tunnel is not reachable then this is the destination IPv4 address of GRE or IP-IP encapsulated packets sent by the delivery service.

- **delivery service**

This is the ID or name of the IES or VPRN service where GRE or IP-IP encapsulated packets are injected and terminated. The delivery service can be the same service where the private tunnel SAP interface resides.

- **DSCP marking in the outer IP header of GRE encapsulated packets**

If this is not configured, then the default is to copy the DSCP from the inner IP header to the outer IP header.

A private tunnel SAP can have only one ip-tunnel sub-object (one GRE or IP-IP tunnel per SAP).

The **show ip tunnel** command displays information about a specific IP tunnel or all configured IP tunnels. The following information is provided for each tunnel:

- service ID that owns the tunnel
- private tunnel SAP that owns the tunnel
- tunnel name, source address
- remote IP address
- backup remote IP address
- local (private) address
- destination (private) address
- delivery service
- dscp

- admin state
- oper state
- type (GRE or IP-IP)

The following is an example of the output of the **show ip tunnel tunnel-name** command.

```
A:config>service>vprn>if>sap>ip-tunnel# show ip tunnel "ipv6-gre"
=====
IP Tunnel Configuration Detail
=====
Service Id       : 1                Sap Id          : tunnel-1.private:1
Tunnel Name      : ipv6-gre
Description      : None
GRE Header       : Yes              Delivery Service : 2
GRE Keys Set     : False
GRE Send Key     : N/A              GRE Receive Key  : N/A
Admin State      : Up                Oper State       : Up
Source Address   : 2001:db8::1:2:3:4
Remote Address   : 3ffe:1::2
Backup Address   : (Not Specified)
Oper Remote Addr : 3ffe:1::2
DSCP             : ef
Reassembly       : inherit
Clear DF Bit     : false            IP MTU           : max
Encap IP MTU    : 1400
Pkt Too Big     : true
Pkt Too Big Numb* : 100             Pkt Too Big Intvl: 10 secs
Oper Flags       : None
Last Oper Changed: 02/09/2015 15:22:38
Host MDA         : 1/2

-----
Target Address Table
-----
Destination IP          IP Resolved Status
-----
172.16.1.2              Yes
2001:db8::2             Yes
-----

=====
IP Tunnel Statistics: ipv6-gre
=====
Errors Rx       : 0                Errors Tx       : 0
Pkts Rx         : 0                Pkts Tx        : 0
Bytes Rx        : 0                Bytes Tx        : 0
Key Ignored Rx  : 0                Too Big Tx     : 0
Seq Ignored Rx  : 0
Vers Unsup. Rx  : 0
Invalid Chksum Rx: 0
Key Mismatch Rx : 0

=====

Fragmentation Statistics
=====
Encapsulation Overhead : 44
Pre-Encapsulation
  Fragmentation Count   : 0
  Last Fragmented Packet Size : 0
Post-Encapsulation
  Fragmentation Count   : 0
```

```

Last Fragmented Packet Size      : 0
=====
=====

```

4.1.1.5 IP fragmentation and reassembly for IP tunnels

An IPsec, GRE or IP-IP tunnel packet that is larger than the IP MTU of some interface in the public network must either be discarded (if the Do Not Fragment (DF) bit is set in the outer IP header) or fragmented. If the tunnel packet is fragmented, then it is up to the destination tunnel endpoint to reassemble the tunnel packet from its fragments. Starting in Release 10, IP reassembly can be enabled for all the IPsec, GRE, and IP-IP tunnels belonging to a tunnel-group. For IP-IP and GRE tunnels, the reassembly option is also configurable on a per-tunnel basis so that some tunnels in the tunnel-group can have reassembly enabled, and others can have the extra processing disabled. When reassembly is disabled for a tunnel, all received fragments belonging to the tunnel are dropped.

To avoid public network fragmentation of IPsec, GRE, or IP-IP packets belonging to a particular tunnel, one possible strategy is to fragment IPv4 payload packets larger than a specified size M at entry into the tunnel (before encapsulation and encryption if applicable). The size M is configurable using the **ip-mtu** command under the **ip-tunnel** or **ipsec-tunnel/tunnel-template** configuration.

If the payload IPv4 packets are all M bytes or less in length then it is guaranteed that all resulting tunnel packets are less than M+N bytes in length, if N is the maximum overhead added by the tunneling protocol. If M+N is less than the smallest interface IP MTU in the public network, fragmentation is avoided. In some cases, some of the IPv4 payload packets entering a tunnel may have their DF bit set. And if needed, the SR OS supports the option (also configurable on a per-tunnel basis) to clear the DF bit in these packets so that they can be fragmented.

The system allows users to configure an **encapsulated-ip-mtu** for a tunnel under an **ip-tunnel** or **ipsec-tunnel/tunnel-template** configuration. This represents the maximum size of the encapsulated tunnel packet. After encapsulation, if the IPv4 or IPv6 tunnel packet size exceeds the configured **encapsulated-ip-mtu**, then the system fragments the packet against the **encapsulated-ip-mtu**.

The following is a description of system behavior about fragmentation:

- **private side**

If the size, before encapsulation, of the IPv4 or IPv6 packet entering the tunnel is larger than the ip-mtu configured under ip-tunnel or ipsec-tunnel/tunnel template:

- **IPv4 payload packet**

If the DF bit is not set in the packet or if the **clear-df-bit** command is configured, then the system fragments the packet against the ip-mtu configured under ip-tunnel or ipsec-tunnel/tunnel-template.

Otherwise, the system drops the packet and sends back an ICMP error Fragmentation required and DF flag set, with the suggested MTU set as the ip-mtu.

- **IPv6 payload packet**

If the packet size >1280 bytes, the system drops the packet and sends back an ICMPv6 Packet Too Big (PTB) message with the suggested MTU set as the ip-mtu.

If the packet size <=1280 bytes, the system forwards the packet into the tunnel.

- **public side**

This applies to both ESP and IKE packets, IPv4 and IPv6.

If the ESP/IKE packet is larger than the encapsulated-ip-mtu, then the system fragments the packet against the encapsulated-ip-mtu, however when the IPv6 ESP/IKE packet is smaller than 1280 bytes, the system does not fragment it, even if it is larger than the encapsulated-ip-mtu.

4.1.1.6 TCP MSS adjustment

The system supports the Transmission Control Protocol (TCP) Maximum Segment Size (MSS) adjustment feature for the following types of tunnels on the ISA:

- IPsec
- IPinIP/GRE
- L2TPv3 (data packet only)

The intent of TCP MSS adjustment is to avoid IP-level fragmentation for TCP traffic encapsulated in a tunnel by updating the MSS option value in the TCP SYN packet with an appropriate value. This feature is useful when there is tunnel encapsulation that is not known by a TCP host, and the extra tunnel encapsulation overhead may cause IP-level fragmentation.

The system supports TCP MSS adjustment on both the public and private sides.

On the public side, when the ISA receives a tunnel packet (such as ESP), after decryption or decapsulation, if the payload packet is a TCP SYN packet, then the ISA replaces the MSS option with a configured value if the configured MSS value is smaller than the received MSS value or when there is no MSS option:

- If **public-tcp-mss-adjust auto** is configured, then:

new MSS value = public_side_MTU – tunnel_overhead – TCP fixed header – IP fixed header

where:

- public_side_MTU = **encapsulated-ip-mtu**

If **encapsulated-ip-mtu** is not configured, which means there is no post-encap fragmentation on ISA, then TCP MSS adjust is disabled.

- TCP fixed header = 20
- IP fixed header = 20 (IPv4) or 40 (IPv6)

- If a specific MSS value such as **public-tcp-mss-adjust new_mss_value** is configured, then the new MSS value sets to the *new_mss_value*.



Note:

- The **public-tcp-mss-adjust auto** command only applies to IPsec and IPinIP/GRE tunnels.
- For an IPsec tunnel, the tunnel_overhead is the maximum overhead of the corresponding CHILD_SA.
- For an IPinIP tunnel, the tunnel_overhead is 0.
- For a GRE tunnel, the tunnel_overhead is length of GRE header.

The private side is similar to the public side. The system processes the received TCP SYN packet on the private side if the TCP MSS adjust is enabled. However, there is no **auto** parameter for **private-tcp-mss-adjust** command.

4.1.1.7 MTU propagation

MTU propagation is an optional feature that allows the system to listen for fragmentation-related ICMP error message received from the public side of the tunnel. These error messages include:

- ICMP Destination Unreachable message "fragmentation needed and DF set" (type 3, code 4)
- ICMPv6 Packet Too Big message (type 2)

The suggested MTU value in the ICMP message is used to derive two MTU values:

- Temporary public MTU (TMTU) are determined as follows:
 - The TMTU starts with a configured **encapsulated-ip-mtu octets** value.
 - If the received MTU is less than 1280 and it is from an ICMPv6 packet, the received value is ignored.
 - If the received MTU is less than 512 and it is from an ICMP packet, the received value is ignored.
 - If the received MTU is greater than or equal to the configured **encapsulated-ip-mtu octets** value, the received value is ignored.
 - If the received MTU is greater than or equal to the current TMTU, the received value is ignored.
 - If the received MTU is less than the current TMTU, it replaces the current TMTU.
 - To prevent attack and rapid change, there is a damp time of 60 seconds after a new TMTU value is set. Within that time frame, all received MTU values are ignored.
 - TMTU has a lifetime timer (configurable with an aging interval). When the lifetime timer expires, the TMTU's value is reset to the **encapsulated-ip-mtu octets** value. The lifetime timer resets whenever a new TMTU value is set.
 - TMTU is a per tunnel value.
- Temporary private MTU (TPMTU) equals TMTU – Tunnel_Encap_Overhead.
 - TPMTU is a per CHILD_SA value.
 - Tunnel_Encap_Overhead is a fixed value for a non-IPsec tunnel-per-tunnel type. For an IPsec tunnel, its value is the maximum overhead based on the **ipsec-transform transform-id** value used by the CHILD_SA.

TMTU and TPMTU are used in the following cases:

- TPMTU is used for fragmenting IP packets received on the private side instead of the configured IP MTU.
- IKEv2 message fragmentation uses TMTU instead of the configured **encapsulated-ip-mtu**.
- IKE IP packet fragmentation uses TMTU instead of the configured **encapsulated-ip-mtu**.
- To derive the TCP MSS value for the TCP MSS adjustment, instead of configured **encapsulated-ip-mtu**.
- ESP packet fragmentation (post-encapsulation fragmentation) does not use TMTU; it only uses the configured **encapsulated-ip-mtu octets** value.

To enable this feature, configure the **propagate-pmtu-v4** and **propagate-pmtu-v6** commands under the **ip-tunnel**, **ipsec-tunnel** or **tunnel-template** contexts.

4.1.2 Operational conditions

A tunnel group that is in use cannot be deleted. In single-active mode, changes to the primary ISA are allowed only when the tunnel group is in a shutdown state. Changing the configured backup ISA (or adding a backup ISA) is allowed at any time unless the ISA is currently active for this tunnel group. When the backup ISA is active, changing the primary ISA is allowed without shutting down the tunnel group.

Changes can be made to the following:

- the mode from multi-active to single-active
- the primary ISA that is in single-active mode
- the active MDA number value that is in multi-active mode
- enabling or disabling the **ipsec-responder-only** configuration

In multi-active mode, if the active member ISA goes down, the system replaces it with a backup ISA. However, if a backup ISA is not available, the tunnel group is placed in an operationally down state. A multi-active tunnel group with MC-IPsec enabled cannot be changed into single active-mode unless it is first removed from the MC-IPsec configuration.

The public interface address can be changed at any time; however, if changed, any static tunnels that were configured to use the public interface address require a configuration changes accordingly. Otherwise, the tunnels are in an operationally down state until their configuration is corrected. The public service cannot be deleted while tunnels are associated.

A tunnel group ID or tag cannot be changed. To remove a tunnel-group instance, it must be in a shutdown state and all IPsec tunnels and IPsec gateways that terminated on the tunnel group must be removed first.

The security policy cannot be changed while an IPsec tunnel is administratively up and using the security policy.

The tunnel local gateway address, peer address, local ID, and public or private service ID parameters cannot be changed while the IPsec gateway or IPsec tunnel is administratively up.

Each IPsec gateway or IPsec tunnel has an administrative state. When the administrative state is down, tunnels cannot be set up.

Each IPsec gateway and IPsec tunnel has an operation state. The operational state can have three possible values:

- **oper-up**
All configuration and related information are valid and fully ready for tunnel setup.
- **oper-down**
Some critical configuration information is missing or not ready. Tunnels cannot be set up.
- **limited**
Not all configuration information is ready to become fully operationally up. When IPsec gateway is in a limited state, it is possible that a new tunnel cannot be established. When the IPsec tunnel is in a limited state, reconnection may fail.

When an IPsec gateway or IPsec tunnel transitions from operationally up to an operationally limited state directly as a result of a hot (non-service affecting) configuration change, established tunnels are not impacted. However, if the IPsec gateway or IPsec tunnel transitions to an operationally down state before it is operationally limited as a result of a service-affecting configuration change, then established tunnels are removed. All operational state transitions are logged.

IPsec gateways or IPsec tunnels can enter the limited state because of the following reasons, among others:

- A Certificate Authority (CA) profile in the configured trust-anchor-profile goes down after the IPsec gateway or IPsec tunnel becomes operationally up.
- An entry in a configured certificate profile goes down after the IPsec gateway or IPsec tunnel becomes operationally up.

4.1.2.1 Dynamic configuration change support for IPsec gateway

All dynamic IPsec tunnels (dynamic LAN-to-LAN tunnels and remote-access tunnels) that terminate on the same IPsec gateway share the same configuration. Use the respective commands in the following contexts to configure an IPsec gateway for an IES or VPRN service:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway
configure service vprn interface sap ipsec-gateway
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw
configure service vprn interface sap ipsec-gw
```

SR OS provides dynamic configuration change capability to modify specific IPsec gateway configurations without impacting existing tunnels.

The following IPsec gateway configurations are dynamically configurable without shutting down the IPsec gateway:

- Changing the pre-shared key, using the following commands:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway pre-shared-key
configure service vprn interface sap ipsec-gateway pre-shared-key
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw pre-shared-key
configure service vprn interface sap ipsec-gw pre-shared-key
```

- Changing the reference of the IKE policy, using the following commands:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway ike-policy
configure service vprn interface sap ipsec-gateway ike-policy
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw ike-policy
configure service vprn interface sap ipsec-gw ike-policy
```

- Changing the reference of the tunnel template, using the following commands:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway default-tunnel-template
configure service vprn interface sap ipsec-gateway default-tunnel-template
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw default-tunnel-template
configure service vprn interface sap ipsec-gw default-tunnel-template
```

- Enabling or changing reference of the RADIUS authentication policy, using the following commands:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway radius authentication-policy
configure service vprn interface sap ipsec-gateway radius authentication-policy
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw radius-authentication-policy
configure service vprn interface sap ipsec-gw radius-authentication-policy
```

- Enable or changing reference of the RADIUS accounting policy, using the following commands:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway radius accounting-policy
configure service vprn interface sap ipsec-gateway radius accounting-policy
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw radius-accounting-policy
configure service vprn interface sap ipsec-gw radius-accounting-policy
```

- Enabling, disabling, or changing reference of the TS negotiation, using the following commands:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway ts-list
configure service vprn interface sap ipsec-gateway ts-list
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw ts-negotiation
configure service vprn interface sap ipsec-gw ts-negotiation
```

- Enabling, disable, or changing reference of the client database, using the command options in the following contexts:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway client-db
configure service vprn interface sap ipsec-gateway client-db
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw client-db
```

```
configure service vprn interface sap ipsec-gw client-db
```

- Changing certificate configuration, using the commands in the followings contexts:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway cert
configure service vprn interface sap ipsec-gateway cert
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw cert
configure service vprn interface sap ipsec-gw cert
```

- Changing DHCPv4-based address assignments, using the commands in the following contexts:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway dhcp-address-assignment dhcpv4
configure service vprn interface sap ipsec-gateway dhcp-address-assignment dhcpv4
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw dhcp
configure service vprn interface sap ipsec-gw dhcp
```

- Changing DHCPv6-based address assignments, using the commands in the following contexts:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway dhcp-address-assignment dhcpv6
configure service vprn interface sap ipsec-gateway dhcp-address-assignment dhcpv6
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw dhcp6
configure service vprn interface sap ipsec-gw dhcp6
```

- Change local address assignment configuration, using the commands in the following contexts:

- **MD-CLI**

```
configure service ies interface sap ipsec-gateway local address-assignment
configure service vprn interface sap ipsec-gateway local address-assignment
```

- **classic CLI**

```
configure service ies interface sap ipsec-gw local-address-assignment
configure service vprn interface sap ipsec-gw local-address-assignment
```

Existing tunnels are not impacted by dynamic configuration changes. The system uses new configurations for new tunnel negotiations. The system continues to use previous configurations that created the tunnels for ongoing operations (such as rekeying) of the existing tunnel.

4.1.3 QoS interactions

The ISA can interact with the queuing functions on the IOM through the ingress and egress QoS provisioning in the IES or IP VPN service where the IPsec session is bound. Multiple IPsec sessions can be assigned to a single IES or VPRN service. In this case, QoS defined at the IES or VPRN service level is applied to the aggregate traffic coming out of, or going into, the set of sessions assigned to that service.

To keep marking relevant in the overall networking design, the following traffic-class processing is supported:

- In the encapsulating direction (private to public), the system copies the traffic class of the payload IP packet header to the outer tunnel IP packet header.
- In the decapsulating direction (public to private), the system can optionally copy the traffic class from the outer tunnel IP packet header to the payload IP packet header using the **copy-traffic-class-upon-decapsulation** command for the template, service, or router IPsec tunnel configuration.

For the tunnel-group ESA VM, if a SAP egress QoS policy is needed on a public or private tunnel SAP, the CIR of all queues configured in the policy should be zero (non-zero CIR is not supported).

4.1.4 OAM interactions

The ISA is IP-addressed by an operator-controlled IP on the public side. That IP address can be used in Ping and Traceroute commands and the ISA can either respond or forward the packets to the CPM.

For static LAN-to-LAN tunnels, in multi-active mode, ping requests to public tunnel addresses would not be answered if the source address is different from the remote address of the static tunnel.

The private side IP address is visible. The status of the interfaces and the tunnels can be viewed using **show** commands.

Traffic that ingresses or egresses an IES or VPRN service associated with specific IPsec tunnels can be mirrored like other traffic.

Mirroring is allowed per interface (public) or IPsec interface (private) side. A filter mirror is allowed for more specific mirroring.

4.1.5 Redundancy

In single-active mode, every tunnel group can be configured with primary and backup ISAs. An ISA can be used as a backup for multiple IPsec groups. The ISAs are cold standby such that upon failure of the primary the standby resumes operation after the tunnels re-negotiate state. While the backup ISA can be shared by multiple tunnel groups only one tunnel group can fail to a single ISA at one time (no double failure support).

In multi-active mode, the active-mda-number value determines the number of ISA MDAs that are active for this tunnel group, and tunnels are spread across all active ISA MDAs. Additional ISA MDAs in this tunnel group are in cold standby.

IPsec also supports dead peer detection (DPD).

BFD can be configured on the private tunnel interfaces associated with GRE tunnels and used by the OSPF, BGP or static routing that is configured inside the tunnel.

SR OS also supports multi-chassis IPsec redundancy, which provides 1:1 stateful protection against ISA failure or chassis failure.

4.1.6 Statistics collection

Input and output octets and packets per service queue are used for billing end customers who are on a metered service plan. Because multiple tunnels can be configured per interface, the statistics can include multiple tunnels. These can be viewed in the CLI and SNMP.

Reporting (syslog, traps) for authentication failures and other IPsec errors are supported, including errors during IKE processing for session setup and errors during encryption or decryption.

A session log indicates the sort of SA setup when there is a possible negotiation. This includes the setup time, teardown time, and negotiated parameters (such as encryption algorithm) as well as identifying the service a particular session is mapped to, and the user associated with the session.

4.1.7 Security

The ISA module provides security utilities for IPsec-related service entities that are assigned to interfaces and SAPs. These entities (such as card, MS-ISA2 module, and IES or VPRN services) must be enabled in order for the security services to process. The module only listens to requests for security services from configured remote endpoints. In the case of a VPN concentrator application, these remote endpoints could come from anywhere on the Internet. In the cases where a point-to-point tunnel is configured, the module listens only to messages from that endpoint.

4.1.7.1 GRE tunnel multicast support

GRE tunnels support unicast and multicast IP packets as payload. From a multicast prospective, a GRE tunnel IP interface (associated with a private tunnel SAP) can be configured as an IGMP interface or as a PIM interface; MLD is not supported. The following multicast features are supported:

- IGMP versions 1, 2 and 3
- IGMP import policies
- IGMP host tracking
- static IGMP membership
- configurable IGMP timers
- IGMP SSM translation
- multicast CAC
- per-interface, per-protocol (IGMP/PIM) multicast group limits
- MVPN support (draft-rosen)
- MVPN support (BGP-MPLS)
- PIM-SM and SSM operation
- PIM BFD support
- configurable PIM timers
- configurable PIM priority
- PIM tracking support
- PIM ECMP (bandwidth or hash-based)

- static multicast route

4.1.7.2 IPv6 over IPv4 GRE tunnel

IPv6 payload packets can be delivered over an IPv4 GRE tunnel. In this scenario the two endpoints of the GRE tunnel have IPv4 addresses and the VPRN or IES SAP interface to the tunnel is an IPv6 only or dual-stack IPv4/IPv6 interface. IPv6 over IPv4 GRE tunneling allows IPv6 islands to be connected over an IPv4 only transport infrastructure.

To configure a tunnel to carry IPv6 payload, the tunnel must be configured with at least one **dest-ip** that contains an IPv6 address (global unicast and/or link local). A tunnel can have up to 16 **dest-ip** addresses (IPv4 and IPv6 together). For a tunnel to come operationally up all the **dest-ip** addresses must be part of locally configured subnets (associated with the private tunnel interface).

To forward IPv6 traffic through a tunnel supporting IPv6 payload, a dynamic routing protocol (such as BGP or OSPFv3) can be configured to run inside the tunnel (by associating the protocol with the private tunnel interface) or else an IPv6 static route next-hop equal to a **dest-ip** of the tunnel can be used.

IPv6 payload packets larger than 1280 bytes (the minimum IPv6 MTU) and also larger than the configured **ip-mtu** value of the tunnel are always discarded. If the **icmp6-generation** and **packet-too-big** commands are configured under the tunnel, then ICMPv6 Packet-Too-Big messages are generated and sent back to the originating host when discards occur because of the private side IP MTU being exceeded.

4.1.8 IKEv2

IKEv2, defined in RFC 4306, *Internet Key Exchange (IKEv2) Protocol*, is the second version of the Internet Key Exchange Protocol. The main driver of IKEv2 is to simplify and optimize IKEv1. An IKE_SA and a CHILD_SA can be created with only four IKEv2 message exchanges. IKEv2 is supported with the following features:

- static LAN-to-LAN tunnel
- dynamic LAN-to-LAN tunnel
- remote-access tunnel
- pre-shared-key authentication, certificate authentication, EAP (remote-access tunnel only)
- liveness check
- IKE_SA rekey
- CHILD_SA rekey (full Traffic-Selector support including protocol and port range)
- extended ESP sequence number

4.1.8.1 IKEv2 traffic selector and TS-list

The SR OS IKEv2 implementation supports the following traffic selectors:

- IPv4/IPv6 address range
- IP protocol ID
- protocol port range

Port range (including OPAQUE ports) is supported for the following protocols:

- TCP
- UDP
- SCTP
- ICMP
- ICMPv6
- MIPv6

With ICMP and ICMPv6, the system treats the most significant 8 bits of the IKEv2 TS port value as the ICMP message type and the least significant 8 bits as ICMP code.

With MIPv6, the system treats the most significant 8 bits of the IKEv2 TS port value as the mobility header type.

With ICMP, ICMPv6, and MIPv6, the port value in TS_i is the value that the tunnel initiator can send, and the port value in TS_r is the value that the tunnel responder can send.

The SR OS supports OPAQUE as a TS port selector. An OPAQUE port means that the corresponding CHILD_SA only accepts packets that are supposed to have port information but do not, such as when a packet is a non-initial fragment.

The system allows users to configure a TS-list for each IPsec gateway, applied to both IKEv2 remote access tunnels and dynamic LAN-to-LAN tunnels. Each TS-list contains a local part and a remote part, with

each part containing up to 32 entries. Each entry can contain address ranges or subnets, protocols, and port range configurations.

The local part of the TS-list represents the traffic selector for the local system, while the remote part is for the remote peer. If a TS-list is applied on an IPsec gateway, and the system is the tunnel responder, then the local part is TSr and the remote part is TSi.

Combinations of address range, protocol, and port range are not allowed to overlap between entries in the same TS-list.

The system performs traffic selector narrowing as follows.

1. For each TS in the received TSi/TSr, independent address, protocol, and port narrowing is performed. The resulting TS-set is the combination of the address, protocol, and range intersections.
2. The collected TS-set is used as the TSi/TSr.

For a remote access tunnel, TSi narrowing results in an intersection between the following three TSis:

- the TSi received from the client
- the remote part configuration of the TS-list
- a generated TS based on the assigned internal address
 - address (the assigned internal address)
 - protocol (any)
 - port range (any)

The following is an example of a dynamic LAN-to-LAN tunnel.

The configured TS-list local part is as follows:

- Entry 1: 10.10.1.0 → 10.10.1.20, udp, port 100 → 200
- Entry 2: 10.20.1.0 → 10.20.1.20, udp, port 300 → 400

The peer proposes the following TSr:

- Entry 1: 10.10.1.1 → 10.10.1.5, udp, port 110 → 150
- Entry 2: 10.10.1.6 → 10.10.1.10, udp, port 180 → 210
- Entry 3: 10.10.1.15 → 10.10.1.28, udp, port 120 → 160
- Entry 4: 10.20.1.15 → 10.20.1.28, tcp, port 250 → 450

The intersections for the proposed entries are as follows:

- Entry 1: 10.10.1.1 → 10.10.1.5, udp, port 110 → 150
- Entry 2: 10.10.1.6 → 10.10.1.10, udp, port 180 → 200
- Entry 3: 10.10.1.15 → 10.10.1.20, udp, port 120 → 160
- Entry 4: 10.20.1.15 → 10.20.1.20, tcp, port 300 → 400

The resulting TSr system return would be:

- 10.10.1.1 → 10.10.1.5, udp, port 110 → 150
- 10.10.1.6 → 10.10.1.10, udp, port 180 → 200
- 10.10.1.15 → 10.10.1.20, udp, port 120 → 160
- 20.20.1.15 → 20.20.1.20, tcp, port 300 → 400

If more than 32 entries are returned, the system rejects the ts-negotiation and returns TS_UNACCEPTABLE to the peer.

For dynamic LAN-to-LAN tunnels, the system can automatically create a reverse route in a private VRF to route clear traffic into the tunnel. The reverse route is created based on the address range part of the narrowed TSi of the CHILD_SA. If there are multiple TSs in the TSi that have overlapping address ranges, the system creates one or more minimal subnet routes that can cover all address ranges in the TSi. If the auto-created reverse route overlaps with an existing reverse route that points to the same tunnel, the system chooses the route with the larger subnet. If the existing route points to a different tunnel, then CHILD_SA creation fails.

For RADIUS authentication, such as **psk-radius**, **cert-radius**, or EAP, the RADIUS server can optionally return a TS-list name via the VSA Alc-IPsec-Ts-Override in the access-accept message, which overrides the TS-list name configured via the CLI.

In the event of a CHILD_SA rekey, if the system is a rekey initiator, it sends the current in-use TS to the peer and expect the peer to return the same TS. If the system is a rekey responder, the system does the same narrowing as was done during CHILD_SA creation.

Configuration of a TS-list can be changed without shutting down the IPsec gateway, although the new TS-list only applies to the subsequent rekey or to the new CHILD_SA creation, and does not affect established CHILD_SAs.

4.1.8.2 IKEv2 fragmentation

In some cases, an IKEv2 message can be large, like an IKE_AUTH message with certificate payload. This is likely to cause the IKEv2 packet to be fragmented into a few smaller IP packets. However, in some deployments, there could be devices or network policing, rate limiting or even dropping UDP fragments. In these cases, the SR OS supports fragmenting IKEv2 messages on the protocol level, as specified in RFC 7383, Internet Key Exchange Protocol Version 2 (IKEv2) Message Fragmentation.

This feature is enabled by configuring the **ikev2-fragment** command in the ike-policy context with an MTU. The specified MTU is the maximum size of IKEv2 packet.

The system only enables IKEv2 fragmentation for a specific tunnel when the **ikev2-fragment** is configured and the peer also announces its support via sending a IKEV2_FRAGMENTATION_SUPPORTED notification.

4.1.9 SHA2 support

According to RFC 4868, *Using HMAC-SHA-256, HMAC-SHA-384, and HMAC-SHA-512 with IPsec*, the following SHA2 variants are supported for authentication or pseudo-random functions:

Use HMAC-SHA-256+ algorithms for data origin authentication and integrity verification in IKEv1/2, ESP:

- AUTH_HMAC_SHA2_256_128
- AUTH_HMAC_SHA2_384_192
- AUTH_HMAC_SHA2_512_256

For use of HMAC-SHA-256+ as a PRF in IKEv1/2:

- PRF_HMAC_SHA2_256
- PRF_HMAC_SHA2_384
- PRF_HMAC_SHA2_512

4.1.10 IPsec client lockout

An optional lockout mechanism can be enabled to block malicious clients and prevent them from using invalid credentials to consume system resources, as well as to prevent malicious users from guessing credentials such as a pre-shared key. This mechanism can be enabled by using the **lockout** command.

If the number of failed authentication attempts from a particular IPsec client exceeds a configured threshold during a specified time interval, the client is blocked for a configurable period of time. If a client is blocked, the system drops all IKE packets from the source IP address and port.

The following authentication failures are counted as failed authentication attempts:

- **IKEv1**
 - **psk**: failed to verify the HASH_I payload in main mode
 - **plain-psk-xauth**:
 - failed to verify the HASH_I payload in main mode
 - RADIUS access-reject received
- **IKEv2**
 - **psk**: failed to verify the AUTH payload in the auth-request packet
 - **psk-radius**:
 - failed to verify the AUTH payload in the auth-request packet
 - RADIUS access-reject received
 - **cert**:
 - failed to verify the AUTH payload in the auth-request packet
 - failed to verify the peer's certification to configured trust-anchors
 - **cert-radius**:
 - failed to verify the AUTH payload in the auth-request packet
 - failed to verify the peer's certification to configured trust-anchors
 - RADIUS access-reject received
 - **eap**: RADIUS access-reject received

Other failures, such as being unable to assign an address, are not counted.

The AUTH failure counter is reset by either a successful authentication before the client is blocked, the expiration of a block timer, or the expiration of the duration timer.

If multiple IPsec clients behind a NAT device share the same public IP address, a limit for the maximum number of clients or ports behind the same IP address can be configured. If the number of ports exceeds the configured limitation, all ports from that IP address are blocked.

The **clear ipsec lockout** command can also be used to manually clear a lockout state for the specified clients.

4.1.11 IPsec tunnel CHILD_SA rekey

SR OS supports CHILD_SA rekeying for both IKEv1 and IKEv2. The following are the behaviors for the rekey:

- **IKEv1 or IKEv2 CHILD_SA rekey initiator**
 - **outbound**

The system immediately switches to the new security association (SA) after a new SA is created.
 - **inbound**

The old SA is kept for three minutes after the new SA is created. Then, it is removed, and upon removal:

 - **IKEv1**

The system does not send a delete message upon removal.
 - **IKEv2**

The systems send a delete message upon removal.
- **IKEv1 or IKEv2 CHILD_SA rekey responder**
 - **outbound**

The system keeps using the old SA for 25 seconds after the new SA is created before switching to the new SA. If a delete message of the old SA is received before 25 seconds, the system removes the old SA and starts using new SA.
 - **inbound**

The old SA is kept for rest of its lifetime. However, if a delete message is received to close the corresponding outbound SA, then the system removes the corresponding inbound SA before its lifetime expires. The system sends a delete message when the old SA lifetime expires.

If the old SA lifetime expires before the 25 seconds or three minutes mentioned above, the old SA is removed upon expiration and the system sends a delete message.

4.1.12 Multiple IKE/ESP transform support

For IPsec tunnels or IPsec gateways, the SR OS allows users to configure up to four IKE transform and four IPsec transform configurations for IKE and ESP traffic.

IKE transform parameters are configured in the **config>ipsec>ike-transform** and referenced in the **ike-policy**, while IPsec transform parameters are configured in the **config>ipsec>ipsec-transform** context and referenced in the tunnel template for dynamic tunnels and under **config>service>vpn>interface>sap>ipsec-tunnel>dynamic-keying** for static tunnels.

IKE transform includes the following configurations:

- IKE encryption algorithm
- IKE authentication algorithm
- Diffie-Hellman group
- IKE SA lifetime

IPsec transform includes the following configurations:

- ESP encryption algorithm

- ESP authentication algorithm
- Diffie-Hellman group for CHILD SA rekey with PFS
- CHILD SA lifetime

If multiple **ike-transform** and **ipsec-transform** parameters are configured for IPsec gateways and IPsec tunnels, the system uses the configured transforms to negotiate with the peer. This negotiation allows IPsec gateways and IPsec tunnels to support peers with different crypto algorithms.

4.2 Using certificates for IPsec tunnel authentication

The SR OS supports X.509v3 certificate authentication for IKEv2 tunnel (LAN-to-LAN tunnel and remote-access tunnel). The SR OS also supports asymmetric authentication. This means the SR OS and the IKEv2 peer can use different methods to authenticate. For example, one side could use pre-shared-key and the other side could use a certificate.

The SR OS supports certificate chain verification. For a static LAN-to-LAN tunnel or ipsec-gw, there is a configurable trust-anchor-profile which specifies the expecting CAs that should be present in the certificate chain before reaching the root CA (self-signed CA) configured in the system.

The SR OS's own key and certificate are also configurable per tunnel or ipsec-gw.

When using certificate authentication, the SR OS uses the subject of the configured certificate as its ID by default.



Note: IPsec application is subject to FIPS restrictions; for more information please see the *7450 ESS, 7750 SR, 7950 XRS, and VSR Basic System Configuration Guide*.

4.2.1 IKEv2 digital signature authentication

RFC 7427 *Signature Authentication in the Internet Key Exchange Version 2 (IKEv2)* defines a new IKEv2 AUTH payload method which not only indicates the type of public key, but also the hash algorithm that used to generate the signature; it also includes a new IKEv2 notification: SIGNATURE_HASH_ALGORITHMS, which is used to signal support of RFC 7427 and a list of support hash algorithms to a peer.

RFC 7427 is the default way to perform certificate authentication for IKEv2. The system negotiates its support with the peer as follows:

- **sending**
 - as tunnel initiator, includes SIGNATURE_HASH_ALGORITHMS in the IKE_SA_INIT request.
 - as tunnel responder, includes SIGNATURE_HASH_ALGORITHMS in IKE_SA_INIT response only if the received IKE_SA_INIT request includes it.
 - includes the SHA1/SHA2-256/SHA2-384/SHA2-512 hash algorithms in SIGNATURE_HASH_ALGORITHMS
- **receiving**
 - If the peer does not include SIGNATURE_HASH_ALGORITHMS in the IKE_SA_INIT packet, then it does not support RFC 7427 and the system uses an RSA Digital Signature for the RSA key (value 1), and DSS Digital Signature (value 3) for the DSA key to generate the AUTH payload.



Note: If the ECDSA key is selected in the cert-profile entry, then the tunnel setup fails in the system.

- If the peer sends SIGNATURE_HASH_ALGORITHMS, then the system uses RFC 7427 and the strongest hash algorithms that is supported by both sides to generate the AUTH payload. If there is no common hash algorithms supported by both sides, the system falls back to RSA Digital Signature (Auth Method value 1) or DSS Digital Signature (Auth Method value 3).



Note: If the selected key is an RSA key, there are specific cases that have a short RSA key with long hash algorithm. The system falls back to RSA Digital Signature for RSA key (value 1) even when both sides send SIGNATURE_HASH_ALGORITHMS and there are common hash algorithms.

To verify the received digital signature of the AUTH payload, the peer must use one of the algorithms in the SIGNATURE_HASH_ALGORITHMS that the system sends. Otherwise, the tunnel setup fails.

The system continues to use CAs in received cert-request payloads to select the cert-profile entry; if the selected entry is an RSA key, then the system needs to decide whether to use PKCS#1-1.5 or RSASS-PSS to generate the signature by using the value set by the `config>ipsec>cert-profile>entry>rsa-signature` command.

4.3 Trust-anchor profile

Since Release 12.0R1, the SR OS supports multiple trust-anchors per ipsec-tunnel/ipsec-gw. Users can configure a trust-anchor-profile that includes up to eight CAs. The system builds a certificate chain by using the certificate in the first certificate payload in the received IKEv2 message. If any of configured trust-anchor CAs in the trust-anchor-profile appears in the chain, then authentication is successful. Otherwise authentication is failed.

The SR OS only supports processing of up to 16 hashes for the trust-anchor list from other products. If the remote end is sending more than 16, and a certificate match is in the > 16 range, the tunnel remains down with authentication failure.

The legacy **trust-anchor** command under ipsec-gw/ipsec-tunnel was deprecated in Release 15.0.R1.

4.4 Cert-profile

Since Release 12.0R1, the SR OS supports sending different certificate/chain according to the received IKEv2 certificate-request payload. This is achieved by configuring a cert-profile which allows up to eight entries. Each entry includes a certificate and a key and optionally also a chain of CA certificates.

The system loads cert/key in cert-profile into memory and build a chain: compare-chain for the certificate configured in each entry of cert-profile upon no shutdown of the cert-profile. These chains are used in IKEv2 certificate authentication. If a chain computation cannot be completed for a configured certificate, then the corresponding compare-chain is empty or only partially computed.

Because there could be multiple entries configured in the cert-profile, the system needs to pick the cert/key in the correct entry that the other side expects to receive. This is achieved by a lookup of the CAs within one cert-request payload or multiple cert-request payloads in the compare-chain and then picking the first

entry that there is a cert-request CA appearing in its chain. If there is no such cert, the system picks the first entry in the cert-profile. The first entry shown in the output below, is the first configured entry in the cert-profile. The entry-id of first entry does not have to be 1.

For example, there are three CA listed in certificate-request payload: CA-1, CA-2 and CA-3, and there are two entries configured in the cert-profile like following:

```
cert-profile "cert-profile-1"
  entry 1
    cert "cert-1"
    key "key-1"
  entry 2
    cert "cert-2"
    key "key-2"
    send-chain
      ca-profile "CA-1"
      ca-profile "CA-2"
```

The system builds two compare-chains: chain-1 for cert-1 and chain-2 for cert-2. Assume CA-2 appears in chain-2, but CA-1 and CA-3 do not appear in either chain-1 or chain-2. Then the system picks entry 2.

After a cert-profile entry is selected, the system generates the AUTH payload by using the configured key in the selected entry. The system also sends the cert in the selected entry as "certificate" payload to the peer.

If a chain is configured in the selected entry, then one certificate payload is needed for each certificate in the configured chain. The first certificate payload in the IKEv2 message is the signing certificate, which is configured by the **cert** command in the chosen cert-profile entry. With the above example, the system sends three certificate payloads: cert-2, CA-1,CA-2.

The following CA chain-related enhancements are supported:

- The no-shut of a ca-profile triggers a re-computation of compute-chain in related cert-profiles. The system also generates a new log-1 to indicate a new compute-chain has been generated; the log includes the ca-profile names on the new chain. Another log-2 is generated if the send-chain in a cert-profile entry is not in compute-chain because of this ca-profile change. Another log is generated if the hash calculation for a certificate under a ca-profile has changed.
- When no-shutting a cert-profile, the system now allows the CAs in the send-chain, not in the compute-chain. The system also generates log-2 as above.
- The system now allows changes of the configuration of send-chain without shutdown of cert-profile.
- When a configure root CA is cross-signed by another CA, multiple overlapping compare-chains for a specific certificate profile entry may occur. Choose one compare-chain by executing the following command to specify the tiebreak CA.

```
configure ipsec cert-profile entry compare-chain-include
```

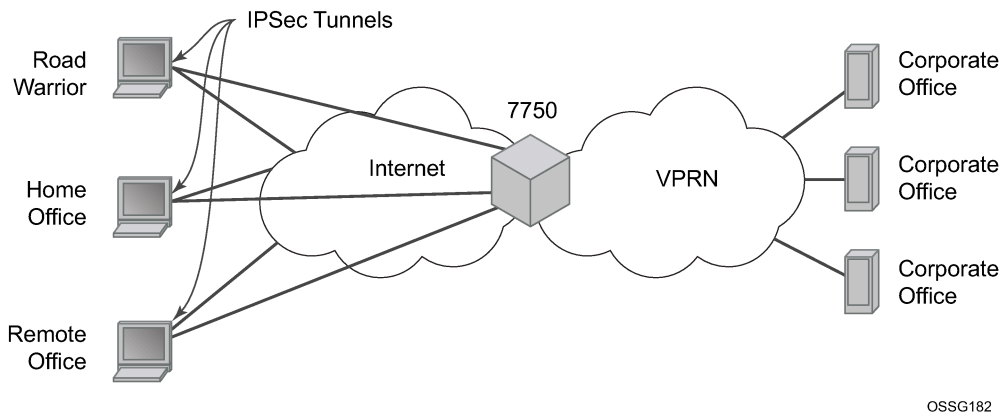
4.5 Video wholesale example

As satellite headend locations can be costly, many municipal and second tier operators cannot justify the investment in their own ground station to offer triple play features. However, it is possible for a larger provider or a cooperative of smaller providers to unite and provide a video headend. Each retail subscriber can purchase content from this single station, and receive it over IP. However, encryption is required so the

signal cannot be understood if intercepted. A high speed encrypted tunnel is preferred over running two layers of double video protection which is cumbersome and computationally intensive.

Figure 40: Video wholesale configuration shows an example of video wholesale configuration.

Figure 40: Video wholesale configuration



OSSG182

4.6 1:1 Multi-chassis IPsec redundancy

This section applies to the 7750 SR, the 7450 ESS, and VSR.

Multi-Chassis IPsec redundancy (MC-IPsec) provides a 1:1 (active/standby) IPsec stateful failover mechanism between two chassis.

The following describes features about MC-IPsec:

- This feature provides protection against ISA failure and chassis failure.
- MC-IPsec is supported for all types of IKEv2 tunnels, including static LAN-to-LAN, dynamic LAN-to-LAN and remote-access tunnels.
- This feature is supported only on multi-active tunnel groups.
- The granularity of failover is per tunnel group, which means a specific tunnel group could failover to the standby chassis independent of other tunnel groups on the master chassis.
- The components included in this feature as shown in the following table.

Table 30: MC-IPsec redundancy feature components

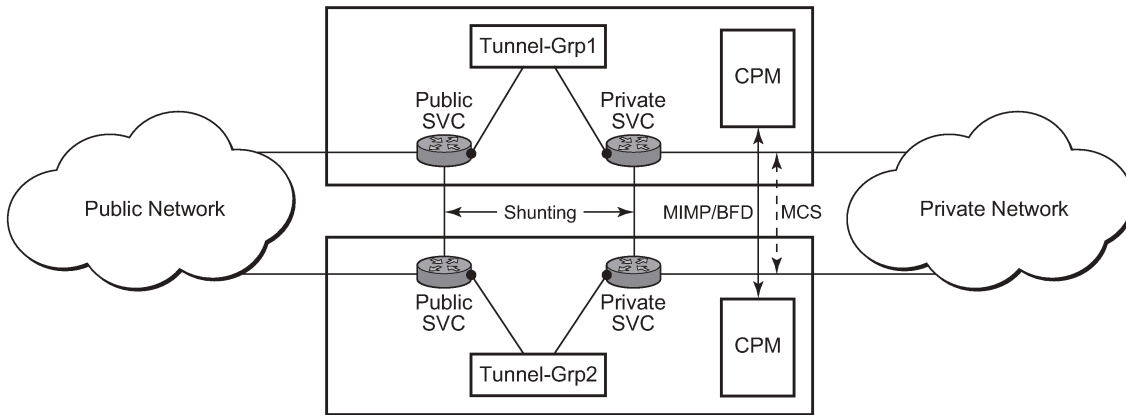
MC-IPsec redundancy feature component	Description
Master election	MC-IPsec Mastership Protocol (MIMP) runs between chassis to elect master, MIMP runs for each tunnel group independently.
Synchronization	Multi-Chassis Synchronization (MCS) synchronizes IPsec states between chassis.
Routing	Routing features include the following:

MC-IPsec redundancy feature component	Description
	<ul style="list-style-type: none"> - MC-IPsec aware routing attracting traffic to the master chassis - shunting support - MC-IPsec aware VRRP

4.6.1 Architecture

The overall MC-IPsec redundancy architecture is displayed in [Figure 41: MC-IPsec architecture](#).

Figure 41: MC-IPsec architecture



al_0099

4.6.2 MIMP

With MIMP enabled, there is a master chassis and a backup chassis. The state of the master or standby is per tunnel-group. For example ([Table 31: Master and backup chassis example](#)), chassis A and B, for tunnel-group 1, A is master, B is standby; for tunnel-group 2, A is standby, B is master.

Table 31: Master and backup chassis example

	Master	Standby
Tunnel Group 1	A	B
Tunnel Group 2	B	A

All IKEv2 negotiation and ESP traffic encryption/decryption only occurs on the master chassis. If the backup chassis receives such traffic, if possible, it shunts them to the master.

There is a mastership election protocol (MIMP) running between the chassis to elect the master. This is an IP-based protocol to avoid any physical topology restrictions.

A central BFD session could be bound to MIMP to achieve fast chassis failure detection.

4.6.2.1 MIMP protocol states

There are five MIMP states:

- discovery
- notEligible
- eligible
- standby
- master

The five MIMP states are described as follows:

- **discovery**

Upon enabling MC-IPsec for the tunnel-group, for example:

1. system starts up
2. no shutdown MC-IPsec peer
3. no shutdown MC-IPsec tunnel-group

Functionally, this means blackhole traffic to the ISA and no shunting.

If the peer is reached before the discovery-interval (configurable) has expired, then the state is changed to whatever the MIMP decides

If the peer is not reached before the discovery-interval has expired, then the state is changed to **eligible** or **notEligible** depending on the oper-status of the tunnel-group.

- **notEligible**

The tunnel-group is operationally down. Functionally, this means blackhole traffic to the ISA and no shunting.

- **eligible**

The peer is not reachable or the associated BFD session is down but the tunnel-group is operationally up. Functionally, this means the ISA processes traffic.

- **standby**

Peer is reachable, elected standby. Functionally, this means blackhole traffic to ISA and shunting if possible.

- **master**

Peer is reachable, elected master. Functionally, this means the ISA processes traffic.

4.6.2.2 Election logic

The following election logic is executed when MIMP packets are exchanged.

1. Calculate master eligibility.

- Set masterEligible to TRUE if the local tunnel group is operationally up, otherwise FALSE.
- Set peerMasterEligible to TRUE if the peer's tunnel group is operationally up, otherwise FALSE.

2. Elect based on eligibility.

- If masterEligible and not peerMasterEligible, elect self master → DONE.

- If not masterEligible and peerMasterEligible, elect peer master → DONE.
 - If not masterEligible and not peerMasterEligible, no master → DONE.
3. Apply stickiness rules (mastership tends not to change).
- If I was "acting master" and peer was not "acting master", then elect self master -> DONE.
 - If I was not "acting master" and peer was "acting master", then elect peer master -> DONE.
- An "acting master" is either in MIMP state "master" or "eligible".
4. Elect based on priority and number of active ISAs.
- If my priority is higher than peer, elect self master → DONE.
 - If peer priority is higher than mine, elect peer master → DONE.
 - If I have more active ISA than peer, elect self master → DONE.
 - If peer has more active ISA than me, elect peer master → DONE.

The tie breaker:

- If the local chassis' MIMP source address is higher than the peer's, elect self master → DONE.
- Elect peer master → DONE.

4.6.2.3 Protection status

Each MC-IPsec-enabled tunnel-group has a "protection status", which could be one of following:

- **notReady**

The tunnel-group is not ready for a switchover because there is either no elected standby to takeover, or there are pending IPsec states which need to be synced. If switchover occurs with this status, then there could be a significant traffic impact.

- **nominal**

The tunnel-group is in a better situation to switchover than notReady. However, traffic still may be impacted.

Protection status serves as an indication for the operator to decide the optimal time to perform a controlled switchover.

The **show redundancy multi-chassis mc-ipsec peer <ip-address> tunnel-group <tunnel-group-id>** command can be used to check current protection status.

4.6.2.4 Other details

- Mastership election is per tunnel-group.
- MIMP is running in the base routing instance.
- MIMP uses the configured value of the **config>redundancy>multi-chassis>peer>source-address** command as the source address. If not configured, then system address is used.
- The priority range is from 0 to 255.
- When an mc-ipsec tunnel-group enters standby from acting master, the tunnel-group is restarted.

- When a tunnel-group enters an admin shutdown state under the mc-ipsec configuration (add a tunnel-group to mc-ipsec configuration, or upon admin shutdown of an mc-ipsec enabled tunnel group):
 - All tunnels in the tunnel-group are deleted/reinstalled to the ISA.
 - All IKE states associated with those tunnels are locally purged from the MS-ISAs.
 - No IKE messages are sent to the IKE peer.These behaviors occur regardless of the presence of a redundant chassis or the state of a redundant chassis.
- With MC-IPsec enabled:
 - auto-establish is blocked.
 - For DPD configuration, only **no dpd** and **dpd** configurations with **reply-only** are allowed.

4.6.3 Routing

4.6.3.1 Routing in public service

A /32 route of the local tunnel address is created automatically for all tunnels on the MC-IPsec enabled tunnel-group.

This /32 route can be exported to a routing protocol by a route policy. The protocol type in route-policy is IPsec.

To attract traffic to the master chassis, a route metric of these /32 routes could be set according to the MIMP state, a metric from the master chassis is better than a metric from the standby chassis. There are three available states that can be used in the **from state** command in the route policy entry configuration:

- IPsec-master-with-peer
 - Corresponding MIMP states: master
- IPsec-master-without-peer
 - Corresponding MIMP states: eligible
- IPsec-non-master
 - Corresponding MIMP states: discovery/standby

However, if the standby chassis receives IPsec traffic, the traffic is shunted to the master chassis by forwarding to a redundant next-hop. The redundant next-hop is an IP next-hop in the public routing instance.

4.6.3.2 Routing in private services

For static LAN-to-LAN tunnels, the static route with the IPsec tunnel as the next-hop could be exported to a routing protocol by a route policy. The protocol type remains **static**. For dynamic LAN-to-LAN tunnels, the reverse-route could be exported to a routing protocol by a route policy. The protocol type is **ipsec**. For remote-access tunnel, the private interface route could be exported to a routing protocol by a route policy.

Similar to routing in public services, the route metric of the above the routes could be set according to the MIMP state. Only a static route with an IPsec tunnel as the nexthop and reverse route has an MIMP state.

If the standby chassis receives IPsec traffic, the traffic is shunted to the master chassis by forwarding to a redundant next-hop. The redundant next-hop is an IP next-hop in a private routing instance.

4.6.3.3 Other details about shunting

Shunting only works when tunnel-group is operationally up.

Shunting is not supported over auto-bind tunnels.

4.6.4 MC-IPsec aware VRRP

In many cases, the public side is a Layer 2 network and VRRP is used to provide link or node protection. However, VRRP and MC-IPsec are two independent processes, each has its own mastership state, which means the VRRP master could be different from MC-IPsec master. This results in shunting traffic unnecessarily.

To address this issue, MC-IPsec aware VRRP is introduced in SR OS Release 10.0R8, which add a new priority event in vrrp-policy: mc-ipsec-non-forwarding. If the configured tunnel-group enters non-forwarding (non-master) state, then the priority of associated VRRP instance is set to the configured value. Delta priority is not supported for this type of event.

4.6.5 Synchronization

To achieve stateful failover, IPsec states are synced between chassis by using the MCS protocol.

- Only successfully created SA after a completed INITIAL EXCHANGES or CREATE_CHILD_SA EXCHANGES is synced.
- Upon switchover, the new standby chassis reboots the tunnel-group.
- The ESP sequence number is not synced except for the high 32 bits of extended sequence numbers.
- The CLI configuration is not synced.

The time must be the same on both chassis (using NTP/SNTP to sync to the same server is an option).

4.6.5.1 Automatic CHILD_SA rekey

Because the ESP sequence number is not synced, a CHILD_SA rekey for each tunnel is initiated by the new active chassis to reset the sequence number upon switchover.

4.6.5.2 Encryption of Synced States

Transport encryption of synced IPsec states can be configured using the **config> redundancy> multi-chassis> peer> sync> transport-encryption> application ipsec** command, with the **ipsec** parameter as the keychain name. This causes the system to encrypt the IPsec states during transportation between the active and standby.

The key used to encrypt states is specified by the referenced keychain, which is defined in the **config> system> security> keychain** context. The **keychain** provides the user the ability to gracefully enable or disable encryption or change the key.

4.6.5.2.1 Gracefully enable encryption steps

Prerequisites

To use a keychain for MCS IPsec, the keychain **aes-128-gcm-16** entry algorithm must be configured.

Use the following steps to gracefully enable encryption. Chassis A and chassis B must both run releases that support the transport-encryption feature. Chassis A is active and chassis B is standby.

Procedure

Step 1. On chassis B, change the configuration to add MCS encryption, using a keychain with two bidirectional entries in the **configure system security keychain direction bi** entry.

- For entry-1, use the following values:
 - For entry-id, use null-key.
 - For begin-time, use now.
 - For tolerance, use forever.
- For entry-2, use the following values:
 - For algorithm, use aes-128-gcm-16.

- For begin-time: t1, add enough time to complete the next step (for example, if the current time is 2019/4/18 10:00 UTC, then add one hour to complete step 2, **begin-time** 2019/4/18 11:00 UTC).

Because both chassis A and chassis B are still using clear transport, A can successfully synchronize states to B.

Step 2. On chassis A, change the configuration to add MCS encryption, using a keychain with two bidirectional entries in the **configure system security keychain direction bi** context:

- For entry-1, use the following values:
 - For entry-id, use null-key.
 - For begin-time, use now.
 - For tolerance, use forever.
- For entry-2, use the following values:
 - For algorithm, use aes-128-gcm-16.
 - For begin-time: t1, use the same value as in step 1 , **begin-time** 2019/4/18 11:00:00 UTC.

Because both A and B can receive either clear or encrypted states, synchronization is successful.

Step 3. After t1, remove entry-1 from both chassis using the **configure system security keychain direction bi no entry 1** command

What to do next

For an example configuration, see [Configuring MCS encryption](#)

4.6.5.2.2 Gracefully change the key steps

About this task

Use the following steps to gracefully change the key. Chassis A and chassis B both have configured **transport-encryption** where chassis A is active and chassis B is standby.

Procedure

Step 1. On Chassis B, change the configuration to add a new keychain, **entry-y**, with a new key:

- entry-x (the old entry)
 - For tolerance, use forever (this step can be performed without administratively disabling the entry).
- entry-y:
 - For algorithm, use aes-128-gcm-16.
 - For begin-time: t1, add enough time to ensure there is enough time to complete the next step (for example, if the current time is 2019/4/18 10:00 UTC, then add one hour to complete step 2, **begin-time** 2019/4/18 11:00 UTC).

At this point, both A and B are still configured to use the old key (**entry-x**) for transport, so A successfully synchronizes states to B.

Step 2. On chassis A, change the configuration to add a new keychain, **entry-y**, with new key:

- entry-x (the old entry)

For tolerance, use forever (this can be performed without shutting down the entry).

- entry-y:
 - For algorithm, use aes-128-gcm-16.
 - For begin-time: t1, use the same value as in step 1. **begin-time** 2019/4/18 11:00 UTC.

Step 3. After t1, both A and B begin to use entry-y. Remove entry-x from both chassis using the **configure system security keychain direction bi no entry x** command.

4.6.6 MC-IPsec Responder only

With MC-IPsec, it is required that MC-IPsec pair can only act as an IKEv2 responder (except for the automatic CHILD_SA rekey upon switchover). To enable this behavior, configure the **configure isa tunnel-group ipsec-responder-only** command.

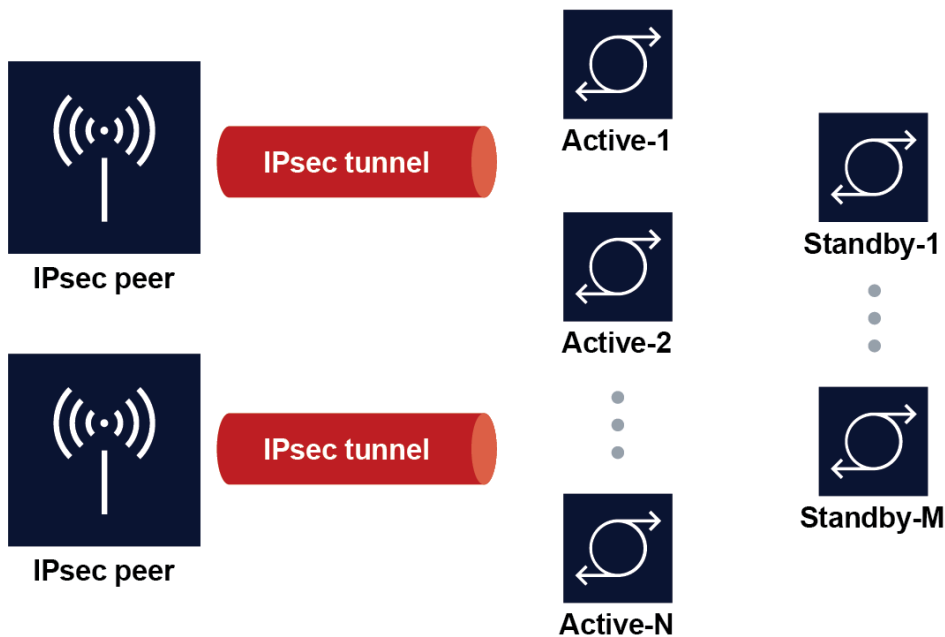
4.7 N:M MC-IPsec redundancy

In addition to 1:1 MC-IPsec, the system also supports N:M multi-chassis IPsec stateful redundancy. The term N:M is used to refer to this feature. N:M provides additional redundancy and potential cost saving compared to 1:1 MC-IPsec.

N:M allows overall N active nodes and tunnel groups to be backed up by the M standby node and tunnel group. The number represented by N can be larger than, equal to, or smaller than the number represented by M.

The following graphic shows a typical N:M configuration.

Figure 42: N active nodes backup by M standby nodes



sw1303

4.7.1 Redundancy domain

The granularity of protection for N:M is per tunnel group. For each to-be-protected tunnel group (a redundancy domain). It is a 1:1 mapping between the tunnel group and domain. Each domain contains up to four nodes. Among the possible nodes are the following:

- One node is elected to be the active node for the tunnel group
- Other nodes are the standby node for the tunnel group

Tunnel states are synchronized between all nodes in the domain.

A node may have multiple tunnel groups, and in turn, participate in multiple redundancy domains. See [Election](#) for information about how the election of nodes is determined.

4.7.2 Redundancy role

Within a specific redundancy domain, each node has a designated role and an operational role:

- **Designated role**
Designated roles are user-configured as the designated active (DA) or designated standby (DS).
- **Operational role**
Operational roles are decided by election protocol during runtime as operational active (OA) or operational standby (OS).

The user can configure each node in the domain as either DA or DS. Any combination of DA:DS is allowed within a domain; however, only one node is elected as OA, other nodes are OS. A designated active role may be elected as OS and a DS may be elected as operationally active.

4.7.3 ISA tunnel-member pool

N:M allows the use of a single set of ISAs to backup multiple tunnel groups. This is achieved by the **tunnel-member-pool** command. The tunnel-member pool includes a set of ISAs, as well as a tunnel group with a DS role that refers to the tunnel-member pool instead of directly referring to the ISA. Multiple DS tunnel groups can refer to the same tunnel-member pool, and in turn, share same set of ISAs.

When a failover occurs, if a DS tunnel group is elected to become the OA, the system then takes the required ISAs out of the underlying tunnel-member pool and uses them to populate the tunnel group and the CPM downloads corresponding tunnel states to the selected ISAs. The number of ISAs taken out of the underlying tunnel member pools equals the value configured by the **active-mda-number** command in the DS tunnel group. If there is enough ISA left in the tunnel member pool for other DS tunnel groups, additional failovers can be supported.

A DA tunnel group refers to the ISA directly, and synchronized states are downloaded to the ISA directly. In the case of a DS tunnel group, synchronized tunnel states are stored first on the CPM and only downloaded to the ISA upon failover. Because of this difference, the traffic impact during a failover is larger than when failing over to a DS tunnel group than to a DA tunnel group.

A DS tunnel group can only refer to a tunnel-member pool, not directly to the ISA.

By default, the synced tunnel states are only downloaded to ISA on switchover. When the optional High Availability (HA) mode for a tunnel-member pool is configured, the CPM pre-downloads the synced tunnel states of all the tunnel groups associated with a pool to every ISA in the pool. After the switchover, the system picks the required number of ISAs out of the pool and activates the states of the switched-over

tunnels. Hence, the switchover to a DS is much faster than the default behavior. However, when compared with failover to DA, this mode is still slower because of the extra time required to activate the states.

With HA, a scale mode is specified, which defines the maximum number of tunnels of all tunnel groups associated with a pool.

4.7.4 Redundancy state and protection status

For each participating domain, an N:M node has a redundancy state with one of following values:

- **Discovery**

The node's initial peer discovery, for example, when the system boots up.

- **notEligible**

The node is not eligible to be elected as OA, for example, when the local ISA fails.

- **Eligible**

The node is unable to reach election consensus with its peers, but the node can function locally as OA, for example, the local tunnel group is up but unable to reach its peers.

- **Standby**

The node is able to reach election consensus with its peers and is elected as OA.

- **Active**

The node is able to reach election consensus with its peers and is elected as OA.

The system also has a protection status for each participating domain that can be Nominal or notReady. In case of a switchover, traffic impact of notReady can be higher than Nominal.

The protection status serves as an indication for to decide the optimal time to perform a controlled switchover.

Use the following command to check the redundancy state and protection status for a domain:

```
show redundancy multi-chassis ipsec-domain
```

4.7.5 Election

The election protocol MIMPv2, for N:M, is used to elect the OA node in a domain. A full meshed MIMPv2 peering is required between all nodes in the domain for each participating domain. The system has a user-configurable priority. The designated role, priority, and router ID impact how an active node is elected.

The election is done through evaluation using the following ordered rules. The evaluation stops when a single winner is left.

1. DA is preferred over DS.
2. With the same designated role, the node with a higher priority is preferred.
3. With the same designated role and priority, the node with the higher router ID is preferred. MIMPv2 uses the router ID as the node identification. This can be different from the MIMPv2 packet's IP address.

A central BFD session can be bound to a MIMPv2 peering to accelerate node failure detection.

A failover is triggered by one of following events:

- node failure
- ISA failure
- manually forcing a switchover. Use the following command to manually force a switchover.

```
tools perform redundancy multi-chassis mc-ipsec force-switchover domain
```

Because a domain can contain up to four nodes, there can be up to three failovers to three different nodes for a tunnel group. This gives more node-level redundancy than the 1:1 MC-IPsec feature.

The following example shows the election.

Example

Domain 1 contains node A, B, C, and D:

Table 32: Election of nodes

Node	Designated role	Operational role	Priority
A	DA	OA	200
B	DA	OS	100
C	DS	OS	200
D	DS	OS	150

If node A fails, then the following occurs.

1. First, node A fails over to node B, because B is DA and C and D are DS. DA is preferred over DS.
2. If node B is not able to recover (if node B also failed), then it fails over to node C because both C and D are DS. Node C has the higher priority and as a result is selected as OA.
3. Lastly, if both node B and C are unable to cover, the only available node D is elected as OA.

4.7.6 Revertive

The revertive function allows a node in an N:M domain to automatically take over as the active router in the domain when it becomes eligible. In this example, the node is called the challenger. To be eligible, the challenger must meet all the following criteria:

1. The domain must be configured as revertive on the challenger.
2. The challenger must continuously have been in a redundancy state, which means on standby for the last 60 seconds.
3. The challenger must continuously have been in a protection status, that is, nominal for the last 60 seconds.
4. The challenger must meet the following criteria:
 - The challenger is DA while the existing OA is DS.
 - The challenger is DS while the existing OA is also DS.
 - With the same designated role, the challenger has a higher priority than the existing OA.

- With the same designated role and priority, the challenger has a higher router ID than the existing OA.

If the 60-second interval for criteria 2 and 3 have elapsed, and the challenger notices that another revertive peer with better properties than in criterion 4 is present and eligible, the challenger continues for another 60 seconds before trying to take over.

4.7.7 DA versus DS

The following table summarizes the differences between DA and DS.

Table 33: DA and DS differences

	DA	DS
ISA member	N/A. Requires a dedicated ISA for each tunnel group.	ISA member pool allows sharing the ISA across multiple tunnel groups.
Sync	Synchronized states are pre-downloaded to the ISA. There is less traffic impact during failover.	Synchronized states are kept on the CPM and downloads to the ISA only upon failover. An optional HA mode allows pre-downloading states to ISA and activating upon switchover.
Election	DA is preferred over DS	
Revertive failover	A DS challenger cannot assume the OA role from a DA node.	
admin ipsec display-key	Supported in case of both OA and OS.	Only supported as OA.

4.7.8 State synchronization

To achieve stateful failover, a fully meshed MCS peering between nodes in a domain is used to synchronize IPsec states across nodes.

The following criteria is considered to synchronize IPsec states across nodes:

- An SA is only successfully created after a completed Initial Exchanges or CREATE_CHILD_SA Exchange is synced.
- Upon switchover, the new OS node resets the tunnel group.
- The ESP sequence number is not synchronized.
- The CLI configuration is not synchronized.
- The system time must be synchronized across all participating nodes (for example, using NTP or SNTP to synchronize the system time).
- Because the ESP sequence number is not synchronized, a CHILD_SA rekey for each tunnel is initiated by the new OA to reset the sequence number on switchover.

- MCS encryption, as described in 1:1 MC-IPsec, can also be used for N:M.

4.7.9 Routing

Like 1:1 MC-IPsec, only an OA in the domain can process IPsec traffic, therefore it is required to always attract traffic to the OA node. This is achieved by MC-aware routing, where the same IPsec routes are advertised from all nodes in a domain. The route advertised by OA has best metric, so traffic is attracted to it. This is achieved by using a route policy that sets different metrics according to the redundancy state of the tunnel group that the routes belong to. The system supports the following states in the route policy:

- **IPsec-master-with-peer**
The corresponding redundancy states are active.
- **IPsec-master-without-peer**
The corresponding redundancy states are eligible.
- **IPsec-non-master**
The corresponding redundancy states are discovery, standby, or notEligible.

4.7.10 VRRP

As with 1:1 MC-IPsec, N:M also supports MC-aware VRRP by using following command:

```
priority-event mc-ipsec-non-forwarding
```

This is equivalent to the following command in routing:

Example: MD-CLI

```
configure policy-options policy-statement entry from state ipsec-non-master
```

Example: classic CLI

```
configure router policy-options policy-statement entry from state ipsec-non-master
```

A VRRP policy can set the priority of the VRRP instance to a configured value after the associated tunnel group enters the redundancy state. Delta priority is not supported for this type of event.

4.7.11 Shunting

Optionally, if an OS node receives IPsec traffic it can shunt the traffic to the OA node to minimize traffic loss.

For each routing instance, N:M allows users to define a multi-chassis shunting profile that specifies, for each N:M peer, the IP interface used to shunt traffic and the corresponding next hop is specified in the following command:

```
multi-chassis-shunt-interface
```

Because shunting is performed by forwarding traffic to the specified next hop over the shunt interface, the shunt interface must be an interface directly connected to its peer a spoke SDP-terminated interface.

4.7.12 Responder only

N:M states that the required node can only act as an IKEv2 responder (except for the automatic CHILD_SA rekey upon switchover). This behavior is enabled with the following command:

```
configure isa tunnel-group ipsec-responder-only
```

4.7.13 Coexisting tunnel groups

The system supports a 1:1 MC-IPsec tunnel group that coexists with an N:M tunnel group. A specific tunnel group cannot be both 1:1 and N:M at the same time.

4.7.14 Provisioning requirements

There are specific provisioning requirement for N:M to work properly, see [IPsec deployment requirements](#) for details.

4.7.15 Configuring N:M

The following configurations are required across all nodes:

- IPsec-specific configurations, such as IKE policy, transform, IPsec tunnel, IPsec gateway, and so on
- N:M redundancy configurations require the following:
 - **configure redundancy multi-chassis peer mc-ipsec domain**
 - Refer to the **domain** parameter with the **configure redundancy multi-chassis peer mc-ipsec domain** command.
 - Enable MCS with the **configure redundancy multi-chassis peer sync ipsec** command.
- For the VRRP routing configuration, use the MC state in the route policy and VRRP policy.
- Shunting configurations require the following:
 - Use either of the following commands to specify the shunt interface name and next-hop address.

```
configure router ipsec multi-chassis-shunt-interface  
configure service vprn ipsec multi-chassis-shunt-interface
```

- Use either of the following commands to specify the shunt interface to use for a specific peer.

```
configure router ipsec multi-chassis-shunting-profile  
configure service vprn ipsec multi-chassis-shunting-profile
```

- Refer to the following command configured under the public or private tunnel interface.

```
multi-chassis-shunting-profile
```

DS tunnel-group-specific configurations require the following:

- The following command specifies a set of ISAs to be included in the pool.

```
configure isa tunnel-member-pool
```

- The tunnel group refers to the tunnel member pool configured with the following command:

```
configure isa tunnel-member-pool
```

The following example shows a DA in an N:M redundancy configuration.

Example: MD-CLI

```
[ex:/configure redundancy multi-chassis]
A:admin@node-2#
  ipsec-domain 1 {
    designated-role active
    priority 250
    tunnel-group 1
  }
  peer 84.84.84.84 {
    sync {
      ipsec true
    }
    mc-ipsec {
      bfd-liveness true
      domain 1 {
      }
    }
  }
  peer 85.85.85.85 {
    sync {
      ipsec true
    }
    mc-ipsec {
      bfd-liveness true
      domain 1 {
      }
    }
  }
}
```

Example: classic CLI

```
*A:admin@node-2>config>redundancy#
-----
  multi-chassis
    ipsec-domain 1 create
      designated-role active
      priority 250
      tunnel-group 1
      no shutdown
    exit
  peer 84.84.84.84 create
    sync
      ipsec
      no shutdown
    exit
  mc-ipsec
    bfd-enable
    domain 1 create
      no shutdown
```



```

        exit
    exit
    no shutdown
exit
peer 85.85.85.85 create
sync
    ipsec
    no shutdown
exit
mc-ipsec
    bfd-enable
    domain 1 create
    no shutdown
    exit
exit
no shutdown
exit
exit

```

The following example shows a DA of a route policy configuration on the private side.

Example: MD-CLI

```

[ex:/configure policy-options]
A:admin@node-2#
policy-statement "export400" {
  entry 10 {
    from {
      state ipsec-master-with-peer
      protocol {
        name [ipsec]
      }
    }
    action {
      action-type accept
      local-preference 255
      community {
        add ["vprn400"]
      }
    }
  }
  entry 20 {
    from {
      state ipsec-non-master
      protocol {
        name [ipsec]
      }
    }
    action {
      action-type accept
      local-preference 70
      community {
        add ["vprn400"]
      }
    }
  }
  entry 30 {
    from {
      state ipsec-master-without-peer
      protocol {
        name [ipsec]
      }
    }
  }
}

```

```

    action {
      action-type accept
      local-preference 100
      community {
        add ["vprn400"]
      }
    }
  }
}

```

Example: classic CLI

```

*A:admin@node-2>config>router>policy-options#
-----
policy-statement "export400"
  entry 10
    from
      protocol ipsec
      state ipsec-master-with-peer
    exit
    action accept
      community add "vprn400"
      local-preference 255
    exit
  exit
  entry 20
    from
      protocol ipsec
      state ipsec-non-master
    exit
    action accept
      community add "vprn400"
      local-preference 70
    exit
  exit
  entry 30
    from
      protocol ipsec
      state ipsec-master-without-peer
    exit
    action accept
      community add "vprn400"
      local-preference 100
    exit
  exit
exit

```

The following example shows a shunting configuration on the private side.

Example: MD-CLI

```

[ex:/configure service vprn "400" ipsec]
A:admin@node-2#
  multi-chassis-shunt-interface "to84" {
    next-hop {
      address 130.100.14.4
    }
  }
  multi-chassis-shunt-interface "to85" {
    next-hop {
      address 130.110.15.5
    }
  }

```

```

}
multi-chassis-shunting-profile "shunt1" {
  peer 84.84.84.84 {
    multi-chassis-shunt-interface "to84"
  }
  peer 85.85.85.85 {
    multi-chassis-shunt-interface "to85"
  }
}
[ex:/configure service vprn "400" interface "priv"]
A:admin@v21# info
  tunnel true
  multi-chassis-shunting-profile "shunt1"

```

Example: classic CLI

```

*A:admin@node-2>config>service>vprn>ipsec#
-----
multi-chassis-shunt-interface "to84" create
  next-hop 130.100.14.4
exit
multi-chassis-shunt-interface "to85" create
  next-hop 130.110.15.5
exit
multi-chassis-shunting-profile "shunt1" create
  peer 84.84.84.84 create
    multi-chassis-shunt-interface "to84"
  exit
  peer 85.85.85.85 create
    multi-chassis-shunt-interface "to85"
  exit
exit

*A:admin@node-2>config>service>vprn#
interface "priv" tunnel create
  ""
  multi-chassis-shunting-profile "shunt1"

```

The following example shows an ISA member pool configuration example for DS tunnel-groups. A single ISA 1/2 is shared by two DS tunnel-groups.

Example: MD-CLI

```

[ex:/configure isa]
A:admin@node-2#
  tunnel-group 1 {
    admin-state enable
    ipsec-responder-only true
    multi-active {
      member-pool "p1"
    }
    reassembly {
      max-wait-time 2000
    }
  }
  tunnel-group 2 {
    admin-state enable
    ipsec-responder-only true
    multi-active {
      member-pool "p1"
    }
  }

```

```

    reassembly {
        max-wait-time 2000
    }
}
tunnel-member-pool "p1" {
    isa 1/2 { }
}

```

Example: classic CLI

```

*A:admin@node-2>config>isa#
  tunnel-member-pool p1 create
    mda 1/2
  exit
  tunnel-group 1 create
    reassembly 2000
    ipsec-responder-only
    multi-active
    member-pool p1
    no shutdown
  exit
  tunnel-group 2 create
    reassembly 2000
    ipsec-responder-only
    multi-active
    member-pool p1
    no shutdown
  exit

```

4.8 IPsec deployment requirements

The following information describes requirements to deploy SR OS IPsec features.

IPsec general

To avoid high CPU loads and some complex cases, the following are the requirements to configure IKEv2 lifetime:

- The IKE_SA lifetime on one side should be approximately twice as large as the other side. The CHILD_SA lifetime on one side should be approximately two or three times larger than the other side.
- With the previous rule, the lifetime of the side with smaller lifetime should not be too small:
 - IKE_SA: >= 86400 seconds
 - CHILD_SA: >= 3600 seconds
- With first rule, on the side with the smaller lifetime, the IKE_SA lifetime should be at least three times larger than CHILD_SA lifetime.
- The IKE protocol is the control plane of IPsec, therefore, the IKE packet should be treated as high QoS priority in the end-to-end path of the public service.

On a public interface, a SAP ingress QoS policy should be configured to ensure the IKE packet is treated as high QoS priority.
- The correct system time is required for certificate authentication to work properly.
- The peer's DPD interval must be larger than 30 seconds and should not send a DPD request if it receives IKE or ESP traffic.

MC-IPsec specific

The following are requirements for MC-IPsec specific deployments:

- The IKEv2 lifetime requirements from [IPsec general](#) should be applied with special care to MC-IPsec deployments.

In an MC-IPsec deployment where the MC-IPsec pair peers with single, non-redundant IKE clients, the IKEv2 lifetime requirements must be applied with the larger lifetimes configured on the MC-IPsec pair.

An MC-IPsec deployment where one MC-IPsec pair peers with another MC-IPsec pair is not recommended. MC-IPsec performs optimally when the multi-chassis pair peers with a single IKE entity. If such a peering (MC-to-MC) is created, the above IKEv2 lifetime requirements should still be followed. However, with one side nominated to be the primary rekey initiator and having the smaller configured lifetimes.

- Responder-only configuration is a mandatory requirement for all types of tunnels on the MC-IPsec pair in the usual deployment scenario of a MC-IPsec pair peering with single, non-redundant IKE clients.
- DPD on the peer side (following the DPD requirement in the above IPsec General section), **dpd interval 300 max-retries 3 reply-only** on the MC-IPsec side.
- Dedicated, redundant, direct physical link between chassis with enough bandwidth for MCS and shunting traffic.

MIMP/MCS and BFD for MC-IPsec traffic must be forwarded over resilient links so that a single IOM/IMM, MDA or port failure does not cause the MIMP to go down. Because this control traffic is forwarded in the base routing instance, the base routing instance links need to spread over multiple ports on multiple IOM/IMMs. Proper QoS configuration is needed to make sure the control traffic gets the highest priority.

- A MC-IPsec switchover when the protection status is not nominal may result in unexpected behavior and traffic loss. A nominal state must be reached on both MC-IPsec chassis before a MC-IPsec switchover is triggered.
- When using VRRP in the public service and a chassis failure occurs, the VRRP/Layer 2 network should re-converge before the MC-IPsec switchover occurs. One way to speed up VRRP switchover is to bind a BFD session to VRRP.
- The system time of the master and standby chassis must be the same. One way to achieve this is for both chassis to sync to an NTP or SNTP server.
- The CLI configuration is not synchronized via MCS so the user must provision the same IPsec-related configurations on the master and standby chassis. This includes using the same IKE policy ID, tunnel template ID, public or private interface name, and so on.
- For an MC-IPsec shunting interface, only one next-hop address is supported; in case of multiple redundant next hops configured using the same shunting interface, the last next hop configured is used.

N:M specific

The following are requirements for N:M-specific deployments.

- All MC-IPsec specific requirements also apply to N:M.
- A fully meshed MIMPv2/MCS peering is required between all participating nodes for a specific domain.
- When adding a node to a domain, ensure that the previous nodes in the domain have fully meshed MIMPv2/MCS peering before adding the new node.

4.9 IKEv2 remote-access tunnel

Since 11.0R6, SR OS supports IKEv2 remote-access tunnel, the difference between a remote-access tunnel and LAN-to-LAN tunnel is remote-access tunnel allows client to request an internal address (and other attributes like DNS address) via IKEv2 configuration payload. The SR OS supports IKEv2 remote-access tunnel with following features:

- authentication methods:
 - pre-shared-key with RADIUS (**psk-radius**) or without RADIUS (**psk**)
 - certificate with RADIUS (**cert-radius**) or without RADIUS (**cert**)
 - EAP/EAP-Only with RADIUS
- internal address assignment via IKEv2 configuration payload
- address assignment support:
 - RADIUS server based
 - local address assignment
- RADIUS accounting to report address usage
- RADIUS disconnect message to remove tunnel
- NAT-Traversal support
- support MC-IPsec

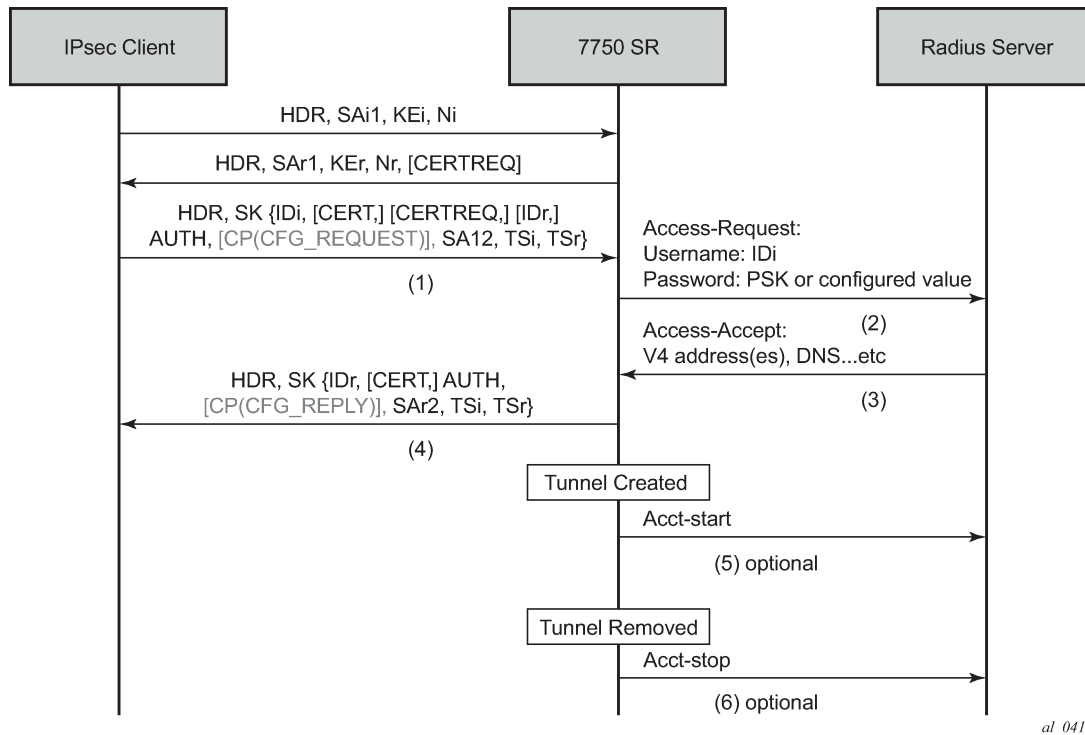
The SR OS only supports address assignments in first CHILD_SA negotiation.

4.9.1 IKEv2 remote access tunnel – RADIUS-based PSK/certificate authentication

If the **auth-method** parameter in the **ike-policy** is configured as **psk-radius** or **cert-radius**, then the system authenticates the client via PSK or certificate accordingly as like a LAN-to-LAN tunnel. The difference being that in the case of **psk-radius** or **cert-radius**, the system also performs a RADIUS authentication or authorization and optionally send RADIUS accounting messages.

[Figure 43: Call flow for psk-radius/cert-radius](#) displays a typical call flow for psk-radius and cert-radius.

Figure 43: Call flow for psk-radius/cert-radius



al_0414

The Access-Request includes the following attributes:

- Username is IDi.
- User-Password is one of following value's hash according to section 5.2 of RFC 2865, *Remote Authentication Dial In User Service (RADIUS)*:
 - client's PSK if the psk-radius is configured (see the CLI)
 - otherwise, a CLI configured key via the **password** command in the radius-authentication-policy; if password is not configured in this case, then system does not include the User-Password attribute in access-request.
- Acct-Session-Id represents the IPsec tunnel session.
 The format is: local_gw_ip-remote_ip:remote_port-time_stamp.
 For example: 172.16.100.1-192.168.5.100:500-1365016423.
- Other RADIUS attributes (dependent on the **config>ipsec>radius-auth-policy> include-radius-attribute** configuration) are as follows:
 - Called-Station-Id (local tunnel address)
 - Calling-station-Id (remote tunnel address:port number)
 - Nas-Identifier (the system name)
 - Nas-Ip-Address (the system IP)
 - Nas-port-id (the public tunnel SAP ID)

If the RADIUS authentication is successful, then the RADIUS server sends an access-accept message back; otherwise, an access-reject message is sent back.

The following are supported attributes in access-accept:

- Alc-IPsec-Serv-Id
- Alc-IPsec-Interface
- Framed-IP-Address
- Framed-IP-Netmask
- Alc-Primary-Dns
- Alc-Secondary-Dns
- Alc-IPsec-Tunnel-Template-Id
- Alc-IPsec-SA-Lifetime
- Alc-IPsec-SA-PFS-Group
- Alc-IPsec-SA-Encr-Algorithm
- Alc-IPsec-SA-Auth-Algorithm
- Alc-IPsec-SA-Replay-Window

After the tunnel is successfully created, the system could optionally (depending on the configuration of the **radius-accounting-policy** under the **ipsec-gw** context), send an accounting-start packet to the RADIUS server, and also send an accounting-stop when the tunnel is removed. The user can also enable the **interim-update** option in the **radius-accounting-policy**.

The following are some attributes included in the acct-start/stop and interim-update:

- Acct-status-type
- Acct-session-id (the same as in the access-request)
- Username

The following attributes are dependent on the **radius-acct-policy>include-radius-attribute** configuration:

- Frame-ip-address: the assigned internal address
- Calling-station-id
- Called-station-id
- Nas-Port-Id
- Nas-Ip-Addr
- Nas-Identifier
- Acct-Session-Time (tunnel session time, only in acct-stop packet)

For a complete list of supported attributes, see the *7450 ESS, 7750 SR, and VSR RADIUS Attributes Reference Guide*.

The system also supports RADIUS disconnect messages to remove an established tunnel, If **accept-coa** (existing command) is enabled in the radius-server configuration, then the system accepts the disconnect-request message (RFC 5176, *Dynamic Authorization Extensions to Remote Authentication Dial In User Service (RADIUS)*), and tear down the specified remote-access tunnel.

```
config>router>radius-server>server#  
[no] accept-coa
```


For security reasons, the system only accepts a disconnect-request when **accept-coa** is configured and the disconnect-request comes from the corresponding server.

The target tunnel is identified by one of following methods:

- Acct-Session-Id
- Nas-Port-Id + Framed-Ip-Addr(Framed-Ipv6-Prefix) + Alc-IPsec-Serv-Id
- User-Name

See the *7450 ESS, 7750 SR, and VSR RADIUS Attributes Reference Guide* for more details about disconnect message support.

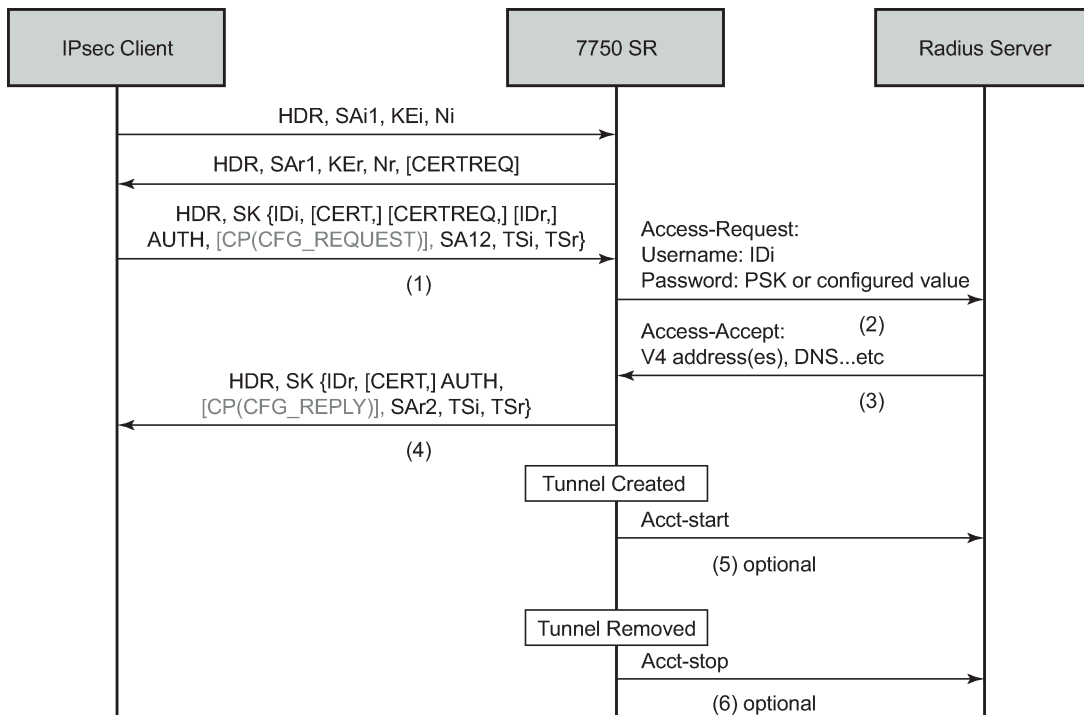
By default, the system only returns what the client has requested in the CFG_REQUEST payload. However, this behavior can be overridden by configuring **relay-unsolicited-cfg-attribute** in the **ike-policy**. With this configuration, the configured attributes returned from the source (such as the RADIUS server) are returned to the client regardless if the client has requested it in the CFG_REQUEST payload.

4.9.1.1 IKEv2 remote-access tunnel – EAP authentication

The SR OS supports EAP authentication for a IKEv2 remote-access tunnel, in which case, the system acts as an authenticator between an IPsec client and a RADIUS server. It transparently forwards EAP messages between the IKEv2 session and RADIUS session. Thus, the actual EAP authentication occurs between the client and the RADIUS server.

Figure 44: Typical call flow of EAP authentication shows a typical call flow of EAP authentication.

Figure 44: Typical call flow of EAP authentication



al_0414

EAP authentication is enabled by configuring **authentication eap**. When enabled, after the received IKE_AUTH request from the client, the system sends an EAP-Response/ID with IDi as the value in the access-request to AAA. AAA returns a method request and the system starts passing through between the client and AAA. (as shown in [Figure 44: Typical call flow of EAP authentication](#)).

The generation of the AUTH payload in the IKE_AUTH response sent by the SR OS (message 4 in flow shown above) is dependent on the **own-auth-method** configuration:

psk

The AUTH payload is present and generated by using PSK.

cert

The AUTH payload is present and generated by the configured public and private key pairs as it does in certificate authentication. Any needed certificates are also sent.

eap-only

Neither AUTH nor CERT payload is present.

The RADIUS attributes in authentication and accounting packets are similar as psk-radius and cert-radius with following differences:

- RADIUS attributes support EAP-Message/Message-Authenticator/State attributes.
- RADIUS attributes support Access-Challenge packet.
- RADIUS attributes support MS-MPPE-Send-Key/ MS-MPPE-Recv-Key in access-accept. These two attributes are required for all EAP methods that generate MSK.

The system provides a method to support EAP and other authentication methods on the same **ipsec-gw** policy. This is enabled by configuring **auto-eap-radius** or **auto-eap** as the **auth-method** in the **ike-policy**.

With **auto-eap-radius**:

- If there is no AUTH payload in an IKE_AUTH request, then the system uses EAP to authenticate the client and also uses **own-auth-method** to generate the AUTH payload.
- If there is an AUTH payload in the IKE_AUTH request, the system uses **auto-eap-own-method** to generate the AUTH payload.
 - If the **auto-eap-method** is **psk**, then the system proceeds as auth-method: psk-radius.
 - If the **auto-eap-method** is **cert**, then the system proceeds as auth-method: cert-radius.
 - If **auto-eap-method** is **psk-or-cert**, then:
 - If the Auth Method field of the AUTH payload is PSK, then the system proceeds as **auth-method:psk-radius**.
 - If the Auth Method field of the AUTH payload is RSA or DSS, then the system proceeds as **auth-method:cert-radius**.

With **auto-eap**:

- If there is no AUTH payload in IKE_AUTH request, then the system uses EAP to authenticate the client and also uses **own-auth-method** to generate AUTH payload.
- If there is an AUTH payload in the IKE_AUTH request:
 - If the **auto-eap-method** is **psk**, then the system proceeds as auth-method: psk.
 - If the **auto-eap-method** is **cert**, then the system proceeds as auth-method: cert.
 - If the **auto-eap-method** is **psk-or-cert**, then:

- If the Auth Method field of the AUTH payload is PSK, then the system proceeds as **auth-method psk**.
- If the Auth Method field of the AUTH payload is RSA or DSS, then the system proceeds as **auth-method cert-auth**.

The system uses auto-eap-own-method to generate the AUTH payload.

4.9.2 IKEv2 remote-access tunnel – authentication without RADIUS

To achieve authentication without RADIUS, the auth-method needs to be configured as psk or cert-auth and local address assignment must be configured under ipsec-gw.

Figure 45: Typical call flow of certificate or PSK authentication without RADIUS shows a typical call flow of certificate or PSK authentication without RADIUS.

Figure 45: Typical call flow of certificate or PSK authentication without RADIUS

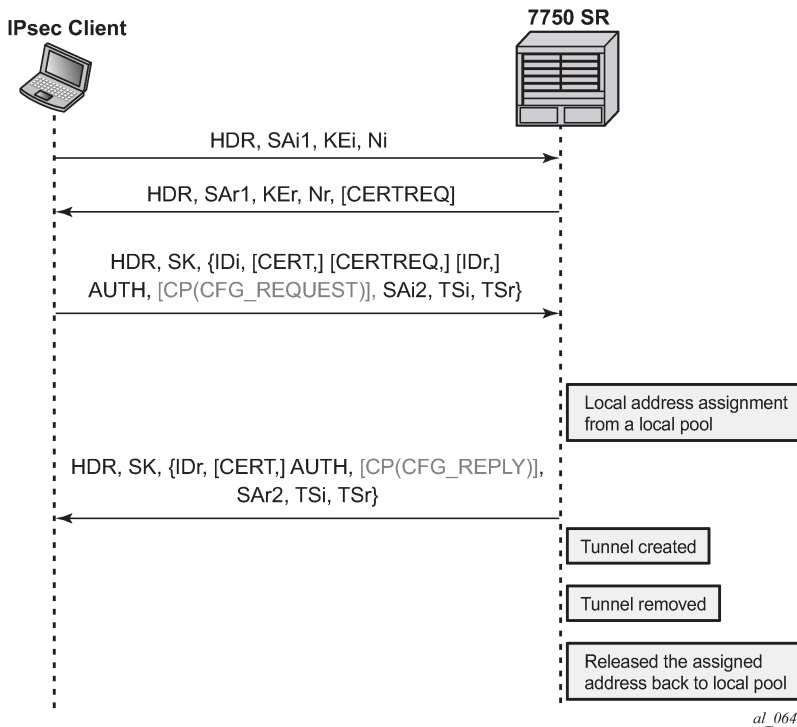
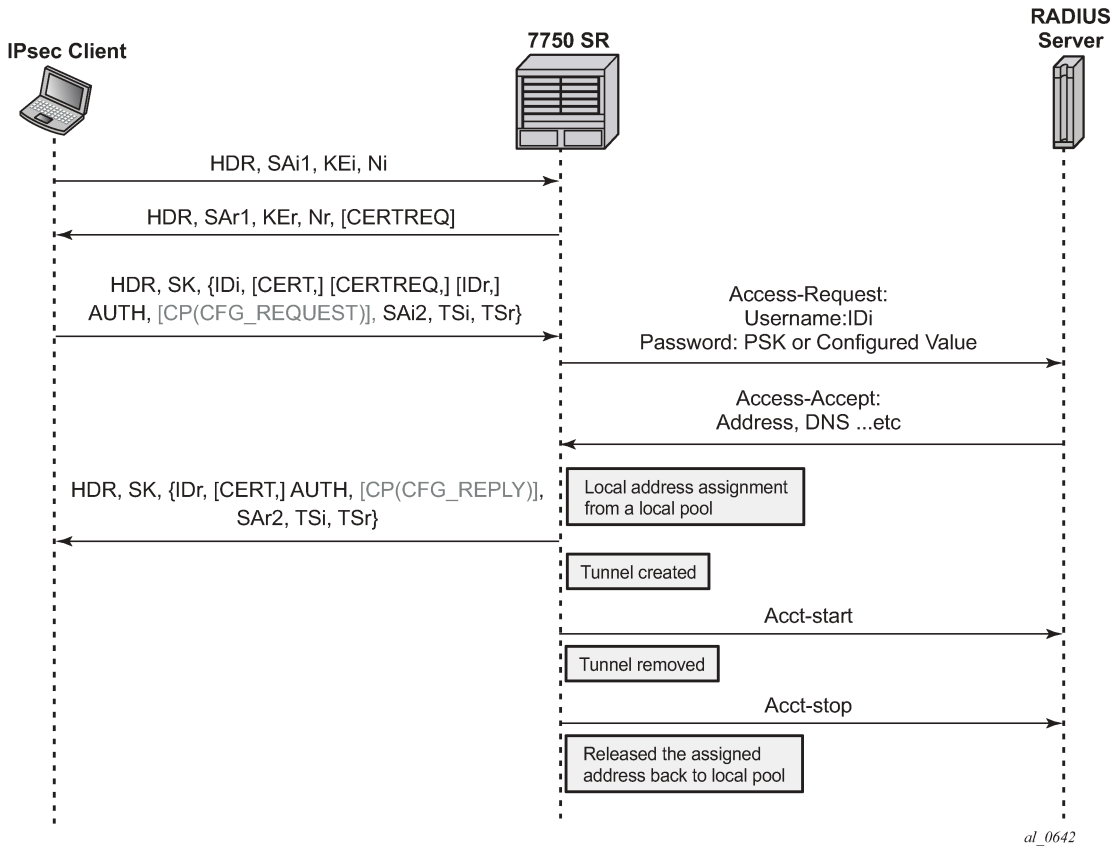


Figure 46: Typical call flow for EAP authentication shows a typical call flow for EAP authentication.

Figure 46: Typical call flow for EAP authentication



In this configuration, the **radius-authentication-policy** and **radius-accounting-policy** in the **ipsec-gw** context are ignored.

RADIUS disconnect messages are supported in this case. Only the following tunnel identification methods are supported:

- Nas-Port-Id + Framed-Ip-Addr(Framed-Ipv6-Prefix) + Alc-IPsec-Serv-Id
- User-Name

4.9.3 IKEv2 remote-access tunnel – address assignment

The SR OS supports the following methods of address assignment for IKEv2 remote-access tunnels:

- RADIUS
- local address assignment (LAA)
- DHCPv4/v6

For RADIUS-based address assignment, the address information is returned in an access-accept packet. This implies that RADIUS-based address assignment requires using an authentication method with RADIUS, such as **psk-radius**, **cert-radius**, or **eap**.

For LAA, the system gets an address from a pool defined in a local DHCPv4/v6 server. When a tunnel is removed, the assigned address is released back to the pool. If the local DHCPv4/v6 server is shut down, all existing tunnels that have an address from the server are removed. If LAA is shut down, the current established tunnel that used LAA stays up.

For DHCP-based address assignment, the system acts as a DHCP client on behalf of the IPsec client and requests an address from an external DHCP server via the standard DHCP exchange. In this case, the system also acts as a DHCP relay agent, which relays all DHCP packets between the DHCP server and the local DHCP client. DHCP renew and rebind are also supported.

4.9.3.1 DHCPv4 address assignment

The client's hardware address field (chaddr) in the DHCPv4 header is generated by the SR OS:

- The first 2 bytes of the MAC address are 02:03.
- The remaining 4 bytes are the hash result of IKEv2 IDi.

The following options are included in the DHCPv4 packets sent by the SR OS:

- Option 82 circuit-id (*private-SAP-id* | *private-interface-name*; for example, tunnel-1.private:100 | priv-int)
- Option 82 remote-id (IKEv2 IDi in text format)
- Option 61 client-id is 1 byte that represents the IKEv2 IDi type plus the IKEv2 IDi in text format. The value of the first byte is as follows:
 - ID_IPV4_ADDR = 1
 - ID_DER_ASN1_DN = 2
 - ID_FQDN = 3
 - ID_RFC822_ADDR = 4
 - ID_IPV6_ADDR = 5

4.9.3.2 DHCPv6 address assignment

Because the system performs a DHCP relay function, all DHCPv6 packets sent or received are encapsulated in DHCPv6 relay-forward and relay-reply messages.

The following items are values of key fields and options in DHCPv6 packets sent by the system:

- Hop-count (0)
- Link address (configurable via the CLI)
- Peer-address (auto-generated based on the IKEv2 IDi)
- Option 1 Client Identifier
 - DUID type (2)
 - Enterprise ID (6527)
 - Value is 1 byte that represents the IKEv2 IDi type plus the IKEv2 IDi in text format. The value of the first byte is the same as that of the first byte in Option 61 for DHCPv4.
- Option 16 Vendor Class
 - Enterprise ID (6527)

- Value (string "SROS IPsec")
- Option 18 Interface ID (*private-SAP-id* | *private-interface-name*; for example, tunnel-1.private:100 | priv-int)
- Option 37 Remote Identifier
 - Enterprise ID (6527)
 - Value (IKEv2 IDi in text format)

4.9.3.3 DHCPv4/v6 usage notes

- Using a local DHCP server on the same chassis for DHCP-based address assignment is not supported. The DHCP server must be external.
- IPsec DHCP Relay uses only the **gi-address** configuration found under the IPsec gateway and does not take into account **gi-address** with **src-ip-addr** configuration below other interfaces.
- The relay-proxy command (**config>service>vprn>if>dhcp>relay-proxy**) must be enabled on an interface that has a gateway IP address as the interface address for the interface to use a DHCPv4 address assignment. The system ignores other DHCP or DHCPv6 configurations on the interface, with the exception of the relay-proxy configuration.
 - If the DHCP server resides in a private service, and the **gi-address** is an address configured on the corresponding tunnel interface, then **relay-proxy** must be enabled on the corresponding private interface.
 - If the DHCP server resides in a routing instance that is different from the private service, then there must be an interface (such as a loopback interface) in the routing instance that has the **gi-address** as the interface address, and **gi-address** must be routable for the DHCP server. Also, **relay-proxy** must be enabled on the interface in the routing instance.

The biggest difference between the LAA and DHCP-based methods is that LAA uses a local API to get an address from a local pool. There is no DHCP packet exchange for LAA, while a DHCP-based method uses standard DHCP packet exchange to request a packet from an external DHCP server.

Because there are three methods for address assignment, the following is the priority order (descending) of sources to choose if more than one source is configured:

- LAA
- DHCP
- RADIUS

There is no fallback between the different sources.

LAA/DHCP can work with an authentication method that does not involve RADIUS, as well as with an authentication method that involves RADIUS. When using LAA/DHCP with an authentication method that involves RADIUS, the following applies:

- LAA/DHCP only happens after RADIUS is successfully authenticated.
- The address information returned by the RADIUS server is ignored (even if LAA/DHCP is configured but is shut down).
- Non-address-related attributes in access-accept messages such as Alc-IPsec-Serv-Id and Alc-IPsec-Tunnel-Template-Id are still accepted.

- RADIUS accounting is supported in this case, but the Framed-IP-Addr/Framed-IPv6-Prefix reported in the acct-request packet is the LAA/DHCP assigned address, not the address returned by the RADIUS server.
- RADIUS disconnect messages are supported.

For MC-IPsec:

- With LAA, the configuration of **config>redundancy>multi-chassis>peer >sync local-dhcp-server** is not needed. This is because the assigned address is synchronized as part of the IPsec tunnel states.
 - Consider the following about DHCP:
 - The DHCP packet exchange process only occurs on the master chassis.
 - The assigned address is synchronized to the standby chassis as part of the IPsec states. The standby chassis does not initiate any DHCP exchanges.
 - The configured DHCP server address (**ipsec-gw>dhcp>server**) should be the same on both chassis.
 - After an MC switchover:
 - The new master does not initiate any DHCP process unless it is time to renew an address or a tunnel goes down.
 - If a new master needs to renew an address or release an address, it sends the DHCP packet to the same DHCP server address that assigned the address on the old master, assuming the external DHCP server is still on, and the renew or release is processed normally.
 - If the new master needs to assign an address for a new tunnel setup, it sends a DHCP discovery or solicit message to all configured DHCP server addresses and then pick the first offer or advertise to finish the DHCP process.
 - For DHCPv4, a gateway IP address is used by the server to forward a response back, so the gateway IP address must be an interface address of the router. For multi-chassis operation, if a DHCPv4 server resides in a private VPRN, there are two options:
 - Configure the same private interface address on both chassis and then use it as the gateway IP address. Configure MC-IPsec-aware routing to make sure that the DHCP response is directed to the master.
 - Configure different private interface addresses with the same subnet on both chassis. The gateway IP address is the private interface address of the local chassis. As well as the private subnet, two /32 private interface address routes from two chassis also need to be advertised so that the DHCP response is routed to the correct chassis.
- If the DHCPv4 server does not reside in a private VPRN, then one method is to configure a loopback interface with a /32 address in the private subnet, and the loopback interface address is used as the gateway IP address. Different addresses must be configured on the master and standby chassis.
- For DHCPv6, unlike DHCPv4, the link-address is not used for the server to forward responses back. The DHCPv6 server sends responses to the source address of the request. This typically is the egress interface address when the system sends out the relay-forward message. For MC-IPsec, no special configuration is required as long as the DHCPv6 server can route relay-reply messages back to the correct chassis.

4.9.4 IPv6 IPsec support

The SR OS provides the following IPv6 support to IPsec functions:

- IPv6 packets as the ESP tunnel payload
- IPv6 as the ESP tunnel encapsulation

4.9.4.1 IPv6 as payload

IPv6 as payload allows IPv6 packets to be forwarded within an IPsec tunnel. Current support includes the following:

- Tunnel type support includes:
 - static LAN-to-LAN tunnel
 - dynamic LAN-to-LAN tunnel
 - remote-access tunnel (only IKEv2 is supported)
- The prefix length of the IPv6 address on a private interface must be /96 or longer.

4.9.4.2 IPv6 as payload: static LAN-to-LAN tunnel

There are three methods to forward IPv6 traffic into static tunnels on the private side:

- The destination address is a configured destination IP (**dest-ip**) under the tunnel context.
 - The **dest-ip** can be either an IPv6 address or an IPv4 address.
 - In the case of IPv6, it must be either an IPv6 global unicast address or an IPv6 link-local address.
 - In the case of IPv4, it can be used to forward IPv4 traffic into the tunnel.
 - In case of unicast address, **dest-ip** must be within the prefix configured on the private interface.
 - Up to 16 destination IPs can be configured per ipsec-tunnel.
- A v6 route with a configured destination IP as the next-hop, this route can be learned from either a static or dynamic from a routing protocol such as BGP.
- An IPv6 static route with an ipsec-tunnel used as the next-hop.

A security policy supports either an IPv4 entry or an IPv6 entry or both for dual-stack.

4.9.4.3 IPv6 as payload: dynamic LAN-to-LAN tunnel

With dynamic LAN-to-LAN tunnels, the system automatically creates a v6 reverse route in the private VPRN based on the received TSi payload with the tunnel as the next hop.

4.9.4.4 IPv6 as payload: remote-access tunnel

The system supports the following IKEv2 IPv6 configuration attributes:

- INTERNAL_IP6_ADDRESS
- INTERNAL_IP6_DNS

The system supports only one internal IPv6 address per tunnel. The following IPv6-related RADIUS attributes are also supported in access-accept:

- Framed-IPv6-Prefix is translated into INTERNAL_IP6_ADDRESS in the configuration payload, which includes two parts. A 16-byte v6 prefix and a one-byte prefix length.
- Alc-Ipv6-Primary-Dns
- Alc-Ipv6-Secondary-Dns

If an internal v6 address has been assigned to the remote-access client, then the Framed-IPv6-Prefix is also included in RADIUS accounting-request packet. The assigned internal v6 address must be within the prefix configured on the corresponding private interface.

If the client request both v4 and v6 address and address source (such as RADIUS or LAA) assign both v4 and v6 address, then both v4 and v6 addresses are assigned to the client via the configuration payload.

4.9.4.5 IPv6 as encapsulation

IPv6 as encapsulation allows IPv4 or IPv6 packets to be forwarded within an IPv6 ESP tunnel, also the IKE protocol can run over IPv6. Current support only includes tunnel type support:

- static LAN-to-LAN tunnel
- dynamic LAN-to-LAN tunnel
- remote-access tunnel (For IKEv1, only v4 over v6 is supported)

For an **ipsec-gw** or **ipsec-tunnel**, only one local gateway address is supported, which could be either an IPv4 or IPv6 address. The SR OS also provides fragmentation and reassembly support for IPv6 ESP/IKE packets.

4.10 MLDv2 over IPsec

The system supports replicating IPv6 multicast traffic into IKEv2 IPsec tunnels based on the MLDv2 report received from the IPsec client.

If a client needs to receive IPv6 multicast traffic over an IPsec tunnel, it includes corresponding multicast address ranges in the traffic selector during CHILD_SA negotiation. After the CHILD_SA is created, the client sends an MLDv2 report to join specific multicast groups over the CHILD_SA. The SeGW terminates the MLDv2 report message and begins replicating requested multicast traffic into the CHILD_SA.

Internally, the system treats each multicast-enabled CHILD_SA as an MLD interface called **ipsec-interface** in various multicast show commands.

This feature supports only IPv6 multicast with MLDv2 and Source Specific Multicast (SSM).

This feature supports only IKEv2 tunnels. For IKEv2 static tunnel, this feature only supports a single child SA per tunnel for multicast traffic.

4.10.1 MLDv2 over IPsec – traffic selector

The negotiated traffic selector of CHILD_SA that transports IPv6 multicast traffic must include the following address ranges:

- an address range for MLDv2 traffic
- an address range for multicast source (IPv6 global unicast address)

- an address range for expected multicast group

The following is an example of a traffic selector of a CHILD_SA that only carries multicast traffic:

- **TSi**
 - FF02::1/128 # destination address of the MLDv2 generic query packet
 - FE80::/64 # IPv6 link local address range, source address of the MLDv2 report packet
 - FF3E::/32 # IPv6 multicast address range, global scope. This is for multicast data traffic and MLDv2 (S,G) specific query
- **TSr**
 - multicast source address range (global unicast)
 - FF02::16/128 # destination address of the MLDv2 report packet
 - FE80::/64 #source address of the MLDv2 query packet

4.10.2 MLDv2 over IPsec – configuration

This feature is enabled by adding private IPsec interface in MLD configuration, for example

```
config>service>vprn
  mld
    interface "priv"
      no shutdown
    exit
  no shutdown
  exit
  interface "priv" tunnel create
    ipv6
      address 2001:dead::1/96
    exit
    sap tunnel-1.private:200 create
  exit
  exit
```

The MLD configuration under the private interface level is ignored by the system.

4.11 Secured interface

A secured interface secures traffic forwarded through a specified IP interface, through one or multiple Secure Interface Tunnels (SI Tunnels) configured under the interface. SI tunnel is conceptually the same as traditional static IPsec tunnels. Some differences are:

- SI tunnels are configured under an IP interface, while static IPsec tunnels are configured under the private tunnel SAP of a tunnel interface.
- With an SI tunnel, the following objects are created automatically with an SI tunnel configuration. There is no need for a separate configuration tunnel configuration:
 - public tunnel SAP
 - public interface
 - private tunnel SAP

- private tunnel interface
- The public service of SI tunnel is the same service of secured interface, which could be either Base router, an IES or an VPRN service.
- The local tunnel address of the SI tunnel must be one of interface addresses of the secure interface. If the secure interface is unnumbered, then it must be one of the interface address of the interface specified by the unnumbered configuration.
- Private service is the same as the public service. The user could also specify a different service.
- On the public side:
 - With a secured interface, by default, all traffic ingress the interface are subject to IPsec processing. If the received traffic is not IPsec traffic (such as ESP and IKE), it is dropped. This behavior can be changed by configuring an **ip-exception** or **ipv6-exception** filter under the interface. All ingress traffic matching the **ip-exception** or **ipv6-exception** filter bypasses IPsec processing and is forwarded through normal routing methods.
 - The system forwards all SI tunnel traffic (after encryption and encapsulation) out through the corresponding secured interface.
 - SSH traffic toward the local system and MPLS/SDP always bypasses IPsec processing.
- On the private side:
 - Like a static IPsec tunnel, traffic is routed into the SI tunnel through a static route or BGP route.
 - When an SI tunnel is operationally down, routes using the next-hop address as the tunnel are unresolved and withdrawn from the route table.
- show, debug, tool, clear, and admin commands that apply to static IPsec tunnels also apply to SI tunnels.
- The following features are not supported with SI tunnels on 7705 SAR-Hm (with cellular exit port):
 - Dest-ip
 - MC-IPsec
 - IPv4 over IPv6
 - IPv6 over IPv6
 - MLDv2 over SI tunnel
- The following features are not supported with SI tunnels on VSR:
 - Dest-ip
 - MC-IPsec
 - MLDv2 over SI tunnel
- SI tunnel are only supported in VSR and 7705 SAR-Hm families.

4.12 IPsec client database

The IPsec client database is a database configured in the (**config>ipsec>client-db**) CLI context, which can be used to authenticate and authorize IKEv2 dynamic LAN-to-LAN tunnels.

Each client database contains one or more client entries. When the system receives a new tunnel request, it performs a look up in the associated database of the IPsec gateway (**ipsec-gw**). If there is match, the

system optionally could use credentials configured in the matched client entry to authenticate the peer. If the authentication succeeds, then, optionally, the matched entry could also return certain IPsec parameters such as the private service ID which can be used for tunnel setup.

If the client database lookup failed to return a match result, then the system can either fall back to the **ipsec-gw** level configuration or fail the tunnel setup. The action to take depends on the CLI configuration.

The system supports one of the following as matching input:

- the peer's tunnel IP address
- the peer's IDi
- a combination of both

The above matching input is defined in the **match-list** context in the client-db configuration. Each client entry contains client matching criteria that corresponds to the match list. The system correlates matching input with the client matching criteria of each client entry in the client-db configuration. The system supports the following matching methods:

- **for the peer's IDi**
 - Any matches any IDi.
 - IPv4/IPv6 prefix matches the peer's address type IDi to a configured prefix. It is considered a match if the IDi falls within the prefix.
 - FQDN matches the peer's FQDN type IDi to a string. This supports a complete string match or a suffix string match.
 - RFC822 matches the peer's RFC 822 type IDi to a string. This supports a complete string match or suffix string match.
- **for the peer's tunnel IP address**
 - Matches the peer's tunnel address to a configured prefix. It is a match if the IDi fall within the prefix.
 - IPv4 Any matches any IPv4 address.
 - IPv6 Any matches any IPv6 address.

Each client entry has a client index (an integer). This is different from a client identification. If there are multiple matched entries in a lookup, the client entry with the smallest client index is used. The client entry supports using a pre-shared key as the credential.

If the credential is not configured in the matched entry, the credential configured under the ipsec-gw context is used.

A client entry could optionally return the following IPsec parameters:

- a private service ID
- a private interface name
- a tunnel-template ID
- a Ts list

The returned parameter overrides the configuration of the ipsec-gw level.

There is only one client-db for each ipsec-gw, but different ipsec-gw configurations can use the same client-db.

Note that the **encapsulated-ip-mtu** command in the client-db returned tunnel-template is not applied to the IKE packet fragmentation. The **encapsulated-ip-mtu** command configured in the **config>ipsec>tunnel-**

template context is used instead. However, the client-db returned encapsulated IP MTU value still applies to the ESP packet fragmentation.



Note:

- A client entry in a shutdown state is skipped while the system performs the matching process.
- If the configuration returned by client-db is invalid, the system fails the tunnel setup.
- The reference of the client-db under the ipsec-gw context can be changed without shutting down ipsec-gw
- Shutting down a referenced client-db without shutting down ipsec-gw is allowed and the established tunnel is not impacted. The system uses the configuration on the ipsec-gw level for new a tunnel request while the client-db is shutdown if a fallback is configured.
- Adding a new client in a referenced client-db without shutting down ipsec-gw or client-db is allowed.
- Removing a client in the referenced client-db without shutting down ipsec-gw or client-db is allowed. However, the shutdown of the client to be removed is required.
- Changing an existing client of a referenced client-db without shutting down ipsec-gw or client-db is allowed. However, the shutdown of the client to be removed is required.

4.13 IPsec transport mode protected IP tunnel

Tunnel-group based GRE tunnels can be protected by applying IPsec transport mode encryption for the GRE tunnel packets. This is achieved by configuring an IPsec transport mode profile under the IP tunnel configuration. When the profile is enabled, the data path flow as follows in the private to public direction:

- The payload packet is received on the private tunnel SAP.
- Optionally perform pre-encapsulation fragmentation based on the payload packet size and **ip-mtu** configuration under the **ip-tunnel** context
- The payload packet is encapsulated into a GRE tunnel packet.
- The IPsec transport mode encryption is applied on the GRE tunnel packet which results in an IPsec ESP packet.
- Optionally perform post-encapsulation fragmentation based on the ESP packet size and the configured **encapsulated-ip-mtu** under the **ip-tunnel** context

In the public to private direction, the data path flows as follows:

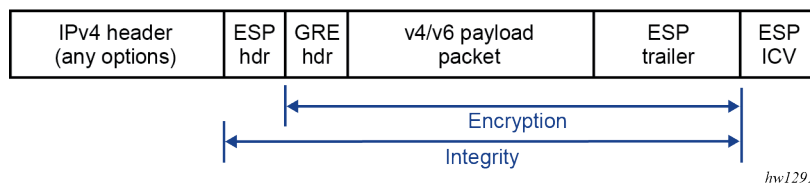
1. The IPsec ESP packet is received on the public tunnel SAP.
2. ESP packet reassembly is performed if it is fragmented.
3. The ESP packet is decrypted, which results as a GRE tunnel packet.
4. The GRE tunnel packet is decapsulated and the payload packet is forwarded out of the private tunnel SAP.

This feature uses IKEv2 to create an IKE_SA and a transport mode CHILD_SA for a specific GRE tunnel. The IKE/IPsec behaves similarly to an IPsec static LAN-to-LAN tunnel, with some transport-mode specific differences.

4.13.1 Packet format

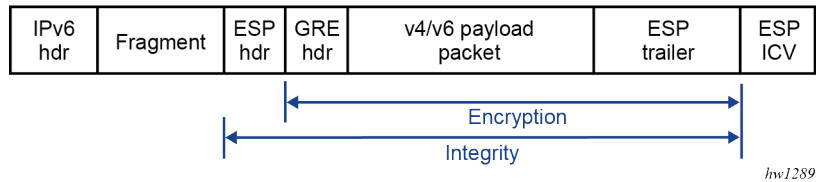
The packet format of GRE with IPsec transport mode follows RFC 4303, *IP Encapsulating Security Payload (ESP)*, Section 3.1.1. The following example shows an IPv4 GRE packet with IPsec transport mode enabled:

Figure 47: IPv4 GRE packet with IPsec transport mode enabled



The following example shows an IPv6 GRE packet with IPsec transport mode enabled:

Figure 48: IPv6 GRE packet with IPsec transport mode enabled



4.13.2 IKEv2 and SA

This feature only requires a single IKE_SA and CHILD_SA per GRE tunnel. Additional CHILD_SA creation requests from a CREATE_CHILD_SA peer exchange are rejected.

This feature uses IKEv2 USE_TRANSPORT_MODE notification to signal the use of the transport mode when the CHILD_SA is created. The system expects the notification to be included in both the request and response messages.

4.13.3 Traffic selector

As a CHILD_SA initiator, the traffic selector is proposed as follows:

- address range (the source and destination GRE tunnel endpoint address, /32 or /128)
- protocol (GRE)
- port range (0 to 65535)

As a CHILD_SA responder, the system uses the same traffic selector to narrow the selection peer's proposal.

4.13.4 Fragmentation and reassembly

ISA only fragments packets in the private to public direction. The following two types of fragmentation are used.

- Fragment before GRE encapsulation is controlled by the **ip-mtu** in the **ip-tunnel** context.
- Fragment after IPsec processing is controlled by the configured **encapsulated-ip-mtu** in the **ip-tunnel** context.

ISA only reassembles received ESP packets on the public side before IPsec decryption. The reassembly behavior is controlled by the **reassembly** command under the tunnel group.

4.13.5 QoS marking

QoS marking refers to the Type of Service field in IPv4 header or the Traffic Class field in the IPv6 header.

- **behavior in the private to the public direction**

If **dscp dscp-name** is configured under the IP tunnel, the configured value is used for the QoS marking in the GRE IP header. If DSCP is not configured, then the payload packet's QoS marking is copied into the GRE IP header.

- **behavior in the public to private direction**

The received ESP packet is decrypted, and as a result, the GRE packet is decapsulated into a payload packet. The QoS marking of the original payload packet is preserved.

4.13.6 Configuring IPsec transport mode protected tunnels

For IPsec transport mode protected tunnels, include the following:

- a GRE tunnel
- IPsec parameters:
 - **ike-policy** *ike-policy-id*
 - **ike-transform** *ike-transform-id*
 - **ipsec-transform** *transform-id*
 - **cert-profile** *profile-name* or **trust-anchor-profile** *name* (the certificate authentication is required)
- an **ipsec-transport-mode-profile** *name* referenced under the GRE tunnel

The following are Classic CLI configuration examples:

```
A:v70>config>ipsec# info
-----
    ike-transform 1 create
      dh-group 20
      ike-auth-algorithm auth-encryption
      ike-encryption-algorithm aes256-gcm16
      ike-prf-algorithm sha384
    exit
    ike-policy 1 create
      ike-version 2
      ike-transform 1
    exit
    ipsec-transform 1 create
      esp-auth-algorithm auth-encryption
      esp-encryption-algorithm aes256-gcm16
      pfs-dh-group 20
    exit
    ipsec-transport-mode-profile "test" create
      dynamic-keying
        ike-policy 1
        pre-shared-key "KrbVPnF6Dg13PM/biw6ErPl5XU7+" hash2
        transform 1
      exit
    exit

A:v70>config>service>vprn# info
-----
    interface "priv" tunnel create
      address 44.44.44.1/24
      sap tunnel-1.private:100 create
        ip-tunnel "t1" create
          dest-ip 44.44.44.2
          gre-header
          source 172.16.100.1
          remote-ip 192.168.1.2
          delivery-service 300
          ipsec-transport-mode-profile "test"
```



```

                no shutdown
            exit
        exit
    exit

```

4.14 Configuring IPsec with CLI

4.14.1 Provisioning a tunnel ISA

A tunnel ISA can only be provisioned on an FP2- or FP3-based IOM. The following output displays a card and ISA configuration.

```

*A:ALA-49>config# info
-----
...
  card 1
    card-type iom4-e
    mda 1
      mda-type me40-1gb-csfp
      no shutdown
    exit
    mda 2
      mda-type isa2-tunnel
      no shutdown
    exit
  no shutdown
...
-----
*A:ALA-49>config#

```

4.14.2 Configuring a tunnel group

The following output displays a tunnel group configuration in the ISA context. The **multi-active** command specifies that there could be multiple active ISAs in the tunnel group, the **mda** command specifies the MDA ID of the ISA in the tunnel group. There could be multiple MDA commands in the tunnel group.

```

*A:ALA-49>config# info
-----
...
  isa
    tunnel-group 1 create
      multi-active
      mda 1/2
      no shutdown
    exit
  exit
...
-----
*A:ALA-49>config#

```

4.14.3 Configuring router interfaces for IPsec

The following output displays an interface "internet" configured using the network port (1/1/1), which provides network connection on the public side.

```
*A:ALA-49>config# info
-----
...
router
  interface "internet"
    address 10.10.7.118/24
    port 1/1/1
  exit
  interface "system"
    address 10.20.1.118/32
  exit
  autonomous-system 123
exit
...
-----
*A:ALA-49>config#
```

4.14.4 Configuring IPsec parameters

The following output displays an IPsec configuration example.

```
config>ipsec
  ike-transform 100 create
    dh-group 14
    ike-auth-algorithm sha256
    ike-encryption-algorithm aes128
    isakmp-lifetime 90000
  exit
  ike-policy 100 create
    ike-version 2
auth-method psk
  ike-transform 100
exit
```

4.14.5 Configuring IPsec in services

The following output displays an IES and VPRN service with IPsec parameters configured.

```
*A:ALA-49>config# info
-----
...
service
  ies 100 customer 1 create
    interface "ipsec-public" create
      address 10.10.10.1/24
      sap tunnel-1.public:1 create
    exit
  exit
  no shutdown
exit
```

```

vprn 200 customer 1 create
  ipsec
    security-policy 1 create
      entry 1 create
        local-ip 172.16.118.0/24
        remote-ip 172.16.91.0/24
      exit
    exit
  exit
route-distinguisher 1:1
interface "ipsec-private" tunnel create
  sap tunnel-1.private:1 create
  ipsec-tunnel "remote-office" create
    security-policy 1
    local-gateway-address 10.10.10.118 peer 10.10.7.91 delivery-service 100
    dynamic-keying
      ike-policy 1
      pre-shared-key "humptydumpty"
      transform 1
    exit
  no shutdown
  exit
exit
interface "corporate-network" create
  address 172.16.118.118/24
  sap 1/1/2 create
  exit
exit
static-route-entry 172.16.91.0/24
  ipsec-tunnel "t1"
  no shutdown
  exit
exit
no shutdown
exit
exit
...
-----
*A:ALA-49>config#

```

4.14.6 Configuring X.509v3 certificate parameters

The following are steps to configure certificate enrollment:

1. Generate a key.

```
admin certificate gen-keypair cf3:/key_plain_rsa2048 size 2048 type rsa
```

2. Generate a certificate request.

```
admin certificate gen-local-cert-req keypair cf3:/key_plain_rsa2048 subject-
rdn "C=US,ST=CA,CN=7750" file 7750_req.cs
```

3. Send the certificate request to CA-1 to sign and get the signed certificate.

4. Import the key.

```
admin certificate import type key input cf3:/
```

```
key_plain_rsa2048 output          key1_rsa2048 format der
```

5. Import the signed certificate.

```
admin certificate import type cert input cf3:/
7750_cert.pem output 7750cert          format pem
```

The following are steps to configure CA certificate/CRL import.

1. Import the CA certificate.

```
admin certificate import type cert input cf3:/
CA_1_cert.pem output ca_cert          format pem
```

2. Import the CA's CRL.

```
admin certificate import type crl input cf3:/
CA_1_crl.pem output ca_crl format          pem
```

The following displays a certificate authentication for IKEv2 static LAN-to-LAN tunnel configuration.

```
config>system>security>pki# info
-----
          ca-profile "alu-root" create
          cert-file "alu_root.cert"
          crl-file "alu_root.crl"
          no shutdown
          exit
-----
config>ipsec# info
-----
          ike-policy 1 create
          ike-version 2
          auth-method cert-auth
          ike-transform 1
          exit
          ipsec-transform 1 create
          exit
          ike-transform 1 create
          exit
          cert-profile "segw" create
          entry 1 create
          cert segw.cert
          key segw.key
          exit
          no shutdown
          exit
          trust-anchor-profile "nokia" create
          trust-anchor "nokia-root"
          exit

config>service>vpn>if>sap
-----
          ipsec-tunnel "t50" create
          security-policy 1
          local-gateway-
address 192.168.55.30 peer 192.168.33.100 delivery-service 300
          dynamic-keying
          ike-policy 1
          transform 1
          cert
          trust-anchor-profile "nokia"
```

```

                                cert-profile "segw"
                                exit
                                exit
                                no shutdown
                                exit

```

The following displays an example of the syntax to import a certificate from the pem format.

```

*A:SR-7/Dut-A# admin certificate import type cert input cf3:/pre-import/R1-
0cert.pem output R1-0cert.der format pem

```

The following displays and example of the syntax to export a certificate to the pem format.

```

*A:SR-7/Dut-A# admin certificate export type cert input R1-0cert.der output cf3:/
R1-0cert.pem format pem

```

4.14.7 Configuring MC-IPsec

4.14.7.1 Configuring MIMP

The following is an MIMP configuration example.

```

config>redundancy>multi-chassis
-----
    peer 10.2.2.2 create
        mc-ipsec
            bfd-enable
            tunnel-group 1 create
                peer-group 2
                priority 120
                no shutdown
            exit
        exit
    no shutdown
exit

```

The peer's tunnel-group ID is not necessarily the same as the local tunnel-group ID. With **bfd-enable**, the BFD parameters are specified under the interface that the MIMP source address resides on, which must be a loopback interface in the base routing instance. The default source address of MIMP is the system address.

The **keep-alive-interval** and **hold-on-neighbor-failure** define the MIMP alive parameter, however, BFD could be used for faster chassis failure detection.

The SR OS also provides a **tools** command to manually trigger the switchover, for example:

```

tools perform redundancy multi-chassis mc-ipsec force-switchover tunnel-group 1

```

4.14.7.2 Configuring multi-chassis synchronization

The following displays an MCS for MC-IPsec configuration.

```

config>redundancy>multi-chassis>

```

```

-----
peer 10.2.2.2 create
  sync
  ipsec
  tunnel-group 1 sync-tag "sync_tag_1" create
  no shutdown
exit

```

The following displays an MCS for MC-IPsec configuration with MCS encryption enabled.

```

config>redundancy>multi-chassis
peer 10.20.1.2 create
  source-address 10.20.1.1
  sync
  ipsec
  tunnel-group 1 sync-tag "ipsec-t-grp-1" create
  transport-encryption
    application "ipsec" keychain "validMultiBi"
  exit
  exit
  no shutdown
exit
mc-ipsec
  bfd-enable
  tunnel-group 1 create
    peer-group 1
    priority 255
    no shutdown
  exit
  exit
  no shutdown
exit
exit

```

The **sync-tag** must match on both chassis for the corresponding tunnel groups.

4.14.7.3 Configuring routing for MC-IPsec

The following configuration is an example using a route policy to export /32 local tunnel address route:

```

config>router>policy-options>
-----
policy-statement "exportOSPF"
  entry 10
    from
      protocol ipsec
      state ipsec-master-with-peer
    exit
    action accept
      metric set 500
    exit
  exit
  entry 20
    from
      protocol ipsec
      state ipsec-non-master
    exit
    action accept
      metric set 1000

```

```

        exit
    exit
    entry 30
    from
        protocol ipsec
        state ipsec-master-without-peer
    exit
    action accept
        metric set 1000
    exit
exit
exit

```

The following configuration shows shunting in public and private services.

Shunting in public service:

```

config>service>ies>
    interface "ipsec-pub" create
        address 172.16.100.254/24
        sap tunnel-1.public:100 create
    exit
    static-tunnel-redundant-next-hop 10.1.1.1
exit

```

Shunting in private service:

```

config>service>vprn>
    interface "ipsec-priv" tunnel create
        ...
        static-tunnel-redundant-next-hop 10.7.7.1
    exit

```

Shunting is enabled by configuring redundant next-hop on a public or private IPsec interface

static-tunnel-redundant-next-hop

shunting nexthop for a static tunnel

dynamic-tunnel-redundant-next-hop

shunting next-hop for a dynamic tunnel

4.14.7.4 Configuring MCS encryption

The following examples show the configuration for MCS encryption.

Keychain example 1

```

*A:vsim4>config>redundancy>multi-chassis>peer>sync>transport-encryption#
-----
    application "ipsec" keychain "mcs"
    exit
-----
config>system>security# info
-----
    keychain "mcs"
        direction
            bi
            entry 2 key "KrbVPnF6Dg13PM/biw6ErIQyX4uD" hash2 algorithm aes-128-gcm-16
            begin-time 2019/03/22 18:50:00 UTC

```

```

    exit
  exit
  exit
  no shutdown
  exit

```

Keychain example 2

```

aes*A:vsim4>config>redundancy>multi-chassis>peer>sync>transport-encryption#
-----
  application "ipsec" keychain "mcs"
  exit

```

```

*A:vsim4>config>system>security# info
-----
  per-peer-queuing
  keychain "mcs"
  direction
    bi
    entry 2 key "KrbVPnF6Dg13PM/biw6ErIQyX4uD" hash2 algorithm aes-128-gcm-16
      begin-time 2019/03/22 18:50:00 UTC
    exit
    entry 3 key "/Ub6BWHC4DsEprLWutGaTcJz1zDPhw==" hash2 algorithm aes-128-gcm-16
      begin-time 2019/03/22 21:30:00 UTC
    exit
    entry null-key
      begin-time 2019/03/22 18:33:23 UTC
      tolerance forever
    exit
  exit
  exit
  no shutdown
  exit

```

4.14.8 Configuring and using CMPv2

CMPv2 server information is configured under the corresponding **ca-profile** using the following key commands:

```

config>system>security>pki>ca-profile
  cmpv2
    url <url-string> [service-id <service-id>]
    response-signing-cert <filename>
    key-list
      key <password> reference <reference-number>

```

The **url** command specifies the HTTP URL of the CMPv2 server, the service specifies the routing instance that the system used to access the CMPv2 server (if omitted, then system uses base routing instance).

The service ID is only needed for inband connections to the server via VPRN services. IES services are not to be referenced by the service ID as any of those are considered base routing instance.

The **response-signing-cert** command specifies a imported certificate that is used to verify the CMP response message if they are protected by signature. If this command is not configured, then CA's certificate is used.

The **keylist** specifies a list of pre-shared-key used for CMPv2 initial registration message protection.

For example:

```
config>system>security>pki>ca-profile>
  cmpv2
    url "http://cmp.example.com/request" service-id 100
    key-list
      key passwordToBeUsed reference "1"
```

All CMPv2 operations are invoked by using the **admin certificate cmpv2** command.

If there is no **key-list** defined under the **cmpv2** configuration, the system defaults to the **cmpv2** transaction input for the command line for authenticating a message without a sender ID. Also, if there is no sender ID in the response message, and there IS a key-list defined, it chooses the lexicographical first entry only, if that fails, it has a fail result for the transaction.

See the command reference section for details about syntax and usage. The system supports optional commands (such as, **always-set-sender-ir**) to support inter-op with CMPv2 servers.

4.14.9 Configuring OCSP

OCSP server information is configured under the corresponding ca-profile:

```
config>system>security>pki>ca-profile>
  oosp
    responder-url <url-string>
    service <service-id>
```

The **responder-url** command specifies the HTTP URL of the OCSP responder. The **service** command specifies the routing instance that system used to access the OCSP responder.

Example:

```
config>system>security>pki>ca-profile>
  oosp
    responder-url "http://ocsp.example.com/request"
    service 100
```

For an ipsec-tunnel or ipsec-gw, the user can configure a primary method, a secondary method and a default result.

```
config>service>ies>if>sap>ipsec-gw>
config>service>vprn>if>sap>ipsec-gw>
config>service>vprn>if>sap>ipsec-tun>dynamic-keying
  cert
    status-verify
      primary {ocsp | crl} secondary {ocsp | crl}
      default-result {revoked | good}
```

Example:

```
config>service>ies>if>sap>ipsec-gw>dynamic-keying
  cert
    status-verify
      primary ocsp secondary crl
```

4.14.10 Configuring IKEv2 remote — access tunnel

The following are configuration tasks for an IKEv2 remote-access tunnel:

- Create an ike-policy with one of the auth-methods that enabled the remote-access tunnel.
- Configure a tunnel-template/ipsec-transform. This is the same as configuring a dynamic LAN-to-LAN tunnel.
- Create a radius-authentication-policy and optionally, a radius-accounting-policy (a radius-server-policy and a radius-server must be preconfigured).
- Configure a private VPRN service and private tunnel interface with an address on the interface. The internal address assigned to the client must come from the subnet on the private interface.
- Configure a public IES/VPRN service and an ipsec-gw under the public tunnel SAP.
- Configure the radius-authentication-policy and radius-accounting-policy (optional) under the ipsec-gw.
- Certificate the related configuration if any certificate related authentication method is used.

The following shows an example using cert-radius:

```

config>system>security>pki# info
-----
      ca-profile "NOKIA-ROOT" create
      cert-file "NOKIA-ROOT.cert"
      crl-file "NOKIA-ROOT.crl"
      no shutdown
      exit
-----
A:SeGW>config>aaa# info
-----
      radius-server-policy "femto-aaa" create
      servers
      router "management"
      server 1 name "svr-1"
      exit
      exit
-----
A:SeGW>config>router# info
-----
      radius-server
      server "svr-
1" address 10.10.10.1 secret "KR35xB3W4aUXtL8o3WzPD." hash2 create
      exit
      exit
-----

config>ipsec# info
-----
      ike-policy 1 create
      ike-version 2
      auth-method cert-radius
      ike-transform 1
      exit
      ipsec-transform 1 create
      exit
      ike-transform 1 create
      exit
      tunnel-template 1 create
      transform 1
      exit

```

```

cert-profile "c1" create
  entry 1 create
    cert SeGW2.cert
    key SeGW2.key
  exit
  no shutdown
exit
trust-anchor-profile "tap-1" create
  trust-anchor "NOKIA-ROOT"
exit
radius-authentication-policy "femto-auth" create
  include-radius-attribute
    calling-station-id
    called-station-id
  exit
  password "DJzlyYKCefyhmnFcFSBuLZovSemMKde" hash2
  radius-server-policy "femto-aaa"
exit
radius-accounting-policy "femto-acct" create
  include-radius-attribute
    calling-station-id
    framed-ip-addr
  exit
  radius-server-policy "femto-aaa"
exit
-----
config>service>ies# info
-----
      interface "pub" create
        address 172.16.100.0/31
        tos-marking-state untrusted
        sap tunnel-1.public:100 create
          ipsec-gw "rw"
            cert
              trust-anchor-profile "tap-1"
              cert-profile "c1"
            exit
          default-secure-service 400 interface "priv"
          default-tunnel-template 1
          ike-policy 1
          local-gateway-address 172.16.100.1
          radius-accounting-policy "femto-acct"
          radius-authentication-policy "femto-auth"
          no shutdown
        exit
      exit
    exit
  no shutdown
-----
A:SeGW>config>service>vprn# info
-----
      route-distinguisher 400:11
      interface "priv" tunnel create
        address 10.20.20.1/24
        sap tunnel-1.private:200 create
        exit
      exit
      interface "l1" create
        address 10.9.9.9/32
        loopback
      exit
    no shutdown
-----

```

4.14.11 Configuring IKEv2 remote — access tunnel with local address assignment

The following are configuration tasks of IKEv2 remote-access tunnel:

- Create an **ike-policy** with any **auth-method**.
- Configure the **tunnel-template** or **ipsec-transform**. This is the same as configuring a dynamic LAN-to-LAN tunnel.
- Configure a private VPRN service and a private tunnel interface with an address on the interface. The internal address assigned to the client must come from the subnet on the private interface.
- Configure a local DHCPv4 or DHCPv6 server with address pool that from which the internal address to be assigned from.
- Configure public IES/VPRN service and **ipsec-gw** under public tunnel SAP.
- Configure the local address assignment under **ipsec-gw**.

The following output shows an example using cert-auth:

```

config>system>security>pki# info
-----
      ca-profile "smallcell-root" create
      cert-file "smallcell-root-ca.cert"
      revocation-check crl-optional
      no shutdown
      exit
-----
config>ipsec# info
-----
      ike-policy 3 create
      ike-version 2
      auth-method cert-auth
      nat-traversal
      ike-transform 1
      exit
      ipsec-transform 1 create
      exit
      ike-transform 1 create
      exit
      cert-profile "segw-mlab" create
      entry 1 create
      cert SeGW-MLAB.cert
      key SeGW-MLAB.key
      exit
      no shutdown
      exit
      trust-anchor-profile "sc-root" create
      trust-anchor "smallcell-root"
      exit
      tunnel-template 1 create
      transform 1
      exit
-----
config>service>ies# info
-----
      interface "pub" create
      address 172.16.100.253/24
      tos-marking-state untrusted
      sap tunnel-1.public:100 create
      ipsec-gw "rw"
      default-secure-service 400 interface "priv"

```

```

default-tunnel-template 1
ike-policy 3
local-address-assignment
  ipv6
    address-source router 400 dhcp-server "d6" pool "1"
  exit
  no shutdown
exit
local-gateway-address 172.16.100.1
cert
  trust-anchor-profile "sc-root"
  cert-profile "segw-mlab"
  status-verify
    default-result good
  exit
exit
local-id type fqdn value segwmobilelab.nokia.com
no shutdown
exit
exit
no shutdown
-----
config>service>vprn# info
-----
  dhcp6
    local-dhcp-server "d6" create
    use-pool-from-client
    pool "1" create
    options
      dns-server 2001:db8:::808:808
    exit
    exclude-prefix 2001:db8:beef::101/128
    prefix 2001:db8::beef::/96 failover access-driven pd wan-host create
    exit
  exit
  no shutdown
  exit
exit
route-distinguisher 400:1
interface "priv" tunnel create
  ipv6
    address 2001:db8::beef::101/96
  exit
  sap tunnel-1.private:200 create
  exit
exit
no shutdown

```

4.14.12 Configuring secured interfaces

The following is an example config for secured interface. In this example, a SI tunnel "t1" is configured under interface "toPeer-1" in Base routing instance, along with an exception filter 100 that allows OSPF packets bypass IPsec processing:

```

config>filter# info
-----
  ip-exception 100 create
  entry 10 create
  match protocol ospf-igp
  exit

```

```
        exit
    exit
-----
config>router# info
-----
#-----
echo "IPsec Configuration"
#-----
    ipsec
        security-policy 1 create
            entry 1 create
                local-ip 100.0.0.20/32
                remote-ip 200.1.1.254/32
            exit
        exit
    exit
#-----
echo "IP Configuration"
#-----
    interface "toPeer-1"
        address 192.168.110.20/24
        port 1/1/3
        ipsec tunnel-group 1 public-sap 300
            ip-exception 100
            ipsec-tunnel "t1" private-sap 300 create
                local-gateway-address 192.168.110.20
                remote-gateway-address 172.16.21.1
                security-policy 1
                dynamic-keying
                    ike-policy 3
                    pre-shared-key "KrbVPnF6Dg13PM/biw6ErD9+g6HZ" hash2
                transform 2
            exit
    exit
```

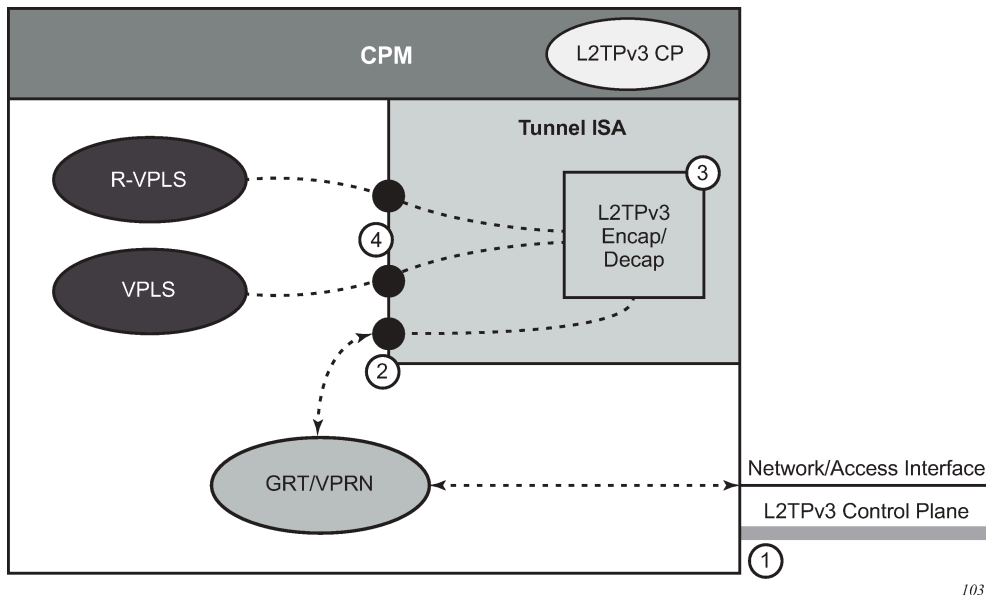
5 L2TPV3 tunnels

5.1 L2TPv3 overview

Layer 2 Tunneling Protocol version 3 (L2TPv3) is a mechanism for the tunneling of Ethernet traffic over an IP network. For this application, the ISA functions as a resource module for the system, performing the L2TPv3 encapsulation and decapsulation functions.

Figure 49: L2TPv3 support for IP transport shows L2TPv3 support for the IP transport model. Table 34: L2TPv3 support for IP transport — tunnel processing steps describes the tunnel processing steps in the figure.

Figure 49: L2TPv3 support for IP transport



1037

Table 34: L2TPv3 support for IP transport — tunnel processing steps

Step number	Description
1	The L2TPv3 control plane can run within either the base routing or VPRN contexts.
2	L2TPv3 encapsulated packets ingress and egress through the public interface, which can be in either the base routing or VPRN contexts.
3	L2TPv3 encapsulation and decapsulation processing is handled within the tunnel ISA.

Step number	Description
4	Unencapsulated packets pass between the tunnel ISA and the associated service via the configured private SAP.

5.2 Control plane

The configuration of the L2TPv3 control plane is similar to that of L2TPv2. A number of the same commands are used for both, but there are new commands specific to L2TPv3. The L2TPv3-specific commands are located in a separate L2TPv3 context in both the general configuration area as well as within the group configuration context.

L2TPv3 control plane parameters can be configured at the global level within the **config>router>l2tp** context, which may include some L2TPv3-specific parameters. This should be used for parameters that are the same for the majority of L2TPv3 tunnels. The same parameters can be configured on a per-tunnel group basis. The tunnel group can be configured within either the base router context or a VPRN service context.

The following example displays an L2TPv3 tunnel group configured within the base routing context:

```
configure
  router
    l2tp
      l2tpv3
        cookie-length 8
        digest-type sha1
        nonce-length 64
        transport-type ip
        exit
      group "base l2tpv3 left" protocol v3draft create
        avp-hiding never
        eth-tunnel
          reconnect-timeout 60
        exit
      l2tpv3
        pw-cap-list ethernet ethernet-vlan
        password "AbkdpF.rY1FgcK4qAYmim sykdmwbAucq" hash2
        exit
        password "rhXAlJTUjuliBn8lVUf KJywztX9cK0Eb/rbWUR/e4ow" hash2
        tunnel "base l2tpv3 tunnel" create
          local-address 172.16.0.100
          peer 192.168.0.100
          no shutdown
        exit
      no shutdown
    exit
```

5.3 Public SAP

The public SAP is the access interface to the L2TPv3 tunnel over which encapsulated traffic is sent to or received from the far end. The IP address bound to this SAP is on the same subnet as the local L2TPv3 tunnel endpoint.

The public SAP must be configured in the same routing context as the L2TPv3 tunnel group configuration. As shown in [Figure 49: L2TPv3 support for IP transport](#), the public SAP can be associated with an IES or VPRN service to connect to the outside or public access network.

The following example displays an L2TPv3 public SAP configured within the base routing context:

```
configure
  service
    ies 10
      interface "l2tp-public-interface" create
        address 172.16.0.1/24
        sap tunnel-1.public:2 create
      exit
    exit
```

5.4 Private SAP

The private SAP is the access interface to the L2TPv3 over which unencapsulated traffic is sent to or received from the far end. The public SAP must be configured within an Ethernet service, such as an Epipe, VPLS, or I-VPLS service.

The private SAP configuration includes the configuration of the following L2TPv3 session parameters:

- VC-ID
- PW-Type
- L2TPv3 tunnel group association

The following example displays an L2TPv3 private SAP configured within the base routing context:

```
configure
  service
    vpls 100 customer 1 create
    ...
    sap tunnel-1.private:100 create
      l2tpv3-session create
        router 2 group "base l2tpv3 left"
        vc-id 100
        pw-type ethernet
        no shutdown
      exit
    no shutdown
  exit
```

6 Video services

6.1 Video services

6.1.1 Video groups

Video services are supported on:

- MS-ISA
- MS-ISA2
- MS-ISM (two MS-ISA2 on a single line card)
- ESA

These MS-ISA and ESA variants are referred to as ISAs in this section.

Both MS-ISA and MS-ISA2 support FCC and RET, VQM, and perfect stream functionality. The ESA currently only supports FCC, RET, and perfect stream . Perfect stream is only supported on multi-complex platforms.

When configured in the router, ISAs are logically grouped into video groups for video services. A video group allows multiple ISAs/ESAs for an application. The system performs load balancing between the ISAs. Nokia recommends not mixing ISA, ISA2, and ESA in the same video group. All ISAs within a video group are always active, there are no standby ISAs.

Video groups provide a redundancy mechanism to guard against hardware failure within a group where the system automatically rebalances tasks to the group excluding the failed MS-ISA or ESA-VM. Video groups also pool the processing capacity of all the group members and increase the application throughput because of the increased packet processing capability of the group. The buffer usage is typically identical for all members of the video group, so increasing the number of members in a group does not increase the scaling numbers for parameters bounded by available buffering, but there is still the increase in performance gained from the pooled packet processor capacity. A video service must be enabled at the video group level before that service can be used.

A video application may restrict the number of ISAs supported in a video group to a smaller number. Please contact your local Nokia representative for more information about video application specifications.



Note: There are restrictions to the number of MS-ISA and ESA-VMs that can be added to a video group. MS-ISA and ESA-VMs must not intermix within a video group. Only the SR-7/12/12e platforms require video groups to be in separate forwarding complexes. For more information, contact your local Nokia representative.

6.1.2 Video SAP

The video group logically interfaces to a service instance with a video Service Access Point (SAP). Like a SAP for connectivity services, the video SAP allows the assignment of an ingress and egress filter policy and QoS policy.

Ingress and egress directions for the filter and QoS policy are named based on the perspective of the router which is the opposite perspective of the ISA. An "egress" policy is one that applies to traffic egressing the router and ingressing the ISA. An "ingress" policy is one that applies to traffic ingressing the router and egressing the video. Although potentially confusing, the labeling of ingress and egress for the ISA policies was chosen so that existing policies for connectivity services can be reused on the ISA unchanged.

If no filter or QoS policy is configured, the default policies are used.

One of the key attributes of a video SAP is a video group association. The video SAP's video group assignment is what determines which video group services on that video SAP. The video groups configuration determines what video services are available.

6.1.3 Video interface

The video interface acts as the following:

- FCC and RET server
- interface for multicast buffering
- interface for the perfect stream

A video interface supports up to 16 IP addresses and is supported on VPRN or IES services.

6.1.3.1 Video interface properties

The video application is supported on various platforms but the video interface has different characteristics on the following platforms:

- **7750 SR-7, 7750 SR-12, 7750 SR-12e**

For 7750 SR-7, 7750 SR-12, 7750 SR-12e, the video interface is a logical interface regardless of whether ISA or ESA is used. For these platforms, the video interface uses direct protocol as the routing policy to export to IGP or BGP protocols. Use the following command to configure the routing policy for a video interface as direct protocol:

- **MD-CLI**

```
configure policy-options policy-statement entry from protocol name direct
```

- **classic CLI**

```
configure router policy-options policy-statement entry from protocol direct
```

- **All other platforms supporting video applications**

The video interface is modeled as a route. Only ESA is supported on these platforms. For these platforms, the video interface uses video protocol as the routing policy to export to IGP or BGP protocols. Use the following command to configure the routing policy for a video interface as video protocol:

- **MD-CLI**

```
configure policy-options policy-statement entry from protocol name video
```

– **classic CLI**

```
configure router policy-options policy-statement from protocol video
```

6.1.4 Multicast information policies

Multicast information policies on the 7750 SR and 7450 ESS serve multiple purposes. In the context of a service with video services, the multicast information policy assigned to the service provides configuration information for the multicast channels and defines video policy elements for a video interface.

This section describes the base elements of a multicast information policy in support of a video service. Specific video service features require additional configuration in the multicast information policy, which is described in the sections dedicated to the video feature.

Multicast information policies are named hierarchically structured policies that are composed of channel bundles that contain channels containing source-overrides.

- A bundle has an assigned name and contains a collection of channels. If attributes are not defined for a bundle they are inherited from a special default bundle named "default". Use commands in the following context to configure bundles for the multicast information policy:

– **MD-CLI**

```
configure multicast-management multicast-info-policy bundle
```

– **classic CLI**

```
configure mcast-management multicast-info-policy bundle
```

- Channels are ranges of IP multicast addresses identified by a start IP multicast address (Gstart) and optional end IP multicast address (Gend), so the channels encompass (*,Gstart) through (*,Gend). A channel attribute is inherited from its bundle unless the attribute is explicitly assigned, in which case the channel attribute takes precedence.
- Source overrides within a channel are IP multicast addresses within the channel that have a specific source IP address (Soverride), so the source override encompasses (Soverride,Gstart) through (Soverride,Gend). A source-override attribute is inherited from its channel, unless the attribute is explicitly assigned in the source-override channel, in which case the source-override channel attribute takes precedence.

For an IP multicast channel (*,G) or (S,G), the most specific policy element in the hierarchy that matches applies to that channel.

A multicast information policy is assigned to a service instance. For video services, the multicast information policy assigned to the service determines the video group for a given IP multicast channel. When a channel is assigned to a video group, the channel is sent to the video group for buffering or processing as appropriate depending on the video services enabled on the video group. If no video group is assigned to a channel, the channel is still distributed within the service instance, but no video services are available for that channel.

In addition to bundles, channels, and source-overrides, multicast information policies also include video policies. Video policies define attributes for the video interfaces within the service instance.

Video policy attributes are specific to the video feature and are covered in detail in the applicable video feature section. Video policies are mentioned here because they are an element of the multicast information policy and provide the link to configuration for a video interface.

6.1.5 Perfect stream protection

To protect against network interruption or reconvergence, it is often more efficient to protect the stream using an alternate transmission path. This can be a separate physical interface, transmission link, system, or even technology.

Perfect stream is also known as duplicate stream in some markets. Perfect stream protection allows an operator to split a single multicast stream (single S,G and common SSRC) into two different transmission paths that may have different transmission characteristics, such as latency or jitter. Instead of selecting one stream for retransmission to the client, the perfect stream protection feature evaluates each stream packet-by-packet, selecting the first, valid packet for retransmission.

A circular buffer is used for perfect stream protection which incorporates both packet-by-packet selection (based on RTP sequence number/timestamp and SSRC) and a re-ordering function whereby any out-of-sequence packets are placed into the buffer in order, therefore creating a corrected, in-order stream.

Playout rate is a function in ingest rate, however because the two streams may be delayed between one-another a few assumptions are made:

- The first arriving packet is always put into the buffer, allowing for the backup medium to wander in terms of latency and jitter.
- Because the source is the same, the rate at which a packet is put into the buffer (from either stream) can be assumed to be the normal bitrate.

The output RTP stream is always maintained in-sequence and the playout speed is always controlled. A moving window calculated against the time stamp smooths jitter that may occur between packets or the two contributing streams. The playout stream is described in [Perfect stream selection](#).

6.1.6 Perfect stream selection

This section describes perfect stream protection and provides details about packet selection on playout.

6.1.6.1 Stream identification

Stream selection is a simple selection algorithm that is applicable to any number of input streams. It is a prerequisite for stream selection that RTPv2 encapsulation be used in UDP.

Each service is identified by multicast source, group/destination address and current synchronization source (SSRC). After the service has been identified, the ISA monitors its ingress for:

- traffic with a DA of the multicast group
- traffic with a DA of the ISA (unicast)

Traffic is further checked as having RTP-in-UDP payload, RTP version 2.

The SSRC of each incoming RTP packet is learned as a unique source. Only one SSRC is supported for each stream; as SSRC may change during abnormal situations (such as encoder failover), it can be updated.

An SSRC can only be updated when a Loss of Transport (LoT) occurs, as other perfect streams (with the original SSRC) may still be operational. When an LoT occurs, the SSRC is deleted, the buffers are purged, and the RTP sequence counters are reset. The SSRC is extracted from the next valid RTP packet and the sequence starts over.

One RTP packet from the perfect stream is selected for insertion into the video ISA buffer. After a packet is selected, the RTP sequence counter is incremented and any further RTP packets received by the ISA with the previous sequence number are discarded.

In summary, perfect stream selection is a FIFO algorithm for RTP packet selection; this is considered optimal because:

- All stream sources are identical. Therefore, for any sequence number, the payload should also be identical.
- Most bit errors should be detected by the CRC-32 algorithm applied to Ethernet, SDH, and so on.

These devices typically discard frames where bit errors occur, with the net result being that the video ISA receives a bit error-free stream (though packet loss can occur).

- The UDP checksum is verified by the video ISA (after input VQM) and any failures result in a silent discard of the packet.
- VQM can be used in conjunction with perfect stream protection.

VQM can be used to monitor the quality of two duplicate ingress multicast streams and the egress multicast stream (after perfect stream selection). This is particularly useful to compare between the ingress and egress multicast. Monitoring the egress multicast after perfect stream selection can provide an insight to the customer viewing experience.

6.1.6.2 Initial sequence identification

When a service is defined and is administratively enabled, the video ISA monitors for valid RTP packets and on first receipt of a valid RTP packet learn the following information:

- SSRC
- sequence number
- timestamp (as timestamp is profile-specific, MPEG2-TS are assumed)

The packet is inserted into the video ISA playout buffer associated with that particular service and playout when directed (playout algorithm).

6.1.6.3 Packet selection

Each valid RTP packet received for a specific service is inserted into the buffer if there is no existing RTP packet that matches the sequence number. When sequence numbers and timestamps discontinue, the video ISA makes a limited attempt to validate the MPEG stream. The video ISA code adopts a philosophy to ensure that sequence numbers and timestamps increment correctly. If packets are non-contiguous, the packet selection algorithm adapts.

Duplicate packets are detected by sequence number or timestamp, unless M-bit resets the timestamp. A packet that is already in the buffer and which has the same sequence number as a received packet (or one recently played out) is discarded. The system monitors the incoming sequence number on each RTP packet. If a packet takes more than 1.5 seconds, it is considered late and is discarded.

In a multiprogram transport stream (MPTS), the timestamp is uniquely set for every RTP packet, as any RTP packet may contain a number of multiplexed elementary streams. As a result, playout is based on the embedded timestamp in each RTP packet. In a single-program transport stream the inverse occurs. Many RTP packets can share the same timestamp as it is referenced from the start of a picture (and a picture can span many RTP packets). As an SPTS does not contain audio, its application is limited to content production and so only MPTS are supported.

Timestamp discontinuities do occur and are normally represented with the Marker bit (M) being set.

Playout time is determined by an internal playout timestamp. The playout timestamp is set independently from the actual timestamp in the packet. The recovered clock can compute the expected timestamp for the very next incoming RTP packet.

When a packet is received, it is first compared to existing packets in the buffer based on sequence number (assuming that a stream may be delayed for hundreds of milliseconds by a backup path yet still be valid); jitter tolerance is only evaluated if this packet is determined to be a new RTP packet eligible for buffer insertion. If jitter tolerance is exceeded, then a timestamp discontinuity is assumed and instead of setting the playout timestamp based on the contained RTP timestamp, the actual received time (offset by playout-buffer) is set for the RTP packet playout timestamp.

In normal operation, the clock is recovered from the timestamp field in the RTP header, is offset by the playout buffer configuration parameter, and is used to schedule playout of the packet. The playout clock is synchronized with the sender by using an adaptive clock recovery algorithm to correct for wander.

Algorithm summary

- For a service marked LoT, if a loss of transport occurred, purge the buffer and reset all counters/timers.
- If the service is UP, check the RTP packet sequence number. Compare to sequence numbers contained in the buffer. If no match then check last played sequence number. If the sequence number of this packet is between last played and last played + 4096 then consider this packet late and discard.
- Check the expected timestamp recovered clock value and compare to RTP timestamp. If the expected timestamp is $(-ve)jitter\ tolerance < timestamp < (+ve)jitter\ tolerance$ then the packet is admitted to the buffer with a playout timestamp per the embedded RTP timestamp. If jitter tolerance is not maintained this marks a discontinuity event. Set playout timestamp to current clock + playout buffer and enqueue.

6.1.6.4 Clock recovery

RFC 2250, *RTP Payload Format for MPEG1/MPEG2 Video*, defines the timestamp format for MPEG2 video streams (which may carry H.264 video): a 90kHz clock referenced to the PCR. Each ingest RTP packet has its timestamp inspected and it is used in an adaptive clock recovery algorithm. Importantly, these adjustments occur on ingress (not on playout). This serves as a long-term, stable, ingress stream recovered clock.

The 90kHz ingress stream recovered clock is adjusted for each service to account for the encoder's reference clock/difference between the clock in the 7750 SR. This input timestamp is derived from the same RTP packet that is inserted into the buffer, and therefore may be subjected to significant jitter. The clock adjustment algorithm must only adjust clock in extremely small increments (in the order of microseconds) over a very long sample period (not bitrate) of at least 30 minutes.

6.1.6.5 Playout

Playout is the process of regenerating the stream based on the playout timestamp.

The playout is an exact offset to the ingress stream-recovered clock and serves as playout time for the video ISA. Because the timestamp is used for buffer playout, CBR, capped VBR, and VBR streams are all supported without pre-configuration. The playout buffer mechanism effectively removes network-induced jitter and restores the output to the rate of the original encoder.

6.1.6.6 Loss of transport

In the circumstance that the playout buffer is emptied an LoT is indicated. The video ISA resets playout timestamp, clock, sequence number, and so on for this event and awaits the next valid RTP packet for this service.

6.1.6.7 Perfect stream in relation to FCC/RET

Perfect stream can enhance the end user experience by removing network induced faults to the multicast stream. The egress of the perfect stream can be fed back into a separate ISA for the FCC/RET application. This ensures that the end user gets the best IPTV viewing experience from the start of a channel change. The perfect stream can be sent to the FCC/RET ISA by using extranet where FCC/RET resides on a separate VPRN instance.

6.1.6.8 Perfect stream in relation to VQM

Perfect stream uses two identical multicast streams traversing two diverse network paths to construct a multicast stream that is less susceptible to network faults. VQM can be configured in conjunction with perfect stream, which allows VQM to monitor the two ingress as well as the egress streams. By using the VQM statistics from all three streams, the operator can assess the quality of the network transporting the multicast stream, the performance of perfect stream, and the viewing quality for the end user. VQM and Perfect Stream can be enabled on the same video-group (the same ISA).

6.1.7 Video quality monitoring

The following terminology is used in this section:

- TNC is an acronym for Technically Non-Conformant.
- QoS is an acronym for Quality of Service.
- POA is an acronym for Program Off Air.
- TNC event (also known as impaired event) refers to a trap/alarm that detects an impairment event and is therefore termed TNC.

An impaired event is said to have occurred if:


- A PAT/PMT syntax error occurs in that second.
- Continuity counter errors were detected.
- PAT /PMT/PCR PIDs were not present in the video stream for a time period equal to or greater than the configured TNC value in the respective alarm.

The default value of the impaired threshold in terms of milliseconds is:

- PAT (100 ms)

- PCR (100 ms)
- PMT (400 ms)
- An unreferenced PID is seen in the video stream which has not been referred in the PMT.
- TNC seconds (also known as impaired seconds) refers to the number of seconds elapsed before an impaired event is detected and a TNC SNMP trap is sent. Use the following command to display TNC events.


```
show video channel analyzer
```

-  **Note:** Although multiple TNC events may occur within a second, the displayed counters only increment once per second.
- QoS event (also known as degraded event) refers to a trap/alarm that detects a degraded event and is therefore termed QoS. A QoS event is said to have occurred if PAT/PMT/PCR PIDs are absent in the video stream for a time period equal to or greater than the configured QoS value in their respective alarms.

The default value of the degraded threshold in terms of milliseconds is:

- PAT (200 ms)
- PCR (200 ms)
- PMT (800 ms)
- QoS seconds (also known as degraded seconds) refer to the number of seconds elapsed before a degraded event was detected and a QoS SNMP trap was sent. Use the following command to display QoS events.

```
show video channel analyzer
```

-  **Note:** Although multiple QoS events may occur within a second, the displayed counters only increment once per second.
- POA event (also known as error event) refers to a trap/alarm that detects an error event and is therefore termed POA.

A POA event has occurred if:

- A synchronization loss error has occurred for that particular second. A synchronization loss is said to have occurred if more than 1 consecutive synchronization byte error is seen in the stream.
- PAT/PMT/PCR PIDs are absent from the video stream for a time period equal to or greater than the configured POA value.

The default value of the degraded threshold in terms of milliseconds is:

- PAT (500 ms)
- PCR (500 ms)
- PMT (2000 ms)
- Traffic loss has occurred for that particular second.
- A transport error indicator or TEI indicator is set in the transport stream packet header for that particular second in the video stream.

- POA seconds (also known as error seconds) refer to the number of seconds elapsed before an error event is detected and a POA SNMP trap is sent. Use the following command to display POA events.

```
show video channel analyzer
```



Note: Although multiple POA events may occur within a second, the displayed counters only increment once per second.

- Good seconds refer to the number of seconds where there are no impaired, degraded, or error events.

PID statistics are described in [Table 35: PID stats — description of fields](#)

Table 35: PID stats — description of fields

Field	Description
PID	Displays the value of the PID.
Is PCR PID	Can be set to Yes or No. If set to Yes, then it indicates that the PID is the PCR PID.
TEI Err Sec	Counts the number of seconds TEI was set for that particular PID.
Absent Err Secs	The number of seconds for which the PID was not seen for a particular interval of time which is decided by the alarms set for the non-Vid PID Absent and Video PID Absent.
PID bitrate	Calculated by counting the number of times the PID occurred in the last second x 188 x 8. <ul style="list-style-type: none"> • 188 = TS packet size • 8 = Number of bits in a byte
CC Err Secs	Number of seconds Continuity Counter errors were seen for that particular PID in the stream.
PID Type	Specifies that the PID is either video, audio, PAT, PMT, or PCR.
MPEG Stream Type	If the PID is video or audio, this field indicates how the video or audio is encoded. For example: <ul style="list-style-type: none"> • For video, H.265, H.264, or MPEG2 (only the decimal equivalent defined by the MPEG standard is displayed and not the string) • For audio, E-AC3, DTS-HD, AC-3, or MPEG-2 (only the decimal equivalent defined by the MPEG standard is displayed and not the string)

To address interval statistics, except for the PID statistics all other statistics described above have interval statistics. Information can be obtained about stream status for the last 1 minute, 5 minute, and 15 minute interval.

MDI - Media Delivery Index (RFC 4445, *A Proposed Media Delivery Index (MDI)*)

- **Delay Factor (RFC 4445)**

The delay factor is a value which indicates the minimum amount of time a STB buffers to resolve network jitter (that is, it is the minimum STB buffer depth in ms). RTP timestamp is used as the definitive time indicator (the notional drain rate).

- **Loss Rate (RFC 4445)**

The Media Loss Rate is the number of media (Transport Stream) packets lost over a specific time interval. This is reported in TS/sec. Each RTP packet lost is assumed to have 7 TS packets lost.



Note: In absence of traffic MDI values are reported as N/A. These stats are reported over current (current second) , 1 minute, 5 minutes and 15 minutes intervals

In many instances IPTV operators are unable to identify the cause of visual impairments which are present in almost every video distribution network because the IPTV network has so many moving parts While head end transport-stream monitoring; full reference video analysis (comparing the source content to the encoded output), and; STB probes allow an operator to establish whether the contribution source, the encoder, or the network is the problem the network is a very complex thing.

Operators can use another measurement point in the network, just before the last mile such that network faults can be characterized as being between the head end and last mile (transport) or in the last-mile itself.

The multicast video quality monitoring solution provides an inspection point for the multicast video stream that is combined with other analysis methods to create a full view of video issues and help troubleshoot the part of the network causing the issue.

Video quality monitoring is one part of a video assurance program and is combined with:

- TS analysis on the encoder output (to detect encoder errors)
- Full-reference PSNR and PQR on the encoder output (to detect over-encoding, noise and other contribution or encoding artifacts)
- STB reporting (such as packet-loss, RET events, packet errors) from the entire STB population
- STB probes performing full-reference monitoring (against test streams)
- STB probes performing channel-change times, estimated PSNR, and so on

Multicast video monitoring within the network can be positioned as complementary to STB reporting and head end analysis, and but should not attempt to perform either of these functions. Because the network node is not capable of decrypting a MPEG transport stream is primarily used to identify correctable and uncorrectable network errors, correlate them with network events (such as routing re-convergence, interface failure, and so on) and provide summary reports and alarms.

For operators who do not have existing STB probes or reporting, a network-based VQM solution can provide insight into quality issues the network may be contributing to, possibly reducing the amount of STB probe investment that is needed. (that is, both probes and the 7750 SR VQM reports many of the same issues in terms of picture quality, fewer probes are needed to test channel change delay, and so on).

The metrics which VQM can report are based on the use of RTP streams which provide per-packet sequencing and an indication of picture type. These two parameters along with measured bitrate allow

VQM to produce estimated MOSv scores for both stream ingress (uncorrected) and stream egress (corrected) outputs.

Reportable metrics include:

- Relevant SCTE-143 error counters
 - PAT
 - PMT
 - PCR
 - transport errors, and so on
- ETSI TR 101 290
 - PID
 - SI repetition
 - degraded blocks/intervals, and so on
- MDI (RFC 4445)
- forwarded and impaired I-/B-/P-frame counts
- GOP length
- video/audio/stream bitrate

These metrics are collected for each multicast group and have relevant command options. Use the following command to generate a detailed numeric metrics report for a specific multicast group.

```
show video channel analyzer address detail
```

Output example: show video channel analyzer address 239.0.1.1 detail

```
=====
Video channel analyzer detail
=====
Channel number : 1
-----
Service Id      : 31905555          Interface Name   : vi-1
Group Address   : 239.0.1.1        Source Address  : 10.10.10.10
MDI Delay Factor : 4              MDI Loss Rate   : 0
Good Secs      : 0

=====
GOP Stats
=====

```

	Min	Max	Avg
GOP Length	10	31	24
Frames/Sec	21	28	25

```
=====
Frame Stats
=====

```

	I-Frame	P-Frame	B-Frame
Good	42	957	0
Bad	0	0	0

```
=====
```

```

Error Stats
=====
                POA Events      QoS Events      TNC Events
-----
PAT Rep          0                1215            0
PMT Rep          0                0               0
PCR Rep          0                0               1
PAT Syntax Err  -                0               -
PMT Syntax Err  -                0               -
Sync Byte Err   -                0               -
Sync Loss       0                -               -
Unref PID       -                -               34600
Traffic Loss    0                -               -
-----
Overall          30285            77             4238
-----
Reoccurring events only increment counter once every second

-----
Number of channels : 1
=====

```

Other information that is captured by VQM but can only be expressed on a per-PID basis. Use the following command to display PID address information per group.

```
show video channel pid address
```

The per-PID information captured by VQM includes the following:

- CC Err secs
- TEI Err secs
- Absent Err secs

Output example: show video channel pid address 239.0.1.1

```

=====
Video Channel PID
=====
Service Id      : 31905555      Interface Name   : vi-1
Group Address   : 239.0.1.1          Source Address  : 10.10.10.10
PID            : 0              PID Type        : pat
MPEG Stream Type : 0          Is PCR PID     : No
Cc Err Secs    : 0              TEI Err Secs   : 0
Absent Err Secs : 0              PID Bitrate    : 0

Service Id      : 31905555      Interface Name   : vi-1
Group Address   : 239.0.1.1          Source Address  : 10.10.10.10
PID            : 100           PID Type        : pmt
MPEG Stream Type : 0          Is PCR PID     : No
Cc Err Secs    : 0              TEI Err Secs   : 0
Absent Err Secs : 0              PID Bitrate    : 0

Service Id      : 31905555      Interface Name   : vi-1
Group Address   : 239.0.1.1          Source Address  : 10.10.10.10
PID            : 101           PID Type        : video
MPEG Stream Type : 27          Is PCR PID     : Yes
Cc Err Secs    : 0              TEI Err Secs   : 0
Absent Err Secs : 0              PID Bitrate    : 5564800

Service Id      : 31905555      Interface Name   : vi-1
Group Address   : 239.0.1.1          Source Address  : 10.10.10.10

```

```

PID           : 201           PID Type      : audio
MPEG Stream Type : 3           Is PCR PID   : No
Cc Err Secs   : 0           TEI Err Secs : 0
Absent Err Secs : 0           PID Bitrate  : 204544

Service Id    : 31905555      Interface Name : vi-1
Group Address : 239.0.1.1     Source Address : 10.10.10.10
PID           : 202           PID Type      : audio
MPEG Stream Type : 3           Is PCR PID   : No
Cc Err Secs   : 0           TEI Err Secs : 0
Absent Err Secs : 0           PID Bitrate  : 132352

Service Id    : 31905555      Interface Name : vi-1
Group Address : 239.0.1.1     Source Address : 10.10.10.10
PID           : 401           PID Type      : other
MPEG Stream Type : 6           Is PCR PID   : No
Cc Err Secs   : 0           TEI Err Secs : 0
Absent Err Secs : 30310       PID Bitrate  : 0

Number of pids for this channel: 6

-----
Number of channels : 1
=====

```

Event alarms are reported by log, syslog, or SNMP. The following is an example of an event trap.

Output example: Event trap

```

1 2017/02/11 18:11:20.42 UTC WARNING: VIDEO #2009 Base Video[1/video-300]
Service Id - 300, Video interface - video-300, Group address - 192.0.2.6, Source address
- 10.20.13.2 Last 10 seconds of analyzer state is good good good good good good good
good good poa

```

A trap is only raised when a POA, QoS, or TNC event occurs within the last 10 seconds. The trap captures events within the last 10 seconds. In the preceding example =, the first nine seconds were "good", which indicates that no events occurred and that every single RTP packet was received. At the 10-second mark, a POA event occurred, which triggered the SNMP trap. This sampling continues every 10 seconds. If an event (POA, QoS, or TNC) continues to be detected, an SNMP trap is raised every 10 seconds.

When the analyzer detects 10 seconds of a "good" condition, another trap is raised to clear the alarm for the multicast (S,G). Subsequent alarms are raised and SNMP traps are triggered only when the analyzer detects another event (POA, QoS, or TNC).

Output example: Event clear trap

```

61 2016/10/18 15:40:56.83 UTC WARNING: VIDEO #2010 Base Video[1/video-300]

Analyzer state is cleared for - Service Id - 300, Video interface - video-300, Group
address - 192.0.2.6, Source address - 10.20.13.2

```

VQM is an optional module available on the input side, or output side of the video ISA. On input, it is applied before perfect stream protection. Conversely when on the output side it is applied only to multicast streams after perfect stream protection.

Because of the large number of channels and the nature of measuring input and output sides, VQM is highly reliant on the use of RTP extensions to provide relevant transmission metrics to the VQM analysis module. In a typical head end, a multicast stream is scrambled to encrypt its video or audio. When this encryption occurs, it is typical for the entire payload of the transport stream (for the nominated PID) to be

completely scrambled. The consequence of such is that the video and audio PES headers, which reveal much about the picture and timing information, are unavailable to the VQM program.

VQM uses intelligent RTP re-wrapping. RTP re-wrapping is a prerequisite for ad insertion and Fast Channel Change (FCC) and involves marking packets before encryption based on the picture type (most importantly, the start of the I frame of IDR frame in H.264).

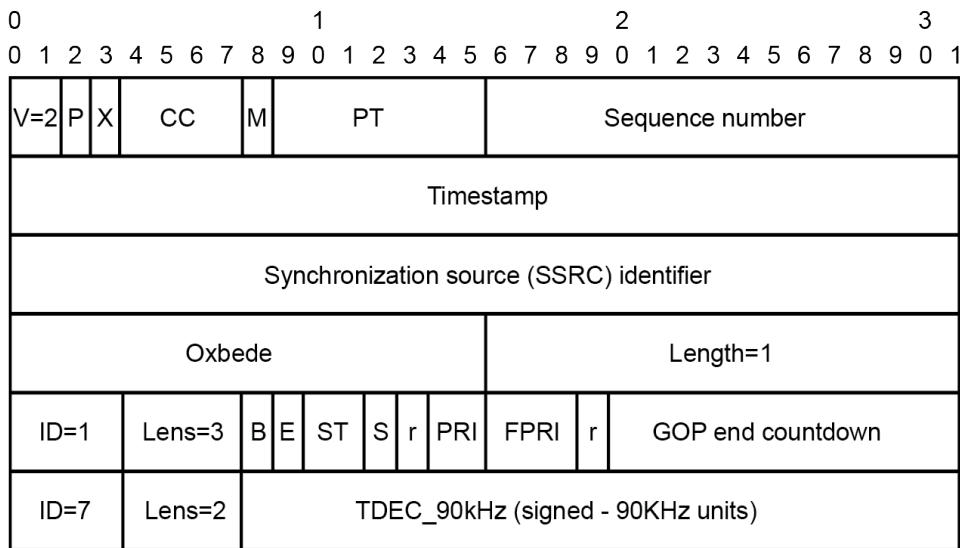
The VSA as currently defined, re-multiplexes each transport stream into a new RTP packet. By doing so it allows the separation of different picture types into their own respective RTP packets, and the separation of audio packets from video packets to allow different synchronization in events of FCC. In effect, it pulls the elementary streams back into their component forms while retaining the syntax and structure of the MPTS.

For information about VSAs, see the *7450 ESS, 7750 SR, 7950 XRS, and VSR System Management Guide*.

For quality analysis, additional information about the picture can be made available before scrambling. The quality analysis performed by the VQM module emphasizes impairments caused by network issues and transport stream syntax, based on the relative proximity of the router to the customer.

When the video ISA is deployed alongside the ALU VSA re-wrapper, a custom RTP header extension is sent with each RTP packet. The RTP header is shown in [Figure 50: RTP header](#).

Figure 50: RTP header



hw1290

Where:

- B (Frame Begin Flag): set if a frame starts in this packet
- E (Frame End Flag): set if a frame finishes in this packet
- ST (Stream Type)
 - 00 video
 - 01 audio
 - 10 data/padd/other
 - 11 Reserved
- S (Switch bit): set to 1 in all RTP packets from the moment the client should do the IGMP join (rewrap does not fill it)
- r: reserved (set to 1)
- PRI: Packet Priority (coarse)
- FPRI: Fine-grained priority

```

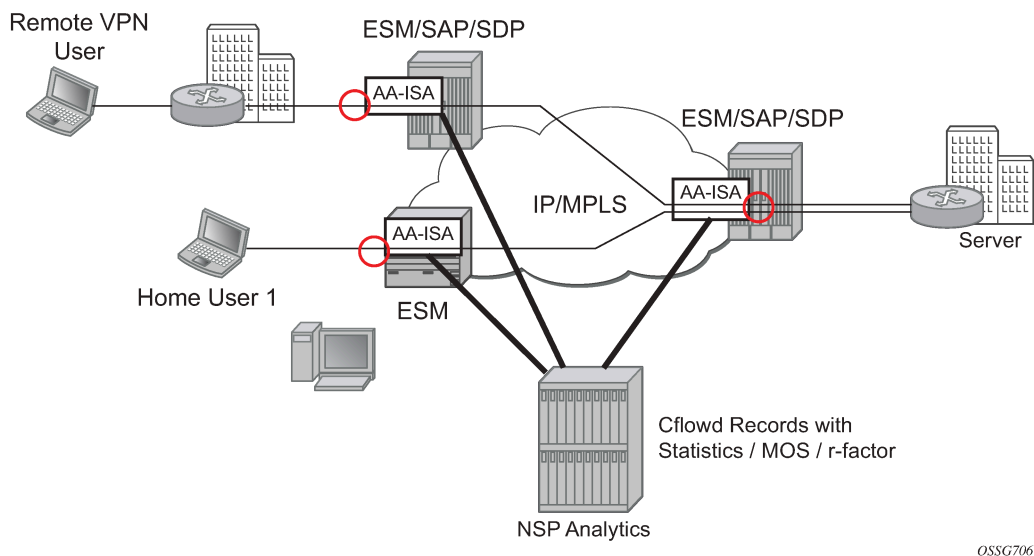
PRI FPRI dec DSCP
-----
3 7 31 AF41 Video IDR frame
3 0 24 AF41 Audio
2 0 16 AF41 Reference frame (P in MPEG2, I or P or some Bs in AVC)
1 7 15 AF42 Non-reference frame (most Bs in MPEG2, some Bs in AVC)
1 5 13 AF42 Trans-GOP frames, undecodable in some circumstances (some Bs in MPEG2)
0 4 4 AF43 Rest of cases (data, secondary videos, etc)
0 1 1 AF43 Padding packets
where AF41=100010, AF42=100100, AF43=100110 (DSCP bits in the IP header)

```

6.1.7.1 VoIP/video/teleconferencing performance measurements

The feature provides ability to measure and provide statistics to allow reporting on voice and video quality for VoIP and teleconferencing (A/V) applications. A sampled deployment is shown in [Figure 51: Voice/video monitoring deployment example](#). Although a distributed model is shown, a hub-and-spoke model, with AA-ISA deployed only on one side of the traffic flow, is also supported.

Figure 51: Voice/video monitoring deployment example



Because of network-based AA, the operator has an ability to monitor voice, video, teleconferencing applications for an AA subscriber regardless of the type of that subscriber (a residential subscriber vs. a user of a business VPN service). AA-ISA monitors UDP/RTP/RTCP/SDP headers for each initiated call/application session (sampling may be provided, although it is expected that a sampling rate is smaller than that of TCP-applications because of the nature of the voice/video applications, which are longer lasting with smaller numbers of sessions or calls per subscriber). AA ISA gathers statistics and computes MOS-scores/R-factor results per each call or application session. At the end of a call (application session closure), AA-ISA sends the statistics and computed scores to a Cflowd collector (the Cflowd infrastructure that was introduced for TCP-performance but modified to carry voice/video specific data is used). The collector summarizes and presents the results to the operator/end user.

6.1.7.2 Mean Opinion Score (MOS) performance measurements solution architecture

AA-ISA integrates a third party MOS software stack to perform VoIP and video MOS measurements. This software provides:

- call quality analysis using optimized ITU-T G.107
- measurements of perceptual effects of burst packet loss and recency using ETSI TS 101 329-5 Annex E Extensions
- measurements and analysis of RTCP XR (RFC3611) VoIP metrics payloads

AA software monitors the associated SDP channel and passes codec information (when available) to the subsystem which monitors VoIP. The video bearer channels traffic generates a wide variety of A/V performance metrics such as:

- call quality metrics
 - listening and conversational quality MOS scores (MOS-LQ, MOS-CQ)
 - listening and conversational quality R-factors (R-LQ, R-CQ)
 - estimated PESQ scores (MOS-PQ)
 - separate R-factors for burst and gap conditions (R-Burst, R-Gap)
 - video MOS-V and audio MOS-A
 - video transmission quality (VSTQ)
- video stream metrics
 - good and impaired I, B, P, SI, SP frame counts
 - automatic detection of GoP structure and other key video stream attributes such as image size, bit rate, codec type
- transport (IP/RTP) metrics
 - packet loss rate, packet discard rate, burst/gap loss rates
 - packet delay variation/jitter
- degradation factors (degradation because of loss, jitter, codec, delay, signal level, noise level, echo, recency)

When a flow terminates, AA software retrieves the flow MOS parameters from the subsystems, formats the info into a Cflowd record and forwards the record to a configured Cflowd collector (RAM).

RAM collects Cflowd records, summarizes these records using route of interest information (source/destinations). In addition, RAM provides the user with statistics (min/max/ avg values) for the different performance parameters that are summarized.

6.2 Retransmission and Fast Channel Change

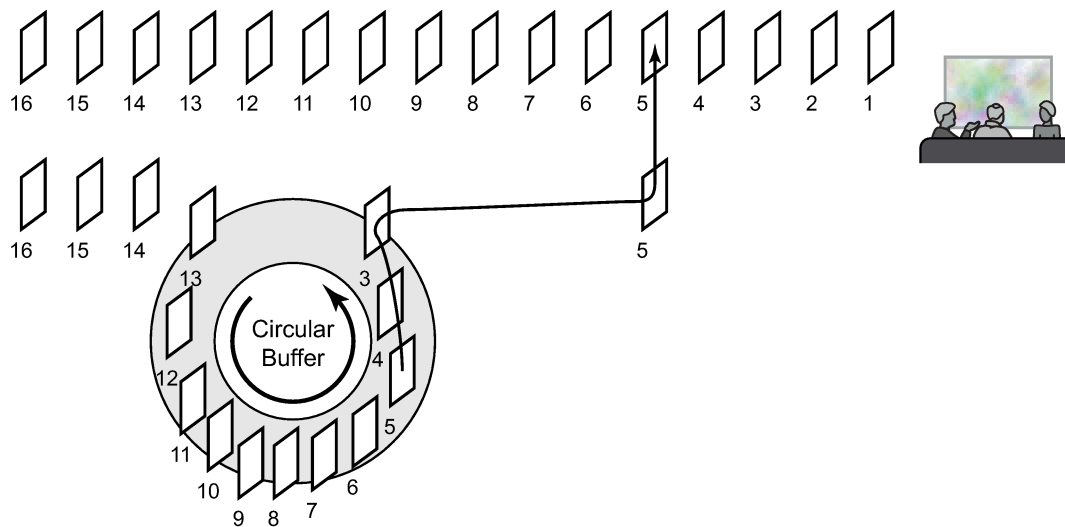
6.2.1 RET and FCC overview

The following sections provide an overview of RET and FCC.

6.2.1.1 Retransmission

Retransmission (RET) for RTP (RFC 3550, *RTP: A Transport Protocol for Real-Time Applications*) is based on a client/server model where the client sends negative acknowledgments (NACKs) using Real-time Transport Control Protocol (RTCP) (RFC 4585, *Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)*) to a RET server when the client detects missing sequence numbers in the RTP stream. The RET server which caches the RTP stream, for example in a circular buffer, detects missing sequence numbers in the replies to the NACKs by resending the missing RTP packets as illustrated in [Figure 52: RET server retransmission of a missing frame](#).

Figure 52: RET server retransmission of a missing frame



OSSG321

The format of the reply must be agreed upon by the RET client and server and can be an exact copy (Payload Type 33 as defined in RFC 3551, *RTP Profile for Audio and Video Conferences with Minimal Control*) or sent with a different Payload Type using an encapsulating RET header format (RFC 4588, *RTP Retransmission Payload Format*).

RET has been defined in standards organizations by the IETF in the above-noted RFCs and Digital Video Broadcasting (DVB) in "Digital Video Broadcasting (DVB); Transport of MPEG-2 TS Based DVB Services over IP Based Networks (DVB-IPTV Phase 1.4)" which describes the STB standards.

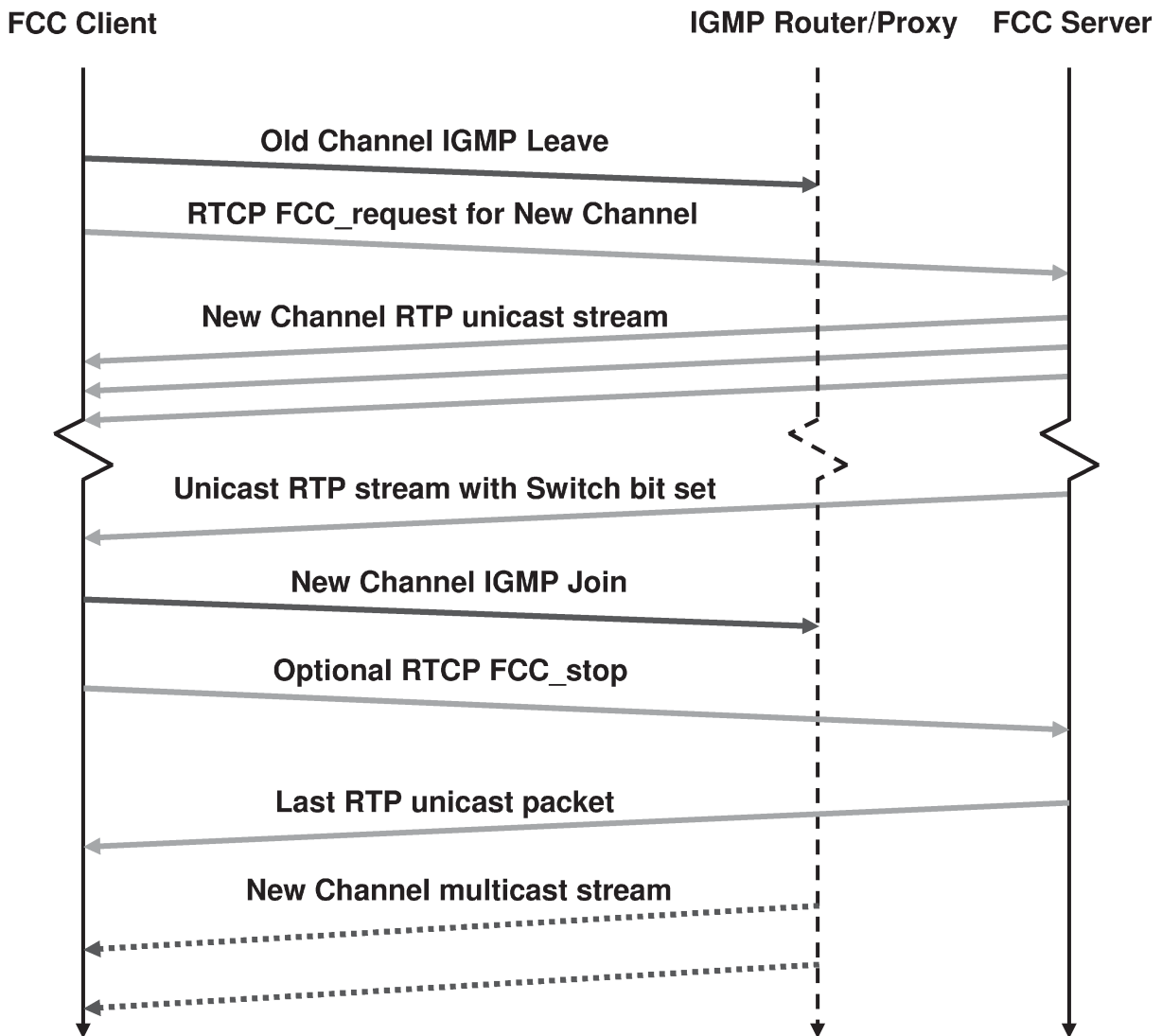
STBs that have a port of the Nokia RET/FCC Client SDK are an example of a standards-compliant RET Client implementation.

6.2.1.2 FCC

FCC is an Nokia method based on a client/server model for providing fast channel changes on multicast IPTV networks distributed over RTP. During a fast channel change, the FCC client initiates a unicast FCC session with the FCC server where the FCC server caches the video stream and sends the channel stream to the FCC client starting at the beginning of a Group of Pictures (GOP). Beginning at a GOP in the past minimizes the visual channel transition on the client/STB, but the FCC unicast stream must be sent at an accelerated rate in the time domain to allow the unicast stream to catch up to the main multicast stream, at which point, the FCC server signals to the client to join the main RTP stream.

[Figure 53: FCC client/server protocol](#) illustrates the FCC client and server communication.

Figure 53: FCC client/server protocol



There are two techniques for compressing the FCC unicast stream in time to allow the unicast session to catch up to the multicast stream: bursting and denting. When bursting, the stream is sent at a rate faster than multicast stream, for example, the stream can be “bursted” at 130% (or 30% over the nominal) multicast rate. “Denting” is a technique where less important video frames are dropped by the FCC server and not sent to the FCC client. Hybrid mode combines bursting and denting.

Bursting is illustrated in [Figure 54: FCC bursting sent faster than nominal rate](#) and denting is illustrated in [Figure 55: FCC denting removing less important frames](#).

Figure 54: FCC bursting sent faster than nominal rate

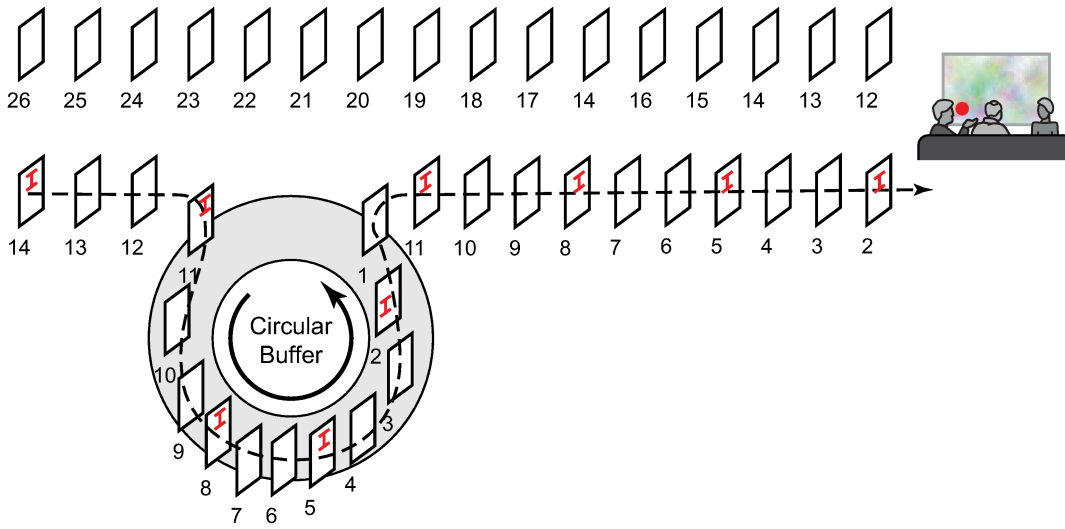
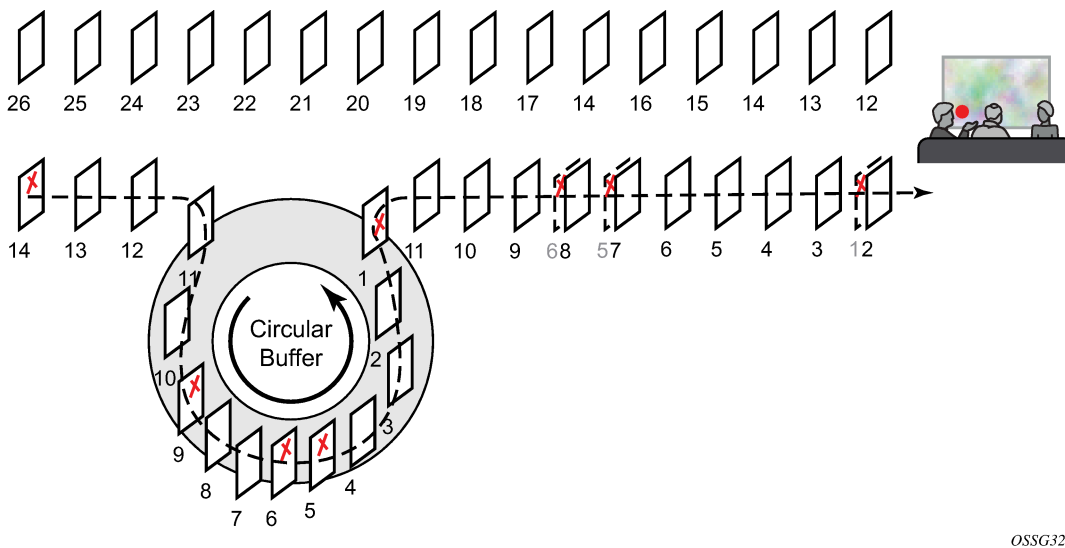


Figure 55: FCC denting removing less important frames



When the unicast session has caught up to the multicast session, the FCC server signals to the FCC client to join the main multicast stream. The FCC server then sends the unicast session at a lower rate called the “handover” rate until the unicast session is terminated.

The FCC server functionality requires the Nokia 5910 Video Services Appliance (VSA) Re-Wrapper which is used to encapsulate and condition the multicast channel streams into RTP, adding important information in the RTP extension header. Also, the ISA FCC server requires an STB FCC client based on the Nokia FCC/RET Client SDK.

6.2.1.2.1 Retransmission server

The ISA RET server is supported within an IES or VPRN service context as applicable to the platform.

Whether the RET server is active for a specific multicast channel is defined in the multicast information policy where channels are defined. The channel configuration for the RET server within the policy is an explicit enable/disable of the local RET server (that is, whether the channel should be buffered), the RET buffer size for the channel in the ISA and a channel type (Picture-in-Picture (PIP), Standard Definition (SD) or High Definition (HD)). The RET buffer should be large enough to account for the round trip delay in the network; typically, a few hundred milliseconds is sufficient.

In an IES or VPRN service, up to 16 IP addresses can be assigned to a video interface.

The video policy within the multicast information policy defines the characteristics for the how the RET server should respond to NACKs received on an IP address. The different characteristics defined in a RET server "profile" are for each channel type (PIP, SD and HD):

- enable/disable for the RET server (that is, whether requests should be serviced or dropped)
- the RET rate (as a percentage of the nominal channel rate)

Typically, RET replies are sent below line rate because most dropped packets occur in the last mile and sending RET replies at a high rate may compound any last mile drop issues.

The IP address(es) of the RET server are defined in the unicast service instance, whereas the UDP port for the RET server is defined in the "default" bundle in the multicast information policy. The same UDP port is used for all RET server IP addresses that use the particular multicast information policy.

The ISA RET server supports the network model where there are separate service instances for unicast and multicast traffic that are cross-connected and multicast replicated downstream in the network. If there are separate multicast and unicast service instances, the unicast and multicast services must use the same multicast information policy.

6.2.1.2.2 FCC server

The ISA FCC server is supported within an IES or VPRN service context as applicable to the platform. VPRN services are not supported on the 7450 ESS.

Whether the FCC server is active for a specific multicast channel is defined in the multicast information policy where channels are defined. The channel configuration for the FCC server within the policy is an explicit enable/disable of the local FCC server (that is, whether the channel should be buffered) and a channel type PIP, SD or HD. When FCC is enabled, three (3) GOPs are stored in the buffer. The channel also defines an optional FCC tuning parameter called the FCC Minimum Duration which is used by the FCC server to determine which GOP to start the FCC unicast session. If there are too few frames of the current GOP stored in the FCC server buffer (based on number of milliseconds of buffering), the FCC server starts the FCC session from the previous GOP.

In an IES or VPRN service, up to 16 IP addresses can be assigned to a video interface.

The Video Policy within the multicast information policy defines the characteristics for the how the FCC server should respond to FCC requests received on an IP address. The different characteristics defined in an FCC server "profile" are for each channel type (PIP, SD and HD):

- enable/disable for the FCC server (for example, should the requests be serviced or dropped)
- the FCC mode (burst, dent or hybrid)
- the burst rate (as a percentage above the nominal channel rate) for PIP, SD and HD channel types

- the multicast handover rate (as a percentage of the nominal channel rate) used by the server after it has signaled the client to join the main multicast channel

Different FCC rates are allowed for each of the channel types because the channel types have different nominal bandwidths. For example, the last mile may only be able to reliably send a 25% burst (above nominal) for HD whereas the equivalent bit rate for SD is a 75% burst. The profiles are designed to provide flexibility.

The IP address of the FCC server is defined in the unicast service instance, whereas the UDP port for the FCC server is defined in the "default" bundle in the multicast information policy. The same UDP port is used for all FCC server IP addresses that use the particular multicast information policy.

The ISA FCC server supports the network model where there are separate service instances for unicast and multicast traffic that are cross-connected and multicast replicated downstream in the network. If there are separate multicast and unicast service instances, the unicast and multicast services must use the same multicast information policy.

6.2.1.2.3 Logging and accounting for RET and FCC

Statistics and accounting are supported, but limited to capturing 256 000 sessions per five-minute period. Configure both FCC and RET timers to five minutes. If either timer is changed, the logging and accounting statistics may not be accurate.

Statistics and accounting are available for the following:

- RET server sessions statistics
- FCC session statistics

6.2.1.3 RET and FCC server concurrency

Even though the previous sections discussed the RET server and FCC server as separate entities, the ISA can support RET and FCC servers at the same service at the same time. Therefore, the configuration commands and operational commands for the services are intermingled. If both the RET server and FCC server are enabled for a specific channel, a single buffer is used for caching of the channel.

A maximum bandwidth limit for all server requests can be defined for a specific "subscriber" which is equated with the source IP address. Before an ISA server processes a request, the ISA calculates the bandwidth to the subscriber required, and drops the request if the subscriber bandwidth limit is exceeded.

The ISA services RET and FCC requests on a first in, first out (FIFO) basis. Before servicing any request, the ISA calculates whether its egress bandwidth can handle the request. If there is insufficient egress bandwidth to handle the service request, the request is dropped. Near the ISA's egress limits, RET requests generally continue to be serviced whereas FCC requests are dropped because RET sessions are generally a fairly small percentage of the nominal rate and FCC sessions are slightly below to above the nominal channel rate.

6.2.1.3.1 Prerequisites and restrictions

This section summarizes some key prerequisites and restrictions for the RET client, RET server and FCC server:

- Both RET and FCC require RTP as the transport stream protocol.

- FCC requires the Nokia 5910 VSA Re-Wrapper.
- FCC requires an implementation of the Nokia 5910 STB Client.
- The multicast information policies must be the same on multicast and unicast services which are cross connected downstream.
- Up to 16 IP addresses can be configured for a Layer 3 service video interface (IES or VPRN) with each supporting a distinct profile.
- There can be a maximum of 32 IP addresses across all Layer 3 service video interfaces per chassis.

6.2.2 Separate timers for FCC and RET

For each fast channel change, a new RTCP session is initialized, whereas for retransmission, the same RTCP session is always reused, with the same source and destination ports. The RTCP session duration is based on a timeout configuration, and a separate timeout option is available for FCC and RET. Use the following commands to configure the timeout:

- **MD-CLI**

```
configure multicast-management multicast-info-policy video-policy video-interface fcc-session-timeout
configure multicast-management multicast-info-policy video-policy video-interface ret-session-timeout
```

- **classic CLI**

```
configure mcast-management multicast-info-policy video-policy video-interface fcc-session-timeout fcc-session-timeout
configure mcast-management multicast-info-policy video-policy video-interface fcc-session-timeout ret-session-timeout
```



Note: Because the FCC RTCP session timeout is generally shorter than the RET session timeout, Nokia recommends configuring the FCC RTCP session timeout to a lower timeout value, accounting for the time required to complete a fast channel change and multicast handoff. This strategy frees RTCP sessions for other subscribers and improves efficiency.

6.2.3 Peak bandwidth and sessions per ISA

Use the following command to display the peak egress bandwidth and sessions per ISA. Each peak value is also marked with a timestamp.

```
show isa video-group
```

There are configurable watermarks for bandwidth and session consumption. The watermarks generate SNMP traps or log events when the ISA has reached the bandwidth or session thresholds. The egress bandwidth and session watermarks are configurable for FCC and RET separately or together. The SNMP trap or log event is generated as the configured watermark is exceeded, and is cleared if the egress rate falls below 10% of the watermark.

For example, if the maximum egress rate is 10 Gb/s and the watermark is set to 90%:

- The alarm is generated when the egress rate exceeds the following:

9 Gb/s (90% of 10 Gb/s)

- The alarm is cleared when the egress rate falls below the following:

0% at 8.1 Gb/s (10% of 9 Gb/s = 0.9 Gb/s, and 9 Gb/s - 0.9 Gb/s = 8.1 Gb/s)

6.2.4 Video support on ESA

The following applies to video support for the ESA:

- The FCC and RET applications are supported.
- VQM is not supported.
- Perfect stream is supported on multi-complex platforms only.
- SR-1/1s/2s/7/12/12e/7s/14s and SR-1-24D are currently supported.
- ISA and ESA must not be used in the same video group.

6.3 IP Video for Live

Traditionally, IPTV is delivered over a broadband network and subscribers require IPTV STB to receive multicast content. Both the STB and its IP gateway must support IP multicast and IGMP. Because native IGMP cannot route beyond the first IP hop, subscriber reach is often restricted.

IP Video for Live further extends the FCC feature. Use the following command to enable IP Video for Live:

- **MD-CLI**

```
configure multicast-management multicast-info-policy video-policy video-interface extended-unicast
```

- **classic CLI**

```
configure mcast-management multicast-info-policy video-policy video-interface extended-unicast
```

Channel requests are triggered by an RTCP request. IPTV content is unicast to each individual subscriber using RTP over IP UDP, which:

- eliminates the need to have both IP multicast support and IGMP support on end user devices
- only requires subscribers to have IP connectivity. IPTV content is routed Over The Top (OTT) to individual subscribers.
- supports all FCC, RET, and IP Video for Live on the same ISA/ESA. Therefore, the ISA/ESA that supports FCC/RET for broadband IPTV multicast can be reused for IP Video for Live.
- allows OTT subscribers to use the same STB and software as the traditional broadband multicast IPTV subscribers. STB can be programmed to automatically use RTCP messages instead of IGMP, if IGMP requests are non-responsive. This allows for unified STB deployment regardless of the end subscriber type.
- delivers video content using UDP with steady bit rate, which does not require network dimensioning for TCP burstiness or buffering

- allows IP Video for Live subscribers to use the FCC and RET function. IP Video for Live uses RTCP control messages to retrieve IPTV content.

Before Release 23.7.R1, the ingress buffer of the multicast packet had a direct impact on the extended unicast streaming to the end subscriber. If the ingress multicast packet had missing packets, or a gap in the RTP sequence, the extended unicast streaming would abruptly stop for the end subscriber. Starting in Release 23.7.R1, the extended unicast feature can tolerate 20 missing consecutive packets from the multicast source. This tolerance also applies to jitter tolerance (20 packets).



Note: The jitter tolerance is based on the number of packets, instead of on the time; the higher the streaming throughput, the lower the jitter tolerance.

Overall, this jitter and loss tolerance helps improve the streaming experience.

6.3.1 Delayed IGMP join for FCC

The FCC solution requires the STB to switch from unicast to multicast after the fast channel change completes. The STB uses IGMP to switch to multicast. In some cases, the multicast process is delayed. Delay may occur because of IGMP processing delay, authentication delay, or network delay. Use the following command, which is used for IP Video for Live, to delay the IGMP for up to 5 minutes:

- **MD-CLI**

```
configure multicast-management multicast-info-policy video-policy video-interface extended-unicast
```

- **classic CLI**

```
configure mcst-management multicast-info-policy video-policy video-interface extended-unicast
```



Note: Nokia recommends enabling this command only in implementations where multicast delivery is delayed beyond 1.5 seconds. By default, extended unicast is disabled to ensure the optimal bandwidth for unicast to multicast handover.

6.4 Configuring video service components with CLI

This section provides information to configure RET/FCC using the command line interface.

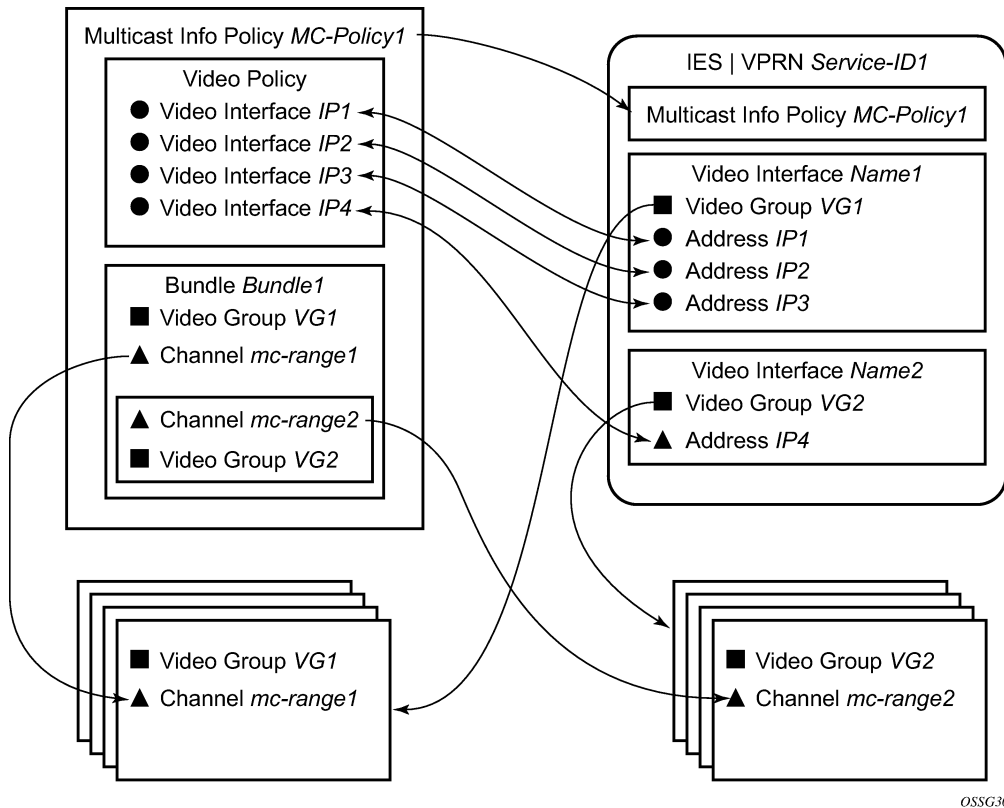
6.4.1 Video services overview

The main entities of video configurations are:

- video group
- multicast information policy
 - a video policy to configure video interface properties
 - multicast bundles and channels to associate bundles/channels with video groups
- within a service, configuring video interfaces and their associations with video groups

Figure 56: Video services configuration elements shows various configuration elements and how they are associated by configuration.

Figure 56: Video services configuration elements



OSSG308

A video interface within a service can have multiple IP address, and their association with the video interfaces within the video policy are based on IP addresses. Support for multiple video interface IP addresses for a specific video interface allows video characteristics (burst rate, retransmission format, and so on) for the channels associated with the video interface to be based on the IP address on which the request is received.

The bundle and channel configuration and the video interface configuration within the service are associated with a specific video group. If the request is received on a video interface for a channel not serviced by the video group associated with the video interface, the request is invalid and is dropped.

Figure 56: Video services configuration elements displays an example of this is a request for mc-range2 received on IP1, IP2 or IP3. A request for mc-range2 would only be valid on IP4.

As with other multicast information policies, the bundle name default is a special bundle and is reserved for setting of default values. If a video parameter is not explicitly set in a bundle/channel, the value set in the default bundle is used.

6.4.1.1 Configuring an ISA module

The ISA hardware has an MDA form factor that is provisioned in the same manner as other MDAs.

Use the following commands to configure an ISA module.

```
configure card card-type
configure card mda mda-type
```

Example: MD-CLI

```
[ex:/configure card 9]
A:admin@node-2# info
  card-type iom4-e-b
  mda 1 {
    mda-type isa2-aa
  }
  mda 2 {
    mda-type isa2-aa
  }
  fp 1 {
  }
```

Example: classic CLI

```
A:node-2>config>card# info
-----
  card-type iom4-e-b
  mda 1
    mda-type isa2-bb
  exit
  mda 2
    mda-type isa2-bb
  exit
  fp 1
  exit
-----
```

6.4.1.2 Configuring a video group

When used for video services, ISA-MSEs are logically grouped into video groups that pool the ISA buffering and processing resources into a single logical entity.

Use the commands in the following context to configure a video group.

```
configure isa video-group
```

The following example shows video-group 1 with a single ISA configured in slot 2/MDA 1.

Example: MD-CLI

```
[ex:/configure isa]
A:admin@node-2# info
  video-group 1 {
    admin-state enable
    description "Video Group 1"
    mda 1/2 { }
  }
```

Example: classic CLI

```
A:node-2>config>isa# info
=====
video-group 1 create
  description "Video Group 1"
  primary 7/2
  no shutdown
exit
=====
```

Within the video group configuration, there are specific video application commands to enable features. These commands are described in the configuration examples for the application. Depending on the video application, more than one primary ISA-MS is allowed increasing the egress capacity of the video group.

ISA-MS in a single video group cannot be on the same IOM. An IOM can accommodate two ISA-MS modules provided that the ISA-MS are members of different video groups.

6.4.1.3 Configuring a video SAP and video interface in a service

Video features in a VPRN service require the user to create the following:

- a video SAP
- a video interface

A video SAP is similar to other SAPs in the system in that QoS and filter policies can be associated with the SAP on ingress (traffic leaving the ISA and ingressing the system) and egress (traffic leaving system and entering the ISA).

The video SAP is associated with a video group. Channels are also associated with a video group which is what establishes the link between what channels can be referenced through the video SAP. The multicast information policy associated with the service is where the channel to video group association is defined.

A single VPRN service can be used to support all video applications. However, it is possible to use two different VPRNs for a FCC/RET service; one VPRN is configured with a FCC/RET server to process FCC/RET requests, and a second VPRN is configured to support only the buffering of the multicast video content. In this case, within the FCC/RET server VPRN, the multicast service of the VPRN (service ID) must be specified.

Another important video services option in a multicast information policy is the administrative bandwidth defined for the channel. Many video applications use the bandwidth to determine if sufficient ISA egress bandwidth exists to service or drop a service request.

Example: Video interface configuration (MD-CLI)

```
[ex:/configure service ies "1"]
A:admin@node-2# info
customer "1"
video-interface "video-100" {
  multicast-service 5
  video-sap {
    video-group-id 4
  }
  address 10.1.1.254/8 { }
  address 192.168.0.254/8 { }
}
```

Example: Video interface configuration (classic CLI)

```
A:node-2>config>service>ies# info
-----
      shutdown
      video-interface "video-100" create
      shutdown
      video-sap 4
      exit
      address 10.1.1.254/8
      address 192.168.0.254/8
      multicast-service 5
      exit
-----
```

6.4.1.4 Basic multicast information policy configuration

Multicast information policies are used by the video applications to define multicast channel attributes and video policies which contains application-specific configuration for a video interface IP address.

It is within the multicast information policy bundles, channels and source-overrides that a video group is assigned to a channel. The video group association is inherited from the more general construct unless it is explicitly disabled.

The administrative bandwidth for channels at the bundle, channel or source-override level is also defined in the multicast information policy. Video applications use the administrative bandwidth here when a channel rate estimate is needed.

A video policy is defined within the multicast information policy for a specific video interface IP address. The IP address for the video policy is the key value that associates it with a specific video interface IP address within a service associated with overall multicast information policy.

See the *7450 ESS, 7750 SR, and VSR Triple Play Service Delivery Architecture Guide* for CLI command descriptions and syntax usage information to configure multicast info policies.

The following example displays a multicast information policy configuration.

Example: MD-CLI

```
[ex:/configure multicast-management]
A:admin@node-2# info
  multicast-info-policy "ies100" {
    bundle "10.5.6.140" {
      admin-bw 8000
      video {
        local-rt-server true
        rt-buffer-size 3000
        video-group 1
      }
      channel start 10.5.6.140 end 10.5.6.140 {
      }
    }
  }
  bundle "10.5.6.241-10.5.6.243" {
    admin-bw 12000
    video {
      rt-buffer-size 4000
      video-group 1
    }
    channel start 10.5.6.241 end 10.5.6.243 {
    }
  }
}
```

```
}
}
```

Example: classic CLI

```
A:node-2>config>mcast-mgmt># info
-----
multicast-info-policy "ies100" create
  bundle "10.5.6.140" create
    admin-bw 8000
    video
      video-group 1
      local-rt-server
      rt-buffer-size 3000
    exit
    channel "10.5.6.140" "10.5.6.140" create
    exit
  exit
  bundle "default" create
  exit
  bundle "5.6.241-5.6.243" create
    admin-bw 12000
    video
      video-group 1
      rt-buffer-size 4000
    exit
    channel "10.5.6.241" "10.5.6.243" create
    exit
  exit
exit
-----
```

6.4.2 Sample configurations

The following example displays the configuration of the VQM with packet selection.

Example: VQM configuration with packet selection (MD-CLI)

```
[ex:/configure multicast-management]
A:admin@node-2# info
  multicast-info-policy "vqm" {
    bundle "ixia" {
      channel start 10.5.5.6 end 10.5.5.7
      admin-bw 2000
      video {
        rt-buffer-size 1000
        video-group 4
        stream-selection {
          source1 192.168.2.1
          intf1 "ineo-ingress1"
          source2 192.168.2.1
          intf2 "ineo-ingress2"
        }
        analyzer {
          alarms {
            cc-error true
            pat-repetition tnc 400 qos 600 poa 700
            pat-syntax
            pcr-repetition tnc 400 qos 600 poa 700
            pid-pmt-unref
            pmt-repetition tnc 2300 qos 2500 poa 2700
          }
        }
      }
    }
  }
}
```

```

        pmt-syntax
        scte-35
        tei-set
        ts-sync-loss
        vid-pid-absent 5000
        non-vid-pid-absent 5000 {
    }
}
}
}
source-override 192.168.2.1 {
}
}
}

[ex:/configure service ies "300"]
A:admin@node-2# info
admin-state enable
description "Default IES description for service ID 300"
customer "1"
vpn-id 300
video-interface "video-300" {
    admin-state enable
    video-sap {
        video-group-id 1
    }
    channel 10.5.5.6 source 192.168.2.1 {
        channel-name "Ineoquest-1"
        zone-channel 10.5.5.6 source 10.20.0.1 {
            adi-channel-name "Ineoquest-1-1"
        }
    }
}
}

[ex:/configure router "Base"]
A:admin@node-2# info
ecmp 2
multicast-info-policy "vqm"
interface "ineo-ingress1" {
    port 3/2/12
    ingress {
        filter {
            ip "100"
        }
    }
    ipv4 {
        primary {
            address 10.200.16.1
            prefix-length 24
        }
    }
}
interface "ineo-ingress2"
port 5/1/1
ingress
filter {
    ip 200
}
}
ipv4 {
    primary {
        address 10.200.17.1
        prefix-length 24
    }
}
}

```

```

}
interface "ixia-egress"
  port 3/2/15
  ipv4 {
    primary {
      address 10.200.15.1
      prefix-length 24
    }
  }
}
interface "system"
  ipv4 {
    primary {
      address 10.20.1.1
      prefix-length 32
    }
  }
}
static-routes {
  route 192.168.2.1/32 route-type multicast {
    next-hop "10.200.16.2" {
      admin-state enable
    }
    next-hop "10.200.17.2" {
      admin-state enable
    }
  }
}
}

[ex:/configure router "Base" igmp]
A:admin@node-2# info
...
interface "video-300-D" {
  static {
    group 10.5.5.6 {
      source 192.168.2.1 { }
    }
  }
}
interface "video-300-D2" {
  static {
    group 10.5.5.6 {
      source 192.168.2.1
    }
  }
}
interface "ixia-egress" {
  static {
    group 10.5.5.6 {
      source 10.20.0.1
    }
  }
}

[ex:/configure router "base" pim]
A:admin@node-2# info
interface "video-300" {
  admin-state enable
}
interface "ineo-ingress1" {
  admin-state enable
  multicast-senders always
}
interface "ineo-ingress2" {

```



```

    admin-state enable
    multicast-senders always
  }
  rp {
    ipv4 {
      bsr-candidate {
        admin-state enable
      }
      rp-candidate {
        admin-state enable
      }
    }
  }
}

[ex:/configure isa]
A:admin@node-2# info
  video-group 4 {
    admin-state enable
    analyzer
    stream-selection true
    mda 3/1 { }
  }

```

Example: VQM configuration with packet selection (classic CLI)

```

A:node-2>config>mcast-mgmt># info
-----
  multicast-info-policy "vqm" create
    bundle "ixia" create
      channel "10.5.5.6" "10.5.5.7" create
        admin-bw 20000
        video
          video-group 4
            rt-buffer-size 1000
            analyzer
            alarms
              cc-error
              pat-repetition tnc 400 qos 600 poa 700
              pat-syntax
              pid-pmt-unref
              pmt-repetition tnc 2300 qos 2500 poa 2700
              pmt-syntax
              vid-pid-absent 5000
              non-vid-pid-absent 5000
              pcr-repetition tnc 400 qos 600 poa 700
              scte-35
              tei-set
              ts-sync-loss
            exit
          exit
          stream-selection source1 192.168.2.1 intf1 "ineo-ingress1"
          source2 192.168.2.1 intf2 "ineo-ingress2"
            exit
            source-override "192.168.2.1" create
            exit
          exit
        exit
      bundle "default" create
      exit
    exit
  -----
A:node-2>config>service# info
-----adm-----

```

```

customer 1 create
  description "Default customer"
exit
ies 300 customer 1 vpn 300 create
  description "Default Ies description for service id 300"
  video-interface "video-300" create
    video-sap 4
    exit
    address 10.20.255.254/16
    channel 10.5.5.6 source 192.168.2.1 channel-name "Ineoquest-1"
    zone-channel 10.5.5.6 source 10.20.0.1 adi-channel-name "Ineoquest-
1-1"
    exit
    adi
    exit
    no shutdown
  exit
  service-name "XYZ Ies 300"
  no shutdown
exit
-----
A:node-2>config>router# info
-----
#-----
echo "IP Configuration"
#-----
interface "ineo-ingress1"
  address 10.200.16.1/24
  port 3/2/12
  ingress
  filter ip 100
  exit
exit
interface "ineo-ingress2"
  address 10.200.17.1/24
  port 5/1/1
  ingress
  filter ip 200
  exit
exit
interface "ixia-egress"
  address 10.200.15.1/24
  port 3/2/15
  exit
interface "system"
  address 10.20.3.1/32
  exit
  ecmp 2
  multicast-info-policy "vqm"
  static-route 192.168.2.1/32
    next-hop 10.200.16.2
    no shutdown
  exit
  next-hop 10.200.17.2
  no shutdown
  exit
exit
#-----
echo "IGMP Configuration"
#-----
igmp
  interface "video-300-D"
  static
  group 10.5.5.6

```

```

        source 192.168.2.1
    exit
    exit
exit
interface "video-300-D2"
    static
        group 10.5.5.6
        source 192.168.2.1
    exit
    exit
exit
interface "ixia-egress"
    static
        group 10.5.5.6
        source 10.20.0.1
    exit
    exit
exit
exit
#-----
echo "PIM Configuration"
#-----
    pim
        rpf-table rtable-m
        interface "video-300"
        exit
        interface "ineo-ingress1"
            multicast-senders always
        exit
        interface "ineo-ingress2"
            multicast-senders always
        exit
        rp
            static
            exit
            bsr-candidate
                shutdown
            exit
            rp-candidate
                shutdown
            exit
        exit
    exit
-----
A:node-2>config>isa# info
-----
    video-group 4 create
        analyzer
        stream-selection
        primary 3/1
        no shutdown
    exit
-----

```

The following output displays configurations of VQM without packet selection.

Example: VQM configuration without packet selection (MD-CLI)

```

[ex:/configure service ies "300"]
A:admin@node-2# info
    admin-state enable
    description "Default IES description for service ID 300"
    customer "1"

```

```

vpn-id 300
interface "linux-egress" {
  sap 3/2/17 {
    description "sap-300-10.10.34.228"
  }
  ipv4 {
    primary {
      address 10.10.34.228
    }
  }
}
interface "linux-ingress" {
  sap 3/2/17 {
    description "sap-300-10.10.33.228"
  }
  ipv4 {
    primary {
      address 10.10.33.228
    }
  }
}
video-interface "video-300" {
  admin-state enable
  video-sap {
    video-group-id 2
  }
  address 10.20.13.1/24 { }
}
channel 10.5.5.6 source 192.168.2.1 {
  channel-name "A2-SP3"
  zone-channel 10.5.5.6 source 10.20.13.2 {
    adi-channel-name "A2-SP3-1"
  }
}
}

[ex:/configure router "Base"]
A:admin@node-2# info
...
multicast-info-policy "A-server"
interface "system" {
  ipv4 {
    primary {
      address 10.20.1.1
      prefix-length 32
    }
  }
}
...
static-routes {
  route 128.251.33.0/24 {
    next-hop "10.10.33.229" {
      admin-state enable
    }
  }
  route 192.168.2.0/24
  next-hop "10.10.33.229" {
    admin-state enable
  }
}
}

[ex:/configure router "Base" igmp]
A:admin@node-2# info

```

```

interface "video-300-D" {
    static {
        group 10.5.5.6 {
            source 192.168.2.1 { }
        }
    }
}
interface "linux-egress" {
    static {
        group 10.5.5.6 {
            source 10.20.13.2
        }
    }
}

[ex:/configure router "base" pim]
A:admin@node-2# info
interface "linux-ingress" {
    admin-state enable
    hello-interval 0
    multicast-senders always
}
interface "linux-egress" {
    admin-state enable
    hello-interval 0
}
apply-to all
}
rp {
    static
    bsr-candidate {
        admin-state enable
    }
    rp-candidate {
        admin-state enable
    }
}
}

[ex:/configure isa]
A:admin@node-2# info
video-group 2 {
    admin-state enable
    analyzer
    mda 2/1 { }
}

[ex:/configure multicast-management]
A:admin@node-2# info
multicast-info-policy "A-server" {
    bundle "LiveTV" {
        channel start 10.5.6.243 end 10.5.6.243"
        admin-bw 3000
        video {
            rt-buffer-size 1000
            video-group 2
        }
    }
    channel start 10.5.5.6 end 10.5.5.6"
    admin-bw 5000
    video {
        rt-buffer-size 1000
        video-group 2
        analyzer {
            alarms {

```



```

        service-name "XYZ Ies 300"
        no shutdown
    exit
-----
A:node-2>config>router# info
-----
#-----
echo "IP Configuration"
#-----
        interface "system"
            address 10.20.1.1/32
        exit
        multicast-info-policy "A-server"
#-----
echo "Static Route Configuration"
#-----
        static-route 128.251.33.0/24
            next-hop 10.10.33.229
            no shutdown
        exit
        static-route 192.168.2.0/24
            next-hop 10.10.33.229
            no shutdown
        exit
    exit
#-----
echo "IGMP Configuration"
#-----
        igmp
            interface "video-300-D"
                static
                    group 10.5.5.6
                    source 192.168.2.1
                exit
            exit
        exit
        interface "linux-egress"
            static
                group 10.5.5.6
                source 10.20.13.2
            exit
        exit
    exit
#-----
echo "PIM Configuration"
#-----
        pim
            interface "linux-ingress"
                hello-interval 0
                multicast-senders always
            exit
            interface "linux-egress"
                hello-interval 0
            exit
        apply-to all
        rp
            static
            exit
        bsr-candidate
            shutdown
        exit
        rp-candidate

```

```

        shutdown
    exit
    exit
exit
-----

A:node-2>config>isa# info
-----
    video-group 2 create
        analyzer
        primary 2/1
        no shutdown
    exit
-----

A:node-2>config>mcast-mgmt># info
-----
    multicast-info-policy "A-server" create
        bundle "LiveTv" create
            channel "10.5.6.243" "10.5.6.243" create
                admin-bw 3000
                video
                    video-group 2
                    rt-buffer-size 1000
                exit
            exit
            channel "10.5.5.6" "10.5.5.6" create
                admin-bw 5000
                video
                    video-group 2
                    rt-buffer-size 1000
                    analyzer
                    alarms
                        cc-error
                        pat-repetition tnc 200 qos 400 poa 600
                        pat-syntax
                        pid-pmt-unref
                        pmt-repetition
                        pmt-syntax
                        vid-pid-absent 1000
                        non-vid-pid-absent 1000
                        pcr-repetition tnc 200 qos 400 poa 600
                        scte-35
                        tei-set
                        ts-sync-loss
                        report-alarm severity tnc
                    exit
                exit
            exit
            source-override "192.168.33.37" create
            exit
        exit
    exit
    bundle "default" create
    exit
    bundle "mp2ts-ads" create
        channel "10.4.5.1" "10.4.5.254" create
            admin-bw 5000
            video
                video-group 2
                rt-buffer-size 1000
            exit
        exit
    exit
exit

```



```
exit
-----
```

6.5 Configuring RET/FCC video components with CLI

This section provides information to configure RET/FCC using the command line interface.

6.5.1 Configuring RET/FCC video features in the CLI

The following sections provide configuration examples for the RET client, RET server and FCC server.

6.5.1.1 Configuring the RET server

This section provides an example configuration for the RET server. The configuration example has the following assumptions:

- a single ISA-MS in slot 2/1 in video group 1
- a single channel 192.0.2.1 within multicast bundle "b1" with an administrative bandwidth of 2700 kb/s defined in **multicast-info-policy** *policy-name*
- a retransmission buffer for the channel set to 300 milliseconds
- RET rate is 5% of nominal
- local RET server address is 10.3.3.3 and destination port is UDP 4096

The first step is to configure the video group, including the RET server, and the ISA-MS hardware. The **local-rt-server** command in the example enables the local RET server on the video group.

Example: Video group with RET server configuration (MD-CLI)

```
[ex:/configure isa]
A:admin@node-2# info
  video-group 1 {
    admin-state enable
    local-rt-server true
    mda 2/1 { }
  }

[ex:/configure card 2 mda 1]
A:admin@node-2# info
  mda-type isa-ms-v
```

Example: Video group with RET server configuration (classic CLI)

```
A:node-2>config>isa# info
-----
  video-group 1 create
    local-rt-server
    primary 2/1
    no shutdown
  exit
-----
```

```
A:node-2>config>card>mda# info
-----
mda-type isa-ms
-----
```

The channel command options for 192.0.2.1 are configured in the multicast management policy. The channel configuration includes the administrative bandwidth and the channel's association with video group 1.

Example: Multicast information policy configuration (MD-CLI)

```
[ex:/configure multicast-management multicast-info-policy "ies100"]
A:admin@node-2# info
  bundle "bl" {
    admin-bw 2700
    video {
      local-rt-port 4096
      rt-buffer-size 300
      video-group 2
    }
    channel start 192.0.2.1 end 192.0.2.1 {
    }
  }
  video-policy {
    video-interface 10.3.3.3 {
      rt-rate 5
      hd {
        local-rt-server true
      }
      pip {
        local-rt-server true
      }
      sd {
        local-rt-server true
      }
    }
  }
}
```

Example: Multicast information policy configuration (classic CLI)

```
A:node-2>config>mcast-mgmt>mcast-info-plcy# info
-----
bundle "default" create
  local-rt-port 4096
exit
bundle "b1" create
  admin-bw 2700
  video
    video-group 2
    rt-buffer-size 300
  exit
  channel "192.0.2.1" "192.0.2.1" create
  exit
exit
video-policy
  video-interface 10.3.3.3 create
    rt-rate 5
    hd
      local-rt-server
    exit
    sd
      local-rt-server
```

```

        exit
        pip local-rt-server
        exit
    exit
exit
-----

```

The **local-rt-port** command in the bundle defines the destination UDP port used to reach the local RET server on the service where the multicast information policy is applied.

In the classic CLI, the RET server port can only be defined in the bundle "default" and applies for all bundles in the policy. If no value is specified, the default is used.

In the example, in the bundle "b1" the **local-rt-server** command enables the RET server for all channels in the bundle, and the **rt-buffer-size** command sets the retransmission buffer for all channels in the bundle to 300 milliseconds.

In the video policy in the example, the **local-rt-server** commands for the video interface 10.3.3.3 enables the RET server on that interface for all channel types "hd" (High Definition), "sd" (Standard Definition) and "pip" (Picture-in-Picture). The **rt-rate** command in the policy indicates that the retransmission rate is 5% of the nominal rate for all channel types; individual rates can be defined if needed.

In the following example, for the RET server in an IES or VPRN service instance and router configuration, the following commands are used to complete the RET server configuration:

1. associate the service with the multicast information policy
2. create the video interface "vi" and assign IP address 10.3.3.3
3. create video SAP and associate it with video group 1
4. create a static IGMP join on video-interface "vi" for the channel 192.0.2.1

In the following example, the services available on the video interface address 10.3.3.3 are defined in the video policy in which the RET server is enabled.

Example: Video interface to service association and IGMP and PIM configuration (MD-CLI)

```

[ex:/configure service ies "100"]
A:admin@node-2# info
  video-interface "vi" {
    admin-state enable
    video-sap {
      video-group-id 1
    }
    address 10.3.3.3/32 { }
  }
}

[ex:/configure router "base"]
A:admin@node-2# info
  multicast-info-policy "ies100"
  igmp
    interface "vi"
      static
        group 192.0.2.1
        starg
      }
    }
  }
}
pim

```

```

        interface "vi" {
    }

```

Example: Video interface to service association and IGMP and PIM configuration (classic CLI)

```

A:node-2>config>service>ies# info
-----
        video-interface "vi" create
            video-sap 1
            exit
            address 10.3.3.3/32
            no shutdown
        exit

A:node-2>config>router# info
-----
...
multicast-info-policy "ies100"
igmp
    interface "vi"
        static
            group 192.0.2.1
            starg
        exit
    exit
pim
    interface "vi"
        exit
    exit
-----

```

6.5.1.2 Configuring the FCC server

This section provides an example configuration for the FCC server. The configuration example has the following assumptions:

- a single ISA-MS in slot 2/1 in video group 1
- a single channel 192.0.2.1 within multicast bundle "b1" with an administrative bandwidth of 8000 kb/s defined using the **multicast-info-policy** command
- FCC mode burst with a rate 130% of nominal for HD, 200% for SD, and disabled for PIP
- local FCC server address 10.3.3.3 and destination port UDP 4098

The first step in the configuration is to configure video group 1 enabling the FCC server and the ISA-MS hardware. The **fcc-server** command in the example enables the FCC server on the video group.

Example: Video group with FCC server configuration (MD-CLI)

```

[ex:/configure isa]
A:admin@node-2# info
    video-group 1 {
        admin-state enable
        fcc-server true
        mda 2/1 { }
    }

```

```
[ex:/configure card 2 mda 1]
A:admin@node-2# info
mda-type isa-ms
```

Example: Video group with FCC server configuration (classic CLI)

```
A:node-2>config>isa# info
-----
video-group 1 create
  fcc-server
  primary 2/1
  no shutdown
exit
-----

A:node-2>config>card>mda# info
-----
mda-type isa-ms
-----
```

The channel parameters for 192.0.2.1 are configured in the **multicast-info-policy** configuration. The channel configuration includes the administrative bandwidth and the channel's association with video group 1.

Example: FCC server configuration for a multicast information policy (MD-CLI)

```
[ex:/configure multicast-management multicast-info-policy "ies100"]
A:admin@node-2# info
bundle "bl" {
  admin-bw 8000
  video {
    fcc-server true
    fcc-channel-type hd
    local-fcc-port 4098
    video-group 1
  }
  channel start 192.0.2.1 end 192.0.2.1 {
  }
}
video-policy {
  video-interface 10.3.3.3 {
    rt-rate 5
    hd {
      fcc-server {
        mode burst
      }
      fcc-burst 30
    }
  }
  pip
}
sd {
  fcc-server {
    mode burst
  }
  fcc-burst 100
}
}
```

Example: FCC server configuration for a multicast information policy (classic CLI)

```

A:node-2>config>mcast-mgmtmcast-info-plcy# info
-----
bundle "default" create
  local-fcc-port 4098
exit
bundle "b1" create
  admin-bw 8000
  video
    video-group 1
    fcc-server
    fcc-channel-type hd
  exit
  channel "192.0.2.1" "192.0.2.1" create
  exit
exit
video-policy
  video-interface 10.3.3.3 create
    rt-rate 5
    hd
      fcc-server mode burst
      fcc-burst 30
    exit
    sd
      fcc-server mode burst
      fcc-burst 100
    exit
    pip
      no fcc-server
    exit
  exit
exit
-----

```

The **local-fcc-port** command in the bundle defines the destination UDP port used to reach the FCC server on the service where the multicast information policy is applied.

In the classic CLI, the FCC server port can only be defined using the **local-fcc-port** command in the bundle "default" and applies for all bundles in the policy. If no value is specified, the default is used.

In the bundle "b1", the **fcc-server** command enables the FCC server for all channels in the bundle, and the **fcc-channel-type hd** command sets the channel type for all channels in the bundle to "hd" (High Definition).

In the preceding video policy example, the **fcc-server** commands for the video interface 10.3.3.3 enables the FCC server on that interface for all channel types "hd" (High Definition), "sd" (Standard Definition). There is no FCC server enabled for the Picture-in-Picture (PIP) channels on the video interface.

The **fcc-burst** command in the policy indicates that the burst rate over the nominal rate for the channel type, that is, HD at 130% (30% over nominal) and SD at 200% (100% over nominal).

In the following example, for the FCC server in an IES or VPRN service instance and router configuration, the following commands are used to complete the FCC server configuration:

1. associate the service with the multicast information policy
2. create the video interface "vi" and assign IP address 10.3.3.3
3. create video SAP and associate it with video group 1
4. create a static IGMP join on video-interface "vi" for the channel 192.0.2.1

The services available on the video interface address are defined in the video policy in which the FCC server was enabled.

Example: Video service association with IGMP and PIM configuration (MD-CLI)

```
[ex:/configure service ies "100"]
A:admin@node-2# info
  video-interface "vi" {
    admin-state enable
    video-sap {
      video-group-id 1
    }
    address 10.4.4.4/32 { }
  }

[ex:/configure router "base"]
A:admin@node-2# info
  ....
  multicast-info-policy "ies100"
  igmp
    interface "vi"
      static
        group 192.0.2.1
        starg
      }
  }
  pim
  interface "vi" {
  }
  ...
```

Example: Video service association with IGMP and PIM configuration (classic CLI)

```
A:node-2>config>service>ies# info
-----
  video-interface "vi" create
  video-sap 1
  exit
  address 10.4.4.4/32
  no shutdown
  exit
-----

A:node-2>config>router# info
-----
  ...
  multicast-info-policy "ies100"
  igmp
    interface "vi"
      static
        group 192.0.2.1
        starg
      exit
    exit
  exit
  pim
  interface "vi"
  exit
  ...
```

6.5.1.3 Logging and accounting collection for video statistics

The following example shows logging and accounting configuration for collecting video statistics. Accounting requires the platforms support no more than 256,000 sessions and the FCC/RET timer must be set to 5 minutes. If these conditions are not met, the accounting record may not be accurate.

Example: Logging and accounting configuration for video statistics (MD-CLI)

```
[ex:/configure log]
A:admin@node-2# info
    file 1
        compact-flash-location {
            primary cf3:
        }
    accounting-policy 1 {
        admin-state disable
        collection-interval 5
        record video
        destination {
            file "1"
        }
    }
}
```

Example: Logging and accounting configuration for video statistics (classic CLI)

```
A:node-2>config>log# info
-----
    file-id 1
        location cf3:
    exit
    accounting-policy 1
        shutdown
        record video
        collection-interval 5
        to file 1
    exit
-----
```

The following example shows how to enable logging and accounting for a service to collect statistics for a particular service and video interface. This example refers to the accounting policy created in the preceding example. When administratively enabled, this starts recording of statistics. The statistics are written in an act-collect directory and administratively enabling the accounting policy moves the recorded file to the act directory.

Example: Enabling video statistics collection for a service (MD-CLI)

```
[ex:/configure service ies "300"]
A:admin@node-2# info
    video-interface "vi" {
        admin-state enable
        accounting-policy 1 {
            admin-state enable
        }
    }
}
```


Example: Enabling video statistics collection for a service (classic CLI)

```
A:node-2>config>service>ies# info
video-interface "vi" create
  accounting-policy "1"
  no shutdown
  exit
no shutdown
exit
```

7 Network Address Translation

7.1 Terminology

- **BNG subscriber**

This is a broader term than the ESM Subscriber, independent of the platform on which the subscriber is instantiated. It includes ESM subscribers on 7750 SR as well as subscribers instantiated on third party BNGs. Some of the NAT functions, such as Subscriber Aware Large Scale NAT44 utilizing standard RADIUS attribute work with subscribers independently of the platform on which they are instantiated.

- **deterministic NAT**

This is a mode of operation where mappings between the NAT subscriber and the outside IP address and port-range are allocated at the time of configuration. Each subscriber is permanently mapped to an outside IP and a dedicated port block. This dedicated port block is referred to as deterministic port block. Logging is not needed as the reverse mapping can be obtained using a known formula. The subscriber's ports can be expanded by allocating a dynamic port block in case that all ports in deterministic port block are exhausted. In such case logging for the dynamic port block allocation/de-allocation is required.

- **Enhanced Subscriber Management (ESM) subscriber**

This is a host or a collection of hosts instantiated in 7750 SR Broadband Network Gateway (BNG). The ESM subscriber represents a household or a business entity for which various services with committed Service Level Agreements (SLA) can be delivered. NAT function is not part of basic ESM functionality.

- **L2-Aware NAT**

In the context of 7750 SR platform combines Enhanced Subscriber Management (ESM) subscriber-id and inside IP address to perform translation into a unique outside IP address and outside port. This is in contrast with classical NAT technique where only inside IP is considered for address translations. Because the subscriber-id alone is sufficient to make the address translation unique, L2-Aware NAT allows many ESM subscribers to share the same inside IP address. The scalability, performance and reliability requirements are the same as in LSN.

- **Large Scale NAT (LSN)**

This refers to a collection of network address translation techniques used in service provider network implemented on a highly scalable, high performance hardware that facilitates various intra and inter-node redundancy mechanisms. The purpose of LSN semantics is to make delineation between high scale and high performance NAT functions found in service provider networks and enterprise NAT that is usually serving much smaller customer base at smaller speeds. The following NAT techniques can be grouped under the LSN name:

- Large Scale NAT44 or Carrier Grade NAT (CGN)
- DS-Lite
- NAT64

Each distinct NAT technique is referred to by its corresponding name (Large Scale NAT44 [or CGN], DS-Lite and NAT64) with the understanding that in the context of 7750 SR platform, they are all part of LSN (and not enterprise based NAT).

Large Scale NAT44 term can be interchangeably used with the term Carrier Grade NAT (CGN) which in its name implies high reliability, high scale and high performance. These are again typical requirements found in service provider (carrier) network.

- **NAT RADIUS accounting**

This is the reporting (or logging) of address translation related events (port-block allocation/de-allocation) via RADIUS accounting facility. NAT RADIUS accounting is facilitated via regular RADIUS accounting messages (start/interim-update/stop) as defined in RFC 2866, *RADIUS Accounting*, with NAT specific VSAs.

- **NAT RADIUS accounting**

This can be used interchangeably with the term NAT RADIUS logging.

- **NAT subscriber**

In NAT terminology, a NAT subscriber is an inside entity whose true identity is hidden from the outside. There are a few types of NAT implementation in 7750 SR and subscribers for each implementation are defined as follows:

- **Large Scale NAT44 (or CGN)**

- The subscriber is an inside IPv4 address.

- **L2-Aware NAT**

- The subscriber is an ESM subscriber which can spawn multiple IPv4 inside addresses.

- **DS-Lite**

- The subscriber in DS-Lite can be identified by the CPE's IPv6 address (B4 element) or an IPv6 prefix. The selection of address or prefix as the representation of a DS-Lite subscriber is configuration dependent.

- **NAT64**

- The subscriber is an IPv6 prefix.

- **non-deterministic NAT**

This is a mode of operation where all outside IP address and port block allocations are made dynamically at the time of subscriber instantiation. Logging in such case is required.

- **port block**

This is collection of ports that is assigned to a subscriber. A deterministic LSN subscriber can have only one deterministic port block that can be extended by multiple dynamic port blocks. Non-deterministic LSN subscriber can be assigned only dynamic port blocks. All port blocks for a LSN subscriber must be allocated from a single outside IP address.

- **port-range**

This is a collection of ports that can spawn multiple port blocks of the same type. For example, deterministic port-range includes all ports that are reserved for deterministic consumption. Similarly dynamic port-range is a total collection of ports that can be allocated in the form of dynamic port blocks. Other types of port-ranges are well-known ports and static port forwards.

7.2 Network Address Translation (NAT) overview

The 7750 SR supports Network Address (and port) Translation (NAPT) to provide continuity of legacy IPv4 services during the migration to native IPv6. By equipping the multi-service ISA (MS ISA) in an IOM4-e, IOM4-e-B, IOM4-e-HS, or in a 7750 SR-1e, 7750 SR-2e, or 7750 SR-3e (IOM-e), the 7750 SR can operate in two different modes, known as:

- Large Scale NAT
- Layer 2-Aware NAT

These two modes both perform source address and port translation as commonly deployed for shared Internet access. The 7750 SR with NAT is used to provide consumer broadband or business Internet customers access to IPv4 Internet resources with a shared pool of IPv4 addresses, such as may occur around the forecast IPv4 exhaustion. During this time it, is expected that native IPv6 services are still growing and a significant amount of Internet content remains IPv4.

7.2.1 Principles of NAT

Network Address Translation devices modify the IP headers of packets between a host and server, changing some or all of the source address, destination address, source port (TCP/UDP), destination port (TCP/UDP), or ICMP query ID (for ping). The 7750 SR in both NAT modes performs Source Network Address and Port Translation (S-NAPT). S-NAPT devices are commonly deployed in residential gateways and enterprise firewalls to allow multiple hosts to share one or more public IPv4 addresses to access the Internet. The common terms of inside and outside in the context of NAT refer to devices inside the NAT (that is behind or masqueraded by the NAT) and outside the NAT, on the public Internet.

TCP/UDP connections use ports for multiplexing, with 65536 ports available for every IP address. Whenever many hosts are trying to share a single public IP address there is a chance of port collision where two different hosts may use the same source port for a connection. The resultant collision is avoided in S-NAPT devices by translating the source port and tracking this in a stateful manner. All S-NAPT devices are stateful in nature and must monitor connection establishment and traffic to maintain translation mappings. The 7750 SR NAT implementation does not use the well-known port range (1 to 1023).

In most circumstances, S-NAPT requires the inside host to establish a connection to the public Internet host or server before a mapping and translation occurs. With the initial outbound IP packet, the S-NAPT knows the inside IP, inside port, remote IP, remote port and protocol. With this information the S-NAPT device can select an IP and port combination (referred to as outside IP and outside port) from its pool of addresses and create a unique mapping for this flow of data.

Any traffic returned from the server uses the outside IP and outside port in the destination IP/port fields – matching the unique NAT mapping. The mapping then provides the inside IP and inside port for translation.

The requirement to create a mapping with inside port and IP, outside port and IP and protocol generally prevents new connections to be established from the outside to the inside as may occur when an inside host needs to be a server.

7.2.2 Traffic load balancing

NAT traffic in SR OS is distributed over ISAs and ESAs within each NAT group. As a result, NAT capacity grows incrementally by adding more ISAs and ESAs to the system while each ISA or ESA participates equally in load sharing.

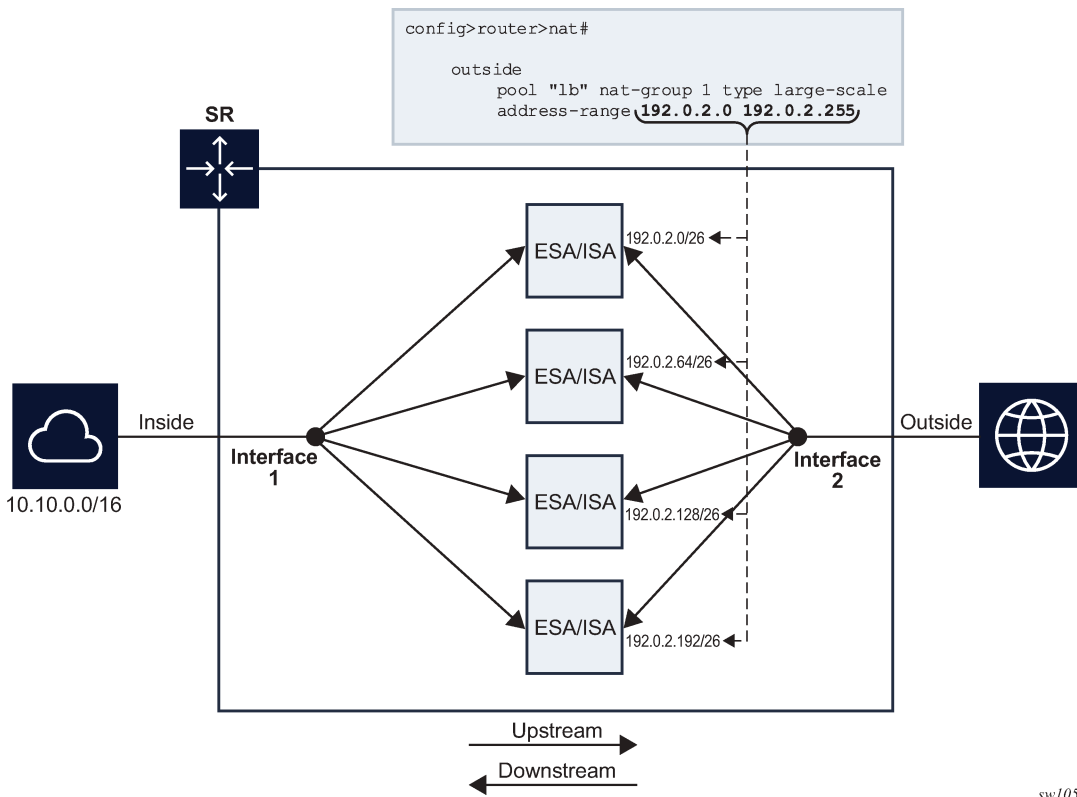
SR OS load balancing mechanisms in CGN (LSN44, DS-Lite, and NAT64) differ in the upstream and downstream directions, they are independent and unaware of each other,

- In the upstream direction, traffic is load balanced based on source IPv4 or IPv6 addresses or subnets.
- In the downstream direction, outside IP address ranges (NAT pool address ranges) are microneted (divided into smaller subnets), and these micronets are assigned to individual ISAs or ESAs in a balanced way. Downstream traffic is assigned to each ISA or ESA based on the micronets.

Figure 57: Load balancing over ISAs and ESAs shows traffic load balancing within SR OS. In the upstream direction, traffic is hashed based on the source IP addresses or subnets from the 10.10.0.0/16 range. A sample of 64000 source IP addresses guarantees equal load distribution.

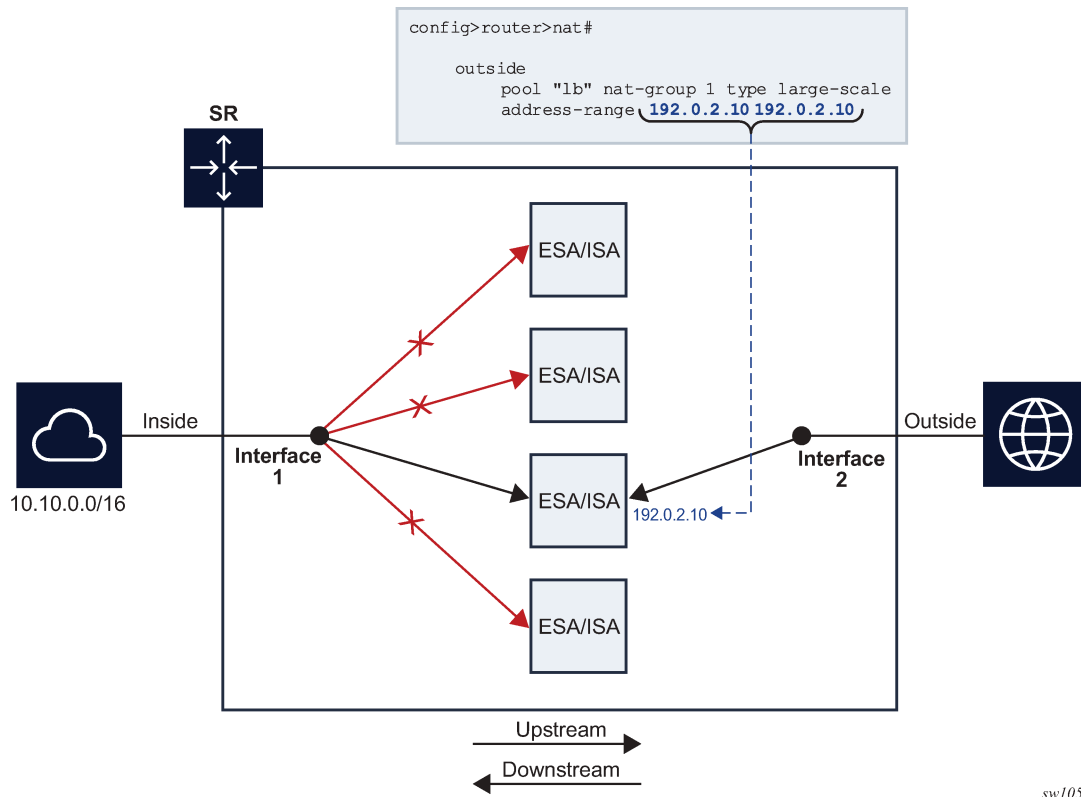
In this example, in the downstream direction, a pool of 256 public addresses is divided into four equal subnets and each subnet is assigned to one ISA or ESA, each ISA or ESA is serving 64 public IP addresses.

Figure 57: Load balancing over ISAs and ESAs



If there are not enough IP addresses on the inside and outside in relation to the number of ISAs and ESAs, unequal load balancing and, in extreme cases, traffic blackholing can occur. **Figure 58: Traffic blackholing** shows an example of an extreme case, where a single IP address is assigned to a pool in a NAT group with four ISAs or ESAs. This single outside IP address can be assigned to a single ISA or ESA that serves downstream traffic. Upstream traffic is unaware of the downstream load distribution, so it sends traffic to all four ISAs and ESAs, and as a result this traffic is dropped at ISAs or ESAs that do not have the public IP address assigned.

Figure 58: Traffic blackholing



The operator is notified when the number of outside IP addresses in a pool is smaller than the number of ISAs or ESAs in the NAT group. The notification is sent in the form of a log.

```

3 2020/04/03 18:48:42.010 CEST MINOR: NAT #2015 Base Resource problem
"The address configuration for pool 'test.' causes one or more ISAs not getting an IP address"
4 2020/04/03 18:48:42.010 CEST MINOR: NAT #2014 Base Resource alarm raised
"The status of the NAT resource problem indication changed to true."
  
```

This configuration is permitted by the CLI, but a message is displayed directly in response to a pool activation.

```

configure router nat outside pool "test" no shutdown
INFO: BB #1221 The address configuration for this pool causes one or more members not getting
an IP address - Router 'Base', pool 'test'
  
```

The load balancing mechanism in L2-Aware NAT relies on a different algorithm than CGN. In L2-Aware NAT, on the inside, traffic is distributed across ISAs and ESAs based on the resource utilization of each ISA or ESA. This load balancing mechanism is control plane driven, contrary to CGN which is forwarding plane driven (hashing is based purely on source IP addresses or subnets). In L2-Aware NAT, an ESM subscriber is directed to an ISA or ESA hosting a large number of subscribers, hosts, and port blocks, as an aggregate. In L2-Aware NAT, traffic is not blackholed when the number of outside IP addresses is smaller than the number of ISAs and ESAs in the pool within a single NAT group. Instead, the outside IP address is assigned to some of the ISAs or ESAs and the ESM host is directed to those. ISAs and ESAs without assigned outside IP addresses remains unused.

7.2.3 Application compatibility

Applications which operate as servers (such as HTTP, SMTP, and so on) or peer-to-peer applications can have difficulty when operating behind an S-NAPT because traffic from the Internet cannot reach the NAT without a mapping in place.

Different methods can be employed to overcome this, including:

- Port forwarding
- STUN support
- Application Layer Gateways (ALG)

The 7750 SR supports all three methods following the best-practice RFC for TCP (RFC 5382, *NAT Behavioral Requirements for TCP*) and UDP (RFC 4787, *Network Address Translation (NAT) Behavioral Requirements for Unicast UDP*). Port Forwarding is supported on the 7750 SR to allow servers which operate on well-known ports <1024 (such as HTTP and SMTP) to request the appropriate outside port for permanent allocation.

STUN is facilitated by the support of Endpoint-Independent Filtering and Endpoint-Independent Mapping (RFC 4787) in the NAT device, allowing STUN-capable applications to detect the NAT and allow inbound P2P connections for that specific application. Many new SIP clients and IM chat applications are STUN capable.

Application Layer Gateways (ALG) allows the NAT to monitor the application running over TCP or UDP and make appropriate changes in the NAT translations to suit. The 7750 SR has an FTP ALG enabled following the recommendation of the IETF BEHAVE RFC for NAT (RFC 5382).

Even with these three mechanisms some applications still experience difficulty operating behind a NAT. As an industry-wide issue, forums like UPnP the IETF, operator and vendor communities are seeking technical alternatives for application developers to traverse NAT (including STUN support). In many cases the alternative of an IPv6-capable application gives better long-term support without the cost or complexity associated with NAT.

7.3 Large Scale NAT

Large Scale NAT represents the most common deployment of S-NAPT in carrier networks today, it is already employed by mobile operators around the world for handset access to the Internet.

A Large Scale NAT is typically deployed in a central network location with two interfaces, the inside toward the customers, and the outside toward the Internet. A Large Scale NAT functions as an IP router and is located between two routed network segments (the ISP network and the Internet).

Traffic can be sent to the Large Scale NAT function on the 7750 SR using IP filters (ACL) applied to SAPs or by installing static routes with a next-hop of the NAT application. These two methods allow for increased flexibility in deploying the Large Scale NAT, especially those environments where IP MPLS VPN are being used in which case the NAT function can be deployed on a single PE and perform NAT for any number of other PE by simply exporting the default route.

The 7750 SR NAT implementation supports NAT in the base routing instance and VPRN, and through NAT traffic may originate in one VPRN (the inside) and leave through another VPRN or the base routing instance (the outside). This technique can be employed to provide customers of IP MPLS VPN with Internet access by introducing a default static route in the customer VPRN, and NATing it into the Internet routing instance.

As Large Scale NAT is deployed between two routed segments, the IP addresses allocated to hosts on the inside must be unique to each host within the VPRN. While RFC1918 private addresses have typically been used for this in enterprise or mobile environments, challenges can occur in fixed residential environments where a subscriber has existing S-NAPT in their residential gateway. In these cases the RFC 1918 private address in the home network may conflict with the address space assigned to the residential gateway WAN interface. Some of these issues are documented in *draft-shirasaki-nat444-isp-shared-addr-02*. Should a conflict occur, many residential gateways fail to forward IP traffic.

7.3.1 Port range blocks

The S-NAPT service on the 7750 SR BNG incorporates a port range block feature to address scalability of a NAT mapping solution. With a single BNG capable of hundreds of thousands of NAT mappings every second, logging each mapping as it is created and destroyed logs for later retrieval (as may be required by law enforcement) could quickly overwhelm the fastest of databases and messaging protocols. Port range blocks address the issue of logging and customer location functions by allocating a block of contiguous outside ports to a single subscriber. Instead of logging each NAT mapping, a single log entry is created when the first mapping is created for a subscriber and a final log entry when the last mapping is destroyed. This can reduce the number of log entries by 5000x or more. An added benefit is that as the range is allocated on the first mapping, external applications or customer location functions may be populated with this data to make real-time subscriber identification, instead of having to query the NAT as to the subscriber identity in real-time and possibly delay applications.

Port range blocks are configurable as part of outside pool configuration, allowing the operator to specify the number of ports allocated to each subscriber when a mapping is created. When a range is allocated to the subscriber, these ports are used for all outbound dynamic mappings and are assigned in a random manner to minimize the predictability of port allocations (*draft-ietf-tsvwg-port-randomization-05*).

Port range blocks also serve another useful function in a Large Scale NAT environment, and that is to manage the fair allocation of the shared IP resources among different subscribers.

When a subscriber exhausts all ports in their block, further mappings are prohibited. As with any enforcement system, some exceptions are allowed and the NAT application can be configured for reserved ports to allow high-priority applications access to outside port resources while exhausted by low priority applications.

7.3.1.1 Reserved ports and priority sessions

Reserved ports allows an operator to configure a small number of ports to be reserved for designated applications should a port range block be exhausted. Such a scenario may occur when a subscriber is unwittingly subjected to a virus or engaged in extreme cases of P2P file transfers. In these situations, instead of blocking all new mappings indiscriminately, the 7750 SR NAT application allows operators to nominate a number of reserved ports and then assign a 7750 SR forwarding class as containing high priority traffic for the NAT application. Whenever traffic reaches the NAT application which matches a priority session forwarding class, reserved ports are consumed to improve the chances of success. Priority sessions could be used by the operator for services such as DNS, web portal, e-mail, VoIP, and so on, to allow these applications even when a subscriber exhausted their ports.

7.3.1.2 Preventing port block starvation

7.3.1.2.1 Dynamic port block starvation in LSN

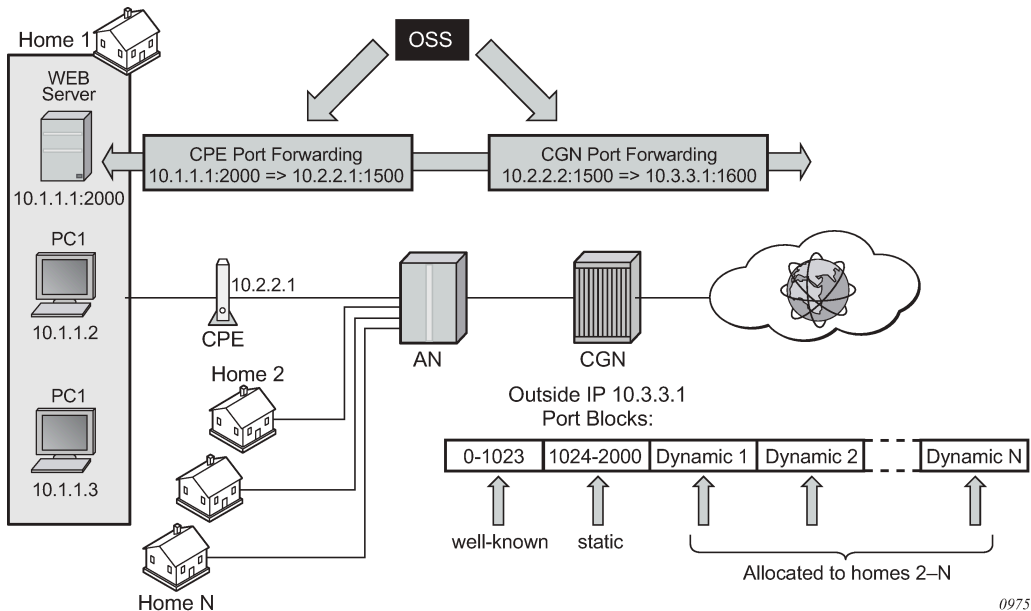
The outside IP address is always shared for the subscriber with a port forward (static or via PCP) and the dynamically allocated port block, insofar as the port from the port forward is in the range >1023. This behavior can lead to starvation of dynamic port blocks for the subscriber. An example for this scenario is shown in [Figure 59: Dynamic port block starvation in LSN](#).

- A static port forward for the WEB server in Home 1 is allocated in the CPE and the CGN. At the time of static port forward creation, no other dynamic port blocks for Home 1 exist (PCs are powered off).
- Assume that the outside IP address for the newly created static port forward in the CGN is 10.3.3.1.
- Over time dynamic port blocks are allocated for a number of other homes that share the same outside IP address, 10.3.3.1. Eventually those dynamic port block allocations exhaust all dynamic port block range for the address 10.3.3.1.
- After the dynamic port blocks are exhausted for outside IP address 10.3.3.1, a new outside IP address (for example, 10.3.3.2) is allocated for additional homes.

Eventually the PCs in Home 1 come to life and they try to connect to the Internet. Because of the dynamic port block exhaustion for the IP address 10.3.3.1 (that is mandated by static port forward – Web Server), the dynamic port block allocation fails and consequently, the PCs are not able to access the Internet. There is no additional attempt within CGN to allocate another outside IP address. In the CGN there is no distinction between the PCs in Home 1 and the Web Server when it comes to source IP address. They both share the same source IP address 10.2.2.1 on the CPE.

The solution for this is to reserve a port block (or blocks) during the static port forward creation for the specific subscriber.

Figure 59: Dynamic port block starvation in LSN



7.3.1.2.2 Dynamic port block reservation

To prevent starvation of dynamic port blocks for the subscribers that use port forwards, a dynamic port block (or blocks) is reserved during the lifetime of the port forward. Those reserved dynamic port blocks are associated with the same subscriber that created the port forward. However, a log would not be generated until the dynamic port block is actually used and mapping within that block are created.

At the time of the port forward creation, the dynamic port block is reserved in the following fashion:

- If the dynamic port block for the subscriber does not exist, then a dynamic port block for the subscriber is reserved. No log for the reserved dynamic port block is generated until the dynamic port block starts being used (mapping created because of the traffic flow).
- If the corresponding dynamic port block already exists, it is reserved even after the last mapping within the last port block had expired.

The reserved dynamic port block (even without any mapping) continues to be associated with the subscriber as long as the port forward for the subscriber is present. The log (syslog or RADIUS) is generated only when there is not active mapping within the dynamic port block and all port forwards for the subscriber are deleted.

Additional considerations with dynamic port block reservation:

- The port block reservation should be triggered only by the first port forward for the subscriber. The subsequent port forwards do not trigger additional dynamic port block reservation.
- Only a single dynamic port block for the subscriber is reserved (that is, no multiple port-block reservations for the subscriber are possible).
- This feature is enabled with the **configure service vprn nat outside pool port-forwarding-dyn-block-reservation** and the **configure router nat outside pool port-forwarding-dyn-block-reservation** commands. This command can be enabled only if the maximum number of configured port blocks per outside IP is greater or equal then the maximum configured number of subscribers per outside IP address. This guarantees that all subscribers (up to the maximum number per outside IP address) configured with port forwards can reserve a dynamic port block.
- If the port-reservation is enabled while the outside pool is operational and subscribers traffic is already present, the following two cases must be considered:
 - The configured number of subscribers per outside IP is less or equal than the configured number of port blocks per outside IP address (this is permitted) but all dynamic port blocks per outside IP address are occupied at the moment when port reservation is enabled. This leaves existing subscribers with port forwards that do not have any dynamic port blocks allocated (orphaned subscribers), unable to reserve dynamic port blocks. In this case the orphaned subscribers must wait until dynamic port blocks allocated to the subscribers without port forwards are freed.
 - The configured number of subscribers per outside IP is greater than the configured number of port blocks per outside IP address. In addition, all dynamic port blocks per outside IP address are allocated. Before the port reservation is even enabled, the subscriber-limit per outside IP address must be lowered (by configuration) so that it is equal or less than the configured number of port blocks per outside IP address. This action causes random deletion of subscribers that do not have any port forwards. Such subscribers are deleted until the number of subscriber falls below the newly configured subscriber limit. Subscribers with static port forwards are not deleted, regardless of the configured subscriber-limit number. When the number of subscribers is within the newly configured subscriber-limit, the port-reservation can take place under the condition that the dynamic port blocks are available. If specific subscribers with port forwards have more than one dynamic port block allocated, the orphaned subscribers must wait for those additional dynamic port blocks to expire and consequently be released.

- This feature is supported on the following applications: CGN, DS-Lite and NAT64.

7.3.2 Pools with flexible port allocations

Pools with flexible port allocations are specialized pools that allow subscribers to configure per-port allocations instead of per port-block allocations for specific use cases. Logging of port allocations and deallocations within such pools is not facilitated. These pools are compatible with ESA-VM and vISA (VSR) and are not applicable to ISA2. These pools cater to users that have a dedicated private pool. Pools with flexible port allocations ensure that the external IP addresses of the pools are associated with a single user entity or a tenant even before the pool provisioning phase, which eliminates the need for logging.

When pools with flexible port allocations are configured, static port forwards are interspersed with dynamically allocated ports. These pool can only be linked via a NAT policy to a source prefix or to a static port forward. This allows NAT processing of traffic solely originating from the configured source prefix or address. Neither of these two mechanisms inherently steers traffic to the ESA-VM/vISA modules. Therefore, traffic for NAT processing is directed to pools with flexible port allocations based on either destination prefix or filter criteria. In this case, explicitly configured NAT policy is not allowed within a destination prefix or a filter . After the traffic is in the ESA-VM/vISA, the pool is selected based on the source prefix or an existing flow created under the static port forward, which does not have a NAT policy explicitly configured.

Execute the following command to configure pools with flexible port allocations.

```
[ex:/configure router "Base" nat outside]
A:admin@node-2# info
  pool "demo" {
    type large-scale
    nat-group 1
    applications {
      flexible-port-allocation true
    }
  }
}
```

The traffic is steered to ESA-VM/vISA based on a destination prefix or a filter criteria.

- **Destination prefix**

Execute the following command to configure a destination prefix .

```
[ex:/configure service vprn nat inside]
A:admin@node-2# info
  large-scale {
    traffic-identification {
      source-prefix-only true
    }
  }
  nat44 {
    destination-prefix 0.0.0.0/0

    source-prefix 10.10.10.0/24 nat-policy "nat-pol-1"
    source-prefix 10.10.11.0/24 nat-policy "nat-pol-1"
    source-prefix 10.10.12.0/25 nat-policy "nat-pol-2"
    source-prefix 10.10.12.128/25 nat-policy "nat-pol-3"
  }
}
```

- **Filter criteria**

– Filter definition

Execute the following command to configure a filter definition .

```
[ex:/configure filter ip-filter "demo-nat" entry 10]
A:admin@node-2# info
  match {
    dst-ip {
      address 0.0.0.0/0
    }
    dst-port {
      eq 2000
    }
  }
  action {
    nat {
    }
  }
}
```

– Applying filter to the ingress interface

Execute the following command to apply a filter to the ingress interface:

```
[ex:/configure service vprn "nat"]
A:admin@node-2# info
  admin-state enable
  service-id 20
  customer "1"
  interface "access" {
    admin-state enable
    ipv4 {
      primary {
        address 172.16.102.1
        prefix-length 24
      }
    }
    sap lag-4:2 {
      ingress {
        filter {
          ip "demo-nat"
        }
      }
    }
  }
  nat {
    inside {
      large-scale {
        traffic-identification {
          source-prefix-only
        }
        nat44 {
          destination-prefix 0.0.0.0/0

          source-prefix 10.10.10.0/24 nat-policy "nat-pol-1"
          source-prefix 10.10.11.0/24 nat-policy "nat-pol-1"
          source-prefix 10.10.12.0/25 nat-policy "nat-pol-2"
          source-prefix 10.10.12.128/25 nat-policy "nat-pol-3"
        }
      }
    }
  }
```



Note: A default NAT policy is not allowed for pools with flexible port allocations. Execute the following command to configure a default NAT policy.

```
[ex:/configure service vprn nat inside]
A:admin@node-2# info
    large-scale {
        nat-policy "cgn44-demo-default"
    }
```

7.3.2.1 Free port limit

The free port limit feature allows the user to configure a limit on free ports per protocol for each external IP address. This avoids rapid port depletion for new subscribers in paired pooling mode, or unnecessary toggling between external IP addresses for existing subscribers in arbitrary pooling mode.

Such port limits ensure that only IP addresses with sufficient free ports, in accordance with the configured limit, are considered for selection and added to the eligible IP address list.



Note: Free port limit does not interfere with port allocation from an outside IP address for subscribers that are already assigned with the IP address. These subscribers can continue to use ports until the exhaustion of the ports on an outside IP address. After all of the ports on an outside IP address are used up, the system maps subscribers, new or those in arbitrary pooling mode, to a new IP address if it has a port count above the configured limit.

Execute the following command to configure free port limits in pools with flexible port allocations.

```
[ex:/configure router "Base" nat outside pool "demo" large-scale]
A:admin@node-2# info
    flexible-port-allocation {
        free-port-limit {
            tcp 1000
            udp 1000
            other 1000
        }
    }
```

7.3.2.2 Restrictions

The following functionalities are not supported for pools with flexible port allocations and is blocked in the CLI:

- referencing this pool in the destination prefix or filter
- destination-based NAT (dNAT)
- PCP
- deterministic NAT
- L2-Aware NAT
- 1:1 NAT
- maximum number of subscribers per-IP address

- no reservation of ports based on QoS settings (port priorities)
- Stateful Inter-Chassis NAT Redundancy (SICR)
- WLAN-GW or L2 aware firewall-specific functionality (dormant pool or V6 translations)
- scaling profile 1 and scaling profile 3

7.3.3 Association between NAT subscribers and IP addresses in a NAT pool

A NAT subscriber can allocate ports on a single outside IP address or multiple IP addresses in a NAT pool. Nokia recommends that NAT subscribers allocate ports from a single outside IP address. If this IP address runs out of ports, the NAT subscriber runs out of ports. In other words, there is no attempt for a new port to be allocated from a different outside IP address. This method of address allocation to a NAT subscriber is referred to as Paired Address Pooling and is the default behavior in SR OS.

The alternative method of port allocation involves port exhaustion on the originally allocated IP address. An attempt is made to allocate ports from another IP addresses that has free ports available. This results in a NAT subscriber be associated with multiple outside IP addresses. This method is referred to as Arbitrary Address Pooling and can be optionally enabled in SR OS. See RFC 7857, Section 4 for more information.

Arbitrary address pooling may offer more efficient allocations of port-blocks across outside IP address in a NAT pool, but it may negatively affect some applications. For example, an application may require two channels for communication, a control channel and a data channel, each on a different source port from the client perspective on the inside of the NAT. The communication channel may be established on the outside address IP1 and outside port X. If port X is the last free port on the IP1, the SR OS attempts to allocate the next port Y for the data channel from a different outside address, IP2. If the application is robust enough to accept communication from the same client on two different IP addresses, there are no issues. However, some applications may not support this scenario and the communication fails.

Arbitrary address pooling implies the following:

- The subscriber limit per outside IP address loses its meaning because the subscriber can now be associated with multiple IP addresses. Hence, the following command cannot be set.

```
configure router nat outside pool subscriber-limit
```

For more information about the **subscriber-limit** command in paired address pooling, see [Managing port block space](#).

- The number of outside IP addresses in a pool must be at least double the number of ESA-VM in a NAT-group hosting the subscriber. Each subscriber is hashed to a single ESA-VM, ISA2, or vISA; therefore at least two outside IP addresses must be available per ESA-VM, ISA2, or vISA for the subscriber to use more than one outside IP address.
- The number of port blocks configured in a NAT policy using the following command is the aggregate limit that a NAT subscriber can be allocated across multiple outside IP addresses.

```
configure service nat nat-policy block-limit
```

- Reserving a port block by SPF configuration (when an SPF is configured before any port blocks are allocated to the subscriber) is not supported. In other words, the following commands are not supported:

MD-CLI

```
configure router nat outside pool port-forwarding dynamic-block-reservation
```

classic CLI

```
configure router nat outside pool port-forwarding-dyn-block-reservation
```

- Arbitrary address pooling is not supported in L2-aware NAT.

Use the following command to show NAT LSN information for the subscriber.

```
show service nat lsn-subscribers subscriber 276824064
```

The asterisk (*) next to the IP address field in the output indicates that additional outside IP addresses are associated with this subscriber in this pool.

Output example

```
=====
NAT LSN subscribers
=====
Subscriber          : [LSN-Host@192.168.1.1]
NAT policy          : nat-policy-lsn-deterministic
Subscriber ID       : 276824064
-----
Type                : classic-lsn-sub
Inside router       : "Base"
Inside IP address prefix : 192.168.1.1/32
ISA NAT group       : 1
ISA NAT group member : 1
Outside router      : 4
Outside IP address   : 192.0.0.1*
```

Use the detailed version of the command to see additional outside IP addresses and port blocks.

7.3.4 Timeouts

Creating a NAT mapping is only one half of the problem – removing a NAT mapping at the appropriate time maximizes the shared port resource. Having ports mapped when an application is no longer active reduces solution scale and may impact the customer experience should they exhaust their port range block. The NAT application provides timeout configuration for TCP, UDP and ICMP.

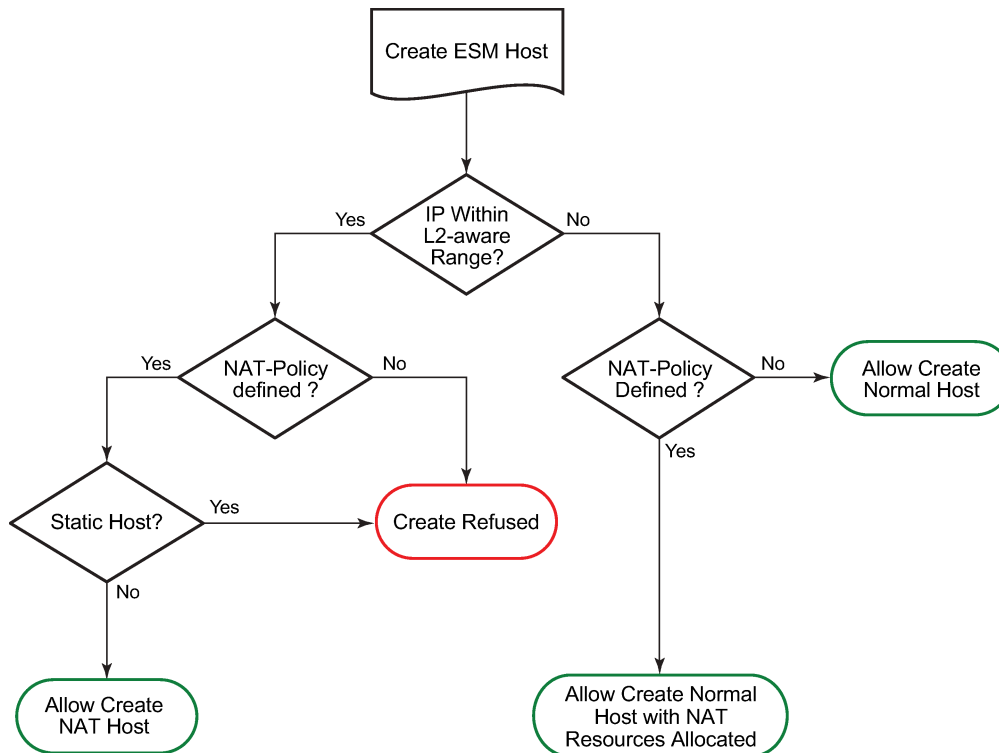
TCP state is tracked for all TCP connections, supporting both three-way handshake and simultaneous TCP SYN connections. Separate and configurable timeouts exist for TCP SYN, TCP transition (between SYN and Open), established and time-wait state. Time-wait assassination is supported and enabled by default to quickly remove TCP mappings in the TIME WAIT state.

UDP does not have the concept of connection state and is subject to a simple inactivity timer. Company-sponsored research into applications and NAT behavior suggested some applications, like the Bittorrent Distributed Hash Protocol (DHT) can make a large number of outbound UDP connections that are unsuccessful. Instead of waiting the default five (5) minutes to time these out, the 7750 SR NAT application supports an udp-initial timeout which defaults to 15 seconds. When the first outbound UDP packet is sent, the 15 second time starts – it is only after subsequent packets (inbound or outbound) that the default UDP timer becomes active, greatly reducing the number of UDP mappings.

7.4 L2-Aware NAT

Figure 60: L2-Aware tree shows the L2-Aware tree.

Figure 60: L2-Aware tree



OSSG711

NAT is supported on DHCP, PPPoE and L2TP. Static and ARP hosts are not supported.

In an effort to address issues of conflicting address space raised in *draft-shirasaki-nat444-isp-shared-addr-02*, an enhancement to Large Scale NAT was co-developed to give every broadband subscriber their own NAT mapping table, yet still share a common outside pool of IPs.

Layer-2 Aware (or subscriber aware) NAT is combined with Enhanced Subscriber Management on the 7750 SR BNG to overcome the issues of colliding address space between home networks and the inside routed network between the customer and Large Scale NAT.

Layer-2 Aware NAT allows every broadband subscriber to be allocated the exact same IPv4 address on their residential gateway WAN link and then proceeds to translate this into a public IP through the NAT application. In doing so, L2-Aware NAT avoids the issues of colliding address space raised in *draft-shirasaki* without any change to the customer gateway or CPE.

Layer-2-Aware NAT is supported on any of the ESM access technologies, including PPPoE, IPoE (DHCP) and L2TP LNS. For IPoE both n:1 (VLAN per service) and 1:1 (VLAN per subscriber) models are supported. A subscriber device operating with L2-Aware NAT needs no modification or enhancement, existing address mechanisms (DHCP or PPP/IPCP) are identical to a public IP service, the 7750 SR BNG simply translates all IPv4 traffic into a pool of IPv4 addresses, allowing many L2-Aware NAT subscribers to share the same IPv4 address.

More information about L2-Aware NAT can be found in *draft-miles-behave-l2nat-00*.

7.4.1 Port block extensions

Similarly to LSN, an L2-Aware NAT subscriber is assigned a single outside IP address per NAT pool, with one or more port blocks tied to the IP address. The outside IP address is shared by multiple subscribers, each with its own unique set of port blocks.

To ensure that a predetermined number of subscribers receive NAT service, an outside IP address and at least one port block on that IP address must be guaranteed. For this reason, the port blocks space in a pool is divided into two partitions:

- **port block space reserved for new L2-Aware NAT subscribers**

Each new subscriber is guaranteed to receive at least one port block, referred to as the initial port block.

- **port block space reserved for the extended port-blocks of existing NAT subscribers**

This port partition can be used by subscribers who exhaust their ports in the initial port block and need additional ports. Pending on the availability and configuration, they are assigned additional port blocks.

Without this type of port space partitioning, the outside IP addresses and the NAT pool may become overtaken by users with heavier port consumption. This denies access to NAT services to a majority of users with lower port consumption.

This division of port space is controlled by limiting the number of subscribers per an outside IP address and configuring the size of the initial port block.

The following shows configuration information relevant to port-block allocation in L2-Aware NAT:

- initial port block size for new subscribers

MD-CLI

```
configure service vprn <service-name> nat outside pool <name>
  port-reservation {
    ports <number>
  }
```

Classic CLI

- The pool name must be **type l2-aware**.
- **port-reservation blocks num** can be set only if **port-block-extension** is not enabled.

```
configure service vprn <id> nat outside pool <name>
  port-reservation blocks <num>
  port-reservation ports <num>
```

- the extended port block size for existing subscribers and the maximum number of subscribers per outside IP address. The size of the initial port blocks and extended port block may differ.

MD-CLI

```
configure service vprn <service-name> nat outside pool <name>
  l2-aware {
    port-block-extension {
      ports <number>
      subscriber-limit <number>
    }
  }
```

Classic CLI

The pool name must be **type I2-aware**.

```
configure service vprn <id> nat outside pool <name>
  port-block-extensions ports <num> subscriber-limit <num>
```

- upper boundary for static port forwards

MD-CLI

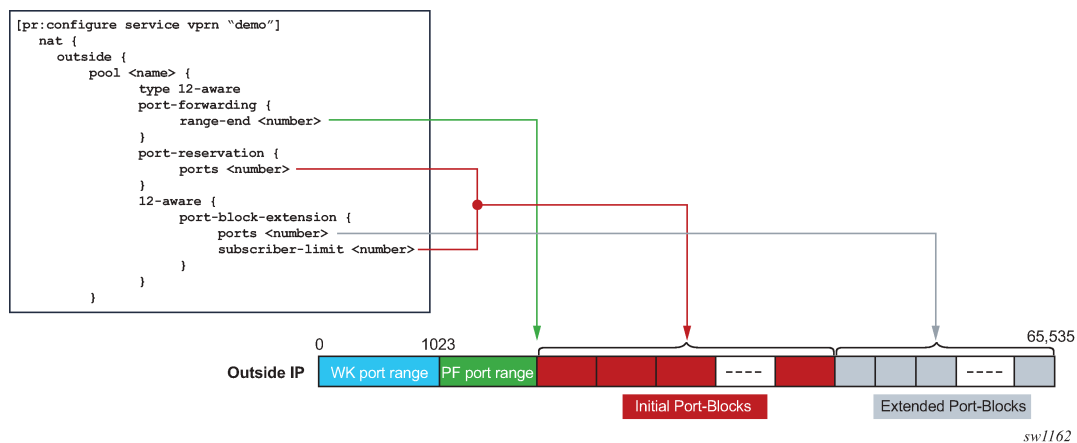
```
[configure service vprn <service-name> nat outside pool <name>]
  port-forwarding {
    range-end <number>
  }
```

Classic CLI

```
configure service vprn <id> nat outside pool <name>
  port-forwarding-range <range-end>
```

Figure 61: Port space partitioning for an outside IP address shows the effects of the commands.

Figure 61: Port space partitioning for an outside IP address



The maximum number of port blocks that can be allocated per subscriber is controlled by the following configuration in the NAT policy.

MD-CLI

```
[configure service nat nat-policy <name>]
  block-limit <number>
```

Classic CLI

```
configure service nat nat-policy <name>
  block-limit <number>
```

7.4.1.1 Managing port block space

Both port partitions, initial and extended, are served on a first-come-first-served basis. The initial port partition guarantees at least one port block for each of the preconfigured number of subscribers per outside IP address (**configure router nat outside pool port-block-extension ports num-ports subscriber-limit number** command). If there are more subscribers in the network than the preconfigured number of NAT subscribers, then this space becomes oversubscribed.

Extended port partitioning, however, does not guarantee that each of the existing NAT subscriber receives additional port blocks. Each subscriber can allocate additional free port blocks only if they are available, up to the maximum combined limit (initial and extended) set in the NAT policy (**configure service nat nat-policy block-limit** command).

For optimized NAT pool management and correct capacity planning, understanding the following parameters in the operator's network is essential:

- IP address compression ratio (the number of subscribers who share one outside IP address)
- subscriber oversubscription ratio (the number of NAT subscribers who are active simultaneously)
- statistical port usage for subscribers (the percentage of subscribers who are heavy, medium, and light port users)
- port block sizes

Based on the previous parameters, an average port block per subscriber can be determined and the following parameters in NAT can be set:

- the subscriber limit per outside IP address configured in the NAT pool
- the size of the initial and extended port blocks configured in the NAT pool
- the maximum number of port blocks per subscriber configured in NAT policy
- the outside IP address range configured in the NAT pool

The following are reasonable guidelines with an example that can serve as an initial configuration for operators who are unsure of their current traffic patterns in terms of port usage for their subscribers.

- An operator has 10,000 subscribers that require NAT, but only 8,000 of them are active simultaneously. This means that the operator can allow oversubscription of outside (NAT) IP address.
- Average port usage:
 - 60% of the subscribers are light port users with less than 1000 ports.
 - 30% of the subscribers are medium port users with less than 2000 ports.
 - 10% of the subscribers are heavy port users with less than 4000 ports.

These assumptions lead to the following calculations:

- $8,000 \text{ active subscribers} \times (0.6 \times 1000 + 0.3 \times 2,000 + 0.1 \times 4,000) = 12,800,000 \text{ total ports.}$
- Consider that one outside IP address can accommodate ~50,000 (64K ports less the static port forwards and well known ports). This yields 256 outside IP addresses (/24) in a pool therefore, $12,800,000 / 50,000 = 256.$
- The compression ratio is 8,000 divided by 256 equals ~32 (32 subscribers share one outside IP address), therefore the subscriber limit equals 32.

Based on this calculation, a reasonable size for the initial port block size is 1000 ports and the extended port block size is 335 ports.

- The maximum number of port blocks per subscriber is set to 10 to accommodate heavy users with 4,000 ports ($1000 + 9 \times 335 = 4015$)

Setting the subscriber limit in a pool to 32, the initial and extended port block sizes to 1000 and 335 respectively, the maximum number of port blocks per subscriber to 10, and configuring a /24 address range in a pool would produce the needed results. This assumes that the subscribers are properly load-balanced over ISAs or ESAs. The following is an example configuration.

```
[configure service vprn "demo vprn" nat outside pool "demo pool"]
  port-reservation {
    ports 1000
  }
  l2-aware {
    port-block-extension {
      ports 335
      subscriber-limit 32
    }
  }
  port-forwarding {
    range-end 15000
  }
}

[configure service nat nat-policy "demo-policy"]
  block-limit 10
```

7.4.2 L2-Aware NAT bypass

L2-Aware NAT bypass refers to the functionality where the entire or partial traffic from a L2-Aware-NAT-enabled ESM1 subscriber circumvents local NAT function. There are three types of bypass supported for L2-Aware NAT in the SR OS:

- full ESM host bypass
- selective ESM host bypass based on an IP filter match
- the entire ESM subscriber bypass because of ISA/ESA failure. This type of bypass is described in [NAT redundancy](#).

7.4.2.1 Full ESM host bypass

In this type of bypass, a subscriber host is implicitly excluded from L2-Aware NAT if its IP address falls outside of the configured subnet in the inside NAT CLI hierarchy under the L2-Aware CLI node.

In the following example, the address under the L2-Aware CLI node (address 10.10.1.254/24) represents the default gateway and a L2-Aware subnet. Hosts with IP addresses within the configured L2-Aware subnet (in this example 10.10.1.0/24) are subjected to L2-Aware NAT (the exception is the default gateway address 10.10.1.254). Hosts outside of this IP range bypass NAT. In this way, a mix of hosts under the same L2-Aware enabled ESM subscriber can coexist, some of which are subject to NAT, and some of which are bypassing NAT.

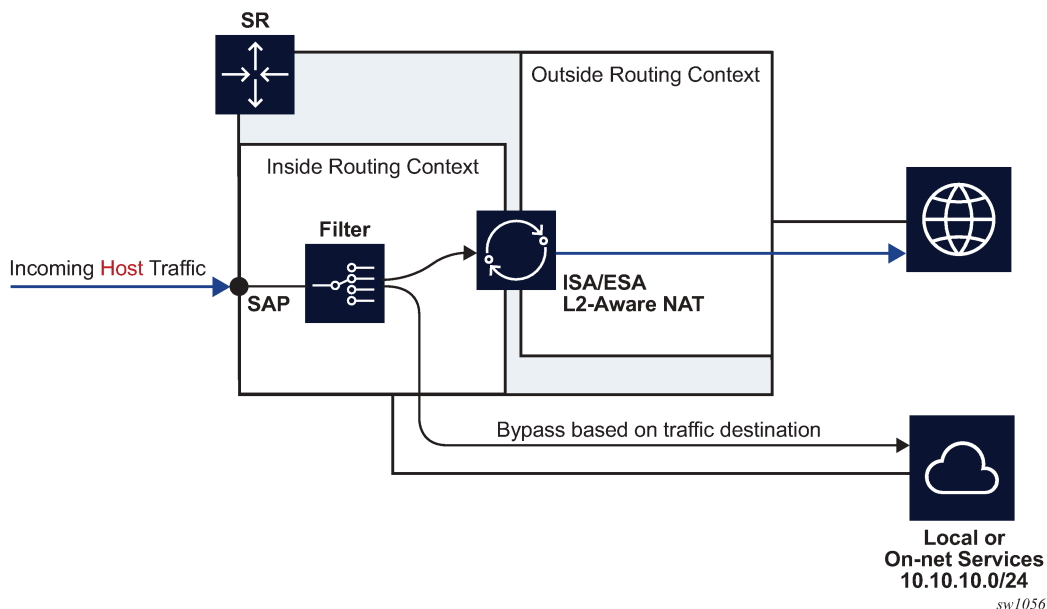
```
configure
  router
    nat
      inside
        l2-aware
          address 10.10.1.254/24
```

7.4.2.2 Selective L2-Aware NAT bypass

In selective L2-Aware NAT bypass, a decision whether to perform NAT is made based on the traffic classifiers (match conditions) defined in an IP filter applied to an ESM host.

A typical use case for selective L2-Aware NAT bypass is based on destinations, where on-net services are needed to be accessed without NAT, while some other off-net destinations, require NAT. Traffic to those on-net services is identified based on the destination IP addresses ([Figure 62: L2-Aware bypass based on traffic destination](#)).

Figure 62: L2-Aware bypass based on traffic destination



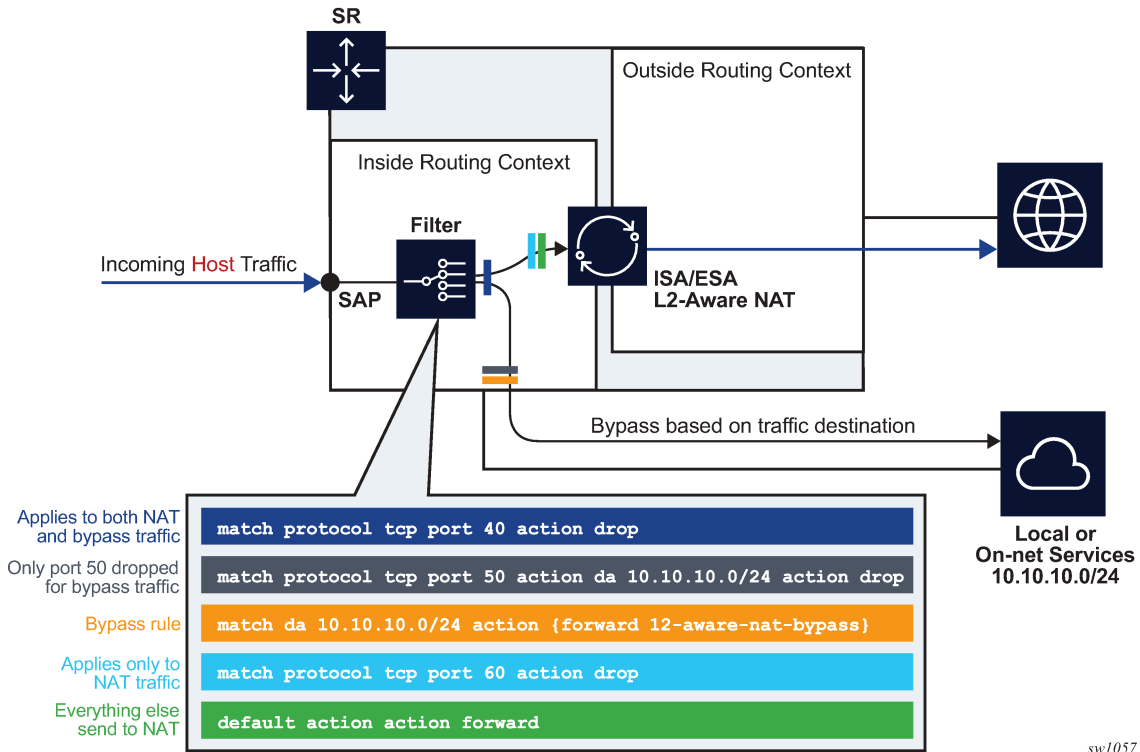
L2-Aware NAT subscribers that are candidates for selective bypass in the SR OS, must be first identified and enabled with the **config>subscr-mgmt>sub-prof>nat-allow-bypass** command:

After the selective L2-Aware NAT bypass is enabled, the determination of whether specific traffic from a host bypasses NAT comes via an IP filter with a newly defined action **l2-aware-nat-bypass**. This new action must be configured in addition to the existing action **accept** (in MD-CLI) or **forward** (in classic CLI). This defined set of actions divert identified traffic away from NAT.

Although most typical use cases require traffic identification based on destination IP addresses, generic match statements in IP filters allow identification of traffic based on any Layer 3 fields.

The filter entries are executed in top-to-bottom order as shown in [Figure 63: Filtering example for L2-Aware NAT bypass](#).

Figure 63: Filtering example for L2-Aware NAT bypass



sw1057

Table 36: Configuration options for selective L2-Aware NAT bypass describes the behavior in relation to the three configuration options that directly influence selective L2-Aware NAT bypass.

Table 36: Configuration options for selective L2-Aware NAT bypass

L2-Aware NAT-enabled host	Selective bypass enabled	IP filter action l2-aware-nat-bypass accept forward	Behavior
Yes	Yes	Yes	Selective bypass is in effect
Yes	Yes	No	The host is enabled for bypass, but without the corresponding IP filter action. Bypass is not in effect and all traffic from the host is NAT'd. After the bypass action is provided via the IP filter, traffic identified in the IP filter is bypassed.
Yes	No	Yes	The host is not enabled for bypass, but the IP filter is configured for bypass. This is an incorrect condition where host traffic is bypassed in the upstream direction but not in the downstream direction. As a result, downstream traffic is dropped.

L2-Aware NAT-enabled host	Selective bypass enabled	IP filter action l2-aware-nat-bypass accept forward	Behavior
Yes	No	No	The host is not enabled for bypass. All host traffic is NAT'd.
No	Yes	Yes	The host is not an L2-Aware NAT host. This is a full bypass case.
No	Yes	No	The host is not an L2-Aware NAT host. This is a full bypass case.
No	No	Yes	The host is not an L2-Aware NAT host. This is full bypass case.
No	No	No	The host is not an L2-Aware NAT host. This is full bypass case.

The following are configuration considerations:

- An ESM-enabled host can be enabled if the following two conditions are met:
 - The subscriber's sub-profile contains the **nat-policy name** command.
 - The host IP address belongs to the subnet configured under the **address** command:

```
configure
  router/service vprn
    nat
      inside
        l2-aware
          address <ip-address/mask>
```

- Selective bypass is enabled if the **configure subscriber-mgmt sub-profile nat-allow-bypass** command is configured under the sub-profile.
- All configuration options are allowed in the CLI and it is up to the operator to consult [Table 36: Configuration options for selective L2-Aware NAT bypass](#) for expected results.

7.4.2.2.1 On-line change of the selective NAT bypass

While traffic is flowing, it is possible to change its path from going through NAT to bypassing NAT. This kind of transition while traffic is flowing is referred to as on-line change.

NAT bypass for L2-Aware NAT subscriber can be influenced with two configuration parameters (**configure subscr-mgmt sub-prof nat-access-mode** and **configure filter ip-filter entry action l2-aware-nat-bypass**). [Table 36: Configuration options for selective L2-Aware NAT bypass](#) lists possible (valid and invalid) combinations of the two.

Enabling and disabling NAT bypass (**config>subscr-mgmt>sub-prof>nat-access-mode**), is supported by changing the subscriber profile for the ESM subscriber with RADIUS/Gx. However, changing the subscriber profile configuration with CLI while the profile is in use is not allowed.

The recommend method to change the IP filter action (**l2-aware-nat-bypass**) is to override the existing IP filter using RADIUS or Gx.

For additional restrictions when using this feature, see Known Limitations section in the Release Notes.

7.4.2.2.2 NAT bypass verification

To verify that NAT bypass is in effect, use the following **show service active-subscribers detail** and **show filter ip** commands.

```
*A:Dut-C# show service active-subscribers detail
=====
Active Subscribers
=====
-----
Subscriber AL_x0ffx6x0x0x2
          (sub_l2-dhcp1)
-----
NAT Policy   : pol-B-1
Outside IP   : 130.0.0.201
Ports        : 1536-1570
NAT Policy   : pol-o-1
Outside IP   : 19.0.0.87 (vprn101)
Ports        : 1024-1055
NAT Policy   : pol-o1-1
Outside IP   : 130.0.0.68 (vprn601)
Ports        : 1152-1349
NAT Policy   : pol-o2-1
Outside IP   : 130.0.0.222 (vprn602)
Ports        : 1152-1349
-----
I. Sched. Policy : N/A
E. Sched. Policy : N/A
E. Agg Rate Limit: Max
E. Min Resv Bw   : 1

I. Policer Ctrl. : N/A
E. Policer Ctrl. : N/A
I. vport-hashing : Disabled
I. sec-sh-hashing: Disabled
Q Frame-Based Ac*: Disabled
Acct. Policy     : N/A
ANCP Pol.        : N/A
Accu-stats-pol  : (Not Specified)
HostTrk Pol.    : N/A
IGMP Policy      : N/A
MLD Policy       : N/A
PIM Policy       : N/A
Sub. MCAC Policy: N/A
NAT Policy       : pol-o-1
Firewall Policy  : N/A
UPnP Policy      : N/A
NAT Prefix List : npl-4
Allow NAT bypass : Yes
Collect Stats    : Disabled
```

This command shows that the ESM subscriber is a L2-Aware NAT subscriber for which bypass is enabled.

The following command provides insight into whether the NAT bypass is in effect:

```
*A:Dut-A>config>filter# show filter ip 10
=====
IP Filter
=====
Filter Id       : 10
Scope           : Template
Type            : Normal
Applied         : Yes
Def. Action     : Drop
```



```

System filter      : Unchained
Radius Ins Pt     : n/a
CrCtl. Ins Pt    : n/a
RadSh. Ins Pt    : n/a
PccRl. Ins Pt    : n/a
Entries          : 1
Description       : (Not Specified)
Filter Name       : 10
-----
Filter Match Criteria : IP
-----
Entry             : 1
Description       : (Not Specified)
Log Id           : n/a
Src. IP          : 0.0.0.0/0
Src. Port        : n/a
Dest. IP         : 0.0.0.0/0
Dest. Port       : n/a
Protocol         : Undefined           Dscp           : Undefined
ICMP Type        : Undefined           ICMP Code      : Undefined
Fragment         : Off                 Src Route Opt  : Off
Sampling         : Off                 Int. Sampling  : On
IP-Option        : 0/0                 Multiple Option: Off
Tcp-flag         : (Not Specified)
Option-pres      : Off
Egress PBR       : Disabled
Primary Action   : Forward
L2 Aware NAT Bypass : Enabled
Ing. Matches     : 0 pkts
Egr. Matches     : 0 pkts
=====

```

7.4.3 L2-Aware NAT destination-based multiple NAT policies

Multiple NAT policies for a L2-Aware subscriber can be selected based on the destination IP address of the packet. This allows the operator to assign different NAT pools and outside routing contexts based on the traffic destinations.

The mapping between the destination IP prefix and the NAT policy is defined in a nat-prefix-list. This nat-prefix-list is applied to the L2-Aware subscriber through a subscriber profile. After the subscriber traffic arrives to the MS-ISA where NAT is performed, an additional lookup based on the destination IP address of the packet is executed to select the specific NAT policy (and consequently the outside NAT pool). Failure to find the specific NAT policy based on the destination IP address lookup results in the selection of the default NAT policy referenced in the subscriber profile.

CLI example:

```

-----
echo "Service Configuration"
#-----
service
  nat
    nat-policy "l2aw nat policy" create
    pool "l2aw-nat-pool" router 1
  exit
  nat-policy "another-l2aw-nat-policy" create
  pool "another-l2aw-nat-pool" router 2
  exit
  nat-policy "default-nat-policy" create

```

```

    pool "default-nat-pool" router Base
    exit

    nat-prefix-list "prefixlist1" application l2-aware-dest-to-policy create
    prefix 192.168.0.0/30 nat-policy "l2aw-nat-pol"
    prefix 192.168.0.64/30 nat-policy "l2aw-nat-pol"
    prefix 192.168.0.128/30 nat-policy "l2aw-nat-pol"
    prefix 192.168.1.0/30 nat-policy "another-l2aw-nat-pol"
    prefix 192.168.1.64/30 nat-policy "another-l2aw-nat-pol"
    prefix 192.168.1.128/30 nat-policy "another-l2aw-nat-pol"
    exit
  exit

#-----
echo "Subscriber-mgmt Configuration"
#-----
  subscriber-mgmt
    sub-profile "sub_profile" create
    nat-policy "def-nat-policy"
    nat-prefix-list "prefixlist1"
  exit

```

As displayed in the example, multiple IP prefixes can be mapped to the same NAT policy.

The NAT prefix list cannot reference the default NAT policy. The default NAT policy is the one that is referenced directly under the subscriber profile.

7.4.3.1 Logging

In L2-Aware NAT with multiple nat-policies, the NAT resources are allocated in each pool associated with the subscriber. This NAT resource allocation is performed at the time when the ESM subscriber is instantiated. Each NAT resource allocation is followed by log generation.

For example, if RADIUS logging is enabled, one Alc-NAT-Port-Range VSA per NAT policy is included in the acct START/STOP message.

[Alc-Nat-Port-Range = "192.168.20.1 1024-1055 router base **nat-pol-1**"

Alc-Nat-Port-Range = "193.168.20.1 1024-1055 router base **nat-pol-2**".

Alc-Nat-Port-Range = "194.168.20.1 1024-1055 router base" **nat-pol-3**.]

7.4.3.1.1 RADIUS logging and NAT-policy change via CoA

Nat-policy change for L2-Aware NAT is supported through a sub-profile change triggered in CoA. However, change of sub-profile alone through CoA does not trigger generation of a new RADIUS accounting message and therefore NAT events related to NAT policy changes are not promptly logged. For this reason, each CoA initiating the sub-profile change in a NAT environment must do one of the following:

- Change the sla-profile.
- Include the Alc-Trigger-Acct-Interim VSA in the CoA messages.

Note that the sla-profile has to be changed and not just refreshed. In other words, replacing the existing sla-profile with the same one does not trigger a new accounting message.

Both of these events trigger an accounting update at the time CoA is processed. This keeps NAT logging current. The information about NAT resources for logging purposes is conveyed in the following RADIUS attributes:

- Alc-Nat-Port-Range-Freed VSA → NAT resources released because of CoA.
- Alc-Nat-Port-Range VSA → NAT resources in use. These can be the existing NAT resources which were not affected by CoA or they can be new NAT resource allocated because of CoA.

NAT logging behavior because of CoA depends on the deployed accounting mode of operation. This is described in [Table 37: NAT-policy change and CoA in L2Aware NAT](#). The **interim-update** keyword must be configured for host/session accounting for Interim-Update messages to be triggered:

```
configure
 subscriber-mgmt
   radius-accounting-policy <name>
     session-accounting interim-update
configure
 subscriber-mgmt
   radius-accounting-policy <name>
     host-accounting interim-update
```

Table Legend:

AATR (Alc-Acct-Triggered-Reason) VSA — This VSA is optionally carried in Interim-Update messages that are triggered by CoA.

ATAI (Alc-Trigger-Acct-Interim) VSA — this VSA can be carried in CoA to trigger Interim-Update message. The string carried in this VSA is reflected in the triggered Interim-Update message.

I-U (Interim-Update Message)

Table 37: NAT-policy change and CoA in L2Aware NAT

	Host or session accounting	Queue-instance accounting	Comments
CoA Sub-prof change + ATAI VSA	Single I-U with: <ul style="list-style-type: none"> • released NAT info • unchanged NAT info • new NAT info • AATR • ATAI 	Single I-U with: <ul style="list-style-type: none"> • released NAT info • unchanged NAT info • new NAT info • AATR • ATAI 	Single I-U message is triggered by CoA.
CoA Sub-profile change + Sla-profile change	First I-U: <ul style="list-style-type: none"> • released NAT info • unchanged NAT info • new NAT info Second I-U: <ul style="list-style-type: none"> • unchanged NAT info • new NAT info 	Acct Stop: <ul style="list-style-type: none"> • released NAT info • unchanged NAT info • new NAT info Acct Start: <ul style="list-style-type: none"> • unchanged NAT info • new NAT info 	Two accounting messages are triggered in succession.
CoA Sub-profile change	—	—	No accounting messages are triggered by CoA. The next regular I-U messages contain:

	Host or session accounting	Queue-instance accounting	Comments
			<ul style="list-style-type: none"> old (released) NAT info unchanged NAT info new NAT info
CoA Sub-profile change+ Sla-profile-change + ATAI VSA	First I-U: <ul style="list-style-type: none"> released NAT info unchanged NAT info new NAT info Second I-U: <ul style="list-style-type: none"> unchanged NAT info new NAT info AATR ATAI 	Acct Stop: <ul style="list-style-type: none"> re-released NAT info unchanged NAT info new NAT info Acct Start: <ul style="list-style-type: none"> unchanged NAT info new NAT info 	Two accounting messages are triggered in succession.

For example, the second CoA row describes the outcome triggered by CoA carrying new sub and sla profiles. In host/session accounting mode this creates two Interim-Update messages. The first Interim-Messages carries information about:

- the released NAT resources at the time when CoA is activated
- existing NAT resources that are not affected by CoA
- new NAT resources allocated at the time when CoA is activated

The second Interim-Update message carries information about the NAT resources that are in use (existing and new) when CoA is activated.

From this, the operator can infer which NAT resources are released by CoA and which NAT resources continue to be in use when CoA is activated.

7.4.3.1.2 Delay between the NAT resource allocation and logging during CoA

Nat-policy change induced by CoA triggers immediate log generation (for example acct STOP or INTERIM-UPDATE) indicating that the nat resources have been released. However, the NAT resources (outside IP addresses and port-blocks) in SR OS are not released for another five seconds. This delay is needed to facilitate proper termination of traffic flow between the NAT user and the outside server during the NAT policy transition. A typical example of this scenario is the following:

- HTTP traffic is redirected to a WEB portal for authentication. Only when the user is authenticated, access to the Internet is granted along with a new NAT policy that provides more NAT resources (larger port-ranges, and so on).
- After the user is authenticated, CoA is used to change the user forwarding properties (HTTP-redirect is removed and the NAT policy is changed). However, CoA must be sent before the authentication acknowledgment (ACK) messages is sent, otherwise the next new HTTP request would be redirected again.

3. Authentication acknowledgment is sent to the NAT user following the CoA which removed the HTTP redirect and instantiated a new NAT policy. Because the original communication between the WEB portal and the NAT user was relying on the original NAT policy, the NAT resources associated with the original NAT policy must be preserved to terminate this communication gracefully. Therefore, the delay of five seconds before the NAT resources are freed.

Similar to other stale dynamic mappings, stale port forwards are released after five seconds. Note that static port forwards are kept on the CPM. New CoAs related to NAT are rejected (NAK'd) in case that the previous change is in progress (during the 5seconds interval until the stale mappings are purged).

7.4.3.2 Static port forwards

Unless the specific NAT policy is provided during Static Port Forward (SPF) creation (SPF creation command), the port forward is created in the pool referenced in the default NAT policy. Nat-policy can be part of the command used to modify or delete SPF. If the NAT policy is not provided, then the behavior is:

- If there is only one match, the port forward is modified or deleted.
- If there is more than one match, modify or delete port forward must specify a NAT policy. Otherwise, the modify or delete action fails.

A match is considered when at least these parameters from the modify or delete command are matched (mandatory parameters in the **spf** command):

- subscriber identification string
- inside IP address
- inside port
- protocol

For a Layer 2-Aware NAT, an alternative AAA interface can be used to specify SPF. An alternative AAA interface and CLI-based port forwards are mutually exclusive. See the 7450 ESS, 7750 SR, and VSR RADIUS Attributes Reference Guide for more details.

7.4.3.3 L2-Aware ping

Similar to the non-L2-Aware **ping** command, understanding how the ICMP Echo Request packets are sourced in L2-Aware ping is crucial for the correct execution of this command and the interpretation of its results. The ICMP Echo Reply packets must be able to reach the source IP address that was used in ICMP Echo Request packets on the SR OS node on which the L2-Aware **ping** command was executed. See [Figure 64: L2-Aware ping](#).

The return packet (the ICMP Echo reply sent by the targeted host) is subject to L2-Aware NAT routing executed in the MS-ISA. The L2-Aware NAT routing process looks at the destination IP address of the upstream packet and then directs the packet to the correct outside routing context. The result of this lookup is a NAT policy that references the NAT pool in an outside routing context. This outside routing context must be the same as the one from which the L2-Aware **ping** command was sourced. Otherwise, the L2-Aware **ping** command fails.

The L2-Aware **ping** command can be run in two modes:

- basic mode (**ping ip-address subscriber subscriber-id**) in which the *subscriber-id* is a required field to differentiate subscriber hosts that assigned the same IP address (although each host has its own instantiation of this IP address)

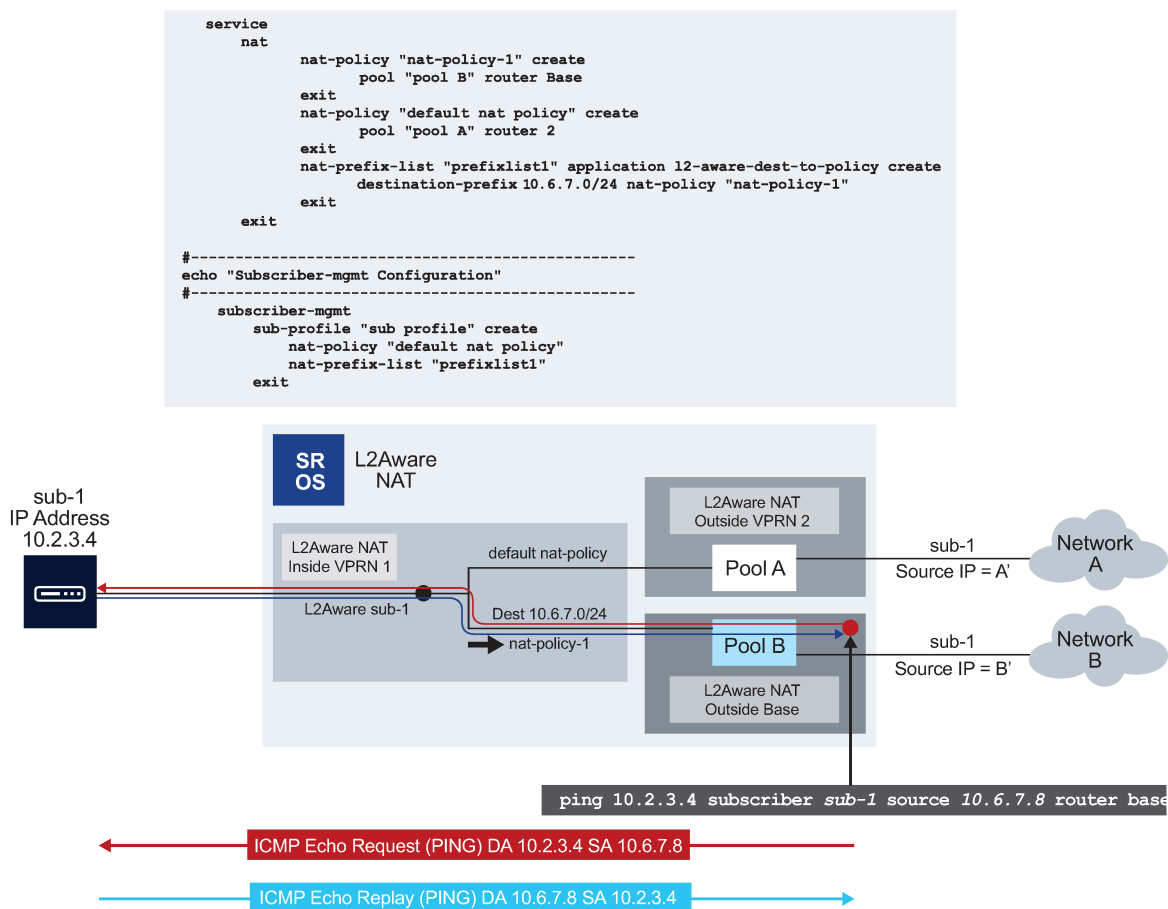
- extended mode where additional parameters can be selected. The two most important being the source IP address (source) and the routing context (router):

ping ip-address subscriber subscriber-id source ip-address router router-id

Figure 64: L2-Aware ping shows the traffic flow for an L2-Aware ping command targeting the subscriber's IP address 10.2.3.4, sourced from the Base routing context using an arbitrary source IP address of 10.6.7.8 (it is not required that this IP address belong to the L2-Aware ping originating node).

When the host 10.2.3.4 replies, the incoming packets with the destination IP address of 10.6.7.8 are matched against the destination-prefix 10.6.7.0/24 referencing the nat-policy-1. nat-policy-1 contains the Pool B which resides in the Base routing context. Hence, the loop is closed and the execution of the L2-Aware ping command is successful.

Figure 64: L2-Aware ping



SIW0097

L2-Aware ping is always sourced from the outside routing context, never from the inside routing context. If the router is not specifically configured as an option in the L2-Aware ping command, the Base routing context is selected by default. If that the Base routing context is not one of the outside routing contexts for the subscriber, the L2-Aware ping command execution fails with the following error message:

"MINOR: OAM #2160 router ID is not an outside router for this subscriber."

7.4.3.4 UPnP

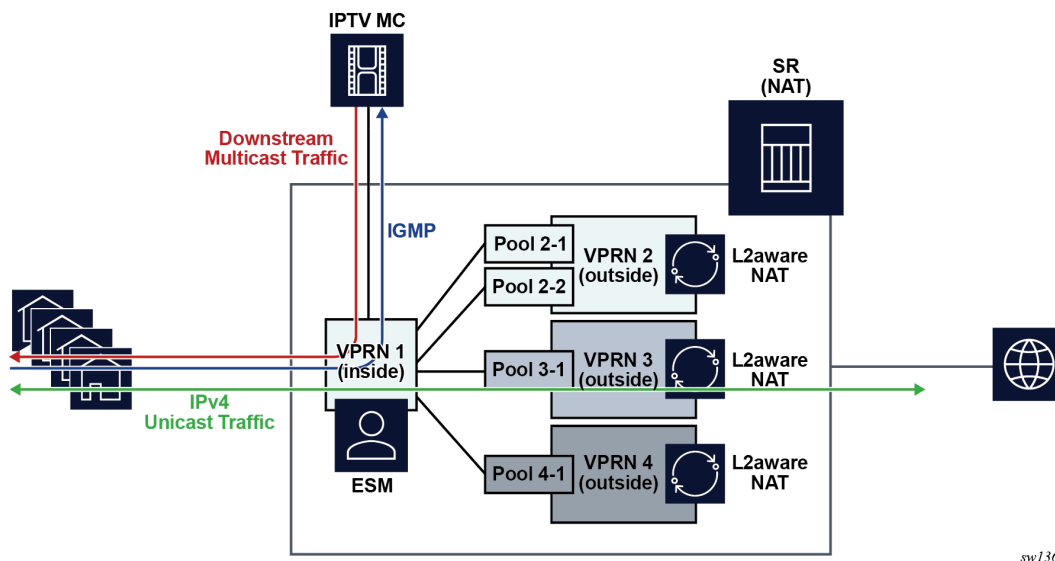
UPnP uses the default NAT policy.

7.4.3.5 L2-Aware NAT and multicast

Multicast traffic through NAT is not supported. However, if the downstream multicast traffic is received in the inside routing context, without going through NAT, the traffic can be forwarded to a L2-Aware host.

The following figure shows an example of downstream multicast traffic on the inside.

Figure 65: Downstream multicast traffic originated on the inside



To enable this type of traffic, the following entities must be configured:

- A NAT policy in the **sub-profile** context. This configuration enables L2-Aware subscribers.
- An IGMP policy in the **sub-profile** context. This enables multicast traffic that is originated on the inside.
- The **force-unique-ip-addresses** command in the **inside l2-aware** context must be enabled. This configuration enforces the uniqueness of IPv4 addresses of L2-Aware subscribers. L2-Aware NAT supports overlapping subscriber IPv4 addresses in the inside routing context.

The following examples show the configuration of L2-Aware NAT and multicast.

Example: MD-CLI

```
A:admin@node-2# configure subscriber-mgmt
  sub-profile "demo-profile" {
    igmp-policy "demo-mcast"
    nat {
      policy "demo-nat-pol"
    }
  }

A:admin@node-2# configure router "Base" nat inside
  l2-aware {
    force-unique-ip-addresses true
```

```

}

/configure service vprn "demo-vprn" nat inside
  l2-aware {
    force-unique-ip-addresses true
  }

```

Example: classic CLI

```

A:node-2>config>subscr-mgmt# info
-----
  sub-profile "demo-profile" create
    nat-policy "demo-nat-pol"
    igmp-policy "demo-mcast"
  exit
  exit

A:node-2>config>service#
  vprn 105 name "demo-vprn" customer 1 create
    nat
      inside
        l2-aware
          address 192.168.100.200/32
        exit
      exit
    exit
  no shutdown
  exit

A:node-2>config>router>nat>inside>l2-aware# info
-----
  force-unique-ip-addresses
-----

*A:node-2>config>service>vprn>nat>inside>l2-aware$ info
-----
  force-unique-ip-addresses
-----

```

7.5 NAT pool addresses and ICMP Echo Request/Reply (ping)

The outside IPv4 addresses in a NAT pool can be configured to answer pings. ICMPv4 Echo Requests are answered with ICMPv4 Echo Replies.

In 1:1 NAT, ICMP Echo Requests are propagated to the host on the inside. The host identified by a NAT binding then answers the ping.

In Network Address Port Translation (NAPT), ICMP Echo Requests are not propagated to the hosts behind the NAT. Instead, the reply is issued by the SR OS from the ESA or ISA.

In NAPT, the behavior is as follows:

- In L2-aware NAT when **port-block-extensions** is disabled, the reply from an outside IP address is generated only when the IP address has at least one host (binding) behind it.

- In L2-aware NAT when **port-block-extensions** is enabled, the reply from an outside IP address is generated regardless if a binding is present.
- In LSN, the reply from an outside IP address is generated regardless if a binding is present.

For security reasons, the ICMP Echo Reply functionality is disabled by default. The following commands enable the behavior:

Classic CLI

```
configure
router
  nat
    outside
      pool <name> nat-group <id> type <large-scale|l2-aware|wlan-gw-anchor>
        [no] icmp-echo-reply
```

MD-CLI

```
configure {
  router "Base" | service vprn <name> {
    nat {
      outside {
        pool <name> {
          icmp-echo-reply <boolean>
        }
      }
    }
  }
}
```

This functionality is on a per pool basis and it can be configured online while the pool is enabled.

7.6 Traffic steering to NAT

Traffic steering to NAT refers to the mechanism by which traffic in the SR line card is redirected to the ISA or VM-ESA for NAT processing. This traffic must be identified first and then redirected to ISA or VM-ESA. The mechanism by which traffic is steered to NAT in an SR node in the upstream direction depends on the NAT type.

For LSN44, the upstream traffic (in the private to public direction) is steered (or redirected) to NAT in an SR node through one of the two mechanisms:

- routing
- filters

Both methods are applied in the inside (private) routing context. Traffic matched through routing or filter criteria is sent to the ISA or VM-ESA for NAT processing and from there to the outside (public) routing context where it exits the node.

In NAT64 and DS-lite, traffic is steered to NAT mainly through routing NAT64 prefix in NAT64 and Address Family Transition Router (AFTR) IPv6 address in DS-lite. However, the routing can be augmented with IPv6 filters to accommodate mapping to multiple NAT pools per subscriber.

In L2-Aware, where NAT is integrated with ESM, traffic is steered to NAT automatically assuming that the subscriber session is associated with NAT during session instantiation phase.

In all NAT types, the downstream traffic arriving in the outside (public) routing context is forwarded to NAT through routing and public pool IPv4 addresses are installed in the routing table with the next hop pointing to the ISA or VM-ESA.

The following sections describe steering logic for LSN44 which is the only NAT type that supports dynamic routing.

7.6.1 Routing approach in LSN44

Routing approach relies on destination IP-based match. A destination IP route leading to NAT can be static (specifically configured) or dynamic (installed through the BGP routing protocol).

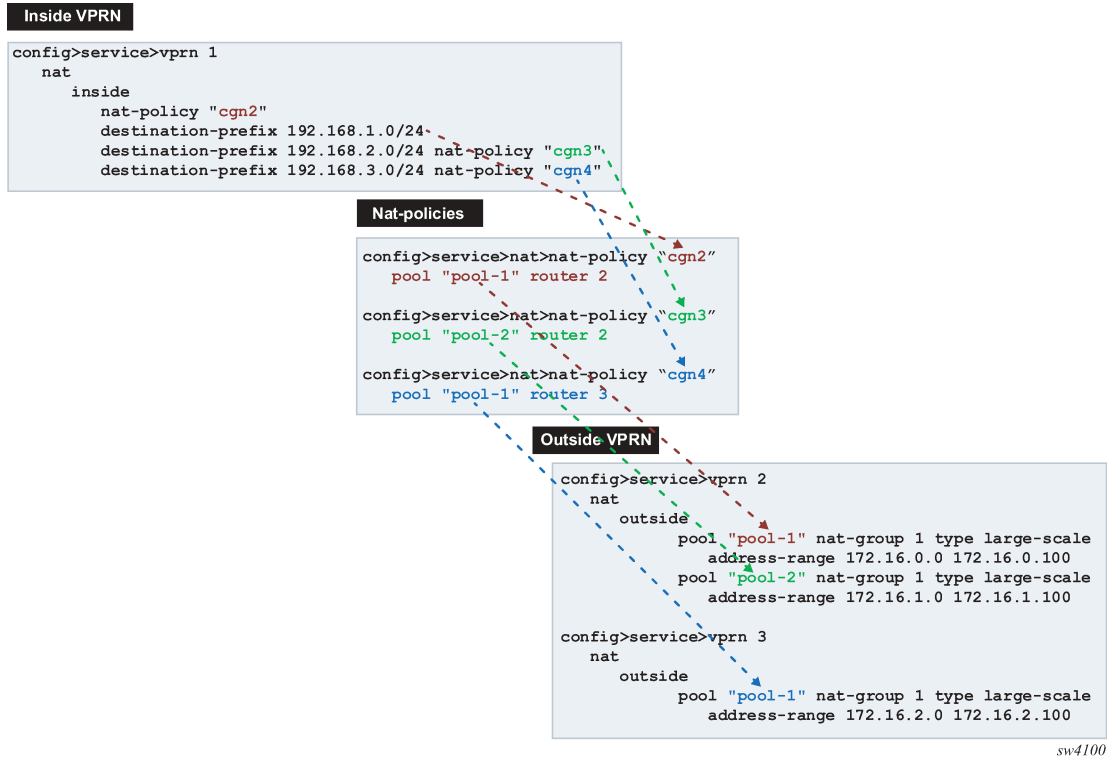
7.6.1.1 Static NAT routes

Static steering to LSN44 is based on the **destination-prefix** command with a statically-configured routing prefix in an inside routing context (VPRN or GRT). This static route points to an ISA or VM-ISA. In transit from the inside routing context to the outside routing context, frames must be redirected through an ISA or VM-ESA where NAT is performed.

If there are multiple ISA or VM-ESAs in a NAT group, an internal LAG per NAT group is used with member ports connected to each ISA or VM-ESA. Upstream traffic is load-balanced between the ISAs or VM-ESAs based on the source IP addresses or prefixes.

The CLI configuration used for static steering to LSN44 is shown in [Figure 66: Basic CLI for NAT](#) where a NAT policy containing a pool name and the outside routing context acts as a bond between the inside routing context and the outside routing context. A destination-prefix without an explicitly configured NAT policy uses a default NAT policy. As shown in [Figure 66: Basic CLI for NAT](#), the prefix 192.168.1.0/24 is mapped to **nat-policy cgn2**.

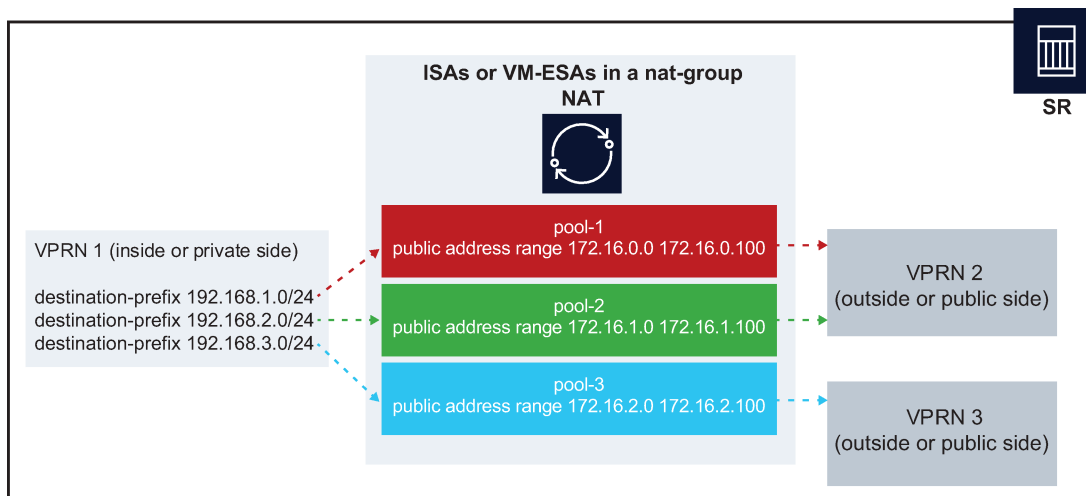
Figure 66: Basic CLI for NAT



sw4100

Figure 67: Logical representation of NAT routing through SR displays the logical configuration.

Figure 67: Logical representation of NAT routing through SR



sw4101

In the routing table, the ISA or VM-ESA next hops for NAT destination prefixes are shown as NAT inside semantics and the listed static routes belong to protocol NAT, as shown in the following example:

```
show router 1 route-table
```

```

=====
Route Table (Service: 1)
=====
Dest Prefix[Flags]          Type  Proto  Age      Pref
  Next Hop[Interface Name]      Metric
-----
192.168.1.0/24             Remote NAT   05d15h21m 0
    NAT inside                0
192.168.2.0/24             Remote NAT   05d15h21m 0
    NAT inside                0
192.168.3.0/24             Remote NAT   05d15h21m 0
    NAT inside                0

```

When forwarding in the downstream direction (in the public to private direction), the outside address ranges (NAT pool ranges 172.16.x.x, as shown in the following example) are subdivided (micronetted) and distributed on a more granular level across the ISAs and VM-ESAs in the same NAT group. The micronetting is necessary so that traffic can be distributed across ISAs or VM-ESAs. The micronets are visible in the routing table and their next hops point to the corresponding ISA or VM-ESA. Micronets are used to attract traffic toward the ISAs or VM-ESAs in the downstream direction. An example of the route table on the outside (public side) is shown in the following output:

```

show router 3 route-table
=====
Route Table (Router: Base)
=====
Dest Prefix[Flags]          Type  Proto  Age      Pref
  Next Hop[Interface Name]      Metric
-----
172.16.2.0/28              Remote NAT   05d15h32m 0
    NAT outside to mda 1/2      0
172.16.2.16/28             Remote NAT   05d15h32m 0
    NAT outside to mda 2/2      0
172.16.2.32/28             Remote NAT   05d15h32m 0
    NAT outside to mda 1/2      0
172.16.2.48/28             Remote NAT   05d15h32m 0
    NAT outside to mda 2/2      0
172.16.2.64/28             Remote NAT   05d15h32m 0
    NAT outside to mda 1/2      0
172.16.2.80/28             Remote NAT   05d15h32m 0
    NAT outside to mda 2/2      0
172.16.2.96/31             Remote NAT   05d15h32m 0
    NAT outside to mda 1/2      0
172.16.2.98/31             Remote NAT   05d15h32m 0
    NAT outside to mda 2/2      0
172.16.2.100/32            Remote NAT   05d15h32m 0
    NAT outside to mda 1/2      0

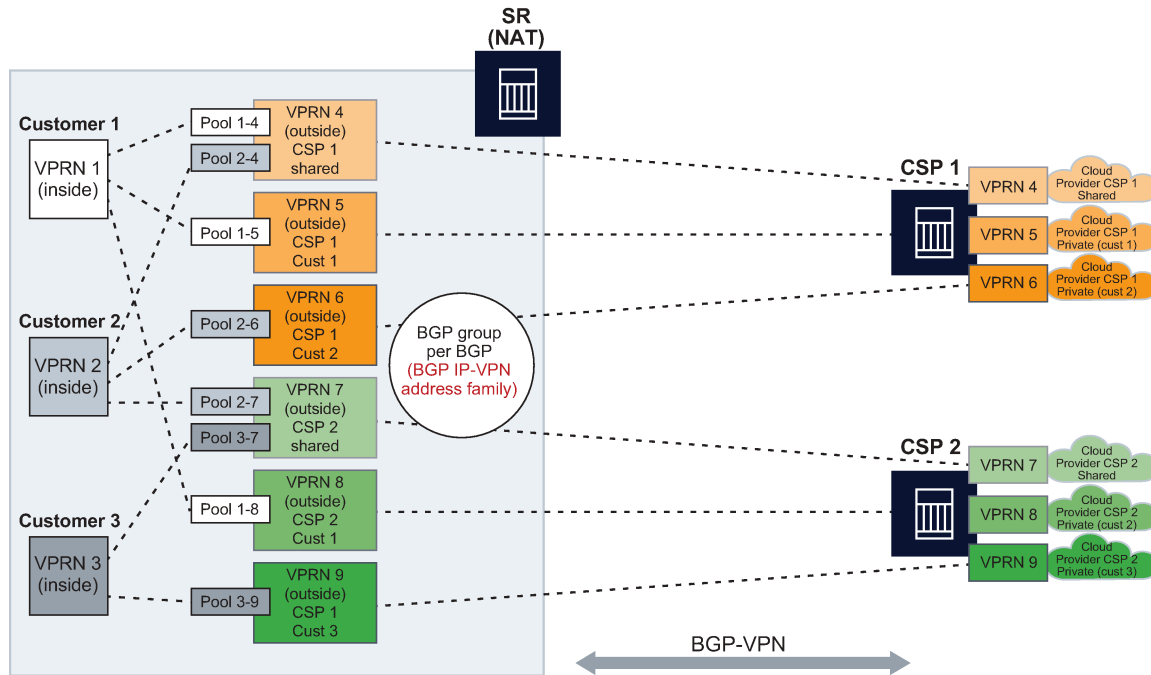
```

7.6.1.2 Dynamic routing to LSN44

LSN44 steering on the inside routing context is supported through routes received from the BGP-VPN protocol. NAT-related BGP-VPN routes are received from BGP-VPN peers in the outside routing contexts and are imported in the inside routing context. This way, the routing information for NAT is dynamically updated without operator intervention. Typical users are operators using NAT who are peering with third parties and frequently re-purposing their IPv4 prefixes based on usage or redundancy. A cloud solution is an example of these types of peering partners.

A typical network design describing this scenario is shown in [Figure 68: Connectivity diagram](#).

Figure 68: Connectivity diagram



Dynamic import of BGP-VPN routes with the next-hop leading to NAT (ISA or VM-ESA) is performed through a routing policy where:

- A route received on the outside by BGP-VPN is matched against the configured route target community and installed in the outside VPRN automatically. Alternatively, the route can be matched against any BGP-VPN supported criteria through an import route policy and installed in the outside routing table.
- Within the NAT context on the inside, the BGP-VPN route is imported through an import route policy, where the route is matched against any BGP-VPN supported criteria. In the same route policy, the route is associated with a NAT policy which determines the NAT pool and outside routing context.

An example configuration for the route policy that is referenced as **import** within a NAT context is shown below.

MD-CLI

```

configure
  policy-options {
    policy-statement <policy-name> {
      entry <id> {
        from {
          : <<any match condition supported for BGP-VPN routes
        }
        action {
          action-type accept
          nat-policy <nat-policy-name>
        }
      }
    }
  }
}
    
```

If the NAT policy under the action is omitted, then a default NAT policy from the inside routing context is used.

The configured route policy is then applied under NAT in the inside routing context.

MD-CLI

```
configure {
  service vprn <name> }
  nat }
  inside }
  large-scale }
  nat-policy <name> <<default nat-policy
  nat44 }
  nat-import <name> }
  {
  {
  {
  {
```

7.6.1.2.1 Deterministic LSN44, non-deterministic LSN44, and 1:1 static LSN44 in a dynamic routing environment

Deterministic and non-deterministic LSN44 can be simultaneously configured in an inside routing context. See [Multiple NAT policies and deterministic NAT](#) for more details.

7.6.1.3 Combination of static and dynamic routes

If the same route is provided by the static configuration and dynamically by a BGP-VPN, only the configured (static) route is installed in the routing table. In other words, a static route has a higher priority than the dynamic route.

7.6.1.4 Scale and logging notes

A NAT route is not installed in the routing table on the inside if one of the following occurs:

- A maximum number of NAT policies per inside VPRN is reached. A NAT policy indirectly represents the route's next hop in the inside routing context (toward the ISA or VM-ESA).
- The next hop in the outside routing context is not available.
- A maximum number of imported NAT routes is reached. There is a maximum numbers of routes that can be imported into the inside routing context from each outside VPRN.
- A maximum number of the dynamic NAT routes per system is reached.

7.6.2 NAT steering through IP filters

Traffic steering to NAT through IP filters is more customizable than steering through routing with traffic identification because of the extensive matching criteria offered by IP filters.

An IP filter can be applied on:

- ingress access and network interfaces on the private side
- network ingress of a VPRN, for example, auto-bind or spoke SDPs
- ingress SLA profile in subscriber management

The filter entry used for NAT steering has an action **nat** which redirects traffic identified through matched criteria toward ISAs and VM-ESAs. The entries in the filter are evaluated the first match ends the filter evaluation.

MD-CLI

```
[pr:/configure filter]
match-list {
  ip-prefix-list "nat-dest" {
    prefix 172.16.0.0/24 { }
  }
}
ip-filter "demo-filter" {
  default-action accept
  entry 10 {
    match {
      protocol udp
      src-port {
        eq 30000
      }
    }
    action {
      accept
    }
  }
  entry 20 {
    match {
      protocol udp
      dst-ip {
        ip-prefix-list "nat-dest"
      }
      src-port {
        range {
          start 40000
          end 50000
        }
      }
    }
    action {
      nat {
      }
    }
  }
  entry 30 {
    match {
      protocol udp
    }
    action {
      drop
    }
  }
}
}
```

In this scenario, any UDP traffic with the source port 3000 as indicated in entry 10, is allowed through the system, bypassing NAT UDP traffic with source ports in the range 40,000 to 50,000 destined for network 172.16.0.0.24 as indicated in entry 20, is NAT'd. The remaining UDP traffic is dropped, according to entry 30.

The remaining non-UDP traffic is allowed through the filter and bypasses NAT as indicated by the default action **accept**.

The following is an example of a filter applied to a network interface on ingress.

MD-CLI

```
[pr:/configure]
router "Base" {
  interface "annex" {
    port 1/x1/2/c10/1:1
    ingress {
      filter {
        ip "demo-filter"
      }
    }
    ipv4 {
      primary {
        address 192.168.12.2
        prefix-length 24
      }
    }
  }
}
```

An example of a filter applied to all ingress network ingress for a specific VPRN.

MD-CLI

```
[pr:/configure service]
vprn "demo" {
  service-id 1
  customer "1"
  network {
    ingress {
      filter {
        ip "demo-filter"
      }
    }
  }
}
```

An example of a filter applied on ingress in an SLA profile.

MD-CLI

```
[pr:/configure service]
subscriber-mgmt {
  sla-profile "demo" {
    ingress {
      ip-filter "demo-filter"
    }
  }
}
```


7.7 L2-Aware support for residential gateway types

L2-Aware NAT functionality is tightly coupled with ESM and therefore, the type of the residential gateway supported in L2-Aware NAT depends on the anti-spoof setting of the ESM subscriber. In this context, the residential gateway types can be:

- **bridged**

Subscriber-hosts behind the residential gateway are individually set up in the BNG and their IP and MAC addresses are known to the BNG during the host setup phase (DHCP/PPPoE).

- **routed with NAT**

The only residential gateway is set up in the BNG. The residential gateway IP and MAC address is known in the BNG during the set up phase. Subscriber hosts behind the residential gateway are not known in the BNG, but instead, they are hidden behind the residential gateway's NAT.

- **routed without NAT**

The residential gateway is set up in the BNG. Hosts behind residential gateway's NAT are not set up in the BNG. The control plane in the BNG is not aware of their IP and MAC addresses. To forward data traffic from these routed hosts in the upstream direction, the anti-spoof in BNG must be set to **nh-mac**. In the downstream direction, a frame route pointing to the residential gateway must be present in the BNG.

Note that DHCP relay on the residential gateway is disabled. If it was enabled, then routed hosts could be set up in BNG with **lease-populate [nbr-of-leases] l2-header [mac ieee-address]** command under the group interface.

Anti-spoof settings in ESM that are relevant to this context include:

- **mac-ip**

Anti-spoof is based on the MAC address and the source IP address of the host. This anti-spoof type is more stringent and secure.

- **nh-mac**

Anti-spoof is based only on the MAC address of the host. This is used in the presence of IP hosts behind the routed RG without NAT. The IP addresses of these hosts are exposed within the data traffic received by BNG, even though those hosts were never explicitly set up in the BNG (using DHCP/PPP). Nh-mac anti-spoof ensures that data traffic from unknown (unknown on the control plane level) IP addresses pass through BNG in the upstream direction. These hosts are behind a known subscriber host, that is, in this case, a routed residential gateway without NAT.

In addition to the anti-spoof setting, an additional CLI command is required in BNG to select the needed residential gateway type:

configure subscr-mgmt sub-prof nat-access-mode {auto | bridged}

The relationship between the anti-spoof setting in ESM, nat-access-mode CLI flag and a compatible residential gateway model is shown in [Table 38: Anti-spoof setting comparisons](#) .

Table 38: Anti-spoof setting comparisons

Model no.	Home model	Anti-spoof	NAT access mode CLI flag	Supported in SR OS	Comments
1	Bridged RG	mac-ip	auto bridged	Yes	All bridged subscriber hosts are eligible for L2-Aware NAT with the most stringent anti-spoof settings. If there is only one host behind the bridged RG, then this model becomes the same as model 3.
2	Bridged RG	nh-mac	bridged	Yes	All bridged subscriber hosts are eligible for L2-Aware NAT. In this model, MAC addresses within the subscriber and SAP must be unique. Even though the anti-spoof in ESM is set to nh-mac , the NAT function still checks the source IP address of the upstream traffic and drops any traffic from spoofed IP addresses (IP source address that do not belong to the bridged hosts, as initially setup in ESM).
3	Routed RG with NAT	mac-ip	auto bridged	Yes	Subscriber hosts behind the residential gateway are hidden behind routed RG's NAT and are not visible in BNG.
4	Routed RG with NAT	nh-mac	auto bridged	Yes	This combination is supported but with inferior anti-spoofing.
5	Routed RG, no NAT	mac-ip	—	No	This combination is not supported. The mac-ip anti-spoof in ESM blocks traffic for the host with an exposed source IP address that resides behind the RG. Those hosts are not set up in the BNG on the control plane level (DHCP/PPPoE is not sent from those hosts).

Model no.	Home model	Anti-spoof	NAT access mode CLI flag	Supported in SR OS	Comments
6	Routed RG, no NAT	nh-mac	auto bridged	Yes	Subscriber hosts with exposed source IP addresses pass the nh-mac anti-spoof check and are eligible for L2-Aware NAT.

7.8 One-to-one (1:1) NAT

In 1:1 NAT, each source IP address is translated in 1:1 fashion to a corresponding outside IP address. However, the source ports are passed transparently without translation.

The mapping between the inside IP addresses and outside IP addresses in 1:1 NAT supports two modes:

- **dynamic**

The operator can specify the outside IP addresses in the pool, but the exact mapping between the inside IP address and the configured outside IP addresses is performed dynamically by the system in a semi-random fashion.

- **static**

The mappings between IP addresses are configurable and they can be explicitly set.

The dynamic version of 1:1 NAT is protocol dependent. Only TCP/UDP/ICMP protocols are allowed to traverse such NAT. All other protocols are discarded, with the exception of PPTP with ALG. In this case, only GRE traffic associated with PPTP is allowed through dynamic 1:1 NAT.

The static version of 1:1 NAT is protocol agnostic. This means that all IP based protocols are allowed to traverse static 1:1 NAT.

The following points are applicable to 1:1 NAT:

- Even though source ports are not being translated, the state maintenance for TCP and UDP traffic is still performed.
- Traffic can be initiated from outside toward any statically mapped IPv4 address.
- 1:1 NAT can be supported simultaneously with NAPT (classic non 1:1 NAT) within the same inside routing context. This is accomplished by configuring two separate NAT pools, one for 1:1 NAT and the other for non 1:1 NAPT.

7.8.1 Static 1:1 NAT

In static 1:1 NAT, inside IP addresses are statically mapped to the outside IP addresses. This way, devices on the outside can predictably initiate traffic to the devices on the inside.

Static configuration is based on the CLI concepts used in deterministic NAT. For example:

```
config
router
nat
inside
deterministic
```

```

prefix 10.0.0.0/24 subscriber-type classic-lsn-sub nat-policy 'one-to-one'
map start 10.0.0.10      end 10.0.0.10      to 1.2.3.4
    map start 10.0.0.15  end 10.0.0.15  to 1.2.3.20
    map start 10.0.0.100 end 10.0.0.100 to 1.2.3.30

```

Static mappings are configured according to the **map** statements. In classic CLI, the map statement can be configured manually by the operator or automatically by the system. In MD-CLI, the map statement must be configured by the operator, but the **tools perform nat deterministic calculate-maps** command can be used to produce system-generated maps. The **calculate-maps** command outputs a set of system-generated map statements. The **map** parameters can then be copied and pasted into an MD-CLI candidate configuration by the operator.

IP addresses from the automatically-generated map statements are sequentially mapped into available outside IP addresses in the pool:

- The first inside IP address is mapped to the first available outside IP address from the pool.
- The second inside IP address is mapped to the second available outside IP address from the pool.

The following mappings apply to the example above:

```

Static mappings
10.0.0.0 - 1.2.3.0
10.0.0.1 - 1.2.3.1
10.0.0.2 - 1.2.3.2
10.0.0.3 - 1.2.3.3
10.0.0.4 - 1.2.3.5
10.0.0.5 - 1.2.3.6
:
10.0.0.9 - 1.2.3.10
10.0.0.10 - 1.2.3.4
10.0.0.11 - 1.2.3.11
10.0.0.12 - 1.2.3.12
:
10.0.0.14 - 1.2.3.14
10.0.0.15 - 1.2.3.20
10.0.0.16 - 1.2.3.15
:
10.0.0.19 - 1.2.3.18
10.0.0.20 - 1.2.3.19
10.0.0.21 - 1.2.3.21
:
10.0.0.28 - 1.2.3.28
10.0.0.29 - 1.2.3.29
10.0.0.30 - 1.2.3.31
:
10.0.0.99 - 1.2.3.100
10.0.0.100 - 1.2.3.30
10.0.0.101 - 1.2.3.101
:
10.0.0.255 - 1.2.3.255

```

7.8.1.1 Protocol agnostic behavior

Although static 1:1 NAT is protocol agnostic, the state maintenance for TCP and UDP traffic is still required to support ALGs. Therefore, the existing scaling limits related to the number of supported flows still apply.

Protocol agnostic behavior in 1:1 NAT is a property of a NAT pool:

```
config
  router / service vprn
    nat
      outside
        pool "one-to-one" nat-group 1 type large-scale applications agnostic create
        address-range 192.168.0.0 192.168.0.10 create
```

The application **agnostic** command is a pool create-time parameter. This command automatically pre-sets the following pool parameters:

```
mode one-to-one
port-forwarding-range 0
port-reservation blocks 1
subscriber-limit 1
deterministic port-reservation 65536.
```

When pre-set, these parameters cannot be changed while the pool is operating in protocol agnostic mode.

The **deterministic port-reservation 65536** command configures the pool to operate in static (or deterministic) mode.

7.8.1.2 Modification of parameters in static 1:1 NAT

Parameters in the static 1:1 NAT can be changed according to the following rules:

- The deterministic pool must be in a **no shutdown** state when a **prefix** or a **map** command in deterministic NAT is added or removed.
- All configured prefixes referencing the pool via the NAT policy must be deleted (unconfigured) before the pool can be shut down.
- Map statements can be modified only when prefix is shutdown state. All existing map statements must be removed before the new ones are created.

7.8.1.3 Load distribution over ISAs in static 1:1 NAT

For best traffic distribution over ISAs, the value of the **classic-lsn-max-subscriber-limit** *max* parameter should be set to 1.

```
config
  router / service vprn X
    nat
      inside
        deterministic
        classic-lsn-max-subscriber-limit <num>
```

This means that traffic is load balanced over ISAs based on inside IP addresses. In static 1:1 NAT this is certainly possible because the subscriber-limit parameter at the pool level is preset to a fixed value of 1.

However, if 1:1 static NAT is simultaneously used with regular (many-to-one) deterministic NAT where the subscriber-limit parameter can be set to a value greater than 1, then the classic-lsn-max-subscriber-limit parameter also has to be set to a value that is greater than 1. The consequence of this is that the traffic is load balanced based on the consecutive blocks of IP addresses (subnets) instead of individual IP addresses. See [Deterministic NAT](#) for information about Deterministic NAT behavior.

7.8.1.4 NAT-policy selection

The traffic match criteria used in the selection of specific NAT policies in static 1:1 NAT (the deterministic part of the configuration) must not overlap with traffic match criteria that is used in the selection of a specific NAT policy used in filters or in destination-prefix statement (these are used for traffic diversion to NAT). Otherwise, traffic is dropped in ISA.

A specific NAT policy in this context refers to a non-default NAT policy, or a NAT policy that is directly referenced in a filter, in a **destination-prefix** command or in a **deterministic prefix** command.

The following example is used to clarify this point:

- Traffic is diverted to NAT using specific **nat-policy pol-2**:

```
service vprn 10
  nat
    inside
      destination-prefix 192.168.0.0/16 nat-policy pol-2
      deterministic
      prefix 10.10.10.0/24 subscriber-type classic-lsn-sub nat-policy pol-1
```

- The deterministic (source) prefix 10.10.10.0/30 is configured to be mapped to **nat-policy pol-1** specifically which points to protocol agnostic 1:1 nat pool.

```
service vprn 10
  nat
    inside
      destination-prefix 192.168.0.0/16 nat-policy pol-2
      deterministic
      prefix 10.10.10.0/30 subscriber-type classic-lsn-sub nat-policy pol-1
```

- Packet received in the ISA has srcIP 10.10.10.1 and destIP 192.168.10.10.
- If no NAT mapping for this traffic exists in the ISA, a NAT policy (and with this, the NAT pool) must be determined to create the mapping. Traffic is diverted to NAT using **nat-policy pol-2**, while the deterministic mapping suggests that **nat-policy pol-1** should be used (this is a different pool from the one referenced in **nat-policy pol-2**). Because of the specific NAT policy conflict, traffic is dropped in the ISA.

To successfully pass traffic between two subnets through NAT while simultaneously using static 1:1 NAT and regular LSN44, a default (**non-specific**) NAT policy can be used for regular LSN44.

For example:

```
service vprn 10
  nat
    inside
      destination-prefix 192.168.0.0/16
      nat-policy pol-2
      deterministic
      prefix 10.10.10.0/30 subscriber-type classic-lsn-sub nat-policy pol-1
```

In this case, the four hosts from the prefix 10.10.10.0/30 are mapped in 1:1 fashion to 4 IP addresses from the pool referenced in the specific nat-policy pol-1, while all other hosts from the 10.10.10.0/24 network are mapped to the NAPT pool referenced by the default nat-policy pol-2. In way, a NAT policy conflict is avoided.

In summary, a specific NAT policy (in filter, **destination-prefix** command or in deterministic **prefix** command) always takes precedence over a default NAT policy. However, traffic that matches classification criteria (in filter, **destination-prefix** command or a deterministic **prefix** command) that leads to multiple specific nat-policies, is dropped.

7.8.1.5 Mapping timeout

Static 1:1 NAT mappings are explicitly configured, and therefore, their lifetime is tied to the configuration.

7.8.1.6 Logging

The logging mechanism for static mapping is the same as in Deterministic NAT. Configuration changes are logged via syslog and enhanced with reverse querying on the system.

7.8.1.7 Restrictions

Static 1:1 NAT is supported only for LSN44 (there is no support for DS-Lite/NAT64 or L2-Aware NAT).

7.8.2 ICMP

In 1:1 NAT, specific ICMP messages contain an additional IP header embedded in the ICMP header. For example, when the ICMP message is sent to the source because of the inability to deliver datagram to its destination, the ICMP generating node includes the original IP header of the packet plus 64bits of the original datagram. This information helps the source node to match the ICMP message to the process associated with this message.

When these messages are received in the downstream direction (on the outside), 1:1 NAT recognizes them and changes the destination IP address not only in the outside header but also in the ICMP header. In other words, a lookup in the downstream direction is performed in the ISA to determine if the packet is ICMP with a specific type. Depending on the outcome, the destination IP address in the ICMP header is changed (reverted to the original source IP address).

Messages carrying the original IP header within ICMP header are:

- Destination Unreachable Messages (Type 3)
- Time Exceeded Message (Type 11)
- Parameter Problem Message (Type 12)
- Source Quench Message (Type 4)

7.9 Deterministic NAT

7.9.1 Overview

In deterministic NAT the subscriber is deterministically mapped into an outside IP address and a port block. The algorithm that performs this deterministic mapping is revertive, which means that a NAT subscriber

can be uniformly derived from the outside IP address and the outside port (and the routing instance). Thus, logging in deterministic NAT is not needed.

The deterministic [subscriber <-> outside-ip, deterministic-port-block] mapping can be automatically extended by a dynamic port-block in case that deterministic port block becomes exhausted of ports. By extending the original deterministic port block of the NAT subscriber by a dynamic port block yields a satisfactory compromise between a deterministic NAT and a non-deterministic NAT. There is no logging as long as the translations are in the domain of the deterministic NAT. After the dynamic port block is allocated for port extension, logging is automatically activated.

NAT subscribers in deterministic NAT are not assigned outside IP address and deterministic port-block on a first come first serve basis. Instead, deterministic mappings are pre-created at the time of configuration regardless of whether the NAT subscriber is active or not. In other words we can say that overbooking of the outside address pool is not supported in deterministic NAT. Consequently, all configured deterministic subscribers (for example, inside IP addresses in LSN44 or IPv6 address/prefix in DS-Lite) are guaranteed access to NAT resources.

7.9.2 Supported deterministic NAT types

The routers support Deterministic LSN44 and Deterministic DS-Lite. The basic deterministic NAT principle is applied equally to both NAT flavors. The difference between the two stem from the difference in interpretation of the subscriber – in LSN44 a subscriber is an IPv4 address, whereas in DS-Lite the subscriber is an IPv6 address or prefix (configuration dependent).

With the exception of **classic-lsn-max-subscriber-limit** and **dslite-max-subscriber-limit** commands in the inside routing context, the deterministic NAT configuration blocks are for the most part common to LSN44 and DS-Lite.

Deterministic DS-Lite section at the end of this section focuses on the features specific to DS-Lite.

7.9.3 Number of subscribers per outside IP and per pool

The outside pools in deterministic NAT can contain an arbitrary number of address ranges, where each address range can contain an arbitrary number of IP addresses (up to the ISA maximum).

The maximum number of NAT subscribers that can be mapped to a single outside IP address is configurable using a **subscriber-limit** command under the pool hierarchy. For Deterministic NAT, this number is restricted to the power of 2 (2^n). The consequence of this is that the number of NAT subscribers must be configuration-wise organized in ranges with the boundary that must be power of 2.

For example, in LSN44 where the NAT subscriber is an IP address, the deterministic subscribers would be configured with prefixes (for example, 10.10.10.0/24 – 256 subscribers) instead of an IP address range that would contain an arbitrary number of addresses (for example, 10.10.10.10 – 10.10.10.50).

On the other hand, in DS-Lite the deterministic subscribers are for the most part already determined by the prefix with the **subscriber-prefix-length** command under the DS-Lite configuration node.

The number of subscribers per outside IP (the **subscriber-limit** command [2^n]) multiplied by the number of IP addresses over all address-range in an outside pool determines the maximum number of subscribers that a deterministic pool can support.

7.9.4 Referencing a pool

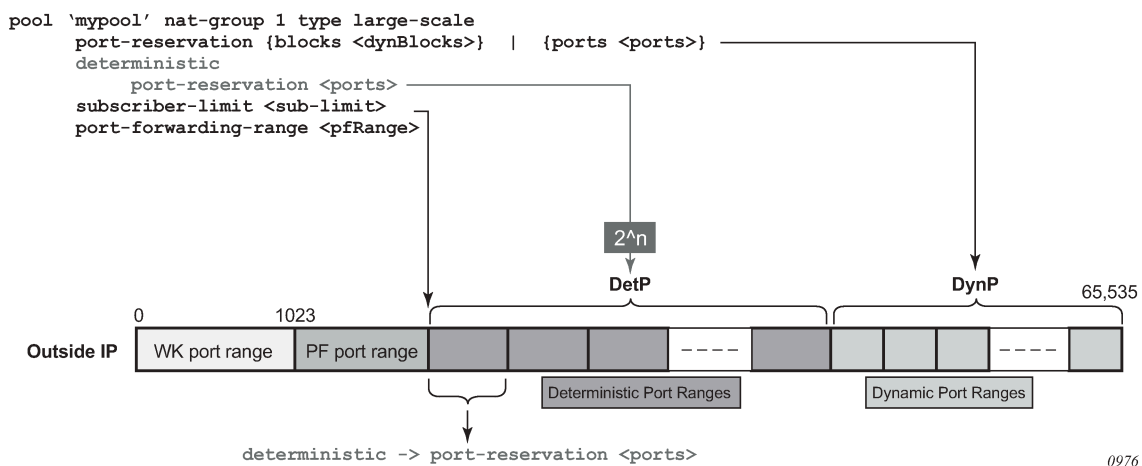
In deterministic NAT, the outside pool can be shared amongst subscribers from multiple routing instances. Also, NAT subscribers from a single routing instance can be selectively mapped to different outside pools.

7.9.5 Outside pool configuration

The number of deterministic mappings that a single outside IP address can sustain is determined through the configuration of the outside pool.

The port allocation per an outside IP is shown in [Figure 69: Outside pool configuration](#).

Figure 69: Outside pool configuration



The well-known ports are predetermined and are in the range 0 — 1023.

The upper limit of the port range for static port forwards (wildcard range) is determined by the existing **port-forwarding-range** command.

The range of ports allocated for deterministic mappings (DetP) is determined by multiplying the number of subscribers per outside IP (**subscriber-limit** command) with the number of ports per deterministic block (**deterministic>port-reservation** command). The number of subscribers per outside IP in deterministic NAT must be power of 2 (2^n).

The remaining ports, extending from the end of the deterministic port range to the end of the total port range (65,535) are used for dynamic port allocation. The size of each dynamic port block is determined with the existing **port-reservation** command.

The **deterministic>port-reservation** command enables deterministic mode of operation for the pool.

Examples:

Three examples follow, with deterministic Large Scale NAT44, where the requirements are:

- 300, 500 or 700 (three separate examples) ports are in each deterministic port block.
- A subscriber (an inside IPv4 address in LSN44) can extend its deterministic ports by a minimum of one dynamic port-block and by a maximum of four dynamic port blocks.
- Each dynamic port-block contains 100 ports.

- Oversubscription of dynamic port blocks is 4:1. This means that 1/4th of inside IP addresses may be starved out of dynamic port blocks in worst case scenario.
- The wildcard (static) port range is 3000 ports.

In the first case, the ideal case is examined where an arbitrary number of subscribers per outside IP address is allocated according to our requirements described above. Then the limitation of the number of subscribers being power of 2 is factored in.

Table 39: Contiguous number of subscribers

Well-known ports ⁴	Static port range ⁴	Number of ports in deterministic block ⁴	Number of deterministic blocks	Number of ports in dynamic block ⁴	Number of dynamic blocks ⁴	Number of inside IP addresses per outside IP address ⁴	Block limit per inside IP address ⁴	Wasted ports
0-1023	1024-4023	300	153	100	153	153	5	312
0-1023	1024-4023	500	102	100	102	102	5	312
0-1023	1024-4023	700	76	100	76	76	5	712

The example in [Table 39: Contiguous number of subscribers](#) shows how port ranges would be carved out in ideal scenario.

The other values are calculated according to the fixed requirements.

port-block-limit includes the deterministic port block plus all dynamic port-blocks.

Next, in [Table 40: Preserving Det/Dyn port ratio with 2^n subscribers](#), a more realistic example with the number of subscribers being equal to 2^n are considered. The ratio between the deterministic ports and the dynamic ports per port-block just like in the example above: 3/1, 5/1 and 7/1 are preserved. In this case, the number of ports per port-block is dictated by the number of subscribers per outside IP address.

Table 40: Preserving Det/Dyn port ratio with 2^n subscribers

Well-known ports ⁴	Static port range ⁴	Number of ports in deterministic block ⁴	Number of deterministic blocks	Number of ports in dynamic block ⁴	Number of dynamic blocks	Number of inside IP addresses per outside IP address ⁴	Block limit per inside IP address ⁴	Wasted ports
0-1023	1024-4023	180	256	60	256	256	5	72
0-1023	1024-4023	400	128	80	128	128	5	72
0-1023	1024-4023	840	64	120	64	64	5	72

⁴ Signifies the fixed parameters (requirements)

The final example ([Table 41: Fixed number of deterministic ports with 2^n subscribers](#)) is similar as [Table 39: Contiguous number of subscribers](#) with the difference that the number of deterministic port blocks fixed are kept, as in the original example (300, 500 and 700).

Table 41: Fixed number of deterministic ports with 2^n subscribers

Well-known ports	Static port range	Number of ports in deterministic block	Number of deterministic blocks	Number of ports in dynamic block	Number of dynamic blocks	Number of inside IP addresses per outside IP	Block limit per inside IP address	Wasted ports
0-1023	1024-4023	300	128	180	128	128	5	72
0-1023	1024-4023	500	64	461	64	64	5	8
0-1023	1024-4023	700	64	261	64	64	5	8

The three examples from above should give us a perspective on the size of deterministic and dynamic port blocks in relation to the number of subscribers (2^n) per outside IP address. Operators should run a similar dimensioning exercise before they start configuring their deterministic NAT.

The CLI for the highlighted case in the [Table 39: Contiguous number of subscribers](#) is displayed:

```

configure
  service
    vprn
      nat
        outside
          pool mypool
            port-reservation ports 180
            deterministic
          port-reservation 300
            subscriber-limit 128
            port-forwarding-range 4023
    
```

Where:

128 subs * 300ports = 38,400 deterministic port range

128 subs * 180ports = 23,040 dynamic port range

Det+dyn available ports = 65,536 – 4024 = 61,512

Det+dyn usable pots = 128*300 + 128 *180 = 61,440 ports

72 ports per outside-ip are wasted.

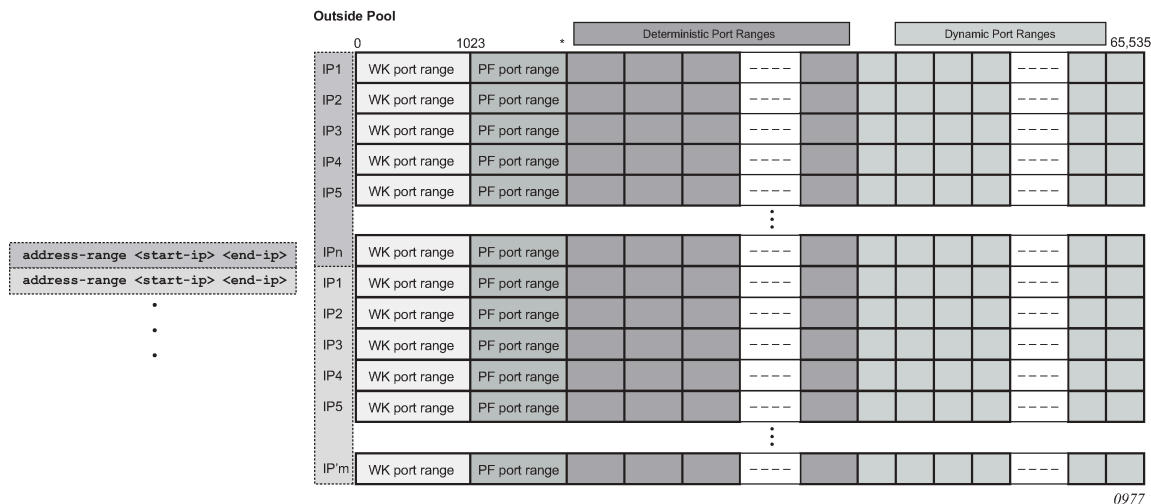
```

configure
  service
    nat
      nat-policy mypolicy
        block-limit 5 # 1 deterministic port block + 4 dynamic port blocks
    
```

This configuration allows 128 subscribers (inside IP addresses in LSN44) for each outside address (compression ratio is 128:1) with each subscriber being assigned up to 1020 ports (300 deterministic and 720 dynamic ports over 4 dynamic port blocks).

The outside IP addresses in the pool and their corresponding port ranges are organized as shown in [Figure 70: Outside address ranges](#).

Figure 70: Outside address ranges



Assuming that the above graph depicts an outside deterministic pool, the number of subscribers that can be accommodated by this deterministic pool is represented by purple squares (number of IP addresses in an outside pool * subscriber-limit). The number of subscribers across all configured prefixes on the inside that are mapped to the same deterministic pool must be less than the outside pool can accommodate. In other words, an outside address pool in deterministic NAT cannot be oversubscribed.

The following is a CLI representation of a deterministic pool definition including the outside IP ranges:

```
pool 'mypool' nat-group 1 type large-scale
  port-reservation {blocks <dynBlocks>} | {ports <ports>}
  deterministic
  port-reservation <ports>
  subscriber-limit <sub-limit>
  port-forwarding-range <pfRange>
  address-range <start-ip-address> <end-ip-address>
  address-range <start-ip-address> <end-ip-address>
```

7.9.6 Mapping rules and the map command in deterministic LSN44

The common building block on the inside in the deterministic LSN44 configuration is a IPv4 prefix. The NAT subscribers (inside IPv4 addresses) from the configured prefix are deterministically mapped to the outside IP addresses and corresponding deterministic port-blocks. Any inside prefix in any routing instance can be mapped to any pool in any routing instance (including the one in which the inside prefix is defined).

The mapping between the inside prefix and the deterministic pool is achieved through a NAT policy that can be referenced per each individual inside IPv4 prefix. IPv4 addresses from the prefixes on the inside are distributed over the IP addresses defined in the outside pool referenced by the NAT policy.

The mapping itself is represented by the **map** command under the prefix hierarchy:

```
router/service vprn
  nat
    inside
      deterministic
        prefix <ip-prefix/length> subscriber-type <nat-sub-type> nat-policy <nat-policy-
name>
        map start <inside-ip-address> end <inside-ip-address> to <outside-ip-address>
```

The purpose of the map statement is to split the number of subscribers within the configured prefix over available sequences of outside IP addresses. The key parameter that governs mappings between the inside IPv4 addresses and outside IPv4 addresses in deterministic LSN44 is defined by the **outside>pool>subscriber-limit** command. This parameter must be power of 2 and it limits the maximum number of NAT subscribers that can be mapped to the same outside IP address.

The follow are rules governing the configuration of the map statement:

In case that the number of subscribers (IP addresses in LSN44) in the **map** statement is larger than the subscriber-limit per outside IP, then the subscribers must be split over a block of consecutive outside IP addresses where the *outside-ip-address* in the map statement represent only the first outside IP address in that block.

The number of subscribers (range of inside IP addresses in LSN44) in the map statement does not have to be a power of 2. Rather it has to be a multiple of a power of two $m * 2^n$, where m is the number of consecutive outside IP addresses to which the subscribers are mapped and the 2^n is the subscriber-limit per outside IP.

An example of the map statement is shown below:

```
router
  nat
    outside
      pool 'my-det-pool' nat-group 1 type large-scale
      subscriber-limit 128
      deterministic
        port-reservation 400
        address-range 192.168.0.0 192.168.0.10

  service vprn 10
  nat
    inside
      deterministic
        prefix 10.0.0.0/24 subscriber-type classic-lsn-sub nat-policy det
        map start 10.0.0.0 end 10.0.0.255 to 192.168.0.1
```

In this case, the configured 10.0.0.0/24 prefix is represented by the range of IP addresses in the map statement (10.0.0.0-10.0.0.255). Because the range of 256 IP addresses in the map statement cannot be mapped into a single outside IP address (subscriber-limit=128), this range must be further implicitly split within the system and mapped into multiple outside IP addresses. The implicit split creates two IP address ranges, each with 128 IP addresses (10.0.0.0/25 and 10.0.0.128/25) so that addresses from each IP range are mapped to one outside IP address. The hosts from the range 10.0.0.0-10.0.0.127 are mapped to the first IP address in the pool (128.251.0.1) as explicitly stated in the map statement (to statement). The hosts from the second range, 10.0.0.128-10.0.0.255 are implicitly mapped to the next consecutive IP address (128.251.0.2).

Alternatively, the **map** statement can be configured as:

```
service vprn 10
```

```

nat
  inside
    deterministic
      prefix 10.0.0.0/24 subscriber-type classic-lsn-sub nat-policy det
        map start 10.0.0.0 end 10.0.0.127 to 192.168.0.1
        map start 10.0.0.128 end 10.0.0.255 to 192.168.0.5

```

In this case the IP address range in the map statement is split into two non-consecutive outside IP addresses. This gives the operator more freedom in configuring the mappings.

However, the following configuration is not supported:

```

service vprn 10
nat
  inside
    deterministic
      prefix 10.0.0.0/24 subscriber-type classic-lsn-sub nat-policy det
        map start 10.0.0.0 end 10.0.0.63 to 192.168.0.1
        map start 10.0.0.64 end 10.0.0.127 to 192.168.0.3
        map start 10.0.0.128 end 10.0.0.255 to 192.168.0.5

```

Considering that the subscriber-limit = 128 (2^n ; where $n=7$), the lower n bits of the start address in the second map statement (map start 10.0.0.64 end 10.0.0.127 to 192.168.0.3) are not 0. This is in violation of the rule #1 that governs the provisioning of the map statement.

Assuming that we use the same pool with 128 subscribers per outside IP address, the following scenario is also not supported (configured prefix in this example is different than in previous example):

```

service vprn 10
nat
  inside
    deterministic

prefix 10.0.0.0/26 subscriber-type classic-lsn-sub nat-policy det
  map start 10.0.0.0 end 10.0.0.63 to 192.168.0.1

prefix 10.0.1.0/26 subscriber-type classic-lsn-sub nat-policy det
  map start 10.0.1.0 end 10.0.1.63 to 192.168.0.1

```

Although the lower n bits in both map statements are 0, both statements are referencing the same outside IP (192.168.0.1). This is violating rule #2 that governs the provisioning of the map statement. Each of the prefixes in this case have to be mapped to a different outside IP address, which leads to underutilization of outside IP addresses (half of the deterministic port-blocks in each of the two outside IP addresses are not used).

In conclusion, considering that the number of subscribers per outside IP (subscriber-limit) must be 2^n , the inside IP addresses from the configured prefix is split on the 2^n boundary so that every deterministic port-block of an outside IP is used. In case that the originally configured prefix contains less subscribers (IP addresses in LSN44) than an outside IP address can accommodate (2^n), all subscribers from such configured prefix are mapped to a single outside IP. Because the outside IP cannot be shared with NAT subscribers from other prefixes, some of the deterministic port-blocks for this particular outside IP address are not used.

Each configured prefix can evaluate into multiple **map** commands. The number of **map** commands depends on the length of the configured prefix, the **subscriber-limit** command and fragmentation of outside address-range within the pool with which the prefix is associated.

In classic CLI, the **map** statement can be configured manually by the operator or automatically by the system. In MD-CLI, the **map** statement must be configured by the operator, but the **tools perform nat**

deterministic calculate-maps command can be used to produce system-generated maps if needed. The **calculate-maps** command outputs a set of system-generated map statements. The **map** parameters can then be copied and pasted into an MD-CLI candidate configuration by the operator.

- If the number of subscribers per configured prefix is greater than the subscriber-limit per outside IP parameter (2^n), then the lowest n bits of the **map start** *inside-ip-address* must be set to 0.
- If the number of subscribers per configured prefix is equal or less than the subscriber-limit per outside IP parameter (2^n), then only one **map** command for this prefix is allowed. In this case there is no restriction on the lower n bits of the **map start** *inside-ip-address*. The range of the inside IP addresses in such map statement represents the prefix itself.
- The *outside-ip-address* in the map statements must be unique amongst all map statements referencing the same pool. In other words, two map statements cannot reference the same *outside-ip-address* in the pool.

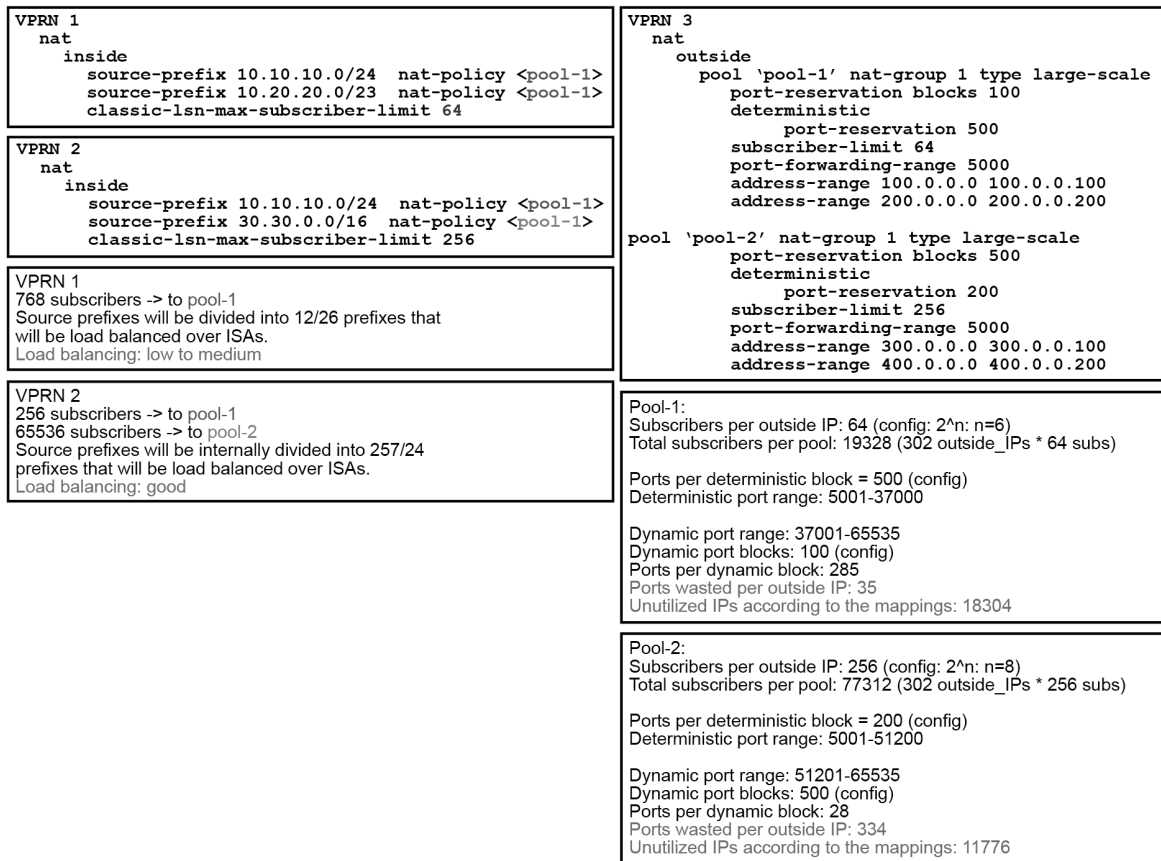
7.9.7 Hashing considerations in deterministic LSN44

Support for multiple MS-ISAs in the nat-group calls for traffic hashing on the inside in the ingress direction. This ensures fair load balancing of the traffic amongst multiple MS-ISAs. While hashing in non-deterministic LSN44 can be performed per source IP address, hashing in deterministic LSN44 is based on subnets instead of individual IP addresses. The length of the hashing subnet is common for all configured prefixes within an inside routing instance. In case that a prefixes from an inside routing instances is referencing multiple pools, the common hashing prefix length is chosen according to the pool with the highest number of subscribers per outside IP address. This ensures that subscribers mapped to the same outside IP address are always hashed to the same MS-ISA.

In general, load distribution based on hashing is dependent on the sample. Large and more diverse sample ensures better load balancing. Therefore the efficiency of load distribution between the MS-ISAs is dependent on the number and diversity of subnets that hashing algorithm is taking into consideration within the inside routing context.

A simple rule for good load balancing is to configure a large number of subscribers relative to the largest subscriber limit in any pool that is referenced from this inside routing instance.

Figure 71: Deterministic LSN44 configuration example



0979

Figure 71: Deterministic LSN44 configuration example shows a case in which prefixes from multiple routing instances are mapped to the same outside pool and at the same time the prefixes from a single inside routing instance are mapped to different pools (Nokia does not support the latter with non-deterministic NAT).



Note: In this example is the inside prefix 10.10.10.0/24 that is present in VPRN 1 and VPRN 2. In both VPRNs, this prefix is mapped to the same pool (pool-1) with the subscriber limit of 64. Four outside IP addresses per prefix per VPRN (eight in total) are allocated to accommodate the mappings for all hosts in prefix 10.10.10.0/24. However, the hashing prefix length in VPRN 1 is based on the subscriber limit of 64 (VPRN 1 references only pool-1) while the hashing prefix length in VPRN 2 is based on the subscriber limit of 256 in pool-2. VPRN 2 references both pools, pool-1 and pool-2, and the larger subscriber limit must be selected. The consequence of this is that the traffic from subnet 10.10.10.0/24 in VPRN 1 can be load balanced over 4 MS-ISA (hashing prefix length is 26) while traffic from the subnet 10.10.10.0/24 in VPRN 2 is always sent to the same MS-ISA (hashing prefix length is 24).

7.9.7.1 Distribution of outside IP addresses across MS-ISAs in an MS-ISA NA group

Distribution of outside IP addresses across the MS-ISAs is dependent on the ingress hashing algorithm. Because traffic from the same subscriber is always pre-hashed to the same MS-ISA, the corresponding outside IP address also must reside on the same ISA. CPM runs the hashing algorithm in advance to determine on which MS-ISA the traffic from particular inside subnet lands and then the corresponding outside IP address (according to deterministic NAT mapping algorithm) is configured in that particular MS-ISA.

7.9.8 Sharing of deterministic NAT pools

Sharing of the deterministic pools between LSN44 and DS-Lite is supported.

7.9.9 Simultaneous support of dynamic and deterministic NAT

Simultaneous support for deterministic and non-deterministic NAT inside of the same routing instance is supported. However, an outside pool can be only deterministic (although expandable by dynamic ports blocks) or non-deterministic at any time.

Ingress hashing for all NAT'd traffic within the VRF, in this case, is performed based on the subnets driven by the classic-lsn-max-subscriber-limit parameter.

7.9.10 Selecting traffic for NAT

Deterministic NAT does not change the way how traffic is selected for the NAT function but instead only defines a predictable way for translating subscribers into outside IP addresses and port-blocks.

Traffic is still diverted to NAT using the existing methods:

- **routing based**

Traffic is forwarded to the NAT function if it matches a configured destination prefix that is part of the routing table. In this case inside and outside routing context must be separated.

- **filter based**

Traffic is forwarded to the NAT function based on any criteria that can be defined inside an IP filter. In this case the inside and outside routing context can be the same.

7.9.11 Inverse mappings

The inverse mapping can be performed with a MIB locally on the node or externally via a script sourced in the router. In both cases, the input parameters are <outside routing instance, outside IP, outside port. The output from the mapping is the subscriber and the inside routing context in which the subscriber resides.

7.9.11.1 MIB approach

Reverse mapping information can be obtained using the following command:

```
tools dump nat deterministic-mapping outside-ip <ipv4-address> router <router-instance>
outside-port <[1..65535]>
<ipv4-address>      : a.b.c.d
<router-instance>  : <router-name> | <service-id>
                    router-name   - "Base"
                    service-id    - [1..2147483647]
```

Example:

```
tools dump nat deterministic-mapping outside-ip 10.0.0.2 router "Base" outside-port 2333
```

Output:

```
Inside router 10 ip 10.0.5.171 -- outside router Base ip 10.0.0.2 port 2333 at Mon Jan 7 10:02:02 PST 2013
```

7.9.11.2 Off-line approach to obtain deterministic mappings

Instead of querying the system directly, there is an option where a Python script can be generated on router and exported to an external node. This Python script contains mapping logic for the configured deterministic NAT in the router. The script can be then queried off-line to obtain mappings in either direction. The external node must have installed Python scripting language with the following modules: getopt, math, os, socket and sys.

The purpose of such off-line approach is to provide fast queries without accessing the router. Exporting the Python script for reverse querying is a manual operation that needs to be repeated every time there is configuration change in deterministic NAT.

The script is exported outside of the box to a remote location (assuming that writing permissions on the external node are correctly set). The remote location is specified with the following command:

```
config service nat deterministic-script location <remote-url>
<remote-url>      - [{ftp:// | tftp://}<login>:<pswd>@<remote-locn>/][<file-path>]
180 chars max
```

The status of the script is shown using the following command:

```
show service nat deterministic-script
=====
Deterministic NAT script data
=====
Location          : ftp://10.10.10.10/pub/det-nat-script/det-nat.py
Save needed       : yes
Last save result  : none
Last save time    : N/A
=====
```

After the script location is specified, the script can be exported to that location with the following command:

```
admin nat save-deterministic-script
```

This needs to be repeated manually every time the configuration affecting deterministic NAT changes.

```
Once the script is exported (saved), the status of the script is changed as well:
show service nat deterministic-script
=====
Deterministic NAT script data
=====
Location          : ftp://10.10.10.10/pub/det-nat-script/det-nat.py
Save needed       : no
Last save result  : success
Last save time    : 2013/01/07 10:33:43
=====
```

The script itself can be run to obtain mapping in forward or backward direction:

```
user@external-server:/home/ftp/pub/det-nat-script$ ./det-nat.py
Usage: det-nat-.py {{DIRECTION PARAMS}} | -h[elp] }
where DIRECTION := { -f[orward] | -b[ackward] }
where PARAMS := { -s[ervice] -a[ddress] -p[ort] }
```

The following displays an example in which source addresses are mapped in the following manner:

```
Router 10, Source-ip: 10.0.5.0-10.0.5.127    to router base, outside-ip 10.0.0.1
Router 10 Source-ip: 10.0.5.128-10.0.5.255  to router base outside-ip 10.0.0.2
```

The forward query for this example is performed as:

```
user@external-server:/home/ftp/pub/det-nat-script$ ./det-nat.py -f -s 10 -a 10.0.5.10
```

Output:

```
subscriber has public ip address 10.0.0.1 from service 0 and is using ports [1324 - 1353]
```

The reverse query for this example is performed as:

```
user@external-server:/home/ftp/pub/det-nat-script$ ./det-nat.py -b -s 0 -a 10.0.0.1 -p 3020
```

Output:

```
subscriber has private ip address 10.0.5.66 from service 10
```

7.9.12 Logging

Every configuration change concerning the deterministic pool is logged and the script (if configured for export) is automatically updated (although not exported). This is needed to keep current track of deterministic mappings. In addition, every time a deterministic port-block is extended by a dynamic block, the dynamic block is logged just as it is today in non-deterministic NAT. The same logic is followed when the dynamic block is de-allocated.

All static port forwards (including PCP) are also logged.

PCP allocates static port forwards from the wildcard-port range.

7.9.13 Deterministic DS-Lite

A subscriber in non-deterministic DS-Lite is defined as v6 prefix, with the prefix length being configured under the DS-Lite NAT node:

```
config>service>vprn>nat>inside>dslite#
  subscriber-prefix-length [32..64 | 128]      (default is 128)
```

All incoming IPv6 traffic with source IPv6 addresses falling under a unique v6 prefix that is configured with **subscriber-prefix-length** command is considered as a single subscriber. As a result, all source IPv4 addresses carried within that IPv6 prefix are mapped to the same outside IPv4 address.

The concept of deterministic DS-Lite is very similar to deterministic LSN44. The DS-Lite subscribers (IPv6 addresses/prefixes) are deterministically mapped to outside IPv4 addresses and corresponding deterministic port-blocks.

Although the subscriber in DS-Lite is considered to be either a B4 element (IPv6 address) or the aggregation of B4 elements (IPv6 prefix determined by the **subscriber-prefix-length** command), only the IPv4 source addresses and ports carried inside of the IPv6 tunnel are actually translated.

The prefix statement for deterministic DS-Lite remains under the same deterministic CLI node as for the deterministic LSN44. However, the prefix statement parameters for deterministic DS-Lite differ from the one for deterministic LSN44 in the following fashion:

- DS-Lite prefix is a v6 prefix (instead of v4). The DS-Lite subscriber whose traffic is mapped to a particular outside IPv4 address and the deterministic port block is deduced from the prefix statement and the subscriber-prefix-length statement.
- Subscriber-type is set to dslite-lsn-sub.

```
config>service>vprn>nat>inside>deterministic#
  prefix <v6-prefix/length> subscriber-type dslite-lsn-sub nat-policy <policy-name>
```

Example:

```
config>service>vprn>nat>inside>deterministic#
  prefix 2001:db8::/56 subscriber-type dslite-lsn-sub nat-policy det-policy

config>service>vprn>nat>inside>dslite#
  subscriber-prefix-length 60
```

In this case, 16 v6 prefixes (from 2001:db8::/60 to 2001:db8:00:F0::/60) are considered DS-Lite subscribers. The source IPv4 addresses/ports inside of the IPv6 tunnels is mapped into respective deterministic port blocks within an outside IPv4 address according to the map statement.

The map statement contains minor modifications as well. It maps DS-Lite subscribers (IPv6 address or prefix) to corresponding outside IPv4 addresses. Continuing on the previous example:

```
map start 2001:db8::/60 end 2001:db8:00:F0::/60 to 192.168.1.1
```

The prefix length (/60) in this case must be the same as configured subscriber-prefix-length. If we assume that the subscriber-limit in the corresponding pool is set to 8 and outside IP address range is 192.168.1.1 to 192.168.1.10, then the actual mapping is the following:

```
2001:db8::/60 to 2001:db8:00:70::/60 to 192.168.1.1
2001:db8:00:80::/60 to 2001:db8:00:F0::/60 to 192.168.1.2
```

7.9.13.1 Hashing considerations in DS-Lite

The ingress hashing and load distribution between the ISAs in Deterministic DS-Lite is governed by the highest number of configured subscribers per outside IP address in any pool referenced within the specific inside routing context.

This limit is configured under:

```
configure
router/service vprn
  nat
    inside
      deterministic
        dslite-max-subscriber-limit <1,2,4,8...32768>
```

While ingress hashing in non-deterministic DS-Lite is governed by the **subscriber-prefix-length** command, in deterministic DS-Lite the ingress hashing is governed by the combination of **dslite-max-subscriber-limit** and **subscriber-prefix-length** commands. This is to ensure that all DS-Lite subscribers that are mapped to a single outside IP address are always sent to the same MS-ISA (on which that outside IPv4 address resides). In essence, as soon as deterministic DS-Lite is enabled, the ingress hashing is performed on an aggregated set of $n = \log_2(\text{dslite-max-subscriber-limit})$ contiguous subscribers. n is the number of bits used to represent the largest number of subscribers within an inside routing context, that is mapped to the same outside IP address in any pool referenced from this inside routing context (referenced through the NAT policy).

After the deterministic DS-Lite is enabled (a **prefix** command under the deterministic CLI node is configured), the ingress hashing influenced by the **dslite-max-subscriber-limit** is in effect for both flavors of DS-Lite (deterministic and non-deterministic) within the inside routing context assuming that both flavors are configured simultaneously.

With introduction of deterministic DS-Lite, the configuration of the subscriber-prefix-length must adhere to the following rule:

The configured value for the subscriber-prefix-length minus the number of bits representing the dslite-max-subscriber-limit value, must be in the range [32 to 64,128]. Or:

```
subscriber-prefix-length - n = [32..64,128]
where n = log2(dslite-max-subscriber-limit)
[or dslite-max-subscriber-limit = 2^n]
```

This can be clarified by the two following examples:

$\text{dslite-max-subscriber-limit} = 64 \text{ — } n=6 \text{ [}\log_2(64) = 6 \text{]} .$

This means that 64 DS-Lite subscribers are mapped to the same outside IP address. Consequently the prefix length of those subscribers must be reduced by 6 bits for hashing purposes (so that chunks of 64 subscribers are always hashed to the same ISA).

According to our rule, the prefix of those subscribers (subscriber-prefix-length) can be only in the range of [38..64], and no longer in the range [32 to 64, 128].

$\text{dslite-max-subscriber-limit} = 1 > n=0 \text{ [}\log_2(1) = 0 \text{]}$

This means that each DS-Lite subscriber is mapped to its own outside IPv4 address. Consequently there is no need for the aggregation of the subscribers for hashing purposes, because each DS-Lite subscriber is mapped to an entire outside IPv4 address (with all ports). Because the subscriber prefix length are not contracted in this case, the prefix length can be configured in the range [32 to 64, 128].

In other words the largest configured prefix length for the deterministic DS-Lite subscriber is $32+n$, where $n = \log_2(\text{dslite-max-subscriber-limit})$. The subscriber prefix length can extend up to 64 bits. Beyond 64 bits for the subscriber prefix length, there is only one value allowed: 128. In the case n must be 0, which means that the mapping between B4 elements (or IPv6 address) and the IPv4 outside addresses is in 1:1 ratio (no sharing of outside IPv4 addresses).

The dependency between the subscriber definition in DS-Lite (based on the subscriber-prefix-length) and the subscriber hashing mechanism on ingress (based on the `dslite-max-subscriber-limit` value), influences the order in which deterministic DS-Lite is configured.

7.9.13.2 Order of configuration steps in deterministic DS-Lite

Configure deterministic DS-Lite in the following order:

1. Configure DS-Lite subscriber-prefix-length.
2. Configure `dslite-max-subscriber-limit`.
3. Configure deterministic prefix (using a NAT policy).
4. Optionally configure map statements under the prefix.
5. Configure DS-Lite AFTR endpoints.
6. Enable (no shutdown) DS-Lite node.

Modifying the `dslite-max-subscriber-limit` requires that all nat-policies be removed from the inside routing context.

To migrate a non-deterministic DS-Lite configuration to a deterministic DS-Lite configuration, the non-deterministic DS-Lite configuration must be first removed from the system. The following steps should be followed:

1. Shutdown DS-Lite node.
2. Remove DS-Lite AFTR endpoints.
3. Remove global NAT policy.
4. Configure/modify DS-Lite subscriber-prefix-length.
5. Configure `dslite-max-subscriber-limit`.
6. Reconfigure global NAT policy.
7. Configure deterministic prefix.
8. Optionally configure one or more manual map statements under the prefix.
9. Reconfigure DS-Lite AFTR endpoints.
10. Enable (no shutdown) DS-Lite node.
11. Configuration Restrictions in Deterministic NAT.
 - **NAT pool**
 - To modify **nat pool** parameters, the **nat pool** must be in a shutdown state.
 - Shutting down the **nat pool** by configuration (**shutdown** command) is not allowed in case that any NAT policy referencing this pool is active. In other words, all configured prefixes referencing the pool via the NAT policy must be deleted system-wide before the pool can be shut down. when the pool is enabled again, all prefixes referencing this pool (with the NAT policy) have to be recreated. For a

large number of prefixes, this can be performed with an offline configuration file executed using the **exec** command.

- **NAT policy**

- All NAT policies (deterministic and non-deterministic) in the same inside routing-instance must point to the same nat-group.
- A NAT policy (be it a global or in a deterministic prefix) must be configured before one can configure an AFTR endpoint.

- **NAT group**

The active-mda-limit in a nat-group cannot be modified as long as a deterministic prefix using that NAT group exists in the configuration (even if that prefix is shutdown). In other words, all deterministic prefixes referencing (with the NAT policy) any pool in that nat-group, must be removed.

- **deterministic mappings (prefix and map statements)**

- Non-deterministic policy must be removed before adding deterministic mappings.
- Modifying, adding or deleting prefix and map statements in deterministic DS-Lite requires that the corresponding nat pool is enabled (in **no-shutdown** state).
- Removing an existing prefix statement requires that the prefix node is in a shutdown state.

```
config>service>vprn>nat>inside>deterministic# info
-----
      classic-lsn-max-subscriber-limit 128
prefix 10.0.5.0/24 subscriber-type classic-lsn-sub nat-policy "det"
      map start 10.0.5.0 end 10.0.5.127 to 192.168.0.7
      map start 10.0.5.128 end 10.0.5.255 to 192.168.0.2
      shutdown

config>service>vprn>nat>inside>deterministic# info
-----
      dslite-max-subscriber-limit 128
prefix 2001:db8:0:1/64 subscriber-type dslite-lsn-sub nat-policy "det"
map start 2001:db8::/64 end 2001:db8::FF:0:0:0/64 to 10.0.0.5
      shutdown

config>service>vprn>nat>inside>ds-lite#
      subscriber-prefix-length 64
      no shutdown
```

Similarly, the map statements can be added or removed only if the prefix node is in a shutdown state.

There are a few rules governing the configuration of the map statement:

- If the number of subscribers per configured prefix is greater than the subscriber-limit per outside IP parameter (2^n), then the lowest n bits of the map start <inside-ip-address> must be set to 0.
- If the number of subscribers per configured prefix is equal or less than the subscriber-limit per outside IP parameter (2^n), then only one **map** command for this prefix is allowed. In this case there is no restriction on the lower n bits of the map start <inside-ip-address>. The range of the inside IP addresses in such map statement represents the prefix itself.

The *outside-ip-address* in the map statements must be unique amongst all map statements referencing the same pool. In other words, two map statements cannot reference the same <outside-ip-address> in a pool.

- **configuration parameters**

- The subscriber-limit in deterministic nat pool must be a power of 2.

- The NAT inside classic-lsn-max-subscriber-limit must be power of 2 and at least as large as the largest subscriber-limit in any deterministic nat pool referenced by this routing instance. To change this parameter, all nat-policies in that inside routing instance must be removed.
- The NAT inside ds-lite-max-subscriber-limit must be power of 2 and at least as large as the largest subscriber-limit in any deterministic nat pool referenced by this routing instance. To change this parameter, all nat-policies in that inside routing instance must be removed.
- In DS-Lite, the [subscriber-prefix-length - log2(dslite-max-subscriber-limit)] value must fall within [32 to 64, 128].
- In DS-Lite, the subscriber-prefix-length can be only modified if the DS-Lite CLI node is in the shutdown state and there are no deterministic DS-Lite prefixes configured.
- **miscellaneous**
 - Deterministic NAT is not supported in combination with 1:1 NAT. Therefore the nat pool cannot be in mode 1:1 when used as deterministic pool. Even if each subscriber is mapped to its own unique outside IP (sub-limit=1, det-port-reservation ports (65535-1023), NAPT (port translation) function is still performed.
 - Wildcard port forwards (including PCP) map to the wildcard port ranges and not the deterministic port range. Consequently logs are generated for static port forwards using PCP.

7.10 Destination Based NAT (DNAT)

Destination NAT (DNAT) in SR OS is supported for LSN44 and L2-Aware NAT. DNAT can be used for traffic steering where the destination IP address of the packet is rewritten. In this fashion traffic can be redirected to an appliance or set of servers that are in control of the operator, without the need for a separate transport service (for example, PBR plus LSP). Applications utilizing traffic steering via DNAT normally require some form of inline traffic processing, such as inline content filtering (parental control, antivirus/spam, firewalling), video caching, and so on.

After the destination IP address of the packet is translated, traffic is naturally routed based on the destination IP address lookup. DNAT translates the destination IP address in the packet while leaving the original destination port untranslated.

Similar to source based NAT (Source Network Address and Port Translation (SNAPT)), the SR OS maintains state of DNAT translations so that the source IP address in the return (downstream) packet is translated back to the original address.

Traffic selection for DNAT processing in MS-ISA is performed via a NAT classifier.

7.10.1 Combination of SNAPT and DNAT

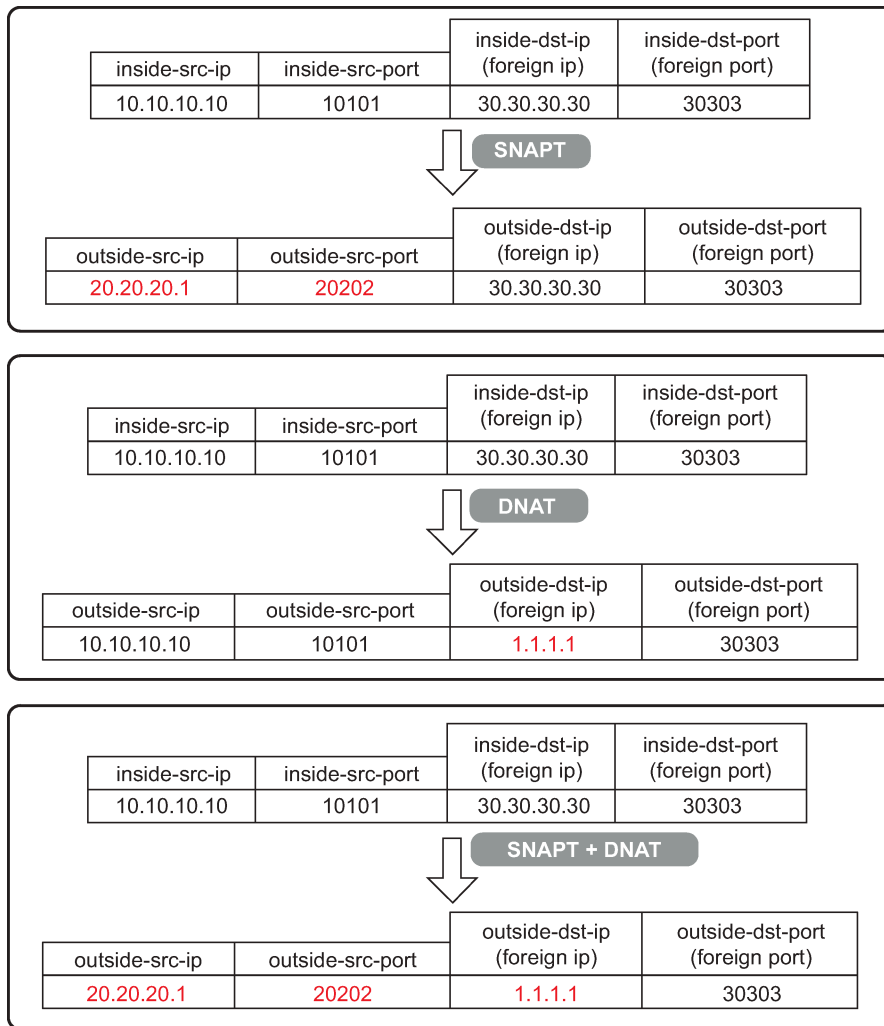
In specific cases SNAPT is required along with DNAT. In other cases only DNAT is required without SNAPT. [Table 42: Supported combinations of SNAPT and DNAT](#) shows the supported combinations of SNAPT and DNAT in SR OS.

Table 42: Supported combinations of SNAPT and DNAT

	SNAPT	DNAT-only	SNAPT + DNAT
LSN44	X	X	X
L2-Aware	X		X

The SNAPT/DNAT address translations are shown in [Figure 72: IP address/port translation modes](#).

Figure 72: IP address/port translation modes



0923

7.10.2 Forwarding model in DNAT

NAT forwarding in SR OS is implemented in two stages:

1. Traffic is first directed toward the MS-ISA. This is performed via a routing lookup, via a filter or via a subscriber-management lookup (L2-Aware NAT). DNAT does not introduce any changes to the steering logic responsible for directing traffic from the I/O line card toward the MS-ISA.
2. When traffic reaches the MS-ISA, translation logic is performed. DNAT functionality incurs an additional lookup in the MS-ISA. This lookup is based on the protocol type and the destination port of the packets, as defined in the nat-classifier.

As part of the NAT state maintenance, the SR OS maintains the following fields for each DNATed flow:

<inside host /port, outside IP/port, *foreign IP address*/port, *destination IP address*/port, protocol (TCP,TCP,ICMP)> Note that the inside host in LSN is inside the IP address and in L2-Aware NAT it is the <inside IP address + subscriber-index>. The subscriber index is carried in session-id of the L2TP.

The *foreign IP address* represents the destination IP address in the original packet, while the *destination IP address* represents the DNAT address (translated destination IP address).

7.10.3 DNAT traffic selection via NAT classifier

Traffic intended for DNAT processing is selected via a nat classifier. The nat classifier has configurable protocol and destination ports. The inclusion of the classifier in the NAT policy is the trigger for performing DNAT. The configuration of the nat classifier determines which of the following is true:

- A specific traffic defined in the match criteria is DNATed while the rest of the traffic is transparently passed through the nat classifier.
- A specific traffic defined in the match criteria is transparently passed through the nat classifier while the rest of the traffic is DNATed.

Classifier cannot drop traffic (no action drop). However, a non-reachable destination IP address in DNAT causes traffic to be black-holed.

7.10.4 Configuring DNAT

DNAT is enabled in the **config>service>nat>nat-policy** context.

```
config>service>nat
nat-policy <nat-policy-name> create
dnat
dnat-only router <router-instance> nat-group <nat-group-id>
nat-classifier <classifier-name>
exit
```

DNAT function is triggered by the presence of the nat classifier (nat-classifier command), referenced in the NAT policy. DNAT-only option is configured in case where SNAPT is not required. This command is necessary to determine the outside routing context and the nat-group when SNAPT is not configured. Pool (relevant to SNAPT) and DNAT-only configuration options within the NAT policy are mutually exclusive.

7.10.4.1 DNAT traffic selection and destination IP address configuration

DNAT traffic selection is performed via a nat-classifier. Nat-classifier is defined under **config>service>nat** hierarchy and is referenced within the **nat-policy**.

```
config>service>nat # nat-classifier <name> create
```

```

default-action {dnat|forward} [ip-addr <ip-address>]
default-dnat-ip-address <ip-addr>
description <description-string>
entry <entry-id> create
action {dnat|forward}[ip-addr <ipv4-address>]
description <description-string>
match protocol {tcp | udp}
match dst-port range start <port-number> end <port-number>
match foreign-ip <ip-address>
exit

```

default-dnat-ip-address is used in all match criteria that contain DNAT action without specific destination IP address. However, the **default-dnat-ip-address** is ignored in cases where IP address is explicitly configured as part of the action within the match criteria.

default-action is applied to all packets that do not satisfy any match criteria.

forward (forwarding action) has no effect on the packets and transparently forwards packets through the nat-classifier.

By default, packets that do not match any matching criteria are transparently passed through the classifier.

7.10.4.2 Micro-netting original source (inside) IP space in DNAT-only case

To forward upstream and downstream traffic for the same NAT binding to the same MS-ISA, the original source IP address space must be known in advance and consequently hashed on the inside ingress toward the MS-ISAs and micro-netted on the outside. This is performed with the following CLI:

```

router | service vprn <id>
nat
inside
classic-lsn-max-subscriber-limit <max>
dnat-only
source-prefix <nat-prefix-list-name>

service nat
nat-prefix-list <name> application dnat-only-subscribers create
prefix <ip-prefix>

```

The **classic-lsn-max-subscriber-limit** parameter was introduced by deterministic NAT and it is reused here. This parameter affects the distribution of the traffic across multiple MS-ISA in the upstream direction traffic. Hashing mechanism based on source IPv4 addresses/prefixes is used to distribute incoming traffic on the inside (private side) across the MS-ISAs. Hashing based on the entire IPv4 address produces the most granular traffic distribution, while hashing based on the IPv4 prefix (determined by prefix length) produces less granular hashing. For further details about this command, consult the CLI command description. The source IP prefix is defined in the nat-prefix-list and then applied under the DNAT-only node in the inside routing context. This instructs the SR OS to create micro-nets in the outside routing context. The number of routes installed in this fashion is limited by the following configuration:

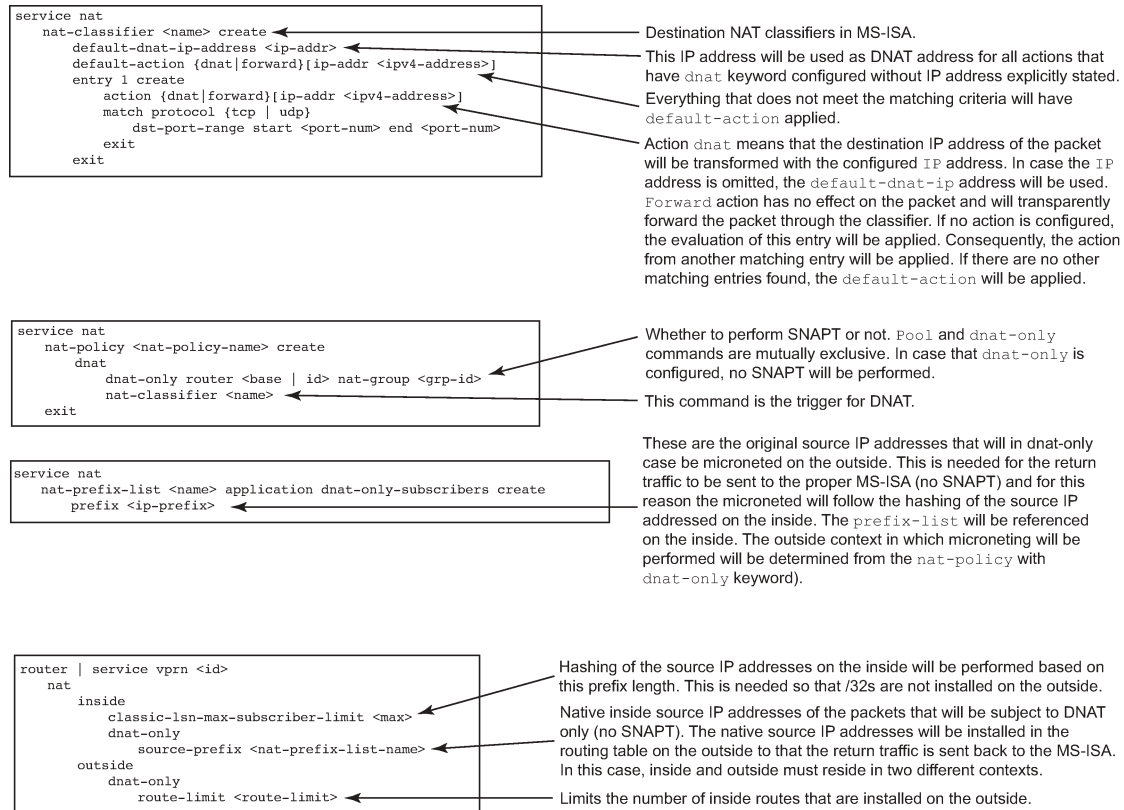
```

router | service vprn <id>
nat
outside
dnat-only
route-limit <route-limit>

```

The configurable range is 1-128K with the default value of 32K. DNAT provisioning concept is shown in [Figure 73: DNAT provisioning model](#).

Figure 73: DNAT provisioning model



0924

7.11 LSN – multiple NAT policies per inside routing context

7.11.1 Restrictions

The following restrictions apply to multiple NAT policies per inside routing context:

- There is no support for L2-Aware NAT.
- DS-Lite and NAT64 diversion to NAT is supported only through IPv6 filters.
- The default NAT policy is counted toward this limit (8).

7.11.2 Multiple NAT policies per inside routing context

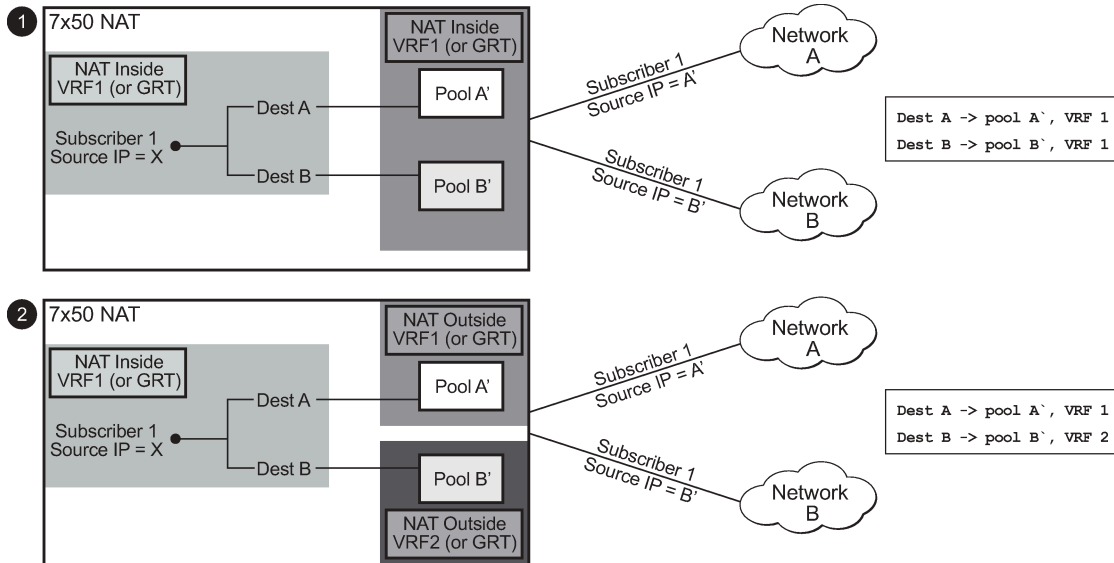
The selection of the NAT pool and the outside routing context is performed through the NAT policy. Multiple NAT policies can be used within an inside routing context. This feature effectively allows selective mapping of the incoming traffic within an inside routing context to different NAT pools (with different mapping properties, such as port-block size, subscriber-limit per pool, address-range, port-forwarding-range,

deterministic vs non-deterministic behavior, port-block watermarks, and so on) and to different outside routing contexts. NAT policies can be configured:

- via filters as part of the **action nat** command
- via routing with the **destination-prefix** command within the inside routing context

The concept of the NAT pool selection mechanism based on the destination of the traffic via routing is shown in [Figure 74: Pool selection based on traffic destination](#).

Figure 74: Pool selection based on traffic destination

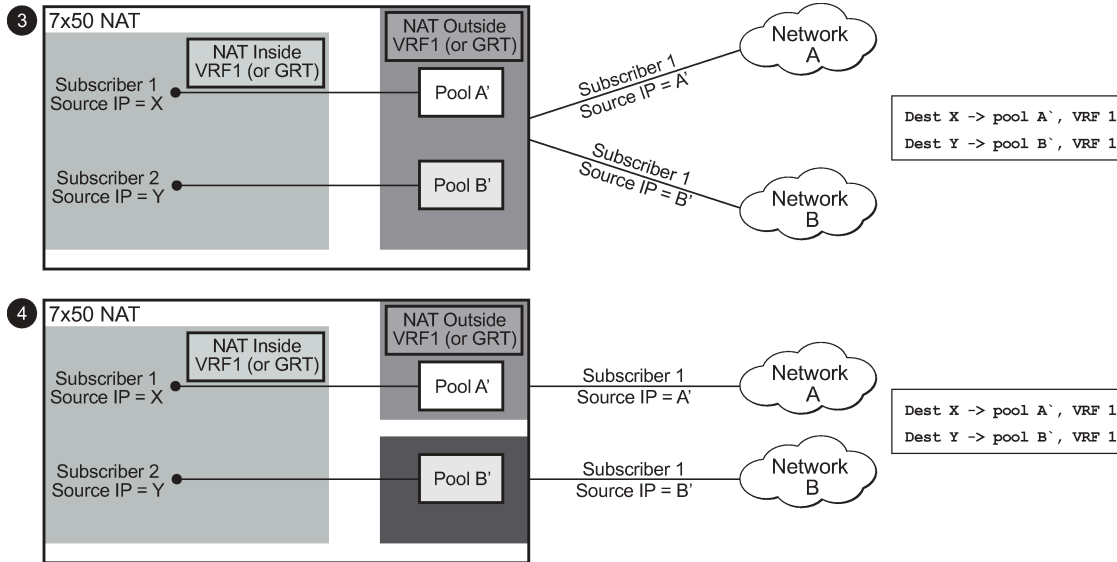


al_0401

Diversion of the traffic to NAT based on the source of the traffic is shown in [Figure 75: NAT pool selection based on the inside source IP address](#).

Only filter-based diversion solution is supported for this case. The filter-based solution can be extended to a 5 tuple matching criteria.

Figure 75: NAT pool selection based on the inside source IP address



al_0402

The following considerations must be taken into account when deploying multiple NAT policies per inside routing context:

- The inside IP address can be mapped into multiple outside IP addresses based on the traffic destination. The relationship between the inside IP and the outside IP is 1:N.
- In case where the source IP address is selected as a matching criteria for a NAT policy (or pool) selection, the inside IP address always stays mapped to the same outside IP address (relationship between the inside IP and outside IP address is, in this case, 1:1)
- Static Port Forwards (SPF); each SPF can be created only in one pool. This means that the pool (or NAT policy) must be an input parameter for SPF creation.

7.11.3 Routing approach for NAT diversion

The routing approach relies on upstream traffic being directed (or diverted) to the NAT function based on the **destination-prefix** command in the **config>service>vprn/router>nat>inside** CLI context. In other words, the upstream traffic is NAT'd only if it matches a preconfigured destination IP prefix. The **destination-prefix** command creates a static route in the routing table of the inside routing context. This static route diverts all traffic with the destination IP address that matches the created entry, toward the MS-ISA. The NAT function itself is performed when the traffic is in the correct context in the MS-ISA.

The CLI for multiple NAT policies per inside routing context with routing based diversion to NAT is the following:

```

service vprn/router
  nat
    inside
      destination-prefix <ip-prefix/length> nat-policy <policy-name>]
      :
      :
  
```

or, for example:

```
service vprn/router
  nat
    inside
      destination-prefix 10.20.10.0/24 nat-policy policy-1
      destination-prefix 10.30.30.0/24 nat-policy policy-1
      destination-prefix 10.40.40.0/24 nat-policy policy-2
```

Different destination prefixes can reference a single NAT policy (policy-1 in this case).

In case that the destination-policy does not directly reference the NAT policy, the default NAT policy is used. The default NAT policy is configured directly in the **vprn/router>nat>inside** context.

After the **destination-prefix** command referencing the NAT policy is configured, an entry in the routing table is created that directs the traffic to the MS-ISA.

7.11.4 Filter-based approach

A filter-based approach diverts traffic to NAT based on the IP matching criteria shown in the CLI below.

```
*A:right-a21>config>filter>ip-filter>entry# match
- match [protocol <protocol-id>]
- no match

<protocol-id>      : protocol numbers - [0..255] (Decimal,
                    Hexadecimal, or Binary representation).
                    Supported IANA IP protocol names -
                    none|crtp|crudp|egp|eigrp|encap|ether-ip|
                    gre|icmp|idrp|igmp|igp|ip|ipv6|ipv6-frag|ipv6-icmp|
                    ipv6-no-nxt|ipv6-opts|ipv6-route|isis|iso-ip|l2tp|
                    ospf-igp|pim|pnni|ptp|rdp|rsvp|sctp|stp|tcp|udp|vrrp
                    * - udp/tcp wildcard

[no] dst-ip        - Configure dest. ip match condition
[no] dst-port     - Configure destination port match condition
[no] port         - Configure port match condition
[no] src-ip       - Configure source ip match condition
[no] src-port     - Configure source port match condition
```

The CLI for the filter-based diversion in conjunction with multiple NAT policies is shown below:

```
filter
  entry
    action nat [nat-policy <nat-policy-name>]
```

The association with the NAT policy is made after the filter is applied to the SAP.

7.11.5 Multiple NAT policies and deterministic NAT

7.11.5.1 Combination of deterministic LSN44, non-deterministic LSN44, and MNP

Deterministic LSN44 is supported in combination with multiple NAT policies (MNP) based on the destination prefix or on a filter term in non-deterministic LSN44. For simplicity, the destination prefix

configuration is used throughout this section instead of the filter terms. However, the combination of deterministic and non-deterministic LSN44 can lead to conflicting scenarios for which the outcomes must be well defined.

The reasons for these conflicting scenarios are the following:

- In private to public (or inside to outside) direction, deterministic NAT uses source IP addresses of the traffic as a match criterion to find the correct NAT pool and outside routing context. This is performed through configuration where each source prefix is associated with a NAT policy.
- In contrast, non-deterministic NAT uses destination IP addresses of the traffic as a match criterion to find the correct NAT pool and outside routing context. This is performed through configuration where each destination prefix is explicitly associated with its own NAT policy.

If both NAT variants are used simultaneously in the same inside NAT routing context, then the conflict resulting from different NAT policies for the same traffic must be resolved. Deterministic and non-deterministic NAT must always use different NAT policies within the same inside routing context.

The rules used to resolve this conflict are:

- Destination prefixes are used to determine the outside routing context from their associated NAT policies. A NAT policy must be explicitly defined for each destination prefix.
- When the outside routing context is determined, the pool selection within that routing context is selected based on the NAT policies associated with the source prefixes.

This means that the outside routing context in both NAT policies must match. If they do not match, traffic is dropped.

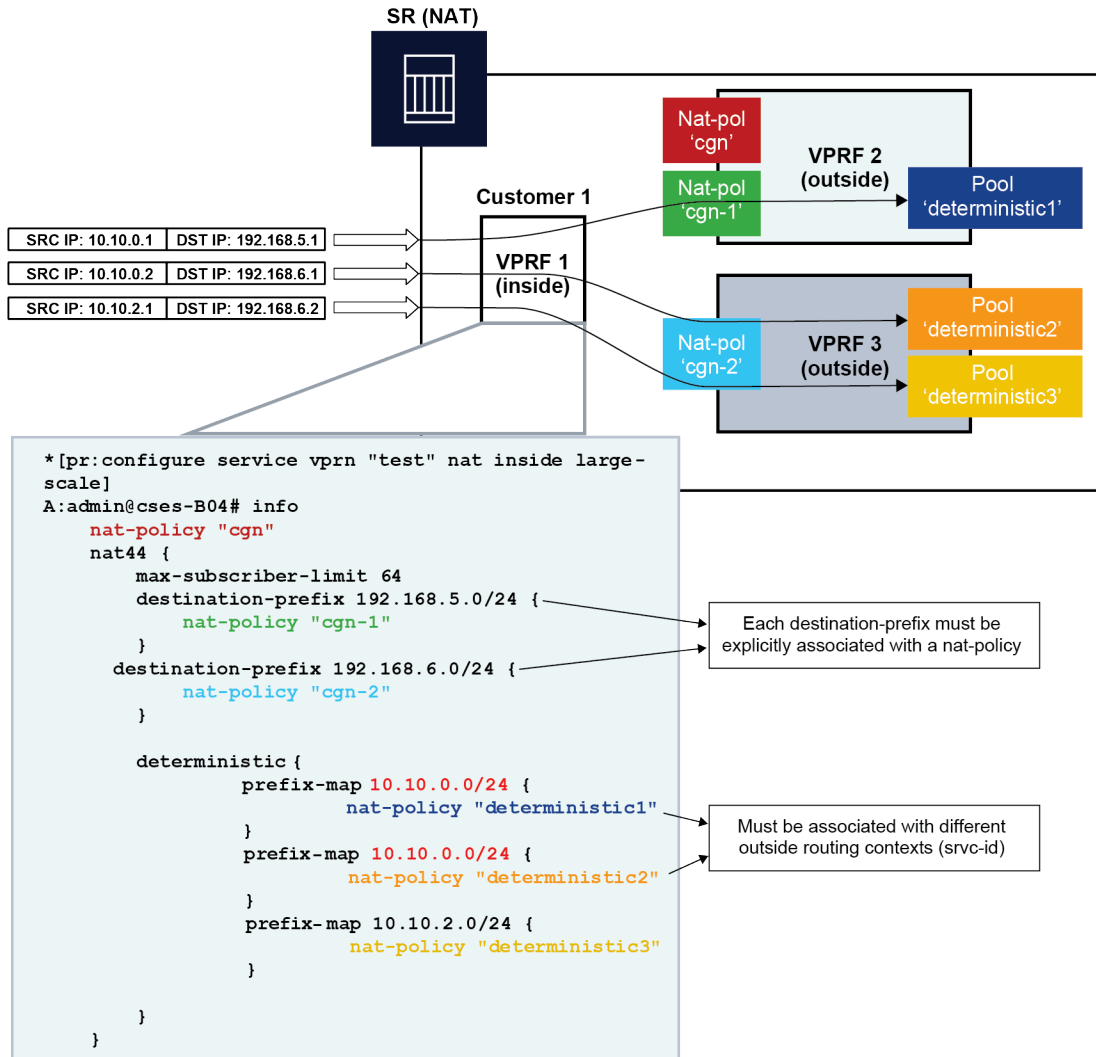
This logic ensures that traffic intended for NAT processing is identified based on the traffic destination (**destination-prefix**) while the pool selection (with outside IP addresses) is determined by the source prefix.

This is useful when non-deterministic NAT and static 1:1 NAT are simultaneously used in the same inside routing context. In this case, a customer to which NAT'd traffic is destined can change and re-purpose its routes that can then be dynamically advertised to the NAT operator. At the same time, the NAT operator who is providing static 1:1 NAT can ensure that its clients are predictably mapped to static outside IP addresses because the source prefixes and their associated NAT policies are not changed, even if the destination prefixes are changed.

An example of a supported scenario is shown in [Figure 76: MNP and deterministic NAT44](#) where all destination prefixes are explicitly associated with NAT policies (no **destination-prefix** is using a default NAT policy) and the source prefixes are explicitly mapped to a different set of NAT policies. The outside routing contexts in the two NAT policies for traffic matching destination and source prefixes simultaneously, must match.

As shown in [Figure 76: MNP and deterministic NAT44](#), traffic from the source network 10.10.0.0/24 is mapped to two different pools in two different outside VPRNs based on the pool names configured in NAT policies associated with the source prefixes. The actual outside VPRN is selected based on the NAT policy associated with the relevant destination prefix.

Figure 76: MNP and deterministic NAT44



sw4102

7.11.6 Multiple NAT policies with DS-Lite and NAT64

DS-Lite and NAT64 diversion to NAT with multiple NAT policies is supported only through IPv6 filters.

Classic CLI

```

configure
  filter
    ipv6-filter
      entry <entry-id> [create]
        action nat nat-type <nat-type> [nat-policy <nat-policy-name>]
        exit
      exit
    exit
  exit

```

Where the **nat-type** parameter can be either **dslite** or **NAT64**.

The DS-Lite AFTR address and NAT64 destination prefix configuration under the corresponding (DS-Lite or NAT64) **router/vprn>nat>inside** context is mandatory. This is even when only filters are needed for traffic diversion to NAT.

For example, every AFTR address and NAT64 prefix that is configured as a match criteria in the filter, must also be duplicated in the **router/vprn>nat>inside** context. However, the opposite is not required.

IPv6 traffic with the destination address outside of the AFTR/NAT64 address/prefix follows normal IPv6 routing path within the 7750 SR.

7.11.7 Default NAT policy

The default **nat-policy** is always mandatory and must be configured under the **router/vprn>nat>inside** context. This default NAT policy can reference any configured pool in the needed ISA group. The pool referenced in the default NAT policy can be then overridden by the NAT policy associated with the destination-prefix in LSN44 or by the NAT policy referenced in the ipv4/ipv6-filter used for NAT diversion in LSN44/DS-Lite/NAT64.

The NAT CLI nodes fail to activate (be brought out of the **no shutdown** state), unless a valid NAT policy is referenced in the **router/vprn>nat>inside** context.

7.11.8 Scaling considerations

Each subscriber using multiple policies is counted as one subscriber for the **inside** resources scaling limits (such as the number of subscribers per MS-ISA), and counted as one subscriber per (subscriber and policy combination) for the **outside** limits (**subscriber-limit** subscribers per IP; **port-reservation** port/block reservations per subscriber).

7.11.9 Multiple NAT policies and SPF configuration considerations

Any Static Port Forward (SPF) can be created only in one pool. This pool, which is referenced through the NAT policy, has to be specified at the SPF creation time, either explicitly through the configuration request or implicitly via defaults.

Explicit requests are submitted either using NSP or the CLI:

```
tools perform nat port-forwarding-action lsn
- lsn create router <router-instance> [b4 <ipv6-address>] [aftr <ipv6-address>] ip <ip-
address> protocol {tcp|udp} [port <port>] lifetime <lifetime> [outside-ip <ipv4-address>]
[outside-port <port>] [nat-policy <nat-policy-name>]
```

In the absence of the NAT policy referenced in the SPF creation request, the default **nat-policy** command under the **vprn/router>nat>inside** context is used.

The consequence of this is that the operator must know the NAT policy in which the SPF is to be created. The SPF cannot be created via PCP outside of the pool referenced by the default NAT policy, because PCP does not provide means to communicate NAT policy name in the SPF creation request.

The static port forward creation and their use by the subscriber types must follow these rules:

- **default NAT policy**

Any subscriber type can use an SPF created in the pool referenced by the default NAT policy.

- **deterministic LSN44 NAT policy**

Only deterministic LSN44 subscribers matching the configured prefix can use the SPF created in the pool referenced by the deterministic LSN44 prefix NAT policy.

- **deterministic DS-Lite NAT policy**

Only deterministic DS-Lite subscribers matching the configured prefix can use the SPF created in the pool referenced by the deterministic DS-Lite prefix NAT policy.

- **LSN44 filter based NAT policy**

Only LSN44 subscribers matching the configured filter entry can use the SPF created in the pool referenced by the non-deterministic LSN44 NAT policy within the filter.

- **DS-Lite filter based NAT policy**

Only DS-Lite subscribers matching the configured filter entry can use the SPF created in the pool referenced by the DS-Lite NAT policy within the filter.

- **NAT64 filter based NAT policy**

Only NAT64 subscribers matching the configured filter entry can use the SPF created in the pool referenced by the NAT64 NAT policy within the filter.

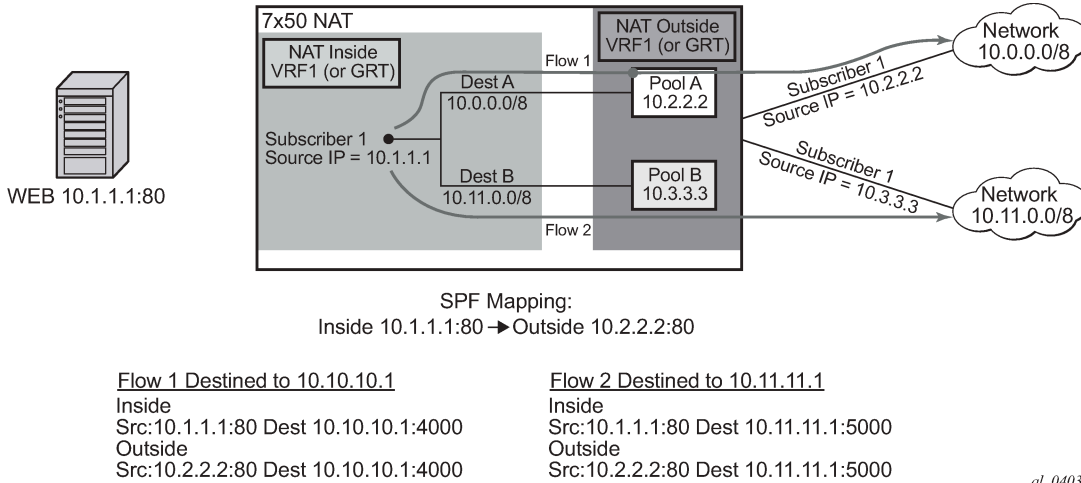
When the last relevant policy for a specific subscriber type is removed from the virtual router, the associated port forwards are automatically deleted.

7.11.9.1 Multiple NAT policies and forwarding considerations

[Figure 77: SPF with multiple NAT policies](#) and [Figure 78: Bypassing NAT policy rule](#) describe specific scenarios that are more theoretical and are less likely to occur in reality. However, they are described here for the purpose of completeness.

[Figure 77: SPF with multiple NAT policies](#) represents the case where traffic from the WEB server 10.1.1.1 is initiated toward the destined network 10.11.0.0/8. Such traffic ends up translated in the Pool B and forwarded to the 10.11.0.0/8 network even though the static port forward has been created in Pool A. In this case, the NAT policy rule (dest 10.11.0.0/8 pool B) determines the pool selection in the upstream direction (even though the SPF for the WEB server already exists in the Pool A).

Figure 77: SPF with multiple NAT policies

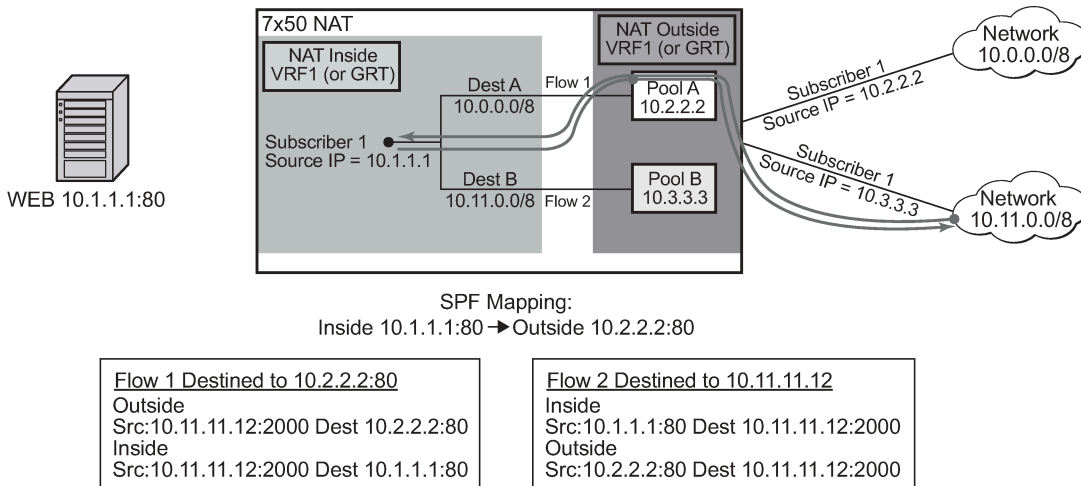


al_0403

The next example in [Figure 78: Bypassing NAT policy rule](#) shows a case where the Flow 1 is initiated from the outside. Because the partial mapping matching this flow already exists (created by SPF) and there is no more specific match (FQF) present, the downstream traffic is mapped according to the SPF (through Pool A to the Web server). At the same time, a more specific entry (FQF) is created (initiated by the very same outside traffic). This FQF now determines the forwarding path for all traffic originating from the inside that is matching this flow. This means that the Flow 2 (reverse of the Flow 1) is not mapped to an IP address from the pool B (as the policy dictates) but instead to the Pool A which has a more specific match.

A more specific match would be in this case fully qualified flows (FQF) that contains information about the foreign host: <host, inside IP/port, outside IP/port, foreign IP address/port, protocol>.

Figure 78: Bypassing NAT policy rule



al_0404

7.12 NAT policy selection in non-deterministic NAT

In deterministic NAT, the NAT policy, and consequently, the NAT pool selection is based on the source IP prefix, as discussed in [Deterministic NAT](#).

The selection of the NAT policy (and then the NAT pool) is based on the source prefix for non-deterministic Large Scale NAT (LSN44). This functionality can be also referred to as NAT traffic identification based on the source prefix.

As described in [Traffic steering to NAT](#), traffic can be redirected to ESA-VM for NAT processing using routing (via the destination prefix) or any criteria defined in the IPv4 filter. The NAT policy selection can be configured in the destination prefix or IPv4 filter. If policy is specified, a default NAT policy configured in the inside routing context is applied to all traffic that does not have an explicitly configured NAT policy in the destination prefix or IPv4 filter.

The source prefix can be used for NAT policy selection as an alternative method to the destination prefix, IPv4 filter, or default NAT policy. After the traffic reaches the ESA-VM, as directed by the destination prefix or IPv4 filter, a source-prefix configuration maps traffic to the NAT policy, and by extension to the NAT pool. Traffic without an explicit mapping between the source prefix and NAT policy is dropped.

NAT policy selection via the source prefix and all other mechanisms for policy selection (such as specific NAT policies in destination-prefix, filter, or a default NAT policy) are mutually exclusive. In other words, traffic steering to ESA-VMs for NAT processing is not performed based on the source prefixes, and instead, relies on existing methods (such as the destination prefix or filter). Source prefixes are only used to associate subscribers with NAT policies, and through the policies with NAT pools.

The following are disabled when the **source-prefix-only** command is enabled:

- default NAT policy
- destination prefixes with an explicit NAT policy
- filter with an explicit NAT policy
- NAT64
- DS-Lite
- destination based NAT
- deterministic DS-lite prefixes
- Stateful Inter-Chassis NAT Redundancy (SICR)

The following example displays a NAT configuration based on the source prefix in non-deterministic NAT.

Example: MD-CLI

```
[ex:/configure router "Base" nat]
A:node-2#
  inside {
    large-scale {
      traffic-identification {
        source-prefix-only true
      }
      nat44 {
        destination-prefix 0.0.0.0/0 {
        }
        source-prefix 10.10.10.0/24 {
          nat-policy "nat-pol-1"
        }
        source-prefix 10.10.11.0/24 {
          nat-policy "nat-pol-1"
        }
      }
    }
  }
```

```
    }
    source-prefix 10.10.12.0/25 {
        nat-policy "nat-pol-2"
    }
    source-prefix 10.10.12.128/25 nat-policy "nat-pol-3" {
        nat-policy "nat-pol-3"
    }
}
}
```

Example: classic CLI

```
A:node-2>config>router>nat>inside
destination-prefix 0.0.0.0/0
traffic-identification source-prefix-only
source-prefix 10.10.10.0/24 nat-policy "nat-pol-1"
source-prefix 10.10.11.0/24 nat-policy "nat-pol-1"
source-prefix 10.10.12.0/25 nat-policy "nat-pol-2"
source-prefix 10.10.12.128/25 nat-policy "nat-pol-3"
```

Configuration notes:

- Aside from the NAT policies linked with the source prefix, other NAT policies are not allowed in this configuration, including the default NAT policy and NAT policies configured under the destination prefix or IPv4 filter.
- The **source-prefix-only** command is mandatory. Only traffic identified with the **source-prefix** command is processed by NAT. Any other traffic that is diverted to NAT and arrives to the ESA-VM is discarded.
- Either the configuration of the **destination-prefix** command or an IPv4 filter is mandatory. This is how traffic is steered toward the ESA-VMs.
- This configuration allows the use of multiple source prefixes.

7.13 Default DMZ Host

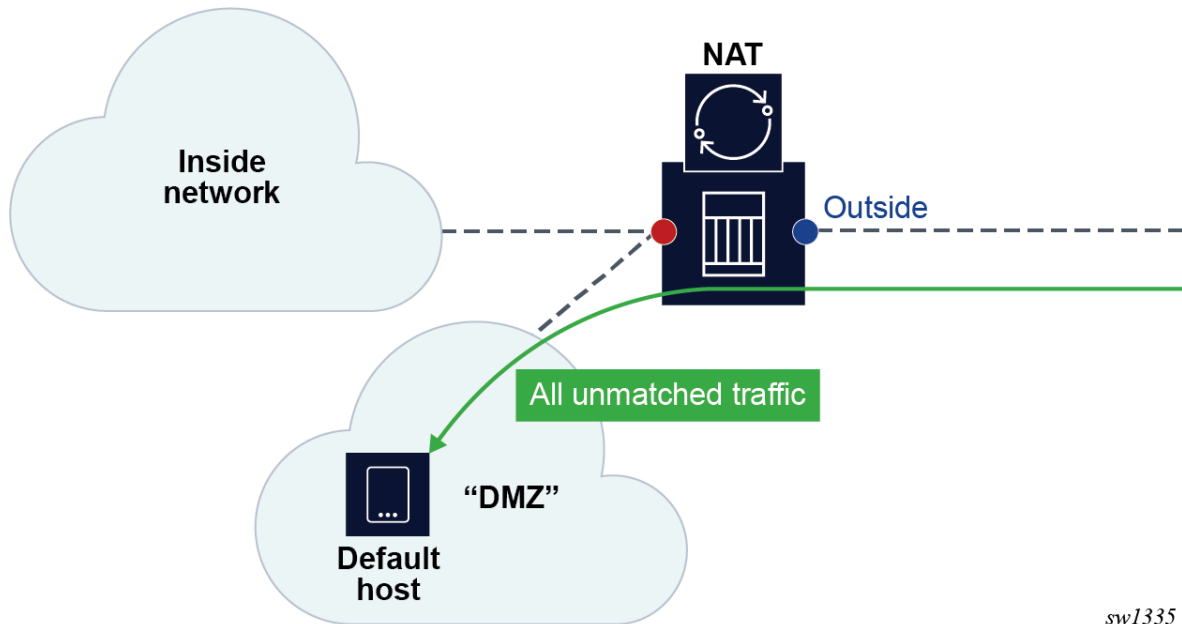
A default demilitarized zone (DMZ) host is a node to which all unmatched traffic from the outside can be redirected. This redirection is achieved by changing the destination IPv4 address in the traffic header to the IPv4 address of the default DMZ host. On the default DMZ host, unmatched traffic can be inspected as part of a threat analysis.

A default DMZ host does not have to be directly connected to the inside NAT segment, but can be located deeper in the network. The default DMZ host does not send any replies to the unknown traffic; and therefore, there is no state maintained for the unknown traffic in NAT. The rate of unmatched traffic sent to the default DMZ host can be restricted by configuration.

In the redirected traffic with swapped destination IPv4 addresses, the Layer 3 and Layer 4 (UDP and TCP) checksums are recalculated.

The following example show a basic default DMZ host configuration.

Figure 79: Default DMZ host



The following example shows a default DMZ host for LSN in a VPRN service configuration.

Example: MD-CLI

```
[ex:/configure service vprn "demo-vprn" nat outside pool "demo-pool" large-scale default-host
configure service vprn "demo-vprn" nat outside pool "demo-pool" large-scale
  default-host {
    ip-address 10.10.10.10
  }
```

Example: classic CLI

```
*A:node-2>config>router#
  nat
    outside
      pool "demo-pool" nat-group 1 type large-scale create
      default-host 10.10.10.10 inside-router-id "Base" rate 10
    exit
  exit
exit
no shutdown
-----
```

The following example shows a L2-aware NAT in a "Base" router configuration.

Example: MD-CLI

```
[ex:/configure router "Base" nat outside pool "demo-pool" l2-aware]
configure router "Base" nat outside pool "demo-pool" l2-aware
  default-host {
    ip-address 10.10.10.10
    inside-router-instance "Base"
    rate-limit 100
```

```
}

```

Example: classic CLI

```
*A:node-2>config#
  router
    nat
      outside
        pool "demo-pool" nat-group 1 type l2-aware create
        default-host 10.10.10.10 inside-router-id "Base" rate 100
      exit
-----
```

7.14 NAT and CoA

RADIUS Change of Authorization (CoA) can be used in subscriber management (ESM) to modify the NAT behavior of the subscriber. This can be performed by:

- replacing a NAT policy in a subscriber profile for the L2-Aware NAT subscriber
- replacing or removing a NAT policy within the IP filter for the ESM subscriber using LSN44, DS-Lite or NAT64
- modifying DNAT parameters directly via CoA for the L2-Aware subscriber

7.14.1 CoA and NAT policies

The behavior for NAT policy changes via CoA for LSN and L2-Aware NAT is summarized in [Table 43: NAT policy changes via CoA](#).

Table 43: NAT policy changes via CoA

Action	Outcome		Remarks
	L2-Aware	LSN	
CoA - replacing NAT policy	<p>Stale flows using the old NAT policy are cleared after 5 seconds.</p> <p>New flows immediately start using a new NAT policy.</p> <p>Restrictions:</p> <p>Allowed only when the previous change is completed (need to wait for a 5 second interval during which the stale mappings caused</p>	<p>Stale flows using the old NAT policy continue to exist and are used for traffic forwarding until they are naturally timed-out or TCP- terminated. The exception to this is when the reference to the NAT policy in the filter was the last one for the inside VRF. In this scenario, the flows from the removed NAT policy are cleared immediately.</p>	<p>A NAT policy change via CoA is performed by changing the sub-profile for the ESM subscriber or by changing the ESM subscriber filter in the LSN case. 1</p> <p>A sub-profile change alone does not trigger accounting messages in L2-Aware NAT and consequently the logging information is lost.</p> <p>To ensure timely RADIUS logging of the NAT policy change in L2-Aware NAT, each CoA must, in addition to the sub-profile change, also do one of the following:</p> <ul style="list-style-type: none"> • Change the sla-profile2.

Action	Outcome		Remarks
	L2-Aware	LSN	
	by previous CoA are purged).	New flows immediately start using new NAT policy.	<ul style="list-style-type: none"> Include the Alc-Trigger-Acct-Interim VSA in the CoA messages. Both of the above events trigger an accounting update at the time when CoA is processed. This keeps NAT logging current.
	(cont.) <ul style="list-style-type: none"> Not allowed if L2-Aware subscriber has multiple hosts and the new prefix-list contains one or more 1:1 NAT policies. Not allowed if the new NAT policy references to a pool in a different NAT group. 		

In non-ESM environments, the NAT policy can be changed by replacing the interface filter via CLI for LSN case.

The SLA profile has to be changed and not just *refreshed*. In other words, replacing the existing SLA profile with the same one does not trigger a new accounting message.

7.14.2 CoA and DNAT

Adding, removing or replacing DNAT parameters in LSN44 can be achieved through NAT policy manipulation in an IP filter for ESM subscriber. The rules for NAT policy manipulation via CoA are given in [Table 43: NAT policy changes via CoA](#). In L2-Aware NAT, CoA can be used to:

- Enable or disable DNAT functionality while leaving the Source Network Address and Port Translation (SNAPT) uninterrupted.
- Modify the default destination IP address in DNAT.

After the DNAT configuration is modified via CoA (enable or disable DNAT or change the default DNAT IP address), the existing flows affected by the change remain active for 5 more seconds while the new flows are created in accordance with the new configuration. After a 5 second timeout, the stale flows are cleared from the system.

The RADIUS attribute used to perform DNAT modifications is a composite attribute with the following format:

Alc-DNAT-Override (234) = "{<DNAT_state> | <DNAT-ip-addr>},[nat-policy]"

where: *DNAT state* = none | disable → and the **DNAT-ip-addr** parameter are mutually exclusive.

DNAT-ip-addr = Provides an implicit enable with the destination IPv4 address in dotted format (a.b.c.d) → and the DNAT-state parameter are mutually exclusive.

nat-policy = nat-policy *name* → This is an optional parameter. If it is not present, then the default NAT policy is assumed.

For example:

Alc-DNAT-Override=none → This negates any previous DNAT related override in the **default** nat-policy. Consequently, the DNAT functionality is set as originally defined in the **default** nat-policy. In case that the 'none' value is received while DNAT is already enabled, a CoA ACK is sent back to the originator.

Alc-DNAT-Override =none,nat-pol-1 → This re-enables DNAT functionality in the specific NAT policy with the name *nat-policy-1*.

Alc-DNAT-Override =none,10.1.1.1 → The DNAT-state and DNAT-ip-addr parameters are mutually exclusive within the same Alc-DNAT-Override attribute. Although a CoA ACK reply is returned to the RADIUS server, an error log message is generated in the SR OS indicating that the attempted override failed.

Alc-DNAT-Override =10.1.1.1 → This changes the default DNAT IP address to 1.1.1.1 in the default NAT policy. In case DNAT was disabled before receiving this CoA, it is implicitly enabled.

Alc-DNAT-Override =10.1.1.1,nat-pol-1 → This changes the default DNAT IP address to 10.1.1.1 in the specific NAT policy named *nat-policy-1*. DNAT is implicitly enabled if it was disabled before receiving this CoA.

The combination of sub-fields with the Alc-DNAT-Override RADIUS attribute and the corresponding actions are shown in [Table 44: CoA and DNAT](#) .

Table 44: CoA and DNAT

DNAT-state	DNAT-ip-addr	NAT policy	DNAT action in L2-Aware NAT
none	-	-	<p>Re-enable DNAT in the default NAT policy.</p> <p>If DNAT was enabled before receiving this CoA, then no specific action is carried out by the SR OS with the exception of sending the CoA ACK back to the CoA server.</p> <p>This negates any previous DNAT-related override in the default nat-policy. Consequently, the DNAT functionality is set as originally defined in the default nat-policy.</p> <p>If the DNAT classifier is not present in the default nat-policy when this CoA is received, an error log message is raised.</p>
none	-	nat-pol-name	<p>Re-enable DNAT in the referenced NAT policy.</p> <p>This negates any previous DNAT related override in the referenced nat-policy. Consequently, the DNAT functionality is set as originally defined in the referenced nat-policy.</p> <p>If the DNAT classifier is not present in the referenced nat-policy when this CoA is received, a CoA ACK reply is returned to the RADIUS server and an error log message is generated in the SR OS indicating that the attempted override has failed.</p>

DNAT-state	DNAT-ip-addr	NAT policy	DNAT action in L2-Aware NAT
none	a.b.c.d	-	<p>These two parameters are mutually exclusive in the same Alc-DNAT-Override attribute.</p> <p>Although a CoA ACK reply is returned to the RADIUS server, an error log message is generated in SR OS indicating that the attempted override has failed.</p>
none	a.b.c.d	nat-pol-name	<p>DNAT-state and DNAT-ip-address parameters are mutually exclusive in the same Alc-DNAT-Override attribute.</p> <p>Although a CoA ACK reply is returned to the RADIUS server, an error log message is generated in SR OS indicating that the attempted override has failed.</p>
disable	-	-	<p>Disable DNAT in the default NAT policy.</p> <p>If the DNAT classifier is not present in the default nat-policy when this CoA is received, a CoA ACK reply is returned to the RADIUS server and an error log message is generated in the SR OS indicating that the attempted override has failed.</p>
disable	-	nat-pol-name	<p>Disable DNAT in the referenced NAT policy.</p> <p>If the DNAT classifier is not present in the referenced nat-policy when this CoA is received, a CoA ACK reply is returned to the RADIUS server and an error log message is generated in SR OS indicating that the attempted override has failed.</p>
disable	a.b.c.d	-	<p>The DNAT-state and DNAT-ip-address parameters are mutually exclusive in the same Alc-DNAT-Override attribute.</p> <p>Although a CoA ACK reply is returned to the RADIUS server, an error log message is generated in SR OS indicating that the attempted override has failed.</p>
disable	a.b.c.d	nat-pol-name	<p>The DNAT-state and DNAT-ip-address parameters are mutually exclusive in the same Alc-DNAT-Override attribute.</p> <p>Although a CoA ACK reply is returned to the RADIUS server, an error log message is generated in the SR OS indicating that the attempted override has failed.</p>
-	a.b.c.d	-	<p>The default destination IP address is changed in the default NAT policy.</p>
-	a.b.c.d	nat-pol-name	<p>The default destination IP address is changed in the referenced NAT policy.</p>

DNAT-state	DNAT-ip-addr	NAT policy	DNAT action in L2-Aware NAT
-	-	- or nat-pol-name	A CoA NAK (error) is generated. Either DNAT-state or DNAT-ip-address parameters must be present in the Alc-DNAT-Override attribute.

If multiple Alc-DNAT-Override attributes with conflicting actions are received in the same CoA or Access-Accept, the action that occurred last takes precedence.

For example, if the following two Alc-DNAT-Override attributes are received in the same CoA, the last one takes effect and consequently DNAT is disabled in the default NAT policy:

Alc-DNAT-Override = "10.1.1.1"

Alc-DNAT-Override = "disable"

7.14.3 Modifying an active NAT prefix list or NAT classifier via CLI

[Table 45: Modifying active NAT prefix list or NAT classifier](#) describes the outcome when the active NAT prefix list or NAT classifier is modified using CLI.

Table 45: Modifying active NAT prefix list or NAT classifier

Action	Outcome		Remarks
	L2-Aware	LSN	
CLI – Modifying prefix in the NAT prefix list	Existing flows are always checked whether they comply with the NAT prefix list that is currently applied in the subscriber profile for the subscriber. If the flows do not comply with the current NAT prefix list, they are cleared after 5 seconds. The new flows immediately start using the updated settings.	Changing the prefix in the NAT prefix list internally re-subnets the outside IP address space.	A NAT prefix list is used with multiple NAT policies in L2-Aware NAT and for downstream Internal subnet in DNAT-only scenario for LSN. The prefix can be modified (added, removed, remapped) at any time in the NAT prefix list, while the NAT policy must be first shut down via CLI.
CLI – Modifying or replacing the NAT classifier	Existing flows are always checked whether they comply with the NAT classifier that is currently applied in the active NAT policy for the subscriber. If the flows do not comply with the current NAT classifier, they are cleared after 5 seconds.	Changing the NAT classifier have the same effect as in L2-Aware NAT; all existing flows using the NAT classifier are checked whether they comply with this classifier or not.	The NAT classifier is used for DNAT. NAT classifier is referenced in the NAT policy.

Action	Outcome		Remarks
	L2-Aware	LSN	
	The new flows immediately start using the updated settings.		
CLI - Removing or adding NAT policy in NAT prefix list	Blocked	Not applicable	
CLI - Removing or adding NAT policy in the subscriber profile	Blocked	Not applicable	
CLI - Removing, adding or replacing NAT prefix list under the rtr/nat/inside/DNAT-only	Not applicable	This action triggers internally re-subnet the source address space according to the new NAT prefix list. However, the current flows in the MS ISA are not affected by this change. In other words, they are not removed if the associated prefix is removed from the prefix list.	

7.15 Watermarks

Watermarks can be configured to monitor the actual usage of sessions, ports, and port blocks.

For each watermark, a high and a low value must be set. When the high threshold value is crossed in the upward direction, an event is generated (SNMP trap), notifying the operator that a NAT resource may be approaching exhaustion. When the low threshold value is crossed in the downward direction, a similar event is generated (clearing the first event), notifying the operator that the resource utilization has dropped below the low threshold value.

Watermarks can be defined on the NAT group, pool, and policy level.

- **NAT group**

Watermarks can be placed to monitor the total number of sessions on an MDA.

- **NAT pool on each NAT group member**

Watermarks can be placed to monitor the port and port-block occupancy in a pool within each NAT group member.

- **NAT policy**

In the policy, the operator can define watermarks on session and port usage. In both cases, the usage per subscriber (for L2-Aware NAT) or per host (for large-scale NAT) is monitored.

7.16 Port forwards

Port forwards allow devices on the public side of NAT (NAT outside) to initiate sessions toward those devices, usually servers, that are hidden behind NAT (NAT inside). Another term for port forwards is NAT pinhole.

A port forward represents a previously created (before any traffic is received from the inside) mapping between a TCP/UDP port on the outside IP address and a TCP/UDP port on the inside IP address assigned to a device behind the NAT. This mapping can be created statically by configuration (such as CLI, MIB, YANG, or NETCONF), or it can be created dynamically with protocols such as PCP or UPnP. Port forwards are supported only in NAT pools in Network Address and Port Translation (NAPT) mode. NAT pools in 1:1 mode do not support configured port forwards because, by default, the pools allow traffic from the outside to the inside and this cannot be disabled. Pools in 1:1 mode (whether protocol agnostic) do not perform port translation; therefore the inside and outside always match.

In UPnP, the forwarded ports are created with the port range of the NAT subscriber, whereas with PCP and Static Port Forwards (SPF), the forwarded ports are allocated from a dedicated port range outside of the port blocks allocated to individual NAT subscribers. There are two ranges dedicated to port forwards in NAT:

- **well-known ports (1 to 1023)**

This range is always enabled and cannot be disabled in NAT pools that support configured port forwards (non 1:1 NAT pools).

- **ports from the ephemeral port range (1025 to 65535)**

Port forwards from the ephemeral port space must be explicitly enabled by configuration. They are allocated from a contiguous block of ports where upper and lower limits are defined. Ports reserved for port forwards allocated in the ephemeral port space are also referred to as wildcard ports.

Port forwarding ranges (well-known ports and wildcard ports) are shared by all NAT subscribers on a specific outside IP address. Port blocks that are individually assigned to the subscriber cannot be allocated from the port forwarding range. The wildcard port forwarding range can be configured only when the pool is administratively disabled.

See the [Port Control Protocol \(PCP\)](#) and [Universal plug and play Internet Gateway Device service](#) sections as well as the SR OS R23.x.Rx Software Release Notes for more information about these protocols and the supported NAT types.

7.16.1 Static port forwards

In the MD-CLI and NETCONF, you must use a **tools** command to manage NAT Static Port Forwards (SPFs).

In the classic CLI, you can manage SPFs through a **tools** or configuration command.

If the **tools** command is configured to manage SPFs and preserve SPFs across reboots, you must use the following command to enable persistency of the SPF. With persistency enabled, SPF configuration is stored on the compact flash. The same command is used for both classic CLI and MD-CLI.

```
configure system persistence nat-port-forwarding
```

Execute the following command to manage port forwarding for Large Scale NAT (LSN):

- **MD-CLI**

```
tools perform nat port-forwarding-action lsn ?
lsn add router <string> [b4 <ipv6 address>] [aftr <ipv6 address>] ip <IP
address> protocol <keyword> [port <number>] [total-port <number>] lifetime
<string> [outside-ip <ipv4 address>] [outside-port <number>] [nat-policy
<string>] [force]
lsn remove router <string> [b4 <ipv6 address>] ip <IP address> protocol
<keyword> port <number> [nat-policy <string>]
lsn modify router <string> [b4 <ipv6 address>] ip <IP address> protocol
<keyword> port <number> lifetime <string> [nat-policy <string>]
```

- **classic CLI**

```
configure service nat port-forwarding lsn
- lsn router <router-instance> [b4 <ipv6-address>] [aftr <ipv6-address>] ip
<ip-address> protocol {tcp|udp} [port <port>] [total-port <total-port>]
[outside-ip <ipv4-address>] [outside-port <port>] [nat-policy
<nat-policy-name>] [force]
```

```
*A:node-2>config>service>nat>fwd# info detail
-----
      lsn router 101 ip 11.11.13.7 protocol udp port 12345 outside-ip
130.0.255.254 outside-port 3171 nat-policy "poll_for_2001-pool-0"
-----
```

For the **tools perform nat port-forwarding-action** command, if you do not explicitly configure the following optional fields, the system selects them automatically:

- port number – number of the source port
- outside IP – IPv4 address for the outside IP address

- outside-port number – number of the outside port
- NAT policy – name of the NAT policy

You can specify a **force** option that is applicable only to LSN pools with flexible port allocations where the dynamic ports in this pool are allocated individually instead of port blocks. The dynamic ports are interleaved with Static Port Forwards (SPFs). This creates increased possibility for a collision between the dynamically-allocated port and the requested SPF during an SPF request.

For instance, if a user requests port X on a public IP address Y, there is a chance that port X is already in use because of the dynamic allocation.

To resolve such conflicts, use the **force** option to ensure that the requested SPF has higher priority, allowing it to preempt an existing dynamically-allocated port. This action overwrites the previous port mapping and deletes all associated sessions.

If you omit the **force** option in such a scenario, the static-port allocation fails. The **force** option can only preempt dynamically-allocated ports and does not affect pre-existing SPFs.

7.16.2 Port Control Protocol (PCP)

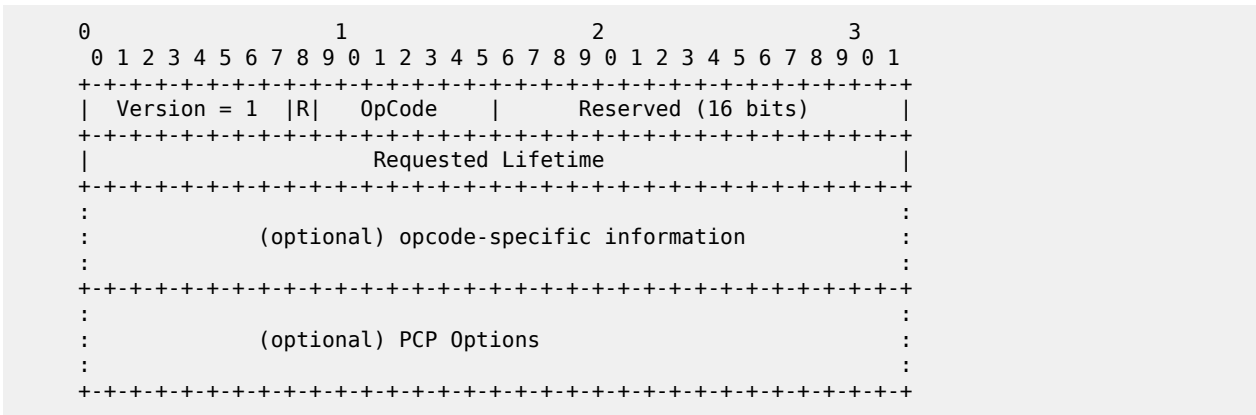
PCP is a protocol that operates between subscribers and the NAT directly. This makes the protocol similar to DHCP or PPP in that the subscriber has a limited but direct control over the NAT behavior.

PCP is designed to allow the configuration of static port-forwards, obtain information about existing port forwards and to obtain the outside IP address from software running in the home network or on the CPE.

PCP runs on each MS-ISA as its own process and make use of the same source-IP hash algorithm as the NAT mappings themselves. The protocol itself is UDP based and is request/response in nature, in some ways, similar to UPnP.

PCP operates on a specified loopback interface in a similar way to the local DHCP server. It operates on UDP and a specified (in CLI) port. As Epoch is used to help recover mappings, a unique PCP service must be configured for each NAT group.

When epoch is lowered, there is no mechanism to inform the clients to refresh their mappings en masse. External synchronization of mappings is possible between two chassis (epoch does not need to be synchronized). If epoch is unsynchronized then the result is clients re-creating their mapping on next communication with the PCP server.



The R-bit (0) indicates request and (1) indicates response. This is a request so (0).

OpCode defined as:

Requested Lifetime: Lifetime 0 means delete.

```

    0           1           2           3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-----+-----+-----+-----+-----+-----+-----+-----+
| Version = 1 |R|  OpCode   |  Reserved   |  Result Code |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                                     Lifetime                                     |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                                     Epoch                                       |
+-----+-----+-----+-----+-----+-----+-----+-----+
:                                                                                   :
:           (optional) OpCode-specific response data                             :
:                                                                                   :
+-----+-----+-----+-----+-----+-----+-----+-----+
:           (optional) Options                                                    :
+-----+-----+-----+-----+-----+-----+-----+-----+

```

As this is a response, R = (1).

The Epoch field increments by 1 every second and can be used by the client to determine if state needs to be restored. On any failure of the PCP server or the NAT to which it is associated Epoch must restart from zero (0).

Result Codes:

- 0 SUCCESS, success.
- 1 UNSUPP_VERSION, unsupported version.
- 2 MALFORMED_REQUEST, a general catch-all error.
- 3 UNSUPP_OPCODE, unsupported OpCode.
- 4 UNSUPP_OPTION, unsupported option. Only if the Option was mandatory.
- 5 MALFORMED_OPTION, malformed option.
- 6 UNSPECIFIED_ERROR, server encountered an error
- 7 MISORDERED_OPTIONS, options not in correct order

Creating a Mapping

Client Sends

```

    0           1           2           3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-----+-----+-----+-----+-----+-----+-----+-----+
| Protocol   |                                     Reserved (24 bits) |
+-----+-----+-----+-----+-----+-----+-----+-----+
|           Internal port           | Suggested external port |
+-----+-----+-----+-----+-----+-----+-----+-----+
:                                                                                   :
: Suggested External IP Address (32 or 128, depending on OpCode):             :
:                                                                                   :
+-----+-----+-----+-----+-----+-----+-----+-----+

```

MAP4 opcode is (1). Protocols: 0 – all; 1 – ICMP; 6 – TCP; 17 – UDP.

MAP4 (1), PEER4 (3) and PREFER_FAILURE are supported. FILTER and THIRD_PARTY are not supported.

7.16.3 PORT_SET option

7.16.3.1 Terminology

The terms internal port and inside port are used interchangeably. They both refer to the original source port before NAT is performed.

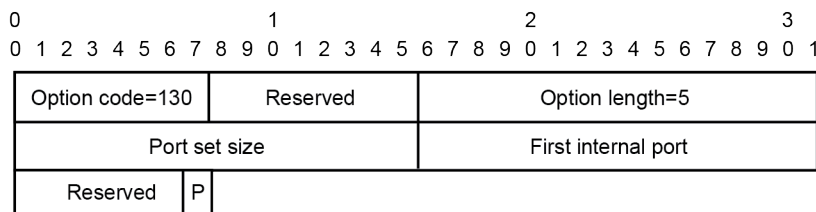
The terms external port and outside port are used interchangeably. They both refer to the translated source port after NAT is performed.

The PCP PORT_SET option is defined in RFC 7753, *Port Control Protocol (PCP) Extension for Port-Set Allocation*, and is used by applications that require a consecutive block of port forwards. The reasons to provide a block of ports in a single request as opposed to multiple requests for single port in a MAP request are described in Section 2 of RFC 7753.

A PCP PORT_SET option indicates to the PCP server (SR OS) that a client needs a block of sequential port forwards. The number of requested port forwards must be greater than one (otherwise, a plain MAP opcode can be used). The ports in the block start at the Internal Port (in MAP opcode) and map to the same number of ports on the outside (or the external side), starting from either the Suggested External Port (in MAP opcode) or from an arbitrary external port. The returned number of ports may be fewer than the requested number, but the number cannot be larger. The allocated port set cannot fall outside of the range defined by the Internal Port (as requested in MAP opcode) plus the PORT_SET Size. Before a port set is assigned to a client, the SR OS always checks if any of the internal ports in the MAP request carrying the PORT_SET option has already been allocated.

The following figure shows an example of PORT_SET option format.

Figure 80: PORT_SET option format example



sw1320

In the MAP request, the needed number of ports is specified in the PORT_SET Size field.

The First Internal Port is set to the same value as the Internal Port field in the MAP opcode.

In the MAP response, the PORT_SET Size represents the number of allocated ports. The First Internal Port represents the first internal port in the port set, which may be different than the Internal Port field in the MAP opcode (see example in Section 6.3 in RFC 7753, overlapping requests). The first external port is the value returned in Assigned External Port field in the MAP opcode.

7.16.3.2 Enabling the PORT_SET option

The PORT_SET option is enabled with the following CLI:

Classic CLI:

```
config>service>nat>pcp-server-policy#
option
```

```
[no] port-set
```

MD-CLI:

```
configure
  service {
    nat {
      pcp-server-policy <name> {
        option {
          port-set <boolean>
        }
      }
    }
  }
}
```

7.16.3.3 Port allocation scheme

The appropriate port set size is determined from the combination of the port set size received via the PORT_SET option and the locally configured policy which may limit the port set size.

If the requested port set is initially not available at the Suggested External Port, an attempt is made to find a new set of ports of the appropriate size. The appropriate PORT_SET size in this context means a combination of the requested PORT_SET size and the local limits set by the operator in the SR OS node.

If the available number of consecutive ports in a set for the specified external IP address is fewer than the requested amount or stated in the policy, the maximum number of available consecutive ports are allocated to the client.

In summary, the SR OS tries to find the biggest available port set (as dictated by the combination of the requested size and local policy), instead of allocating random port sets available at the suggested external port.

7.16.3.4 Limits and quotas

The maximum number of port forwards per subscriber can be limited with the following CLI:

Classic CLI:

```
config>service>nat>nat-policy# port-limits forwarding
- forwarding <limit>
- no forwarding
```

MD-CLI:

```
configure
  service {
    nat {
      nat-policy <name> {
        port-limits {
        }
      }
    }
  }
}
```

This limit is the total number of port forwards, regardless of the methods by which the port forwards are requested, either PCP MAP, PORT_SET, or static port forwards.

7.16.3.5 Port overlaps

PCP clients should not request overlapping ports. Requesting overlapping ports produces an erroneous condition. If this condition occurs, then the request is considered as a refresh of the existing ports. This is described in RFC 7753, Section 4.4.1.

7.16.3.6 Port allocation example

The following example depicts PCP port allocation with the PORT_SET option in action. For additional examples, see RFC 7753.

In this example, the sequence on the top of [Figure 81: PORT_SET example](#) represents the state of the port forwards for an external IP address before the PCP Request with the PORT_SET option is received. The port 1032 (in red) has already been mapped to a source port for the same client or to a different client.

A PCP client requesting an overlapping set of external ports (while internal ports are different) triggers the following action in SR OS (PCP server):

- The SR OS checks if the existing mapping is overlapping with the one for this client. It checks to see if the occupied external port is already mapped to one of the requested internal ports from the 20000 to 20009 range for the same client.
- If such overlap between internal and external ports is detected (for example, 20001 is mapped to 1032), then this is considered a refresh, and the response is:

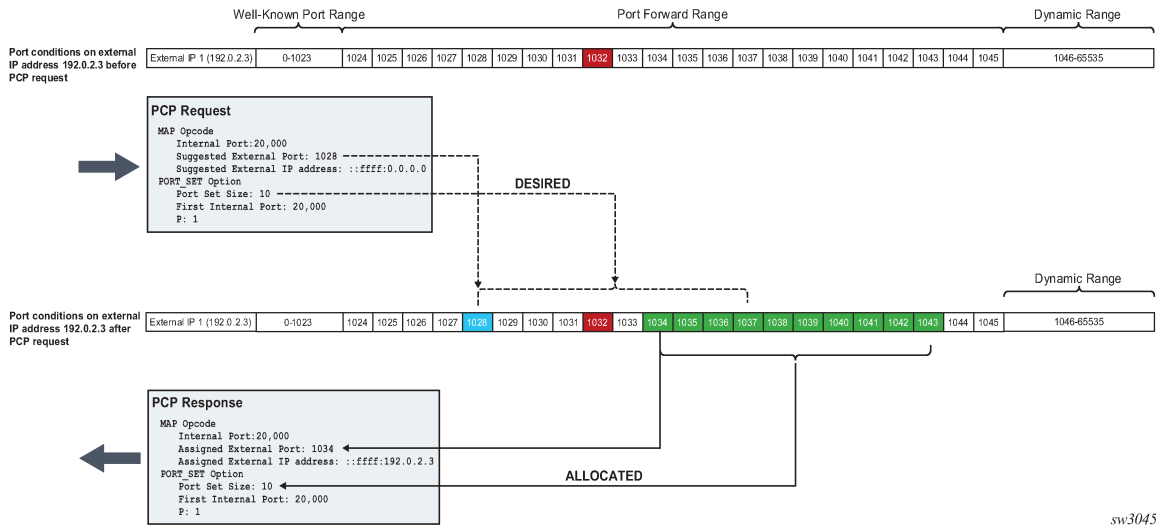
```
opcode: 1(MAP)
internal port: 20,000
assigned external port: 1032
assigned external ip address:192.0.2.3
```

Every overlapping pair of internal and external ports are individually acknowledged. No new mapping is allocated.

- If there is no overlap (when the external port is mapped to a port outside of the 20000 to 20009 range), and no limit is reached, then the SR OS honors the request where it finds any set containing consecutive 10 ports.

In [Figure 81: PORT_SET example](#), there is no overlap (external port 1032 belongs to another client or the same client outside of the 20000 to 20009 range) and consequently a block of 10 ports (following parity) is allocated to the client.

Figure 81: PORT_SET example



7.16.3.7 Operational considerations

Consider the following points when the when PCP PORT_SET option is enabled. Points listed below are in accordance with the RFC 7753:

- The SR OS attempts to allocate the first external port from the suggested external port. If the port is not available, another port is selected.
- Parity is honored.
- The PORT_SET size value 0xffff in the request indicates that the client is willing to accept as many ports as the SR OS can offer.
- If the requested PORT-SET request cannot be fully (or exactly) satisfied because of the unavailability of the requested consecutive port range, the SR OS tries to find and allocate the next largest port set.
- If the SR OS, because of the lack of contiguous port ranges, allocates only a single port, the PORT_SET option is not present in the response.
- If the SR OS receives a PCP MAP request, with or without a PORT_SET option that tries to map one or more internal ports or port sets belonging to already existing mappings, then the request is considered to be a refresh request. Each of the matching port or port set mappings is processed independently, as if a separate refresh request was received. Consequently, the SR OS sends a Mapping Update message for each of the mappings.
- If multiple PORT_SET options are present in a single PCP MAP request, a MALFORMED_OPTION error is returned.
- If the PORT_SET size is zero, a MALFORMED_OPTION error is returned.
- If a Prefer Failure option is present, a MALFORMED_OPTION error is returned.
- When a PCP request contains both the PORT_SET and Port Reservation Port (N/N+1) options, only the PORT_SET is honored.
- PCP (with PORT_SET option or otherwise) should not be configured simultaneously with static port forwards in the same NAT pool. PCP allows for dynamic refreshment of port forwards while static port

forwards do not have this capability. As a result, PCP port refresh of a port forward allocated statically may lead to unwanted behavior.

- If the PORT_SET capability is added to or removed from an operationally up PCP server on an SR OS node, the server resets its Epoch time and sends a Version 2 ANNOUNCE message as described in the PCP specification.

7.16.4 Universal plug and play Internet Gateway Device service

Universal Plug and Play (UPnP), which is a set of specifications defined by the UPnP forum. One specification is called Internet Gateway Device (IGD) which defines a protocol for clients to automatically configure port mappings on a NAT device. Today, many gaming, P2P, VoIP applications support the UPnP IGD protocol. The SR OS supports the following UPnP version 1 **InternetGatewayDevice version 1** features:

- Supports only L2-Aware NAT hosts.
- Distributed subscriber management is not supported.
- The UPnP server runs on NAT ISA and only serves the local L2-Aware NAT hosts on the same ISA.
- The UPnP server can be enabled per subscriber by configuring a **upnp-policy** in the sub-profile.
- UPnP discovery is supported.
- UPnP eventing is not supported.
- The following IGD devices and services are supported:
 - WANDevice
 - WANConnectionDevice
 - WANIPConnection service
- For WANIPConnection services:
 - Optional state variables in a WANIPConnection service are not supported.
 - Optional actions in a WANIPConnection services are not supported.
 - Wildcard ExternalPort is not supported.
 - Only supports wildcard RemoteHost.
 - Up to 64 bytes of port mapping description are supported.
 - The SR OS supports a vendor specific action **X_ClearPortMapping**. This clears all port mappings of the subscriber belonging to the requesting host. This action has no in or out arguments.
- If the NewExternalPort in an addPortMapping request is same as the external port of one existing UPnP port mapping:
 - If NewInternalClient is different from InternalClient of existing mapping, then the system rejects the request.
 - If NewInternalClient is same as InternalClient of existing mapping:
 - With strict-mode on, if the source IP address of the request is same as InternalClient of existing mapping, then the request is accepted; otherwise the request is rejected.
 - With strict-mode off, the request is accepted.

- The system also supports the Alc-UPnP-Sub-Override-Policy RADIUS VSA which can be included in access-accept or CoA request. It can be used to override the **upnp-policy** configured in sub-profile or disable UPnP for the subscriber. See RADIUS reference guide for detail usage.

7.16.4.1 Configuring UPnP IGD service

Procedure

Step 1. Configure L2-Aware NAT.

Step 2. Create a **upnp-policy**:

Step 3. Configure the **upnp-policy** as created in step 2 in the subscriber profile:

```
config>service
  upnp
    upnp-policy "test" create
    no description
    http-listening-port 5000
    mapping-limit 100
    no strict-mode
  exit
```

```
config>subscr-mgmt
  sub-profile "l2nat-upnp" create
  nat-policy "l2"
  upnp-policy "test"
  exit
```

7.17 NAT Point-to-Point Tunneling Protocol (PPTP) ALG

PPTP is defined in RFC 2637, *Point-to-Point Tunneling Protocol (PPTP)*, and is used to provide VPN connection for home/mobile users to gain secure access to the enterprise network. Encrypted payload is transported over GRE tunnel that is negotiated over TCP control channel. In order for PPTP traffic to pass through NAT, the NAT device must correlate the TCP control channel with the corresponding GRE tunnel. This mechanism is referred to as PPTP ALG.

7.17.1 PPTP protocol

There are two components of PPTP:

- TCP control connection between the two endpoints
- an IP tunnel operating between the same endpoints. These are used to transport GRE encapsulated PPP packets for user sessions between the endpoints. PPTP uses an extended version of GRE to carry user PPP packets.

The control connection is established from the PPTP clients (for example, home users behind the NAT) to the PPTP server which is located on the outside of the NAT. Each session that carries data between the two endpoints can be referred as call. Multiple sessions (or calls) can carry data in a multiplexed fashion over a tunnel. The tunnel protocol is defined by a modified version of GRE. Call ID in the GRE header is

used to multiplex sessions over the tunnel. The Call-ID is negotiated during the session/call establishment phase.

7.17.1.1 Supported control messages

- **control connection management**

The following messages are used to maintain the control connection:

- Start-Control-Connection-Request
- Start-Control-Connection-Reply
- Stop-Control-Connection-Request
- Stop-Control-Connection-Reply
- Echo-Request
- Echo-Reply

The remaining control message types are sent over the established TCP session to open/maintain sessions and to convey information about the link state:

- **call management**

Call management messages are used to establish/terminate a session/call and to exchange information about the multiplexing field (Call-id). Call-IDs must be captured and translated by the NAT. The call management messages are:

- Outgoing-Call-Request (contains Call ID)
- Outgoing-Call-Reply (contains Call ID and peer's Call-ID)
- Call-Clear-Request (contains Call ID)
- Call-Disconnect-Notify (contains Call ID)

- **error reporting**

This message is sent by the client to indicate WAN error conditions that occur on the interface supporting PPP.

Wan-Error-Notify contains Call ID and Peer's Call ID.

- **PPP session control**

This message is sent in both directions to setup PPP-negotiated options.

Set-Link-Info contains Call ID and Peer's Call ID.

After Call-ID is negotiated by both endpoints, it is inserted in GRE header and used as multiplexing field in the tunnel that carries data traffic.

7.17.1.2 GRE tunnel

A GRE tunnel is used to transport data between two PPTP endpoints. The packet transmitted over this tunnel has the general structure shown in the following figure.

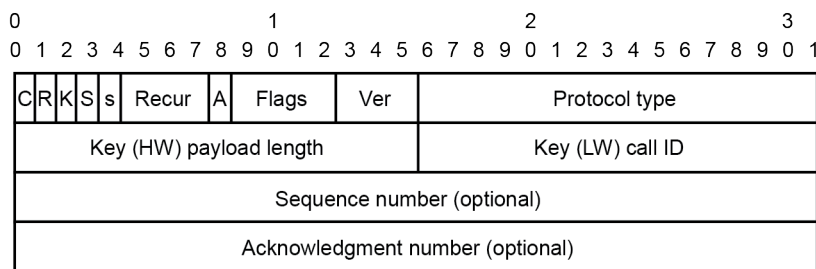
Figure 82: Structure of a packet transmitted between two PPTP endpoints over a GRE tunnel

Media header	Ethernet header, for example
IP header	Tunnel endpoints
GRE header	See following example
PPP packet	Packet payload including PPP header

sw1321

The following figure shows an example GRE header containing the Call ID of the peer for the session for which the GRE packet belongs.

Figure 83: GRE header example



sw1322

7.17.2 PPTP ALG operation

PPTP ALG is aware of the control session (Start Control Connection Request/Replay) and consequently it captures the Call ID field in all PPTP messages that carry that field. In addition to translating inside IP and TCP port, the PPTP ALG process data beyond the TCP header to extract the Call ID field and translate it inside of the Outgoing Call Request messages initiated from the inside of the NAT.

The GRE packets with corresponding Call IDs are translated through the NAT as follows:

- The inside source IP address is replaced by the outside IP address and the opposite is true for traffic in the opposite direction. This is standard IP address translation technique. The key is to keep the outside IP address of the control packets and corresponding data packets (GRE tunnel) the same.
- The Call-ID in the GRE packets in the direction of outside to inside is translated by the NAT according to the mappings that were created during session negotiation.

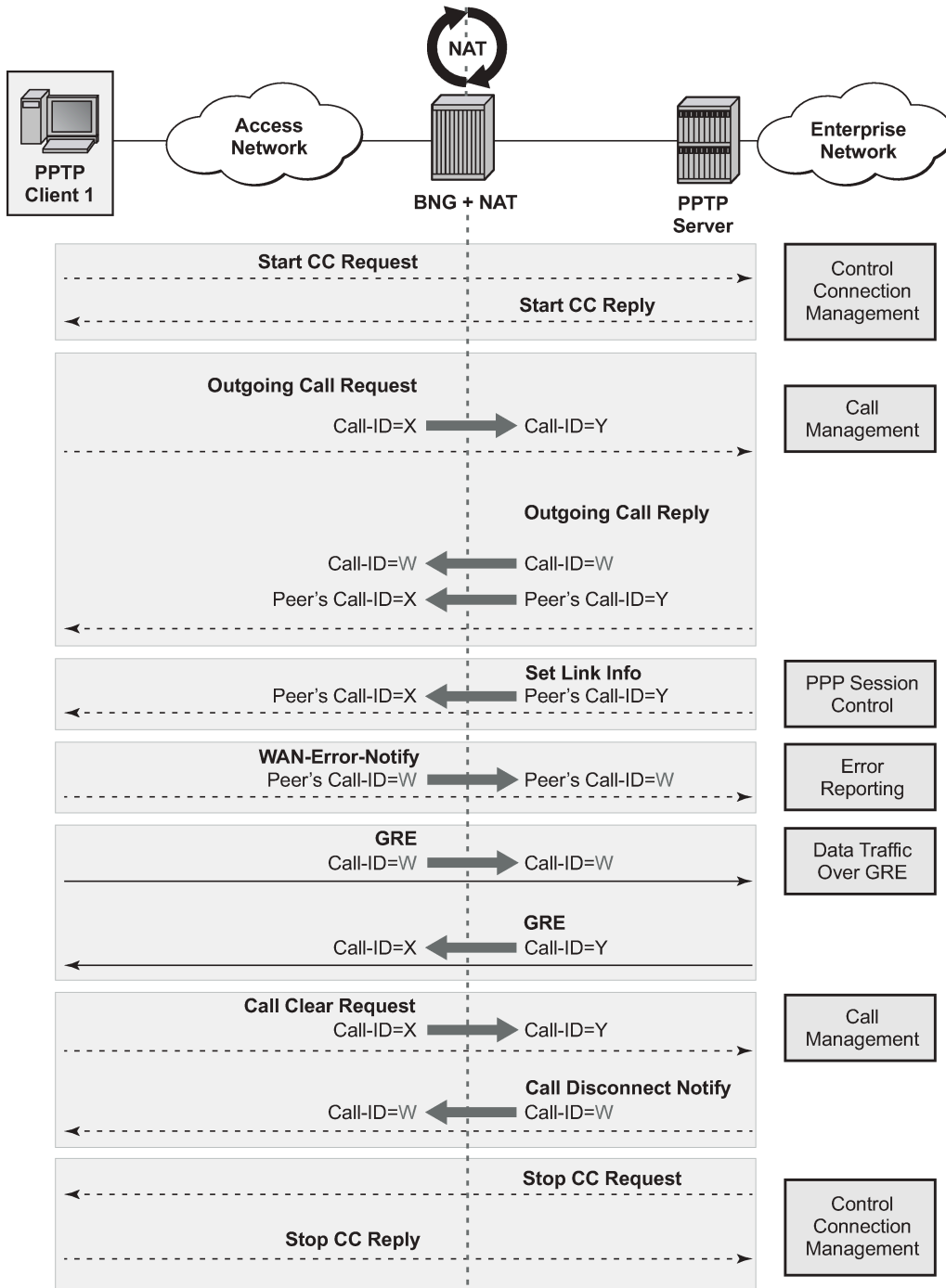
In addition, the following applies:

- GRE packets are translated and passed through the NAT only if they can be matched to an existing PPTP call for which the mapping already exists.
- Translation of the Call-IDs advertised by the PPTP server in the Outgoing Call Reply control message (this message is sent from the outside of the NAT to the inside) are not translated. Subsequently the Call ID in such messages are transparently passed through the NAT. There is no need to translate those Call IDs as their uniqueness between the two endpoints are guaranteed by the selection algorithm of the PPTP server. This can be thought of as destination TCP/UDP ports. They are not translated in the NAT. Instead only the source ports are translated.
- PPTP session initiation in the outside to inside direction through the NAT is not supported.

- Call-ID's are allocated and used in the same fashion as the outside TCP/UDP ports (random with parity). They are taken from the same port range as ICMP ports.

The basic principle of PPTP NAT ALG is shown in [Figure 84: NAT PPTP operation](#).

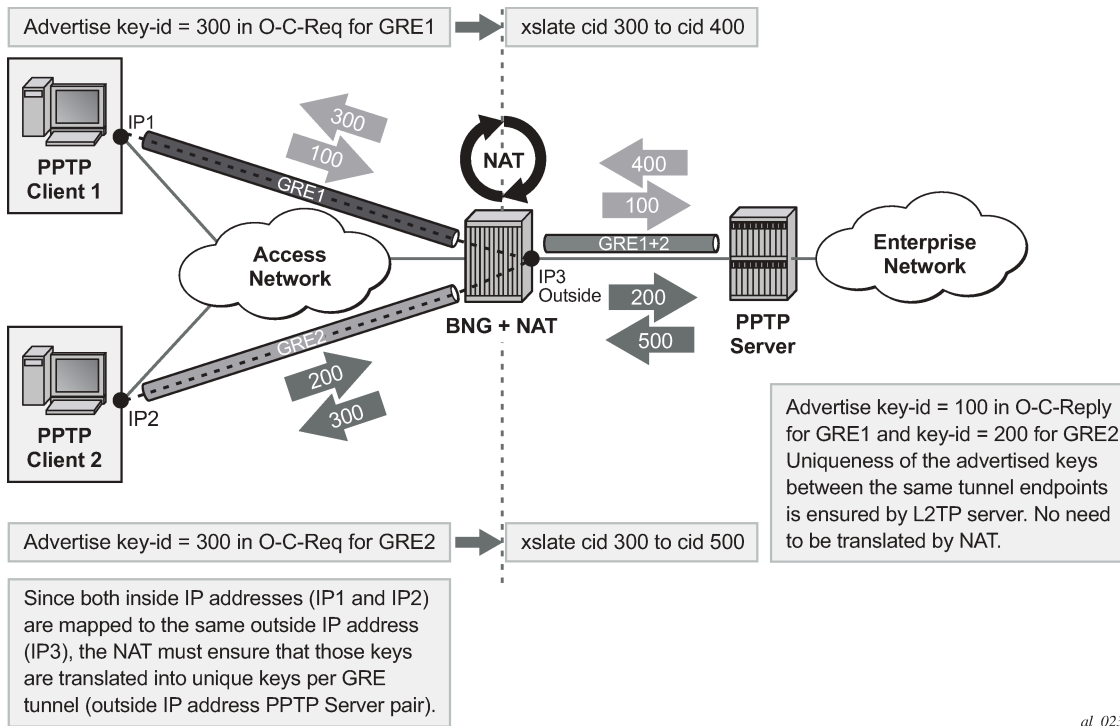
Figure 84: NAT PPTP operation



al_0238

The scenario where multiple clients behind the NAT are terminated to the same PPTP server is shown in [Figure 85: Merging of endpoints in NAT](#). In this case, it is possible that the source IP addresses of the two PPTP clients are mapped to the same outside address of the NAT. Because the endpoints of the GRE tunnel from the NAT to the PPTP server are the same for both PPTP clients (although their real source IP addresses are different), the NAT must ensure the uniqueness of the Call-IDs in the outbound data connection. This is where Call-ID translation in the NAT becomes crucial.

Figure 85: Merging of endpoints in NAT



7.17.3 Multiple sessions initiated from the same PPTP client node

The routers supports a deployment scenario where multiple calls (or tunnels) are established from a single PPTP node within a single control connection. In this case, there is only one set of Start-Control-Connection-Req/Reply messages (one control channel) and multiple sets of Outgoing-Call-Req/Reply messages.

7.17.4 Selection of call IDs in NAT

Call-Id are taken from the same pool as the ICMP port ranges. Port-ranges and Call-IDs are both 16-bit values. Call-id selection mechanism is the same as the outside TCP/UDP port selection mechanism (random with parity).

7.18 Modifying active NAT prefix list or NAT classifier via CLI

[Table 46: Modifying active NAT prefix list or NAT classifier](#) describes the outcome when the active NAT prefix list or NAT classifier is modified using CLI.

Table 46: Modifying active NAT prefix list or NAT classifier

Action	Outcome		Remarks
	L2-Aware	LSN	
CLI – Modifying prefix in the NAT prefix list	Existing flows are always checked whether they comply with the NAT prefix list that is currently applied in the sub-profile for the subscriber. If the flows do not comply with the current NAT prefix list, they are cleared after 5 seconds. The new flows immediately start using the updated settings.	Changing the prefix in the NAT prefix list internally re-subnets the outside IP address space.	Nat-prefix list is used with multiple NAT policies in L2-Aware NAT and for downstream internal subnet in dNAT-only scenario for LSN. Prefix can be modified (added, removed, remapped) at any time in the NAT prefix list, while the NAT policy must be first shut down via CLI.
CLI – Modifying the NAT classifier	Existing flows are always checked whether they comply with the NAT classifier that is currently applied in the active NAT policy for the subscriber. If the flows do not comply with the current NAT classifier, they are cleared after 5 seconds. The new flows immediately start using the updated settings.	Changing the NAT classifier has the same effect as in L2-Aware NAT; all existing flows using the NAT classifier are checked to see whether they comply with this classifier.	The NAT classifier is used for dNAT. The NAT classifier is referenced in the NAT policy.
CLI - Removing/ adding NAT policy in nat-prefix-list	Blocked	Not Applicable	—
CLI - Removing/ adding/ replacing NAT policy in sub-profile	Blocked	Not Applicable	—
CLI - Removing/ adding/ replacing NAT prefix-list under the	Not Applicable	Internally re-subnet, no effect on the flows	—

Action	Outcome		Remarks
	L2-Aware	LSN	
rtr/nat/inside/ dnat-only			

7.19 NAT logging

LSN logging is extremely important to the Service Providers (SP) who are required by the government agencies to track source of suspicious Internet activities back to the users that are hidden behind the LSN device.

The 7750 SR supports several modes of logging for LSN applications. Choosing the right logging model depends on the required scale, simplicity of deployment and granularity of the logged data.

For most purposes logging of allocation/de-allocation of outside port-blocks and outside IP address along with the corresponding LSN subscriber and inside service-id is sufficient.

In some cases, port-block based logging is not satisfactory and per flow logging is required.

7.19.1 Syslog/SNMP/local-file logging

The simplest form of LSN and L2-Aware NAT logging is via logging facility in the 7750 SR, commonly called logger. Each port-block allocation/de-allocation event is recorded and send to the system logging facility (logger). Such an event can be:

- recorded in the system memory as part of regular logs
- written to a local file
- sent to an external server by a syslog facility
- sent to a SNMT trap destination

In this mode of logging, all applications in the system share the same logger.

Syslog/SNMP/Local-File logging on LSN and NAT RADIUS-based logging are mutually exclusive.

Syslog/SNMP/local-file logging must be separately enabled for LSN and L2-Aware NAT in log even-control. The following displays relevant MIB events:

```
2012 tmnxNatPlBlockAllocationLsn
2013 tmnxNatPlBlockAllocationL2Aw
```

7.19.1.1 Filtering LSN events to system memory

In this example a single port-block [1884-1888] is allocated/de-allocated for the inside IP address 10.5.5.5 which is mapped to the outside IP address 198.51.100.1. Consequently the event is logged in the memory as.

```
2 2012/07/12 16:40:58.23 WEST MINOR: NAT #2012 Base NAT
"{2} Free 198.51.100.1 [1884-1888] -- vprn10 10.5.5.5 at 2012/07/12 16:40:58"

1 2012/07/12 16:39:55.15 WEST MINOR: NAT #2012 Base NAT
"{1} Map 198.51.100.1 [1884-1888] -- vprn10 10.5.5.5 at 2012/07/12 16:39:55"
```

When the needed LSN events are enabled for logging via event-control configuration, they can be logged to memory via standard log-id 99 or be filtered with a custom log-id, such as in this example (log-id 5):

Configuration:

```
*A:left-a20>config>log# info
-----
filter 1
  default-action drop
  entry 1
    action forward
    match
      application eq "nat"
      numbr eq 2012
    exit
  exit
exit
event-control "nat" 2001 suppress
event-control "nat" 2002 suppress
event-control "nat" 2003 suppress
event-control "nat" 2004 suppress
event-control "nat" 2005 suppress
event-control "nat" 2006 suppress
event-control "nat" 2007 suppress
event-control "nat" 2008 suppress
event-control "nat" 2009 suppress
event-control "nat" 2010 suppress
event-control "nat" 2011 suppress
event-control "nat" 2012 generate
event-control "nat" 2014 suppress
event-control "nat" 2015 suppress
event-control "nat" 2017 suppress
syslog 10
exit
log-id 5
  filter 1
    from main
    to memory
  exit
-----

*A:left-a20# show log event-control "nat"
=====
Log Events
=====
Application
ID#   Event Name                               P  g/s   Logged   Dropped
-----
2001  tmnxNatPLL2AwBlockUsageHigh             WA thr    0       0
2002  tmnxNatIsaMemberSessionUsageHigh        WA thr    0       0
2003  tmnxNatPLLSnMemberBlockUsageHigh        WA thr    0       0
2007  tmnxNatL2AwSubIcmpPortUsageHigh         WA thr    0       0
2008  tmnxNatL2AwSubUdpPortUsageHigh          WA thr    0       0
2009  tmnxNatL2AwSubTcpPortUsageHigh          WA thr    0       0
2010  tmnxNatL2AwSubSessionUsageHigh          WA thr    0       0
2012  tmnxNatPLBlockAllocationLsn             MI sup    0       0
2013  tmnxNatPLBlockAllocationL2Aw            MI sup    0       0
2014  tmnxNatResourceProblemDetected          MI thr    0       0
2015  tmnxNatResourceProblemCause             MI thr    0       0
2016  tmnxNatPLAddrFree                       MI sup    0       0
2017  tmnxNatPLLSnRedActiveChanged            WA thr    0       0
2018  tmnxNatPcpSrvStateChanged               MI thr    0       0
2020  tmnxNatMdaActive                        MI thr    0       0
2021  tmnxNatLsnSubBlksFree                   MI sup    0       0
2022  tmnxNatDetPlyChanged                     MI thr    0       0
2023  tmnxNatMdaDetectsLoadSharingErr         MI thr    0       0
2024  tmnxNatIsaGrpOperStateChanged           MI thr    0       0
```

```

2025 tmnxNatIsaGrpIsDegraded      MI thr      0      0
2026 tmnxNatLsnSubIcmpPortUsgHigh WA thr      0      0
2027 tmnxNatLsnSubUdpPortUsgHigh WA thr      0      0
2028 tmnxNatLsnSubTcpPortUsgHigh WA thr      0      0
2029 tmnxNatLsnSubSessionUsgHigh WA thr      0      0
2030 tmnxNatInAddrPrefixBlksFree  MI sup      0      0
2031 tmnxNatFwd2EntryAdded         MI sup      0      0
2032 tmnxNatDetPlcy0perStateChanged MI thr      0      0
2033 tmnxNatDetMap0perStateChanged MI thr      0      0
2034 tmnxNatFwd20perStateChanged   WA thr      0      0
=====

```

The event description is shown below:

tmnxNatPLL2AwBlockUsageHigh

The tmnxNatPLL2AwBlockUsageHigh notification is sent when the block usage of a Layer-2-Aware NAT address pool reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatIsaMemberSessionUsageHigh

The tmnxNatIsaMemberSessionUsageHigh notification is sent when the session usage of a NAT ISA group member reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatPLLsnMemberBlockUsageHigh

The tmnxNatPLLsnMemberBlockUsageHigh notification is sent when the block usage of a Large Scale NAT address pool reaches its high watermark ('true') or when it reaches its low watermark again ('false') on a particular member MDA of its ISA group.

tmnxNatLsnSubIcmpPortUsageHigh

The tmnxNatLsnSubIcmpPortUsageHigh notification is sent when the ICMP port usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatLsnSubUdpPortUsageHigh

The tmnxNatLsnSubUdpPortUsageHigh notification is sent when the UDP port usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatLsnSubTcpPortUsageHigh

The tmnxNatLsnSubTcpPortUsageHigh notification is sent when the TCP port usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatL2AwSubIcmpPortUsageHigh

The tmnxNatL2AwSubIcmpPortUsageHigh notification is sent when the ICMP port usage of a Layer-2-Aware NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatL2AwSubUdpPortUsageHigh

The tmnxNatL2AwSubUdpPortUsageHigh notification is sent when the UDP port usage of a Layer-2-Aware NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatL2AwSubTcpPortUsageHigh

The **tmnxNatL2AwSubTcpPortUsageHigh** notification is sent when the TCP port usage of a Layer-2-Aware NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatL2AwSubSessionUsageHigh

The **tmnxNatL2AwSubSessionUsageHigh** notification is sent when the session usage of a Layer-2-Aware NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatLsnSubSessionUsageHigh

The **tmnxNatLsnSubSessionUsageHigh** notification is sent when the session usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

tmnxNatPlBlockAllocationLsn

The **tmnxNatPlBlockAllocationLsn** notification is sent when an outside IP address and a range of ports is allocated to a NAT subscriber associated with a Large Scale NAT (LSN) pool, and when this allocation expires.

tmnxNatPlBlockAllocationL2Aw

The **tmnxNatPlBlockAllocationL2Aw** notification is sent when an outside IP address and a range of ports is allocated to a NAT subscriber associated with a Layer-2-Aware NAT pool, and when this allocation expires.

tmnxNatResourceProblemDetected

The **tmnxNatResourceProblemDetected** notification is sent when the value of the object **tmnxNatResourceProblem** changes.

tmnxNatResourceProblemCause

The **tmnxNatResourceProblemCause** notification is to describe the cause of a NAT resource problem.

tmnxNatPlAddrFree

The **tmnxNatPlAddrFree** notification is sent when a range of outside IP addresses becomes free at once.

tmnxNatPlLsnRedActiveChanged

The **tmnxNatPlLsnRedActiveChanged** notification is related to NAT Redundancy sent when the value of the object **tmnxNatPlLsnRedActive** changes. The cause is explained in the **tmnxNatNotifyDescription** which is a printable character string.

tmnxNatMdaActive

The **tmnxNatMdaActive** notification is sent when the value of the object **tmnxNatIsaMdaStatOperState** changes from 'primary' to any other value, or the other way around. The value 'primary' means that the MDA is active in the group.

tmnxNatLsnSubBlksFree

The **tmnxNatLsnSubBlksFree** notification is sent when all port blocks allocated to a Large Scale NAT (LSN) subscriber are released.

The NAT subscriber is identified with its subscriber ID **tmnxNatNotifyLsnSubId**.

To further facilitate the identification of the NAT subscriber,

its type `tmnxNatNotifySubscriberType`, inside IP address `tmnxNatNotifyInsideAddr` and inside virtual router instance `tmnxNatNotifyInsideVRtrID` are provided.

The values of `tmnxNatNotifyMdaChassisIndex`, `tmnxNatNotifyMdaCardSlotNum` and `tmnxNatNotifyMdaSlotNum` identify the ISA MDA where the blocks were processed.

All notifications of this type are sequentially numbered with the `tmnxNatNotifyPlSeqNum`.

The value of `tmnxNatNotifyNumber` is the numerical identifier of the NAT policy used for this allocation; it can be used for correlation with the `tmnxNatPlBlockAllocationLsn` notification; the value zero means that this notification can be correlated with all the `tmnxNatPlBlockAllocationLsn` notifications of the subscriber.

`tmnxNatDetPlcyChanged`

The `tmnxNatDetPlcyChanged` notification is sent when something changed in the Deterministic NAT map.

[CAUSE] Such a change may be caused by a modification of the `tmnxNatDetPlcyTable` or the `tmnxNatDetMapTable`.

[EFFECT] Traffic flows of one or more given subscribers, subject to NAT, may be assigned different outside IP address and/or outside port.

[RECOVERY] Managers that rely on the offline representation of the Deterministic NAT map should get an updated copy.

`tmnxNatMdaDetectsLoadSharingErr`

The `tmnxNatMdaDetectsLoadSharingErr` notification is sent periodically at most every 10 seconds while a NAT ISA MDA detects that it is receiving packets erroneously, due to incorrect load-balancing by the ingress IOM.

The value of `tmnxNatNotifyCounter` is the incremental count of dropped packets since the previous notification sent by the same MDA.

[CAUSE] The ingress IOM hardware does not support a particular NAT function's load-balancing, for example an IOM-2 does not support deterministic NAT.

[EFFECT] The MDA drops all incorrectly load-balanced traffic.

[RECOVERY] Upgrade the ingress IOM, or change the configuration.

`tmnxNatIsaGrpOperStateChanged`

The `tmnxNatIsaGrpOperStateChanged` notification is sent when the value of the object `tmnxNatIsaGrpOperState` changes.

`tmnxNatIsaGrpIsDegraded`

The `tmnxNatIsaGrpIsDegraded` notification is sent when the value of the object `tmnxNatIsaGrpDegraded` changes.

`tmnxNatLsnSubIcmpPortUsgHigh`

The `tmnxNatLsnSubIcmpPortUsgHigh` notification is sent when the ICMP port usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

The subscriber is identified with its inside IP address or prefix `tmnxNatNotifyInsideAddr` in the inside virtual router instance `tmnxNatNotifyInsideVRtrID`.

tmnxNatLsnSubUdpPortUsghigh

The **tmnxNatLsnSubUdpPortUsghigh** notification is sent when the UDP port usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

The subscriber is identified with its inside IP address or prefix **tmnxNatNotifyInsideAddr** in the inside virtual router instance **tmnxNatNotifyInsideVRtrID**.

tmnxNatLsnSubTcpPortUsghigh

The **tmnxNatLsnSubTcpPortUsghigh** notification is sent when the TCP port usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

The subscriber is identified with its inside IP address or prefix **tmnxNatNotifyInsideAddr** in the inside virtual router instance **tmnxNatNotifyInsideVRtrID**.

tmnxNatLsnSubSessionUsghigh

The **tmnxNatLsnSubSessionUsghigh** notification is sent when the session usage of a Large Scale NAT subscriber reaches its high watermark ('true') or when it reaches its low watermark again ('false').

The subscriber is identified with its inside IP address or prefix **tmnxNatNotifyInsideAddr** in the inside virtual router instance **tmnxNatNotifyInsideVRtrID**.

tmnxNatInAddrPrefixBlksFree

The **tmnxNatInAddrPrefixBlksFree** notification is sent when all port blocks allocated to one or more subscribers associated with a particular set of inside addresses are released by this system.

The type of subscriber(s) is indicated by **tmnxNatNotifySubscriberType**.

The set of inside IP addresses is associated with the virtual router instance indicated by **tmnxNatNotifyInsideVRtrID** and is of the type indicated by **tmnxNatNotifyInsideAddrType**

The set of inside IP addresses consists of the address prefix indicated with **tmnxNatNotifyInsideAddr** and **tmnxNatNotifyInsideAddrPrefixLen** unless these objects are empty and zero; if **tmnxNatNotifyInsideAddr** is empty and **tmnxNatNotifyInsideAddrPrefixLen** is zero, the set contains all IP addresses of the indicated type.

The values of **tmnxNatNotifyMdaChassisIndex**, **tmnxNatNotifyMdaCardSlotNum** and **tmnxNatNotifyMdaSlotNum** identify the ISA MDA where the blocks were processed.

All notifications of this type are sequentially numbered with the **tmnxNatNotifyPlSeqNum**.

This type of notification is typically the consequence of one or more configuration changes; the nature of these changes is indicated in the **tmnxNatNotifyDescription**.

tmnxNatFwd2EntryAdded

[CAUSE] The **tmnxNatFwd2EntryAdded** notification is sent when a row is added to or removed from the **tmnxNatFwd2Table** by other means than operations on the **tmnxNatFwdAction**;

a conceptual row can be added to or removed from the table by operations on the `tmnxNatFwdAction` object group or otherwise, by means of the PCP protocol or automatically by the system, for example when a subscriber profile is changed.

When the row is added, the value of the object `tmnxNatNotifyTruthValue` is 'true'; when the row is removed, it is 'false'.

[EFFECT] The specified NAT subscriber can start receiving inbound traffic flows.

[RECOVERY] No recovery required; this notification is the result of an operator or protocol action.

`tmnxNatDetPlcyOperStateChanged`

[CAUSE] The `tmnxNatDetPlcyOperStateChanged` notification is sent when the value of the object `tmnxNatDetPlcyOperState` changes. The cause is explained in the `tmnxNatNotifyDescription`.

`tmnxNatDetMapOperStateChanged`

[CAUSE] The `tmnxNatDetMapOperStateChanged` notification is sent when the value of the object `tmnxNatDetMapOperState` changes. The cause is explained in the `tmnxNatNotifyDescription`.

`tmnxNatFwd2OperStateChanged`

[CAUSE] The `tmnxNatFwd2OperStateChanged` notification is sent when the value of the object `tmnxNatFwd2OperState` changes. This is related to the state of the ISA MDA where the forwarding entry is located, or the availability of resources on that MDA.

In the case of Layer-2-Aware NAT subscribers, the `tmnxNatFwd2OperState` is 'down' while the subscriber is not instantiated. This would typically be a transient situation.

[EFFECT] The corresponding inward bound packets are dropped while the operational status is 'down'.

[RECOVERY] If the ISA MDA reboots successfully, or another ISA MDA takes over, no recovery is required. If more resources become available on the ISA MDA, no recovery is required.

7.19.1.2 NAT logging to a local file

In this case, the destination of log-id 5 in the following example would be a local file instead of memory:

```
*A:left-a20>config>log# info
-----
file-id 5
  description "nat logging"
  location cf1:
  rollover 15 retention 12
exit

log-id 5
  filter 1
  from main
  to file 5
exit
```

The events are logged to a local file on the compact flash cf1 in a file under the /log directory.



Note: Logging to the compact flash (CF) represents a single point of failure. Performance (logs per second) of logging onto the CF is limited in comparison to other logging methods (RADIUS, Syslog, and IPFIX). Failure to generate logs because of a failed CF or performance limitation results in dropped NAT traffic. For this reason, local NAT logging in the SR OS is recommended only in a lab environment.

7.19.2 SNMP trap logging

In case of SNMP logging to a remote node, the log destination should be set to SNMP destination. Allocation de-allocation of each port block triggers sending a SNMP trap message to the trap destination.

```
*A:left-a20>config>log# info
-----
filter 1
  default-action drop
  entry 1
    action forward
    match
      application eq "nat"
      number eq 2012
    exit
  exit
exit

snmp-trap-group 6
  trap-target "nat" address 192.168.1.10 port 9001 snmpv2c notify-community "private"
exit
log-id 6
  filter 1
  from main
  to snmp
exit
```

```

⊕ Internet Protocol Version 4, Src: 1.1.1.1 (1.1.1.1), Dst: 114.0.1.10 (114.0.1.10)
⊖ User Datagram Protocol, Src Port: snmptrap (162), Dst Port: etl servicemgr (9001)
    Source port: snmptrap (162)
    Destination port: etl servicemgr (9001)
    Length: 358
    ⊕ Checksum: 0x0e2c [correct]
⊖ Simple Network Management Protocol
    version: v2c (1)
    community: private
    ⊖ data: snmpv2-trap (7)
        ⊖ snmpv2-trap
            request-id: 1
            error-status: noError (0)
            error-index: 0
            ⊖ variable-bindings: 14 items
                ⊕ 1.3.6.1.2.1.1.3.0: 19054240
                ⊕ 1.3.6.1.6.3.1.1.4.1.0: 1.3.6.1.4.1.6527.3.1.3.65.0.12 (iso.3.6.1.4.1.6527.3.1.3.65.0.12)
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.2.0:
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.4.0:
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.5.0: 50000001
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.8.0: 1894
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.9.0: 1898
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.10.0: 07dc070d00321b002b0000
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.13.0: 1
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.3.0:
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.6.0:
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.7.0: 1a000038
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.11.0:
                ⊕ 1.3.6.1.4.1.6527.3.1.2.65.2.17.0: 5

```

7.19.3 NAT syslog

NAT logs can be sent to a syslog remote facility. A separate syslog message is generated for every port-block allocation/de-allocation.

```

*A:left-a20>config>log#info
-----
...
    filter 1
        default-action drop
        entry 1
            action forward
            match
                application eq "nat"
                number eq 2012
            exit
        exit
    exit
    syslog 7
        address 192.168.1.10
    exit

    log-id 7
        filter 1
            from main
            to syslog 7
        exit
-----

```

```

Internet Protocol Version 4, Src: 1.1.1.1 (1.1.1.1), Dst: 114.0.1.10 (114.0.1.10)
User Datagram Protocol, Src Port: syslog (514), Dst Port: syslog (514)
  Source port: syslog (514)
  Destination port: syslog (514)
  Length: 184
  Checksum: 0x3539 [correct]
    [Good Checksum: True]
    [Bad Checksum: False]
Syslog message: LOCAL7.INFO: Jul 13 15:04:53 1.1.1.1 TMNX: 35 Base NAT-INDETERMINATE-tmNxNatPBlockAllocationLsn-2012 [NAT]: {45} Map 80.0.0.1 [1994-1998] -- vprn10 26.0.0.56 at 2012/07/13 15:04:53
  1011 1... = Facility: LOCAL7 - reserved for local use (23)
  ....110 = Level: INFO - informational (6)
  Message: Jul 13 15:04:53 1.1.1.1 TMNX: 35 Base NAT-INDETERMINATE-tmNxNatPBlockAllocationLsn-2012 [NAT]: {45} Map 80.0.0.1 [1994-1998] -- vprn10 26.0.0.56 at 2012/07/13 08:04:53\n

```

Severity level for this event can be changed via CLI:

```

*A:left-a20# configure log event-control "nat" 2012 generate
<severity-level>
cleared indeterminate critical major minor warning

```

7.19.4 LSN RADIUS logging

LSN RADIUS logging (or accounting) is based on RADIUS accounting messages defined in RFC 2866. It requires a user to have RADIUS accounting infrastructure in place. For that reason, LSN RADIUS logging and LSN RADIUS accounting terms can be used interchangeably.

This mode of logging operation is introduced so that the shared logging infrastructure in 7750 SR can be offloaded by disabling syslog/SNMP/local-file LSN logging. The result is increased performance and higher scale, particularly in cases when multiple BB-ISA cards within the same system are deployed to perform aggregated LSN functions.

An additional benefit of LSN RADIUS logging over syslog/SNMP/local-file logging is reliable transport. Although RADIUS accounting relies on unreliable UDP transport, each accounting message from the RADIUS client must be acknowledged on the application level by the receiving end (accounting server).

Each port-block allocation or deallocation is reported to an external accounting (logging) server in the form of START, INTERIM-UPDATE, or STOP messages. The type of accounting messages generated depends on the mode of operation. The modes of operation are as follows:

- **START and STOP per port-block**

An accounting START is generated when a new port-block for the LSN subscriber is allocated. Similarly, the accounting STOP is generated when the port-block is released. Each accounting START and STOP pair of messages that are triggered by port-block allocation or deallocation within the same subscriber have the same Acct-Multi-Session-Id (subscriber significant) but a different Acct-Session-Id (port-block significant). This mode of operation is enabled by the inclusion of Acct-Multi-Session-Id within the NAT accounting policy.

- **START and STOP per subscriber**

An accounting START is generated when the first port block for the NAT subscriber is allocated. Each consecutive port-block allocation or deallocation triggers an INTERIM-UPDATE message with the same Acct-Session-Id (subscriber significant). The termination cause attribute in accounting STOP messages indicates the reason for port-block deallocation. Deallocation of the last port block for the LSN subscriber triggers an accounting STOP message. There is no Acct-Multi-Session-Id present in this mode of operation.

The accounting messages are generated and reported directly from the BB-ISA card, therefore bypassing accounting infrastructure residing on the Control Plane Module (CPM).

LSN RADIUS logging is enabled per NAT group. To achieve the required scale, each BB-ISA card in the NAT group with LSN RADIUS logging enabled runs a RADIUS client with its own unique source IP address. Accounting messages can be distributed to up to five accounting servers that can be accessed in round-robin fashion. Alternatively, in direct access mode, only one accounting server in the list is used. When this server fails, the next one in the list is used.

Perform the following steps to enable LSN RADIUS logging:

1. Configure the LSN RADIUS policy. The policy defines the following:

- accounting destination
- inclusion of RADIUS attributes that are sent in accounting messages to the destination
- source IP addresses per BB-ISA card (RADIUS client) in the NAT group



Note: The accounting policy applies to both LSN and WLAN-GW. Some attributes are only applicable to NAT, some are only applicable WLAN-GW, and some are applicable to both.

2. Apply this policy to the NAT group. This automatically enables RADIUS accounting on every BB-ISA card in the group, provided that each BB-ISA card has an IP address.

The following example shows the LSN RADIUS accounting policy configuration.

Example: LSN RADIUS accounting policy configuration (MD-CLI)

```
[ex:/configure aaa radius isa-policy "1"]
A:admin@node-2# info detail
## apply-groups
## apply-groups-exclude
description "RADIUS accounting policy for NAT"
## password
nas-ip-address-origin system-ip
## python-policy
accounting {
    include-attributes {
        acct-delay-time false
        acct-triggered-reason false
        called-station-id true
        calling-station-id false
        circuit-id false
        class false
        dhcp-options false
        dhcp-vendor-class-id false
        frame-counters true
        framed-ip-address true
        framed-ip-netmask false
        framed-ipv6-prefix false
        hardware-timestamp true
        ipv6-address false
        mac-address false
        multi-session-id true
        nas-identifier true
        nas-ip-address false
        nas-ipv6-address false
        nas-port false
        nas-port-id false
        nas-port-type false
        nat-inside-service-id true
        nat-outside-ip-address true
        nat-outside-service-id true
        nat-port-range-block true
        nat-subscriber-string false
    }
}
```



```

    octet-counters true
    proxied-subscriber-data false
    release-reason true
    remote-id false
    rssi false
    session-time true
    subscriber-id false
    toserver-dhcp6-options false
    ue-creation-type false
    user-name true
    wlan-ssid-vlan false
    xconnect-tunnel-local-ipv6-address false
    xconnect-tunnel-remote-ipv6-address false
    xconnect-tunnel-service false
    xconnect-tunnel-type false
    xconnect-tunnel-home-address false
    millisecond-event-timestamp false
    credit-control-quota false
  }
}
...
servers {
  source-address-range 192.168.1.20
  timeout 5
  total-tries 3
  router-instance "Base"
  access-algorithm direct
  ipv6 {
    mtu 9000
    ## source-prefix
  }
  server 1 {
    ## apply-groups
    ## apply-groups-exclude
    admin-state disable
    ip-address 192.168.1.10
    secret "ZVo7IYMjSxbdWlocPvLeFh5a8Xa1DVL0c2uzKvMGRnIKJo37JJjKLeoTXIPkD7h
QljmD3aC8ZdQ0Slw=" hash2
    purpose {
      accounting {
        udp-port 1813
      }
      ## authentication
      ## coa
    }
  }
}
}

```

Example: LSN RADIUS accounting policy configuration (classic CLI)

```

*A:node-2>config>aaa>isa-radius-plcy$ info detail
-----
description "RADIUS accounting policy for NAT"
nas-ip-address-origin system-ip
no password
no periodic-update
user-name-format mac mac-format alu
acct-include-attributes
  no acct-delay-time
  no acct-trigger-reason
called-station-id

```

```
no calling-station-id
no circuit-id
no class
no credit-control-quota
no dhcp-options
no dhcp-vendor-class-id
no dhcp6-options
frame-counters
framed-ip-addr
no framed-ip-netmask
no framed-ipv6-prefix
hardware-timestamp
inside-service-id
no ipv6-address
no mac-address
no millisecond-event-timestamp
multi-session-id
nas-identifier
no nas-ip-address
no nas-ipv6-address
no nas-port
no nas-port-id
no nas-port-type
no nat-subscriber-string
octet-counters
outside-ip
outside-service-id
port-range-block
release-reason
no remote-id
session-time
no subscriber-data
no subscriber-id
no ue-creation-type
user-name
no wifi-rssi
no wifi-ssid-vlan
no xconnect-tunnel-home-address
no xconnect-tunnel-local-ipv6-address
no xconnect-tunnel-remote-ipv6-address
no xconnect-tunnel-service
no xconnect-tunnel-type
exit

...

servers
access-algorithm direct
retry 3
router "Base"
source-address-range 192.168.1.20
timeout sec 5
ipv6
    mtu 9000
    no source-prefix
exit
server 1 create
shutdown
accounting port 1813
no authentication
no coa
ip-address 192.168.1.10
secret "ZVo7IYMjSxbdWlocPvLeFh5a8Xa1DVL0c2uzKvMGRnIKJo37JJj
KleoTXIPkd7hQljmD3aC8ZdQ0Slw=" hash2
```

```
exit
exit
-----
```



Note: The NAT subscriber string and subscriber data attributes are only relevant when subscriber-aware NAT is enabled.

Use the following command to assign one unique IPv4 address to each BB-ISA card from the range of IPv4 addresses configured:

- **MD-CLI**

```
configure aaa radius isa-policy servers source-address-range
```

- **classic CLI**

```
configure aaa isa-radius-policy servers source-address-range
```



Note: This IPv4 address must be accessible from the accounting server.

The IP addresses are consecutively assigned to each BB-ISA, starting from the IP address configured by this command. The number of IP addresses allocated internally by the system corresponds to the number of BB-ISAs in the system.

Each BB-ISA is provisioned automatically with the first free IP address available, starting from the IP address that is configured using the **source-address-range** command. When a BB-ISA is removed from the system (or NAT group), it releases that IP address to be available to the next BB-ISA that comes online within the NAT group.

It is important to be mindful of the internally allocated IP addresses, because they are not explicitly configured in the system (other than the first IP address configured using the **source-address-range** command). However, those internally-assigned IP addresses can be seen using show commands in the routing table.

Use the following command to show route-table information.

```
show router route-table
```

The following example shows there is one BB-ISA card in the NAT group 1. Its source IP address is 192.168.1.120.

Output example: NAT group with one BB-ISA card

```
=====
Route Table (Router: Base)
=====
Dest Prefix[Flags]  Type    Proto  Age           Pref  Next Hop[Interface Name]  Metric
-----
80.0.0.1/32        Remote  NAT    02d18h24m    0     NAT outside: group 1 member 1  0
192.168.1.0/28     Local   Local  02d20h25m    0     radius                        0
192.168.1.20/32   Remote  NAT    00h38m29s    0     NAT outside: group 1 member 1  0
=====
```

To communicate with IPv6 servers, provide a /64 source prefix for all ISAs and ESA VMs to use. Use the following command to configure this prefix:

- **MD-CLI**

```
configure aaa radius isa-policy servers ipv6 source-prefix
```

- **classic CLI**

```
configure aaa isa-radius-policy servers ipv6 source-prefix
```

Each ISA or ESA uses one /128 address from this prefix as a source address when communicating with an IPv6 RADIUS server. Because this is a /64 prefix, there is no risk of the ISA or ESA VMs running out of allocated addresses or using addresses not assigned to them, as there is with IPv4.

There is no path MTU discovery when communicating with an IPv6 server. On IPv6 servers, only the originating host is allowed to fragment a packet. In this case, the host is the ISA or ESA VM. This means that any MTU applied to the IOM does not fragment the packet, as it does with IPv4. However, an MTU for IPv6 fragmentation can be configured manually. Use the following command to configure an MTU for IPv6 fragmentation:

- **MD-CLI**

```
configure aaa radius isa-policy servers ipv6 mtu
```

- **classic CLI**

```
configure aaa isa-radius-policy servers ipv6 mtu
```

When the ISA or ESA VM generates a packet larger than the MTU, IPv6 fragmentation is applied to the packet.



Note: Fragmentation and reassembly adds processing overhead. To avoid this, increase the MTU or reduce the message size (by including less RADIUS attributes, for example) when possible.

It is possible to load-balance accounting messages over multiple logging servers by configuring the access-algorithm to round-robin mode. After the LSN RADIUS accounting policy is defined, it must be applied to a NAT group.

The following example shows a NAT group with an LSN RADIUS accounting policy applied to it.

Example: NAT group with an LSN RADIUS accounting policy (MD-CLI)

```
*[ex:/configure isa nat-group 1]
A:admin@node-2# info
  admin-state enable
  radius-accounting-policy "nat-acct-basic"
  redundancy {
    active-mda-limit 1
  }
  mda 1/2 { }
```

Example: NAT group with an LSN RADIUS accounting policy (classic CLI)

```
*A:node-2>config>isa>nat-group# info
-----
  active-mda-limit 1
  radius-accounting-policy "nat-acct-basic"
  mda 1/2
```

```
no shutdown
```

The RADIUS accounting messages for when a Large Scale NAT44 subscriber has allocated two port blocks in a logging mode where accounting START or STOP is generated per port block are as follows.

```
Fri Jul 13 09:55:15 2012
NAS-IP-Address = 10.1.1.1
NAS-Identifier = "left-a20"
NAS-Port = 37814272
Acct-Status-Type = Start
Acct-Multi-Session-Id = "500052cd2edcaeb97c2dad3d7c2dad3d"
Acct-Session-Id = "500052cd2edcaeb96206475d7c2dad3d"
Called-Station-Id = "00-00-00-00-01-01"
User-Name = "LSN44@10.0.0.58"
Alc-Serv-Id = 10
Framed-IP-Address = 10.0.0.58
Alc-Nat-Outside-Ip-Addr = 198.51.100.1
Alc-Nat-Port-Range = "198.51.100.1 2024-2028 router base"
Acct-Input-Packets = 0
Acct-Output-Packets = 0
Acct-Input-Octets = 0
Acct-Output-Octets = 0
Acct-Input-Gigawords = 0
Acct-Output-Gigawords = 0
Acct-Session-Time = 0
Event-Timestamp = "Jul 13 2012 09:54:37 PDT"
Acct-Unique-Session-Id = "21c45a8b92709fb8"
Timestamp = 1342198515
Request-Authenticator = Verified

Fri Jul 13 09:55:16 2012
NAS-IP-Address = 10.1.1.1
NAS-Identifier = "left-a20"
NAS-Port = 37814272
Acct-Status-Type = Start
Acct-Multi-Session-Id = "500052cd2edcaeb97c2dad3d7c2dad3d"
Acct-Session-Id = "500052cd2edcaeb9620647297c2dad3d"
Called-Station-Id = "00-00-00-00-01-01"
User-Name = "LSN44@10.0.0.58"
Alc-Serv-Id = 10
Framed-IP-Address = 10.0.0.58
Alc-Nat-Outside-Ip-Addr = 198.51.100.1
Alc-Nat-Port-Range = "198.51.100.1 2029-2033 router base"
Acct-Input-Packets = 0
Acct-Output-Packets = 5
Acct-Input-Octets = 0
Acct-Output-Octets = 370
Acct-Input-Gigawords = 0
Acct-Output-Gigawords = 0
Acct-Session-Time = 1
Event-Timestamp = "Jul 13 2012 09:54:38 PDT"
Acct-Unique-Session-Id = "baf26e8a35e31020"
Timestamp = 1342198516
Request-Authenticator = Verified
```

The RADIUS accounting messages for when a Large Scale NAT44 subscriber has deallocated two port blocks in a logging mode where accounting START or STOP is generated per port block are as follows.

```
Fri Jul 13 09:56:18 2012
NAS-IP-Address = 10.1.1.1
NAS-Identifier = "left-a20"
NAS-Port = 37814272
```

```

Acct-Status-Type = Stop
Acct-Multi-Session-Id = "500052cd2edcaeb97c2dad3d7c2dad3d"
Acct-Session-Id = "500052cd2edcaeb96206475d7c2dad3d"
Called-Station-Id = "00-00-00-00-01-01"
User-Name = "LSN44@10.0.0.58"
Alc-Serv-Id = 10
Framed-IP-Address = 10.0.0.58
Alc-Nat-Outside-IP-Addr = 198.51.100.1
Alc-Nat-Port-Range = "198.51.100.1 2024-2028 router base"
Acct-Terminate-Cause = Port-Unneeded
Acct-Input-Packets = 0
Acct-Output-Packets = 25
Acct-Input-Octets = 0
Acct-Output-Octets = 1850
Acct-Input-Gigawords = 0
Acct-Output-Gigawords = 0
Acct-Session-Time = 64
Event-Timestamp = "Jul 13 2012 09:55:41 PDT"
Acct-Unique-Session-Id = "21c45a8b92709fb8"
Timestamp = 1342198578
Request-Authenticator = Verified

Fri Jul 13 09:56:20 2012
NAS-IP-Address = 10.1.1.1
NAS-Identifier = "left-a20"
NAS-Port = 37814272
Acct-Status-Type = Stop
Acct-Multi-Session-Id = "500052cd2edcaeb97c2dad3d7c2dad3d"
Acct-Session-Id = "500052cd2edcaeb9620647297c2dad3d"
Called-Station-Id = "00-00-00-00-01-01"
User-Name = "LSN44@10.0.0.58"
Alc-Serv-Id = 10
Framed-IP-Address = 10.0.0.58
Alc-Nat-Outside-IP-Addr = 198.51.100.1
Alc-Nat-Port-Range = "198.51.100.1 2029-2033 router base"
Acct-Terminate-Cause = Host-Request
Acct-Input-Packets = 0
Acct-Output-Packets = 25
Acct-Input-Octets = 0
Acct-Output-Octets = 1850
Acct-Input-Gigawords = 0
Acct-Output-Gigawords = 0
Acct-Session-Time = 65
Event-Timestamp = "Jul 13 2012 09:55:42 PDT"
Acct-Unique-Session-Id = "baf26e8a35e31020"
Timestamp = 1342198580
Request-Authenticator = Verified

```

Including the `Acct-Multi-Session-Id` attribute in the NAT accounting policy enables generating START and STOP messages for each allocation or deallocation of a port block within the subscriber. Otherwise, only the first and last port block for the subscriber would generate a pair of START and STOP messages. All port blocks in between would generate INTERIM-UPDATE messages.

The `User-Name` attribute in accounting messages is set to `app-name@inside-ip-address`. The `app-name` can be any of the following:

- LSN44
- DS-Lite
- NAT64

7.19.4.1 Periodic RADIUS logging

Currently-allocated NAT resources (such as a public IP address and a port block for a NAT subscriber) can be periodically refreshed via Interim-Update (I-U) accounting messages. This functionality is enabled by the periodic RADIUS logging facility. Its primary purpose is to keep logging information preserved for long-lived sessions in environments where NAT logs are periodically and deliberately deleted from the service provider's network. This is typically the case in countries where privacy laws impose a limit on the amount of time that the information about customer's traffic can be retained/stored in service provider's network.

Periodic RADIUS logging for NAT is enabled by the following command:

```
configure
aaa
  isa-radius-policy <name> create
    [no] periodic-update interval <hours> [rate-limit <r>]
```

The configurable interval dictates the frequency of I-U messages that are generated for each currently allocated NAT resource (such as a public IP address and a port block).

By default, the I-U messages are sent in rapid succession for a subscriber without any intentional delay inserted by SR OS. For example, a NAT subscriber with 8 NAT policies, each configured with 40 port ranges generates 320 consecutive I-U messages at the expiration of the configured interval. This can create a surge in I-U message generation in cases where intervals are synchronized for multiple NAT subscribers. This can have adverse effects on the logging behavior. For example, the logging server can drop messages because of its inability to process the high rate of incoming I-U messages.

To prevent this, the rate of I-U message generation can be controlled by a **rate-limit** CLI parameter.

The periodic logging is applicable to both modes of RADIUS logging in NAT:

- **Acct-Multi-Session-Id AVP is enabled**

In this case, accounting START/STOP messages are generated for each NAT resource (such as a public IP address and a port block) allocation/de-allocation. Acct-multi-session-id and acct-session-id messages in the periodic I-U messages for the currently allocated NAT resource are inherited from the acct START messages related to the same NAT resource.

- **Acct-Multi-Session-Id AVP is disabled**

In this case, the acct START is generated for the first allocated NAT resource for the subscriber (a public IP address and a port block) and the acct STOP message is generated when the last NAT resource for the subscriber is released. All of the in-between port block allocations for the same subscriber trigger I-U messages with the same acct-session-id as the one contained in the acct START message. To differentiate between the port-block allocations, releases and updates within the I-U messages for the same NAT subscriber, the Alc-Acct-Triggered-Reason AVP is included in every periodic I-U message. Sending the Alc-Acct-Triggered-Reason AVP is configuration dependent (enabled in the **isa-radius-policy>acct-include-attributes** context). The supported values for Alc-Acct-Triggered-Reason AVP in I-U messages are:

- Alc-Acct-Triggered-Reason=Nat-FREE (19) Generated when the port-block is released.
- Alc-Acct-Triggered-Reason=Nat-MAP (20) Generated when the port-block is allocated.
- Alc-Acct-Triggered-Reason = Nat-UPDATE (21) Generated during periodically scheduled I-U update.

The log for each port-block periodic update is carried in a separate I-U message.

7.19.4.1.1 Message pacing

Periodic I-U message output can be paced to avoid congestion at the logging server. Pacing is controlled by the **rate-limit** option of the **periodic-update** command. As an example, consider the following hypothetical case:

- 1 million NAT subscribers came up within 1 hour (16,666 NAT-subs per minute).
- On average, each NAT subscriber allocates two port blocks.
- This means that 2 million logs are sent to the logging server.
- If the **rate-limit** value is set to 100 (messages per second), on average it would take over 5 hours to send all those messages at the given rate.
- In this case, it would be prudent to set the interval value to at least 6 hours, or increase the **rate-limit** value so there is no time overlap between the old and new logs.

In the case of an MS-ISA switchover or a NAT multi-chassis redundancy switchover, there is a chance that a large number of subscribers become active at approximately the same time on the newly active MS-ISA (or chassis). This causes a large number of logs to be sent in a relatively short amount of time, which may overwhelm the logging server. The **rate-limit** parameter is designed to help in such situations.

7.19.4.2 RADIUS buffer management on ISA or ESA-VM

RADIUS accounting messages (Accounting-Request) sent from an SR OS node to a RADIUS server are acknowledged by an Accounting-Response message generated by the RADIUS server. This acknowledgment confirms that the RADIUS server has successfully received and recorded the client information, such as NAT logs. This communication between the SR OS node and the RADIUS server occurs over UDP transport, as defined in RFC 2866.

If there is a lack of acknowledgments to the SR OS because of RADIUS server overload, server failure, or network failure, the SR OS node backs off RADIUS messages and retransmits them in line with the configured ISA RADIUS policy. After several retransmissions, as specified by the ISA RADIUS policy, the message is discarded.

Each ISA or ESA-VM maintains a buffer capable of storing 32,000 outstanding transactions toward the RADIUS server. A slow or unresponsive RADIUS server can result in buffer exhaustion.

When the buffer is full and unable to accept additional messages, the new NAT port blocks are not allocated or existing ones are not released.

The Acct Tx Timeouts counter in the following output example shows the number of dropped RADIUS messages because of these timeouts. This command was executed while the RADIUS server was unresponsive.

```
show aaa isa-radius-policy "AcctPolicy1"

=====
ISA RADIUS policy "AcctPolicy1"
=====
Description          : AcctPolicy1 associated with nat-grp 1
Include attributes acct : framed-ip-addr
                       : nas-identifier
                       : nat-subscriber-string
                       : user-name
                       : inside-service-id
                       : outside-service-id
                       : outside-ip
```



```

port-range-block
hardware-timestamp
release-reason
multi-session-id
frame-counters
octet-counters
session-time
called-station-id
subscriber-data
framed-ip-netmask
circuit-id
remote-id
dhcp-options
dhcp-vendor-class-id
mac-address
nas-port-id
nas-port-type
calling-station-id
subscriber-id
acct-trigger-reason
ue-creation-type
wifi-rssi
acct-delay-time
wifi-ssid-vlan
nas-ip-address
nas-port
class
ipv6-address
framed-ipv6-prefix
dhcp6-options
Include attributes auth      : nas-ip-address
                             : nas-ipv6-address
User name format             : mac
User name MAC format        : alu
NAS-IP-Address              : system-ip
Python policy                : (Not Specified)
Periodic update
  interval (hours)          : (Not Specified)
  rate limit (messages/s)   : (Not Specified)
-----
RADIUS server settings
-----
Router                       : 2000
Source address start         : 192.0.2.100
Source address end           : 192.0.2.100
IPv6 Source prefix           : (Not Specified)
IPv6 MTU                     : 9000
Access algorithm             : direct
Retry                        : 10
Timeout (s)                  : 60
Last management change       : 05/04/2023 11:36:05
=====
Servers for "AcctPolicy1"
=====
Index Address                Acct-port Auth-port CoA-port
-----
1    10.10.17.7              1813     0         0
=====
Status for ISA RADIUS server policy "AcctPolicy1"
=====

```

```

Server 1, group 1, member 1
-----
Purposes Up                : (None)
Purposes Hold down        : (None)
Source IP address         : 192.0.2.100
Acct Tx Requests          : 655451
Acct Tx Retries           : 589905
Acct Tx Timeouts          : 65544
Acct Rx Replies           : 0
Auth Tx Requests          : 0
Auth Tx Retries           : 0
Auth Tx Timeouts          : 0
Auth Rx Replies           : 0
CoA Rx Requests          : 0
=====

```

7.19.5 Summarization logs and bulk operations

Bulk operations, such as removing a NAT policy or shutting down a NAT pool, can trigger a cascade of events, such as release of NAT subscribers associated with the NAT policy or a NAT pool. To avoid excessive logging during those operations, summarization logs are used. These logs carry relational information that connect multiple events and are categorized under event log 99 on the CPM. Configurable destinations for these logs include SNMP notification (trap), syslog (sent in syslog format to the syslog collector), memory (sent to memory buffer), local file, and NETCONF.

Tracking NAT subscribers from the logs becomes more complicated if they were terminated because of bulk operations. A MAP log is generated when NAT resources for the subscriber are allocated; a FREE log is generated when NAT resources for the subscriber are released. Typically, individual MAP logs are paired with corresponding FREE logs to determine the identity and activity duration for the subscriber. However, during bulk operations, individual FREE logs are substituted with a summarized log containing relational information. In such cases, identifying NAT subscriber mappings may necessitate examining multiple logging sources, such as a combination of RADIUS and summarization logs.

To simplify log summarization, a policy ID is added as a connecting parameter in all logs. The policy ID follows the format:

```
plcy-id XX
```

Where: *XX* is a unique number representing a NAT policy and assigned by the router for each inside routing context, as shown in the following example:

```

670 2023/05/31 12:55:00.952 UTC MINOR: NAT #2012 vprn601 NAT
"{986} Map 10.10.10.1 [4001-4279] MDA 5/1 -- 1166016512 classic-lsn-sub
%203 vprn101 192.0.2.1 at 2023/05/31 12:55:00"

```

When an active NAT policy is removed from the configuration within an inside routing context, all NAT subscribers associated with that NAT policy in that context are removed from the system. Instead of generating individual FREE logs for each subscriber, a single summarized log is generated. This summarized log entry contains only the policy ID of the removed NAT policy and the inside service ID. To determine which NAT resources were released, the user must match the policy ID and the service ID in the summarization log with those in all MAP logs that lack a pairing explicit FREE log.

A summarization log is always created on the CPM, regardless of whether RADIUS logging is enabled.

A summarization log is generated on the CPM under the following circumstances:

- **NAT policy removal**

If there is a single NAT policy for each inside routing context, the summarization log contains the inside service ID (vprn or base). To identify the terminated NAT mappings for subscribers, search all individual MAP logs matching the service ID from the summarization log.

When there are multiple NAT policies per inside routing context, the summarization log contains the inside service ID and policy ID. Search individual logs based on policy ID and inside service ID to identify subscribers affected by the NAT policy removal.

- **pool shutdown**

The router sends a summarization log with the outside service ID and all IP address ranges in the pool. Match individual logs based on outside IP address and outside service ID to identify released subscribers.

- **IP address range removal from the pool**

The summarization log includes the outside service ID and the removed IP address range. Match individual logs based on the outside IP addresses in the range and the outside service ID to identify the released subscribers.

- **Non deterministic source prefix removal**

The summarization log includes the removed source prefix, policy ID, and inside service ID.

- **Last AFTR address removal**

The summarization log includes the inside service ID.

- **DS-Lite or NAT64 node shutdown**

The summarization log includes the inside service ID.

- **Deterministic NAT prefix creation or removal**

The summarization log includes the inside service ID.

Summarization logs are enabled by event controls 2021 (tmnxNatLsnSubBlksFree), 2016 (tmnxNatPIAddrFree), and 2030 (tmnxNatInAddrPrefixBlksFree). These events are suppressed by default. Event control 2021 also reports when all port blocks for a subscriber are freed.

7.19.5.1 Summarization logs and RADIUS logging

RADIUS logging does not generate summarization logs because RADIUS accounting messages (start, interim-updates, and stop messages) are generated for each subscriber individually. Therefore, using RADIUS logging to also send summarization logs to every subscriber would be ineffective.

Instead, during RADIUS logging bulk operations, summarization logs are generated exclusively on the CPM using event logs. One exception is when a NAT accounting policy is removed, in which case a RADIUS acct-off message is sent without an accompanying summarization log. For bulk operations with RADIUS logging, operators must rely on both RADIUS logging and summarization logs on the CPM.

For example, if a RADIUS log sequence indicates a mapping for <inside IP 1, outside IP 1, port-block 1>, and later a mapping log for <inside IP 2, outside IP 1, port-block 1>, it suggests that the FREE log for <inside IP 1, outside IP 1, port-block 1> is missing. This could mean that either the FREE log for <inside IP 1, outside IP 1, port-block 1> was lost, or a policy, pool, and address range were removed from the configuration. In the latter case, the operator should check the CPM log for the summarization message.

7.19.6 Integrated L2-Aware NAT RADIUS logging and BNG accounting

In L2-Aware NAT, the logging of NAT resources is integrated with ESM RADIUS accounting. The NAT-related resources reporting is described in [Table 47: Integrated ESM and NAT accounting](#).

Accounting START messages carry only the RADIUS Event-Timestamp (type 55), which correctly reflects the creation of the initial port block and outside IP address for L2-Aware NAT. The initial port block and outside IP address allocation in the ISA or ESA for a L2-Aware subscriber is triggered by the control plane (CPM) when the first session or host is created. This means that the initial port block and outside IP address creation in the ISA or ESA is not triggered by data traffic. However, data traffic triggers the creation of extend port blocks.

Interim-Updates and STOP accounting messages carry two timestamps. This is because the RADIUS accounting message is generated by the CPM at the time indicated by the Event-Timestamp, which may not accurately reflect the time of the extended port block allocation or de-allocation that occurs on ISA or ESA.

- **RADIUS Event-Timestamp (type 55) with a 1 second resolution**

This timestamp is updated by the CPM with the time that the Interim-Update message is generated.

- **Nokia Alc-ISA-Event-Timestamp (type 86)**

This is updated only when an event on the ISA or ESA occurs, for example, an extension port block is allocated or de-allocated. The format and resolution of this timestamp is the same as the format of the Event-Timestamp.

A summary of integrated ESM and NAT RADIUS logging is shown in [Table 47: Integrated ESM and NAT accounting](#). Only RADIUS attributes relevant to NAT are shown.

Table 47: Integrated ESM and NAT accounting

ESM and NAT integrated RADIUS accounting/logging			
Acct msg type	Queue-instance (SLA-profile instance) accounting	Session or host accounting	Comments
Start	<p>An Acct START message is generated for every SLA profile instantiation and every accounting START message contains NAT-related information carried in Alc-Nat-Port-Range (26.6527.121) which includes the outside IP address, newly allocated initial port block, outside router ID, and NAT policy ID.</p> <p>If there are multiple SLA profile instances per a NAT-enabled ESM subscriber, this information is repeated for all additional SLA profile instances.</p>	<p>Acct START is generated for every new session or host of a NAT-enabled subscriber.</p> <p>This message carries:</p> <ul style="list-style-type: none"> • the outside IP address and the initial port for the first session or the host for the subscriber • the outside IP address, the initial port block, and extended port blocks for any existing sessions or hosts of the subscriber <p>The NAT related information is carried in the following RADIUS attribute: Alc-Nat-Port-Range(26.6527.121)</p> <p>This attribute includes the outside IP address, port blocks, outside router ID, and NAT policy. There is</p>	<p>The initial port block and outside IP address are always advertised in accounting START messages, regardless of whether there is a single session, host, multiple sessions, hosts per subscriber, or the sessions or hosts are NAT-enabled.</p>

ESM and NAT integrated RADIUS accounting/logging			
Acct msg type	Queue-instance (Sla-profile instance) accounting	Session or host accounting	Comments
		no distinction between NAT-enabled sessions or hosts and non NAT-enabled sessions of hosts (that is, non NAT-enabled sessions or hosts also carry NAT information) for a NAT enabled subscriber.	
Regular Interim-Update	<p>The message reports existing in-use NAT resources (the cumulative update) for each SLA profile instance:</p> <p>Alc-Nat-Port-Range (26.6527.121)</p> <p>The outside IP address, all existing port blocks, outside router ID, and NAT policy ID.</p> <p>Alc-ISA-Event-Timestamp(241.26.6527.86)</p> <p>The time of the last extended port block allocation or de-allocation on the ISA or ESA.</p> <p>Event-Timestamp (55)</p> <p>The time when the RADIUS message is generated on the CPM.</p> <p>This is repeated for all NAT-enabled sessions of an ESM subscriber.</p>	<p>This message reports the existing in-use NAT resources (the cumulative update) for each session:</p> <p>Alc-Nat-Port-Range (26.6527.121)</p> <p>The outside IP address, all existing port blocks, outside router ID, and NAT policy.</p> <p>Alc-ISA-Event-Timestamp(241.26.6527.86)</p> <p>The time of the last extended port block allocation or de-allocation on the ISA or ESA.</p> <p>Event-Timestamp (55)</p> <p>The time when the RADIUS message is generated on the CPM</p> <p>This is repeated for all NAT-enabled sessions or hosts of an ESM subscriber.</p>	—
Triggered Interim-Update	<p>This message carries differential updates tracking changes only for extended port blocks of the existing subscriber. The initial port-block is not advertised in the triggered Interim-Update but instead it is only advertised in the accounting START (map) or STOP (free) message.</p> <p>Alc-Nat-Port-Range (26.6527.121)</p> <p>The outside IP address, newly allocated or de-allocated extended port block, outside router ID, and NAT policy ID.</p> <p>Alc-Acct-Triggered-Reason (26.6527.163)</p> <ul style="list-style-type: none"> NAT-MAP (20) 	<p>This message carries differential updates tracking changes only for extended port blocks of the existing subscriber. The initial port block is never advertised in the triggered Interim-Update but is only advertised in accounting START (map) or STOP (free) message.</p> <p>Alc-Nat-Port-Range (26.6527.121)</p> <p>The outside IP address, newly allocated or de-allocated extended port block, outside router ID, and NAT policy ID.</p> <p>Alc-Acct-Triggered-Reason (26.6527.163)</p> <ul style="list-style-type: none"> NAT-MAP (20) 	<p>If the last session of the subscriber is terminated, and at the same time this session has extended port blocks in use, two consecutive RADIUS accounting messages are sent (regardless of the accounting model):</p> <ul style="list-style-type: none"> a triggered I-U message with extended PBs

ESM and NAT integrated RADIUS accounting/logging			
Acct msg type	Queue-instance (SLA-profile instance) accounting	Session or host accounting	Comments
	<ul style="list-style-type: none"> NAT-FREE (19) <p>These are the reasons for this triggered Interim-Update message. An extended port block is allocated (MAP) or de-allocated (FREE).</p> <p>Alc-ISA-Event-Timestamp (241.26.6527.86)</p> <p>The time of the extended port block allocation or de-allocation on the ISA or ESA.</p> <p>Event-Timestamp (55)</p> <p>The time when the RADIUS message is generated on the CPM.</p> <p>This is repeated for each SLA-profile instance (queuing instance).</p>	<ul style="list-style-type: none"> NAT-FREE (19) <p>The reason for this triggered Interim-Update message. An extended port block is allocated (MAP) or de-allocated (FREE).</p> <p>Alc-ISA-Event-Timestamp (241.26.6527.86)</p> <p>The time of the extended port block allocation or de-allocation on the ISA or ESA.</p> <p>Event-Timestamp (55)</p> <p>The time when the RADIUS message is generated on the CPM.</p> <p>This is repeated for all sessions or hosts of an ESM subscriber.</p>	<ul style="list-style-type: none"> a STOP message for the last session termination for the subscriber. This STOP message contains the initial PB (and outside IP address). <p>A subscriber termination is infrequent event.</p>
Stop	<p>Accounting STOP messages are sent when an SLA profile instance (queuing-instance) is terminated (the last session associated with it is terminated).</p> <p>If the terminated SLA-profile instance (queuing instance) is the last for the subscriber, the accounting STOP message only carry the initial port block (and outside IP address). Any extended port blocks that were released are be reported in immediately preceding triggered Interim-Update message.</p> <p>If the terminated SLA-profile instance (queuing instance) is not the last for the subscriber, the accounting STOP message carries the initial port-block (and outside IP address) and any extended port blocks that are still allocated for the subscriber, but not used any more by this terminated SLA-profile instance.</p> <p>Alc-Nat-Port-Range (26.6527.121)</p> <p>The outside IP address, initial port block, outside router ID, and NAT policy ID.</p>	<p>Accounting STOP message is sent when a session or host of a NAT enabled subscriber is terminated.</p> <p>If the terminated session of host is the last for the subscriber, the accounting STOP message carries only the initial port-block (and outside IP address). Any extended port blocks that were released are reported in immediately preceding triggered Interim-Update messages.</p> <p>If the terminated session or host is not the last for the subscriber, the accounting STOP message carries the initial port-block (and outside IP address) and any extended port blocks that are still allocated for the subscriber, but not used any more by this terminated session or host.</p> <p>Alc-Nat-Port-Range (26.6527.121)</p> <p>The outside IP address, initial port block, outside router ID, and NAT policy ID.</p> <p>Alc-ISA-Event-Timestamp (241.26.6527.86)</p>	<p>Each accounting stream (START, I-U, STOP) is treated as a separate entity and it contains NAT information that can overlap with other accounting streams (for the queuing instance or a session) of the same subscriber.</p> <p>A complete NAT information is always conveyed in an accounting stream, for example, for every PB allocation a matching de-allocation can be found on the same stream. In other words, there are no known cases where a PB allocation is reported on one accounting stream,</p>

ESM and NAT integrated RADIUS accounting/logging			
Acct msg type	Queue-instance (SLA-profile instance) accounting	Session or host accounting	Comments
	<p>Alc-ISA-Event-Timestamp (241.26.6527.86)</p> <p>The time of the last extended port block allocation or de-allocation on the ISA or ESA.</p> <p>Event-Timestamp (55)</p> <p>The time when the RADIUS message is generated on the CPM. This information is generated for every SLA-profile instance (queuing instance) termination, meaning that the information is repeated if the subscriber has multiple SLA-profile instances.</p>	<p>The time of the last extended port block allocation or de-allocation on the ISA or ESA.</p> <p>Event-Timestamp (55)</p> <p>The time when the RADIUS message is generated on the CPM</p> <p>This information is generated upon termination of every session or host of an L2-Aware subscriber.</p>	<p>but de-allocation is reported on another.</p>

The following are examples showing only relevant NAT related attributes:

- A session is created for a L2-Aware NAT subscriber. At the time of session instantiation, a RADIUS accounting START messages is generated.

```
Alc-Nat-Port-Range = "192.168.20.2 2001-2024 router base l2-aware"
Event-Timestamp = T1
```

The outside IP address 192.168.20.2 and initial port block [2001-2004] are allocated at time T1.

- New extended port block is allocated. Differential data is carried in a triggered Interim-Update message.

```
Alc-Nat-Port-Range = "192.168.20.2 3000-3023 router base l2-aware"
Alc-Acct-Triggered-Reason = Nat-Map (20)
Event-Timestamp = T3
Alc-ISA-Event-Timestamp = T2
```

Only the new allocated port blocks are present in this update with the triggered reason Nat-Map (20).

This port block was allocated on ISA or ESA at time T2 which may be different than time T3 at which the Interim-Update is sent to the RADIUS server.

This difference may be small if there is no congestion in the system. It may be larger if there is congestion in the system while the notifications from the ISA or ESA are queued internally in the system waiting to be transported to the CPM which is backlogged. A reason for CPM backlog can be from a high volume of RADIUS messages that are sent to the RADIUS servers.

- Periodic Interim-Update messages are triggered at regular intervals and carries cumulative (or absolute) data.

```
Alc-Nat-Port-Range = "192.168.20.2 2001-2024, 3000-3023 router base l2-aware"
Event-Timestamp = T4
```

```
Alc-ISA-Event-Timestamp = T2
```

This update carries both previously allocated port blocks, the initial port block and the extended port block.

T4 in Event-Timestamp reflects the time when the message is generated, while the Alc-ISA-Event-Timestamp is unchanged from the previous update because no new event occurred on the ISA or ESA.

- An existing extended port block is de-allocated. Differential data is carried in triggered Interim-Update message.

```
Alc-Acct-Triggered-Reason = Nat-Free (19)
Alc-Nat-Port-Range = "192.168.20.2 3000-3023 router base l2-aware"
Event-Timestamp = T6
Alc-ISA-Event-Timestamp = T5
```

Only the de-allocated port block is present in this update with the triggered reason NAT-Free (19).

This port block was de-allocated on ISA or ESA at time T5 which may be different than time T6 at which the Interim-Update is sent to the RADIUS server.

- At session termination, a RADIUS accounting STOP message with initial port block is generated.

```
Alc-Nat-Port-Range = "192.168.20.2 2001-2024 router base l2-aware"
Event-Timestamp = T7
Alc-ISA-Event-Timestamp = T5
```

This final update for the session carries the initial port block that is no longer used by the terminated session, host or queuing instance. Although this session is terminated, the initial port block can be used by another sessions still present under the same L2-Aware NAT subscriber.

T7 in Event-Timestamp reflects the time when the message is generated, while the Alc-ISA-Event-Timestamp is always the same as in the previous triggered accounting Interim-Update message.

7.19.6.1 Enabling RADIUS logging for L2-Aware NAT subscribers

This method of logging is enabled with a RADIUS accounting policy as follows:

```
[configure subscriber-mgmt radius-accounting-policy <name>]
  include-radius-attribute {
    nat-port-range true
    acct-triggered-reason true
  }
```

For session or host type accounting, the generation of periodic Interim-Update messages must be enabled as follows:

```
[configure subscriber-mgmt radius-accounting-policy <name>]
  session-accounting {
    interim-update true
  }
```


7.19.6.2 Timestamp interpretation

Extended port block functionality in L2-Aware NAT contains an additional time stamp into the logging framework. In addition to the standardized Event-Timestamp that is carried in every RADIUS accounting message, a NAT-related timestamp is included. This additional timestamp is introduced in the accounting stream after the first extended port block for the subscriber is allocated and then it is present in every accounting message in the stream. It represents the time of the last extended port block allocation or de-allocation as recoded by the ISA or ESA.

The two timestamps should be interpreted as the following:

- The Standard Event-Timestamp (55) attribute records the time when the accounting message was generated on the CPM.
- The Alc-ISA-Event-Timestamp (241.26.6527.86) attribute records the time of the last NAT related event (extended port block allocation or de-allocation).

For example, a periodic I-U message below indicates that at the time 1000, a subscriber has two ports blocks allocated: [2001-2024] and [3000-3023]. The last change related to extended port blocks was at time 500.

The following are Periodic Interim-Updates with the NAT-related attributes:

```
Alc-Nat-Port-Range = "192.168.20.2 2001-2024,3000-3023 router base l2-aware"  
Event-Timestamp = 1000  
Alc-ISA-Event-Timestamp = 500
```

In the following scenario where:

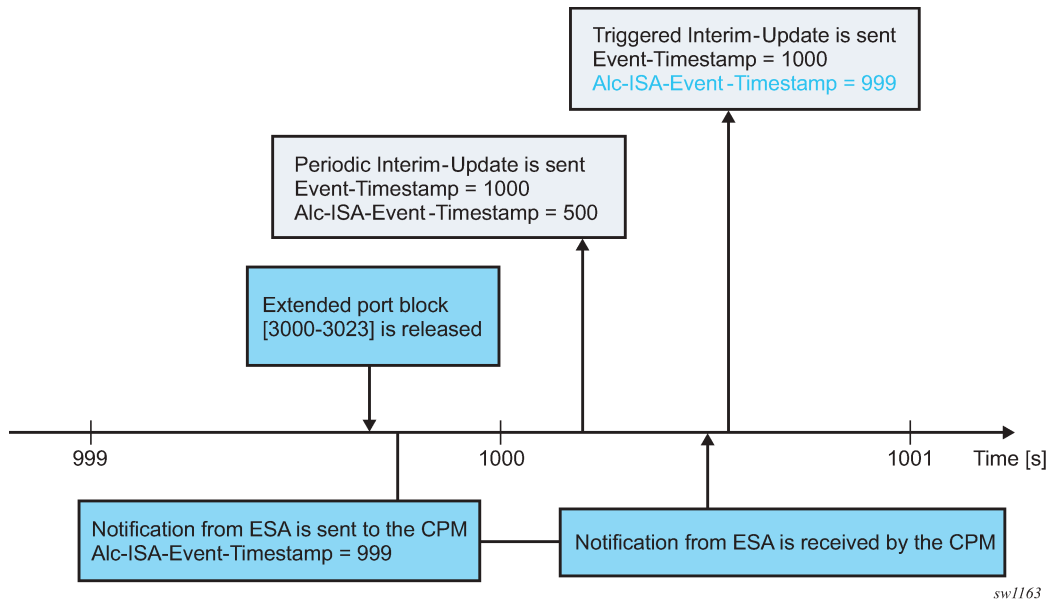
- the extended port block [3000-3023] was released a few milliseconds before the previous periodic Interim-Update message was sent
- notification from the ISA or ESA about this event has not reached the CPM in time for the event to be included in periodic Interim-Update

Then, a triggered Interim-Update would immediately follow the above periodic Interim-Update with relevant NAT related attributes:

```
Alc-Nat-Port-Range = "192.168.20.2 3000-3023 router base l2-aware"  
Alc-Acct-Triggered-Reason = Nat-Free  
Event-Timestamp = 1000  
Alc-ISA-Event-Timestamp = 999
```

Both messages have the same Event-Timestamp of 1000 because the resolution of this timestamp is 1 second. However, the port block [3000-3023] was released at time 999 as indicated by the Alc-ISA-Event-Timestamp. This scenario is shown in [Figure 86: Alc-ISA-Event-Timestamp](#).

Figure 86: Alc-ISA-Event-Timestamp



7.19.6.3 High logging rates

A system with on-demand port block allocations is dynamic with a possibility of generating a high volume of logs. Transporting NAT logs through ESM accounting relies on a generic RADIUS accounting infrastructure implemented in the SR which supports multiple RADIUS servers and failover mechanisms. In cases where the rate of accounting message exceeds the capacity of the entire accounting system, the queue of the accounting message toward the RADIUS servers in SR start to fill up. This can be caused by the internal condition in the SR by slow or even unresponsive RADIUS servers. Considering that NAT is only a contributor of the accounting messages in a larger accounting framework that includes ESM, the rate of the allocation and de-allocations of extended port blocks is internally limited. Although this does not prevent loss of accounting messages in the overloaded accounting system (for example, caused by slow RADIUS servers), it helps to reduce the chances that the system become overloaded.

7.19.6.4 Intra-chassis redundancy

Initial port blocks are preserved during an ISA or ESA switchover and are not affected by the switchover. However, extended port blocks are released during the switchover. Consequently, their release is reported in triggered Interim-Update messages.

7.19.7 LSN and L2-Aware NAT flow logging

LSN and L2-Aware NAT flow logging allows each BB-ISA card to export the creation and deletion of NAT flows to an external server. A NAT flow, or a Fully Qualified Flow, consists of the following parameters: inside IP, inside port, outside IP, outside port, foreign IP, foreign port, and protocol (UDP, TCP, ICMP).

```
Owner           : LSN-Host@10.10.10.101
Router          : 1
Policy          : mnp
FlowType        : UDP
Inside IP Addr  : 10.10.10.101
Inside Port     : 20001
Outside IP Addr : 192.168.20.28
Outside Port    : 2001
Foreign IP Addr : 192.168.5.4
Foreign Port    : 20001
Dest IP Addr    : 192.168.5.4
Nat Group       : 1
Nat Group Member : 1
```

The foreign IP address is the original IPv4 destination address as received by NAT on the inside. The destination IP address is the translated foreign IP address if that destination NAT is active (the destination NAT translates the destination IPv4 address of the packet).

Additional information, such as the inside or outside service ID and subscriber string, can be added to a flow record.

Flow logging can be deployed as an alternative to port-range logging or can be complementary (providing a more granular log for offline reporting or compliance). Certain operators have legal and compliance requirements that require extremely detailed logs, created per flow, to be exportable from the NAT node.

Because the setup rate of new flows is excessive, logging to an internal facility (like compact flash) is not possible except in debugging mode (which must specify match criteria down to the inside IP and service level).

Flow logging can be enabled on a per-NAT policy basis and, consequently, it is initiated from each BB-ISA card. The flow records can be exported to an external collector in an IPFIX format or a syslog format, both of which use UDP as the transport protocol. These UDP streams are stateless because of the significant volume of transactions. However they do contain sequence numbers so packet loss can be identified. They egress the chassis at the **fc nc**.

IPFIX and SYSLOG flow logging are configured using respective flow logging policies (such as **ipfix-export-policy** and **syslog-export-policy**). Each flow logging policy supports two destinations (collectors). One **ipfix-export-policy** and one **syslog-export-policy** can be used simultaneously in any one NAT group.

7.19.7.1 IPFIX flow logging

IPFIX defines two different types of messages that are sent from the IPFIX exporter (SR OS NAT node). The first contains a template set which is an IPFIX message that defines fields for subsequent IPFIX messages but contains no actual data of its own. The second IPFIX message type contains data sets. Here, the data is passed using the previous template set message to define the fields. This means an IPFIX message is not passed as sets of TLV, but instead, data is encoded with a scheme defined through the template set message.

While an IPFIX message can contain both a template set and data set, the SR OS node sends template set messages periodically without any data, whereas the data set messages are sent on demand and as required. When IPFIX is used over UDP, the default retransmission frequency of the template set messages defaults to 10 minutes. The interval for retransmission is configurable in CLI with a minimum interval of one minute and a maximum interval of 10 minutes. When the exporter first initializes, or when a configuration change occurs, the template set is sent out three times, one second apart. Templates are sent before any data sets, assuming that the collector is enabled, so that an IPFIX collector can establish the data template set.

Although the UDP transport is unreliable, the IPFIX sequence number is a 32-bit number that contains the total number of IPFIX data records sent for the UDP transport session before the receipt of the new IPFIX message. The sequence number starts with 0 and rolls over when it reaches 4,294,967,268.

The default packet size is 1500 bytes unless another value has been defined in the configuration (the range is 512 bytes through 9212 bytes inclusively). Traffic is originated from a random high port to the collector on port 4739. Multiple create and delete flow records are stuffed into a single IPFIX packet (although the mappings created are not delayed) until stuffing an additional data record would exceed MTU or a timer expires. The timer is not configurable and is set to 250 milliseconds (that is, should any mapping occur, a packet is sent within 250 milliseconds of that mapping being created).

Each collector has a 50-packet buffering space. If, because of excessive logging, the buffering space becomes unavailable, new flows are denied and the deletion of flows is delayed until buffering space becomes available.

Two collector nodes can be defined in the same IPFIX export policy for redundancy purposes.

7.19.7.2 Template formats

The SR OS supports two data formats. Their selection is controlled through CLI:

```
configure
service
  ipfix
    ipfix-export-policy <name> [create]
    template-format {format1|format2}
```

The difference between the two formats is related to the fields conveying information about the translated source IP addresses and ports (outside IP addresses and ports).

Format1 carries information about translated (outside) IP address in the sourceIPv4Address information element while in format2 this information element is replaced by the postNATSourceIPv4Address. Further, format1 does not convey any information about the translated source port (post-NAT) while a new information element postNATsourceTransportPort is introduced in format2 to carry this information.

Both formats use proprietary information element AluNatSubString carrying the original source IP address, before NAT is performed.

The template and data sets are formatted according to RFC 5101, *Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of IP Traffic Flow Information*.

Standardized data fields are defined in RFC 5102, *Information Model for IP Flow Information Export*, and in IANA registry <https://www.iana.org/assignments/ipfix/ipfix.xhtml#ipfix-information-elements>.

In addition to standardized data fields, IPFIX supports vendor-proprietary data fields which contains an Enterprise Number specific to each vendor.

The supported information elements and their description for each format is provided in [Table 48: IPFIX fields and formats](#). EN in [Table 48: IPFIX fields and formats](#) stands for Enterprise Number (0 = IETF, 637 = Nokia) and IE-Id represents Information Element Identifier.

Table 48: IPFIX fields and formats

Field	EN, IE-Id	Format 1	Format 2
flowId	0, 148	A unique (per-observation domain ID) ID for this flow. Used for tracking purposes only (opaque value). The flow ID in a create and a delete mapping record must be the same for a specific NAT mapping.	A unique (per-observation domain ID) ID for this flow. Used for tracking purposes only (opaque value). The flow ID in a create and a delete mapping record must be the same for a specific NAT mapping.
sourceIPv4Address	0, 8	The outside (translated) IP address used in the NAT mapping. In format2, this is replaced by postNATSourceIPv4Address.	N/A
postNATSourceIPv4Address	0, 255	N/A	The outside (translated) IP address used in the NAT mapping. This replaces the sourceIPv4Address field from format1.
destinationIPv4Address	0, 12	The foreign or remote IP address used in the NAT mapping.	The foreign or remote IP address used in the NAT mapping.
sourceTransportPort	0, 7	The outside (translated) source port used in the NAT mapping.	This is the original source port (before NAT translation) on the inside
postNAPTsourceTrasportPort	0, 227	N/A	The outside (translated) source port used in the NAT mapping
destinationTransportPort	0, 11	The destination port used in the NAT mapping.	The destination port used in the NAT mapping.
flowStartMilliseconds	0, 152	The timestamp of when the flow was created (chassis NTP derived) in milliseconds from epoch.	The timestamp of when the flow was created (chassis NTP derived) in milliseconds from epoch.
flowEndMilliseconds	0, 153	The timestamp of when the flow was destroyed (chassis NTP	The timestamp of when the flow was destroyed (chassis

Field	EN, IE-Id	Format 1	Format 2
		derived) in milliseconds from epoch.	NTP derived) in milliseconds from epoch.
protocolIdentifier	0, 4	Protocol (UDP, TCP, ICMP)	Protocol (UDP, TCP, ICMP)
flowEndReason	0, 136	<p>The reasons for flow termination. The following Flow End Reasons are supported:</p> <ul style="list-style-type: none"> 0x01: Idle Timeout. A mapping expired (because of UDP or TCP timeout) 0x03: end of Flow Detected. A mapping closed (only used for TCP after a FIN or RST). 0x04: forced end. Collects all other reasons included administrative or failure case. 	<p>The reasons for flow termination. The following Flow End Reasons are supported:</p> <ul style="list-style-type: none"> 0x01: Idle Timeout. A mapping expired (because of UDP or TCP timeout) 0x03: end of Flow Detected. A mapping closed (only used for TCP after a FIN or RST). 0x04: forced end. Collects all other reasons included administrative or failure case.
paddingOctets	0, 210	Padding	N/A
aluInsideServiceId	637, 91	The 16-bit service ID representing the inside service ID. This field is not applicable in L2-Aware NAT and is set to NULL in this case.	The 16-bit service ID representing the inside service ID. This field is not applicable in L2-Aware NAT and is set to NULL in this case.
aluOutsideServiceId	637, 92	The 16-bit service ID representing the outside service ID.	The 16-bit service ID representing the outside service ID.
aluNatSubString	637, 93	<p>A variable 8B aligned string that represents the NAT subscriber construct (as currently used in the tools>dump>service>nat>session commands). The original IP source address, before NAT is performed is included in this string.</p> <p>For example: LSN-Host@10.10.10.101</p>	<p>A variable 8B aligned string that represents the NAT subscriber construct (as currently used in the tools>dump>service>nat>session commands). The original IP source address, before NAT is performed is included in this string.</p> <p>For example: LSN-Host@10.10.10.101</p>

7.19.7.3 Template format 1 and format 2

Table 49: Template format 1 and Table 50: Template format 2 show information elements that are present in data sets during flow creation and deletion for the two formats. Each template set carries a unique template ID that is used to match the corresponding data set that carries the same ID in the set header (Set Id field).

Table 49: Template format 1

Flow creation template set		Flow deletion templates set	
Description	Size (B)	Description	Size (B)
flowId	8	flowId	8
sourceIPv4Address	4	sourceIPv4Address	4
destinationIPv4Address	4	destinationIPv4Address	4
sourceTransportPort	2	sourceTransportPort	2
destinationTransportPort	2	destinationTransportPort	2
flowStartMilliseconds	8	flowEndMilliseconds	8
protocolIdentifier	1	protocolIdentifier	1
paddingOctets	1	flowEndReason	1
aluInsideServiceID	2	aluInsideServiceID	2
aluOutsideServiceID	2	aluOutsideServiceID	2
aluNatSubString	var	aluNatSubString	var

Table 50: Template format 2

Flow creation template set		Flow deletion templates set	
Description	Size (B)	Description	Size (B)
flowId	8	flowId	8
postNATSourceIPv4Address	4	postNATSourceIPv4Address	4
destinationIPv4Address	4	destinationIPv4Address	4
sourceTransportPort	2	sourceTransportPort	2
postNAPTSourceTransportPort	2	postNAPTSourceTransportPort	2
destinationTransportPort	2	destinationTransportPort	2
flowStartMilliseconds	8	protocolIdentifier	1
protocolIdentifier	1	flowEndReason	1

Flow creation template set		Flow deletion templates set	
Description	Size (B)	Description	Size (B)
paddingOctets	1	flowEndMilliseconds	8
aluInsideServiceID	2	aluInsideServiceID	2
aluOutsideServiceID	2	aluOutsideServiceID	2
aluNatSubString	var	aluNatSubString	var

7.19.7.4 Configuration example

Large Scale NAT44 Flow Logging with format2:

1. A collector node along with other local transport parameters must be defined through an IPFIX export policy.

```
*A:BNGL>config>service>ipfix# info detail
-----
ipfix-export-policy "flow-logging" create
  no description
  template-format format2
  collector router "Base" ip 192.168.115.1 create
    mtu 1500
    source-address 192.0.2.2
    template-refresh-timeout min 5
    no shutdown
  exit
exit
```

To export flow records using UDP stream, the BB-ISA card must be configured with an appropriate IPv4 address within a designated VPRN. This address (/32) acts as the source for sending all IPFIX records and is shared by all ISA.

2. After the IPFIX export policy is defined, it must be applied within the NAT policy:

```
*A:BNGL>config>service>nat# info
-----
nat-policy "mnp" create
  pool "mnp" router Base
  ipfix-export-policy "flow-logging"
exit
```

Flow creation and flow deletion templates for format2, as captured in Wireshark, are shown in [Figure 87: Format2 templates](#).

Figure 87: Format2 templates

```

Cisco NetFlow/IPFIX
Version: 10
Length: 148
Timestamp: Oct 31, 2017 04:33:36.000000000 Central Daylight Time
FlowSequence: 0
Observation Domain Id: 1179650
Set 1
  FlowSet Id: Data Template (v10 [IPFIX]) (2)
  FlowSet Length: 64
  Template (Id = 256, Count = 11)
    Template Id: 256
    Field Count: 11
    Field (1/11): flowId
    Field (2/11): postNATSourceIPv4Address
    Field (3/11): IP_DST_ADDR
    Field (4/11): L4_SRC_PORT
    Field (5/11): postNAPTSourceTransportPort
    Field (6/11): L4_DST_PORT
    Field (7/11): flowStartMilliseconds
    Field (8/11): PROTOCOL
    Field (9/11): 91 [pen: Alcatel-Lucent (previously was 'Alcatel Data Network')]
    Field (10/11): 92 [pen: Alcatel-Lucent (previously was 'Alcatel Data Network')]
    Field (11/11): 93 [pen: Alcatel-Lucent (previously was 'Alcatel Data Network')]
Set 2
  FlowSet Id: Data Template (v10 [IPFIX]) (2)
  FlowSet Length: 68
  Template (Id = 257, Count = 12)
    Template Id: 257
    Field Count: 12
    Field (1/12): flowId
    Field (2/12): postNATSourceIPv4Address
    Field (3/12): IP_DST_ADDR
    Field (4/12): L4_SRC_PORT
    Field (5/12): postNAPTSourceTransportPort
    Field (6/12): L4_DST_PORT
    Field (7/12): PROTOCOL
    Field (8/12): flowEndReason
    Field (9/12): flowEndMilliseconds
    Field (10/12): 91 [pen: Alcatel-Lucent (previously was 'Alcatel Data Network')]
    Field (11/12): 92 [pen: Alcatel-Lucent (previously was 'Alcatel Data Network')]
    Field (12/12): 93 [pen: Alcatel-Lucent (previously was 'Alcatel Data Network')]

```

IPFIX flow creation data set, as captured in Wireshark, is shown in [Figure 88: Flow creation data set](#).

Figure 88: Flow creation data set

```

Cisco NetFlow/IPFIX
Version: 10
Length: 74
Timestamp: Oct 30, 2017 20:33:45.000000000 Central Daylight Time
FlowSequence: 0
Observation Domain Id: 1179650
Set 1
  FlowSet Id: (Data) (256)
  FlowSet Length: 58
  Flow 1
    Flow Id: 216172784529768566
    Post NAT Source IPv4 Address: 192.168.20.24 (192.168.20.24)
    DstAddr: 192.168.5.4 (192.168.5.4)
    SrcPort: 20001
    Post NAPT Source Transport Port: 2003
    DstPort: 20001
    Protocol: 17
    Enterprise Private entry: (Alcatel-Lucent (previously was 'Alcatel Data Network')) Type 91: value (hex bytes): 00 01
    Enterprise Private entry: (Alcatel-Lucent (previously was 'Alcatel Data Network')) Type 92: value (hex bytes): 00 00
    Enterprise Private entry: (Alcatel-Lucent (previously was 'Alcatel Data Network')) Type 93: value (hex bytes): 4c 53
    string_len_short: 18
    StartTime: Oct 31, 2017 04:33:45.460000000 Central Daylight Time

```

IPFIX flow destruction data set, as captured in Wireshark, is shown in [Figure 89: Flow destruction](#).

Figure 89: Flow destruction

```

Cisco NetFlow/IPFIX
  Version: 10
  Length: 75
  Timestamp: Oct 30, 2017 20:35:53.000000000 Central Daylight Time
  FlowSequence: 1
  Observation Domain Id: 1179650
  Set 1
    FlowSet Id: (Data) (257)
    FlowSet Length: 59
    Flow 1
      Flow Id: 216172784529768566
      Post NAT Source IPv4 Address: 192.168.20.24 (192.168.20.24)
      DstAddr: 192.168.5.4 (192.168.5.4)
      SrcPort: 20001
      Post NAT Source Transport Port: 2003
      DstPort: 20001
      Protocol: 17
      Flow End Reason: Idle timeout (1)
      Enterprise Private entry: (Alcatel-Lucent (previously was 'Alcatel Data Network')) Type 91: value (hex bytes): 00 01
      Enterprise Private entry: (Alcatel-Lucent (previously was 'Alcatel Data Network')) Type 92: value (hex bytes): 00 00
    [Enterprise Private entry: (Alcatel-Lucent (previously was 'Alcatel Data Network')) Type 93: value (hex bytes): 4c 53
      string_len_short: 18
    EndTime: Jan 4, 1973 14:48:48.103230848 Central Standard Time

```

7.19.7.5 Syslog flow logging

The format of syslog messages for NAT flow logging in SR OS adheres to RFC 3164, *The BSD Syslog Protocol*:

<PRI> <HEADER><MSG>

where:

- <PRI> (the "<" and ">" are included in the syslog message) is the configured facility*8+severity (as described in the *7450 ESS*, *7750 SR*, *7950 XRS*, and *VSR System Management Guide* and RFC 3164).
- <HEADER> defines the MMM DD HH:MM:SS <hostname>. Two characters always appear for the day (DD) field. Single-digit days are preceded with a space character. Time is recorded as local time (and not UTC). The time zone designator is not shown in this example, but each event has its own timestamp where the time-zone designator is shown.
- <MSG> defines the <log-prefix>: <seq> <application [<subject>]: <message>\n

where:

- <log-prefix> is an optional 32-character string of text (default = 'TMNX') as configured in the **log-prefix** command.
- <seq> is the log event sequence number (always preceded by a colon and a space char).
- The [<subject>] field may be empty resulting in []:
- <message> display a custom message relevant to the log event.
- \n is the standard ASCII new line character (hex 0A).

[Table 51: Syslog message fields for NAT flow logging](#) shows the syslog message fields for NAT flow logging.

Table 51: Syslog message fields for NAT flow logging

Field name	Value	Comments
PRI • severity • facility	• Default: 6 • Default: 16	• Configurable • Configurable
Timestamps	MMM DD HH:MM:SS	
<hostname>		The IP address of the SR OS system that is generating the message.
<log-prefix>		Configurable. This can be used as a field to differentiate between the vendors. For example, NOK(ia) in log-prefix indicates that this is a log format from a Nokia node so the operator can apply parsing logic accordingly.
<seq>		Sequence numbers can be used for tracking if loss in transit occurs.
<application>	NAT	The application that generated the log.
[<subject>]:	MDA ID	The BB-ISA on which the event occurred.
<message>		This is a custom part with specific information related to the event itself.

The message portion contains information relevant to the respective log event, even if this information is already repeated outside of the message (for example, timestamp). The fields in the message part are separated by a single whitespace for easier parsing and are placed in the order shown [Table 52: Message fields](#).

Table 52: Message fields

Field name	Value	Presence	Comments
NAT type	LSN44 NAT64	M(andatory)	
Event name	SADD	M	SADD – session added event

Field name	Value	Presence	Comments
	SDEL		SDEL – session deleted event
Timestamp	<TimeStamp>: <Year> <Mon> <Day> <hh:mm:ss:cs> <TZ>, Year is 4-digit, Mon is 3-letter abbreviation, TZ is a 1-5 character time-zone designator.	M	Because events can be combined in the same syslog message, each event is uniquely timestamped with the local time (not UTC), including the time zone designator. During daylight saving's time (summer), the time zone designator is replaced by the DST designator, which is configurable.
Protocol ID	1, 6, 17	M	ICMP, UDP, TCMP
Inside router	0 to 2147483650	M	0 represents Base 1 to 2147483650 represents VPRNs
Source IP address	IPv4 address in LSN44 and IPv6 address in NAT64	M	
Source port or ICMP identifier	0 to 65535	M	
Outside router	0 to 2147483650	M	
Outside (post NAT) IP address	IPv4 address	M	
Outside (post NAT) port or ICMP identifier	0 to 65535	M	
Foreign IP address	IPv4 address	O(ptional)	This is the original destination IPv4 address.
Foreign port or ICMP identifier	0 to 65535	O	
Destination IP address	IPv4 address	O	It represents the translated destination IP address.
Nat-policy	<name>	O	
Sub-ID	<sub-name>	O	"-" if requested by the configuration (includes the sub-id statement) but the sub-aware NAT is not

Field name	Value	Presence	Comments
			enabled. Otherwise, the sub-ID in the sub-aware NAT.

7.19.7.5.1 Sequence numbers

Each syslog message contains a sequence number. The sequence numbers are independently generated by each BB-ISA per collector, and they are monotonically increased by 1. The MDA ID carried in the syslog message is used to differentiate between the overlapping sequence numbers generated by different BB-ISAs in the same NAT group.

7.19.7.5.2 Timestamp

Each flow creation or deletion event is timestamped individually using the local time in the system. The event timestamp (including the time-zone designator) is carried in the message part of the log. This timestamp is carried in addition to the syslog timestamp which is generated at the time of syslog message generation and carried in the <HEADER> part of the syslog messages.

7.19.7.5.3 Event aggregation

By default, flow logging events are transported to the collector as fast as they are generated. This does not imply that each event is transported individually, instead a few events can be still aggregated in a single message. However, this aggregation is not user controllable and it depends on the current condition in the system (events that are generated at approximately the same time).

To further optimize transport of logging events to the collector, the events can be aggregated in a controlled fashion. The flow logging events can be aggregated based on:

- expiry of a configurable timer
- transport message size (logs are collected until the size of the syslog message reaches the MTU size)

Whichever of the two conditions is met first triggers the generation of a syslog frame carrying multiple events. The separating character between the logs in a syslog message is "|" surrounded by a whitespace on each side.

```
<186>Jan 11 18:51:22 135.221.38.108 NOK: 47 NAT [MDA 1/1]: NAT44 SADD 2017 Jan 11 18:51:22:50
PST 6 0 10.10.10.1 3000 20 11.11.11.11 5000 12.12.12.12 8000 pol-name-1 sub-1 | NAT44 SADD 2017
Jan 11 18:51:22:60 PST 6 0 10.10.10.2 4000 20 11.11.11.11 6000 13.13.13.13 9000 pol-name-1 sub-1\n
```

7.19.7.5.4 Syslog transmission rate limit and overload conditions

The transmission rate of syslog messages can be limited by configuration. The rate limit is enforced in packets-per-second (pps). When the rate limit is exceeded, NAT flow logs are buffered. An overload condition is characterized by exhaustion of this buffer space. This condition can occur because of imposed rate limit or the software speed limit. Once the buffer space is exhausted, new flow creation is denied, and the teardown of the existing flows are delayed until the buffer space becomes available. Rate limit determines how fast the buffers are freed (by sending packets to the collector).

7.20 DS-Lite and NAT64 fragmentation

7.20.1 Overview

Fragmentation functionality is invoked when the size of a fragmentation eligible packet exceeds the size of the MTU of the egress interface or tunnel. Packets eligible for fragmentation are:

- IPv4 packets or fragments with the DF bit in the IPv4 header cleared. Fragmentation can be performed on any routing node between the source and the destination of the packet.
- IPv6 packets on the source node. Fragmentation of IPv6 packet on the transient routing nodes is not allowed.

The best practice is to avoid fragmentation in the network by ensuring adequate MTU size on the transient or source nodes. Drawbacks of the fragmentation are:

- increased processing and memory demands to the network nodes (especially during reassembly process)
- increased byte overhead
- increased latency

Fragmentation can be particularly deceiving in a tunneled environment whereby the tunnel encapsulation adds extra overhead to the original packet. This extra overhead could tip the size of the resulting packet over the egress MTU limit.

Fragmentation could be one solution in cases where the restriction in the mtu size on the packet's path from source to the destination cannot be avoided. Routers support IPv6 fragmentation in DS-Lite and NAT64 with some enriched capabilities, such as optional packet IPv6 fragmentation even in cases where DF-bit in corresponding IPv4 packet is set.

In general, the lengths of the fragments must be chosen such that resulting fragment packets fit within the MTU of the path to the packets destinations.

In downstream direction fragmentation can be implemented in two ways:

- IPv4 packet can be fragmented in the carrier IOM before it reaches ISA for any NAT function.
- IPv6 packet can be fragmented in the ISA, after the IPv4 packet is IPv6 encapsulated in DS-Lite or IPv6 translated in NAT64.

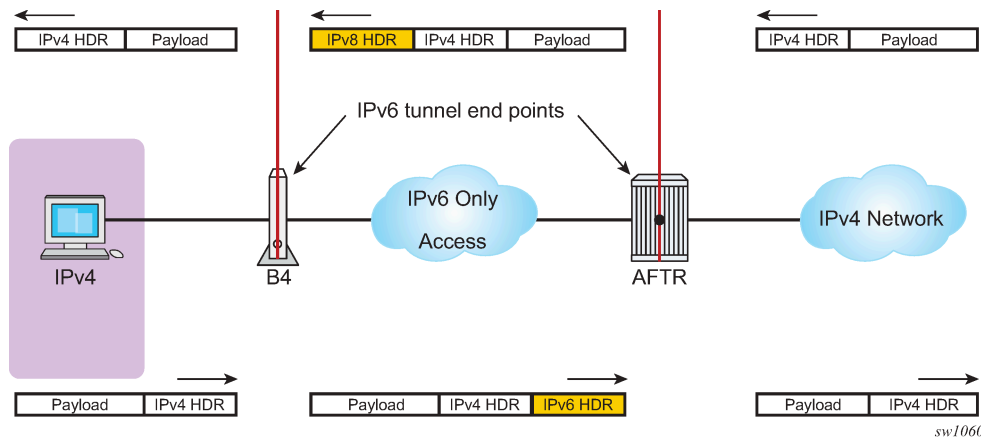
In upstream direction, IPv4 packets can be fragmented after they are decapsulated in DS-Lite or translated in NAT64. The fragmentation occurs in the IOM.

7.20.2 IPv6 fragmentation in DS-Lite

In the downstream direction, the IPv6 packet carrying IPv4 packet (IPv4-in-IPv6) is fragmented in the ISA in case the configured DS-Lite tunnel-mtu is smaller than the size of the IPv4 packet that is to be tunneled inside of the IPv6 packet. The maximum IPv6 fragment size is 48bytes larger than the value set by the tunnel-mtu. The additional 48 bytes is added by the IPv6 header fields: 40 bytes for the basic IPv6 header plus 8 bytes for extended IPv6 fragmentation header. NAT implementation in the routers does not insert any extension IPv6 headers other than fragmentation header.

[Figure 90: DS-Lite](#) shows DS-Lite IPv6 fragmentation.

Figure 90: DS-Lite



If the IPv4 packet is larger than the value set by the tunnel-mtu, the fragmentation action depends on the configuration options and the DF bit setting in the header of the received IPv4 header:

- The IPv4 packet can be dropped regardless of the DF bit setting. IPv6 fragmentation is disabled.
- The IPv4 packet can be encapsulated in IPv6 packet and then the IPv6 can be fragmented regardless of the DF bit setting in the IPv4 tunneled packet. The IPv6 fragment payload is limited to the value set by the tunnel-mtu.
- The IPv4 packet can be encapsulated in IPv6 packet and then the IPv6 can be fragmented only if the DF bit is cleared. The IPv6 fragment payload is limited to the value set by the tunnel-mtu.

If the IPv4 packet is dropped because of fragmentation not being allowed, an ICMPv4 Datagram Too Big message is returned to the source. This message carries the information about the size of the MTU that is supported, by notifying the source to reduce its MTU size to the requested value (tunnel-mtu).

The maximum number of supported fragments per IPv6 packet is 8. Considering that the minimum standard based size for IPv6 packet is 1280 bytes, 8 fragments is enough to cover jumbo Ethernet frames.

```
configure
  [router] | [service vprn]
    nat
      inside
        dual-stack-lite
        address <IPv6 Addr>
          tunnel-mtu bytes
          ip-fragmentation {disabled | fragment-ipv6 | fragment-ipv6-unless-ipv4-df-
set}
```

7.20.3 NAT64

Downstream fragmentation in NAT64 works in similar fashion. The difference between DS-Lite is that in NAT64 the configured ipv6-mtu represents the mtu size of the IPv6 packet (as opposed to payload of the IPv6 tunnel in DS-Lite). In addition, IPv4 packet in NAT64 is not tunneled but instead IPv4 or IPv6 headers are translated. Consequently, the fragmented IPv6 packet size is 28 bytes larger than the translated IPv4 packet 20 bytes difference in basic IP header sizes (40-bytes IPv6 header versus a 20-byte IPv4 header) plus 8 bytes for extended fragmentation IPv6 header. The only extended IPv6 header that NAT64 generates is the fragmentation header.

If the IPv4 packet is dropped because of the fragmentation not being allowed, the returned ICMP message contains MTU size of ipv6-mtu minus 28 bytes.

Otherwise the fragmentation options are the same as in DS-Lite.

```
configure
 [router] | [service vprn]
   nat
     inside
       nat64
         ipv6-mtu bytes
         ip-fragmentation {disabled | fragment-ipv6 | fragment-ipv6-unless-ipv4-df-set}
```

7.21 DS-Lite reassembly

In a tunneled environment such as DS-Lite, a fragmented packet must be reassembled in the end node before it is decapsulated. DS-Lite reassembly is implemented in-line, which means that the reassembly function runs in the same MS-ISA where native DS-Lite processing occurs. The presence of the Fragment Extension header in the IPv6 header signals the need for reassembly for all traffic destined for the AFTR in the upstream direction.

Fragments of a frame can be buffered for up to two seconds in an MS-ISA to wait for all the fragments of the original frame to arrive can be reassembled.

DS-Lite reassembly is performed only in the upstream direction.

7.21.1 Interpreting fragmentation statistics

Fragmentation statistics in DS-Lite can be observed by issuing the following command:

show isa nat-group <grp-id> mda <slot-id/mda-id> statistics

The command output displays only relevant DS-Lite fragmentation counters are shown. The remaining counters are removed from the output for easier reading.

```
show isa nat-group 1 mda 1/2 statistics
=====
ISA NAT Group 1 MDA 1/2
=====
--snip--
too many fragments for IP packet           : 0
too many fragmented packets                : 0
too many fragment holes                    : 0
too many frags buffered                    : 0
fragment list expired                      : 0
Reassembly Failures                       : 0
Fragments RX DSL                           : 0
Fragments RX DORMANT                       : 0
Fragments TX DSL                           : 0
Fragments TX DORMANT                       : 0
Fragments RX OUT                           : 0
Fragments TX OUT                           : 0
too many frag. lists for flow              : 0
frag. list cleanup in progress             : 0
--snip--
=====
```


To interpret these counters, familiarity with the following terms in the context of fragmentation is required:

- **packet**
An IP packet that is split into fragments (multiple frames) because of its original size being larger than the MTU configured on any node servicing this packet.
- **fragment**
A fragment comprises the frames that make up a packet. Multiple fragments (frames on the wire) are eventually reassembled into the original packet.
- **fragment list**
MS-ISA maintains a list of fragments belonging to a single packet. Each list represents a single fragmented packet and the list contains multiple fragments.
- **fragment hole**
A hole refers to a fragment or a group of consecutive fragments in a fragment list. Fragments of a packet are sequentially numbered from first to last and they must be reassembled in the same order in which they are fragmented. For example, if a packet contains 9 fragments but only fragments [1,3,5,9] are received by MS-ISA, then there are 3 holes in this list [2],[3,4],[6,7,8].
- **flow**
This is identified by 5 tuple <src IP, dst IP, src Port, dst Port, protocol>. Flows can have many packets and each packet of a flow can be fragmented.

[Table 53: Counter names and descriptions](#) describes counter names.

Table 53: Counter names and descriptions

Counter name	Description
too many fragments for IP packet	This counter increments if there are more than 20 fragments of a received single packet (the maximum number of fragments per packet is 20). In this case, all fragments of the packet are dropped.
too many fragmented packets	This counter increases if the maximum number of fragmented packets per MS-ISA is reached. See the MS-ISA Scaling Guides for the maximum number of fragmented packets per MS-ISA (specifically, the max num of frag lists parameter).
too many fragment holes	This counter increments if there are more than 11 holes tracked for a single packet. In this case, all fragments are dropped.
too many frags buffered	This counter increases if the maximum number of fragments that can be stored on MS-ISA is reached.
fragment list expired	This counter increases if all fragments of single packets are not received within two seconds. In this case, all fragments of this packet are dropped.
lists for flow	This counter increases when more than five fragmented packets per flow are being maintained simultaneously in

Counter name	Description
	MS-ISA. In this case, the fragments of the sixth packets are dropped.
Reassembly Failures	This counter increases when the reassembly of a packets (when all fragments are received) failed. This can attributed to the size of the first DS-Lite IPv6 fragment is smaller than 1280B, or the total reassembled packet is too big (greater than 9212B).
Fragments RX DSL	This counter increments only when a DS-Lite packet/fragment is received in the upstream direction that contains an IPv4 fragment. In other words, this counter is relevant only to IPv4 fragments inside of the DS-Lite packet/fragment that is received from the subscriber, and is not affected by DS-Lite fragments.
Fragments TX DSL	This counter increments only in case that DS-Lite packet/fragment sent in the downstream direction contains an encapsulated IPv4 fragment (which is received from the public side). In other words, this counter is relevant only to IPv4 fragments inside of the DS-Lite packet/fragment that is sent toward the subscriber, and is not affected by DS-Lite fragments.
Fragments RX OUT	This counter increments when an IPv4 fragment is received in the downstream direction, toward the subscriber.
Fragments TX OUT	This counter increments when an IPv4 fragment is transmitted in the upstream direction (public side).

7.21.2 Support for small first fragments

RFC 8200, *Internet Protocol, Version 6 (IPv6) Specification*, recommends the minimum MTU size in IPv6 should be 1280 bytes. However, some devices sourcing IPv6 traffic do not follow this recommendation and fragments packets to a size smaller than 1280 bytes. To accommodate these devices, the DS-Lite implementation in SR OS can process first fragments (with the fragment offset equaling 0 in the IPv6 fragmentation header) that are smaller than 1280 bytes. The SR OS can reassemble such packets in the upstream direction, as well as fragment them in the downstream direction.

7.21.2.1 Upstream reassembly with small first IPv6 fragments less than 1280 bytes

By default, DS-Lite implementation in SR OS drops the first fragment in the upstream direction if it is smaller than 1280 bytes.

For a router instance, use the following command to configure the minimum MTU size for the first fragment in the upstream direction and enable the processing of first IPv6 fragments smaller than 1280 bytes:

- **MD-CLI**

```
configure router nat inside large-scale dual-stack-lite endpoint min-first-fragment-size-rx
```

- **classic CLI**

```
configure router nat inside dual-stack-lite address min-first-fragment-size-rx
```

For a VPRN service, use the following command to configure the minimum MTU size for the first fragment in the upstream direction and enable the processing of first IPv6 fragments smaller than 1280 bytes:

- **MD-CLI**

```
configure service vprn nat inside large-scale dual-stack-lite endpoint min-first-fragment-size-rx
```

- **classic CLI**

```
configure service vprn nat inside dual-stack-lite address min-first-fragment-size-rx
```

7.21.2.2 Downstream fragmentation with small first IPv6 fragment

Use the following command to set the size of the frames in the downstream direction for a router instance:

- **MD-CLI**

```
configure router nat inside large-scale dual-stack-lite endpoint tunnel-mtu
```

- **classic CLI**

```
configure router nat inside dual-stack-lite address tunnel-mtu
```

Use the following command to set the size of the frames in the downstream direction for a VPRN service:

- **MD-CLI**

```
configure service vprn nat inside large-scale dual-stack-lite endpoint tunnel-mtu
```

- **classic CLI**

```
configure service vprn nat inside dual-stack-lite address tunnel-mtu
```

The **tunnel-mtu** value represents the size of the IPv4 payload which is encapsulated in the IPv6 packet with an additional 48-byte header. These IPv6 packets can then be fragmented when fragmentation is enabled.

7.22 Histogram

The distribution of the following resources in a NAT pool is tracked in the form of a histogram.

- **Ports and subscribers**

The distribution of outside ports in a NAT pool is tracked for an aggregate number of subscribers. The output of the **histogram** command can reveal the number of subscribers in a pool that are heavy port users, or it can reveal the average number of ports used by most subscribers

- **Port blocks and subscribers in a NAT pool**

The distribution of port blocks in an L2-Aware NAT pool is tracked for an aggregate number of subscribers. The output of the **histogram** command can reveal how subscribers are using port blocks in the aggregate.

- **Subscribers and IP addresses**

The distribution of subscribers across IP addresses is tracked. The output of the **histogram** command is used to determine if any substantial imbalances exist.

- **Extended port blocks and outside IP addresses in a NAT pool**

The distribution of extended port blocks in the NAT pool is tracked in relation to an aggregate number of outside IP addresses. The output of the **histogram** command can reveal how extended port blocks are distributed over IP addresses in an aggregate. This is applicable only for an L2-Aware NAT pool with extended port blocks enabled, or a deterministic LSN pool.

The operator can use the displayed information to adjust the port block size per subscriber, the amount of port blocks per subscriber, or see port usage trends over time. Consequently, the operator can adjust the configuration as the port demand per subscriber increases or decreases over time. For example, an operator may find that the port usage in a pool increased over a period of time. Accordingly, the operator can plan to increase the number of ports per port block.

Execute the following show commands to display the **histogram** output.

Use the following command to show route-table information

- **Ports and subscribers per NAT pool (L2-Aware or LSN)**

Use the following command to show ports and subscribers per NAT pool (L2-Aware or LSN)

```
show router nat pool <name> histogram ports - ports bucket-size <size> num-buckets <number>
```

Output example: The output is organized in port buckets with the number of subscribers in each bucket.

```
show router nat pool <name> histogram ports
- ports bucket-size <size> num-buckets <number>
<size>                : [1..65536]
<number>              : [2..50]
```

```
show router nat pool "pool-1" histogram ports bucket-size 200 num-buckets 10
=====
Usage histogram NAT pool "pool-1" router "Base"
=====
```

Num-ports	Sub-TCP	Sub-UDP	Sub-ICMP
1-199	17170	0	0

```

200-399      8707      0      0
400-599      2406      0      0
600-799      635       0      0
800-999      322       0      0
1000-1199    0           0      0
1200-1399    0           0      0
1400-1599    0           0      0
1600-1799    0           0      0
1800-        0           0      0
-----
No. of entries: 10
=====
    
```

- **Port blocks and subscribers per NAT pool (L2-Aware and LSN)**

Use the following command to show port blocks and subscribers per NAT pool (L2-Aware and LSN)

```
show router nat pool <name> histogram port-blocks
```

Output example: The output is organized by the increasing number of port blocks in a NAT pool with the number of subscribers using the number of port blocks indicated in each line.

```

show router nat pool <name> histogram port-blocks

show router nat pool "l2a" histogram port-blocks
=====
Usage histogram NAT pool "l2a" router "Base" port blocks per subscriber
=====
Num port-blocks      Num subscribers
-----
1                    17398
2                    8550
3                    2352
4                    940
5                    0
6                    0
7                    0
8                    0
9                    0
10                   0
-----
No. of entries: 10
=====
    
```

- **Subscribers and IP addresses per NAT pool (LSN)**

Use the following command to show subscribers and IP addresses per NAT pool (LSN)

```

show router nat pool demo histogram subscribers-per-ip
- subscribers-per-ip bucket-size <size> num-buckets <number>

<size>              : [1..65536]
<number>            : [2..100]
    
```

Output example: The output is organized in buckets where each bucket shows how the subscribers are spread over the preferred outside IP addresses. For example, the output of the below command

shows that each of the 513 IP addresses in the pool have 120 to 129 subscribers. This is a fairly even distribution of subscribers over IP addresses and the favorable output of this command.

```
*A:Dut-C# show router 5 nat pool "demo" histogram subscribers-per-ip bucket-size 10 num-buckets 50
```

```
=====
Usage histogram NAT pool "demo" router 5 subscribers per IP address
=====
```

```
Num subscribers      Num IP addresses
-----
```

1-9	0
10-19	0
20-29	0
30-39	0
40-49	0
50-59	0
60-69	0
70-79	0
80-89	0
90-99	0
100-109	0
110-119	0
120-129	513
130-139	0
140-149	0
150-159	0
160-169	0
170-179	0
180-189	0
190-199	0
200-209	0
210-219	0
220-229	0
230-239	0
240-249	0
250-259	0
260-269	0
270-279	0
280-289	0
290-299	0
300-309	0
310-319	0
320-329	0
330-339	0
340-349	0
350-359	0
360-369	0
370-379	0
380-389	0
390-399	0
400-409	0
410-419	0
420-429	0
430-439	0
440-449	0
450-459	0
460-469	0
470-479	0
480-489	0
490-	0

- **Extended port blocks in a NAT pool and outside IP addresses (L2-Aware NAT and deterministic LSN)**

Use the following command to show subscribers and IP addresses per NAT pool (LSN).

```
show router nat pool demo histogram extended-port-blocks-per-ip
- extended-port-blocks-per-ip bucket-size <size> num-buckets <number>
<size>                : [1..65536]
<number>              : [2..50]
```

Output example: The output is organized in extended port-block buckets in a NAT pool with the number of outside IP addresses in each bucket.

```
show router nat pool demo histogram extended-port-blocks-per-ip
- extended-port-blocks-per-ip bucket-size <size> num-buckets <number>
<size>                : [1..65536]
<number>              : [2..50]
show router nat pool "l2a" histogram extended-port-blocks-per-ip bucket-size 1 num-buckets
10
=====
Usage histogram NAT pool "l2a" router "Base" extended port blocks per IP address
=====
Num extended-port-blocks      Num IP addresses
-----
-
1-1                          1039
2-2                          6182
3-3                          777
4-4                          194
5-5                          0
6-6                          0
7-7                          0
8-8                          0
9-                            0
-----
No. of entries: 10
=====
```

The output of each command can be periodically exported to an external destination with the **cron** command. The following displays an example of the output:

```
/configure system
script-control {
  script "nat_histogram" owner "TiMOSCLI" {
    admin-state enable
    location "ftp://*:*@138.203.8.62/nat-histogram.txt"
  }
  script-policy "dump_nat_histogram" owner "TiMOS CLI" {
    admin-state enable
    results "ftp://*:*@138.203.8.62/nat_histogram_results.txt"
    script {
      name "nat_histogram"
    }
  }
}
cron {
  schedule "nat_histogram_schedule" owner "TiMOS CLI" {
    admin-state enable
    interval 600
    script-policy {
      name "dump_nat_histogram"
    }
  }
}
```



```
}

```

The `nat-histogram.txt` file contains the command execution line. For example:

```
show router nat pool "pool-1" histogram ports bucket-size 200 num-buckets 10

```

This command is executed every 10 minutes (600 seconds) and the output of the command is written into a set of files on an external TFTP server as displayed in the following example.

```
[root@ftp]# ls nat_histogram_results.txt*
nat_histogram_results.txt_20130117-153548.out
nat_histogram_results.txt_20130117-153648.out
nat_histogram_results.txt_20130117-153748.out
nat_histogram_results.txt_20130117-153848.out
nat_histogram_results.txt_20130117-153948.out
nat_histogram_results.txt_20130117-154048.out
[root@ftp]#

```

7.23 NAT redundancy

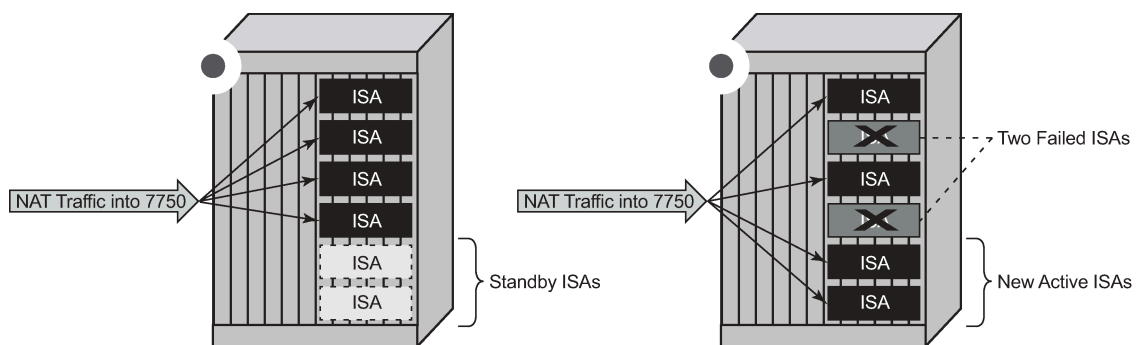
NAT ISA redundancy helps protect against Integrated Service Adapter (ISA) failures. This protection mechanism relies on the CPM maintaining configuration copy of each ISA. In case that an ISA fails, the CPM restores the NAT configuration from the failed ISA to the remaining ISAs in the system. NAT configuration copy of each ISA, as maintained by CPM, is concerned with configuration of outside IP address and port forwards on each ISA. However, CPM does not maintain the state of dynamically created translations on each ISA. This causes interruption in traffic until the translation are re-initiated by the devices behind the NAT.

Two modes of operation are supported:

- **Active-Standby**

In this mode of operation, any number of standby ISAs can be allocated for protection purposes. When there are no failures in the router, standby ISAs are idle, in a state ready to accept traffic from failed ISA. Mapping between the failed ISA and the standby ISA is always 1:1. This means that one standby ISA entirely replaces one failed ISA. In this respect, NAT bandwidth from the failed ISA is reserved and restored upon failure. This model is shown in [Figure 91: Active-Standby intra-chassis redundancy model](#).

Figure 91: Active-Standby intra-chassis redundancy model



al_0789

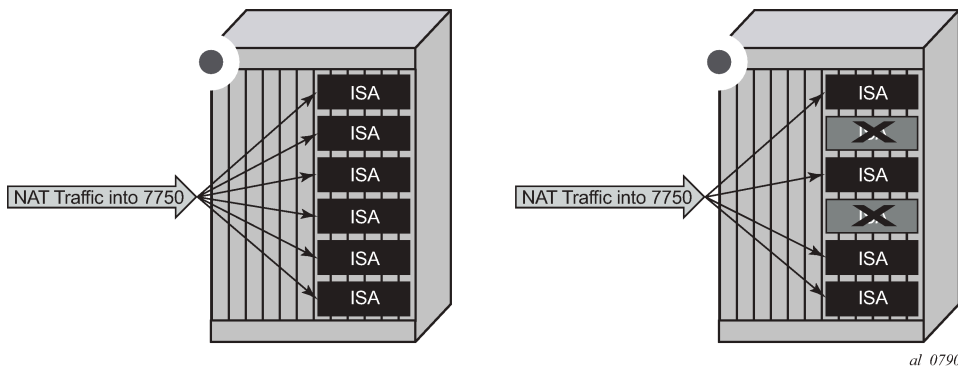
- **Active-Active**

In this mode all ISAs in the system are active. When an ISA fails, its load is distributed across the remaining active ISA. In this mode of operation there is no bandwidth reservation across active ISA. Each ISA can operate at full speed at any time. However, memory resources necessary to setup new translations from the failed ISAs are reserved. The reserved resources are:

- subscribers (inside IPv4 addresses for LSN44, IPv6 prefixes for DS-Lite/NAT64 and L2-Aware subscriber)
- outside IPv4 addresses
- outside port-ranges

By reserving memory resources it can be assured that failed traffic can be recovered by remaining ISAs, potentially with some bandwidth reduction in case that remaining ISAs operated at full or close to full speed before the failure occurred. Active-active ISA redundancy model is shown in [Figure 92: Active-Active intra-chassis redundancy model](#).

Figure 92: Active-Active intra-chassis redundancy model



In case of an ISA failure, the member-id of the member ISA that failed is contained in the FREE log. This info is used to find the corresponding MAP log which also contains the member-id field.

In case of RADIUS logging, CPM summarization trap is generated (because RADIUS log is sent from the ISA – which is failed).

7.23.1 NAT stateless dual-homing

Multi-chassis stateless NAT redundancy is based on a switchover of the NAT pool that can assume active (master) or standby state. The inside/outside routes that attract traffic to the NAT pool are always advertised from the active node (the node on which the pool is active).

This dual-homed redundancy based on the pool mastership state works well in scenarios where each inside routing context is configured with a single NAT policy (NAT'd traffic within this inside routing context is mapped to a single NAT pool).

However, in cases where the inside traffic is mapped to multiple pools (with deterministic NAT and in case when multiple NAT policies are configured per inside routing context), the basic per pool multi-chassis redundancy mode can cause the inside traffic within the same routing instance to fail because some pools referenced from the routing instance may be active on one node while other pools may be active on the other node.

Imagine a case where traffic ingressing the same inside routing instance is mapped as follows (this mapping can be achieved via filters):

- Source ip-address A → Pool 1 (nat-policy 1) active on Node 1
- Source ip-address B → Pool 2 (nat-policy 2) active on Node 2

Traffic for the same destination is normally attracted only to one NAT node (the destination route is advertised only from a single NAT node). Assume that this node is Node 1 in the example. After the traffic arrives to the NAT node, it is mapped to the corresponding pool according to the mapping criteria (routing based or filter based). But if active pools are not co-located, traffic destined for the pool that is active on the neighboring node would fail. In our example traffic from the source ip-address B would arrive to the Node 1, while the corresponding Pool 2 is inactive on that node. Consequently the traffic forwarding would fail.

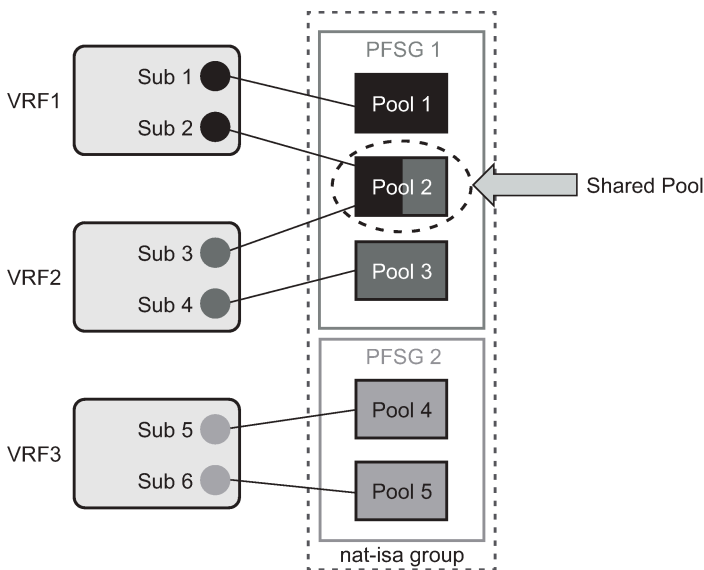
To remedy this situation, a group of pools targeted from the same inside routing context must be active on the same node simultaneously. In other words, the active pools referenced from the same inside routing instance must be co-located. This group of pools is referred to as Pool Fate Sharing Group (PFSG). The PFSG is defined as a group of all NAT pools referenced by inside routing contexts whereby at least one of those pools is shared by those inside routing contexts. This is shown in [Figure 91: Active-Standby intra-chassis redundancy model](#).

Even though only Pool 2 is shared between subscribers in VRF 1 and VRF 2, the remaining pools in VRF 1 and VRF 2 must be made part of PFSG 1 as well.

This ensures that the inside traffic is always mapped to pools that are active in a single box.

[Figure 93: Pool fate sharing group](#) shows the pool fate sharing group.

Figure 93: Pool fate sharing group



al_0409

There is always one lead pool in PFSG. The Lead pool is the only pool that is exporting/monitoring routes. Other pools in the PFSG are referencing the lead pool and they inherit its (activity) state. If any of the pools in PFSG fails, all the pools in the PFSG switch the activity, or in another words they share the fate of the lead pool (active/standby/disabled).

There is one lead pool per PFSG per node in a dual-homed environment. Each lead pool in a PFSG has its own export route that must match the monitoring route of the lead pool in the corresponding PFSG on the peering node.

PFSG is implicitly enabled by configuring multiple pools to follow the same lead pool.

7.23.1.1 Configuration considerations

Attracting traffic to the active NAT node (from inside and outside) is based on the routing.

On the outside, the active pool address range is advertised. On the inside, the destination prefix or steering route (in case of filter based diversion to the NAT function) is advertised by the node with the active pool.

The advertisement of the routes is driven by the activity of the pools in the pool fate sharing group:

```
configure
  router/service vprn
    nat
      outside
        pool <name>
          redundancy
            export <ip-prefix/length>
            monitor <ip-prefix/length>[no] shutdown
            follow router <rtr-id> pool <master-pool>
```

For example:

```
router/service vprn
  nat
    outside
      pool "nat0-pool" nat-group 1 type large-scale create
      port-reservation ports 252
      redundancy
        follow router 500 pool "nat500-pool"
      exit
      address-range 192.168.12.0 192.168.12.10 create
      exit
      no shutdown
    exit
  exit
exit
```

A pool can be one of the following:

- a leading pool (configure export- and monitor-route and put in no shutdown)
- A following pool (configure follow)

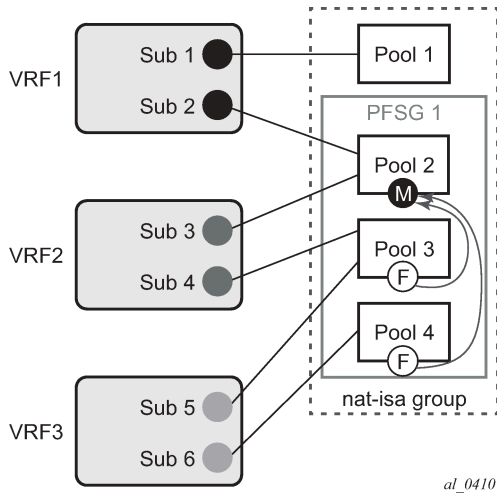
Both sets of options are therefore mutually exclusive.

A leading pool redundancy is only enabled when the redundancy node is in no shutdown. For a following pool, the administrate has no effect, and the redundancy is only enabled when the leading pool is enabled.

Before a lead pool is enabled, consistency checks are performed to make sure that PFSG is properly configured and that the all pools in the PFSG belong to the same NAT ISA group. PFSG is implicitly enabled by configuring multiple pools to follow the same lead pool. To ensure the effective functioning of a PFSG, it is essential that all participating NAT pools are enabled. Modifications to the PFSG, such as adding or removing pools, can only be executed when the primary pool is temporarily disabled.

For example in the case shown in [Figure 94: Consistency check](#), the consistency check would fail because pool 1 is not part of the PFSG 1 (where it should be).

Figure 94: Consistency check



7.23.1.2 Troubleshooting commands

The following command displays the state of the leading pool (dual-homing section toward the bottom of the command output):

```
*A:Dut-B# show router 500 nat pool "nat500-pool"
=====
NAT Pool nat500-pool
=====
Description                : (Not Specified)
ISA NAT Group               : 1
Pool type                   : largeScale
Admin state                 : inService
Mode                        : auto (napt)
Port forwarding dyn blocks reserved : 0
Port forwarding range       : 1 - 1023
Port reservation           : 2300 blocks
Block usage High Watermark (%) : (Not Specified)
Block usage Low Watermark (%)  : (Not Specified)
Subscriber limit per IP address : 65535
Active                      : true
Deterministic port reservation : (Not Specified)
Last Mgmt Change            : 02/17/2014 09:41:43
=====
NAT address ranges of pool nat500-pool
=====
Range                        Drain Num-blk
-----
192.168.1.0 - 192.168.1.255 0
-----
No. of ranges: 1
=====
```

```

NAT members of pool nat500-pool ISA NAT group 1
=====
Member                               Block-Usage-% Hi
-----
1                                     < 1          N
2                                     < 1          N
3                                     < 1          N
4                                     < 1          N
5                                     < 1          N
6                                     < 1          N
-----
No. of members: 6
=====
Dual-Homing
=====
Type                                 : Leader
Export route                         : 10.0.0.3/32
Monitor route                        : 10.0.0.2/32
Admin state                           : inService
Dual-Homing State                    : Active
=====
Dual-Homing fate-share-group
=====
Router      Pool                      Type
-----
Base        nat0-pool                 Follower
vprn500     nat500-pool                Leader
vprn501     nat501-pool                 Follower
vprn502     nat502-pool                 Follower
-----
No. of pools: 4
=====

```

The following command displays the state of the follower pool (dual-homing section toward the bottom of the command output):

```

*A:Dut-B# show router 501 nat pool "nat501-pool"
=====
NAT Pool nat501-pool
=====
Description                               : (Not Specified)
ISA NAT Group                             : 1
Pool type                                  : largeScale
Admin state                                : inService
Mode                                        : auto (napt)
Port forwarding dyn blocks reserved        : 0
Port forwarding range                      : 1 - 1023
Port reservation                           : 2300 blocks
Block usage High Watermark (%)             : (Not Specified)
Block usage Low Watermark (%)              : (Not Specified)
Subscriber limit per IP address            : 65535
Active                                     : true
Deterministic port reservation             : (Not Specified)
Last Mgmt Change                           : 02/17/2014 09:41:43
=====
NAT address ranges of pool nat501-pool
=====
Range                                     Drain Num-blk
-----

```

```

192.168.2.0 - 192.168.2.255          0
192.168.3.0 - 192.168.3.255          0
-----
No. of ranges: 2
=====
NAT members of pool nat501-pool ISA NAT group 1
=====
Member                                Block-Usage-% Hi
-----
1                                     < 1           N
2                                     < 1           N
3                                     < 1           N
4                                     < 1           N
5                                     < 1           N
6                                     < 1           N
-----
No. of members: 6
=====
Dual-Homing
=====
Type                                  : Follower
Follow-pool                            : "nat500-pool" router 500
Dual-Homing State                       : Active
=====
Dual-Homing fate-share-group
=====
Router      Pool                                Type
-----
Base        nat0-pool                            Follower
vprn500     nat500-pool                            Leader
vprn501     nat501-pool                            Follower
vprn502     nat502-pool                            Follower
-----
No. of pools: 4
=====

```

The following command lists all the pools that are configured along with the NAT inside/outside routing context.

```

*A:Dut-B# show service nat overview
=====
NAT overview
=====
Inside/  Policy/  Type
Outside  Pool
-----
vprn550  lsn-policy_unused  default
Base     nat0-pool

vprn550  lsn-policy_nat1    destination prefix
vprn500  nat500-pool

vprn550  lsn-policy_nat2    destination prefix
vprn501  nat501-pool

vprn551  lsn-policy_unused  default
Base     nat0-pool

vprn551  lsn-policy_nat3    destination prefix

```

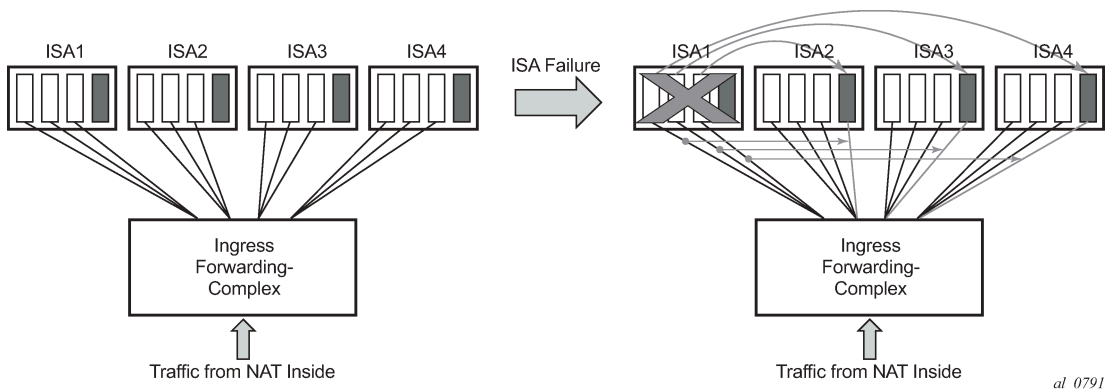
vprn501	nat501-pool	
vprn551	lsn-policy-nat4	destination prefix
vprn502	nat502-pool	
vprn552	lsn-policy_unused	default
Base	nat0-pool	
vprn552	lsn-policy-nat5	destination prefix
vprn502	nat502-pool	
=====		

7.23.2 Active-active ISA redundancy model

In active-active ISA redundancy, each ISA is subdivided into multiple logical ISAs. These logical sub-entities are referred to as members. NAT configuration of each member is saved in the CPM. In case that any one ISA fails, its members are downloaded by the CPM to the remaining active ISAs. Memory resources on each ISA are reserved to accommodate additional traffic from the failed ISAs. The amount of resources reserved per ISA depends on the number of ISAs in the system and the number of simultaneously supported ISA failures. The number of simultaneous ISA failures per system is configurable. Memory reservation affects NAT scale per ISA.

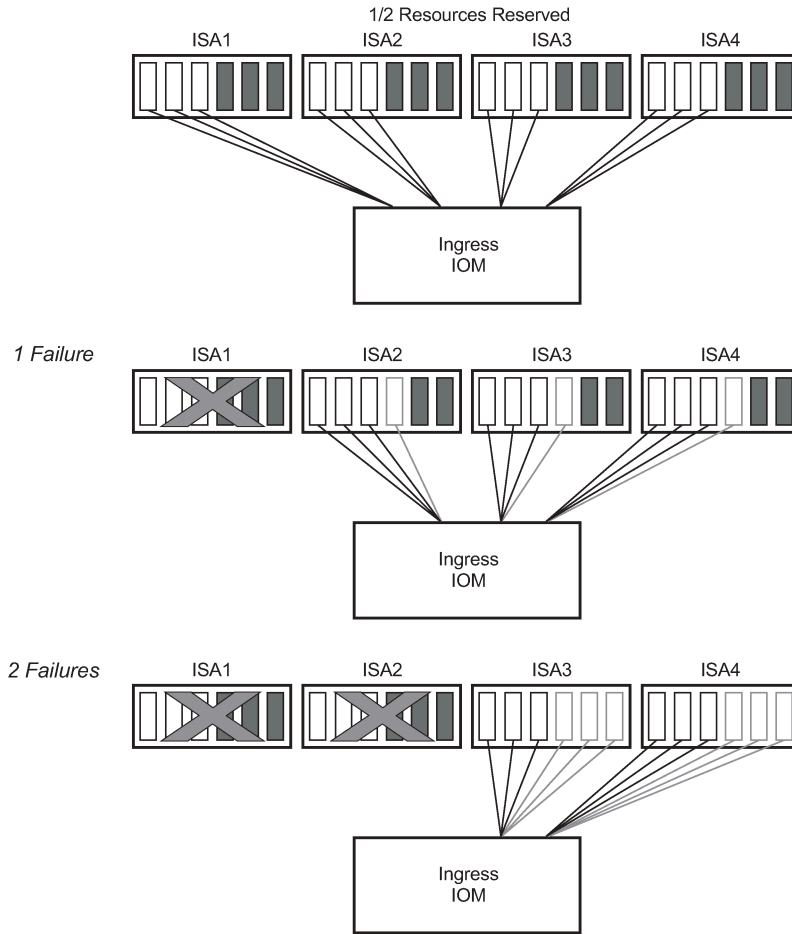
Traffic received on the inside is forwarded by the ingress forwarding complex to a predetermined member ISAs for further NAT processing. Each ingress forwarding complex maintains an internal link per member. The number of these internal links, along with other factors, determine the maximum number of members per system and with this, the granularity of traffic distribution over remaining ISAs in case of an ISA failure. The segmentation of ISAs into members for a single failure scenario is shown in [Figure 95: Load distribution in active-active intra-chassis redundancy model](#). The protection mechanism in this example is designed to cover for one physical ISA failure. Each ISA is divided into four members. Three of those carry traffic during normal operation, while the fourth one has resources reserved to accommodate traffic from one of the members in case of failure. When an ISA failure occurs, three members are delegated to the remaining ISAs. Each member from the failed ISA is mapped to a corresponding reserved member on the remaining ISAs.

Figure 95: Load distribution in active-active intra-chassis redundancy model



Active-active ISA redundancy model supports multiple failures simultaneously. The protection mechanism shown in [Figure 96: Multiple failures](#) is designed to protect against two simultaneous ISA failures. As the previous case, each ISA is divided into six members, three of which are carrying traffic under normal circumstances while the remaining three members have reserved memory resources.

Figure 96: Multiple failures



al_0792

Table 54: Load distribution in active-active ISA redundancy model supporting single ISA failure shows resource utilization for a single ISA failure in relation to the total number of ISAs in the system. The resource utilization affects only scale of each ISA. However, bandwidth per ISA is not reserved and each ISA can operate at full speed at any time (with or without failures).

Table 54: Load distribution in active-active ISA redundancy model supporting single ISA failure

Number of physical ISAs per system	Number of member ISAs per physical ISA (active/reserved)	Resource utilization per system in non-failed condition	Resource utilization per system with one failed ISA
2	1A 1R	50%	100%
3	2A 1R	67%	100%
4	3A 1R	75%	100%
5	3A 1R	75%	95%

Number of physical ISAs per system	Number of member ISAs per physical ISA (active/reserved)	Resource utilization per system in non-failed condition	Resource utilization per system with one failed ISA
6	2A 1R	66%	83%
7	2A 1R	66%	80%
8	2A 1R	66%	79%
9	1A 1R	50%	61%
10	1A 1R	50%	60%
11	1A 1R	50%	59%
12	1A 1R	50%	58%
13	1A 1R	50%	58%
14	1A 1R	50%	57%

7.23.2.1 Startup conditions

During the first five minutes of system startup or nat-group activation, the system behaves as if all ISAs are operational. Consequently, ISAs are segmented in its members according to the configured maximum number of supported failures.

Upon expiration of this initial five minute interval, the system is re-evaluated. In case that one or more ISAs are found in faulty state during re-evaluation, the members of the failed ISAs are distributed to the remaining ISAs that are operational.

7.23.2.2 Recovery

After a failed ISA is recovered, the system automatically accepts it and traffic is assigned to it. Traffic that is moved to the recovered ISA is interrupted.

7.23.2.3 Adding additional ISAs in the ISA group

Adding additional ISAs in an operational nat-group requires reconfiguration of the active mda-limit for the nat-group (or the failed mda-limit for that matter). This is only possible when nat-group is in an administratively shutdown state.

7.23.3 L2-Aware bypass

L2-Aware bypass provides the basis for traffic continuity if an MS-ISA fails. With L2-Aware bypass functionality disabled and without an intra-chassis redundancy scheme deployed (such as active/active or active/standby), the traffic to be processed by the failed MS-ISA is blackholed. This means that traffic continues to be sent to the failed MS-ISA. By enabling L2-Aware bypass, instead of being blackholed,

the traffic is routed outside of the SR OS node without being NAT'd in accordance to the routing table in the inside routing context. The intent is that non-NAT'd traffic is intercepted by a central NAT node that performs the NAT function. This way, traffic served by the failed MS-ISA continues to be NAT'd by a central NAT node. The central NAT node provides redundancy for multiple SR OS nodes, therefore reducing the need to equip each individual SR OS node with multiple MS-ISAs which are normally used in an active/active or active/standby intra-chassis redundancy mode.

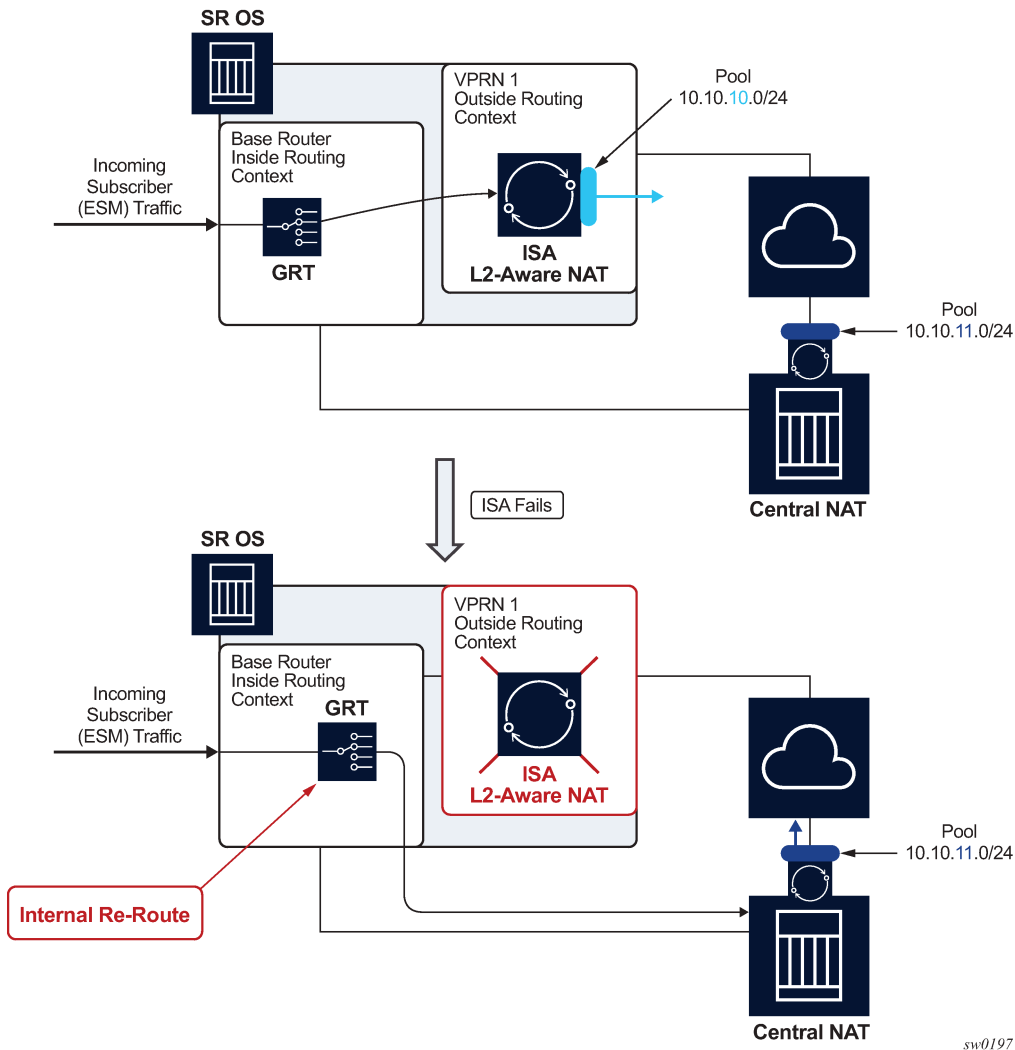
This concept is shown in [Figure 97: L2-Aware bypass](#) . The example shows the base router as an inside routing context where the global routing table (GRT) is used to decide where to send traffic if an MS-ISA is unavailable. Apart from this example, the inside routing context is not limited to the base router but instead can be an VPRN instance.

L2-Aware bypass is considered as an optional redundancy model in L2-Aware NAT which is mutually exclusive with the other two MS-ISA redundancy modes (active/active and active/standby).

L2-Aware bypass is enabled with the following CLI:

```
configure
isa
  nat-group <id>
  redundancy {active-active|active-standby|l2aware-bypass}
```

Figure 97: L2-Aware bypass



sw0197

7.23.3.1 Sharing IP addresses in L2-Aware NAT

L2-Aware NAT allows overlap of inside (private) IP addresses between Enhanced Subscriber Management (ESM) (or L2-Aware NAT) subscribers. For example, IP addresses assigned to hosts within, for example, subscriber SUB-1, can be identical to IP addresses assigned to hosts within, for example, subscriber SUB-2. This is possible by the subscriber-ID field (which must be unique in the system) that is a part of the NAT translation key. This way the return traffic (in downstream direction) belonging to different ESM subscribers with overlapping IP addresses can still be differentiated by a unique ESM subscriber-id field that is used in reverse NAT translation.

L2-Aware bypass functionality with a failed MS-ISA breaks the logic because traffic is not translated (NAT'd) in SR OS node, and therefore, the return traffic does not take subscriber-id field into forwarding consideration. For this reason, the overlap of inside (private) IP addresses between ESM subscribers is not supported by the L2-Aware bypass functionality for the routed traffic within the same inside routing

context. In other words, private IP addresses must be unique across the subscribers within a specified inside routing context.

7.23.3.2 Recovery

Upon the recovery of the failed MS-ISA, all existing subscribers that are affected by the bypass continue to use the bypass. However, all new subscribers that come online after the recovery, are automatically L2-Aware NAT'd (and therefore do not use the bypass).

Restoring bypassed subscribers to the L2-Aware NAT after the recovery, requires manual intervention by the operator. This is accomplished by executing the **tools perform nat recover-l2aw-bypass mda mda-id** command:

In L2-Aware NAT, the <subscriber to outside IP, port-block> mappings are allocated during the subscriber attachment phase (when the subscriber comes online) and are maintained in the CPM. Therefore, they are preserved in the CPM during MS-ISA failure. This means that the original mappings for the recovered subscribers continue to be used when the MS-ISA is recovered.

Be aware that only the partial mappings <subscriber to outside IP, port-block> are preserved. This does not include preservation of NAT sessions sometimes referred as fully qualified flows. NAT sessions are maintained in MS-ISA and they are lost during MS-ISA failures. Hence, this model provides stateless failover to an external NAT node.

The operator is notified about the MS-ISA failure by a log message or an SNMP trap. An example of such log is:

```
9 2017/06/07 11:32:49.748 UTC MINOR: NAT #2020 Base NAT
"The NAT MDA 5/1 is now inactive in group 2."
```

7.23.3.3 Default bypass during reboot or MS-ISA provisioning

If enabled, L2-Aware bypass takes effect automatically if an MS-ISA does not become operational within 10 minutes of provisioning (configuring) or after a system bootup.

7.23.3.4 Logging

Because partial mappings in L2-Aware NAT <subscriber to outside IP addresses, port block> are preserved in the CPM during an MS-ISA failure, no logging is generated for existing ESM/NAT subscribers when the MS-ISA fails or is recovered.

7.23.4 Stateful inter-chassis NAT redundancy

Stateful inter-chassis NAT redundancy provides seamless NAT failover between the two redundant SR OS nodes. A pair of redundant nodes operates in active or standby modes per NAT group. If traffic distribution between the nodes is needed, then up to four NAT groups per node can be deployed, and with each NAT group having its own set of ISAs. In this scenario, traffic between the nodes can be load-balanced per NAT group.

[Figure 98: CGN stateful inter-chassis redundancy](#) shows a scenario where inside routes are advertised from the node with the active NAT group that ensures traffic is symmetric. This means that upstream and

downstream traffic is fully flowing through the same node. Although this scenario represents the majority of use cases, it is allowed for the upstream traffic to arrive on the node with a standby NAT group and be shunted over to the node with active nat-group over a link that interconnects the two nodes. This scenario is shown in [Figure 99: Asymmetric traffic](#).

The redundant pair of NAT nodes protect against the following:

- node failure
- ISA failure
- link failure
- path failure (BFD, VRRP)

Figure 98: CGN stateful inter-chassis redundancy

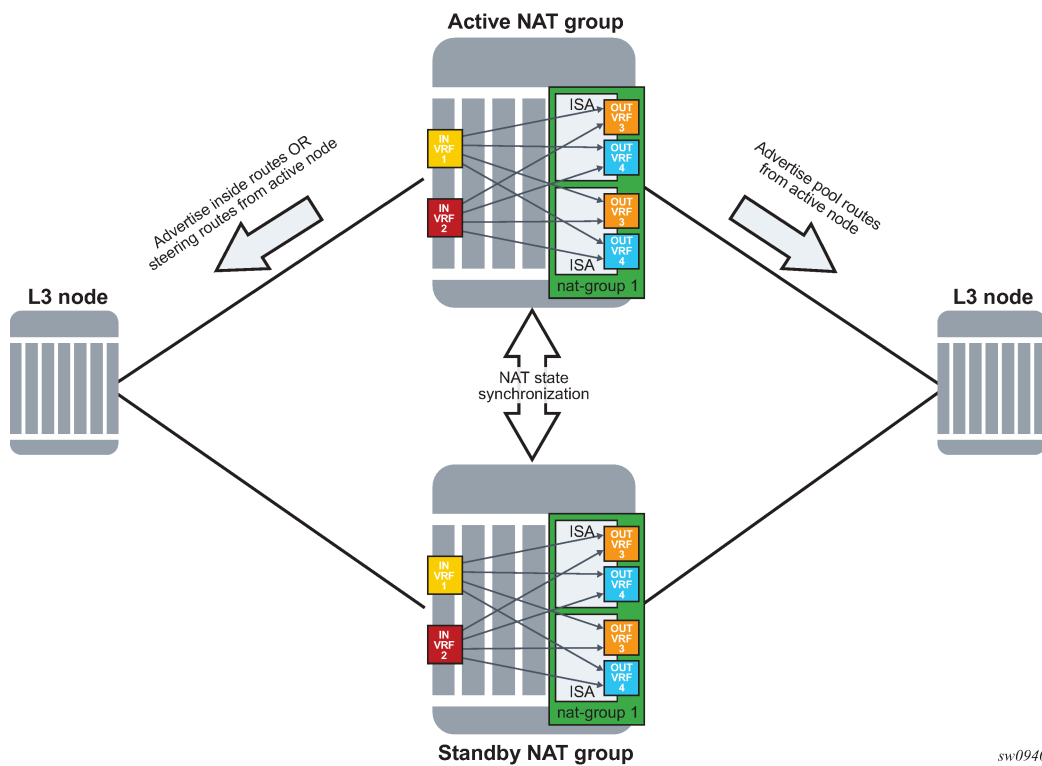
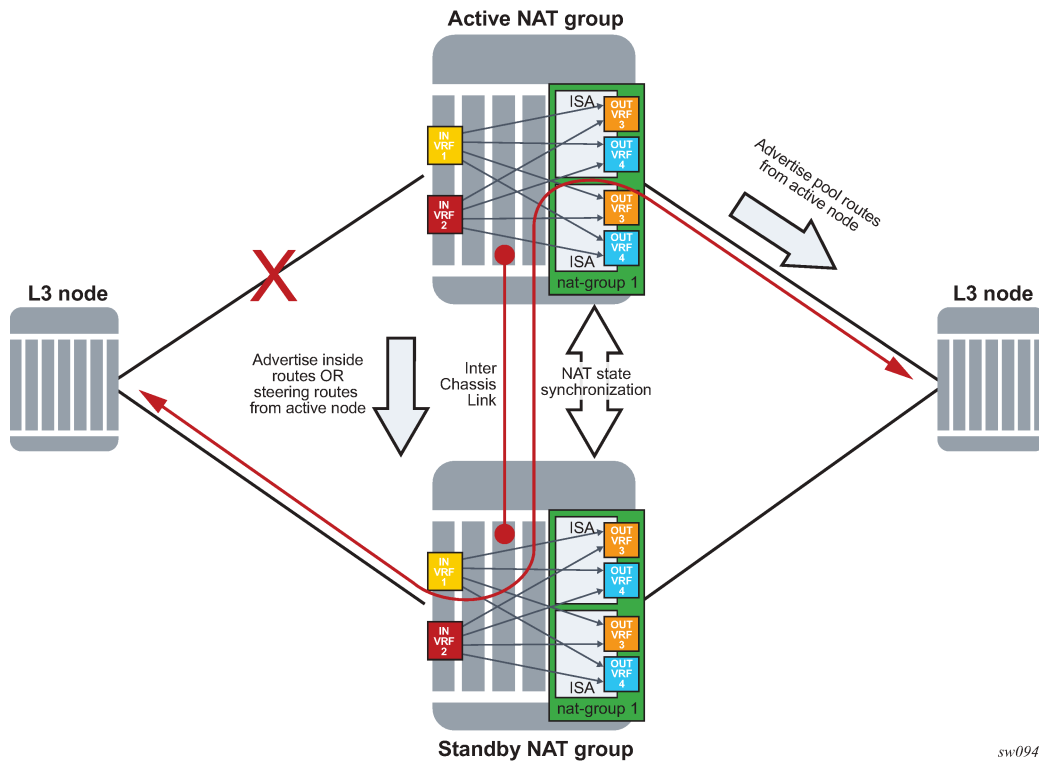


Figure 99: Asymmetric traffic

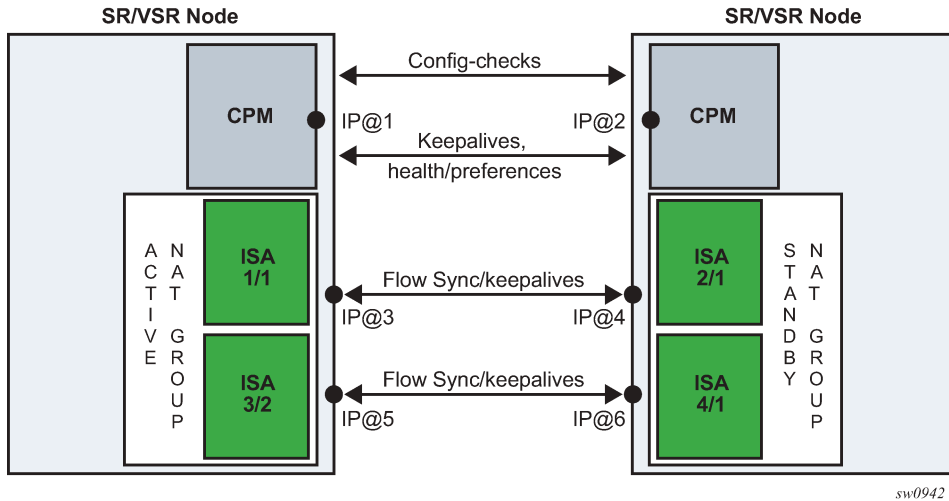


sw0941

The basic premise of stateful inter-chassis NAT redundancy is as follows:

- ISAs/ESAs within a NAT group can be either active or standby. This means that in a pair of redundant NAT nodes, only one node attracts traffic per NAT group while the NAT group on the peering node is in a standby state with synchronized flows (or sessions).
- The activity (active or standby) of a NAT group is determined by an internal health parameter that represents the node's ability to perform NAT at full capacity. The value of the health parameter can change dynamically and is based on the events that are related to various failures against which the stateful inter-chassis NAT redundancy offers protection. The health value is communicated between the redundant (or peering) nodes on the CPM level.
- In case of equal health, an operator can, through configuration, influence the activity of a NAT group.
- The activity of a NAT group is characterized by the advertisement of the NAT- related routes on the inside (steering and destination-prefix) and the outside (pool routes).
- Only TCP/UDP/ICMP flows that are older than the preconfigured amount of time are synchronized.
- Flow synchronization is performed directly between the ISAs, bypassing the CPM. The CPM's main role is to determine the activity of a NAT group with health related parameters as shown in [Figure 100: NAT synchronization](#).
- Flows are created on the active NAT group and are synchronized from the active NAT group to the standby NAT group.
- Immediately following the switchover, all flows on the newly standby NAT group are resynchronized from the newly active NAT group.
- If the link between the chassis is lost, then each chassis operates in a standalone mode.

Figure 100: NAT synchronization



A reliable and redundant link should always be available between the two redundant NAT nodes. This link is referred to as Inter-Chassis Link (ICL) and is used for:

- CPM communication (activity negotiation and presence detection)
- ISA-to-ISA communication (flow synchronization)
- transient forwarding of data traffic during switchovers

7.23.4.1 Health status and failure events

A health value determines the activity of a NAT group within a pair of redundant nodes. The health value of a NAT group is internally calculated. The system can automatically decrease this value depending on the events that can negatively affect the system’s ability to perform NAT at a needed capacity.

A NAT group with a higher health value becomes active.

[Table 55: Activity states at equal health](#) shows activity states, if paired NAT groups have equal health values on both nodes. Preferred is a configuration parameter that influences the activity state for a pair of NAT groups with equal health value (typical use case would be load balancing per NAT group).

Table 55: Activity states at equal health

Node 1	Node 2	Active node	Comments
no preferred configured	no preferred configured	Whichever node becomes active first, remains active	If both nodes are becoming active simultaneously, the node with the highest system chassis MAC address becomes a controller node that decides which node becomes active node and which standby based on the health and preference values.

Node 1	Node 2	Active node	Comments
			When the health and preference are equal, the controller node does not preempt (trigger a switchover) an already active node.
preferred configured	no preferred configured	Node 1	Node 1 always preempts Node 2 (if the health values are equal)
no preferred configured	preferred configured	Node 2	Node 2 always preempts Node 1 (if the health values are equal)
preferred configured	preferred configured	Whichever node becomes active first, remains active	Same as for no preference on both nodes

The *health* parameter is initially set to a value of 1000 under the following circumstances:

- The number of active ISAs in a NAT group is matching the configured value for the **active-mda-limit** configuration parameter.
- There are no port failures that are being monitored.
- There are no failures within the operation groups that are being monitored.

The above circumstances imply that the system is fully operational with no failures that would affect NAT operation.

However, the health value can be influenced by the events that can affect NAT operation, and that are outside of ISA-related failures, for example, unhealthy ports and paths that lead traffic in and out of the NAT node. Such events are explicitly tracked or monitored for the purpose of dynamically adjusting the health value and therefore influencing the activity of the NAT groups.

Stateful inter-chassis NAT redundancy protects against the following failures:

- Nodal failure; if the active node fails, the standby node is notified of such an event by the lack of received keepalives.
- ISA failure; a NAT group must have exactly the **active-mda-limit** number of ISAs that are operationally up to participate in stateful inter-chassis NAT redundancy. If the number of operational ISAs falls below the configured limit, then the health of the NAT group drops to 0.
- Ports on the node can be monitored, and their operational state can trigger a change of the health value.
- BFD sessions on the node can be monitored using the oper-groups and their state change can trigger a change of the health value.
- VRRP instances under interface configurations can be monitored using the oper-groups and their state change can trigger a change of the health value.
- SAPs on the node can be monitored using the oper-groups and their operational state can trigger a change of the health value.

Port and oper-group state change influences the reachability of the NAT node and consequently this affects network-wide NAT operation. If that port or path capacity in and out of the NAT node drops below a specific level, a switchover to a healthier NAT node may be needed.

Port states can be tracked or monitored on the private side (inside) and on the public side (outside) of NAT.

Oper-groups are constructs that are tracking states of BFD enabled interfaces, SAPs, and VRRP instances.

BFD sessions targeted to the next hop can traverse intermediate Layer 2 nodes and can have longer reach than port tracking.

Another benefit of monitoring ports and paths is that it can help reduce the amount of traffic on the inter-chassis communication link (ICL) if that active node loses direct connection to the node downstream or upstream from it. The link for inter-chassis control communication (ICL) must always be present (for synchronization purposes). However, this link does not need to be designed for heavy traffic loads during extended periods of time occurs if traffic bearing ports are not colocated with the active node. However, this link is used for shorter transient periods that are caused by switchovers.

7.23.4.2 Route advertisements

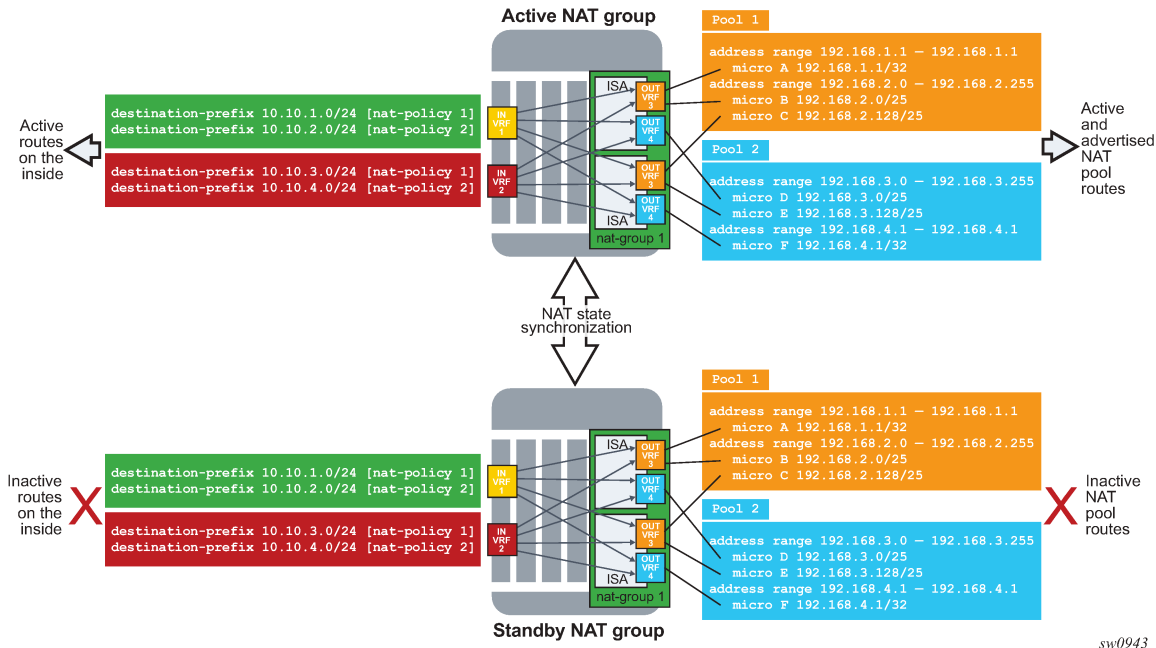
NAT-related routes are active only on the node with the active NAT group. On the outside, these are the pool routes that attract traffic in the downstream direction. On the inside, the destination-prefix (which can be configured per NAT policy) represents a route that diverts traffic to NAT in the upstream direction. In case of a filter-based diversion to NAT, a steering route is advertised only from the active node. The following are the steering-route contexts.

config>router>nat>inside>redundancy>steering-route ip-prefix/length

config>service>vprn>nat>inside>redundancy>steering-route ip-prefix/length

The route advertisement concept is shown in [Figure 101: Route advertisement](#).

Figure 101: Route advertisement



7.23.4.3 Flow synchronization

The goal of flow synchronization is to minimize service interruption after a switchover. Minimum service interruption in this context allows some packet drop during a switchover, but the user sessions are preserved without requiring a user to restart them. Flows are synchronized directly between the ISAs at the time of creation and deletion, and always only from the active side to the standby side.

The amount of traffic carried across the inter-chassis link is proportional to the number of flows that are synchronized, and it also depends on the NAT type (NAT44, DS-Lite or NAT64). The size of a flow record on the wire is around 100 bytes and a number of flow records can be packed into a single frame whose MTU is adjustable. With approximately 90 bytes of header overhead per frame, and the number of synchronized flows, it is relatively easy to estimate the required bandwidth of the inter-chassis links. The minimum recommended bandwidth of the inter-chassis link is 10 Gb/s.

Excluding short-lived flows from synchronization could further reduce the necessary bandwidth for synchronization. The flow replication threshold is configurable:

```
configure
  isa
    nat-group <id>
      active-mda-limit <limit>
      inter-chassis-redundancy
      replication-threshold <seconds> [0..300]
```

For the replication-threshold, Nokia recommends choosing a value that is larger than any of the typical short-lived flows that are not closed by the protocol itself but rely on the timeout value for its expiration (udp-dns, udp-initial, or icmp-query).

After a switchover, a resynchronization of flows occurs. The new standby ISA starts clearing all its flows and the synchronization process restarts. This means that an attempt is made to resynchronize flows from the currently active side to the standby side. During this process, the active ISAs continue to forward traffic and create new flows.

While the flow synchronization is in progress, a switchover is not allowed unless the health on the active side drops to 0, which means that one or more ISAs on the active side have failed.

7.23.4.3.1 Loss of synchronization

After the loss of synchronization, an ISA transitions into a "com-sync" state. Some of the events that can cause loss of synchronization on the ISA level are:

- Misconfigurations such as:
 - pools not matching on both nodes (outside IPs do not match between the ISAs)
 - NAT policies not matching on both nodes
- ISA-to-ISA timeout. If an ACK for any flow synchronization frame is not received within one second, the system transitions to a non-synchronized state.

When the synchronization is lost, the standby ISA starts clearing all the flows and the synchronization process restarts. This means that an attempt is made to resynchronize flows from the currently active side to the standby side. During this process, the active ISAs continue to forward traffic and create new flows.

7.23.4.3.2 Flow timeout on the standby node

When a flow is synchronized to a standby ISA, its record is present in the standby ISA until it is deleted explicitly by the active node with a delete synchronization message. There are no additional updates sent for that flow from the active node between its creation and deletion time.

7.23.4.3.3 Flow timeout following the switchover

When an ISA transitions from a standby to an active state, the timeout for the synchronized flows on the newly active ISA is set to a percentage of the flow timeout value that is configured in the NAT policy. Timeout of the flow refers to the clearing of its state in NAT after a period of traffic inactivity. Flow timeouts are configured in the NAT policy. This adjustment of the flow timeout on the newly active ISA is necessary because the standby ISA, although aware of the flows, is not aware of the flow forwarding status at the time of the switchover. In other words, the flow on the active node just before the switchover may have been inactive and close to the timer expiry. In this case, it may not be desirable to extend the life time of such a flow on the newly active ISA to its initially configured value.

This flow timeout after the switchover is set by the following command:

```
config>isa>nat-group>inter-chassis-redundancy>flow-timeout-on-switchover percent [1..50]
```

A value of 50 for the **flow-timeout-on-switchover** means that synchronized flows on the newly active ISA inherits half of the value set in the NAT policy.

7.23.4.4 Rapid consecutive switchovers

Excessive switchovers can be caused by unstable network elements. Events causing instability should be dampened at the source (ports may support event dampening). Event dampening control is not configurable under NAT. However, a dampening mechanism is built into stateful inter-chassis NAT redundancy by not allowing a switchover while synchronization is in progress.

Considering that a full re-synchronization is triggered after every switchover, the next switchover can occur only after the amount of time needed for full synchronization of flows. The new standby starts deleting flows and after this is completed, all flows from the newly active ISA are copied over to the standby.

7.23.4.5 ISA-to-ISA communication

Flow synchronization occurs directly between the ISAs, bypassing CPMs. Each ISA has its own IP address through which it communicates with its counterpart on the peering node.

Although each ISA has its own IP address, only one IP address in a NAT group is configured. This IP address is used by one of the ISAs and the rest of the ISAs are assigned consecutive IP addresses automatically by the system.

7.23.4.6 Preemption

Preemption is a feature that allows one node to relinquish activity to a peering node when the health values are equal. This functionality can be enabled through configuration (**preferred** CLI command). If one of the nodes in the redundant pair is enabled for preemption (has a **preferred** keyword configured), then this node, upon a boot up, takes over the activity from the existing active peering node with the same health value.

In a setup with two nodes with the same health value and no **preferred** command is configured on either node, the currently active node remains active and does not relinquish its activity status to a just booted up node unless the health value is changed.

If the health values are different between the peering nodes, preemption is automatically disabled, and the activity is driven solely by the health value.

A freshly booted node that is about to take over the activity from an existing node (because of preemption), waits for all the flows to be synchronized and transferred from the currently active node. Only after this transfer is complete, does switchover occur.

7.23.4.7 Message delivery prioritization

Control messages originated by the CPM that pertain to NAT redundancy are treated with the highest priority and are marked by the system with a DSCP value of Network Control NC1 (48d). Those messages are crucial to stability of NAT redundancy (otherwise inadvertent switchovers could occur).

Flow synchronization control messages originated by ISAs are marked with EF (46d) for DSCP and 5 for dot1p, which are lower priorities than CPM-originated messages. The sync (ISA keepalive) messages are sent with DSCP NC1 and dot1p 6.

Flow synchronization messages are more tolerant to delays and loss.

Both of these types of control messages should be higher priority than any customer -originated traffic that is expected to cross the ICL. Customer traffic can be marked in SR OS node by the appropriate QoS configuration on egress interfaces.

7.23.4.8 Subscriber-aware NAT

Subscriber-aware information must be supplied to both nodes (active and standby) through RADIUS accounting.

With Nokia BNG, the following configuration ensures that subscriber information is sent through RADIUS to more than one target:

```
*A:BNG1>config>subscr-mgmt>sub-prof# info
-----
    radius-accounting
      policy "acct-nat-1" "acct-nat-2"
    exit
```

Accounting policies contain respective NAT node destinations and consequently the subscriber information is sent to both NAT nodes.

7.23.4.9 Matching configuration on redundant pair of nodes

A pair of nodes participating in stateful NAT inter-chassis redundancy must have matching NAT configurations, the inside service-id and outside service-id. That is, parameters other than the configuration items referring to local objects, such as ports and interfaces, should be configured with the same values on both nodes.

For example,

```
config isa nat-group inter-chassis-redundancy keepalive interval dropcount count
```

must be the same on both nodes.

However, the statement

config isa nat-group inter-chassis-redundancy monitor-port *port-id*

does not need to match because each node can monitor its own set of unique ports, or not monitor ports at all.

Detection of configuration mismatch is logged in the system and operators are encouraged to check the logs periodically for any misaligned statements.

7.23.4.10 Online configuration changes

Certain NAT configuration changes can be performed online, which allows NAT to continue running in a redundant configuration.

In classic CLI, the user cannot perform the following online:

- delete flows
- block subscribers

In MD-CLI, these changes are allowed; however, during the commit, the inter-chassis NAT redundancy is temporarily disabled and then re-enabled after the commit is completed. Examples of configuration changes that require halting synchronization, include manipulation of NAT pools with active flows or subscribers, or removing NAT policies for active subscribers.

For configuration changes that can be performed online without halting the synchronization, there is a period during which the configuration between the nodes are misaligned (while the user is performing configuration changes). During these periods, the system continuously tries to synchronize the configuration. This perpetual attempt to synchronize flows is demanding for processing power and Nokia recommends it should be avoided.

To properly perform NAT configuration changes in a redundant configuration, the user must temporarily disable synchronization of the flows between the nodes with the following configuration. Perform the changes to the NAT configuration in an inter-chassis redundant setup.

Example: MD-CLI

```
A:admin@node-2#
  isa
    nat-group 1 {
      redundancy {
        inter-chassis {
          sync false
        }
      }
    }
  }
```

Example: classic CLI

```
-----
  nat-group 1 create
  no shutdown
  redundancy inter-chassis
  active-mda-limit 1
  exit
exit
```

In this configuration, the nodes continue to operate in active or standby mode but the newly added and deleted flows on the active node are not synchronized. During this non-synchronizing period, the user can make any NAT changes that are possible in a standalone node. When the configuration changes are performed, the **sync** command must be reversed on both nodes. At that point, the two nodes resynchronize.

7.23.4.10.1 Setting up inter-chassis redundancy

About this task

Perform the following steps to change the NAT configuration in an inter-chassis redundant setup.

Procedure

Step 1. Disable the committed synchronization of flows between the ISAs or ESAs on both nodes.

Example

MD-CLI

```
A:node-2>config>isa>nat-group>redundancy inter-chassis sync false
A:node-2>config>isa>nat-group commit
```

Example

classic CLI

```
A:node-2>config>isa>nat-group>inter-chassis-redundancy# no sync
```

This **sync** command causes the nodes to behave as if the flow synchronizations are not configured, which allows the online configuration changes.

The order in which the nodes (active or standby) are configured is irrelevant.

Step 2. Perform the configuration changes on both nodes.

Step 3. Re-enable synchronization of flows on both nodes.

Example

MD-CLI

```
A:node-2>config>isa>nat-group>redundancy inter-chassis sync true
A:node-2>config>isa>nat-group commit
```

Example

classic CLI

```
A:node-2>config>isa>nat-group>inter-chassis-redundancy# sync
```

7.23.4.11 Scenario with monitoring ports

This example relies on the following assumptions in [Figure 102: Port monitoring scenario](#):

- Load sharing over redundant chassis is achieved through two nat-groups that are, under normal conditions (no failures), active on respective chassis:
 - nat-group 1 is active on NAT node 1.
 - nat-group 2 is active on NAT node 2.
- Two 100G links on the network/public/outside side are shared between the two NAT groups on each node (Internet access). These links are redundant, and failure of one link does not have a negative effect on the traffic.
- Each NAT group has five 10G ports connected on the subscriber/private/inside side. Planned traffic load over those links is between 30G and 40G, which means that one link can be safely lost, without affecting traffic in the NAT group.

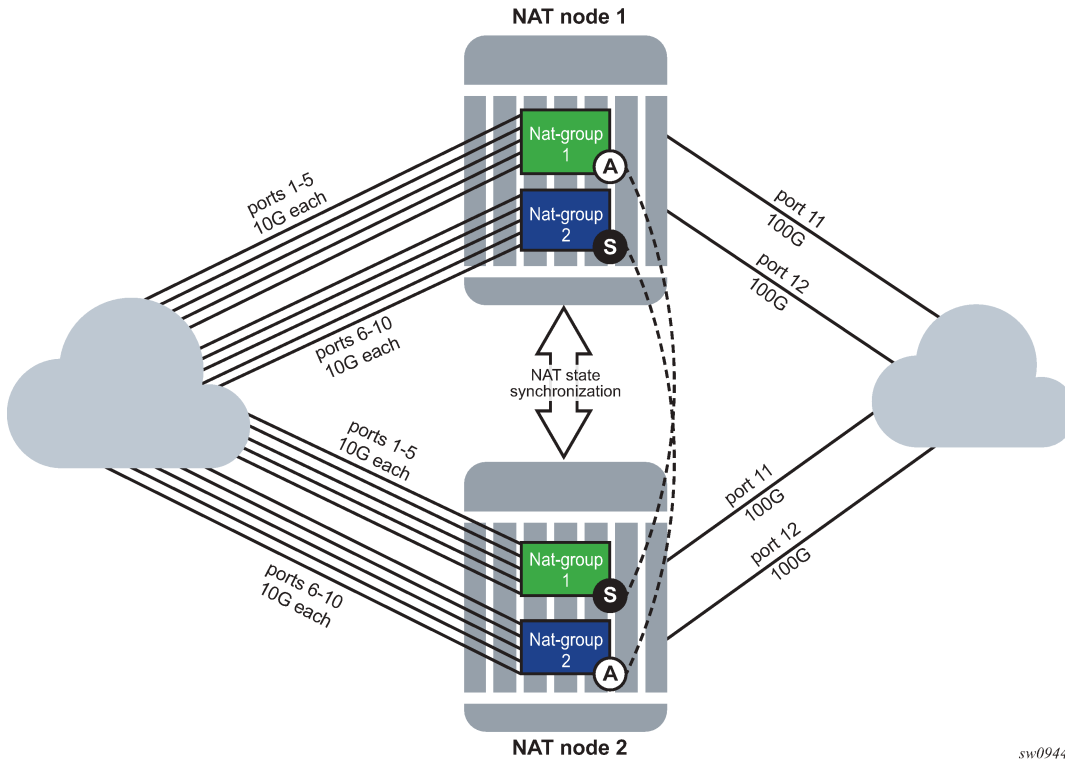
The operator's rules for managing failures are the following:

- The scheme protects against two access link failures per NAT group and one network link failure, simultaneously.
- The scheme protects against three access link failures per NAT group, simultaneously. However, in this case, there cannot be any network link failures.
- In the two above scenarios, if both network links fail on the same node (while on the other node at least one is available), the node with two failed links becomes standby.

According to those rules, the following configuration can be applied:

```
configure
  isa
    nat-group 1
      active-mda-limit 5
      inter-chassis-redundancy
        monitor-port port-1 health-drop 6
        monitor-port port-2 health-drop 6
        monitor-port port-3 health-drop 6
        monitor-port port-4 health-drop 6
        monitor-port port-5 health-drop 6
        monitor-port port-11 health-drop 10
        monitor-port port-12 health-drop 10
```


Figure 102: Port monitoring scenario



sw0944

The results for a randomly selected number of failure combinations (out of 360 valid combinations) is shown in [Table 56: Randomly selected number of failure combinations](#).

“N” indicates that the priority is equal, and unless preemption is enabled, the node that becomes active first, remains active.

Table 56: Randomly selected number of failure combinations

Node	Number of failures in nat-group 1 (10G ports)	Number of failures in nat-group 2 (10G ports)	Number of failures on shared network side (100G ports)	Health of nat-group 1	Health of nat-group 2	State of nat-group 1 (active/standby)	State of nat-group 2 (active/standby)
1	0	0	0	1000	1000	A	A
2	1	0	1	984	990	S	S
1	0	0	1	990	990	A	A
2	2	1	1	978	984	S	S
1	0	1	0	1000	994	N	S
2	0	0	0	1000	1000	N	A
1	1	1	0	1000	994	A	S

Node	Number of failures in nat-group 1 (10G ports)	Number of failures in nat-group 2 (10G ports)	Number of failures on shared network side (100G ports)	Health of nat-group 1	Health of nat-group 2	State of nat-group 1 (active/standby)	State of nat-group 2 (active/standby)
2	2	0	0	988	1000	S	A
1	1	0	1	984	990	S	A
2	0	2	1	990	978	A	S
1	1	1	1	984	984	A	N
2	2	1	1	978	984	S	N
1	1	2	0	994	988	N	S
2	1	0	0	994	1000	N	A
1	1	2	1	984	978	S	S
2	0	1	1	990	984	A	A
1	1	2	1	984	978	A	S
2	2	0	1	978	990	S	A
1	2	2	0	988	988	S	A
2	0	2	1	990	978	A	S

7.23.4.12 Configuring stateful inter-chassis NAT redundancy

Stateful inter-chassis NAT redundancy requires synchronization on the CPM level and on the ISA and ESA levels.

CPM level synchronization is required to primarily exchange health information and keepalives between the nodes for the purpose of determining active and standby NAT groups between the two peers (nodes). Each peer is identified by a single IP address. The level of traffic exchanged between the peers for CPM synchronization is low.

The ISA or ESA level synchronization is required to synchronize flows between the ISA or ESAs. Each ISA or ESA becomes a peer and is identified by its own IP address. The level of traffic exchanged between ISA or ESA for synchronization purposes depends on the configuration and the amount of NAT traffic.

Basic configuration steps are described below with command syntax examples. Some of the steps are optional and can assume default values. For more information about each command, see the *7450 ESS*, *7750 SR*, *7950 XRS*, and *VSR Classic CLI Command Reference Guide*.

1. Use the command options in the following context to configure a synchronization peer on the CPM level.

```
configure redundancy multi-chassis peer sync nat nat-group
```

The health of the NAT group is exchanged between the chassis and the node that is elected as active for the NAT group. The other node becomes the standby for the same NAT group.

2. Use the command options in the following context to configure keepalives between the nodes (CPMs):

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis keepalive
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy keepalive
```

3. Use the following command to configure the minimum duration of the flow before it is synchronized:

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis replication-threshold
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy replication-threshold
```

The user may choose to synchronize only long-lived flows.

4. Use the following command to configure a timeout of the flow after a switchover:

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis flow-timeout-on-switchover
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy flow-timeout-on-switchover
```

Independent of stateful redundancy, and depending on the type of traffic, each flow has a timeout value that determines its expiration time if there is inactivity. The initial flow timeouts are configured in a NAT policy. After a switchover, this timeout can be reset to the percentage of the originally-configured value. This can be useful because some of the flows switched over may already have been in an inactive state before the switchover.

5. Use the following command to configure the IP-MTU size of the packets carrying flow synchronization information between the ISA or ESAs:

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis ip-mtu
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy ip-mtu
```

6. Use the following commands to configure the IP address of the first ISA or ESA in a NAT group on local and remote nodes:

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis local-ip-range-start  
configure isa nat-group redundancy inter-chassis remote-ip-range-start
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy local-ip-range-start  
configure isa nat-group inter-chassis-redundancy remote-ip-range-start
```

The IP addresses for the remaining ISA or ESA are assigned automatically consecutively. These are peering addresses between the ISA and ESAs over which the flows are synchronized. Traffic from the first IP address on the local node is sent to the first IP address on the remote node.

7. Use the command options in the following contexts to configure the monitoring status of the ports and other objects, such as SAPs, BFD sessions, or VRRP sessions in the system:

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis monitor-oper-group  
configure isa nat-group redundancy inter-chassis monitor-port
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy monitor-oper-group  
configure isa nat-group inter-chassis-redundancy monitor-port
```

The status of these objects can affect the health of the system and can trigger a switchover.

8. Use the following command to select the activity preference for a NAT group:

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis preferred true
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy preferred
```

9. Use the following command to reference a routing instance through which ISA or ESAs on redundant nodes exchange synchronization information:

- **MD-CLI**

```
configure isa nat-group redundancy inter-chassis router-instance
```

- **classic CLI**

```
configure isa nat-group inter-chassis-redundancy router
```

The following are relevant command options for the **show isa nat-group** command.

```
show isa nat-group inter-chassis-redundancy
show isa nat-group inter-chassis-redundancy statistics
show isa nat-group member inter-chassis-redundancy
show isa nat-group member inter-chassis-redundancy statistics
```

7.24 ISA feature interactions

This section describes the interaction between MS-ISA applications and other system features.

7.24.1 MS-ISA use with service mirrors

All MS-ISA uses include support for service mirroring running with no feature interactions or impacts. For example, any service diverted to AA, IPsec, NAT, LNS, or supported combinations of MS-ISA application also supports service mirroring simultaneously.

7.24.2 Network Address Translation

7.24.3 Subscriber aware Large Scale NAT44

Subscriber aware Large Scale NAT44 attempts to combine the positive attributes of Large Scale NAT44 and L2-Aware NAT, namely:

- the ability for some traffic to bypass the NAT function, such as IPTV traffic and VoIP traffic whenever a unique IP address per subscriber is used (for example, not L2-Aware NAT where all subs share the same IP). This can be achieved using existing Large Scale NAT44 mechanisms (ingress IP-filters)
- the use of RADIUS Acct for logging of port-ranges, including multiple port-range blocks
- the use of subscriber-identification/RADIUS username to identify the customer to simplify management of Large Scale NAT44 subscribers

Subscriber awareness in Large Scale NAT44 facilitates the release of NAT resources immediately after the BNG subscriber is terminated, without having to wait for the last flow of the subscriber to expire on its own (TCP timeout is 4hours by default).

The subscriber aware Large Scale NAT44 function leverages RADIUS accounting proxy built-in to the 7750 SR. The RADIUS accounting proxy allows the 7750 SR to inform Large Scale NAT44 application about individual BNG subscribers from the RADIUS accounting messages generated by a remote BNG and use this information in the management of Large Scale NAT44 subscribers. The combination of the two allows, for example, the 7750 SR running as a Large Scale NAT44 to make the correlation between the BNG subscriber (represented in the Large Scale NAT44 by the Inside IP Address) and RADIUS attributes such as User-Name, Alc-Sub-Ident-String, Calling-Station-Id or Class. These attributes can subsequently be used for either management of the Large Scale NAT44 subscriber, or in the NAT RADIUS Accounting messages generated by the 7750 SR Large Scale NAT44 application. Doing so simplifies both the administration of the Large Scale NAT44 and the logging function for port-range blocks.

As BNG subscribers authenticate and come online, the RADIUS accounting messages are 'snooped' through RADIUS accounting proxy which creates a cache of attributes from the BNG subscriber. BNG subscribers are correlated with the NAT subscriber by framed-ip address, and one of the following attributes that must be present in the accounting messages generated by BNG:

- username
- subscriber ID
- RADIUS Class attribute
- Calling-Station-id
- IMSI
- IMEI

Framed-ip address must also be present in the accounting messages generated by BNG.

Large Scale NAT44 Subscriber Aware application receives a number of cached attributes which are used for appropriate management of Large Scale NAT44 subscribers, for example:

- Delete the Large Scale NAT44 subscriber when the BNG subscriber is terminated.
- Report attributes in Large Scale NAT44 accounting messages according to configuration options.

Creation and removal of RADIUS accounting proxy cache entries related to BNG subscriber is triggered by the receipt of accounting start/stop messages sourced by the BNG subscriber. Modification of entries can be triggered by interim-update messages carrying updated attributes. Cached entries can also be purged via CLI.

In addition to passing one of the above attributes in Large Scale NAT44 RADIUS accounting messages, a set of opaque BNG subscriber RADIUS attributes can optionally be passed in Large Scale NAT44 RADIUS accounting messages. Up to 128B of these opaque attributes are accepted. The remaining attributes are truncated.

Large Scale NAT44 subscriber instantiation can optionally be denied in case that corresponding BNG subscriber cannot be identified in Large Scale NAT44 through RADIUS accounting proxy.

Configuration guidelines:

Configure RADIUS accounting proxy functionality in a routing instance that receives accounting messages from the remote or local BNG. Optionally forward received accounting message received by RADIUS accounting proxy to the final accounting destination (accounting server).

Point the BNG RADIUS accounting destination to the RADIUS accounting proxy – this way RADIUS accounting proxy receives and 'snoop' BNG RADIUS accounting data.

BNG subscriber can be associated with two accounting policies, therefore pointing to two different accounting destinations. For example, one to the RADIUS accounting proxy, the other one to the real accounting server.

Configure subscriber aware Large Scale NAT44. From Large Scale NAT44 Subscriber Aware application reference the RADIUS Proxy accounting server and define the string that is used to correlate BNG subscriber with the Large Scale NAT44 subscriber.

Optionally enable NAT RADIUS accounting that includes BNG subscriber relevant data.

(1)

```
*A:left-a20>config>service>vprn#
  radius-proxy
    server "proxy-acct" purpose accounting create
      default-accounting-server-policy "lsn-policy"
      description "two side server -interface:client ; default-plcy:real server"
      interface "rad-proxy-loopback"
      secret "TEg1UEZzemRMyZXD1HvvQGkeGfoQ58MF" hash2
      no shutdown
    exit
  exit
```

RADIUS accounting proxy listens to accounting messages on interface 'rad-proxy-loopback'.

The name 'proxy-acct' as defined by the **server** command is used to reference this proxy accounting server from Large Scale NAT44.

Received accounting messages can be relayed further from RADIUS accounting proxy to the accounting server which can be indirectly referenced in the default-accounting-policy 'lsn-policy'.

```
The lsn-policy is defined as:
*A:left-a20>config>aaa#
  radius-server-policy "lsn-policy" create
  servers
    router "Base"
      source-address 192.168.1.12
      server 1 name "192"
    exit
  exit
```

This lsn-policy can then reference an external RADIUS accounting server with its own security credentials. This external accounting server can be configured in any routing instance.

```
*A:left-a20>config>router>radius-server# info
-----
  server "192" address 192.168.1.10 secret "KRr7H.K3i0z90/hj2BUSmdJUdl.zWrkE" hash2
  port 1813 create
  description "real radius or acct server"
  exit
```

(2)

Two RADIUS accounting policies can be configured in BNG, one to the real RADIUS server, the other one to the RADIUS accounting proxy.

```
*A:left-a20>config>subscr-mgmt>sub-prof# info
-----
  radius-accounting-policy "real-acct-srvr" duplicate "lsn"
  egress
  agg-rate-limit 10000
```

```

exit
-----
*A:left-a20>config>subscr-mgmt>acct-plcy# info

```

```

-----
description "lsn radius-acct-policy"
update-interval 5
include-radius-attribute
  acct-authentic
  acct-delay-time
  called-station-id
  calling-station-id remote-id
  circuit-id
  framed-interface-id
  framed-ip-addr
  framed-ip-netmask
  mac-address
  nas-identifier
  nas-port-id
  nas-port-type
  nat-port-range
  remote-id
  sla-profile
  sub-profile
  subscriber-id
  user-name
  alc-acct-triggered-reason
exit
session-id-format number
radius-accounting-server
  router 10 (service id where proxy radius is configured)
  server 1 address 10.5.5.5 secret "cVi1sidvgH28Pd9QoN1fLE" hash2
(radius proxy IP address is 10.5.5.5 on interface "rad-proxy-loopback"; the 'secret' is the
same as configured on RADIUS accounting proxy)
exit

```

(3)

Sub-aware Large Scale NAT44 references the RADIUS accounting proxy server 'proxy-acct' and defines the calling-station-id attribute from the BNG subscriber as the matching attribute:

```

*A:left-a20>config>service>vprn>nat>inside# info
-----
nat-policy "nat-base"
  destination-prefix 10.0.0.0/16
  subscriber-identification
    attribute vendor "standard" attribute-type "station-id"
  description "sub-aware CGN"
  radius-proxy-server router 10 name "proxy-acct"
  no shutdown
exit
-----

```

(4)

Optionally RADIUS NAT accounting can be enabled:

```

*A:left-a20>config>isa>nat-group# info

```



```

-----
    active-mda-limit 1
    radius-accounting-policy "nat-acct-basic"
    mda 1/2
    no shutdown

*A:left-a20>config>aaa>isa-radius-plcy# info detail
-----
    description "radius accounting policy for NAT"
    include-radius-attribute
        framed-ip-addr
        nas-identifier
        nat-subscriber-string
        user-name
        inside-service-id
        outside-service-id
        outside-ip
        port-range-block
        hardware-timestamp
        release-reason
        multi-session-id
        frame-counters
        octet-counters
        session-time
        called-station-id
        subscriber-data
    exit
    radius-accounting-server
        access-algorithm direct
        retry 3
        router "Base"
        source-address-range 192.168.1.20 192.168.1.20
        timeout sec 5
        server 1 address 192.168.1.10 secret "KlWIBi08CxTyM/YXaU2gQit
0u8Ggf5D70j5hjese27A" hash2 port 1813
    exit
-----

```

Such setup would produce a stream of following Large Scale NAT44 RADIUS accounting messages:

```

Mon Jul 16 10:59:27 2012
NAS-IP-Address = 10.1.1.1
NAS-Identifier = "left-a20"
NAS-Port = 37814272
Acct-Status-Type = Start
Acct-Multi-Session-Id = "500456500365a4de7c29a9a07c29a9a0"
Acct-Session-Id = "500456500365a4de6201d7b87c29a9a0"
Called-Station-Id = "00-00-00-00-01-01"
User-Name = "remote0"
Calling-Station-Id = "remote0"
Alc-Serv-Id = 10
Framed-IP-Address = 10.0.0.7
Alc-Nat-Outside-IP-Addr = 198.51.100.1
Alc-Nat-Port-Range = "198.51.100.1 1054-1058 router base"
Acct-Input-Packets = 0
Acct-Output-Packets = 0
Acct-Input-Octets = 0
Acct-Output-Octets = 0
Acct-Input-Gigawords = 0
Acct-Output-Gigawords = 0
Acct-Session-Time = 0
Event-Timestamp = "Jul 16 2012 10:58:40 PDT"

```

```
NAS-IP-Address = 10.1.1.1
User-Name = "cgn_1_ipoe"
Framed-IP-Netmask = 255.255.255.0
Class = 0x63676e2d636c6173732d7375622d6177617265
NAS-Identifier = "left-a20"
Acct-Session-Id = "D896FF0000000550045640"
Event-Timestamp = "Jul 16 2012 10:58:24 PDT"
NAS-Port-Type = Ethernet
NAS-Port-Id = "1/1/5:5.10"
Acct-Delay-Time = 0
Acct-Authentic = RADIUS
Acct-Unique-Session-Id = "10f8bce6e5e7eb41"
Timestamp = 1342461567
Request-Authenticator = Verified
```

Mon Jul 16 11:03:56 2012

```
NAS-IP-Address = 10.1.1.1
NAS-Identifier = "left-a20"
NAS-Port = 37814272
Acct-Status-Type = Interim-Update
Acct-Multi-Session-Id = "500456500365a4de7c29a9a07c29a9a0"
Acct-Session-Id = "500456500365a4de6201d7b87c29a9a0"
Called-Station-Id = "00-00-00-00-01-01"
User-Name = "remote0"
Calling-Station-Id = "remote0"
Alc-Serv-Id = 10
Framed-IP-Address = 10.0.0.7
Alc-Nat-Outside-IP-Addr = 198.51.100.1
Alc-Nat-Port-Range = "198.51.100.1 1054-1058 router base"
Acct-Input-Packets = 0
Acct-Output-Packets = 1168
Acct-Input-Octets = 0
Acct-Output-Octets = 86432
Acct-Input-Gigawords = 0
Acct-Output-Gigawords = 0
Acct-Session-Time = 264
Event-Timestamp = "Jul 16 2012 11:03:04 PDT"
Acct-Delay-Time = 5
NAS-IP-Address = 10.1.1.1
User-Name = "cgn_1_ipoe"
Framed-IP-Netmask = 255.255.255.0
Class = 0x63676e2d636c6173732d7375622d6177617265
NAS-Identifier = "left-a20"
Acct-Session-Id = "D896FF0000000550045640"
Acct-Session-Time = 279
Event-Timestamp = "Jul 16 2012 11:03:04 PDT"
NAS-Port-Type = Ethernet
NAS-Port-Id = "1/1/5:5.10"
Acct-Delay-Time = 0
Acct-Authentic = RADIUS
Acct-Unique-Session-Id = "10f8bce6e5e7eb41"
Timestamp = 1342461836
Request-Authenticator = Verified
```

Mon Jul 16 11:04:34 2012

```
NAS-IP-Address = 10.1.1.1
NAS-Identifier = "left-a20"
NAS-Port = 37814272
Acct-Status-Type = Stop
Acct-Multi-Session-Id = "500456500365a4de7c29a9a07c29a9a0"
Acct-Session-Id = "500456500365a4de6201d7b87c29a9a0"
Called-Station-Id = "00-00-00-00-01-01"
User-Name = "remote0"
Calling-Station-Id = "remote0"
```

```

Alc-Serv-Id = 10
Framed-IP-Address = 10.0.0.7
Alc-Nat-Outside-IP-Addr = 198.51.100.1
Alc-Nat-Port-Range = "198.51.100.1 1054-1058 router base"
Acct-Terminate-Cause = Host-Request
Acct-Input-Packets = 0
Acct-Output-Packets = 1321
Acct-Input-Octets = 0
Acct-Output-Octets = 97754
Acct-Input-Gigawords = 0
Acct-Output-Gigawords = 0
Acct-Session-Time = 307
Event-Timestamp = "Jul 16 2012 11:03:47 PDT"
NAS-IP-Address = 10.1.1.1
User-Name = "cgn_1_ipoe"
Framed-IP-Netmask = 255.255.255.0
Class = 0x63676e2d636c6173732d7375622d6177617265
NAS-Identifier = "left-a20"
Acct-Session-Id = "D896FF0000000550045640"
Acct-Session-Time = 279
Event-Timestamp = "Jul 16 2012 11:03:04 PDT"
NAS-Port-Type = Ethernet
NAS-Port-Id = "1/1/5:5.10"
Acct-Delay-Time = 0
Acct-Authentic = RADIUS
Acct-Unique-Session-Id = "10f8bce6e5e7eb41"
Timestamp = 1342461874
Request-Authenticator = Verified

```

The matching accounting stream generated on the BNG is shown below:

```

Mon Jul 16 10:59:11 2012
Acct-Status-Type = Start
NAS-IP-Address = 10.1.1.1
User-Name = "cgn_1_ipoe"
Framed-IP-Address = 10.0.0.7
Framed-IP-Netmask = 255.255.255.0
Class = 0x63676e2d636c6173732d7375622d6177617265
Calling-Station-Id = "remote0"
NAS-Identifier = "left-a20"
Acct-Session-Id = "D896FF0000000550045640"
Event-Timestamp = "Jul 16 2012 10:58:24 PDT"
NAS-Port-Type = Ethernet
NAS-Port-Id = "1/1/5:5.10"
ADSL-Agent-Circuit-Id = "cgn_1_ipoe"
ADSL-Agent-Remote-Id = "remote0"
Alc-Subsc-ID-Str = "CGN1"
Alc-Subsc-Prof-Str = "nat"
Alc-SLA-Prof-Str = "tp_sla_prem"
Alc-Client-Hardware-Addr = "2001:db8:65:05:10:01"
Acct-Delay-Time = 0
Acct-Authentic = RADIUS
Acct-Unique-Session-Id = "9c1723d05e87c043"
Timestamp = 1342461551
Request-Authenticator = Verified

Mon Jul 16 11:03:51 2012
Acct-Status-Type = Interim-Update
NAS-IP-Address = 10.1.1.1
User-Name = "cgn_1_ipoe"
Framed-IP-Address = 10.0.0.7
Framed-IP-Netmask = 255.255.255.0

```

```
Class = 0x63676e2d636c6173732d7375622d6177617265
Calling-Station-Id = "remote0"
NAS-Identifier = "left-a20"
Acct-Session-Id = "D896FF0000000550045640"
Acct-Session-Time = 279
Event-Timestamp = "Jul 16 2012 11:03:04 PDT"
NAS-Port-Type = Ethernet
NAS-Port-Id = "1/1/5:5.10"
ADSL-Agent-Circuit-Id = "cgn_1_ipoe"
ADSL-Agent-Remote-Id = "remote0"
Alc-Subsc-ID-Str = "CGN1"
Alc-Subsc-Prof-Str = "nat"
Alc-SLA-Prof-Str = "tp_sla_prem"
Alc-Client-Hardware-Addr = "2001:db8:65:05:10:01"
Acct-Delay-Time = 0
Acct-Authentic = RADIUS
Alcatel-IPD-Attr-163 = 0x00000001
Alc-Acct-I-Inprof-Octets-64 = 0x00010000000000000000
Alc-Acct-I-Outprof-Octets-64 = 0x000100000000000020468
Alc-Acct-I-Inprof-Pkts-64 = 0x00010000000000000000
Alc-Acct-I-Outprof-Pkts-64 = 0x00010000000000000052a
Alc-Acct-I-Inprof-Octets-64 = 0x00030000000000000000
Alc-Acct-I-Outprof-Octets-64 = 0x00030000000000000000
Alc-Acct-I-Inprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-I-Outprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-I-Inprof-Octets-64 = 0x00050000000000000000
Alc-Acct-I-Outprof-Octets-64 = 0x00050000000000000000
Alc-Acct-I-Inprof-Pkts-64 = 0x00050000000000000000
Alc-Acct-I-Outprof-Pkts-64 = 0x00050000000000000000
Alc-Acct-0-Inprof-Octets-64 = 0x00010000000000000000
Alc-Acct-0-Outprof-Octets-64 = 0x00010000000000003154
Alc-Acct-0-Inprof-Pkts-64 = 0x00010000000000000000
Alc-Acct-0-Outprof-Pkts-64 = 0x00010000000000000009a
Alc-Acct-0-Inprof-Octets-64 = 0x00030000000000000000
Alc-Acct-0-Outprof-Octets-64 = 0x00030000000000000000
Alc-Acct-0-Inprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-0-Outprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-0-Inprof-Octets-64 = 0x00050000000000000000
Alc-Acct-0-Outprof-Octets-64 = 0x00050000000000000000
Alc-Acct-0-Inprof-Pkts-64 = 0x00050000000000000000
Alc-Acct-0-Outprof-Pkts-64 = 0x00050000000000000000
Acct-Unique-Session-Id = "9c1723d05e87c043"
Timestamp = 1342461831
Request-Authenticator = Verified
```

Mon Jul 16 11:04:34 2012

```
Acct-Status-Type = Stop
NAS-IP-Address = 10.1.1.1
User-Name = "cgn_1_ipoe"
Framed-IP-Address = 10.0.0.7
Framed-IP-Netmask = 255.255.255.0
Class = 0x63676e2d636c6173732d7375622d6177617265
Calling-Station-Id = "remote0"
NAS-Identifier = "left-a20"
Acct-Session-Id = "D896FF0000000550045640"
Acct-Session-Time = 322
Acct-Terminate-Cause = User-Request
Event-Timestamp = "Jul 16 2012 11:03:47 PDT"
NAS-Port-Type = Ethernet
NAS-Port-Id = "1/1/5:5.10"
ADSL-Agent-Circuit-Id = "cgn_1_ipoe"
ADSL-Agent-Remote-Id = "remote0"
Alc-Subsc-ID-Str = "CGN1"
Alc-Subsc-Prof-Str = "nat"
```

```

Alc-SLA-Prof-Str = "tp_sla_prem"
Alc-Client-Hardware-Addr = "2001:db8:65:05:10:01"
Acct-Delay-Time = 0
Acct-Authentic = RADIUS
Alc-Acct-I-Inprof-Octets-64 = 0x00010000000000000000
Alc-Acct-I-Outprof-Octets-64 = 0x0001000000000000248c4
Alc-Acct-I-Inprof-Pkts-64 = 0x00010000000000000000
Alc-Acct-I-Outprof-Pkts-64 = 0x0001000000000000005d9
Alc-Acct-I-Inprof-Octets-64 = 0x00030000000000000000
Alc-Acct-I-Outprof-Octets-64 = 0x00030000000000000000
Alc-Acct-I-Inprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-I-Outprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-I-Inprof-Octets-64 = 0x00050000000000000000
Alc-Acct-I-Outprof-Octets-64 = 0x00050000000000000000
Alc-Acct-I-Inprof-Pkts-64 = 0x00050000000000000000
Alc-Acct-I-Outprof-Pkts-64 = 0x00050000000000000000
Alc-Acct-0-Inprof-Octets-64 = 0x00010000000000000000
Alc-Acct-0-Outprof-Octets-64 = 0x00010000000000003860
Alc-Acct-0-Inprof-Pkts-64 = 0x00010000000000000000
Alc-Acct-0-Outprof-Pkts-64 = 0x000100000000000000b0
Alc-Acct-0-Inprof-Octets-64 = 0x00030000000000000000
Alc-Acct-0-Outprof-Octets-64 = 0x00030000000000000000
Alc-Acct-0-Inprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-0-Outprof-Pkts-64 = 0x00030000000000000000
Alc-Acct-0-Inprof-Octets-64 = 0x00050000000000000000
Alc-Acct-0-Outprof-Octets-64 = 0x00050000000000000000
Alc-Acct-0-Inprof-Pkts-64 = 0x00050000000000000000
Alc-Acct-0-Outprof-Pkts-64 = 0x00050000000000000000
Acct-Unique-Session-Id = "9c1723d05e87c043"
Timestamp = 1342461874
Request-Authenticator = Verified

```

7.25 Mapping of Address and Port using Translation (MAP-T)

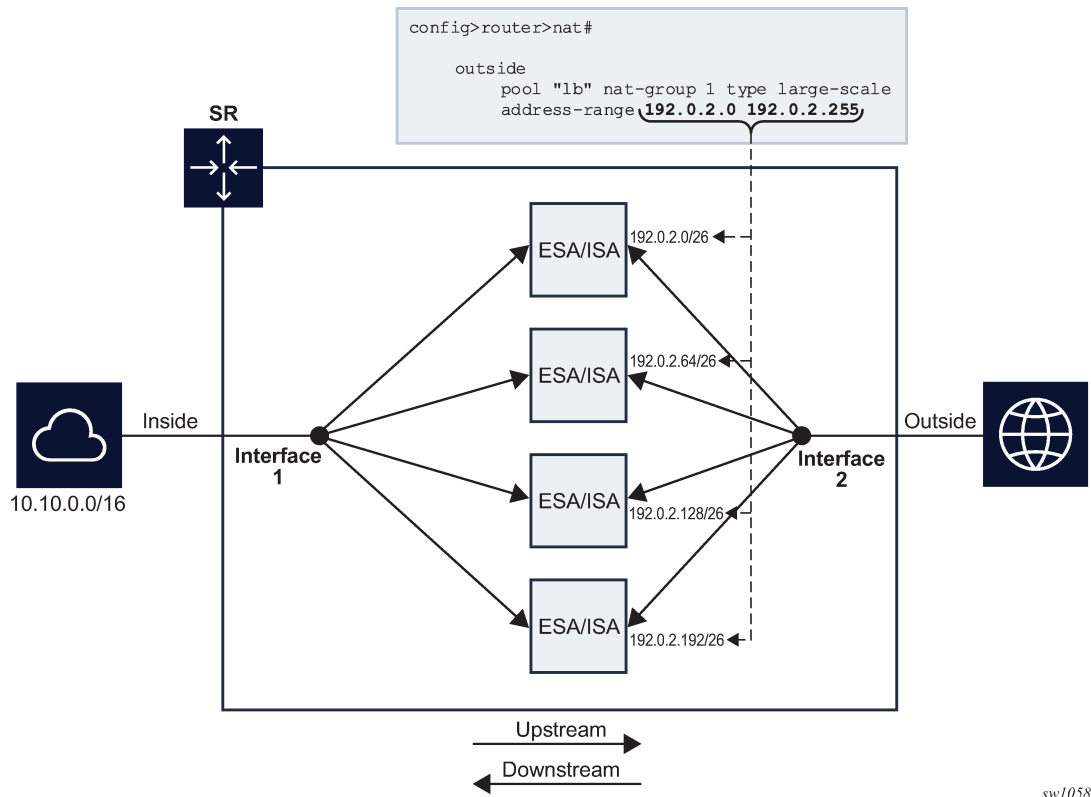


Note: The MAP-T feature and commands described in this section apply to the Nokia Virtualized Service Router (VSR) only.

MAP-T is a NAT technique defined in RFC 7599. Its key advantage is the decentralization of stateful NAT while enabling the sharing of public IPv4 addresses among the customer edge (CE) devices. In a nutshell, the CE performs the stateful NAT44 function and translates the resulting IPv4 packet into an IPv6 packet. The IPv6 packet is transported over the IPv6 network to the Border Router (BR), which then translates the IPv6 packet to IPv4 and sends it into the public domain.

As multiple CEs can share a single public IPv4 address, MAP-T must rely on an algorithm (A+P algorithm running on the CEs and BR) to ensure that each CE is assigned a unique port-range on a shared IPv4 public address. In this way, each CE can be uniquely identified at the BR by a combination of the shared IPv4 public address and unique port-range. A set of CEs and BR that share a common set of MAP algorithm rules constitutes a MAP domain. [Figure 103: MAP-T network level view](#) shows a network-level view of Map-T.

Figure 103: MAP-T network level view



MAP-T offers the following advantages mainly as a result of its stateless BR operation:

- **improved scaling**
State maintenance is decentralized, which enables better scaling.
- **simplified redundancy**
There are no sessions synchronized between redundant BRs and this translates to simplified redundancy.
- **reduced logging**
As there are no NAT resources in the BR that require logging, only configuration changes in the BR are logged, which reduces the volume of logging data.
- **simpler communication**
MAP-T simplifies user-to-user communication.
- **higher throughput**
MAP-T offers higher throughput than a stateful solution, with less processing required in the BR.

Mapping of address and port (MAP) is a generic function, regardless of the underlying transport mechanism (MAP-T or MAP-E) used. Each MAP CE is assigned as follows:

- **a shared public IPv4 address with a unique port-range on the shared IPv4 address**
Although a shared IPv4 address is used in most cases, the CE is sometimes assigned a unique IPv4 address or even an IPv4 prefix. This information is used for stateful NAT44 at the CE.

- **an IPv6 prefix (IA-PD)**

A "subnet" from the IPv6 prefix is allocated to the CE as a MAP prefix. The MAP prefix is used to encode public IPv4 information and identify the CE in a MAP domain. The remainder of the IA-PD is used on the LAN side of the CE.

- **an IPv6 address (IA-NA)**

The IPv6 address is independent of MAP and is a regular IPv6 address on the WAN side. The address is used for native end-to-end IPv6 communication (it can participate in forming routing adjacencies and other tasks).

The CE and BR perform the following functions in the MAP-T domain:

- **CE upstream direction (IPv4→IPv6)**

- Perform stateful NAT44 function (private→public).
- Translate the public IPv4 address and port into an assigned IPv6 MAP source address.
- Send the IPv6 packet with encoded IPv4 information toward the BR.

- **BR upstream direction (IPv6→IPv4)**

- Perform an anti-spoof check on the received IPv6 packet to ensure that it is coming from a trusted source (CE).

Anti-spoofing is achieved by checking the source IPv6 MAP address against the configured MAP rules and making sure that the correct public IPv4 address and port-range of the CE are encoded in the CE's source IPv6 MAP address.

- Translate the IPv6 packet into an IPv4 packet and forward it into the public domain.

- **BR downstream direction (IPv6←-IPv4)**

- Translate the IPv4 packet into an IPv6 packet according to MAP rules.

The IPv4 destination address of the received packet is translated into an IPv6 MAP address of the CE.

- Send the IPv6 packet toward the CE.

- **CE downstream direction (IPv4 ←- IPv6)**

- Perform the anti-spoofing function using the destination IPv6 address to verify that the packet is destined for the CE.

MAP rules are used to verify that the public IPv4 address and the port-range of the CE is encoded in the IPv6 destination IP address of the received packet (IPv6 MAP address of the CE).

- Translate the IPv6 packet into an IPv4 packet.
- Perform the NAT44 function (public→private).
- Forward the packet into the private IPv4 network.

Each device (CE and BR) is also responsible for fragmentation handling and ICMP error reporting (MTU too small, TTL expired, and so on).

7.25.1 MAP-T rules

MAP-T rules control the address translation in a MAP-T domain. The mapping rules can be delivered to the devices in the MAP domain using RADIUS or DHCP, or be statically provisioned.

The MAP-T rules are:

- **Basic Mapping Rule (BMR)**

The BMR is used to translate the public IPv4 address and port-range assigned to the CE into the IPv6 MAP address. It is composed of the following parameters:

- rule IPv6 prefix (including prefix length)
- rule IPv4 prefix (including prefix length)
- rule Embedded Address bits (EA-bits) define the portion of the IA-PD that encodes the IPv4 suffix and port-range
- rule Port Parameters (optional)

- **Forwarding Mapping Rule (FMR)**

The FMR is used for forwarding within the MAP domain. FMRs are instantiated in the BR so that the BR can forward traffic to the CEs. FMRs can also be instantiated in CEs to forward traffic directly between CEs, effectively bypassing BR. The FMR is composed of the same set of parameters as the BMR:

- rule IPv6 prefix (including prefix length)
- rule IPv4 prefix (including prefix length)
- rule Embedded Address (EA) bits that define the portion of the IA-PD that encodes the IPv4 suffix and port-range
- rule Port Parameters (optional)

- **Default Mapping Rule (DMR)**

The DMR is used to forward traffic outside the MAP domain. This rule contains the IPv6 prefix of the BR in MAP-T and it is used as the default route.

7.25.2 A+P mapping algorithm

The public IPv4 address and the port-range information of the CE is encoded in its assigned IPv6 delegated prefix (IA-PD). The BMR holds the key to decode this information from the IA-PD of the CE. The BMR identifies the portion of bits of the IA-PD that contain the suffix of the IPv4 address and the port-set ID (PSID). These bits are called the EA bits. The PSID represents the port-range assigned to the CE.

The public IPv4 address of the CE is constructed by concatenating the IPv4 prefix carried in the BMR and the suffix, which is extracted from the EA bits within the IA-PD. The port-range is identified by the remaining EA bits (PSID portion). The EA bits uniquely identify the CE within the IPv6 network in a MAP domain.

The *psid-offset* value must be set to a value greater than 0. It represents ports that are omitted from the mapping (for example, well-known ports).

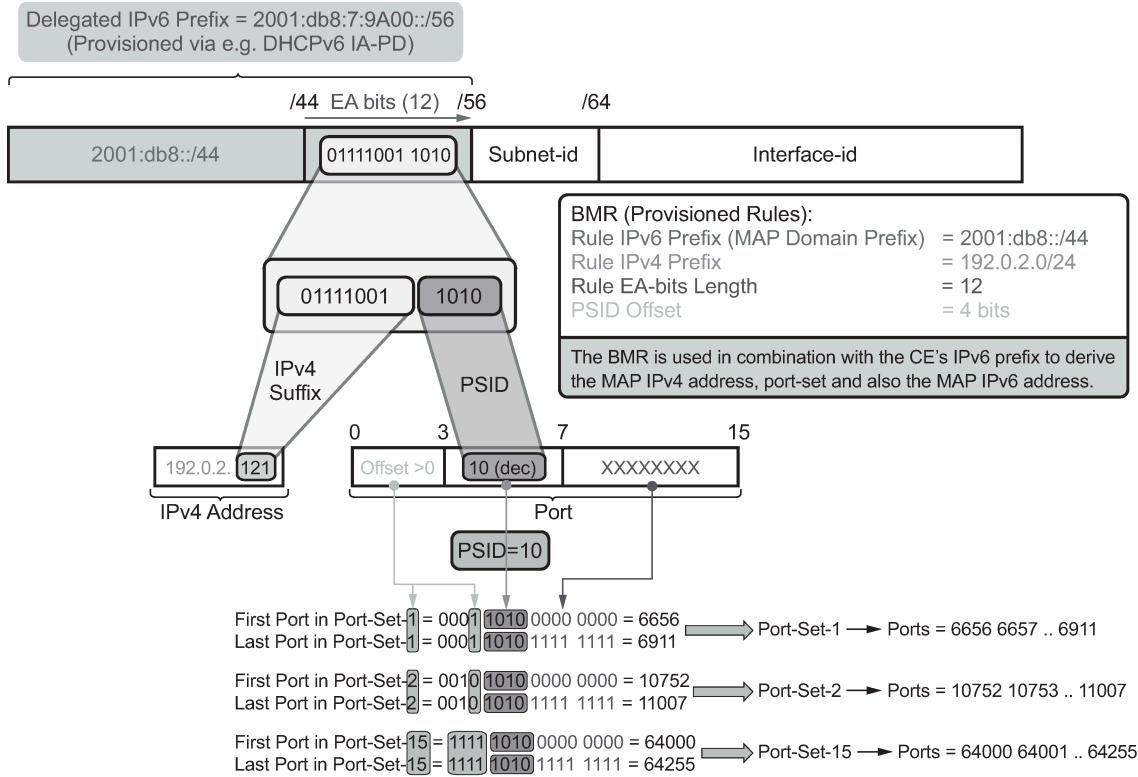
An IPv4 address and port on the private side of the CE must be statefully translated to a public IPv4 address, and the port within the assigned port must be set on the public side of the CE. This ensures that the BR, based on the same MAP rules, can extrapolate the IPv4 source of the packet for verification (anti-spoofing) purposes in the upstream direction, and conversely, to determine the destination IPv6 MAP address (CE address) in the downstream direction (based on the destination IPv4+port).

The IPv6 MAP address is constructed by setting the subnet-id in the delegated IPv6 prefix to 0. In this way, the subnet-id of 0 is reserved for MAP function. The remaining subnets can be delegated on the LAN side of the CE.

The interface-id is set to the IPv4 public address and PSID. This is described in RFC 7599, §6.

In this way, the IPv4 and IPv6 addresses of the CE are defined and easily converted to each other based on the BMR and the port information in the packet. [Figure 104: A+P mapping](#) shows the A+P mapping algorithm.

Figure 104: A+P mapping

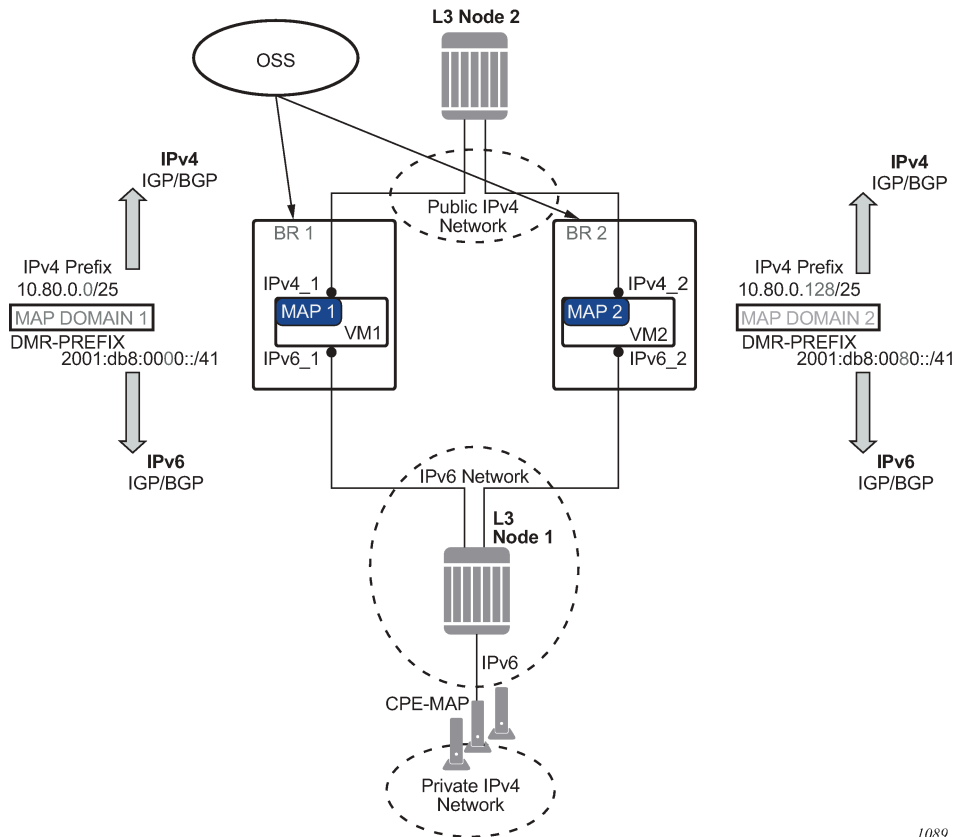


1086

7.25.3 Routing considerations

[Figure 105: MAP-T deployment scenario](#) shows a MAP-T deployment scenario.

Figure 105: MAP-T deployment scenario



1089

The routes related to MAP-T are:

- **IPv4 prefixes from the MAP-rules**

These routes are advertised in the upstream direction.

- **DMR**

This is the BR prefix for a specific domain. This route is advertised in the downstream direction.

Routes related to MAP-T are advertised through IPv4 and IPv6 routing protocols. MAP-T routes in the VSR are owned by "protocol nat" with a metric of 50. This can be used to configure an export routing policy when advertising MAP-T routes in IGP or BGP.

Multiple MAP-T domains can be supported in the same routing context.



Note: IPv6 IA-PD end-user prefixes are carved out of the IPv6 rule prefix. Aside from MAP-T, IA-PD is used for native IPv6 end-to-end traffic outside of MAP-T. Although the IPv6 rule prefix is not marked as a NAT route in the routing table, it is nonetheless advertised in the upstream direction.

7.25.4 Forwarding considerations in the BR

In the upstream direction, when the BR receives an IPv6 packet destined for the BR prefix, a source-based IPv6 address lookup (anti-spoofing) is performed to verify that the packet is sent by the credible CE.

In the downstream direction, a destination-based IPv4 lookup is performed. This leads to the MAP-T rule entry, which provides the information necessary to derive the IPv6 address of the destination CE.

The MAP-T forwarding function in the VSR is also responsible for:

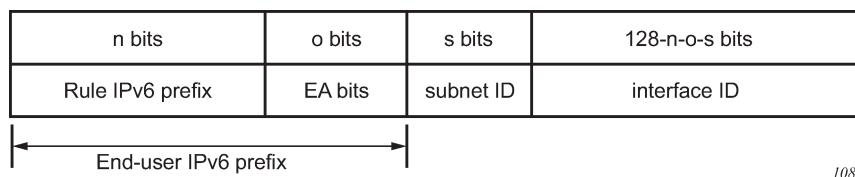
- address conversion between IPv4 and IPv6 based on the BMR rule
- header translation between IPv4 and IPv6, as described in RFC 6145

In address-sharing scenarios, address translation is performed for TCP/UDP and a subset of ICMP traffic; everything else is dropped. In contrast, 1:1 and prefix-sharing scenarios are protocol agnostic.

7.25.4.1 IPv6 addresses

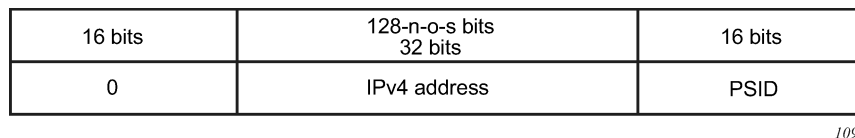
An IPv6 address of the MAP-T node is constructed according to RFC 7597, §5.2 and RFC 7599, §6. [Figure 106: IPv6 address construction](#) shows the IPv6 address of the MAP-T node.

Figure 106: IPv6 address construction



The subnet ID for a MAP node (CE) is set to 0. [Figure 107: Node interface](#) shows the node interface (PSID is left-padded with zeros to create a 16-bit field and the IPv4 address is the public IPv4 address assigned to the CE).

Figure 107: Node interface



This constructed IPv6 address represents the source IPv6 address of traffic sent from the CE to BR (upstream direction), and the destination IPv6 address in the opposite direction (downstream traffic sent from the BR to the CE).

The source IPv6 address in the downstream direction is a combination of the BR IPv6 prefix and the source IPv4 address (per RFC 7599, §5.1) received in the original packet.

The destination IPv6 address in the upstream direction is a combination of the BR IPv6 prefix and the IPv4 destination address (RFC 7599, §5.1) in the original packet.

7.25.4.2 1:1 translations and IPv4 prefix translations

1:1 translations refer to the case in which each CE is assigned a distinct public IPv4 address; that is, there is no public IPv4 address sharing between the CEs. In this case, the PSID field is 0 and the sum of lengths for the IPv4 rule prefix and EA bits is 32. In other words, all the EA bits represent the IPv4 suffix. The public IPv4 address of the CE is created by concatenating the Rule IPv4 Prefix and the EA bits.

IPv4 Prefix translations refer to the case where an IPv4 prefix is assigned to a CE. In this case, the PSID field is 0 and the sum of the lengths for IPv4 rule prefix and EA bits is less than 32.

In both preceding cases, the translations are protocol agnostic; all protocols, not just TCP/UDP or ICMP, is translated.

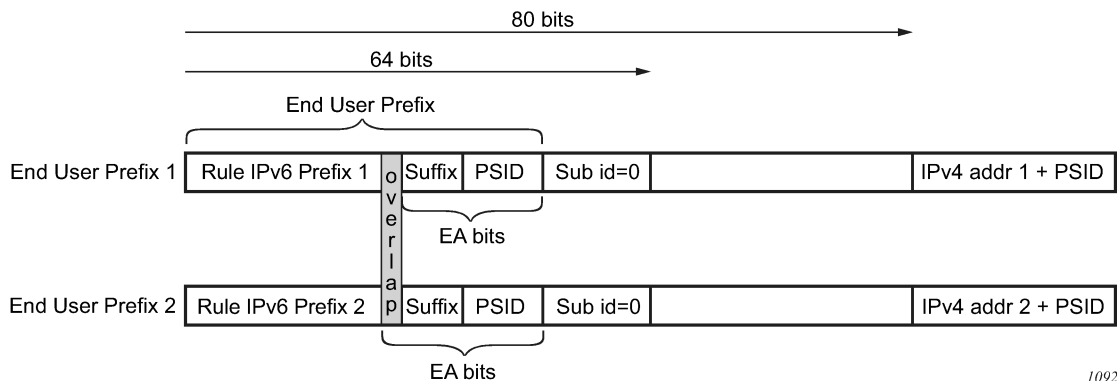
7.25.4.3 Hub-and-spoke topology

The BR supports hub-and-spoke topology, which means that the BR facilitates communication between MAP-T CEs.

7.25.4.4 Rule prefix overlap

Rule prefix overlap is not supported because it can cause lookup ambiguity. [Figure 108: IPv6 rule prefix overlap](#) shows a rule prefix overlap example.

Figure 108: IPv6 rule prefix overlap



In the case where rule IPv6 prefix 1 is a subset of rule IPv6 prefix 2, the overlapping bits between the EA-bits in end user prefix 2 and the overlapping bits in rule prefix 1 (represented by the shaded sections in [Figure 108: IPv6 rule prefix overlap](#)) could render end-user prefixes 1 and 2 indistinguishable (everything else being the same) when anti-spoof lookup is performed in the upstream direction. This could result in an incorrect anti-spoofing lookup.

A similar logic can be applied to overlapping IPv4 prefixes in the downstream direction, where the longest prefix match always leads to the same CE, while the shortest match (leading to a different CE) is not evaluated.

7.25.5 BMR rules implementation example

This section examines an example MAP-T deployment with three MAP rules. The deployment assumes the following:

- There are about 12,000 private IPv4 addresses that need to be translated via MAP-T.
- Each such address should have approximately 4000 ports available per CE. Therefore, the IP address sharing ratio is 16:1; that is, 16 CEs share the same public IP address.

- The public IPv4 addresses that are available to the operator for this translation are from three /24 subnets (10.11.11.0/24, 10.12.12.0/24 and 10.13.13.0/24).
- All users (or CEs) are assigned a /60 IA-PD.

The 12,000 private IPv4 addresses (CEs) in a 16:1 sharing scenario can be covered using three /24 subnets as follows:

$$(3 * 2^8 * 16 = 12,288)$$

The IPv4 rule prefix and EA bits length per rule in this scenario are:

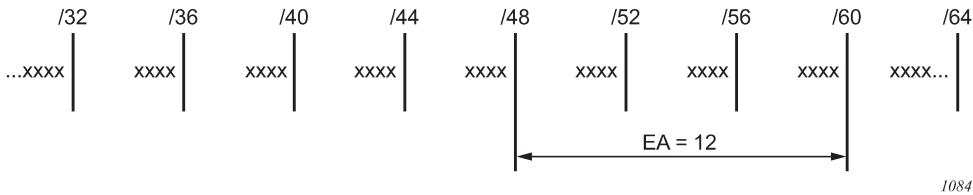
- 10.11.11.0/24 EA length: 12 bits (8 bits for the IPv4 suffix and 4 bits for PSID)
- 10.12.12.0/24 EA length: 12 bits (8 bits for the IPv4 suffix and 4 bits for PSID)
- 10.13.13.0/24 EA length: 12 bits (8 bits for the IPv4 suffix and 4 bits for PSID)

The first 6 bits of the 16 bit port-range are set to 000000 and are reserved for *psid-offset* (ports 0-1023 are reserved); therefore, the user-allocated port space is calculated as follows:

$$4000 - 64 = 4032 \text{ ports}$$

The IPv6 rule prefix is the next parameter in the MAP rule. [Figure 109: Determining the rule IPv6 prefix](#) shows the relevant bits in the IPv6 address: only bits /32 to /64 are considered; the irrelevant bits of the IPv6 addresses are ignored in this example.

Figure 109: Determining the rule IPv6 prefix



The following three rules are created in this example:

- Rule 1 covers subnet 1.
- Rule 2 covers subnet 2.
- Rule 3 covers subnet 3.

In each of the three cases, the EA bits extend from the PD length (/60) to the IPv6 rule prefix length (/48).

The IPv6 rule prefix length is determined for each of the three rules. However, the IPv6 rule prefixes must not overlap, see section [Rule prefix overlap](#) for more information. Non-overlapping IPv6 rule prefixes ensure that each CE is assigned a unique IA-PD. [Table 57: IPv6 rule prefixes](#) describes the rules.

Table 57: IPv6 rule prefixes

	Rule 1	Rule 1	Rule 1
IPv6 rule prefix	2001:db8:0000::/48-	2001:db8:0001::/48-	2001:db8:0002::/48-
IPv4 rule prefix	10.11.11.0/24	10.12.12.0/24	10.13.13.0/24
EA bits	12	12	12
Paid-Offset	6	6	6

The final step is to ensure that the DHCPv6 server hands out correct end-user prefixes (IA-PD), and the rules are also delegated.

In this example, each /48 IPv6 rule prefix supports 4,000 MAP-T CEs, where each CE can further delegate 15 IPv6 "subnets" on the LAN side and each CE is allocated about 4,000 ports to use in stateful NAT44.



Note: The VSR-BR supports only IPv6 rule prefixes of the same length within a domain. To accommodate a different prefix length assignment for IA-PD (for example /56), create another domain with a different IPv6 rule prefix (/44 instead of /48).

7.25.6 ICMP

The following ICMPv4 messages are supported in MAP-T on the VSR; other types of ICMP messages are not supported:

- **ICMP query messages**

These messages contain an identifier field in the ICMP header, which is referred to as the "query identifier" or "query-id" and it is used in MAP-T in the same way as the L4 ports are used in TCP or UDP. ICMP Echo Req/Rep (PING) and traceroute are examples that rely on ICMP Query messages.

- **ICMP error messages**

These messages contain the embedded original datagram that triggers the ICMP error message. The ICMP error messages do not contain the **query-id** field.

The ICMP Query messages and ICMP Error messages are supported regardless of whether they are just passing through a VSR (transit messages), or are terminated or generated in or from a VSR.

The NAT-related ICMPv4 behavior is described in RFC 5508. The following NAT messages are supported in the MAP-T VSR (RFC 5508, §7, Requirement 10a):

- ICMPv4 Error Message: Destination Unreachable Message (Type 3)
- ICMPv4 Error Message: Time Exceeded Message (Type 11)
- ICMPv4 Query Message: Echo and Echo Reply Messages (Type 8 and Type 0)

7.25.7 Fragmentation

The IPv6 header of the IPv4-translated packet in MAP-T can be up to 28 bytes larger than the IPv4 header (40-byte IPv6 header plus 8-byte fragmentation header versus 20-byte IPv4 header). In the case where the IPv4-to-IPv6 translated packet is larger than the IPv6 MTU, the original IPv4 packet is fragmented so that the size of the translated IPv6 packet is within IPv6 MTU. IPv6 packets are never fragmented, although they may contain the fragmentation header that carries fragmentation information related to the original IPv4 packet/fragment.

The IPv6 MTU in the VSR is configurable for each MAP-T domain. The L2 header is excluded from the IPv6 MTU.

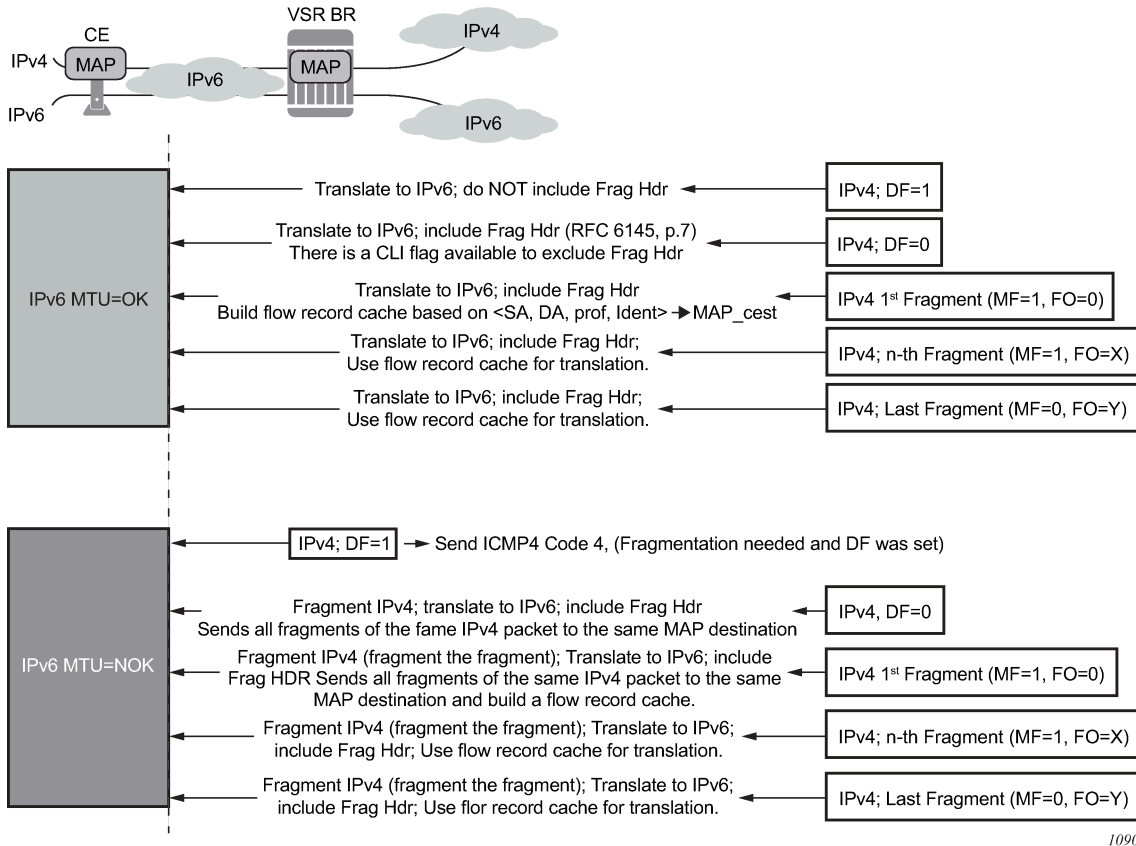
7.25.7.1 Fragmentation in the downstream direction

All fragments of the same IPv4 packet are translated and sent toward the same CE. As the second and consecutive fragments do not contain any port information, the translation is performed based on the <SA, DA, Prot, Ident> cached flow records extracted from the IPv4 header.

Note that the VSR may further fragment an IPv4 fragment that it has received to fit it within the IPv6 MTU.

Figure 110: Fragmentation in the downstream direction shows downstream fragmentation scenarios.

Figure 110: Fragmentation in the downstream direction

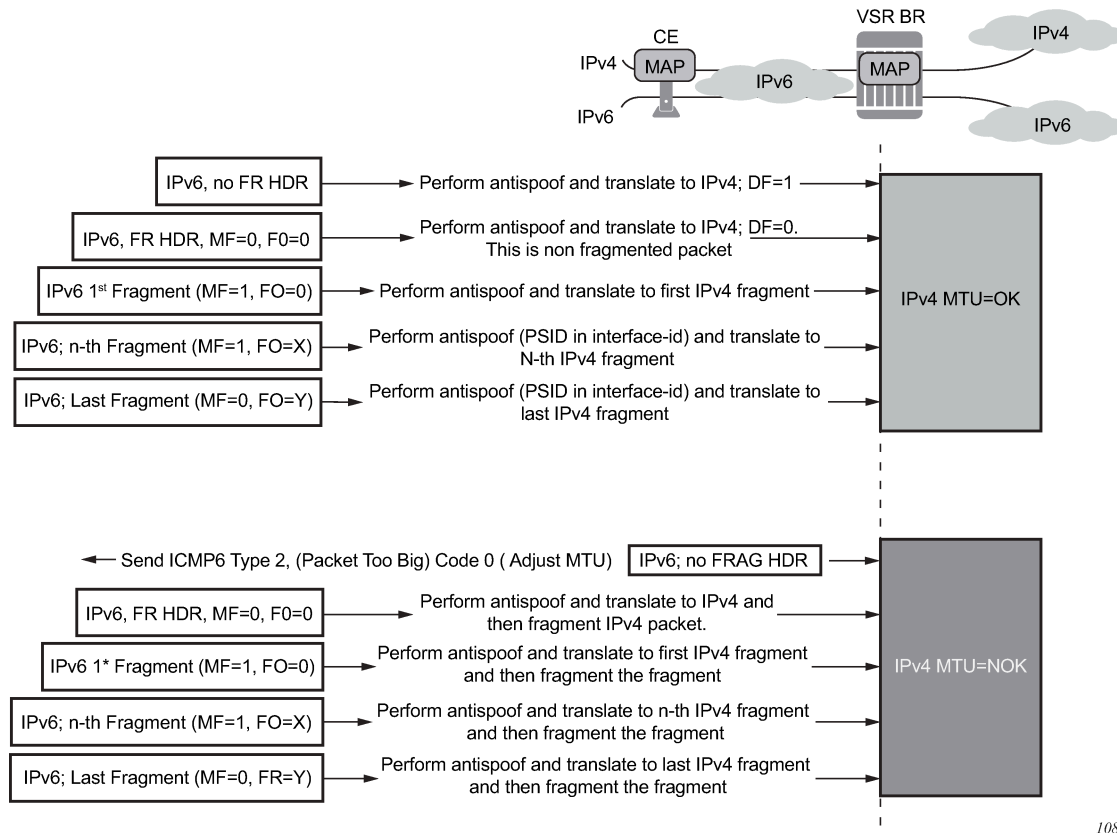


1090

7.25.7.2 Fragmentation in the upstream direction

In the upstream direction, the received IPv6 fragments are artifacts of the IPv4 packets being fragmented on the CP side, before they are translated into IPv6. No flow caching is performed in the upstream direction. The BR performs an anti-spoof for each fragment and if the anti-spoof is successful, the fragment is translated to IPv4. Figure 111: Fragmentation in the upstream direction shows the upstream fragmentation scenario.

Figure 111: Fragmentation in the upstream direction



1088

7.25.7.3 Fragmentation statistics

Fragmentation statistics can be cleared using the **clear nat map frag-stats** command. The following fragmentation statistics are available:

- **Rx Resolved Frags**

This counter shows fragments that were resolved and never buffered; for example:

- first fragments (MF=1, FO=0), which are always resolved by definition
- non-first fragments with matching flow records

- **Rx Unresolved Frags**

This counter shows the number of packets that were queued in the system since the last clear command was invoked. For example, packets with out-of-order fragments without a matching flow record (missing 1st fragment) can be eventually resolved and forwarded, or discarded (for example, because of timeout).

- **Tx Frags**

This counter shows the fragments that were transmitted (**Rx Resolved** and **Rx Unresolved** that were eventually resolved) out of fragmentation logic within the VSR. There is no guarantee that the fragments are transmitted out of the system as they may be dropped on egress because of congestion or restrictions imposed by the configured filter.

- **Dropped Frags**

This counter represents the fragments that are dropped because of fragmentation issues such as timeout, buffer full, and so on.

- **Buffers in Use**

This counter represents the amount of buffered fragments expressed as a percentage of the maximum buffer space that can be used for fragmentation.

- **Max Buffers**

This is a non-cumulative counter that represents the maximum number of buffers allocated since the last **clear** command. The counter captures the highest value of the **buffers-in-use** counter since the last **clear** command. The unit of this counter is the percentage of the total buffer space that can be used by fragmentation.

- **Created Flows**

This is a cumulative counter that represents the total number of flow records since the last **clear** command was invoked. It only counts the first fragment. It represents the number of fragmented packets that were processed by the system since the last **clear** command. This counter does not indicate the number of flows (packets whose fragments were transmitted fully) that were actually transmitted.

- **Flows in Use**

The counter gives an approximation of the number of flow records currently in use and the number of fragmented packets being processed at the time the counter was invoked, as a percentage.

- **Max Flows**

A non-cumulative counter that represents the maximum number of flow records reached since the last **clear** command. The counter shows the highest value of the **flows-in-use** counter since the last **clear** command, as a percentage.

- **Flow Collisions**

This counter represents the number of overlapping first fragments. For example, in the case where a flow record already exists and another first fragment for the flow is received.

- **Exceeded Max Timeouts**

This counter shows the number of fragments that have timed out since the last **clear** command. The represented fragments are:

- Rx unresolved (buffered) fragments that have timed out because of a missing first fragment
- deleted flow-records because they have not received all fragments within the timeout period

- **Exceeded Max Flows**

This counter represents the number of times that the flows in the system has exceeded the maximum supported value.

- **Exceeded Max Buffers**

This counter represents the number of times that the buffers in the system has exceeded the maximum supported value.

- **Exceeded Max Buffers Per Flow**

This counter represents the number of times that the fragment counter per flow has exceeded its limit.

7.25.8 Maximum Segment Size (MSS) adjust

The MSS Adjust feature is used to prevent fragmentation of TCP traffic. The TCP synchronize/start (SYN) packets are intercepted and their MSS value inspected to ensure that it conforms with the configured MSS value. If the inspected value is greater than the value configured in the VSR BR, the MSS value in the packet is lowered to match the configured value before the TCP SYN packet is forwarded.

As the end nodes governing the MSS value are IPv4 nodes, this feature is supported for IPv4 packets only. An MSS adjust is performed in both the upstream and downstream directions.

7.25.9 Statistics collection

The VSR BR maintains a count of the forwarded and dropped packets/octets per MAP-T-domain per direction. The statistics are collected on ingress (upstream v6 and downstream v4) and stored in 64-bit counters.

7.25.10 Logging

As with any NAT operation where the identity of the user is hidden behind the NAT identity, logging of the NAT translation information is required. In the MAP-T domain, NAT logging is based on configuration changes because the user identity can be derived from the configured rules.

A system can have a large number of rules and each configured MAP rule generates a separate log. As a result, the amount of logs generated can be substantial. Logging is explicitly enabled using a log event.

A NAT log contains information about the following:

- MAP type (map-t)
- map-domain name
- map-rule name
- v6 rule-prefix
- v4 rule-prefix
- EA bits
- psid-offset bits
- associated routing context for the MAP-T rule
- timestamp

A MAP rule log is generated when both of the following conditions are met:

- a MAP rule is activated and deactivated in the system (administratively **shutdown/no shutdown**, corresponding MAP domain is associated/dissociated from the routing context, corresponding MAP domain is **shutdown /no shutdown**, and so on)
- event **tmnxNatMapRuleChange (id=2036)** has been enabled in event-control

Example:

```
551 2016/04/22 14:56:35.44 UTC MINOR: NAT #2036 vprn220 NAT MAP
"map-type=map-t map-domain=domain-name-1 rule-name=rule-name-1 rule-prefix=2001:db8::/44
ipv4-prefix=192.168.10.0/24 ea-length=12 psid-offset=6 enabled router=vprn220 at 2016/04/22
14:56:35"
```

7.25.11 Licensing

A valid MAP-T license is required to enable the MAP-T functionality in the VSR BR. A MAP-T domain can only be instantiated with the appropriate license, which enables the following CLI command:

```
configure
  service
    vprn <id> customer <cust-id> create
      map-domain <domain-name>
```

7.25.12 Configuration

The MAP-T configuration consists of defining MAP-T parameters within a template. The MAP-T domain is then instantiated by applying (referencing) this template within a routing (router or VPRN) context.

Defining a MAP domain template:

```
configure
  service
    nat
      map-domain <domain-name> create
        [no] shutdown
        dmr-prefix <ipv6-prefix>
        tcp-mss-adjust <segment-size>
        mtu <mtu-size>
        ip-fragmentation
          [no] v6-frag-header
        mapping-rule <rule-name>
          [no] shutdown
          rule-prefix <ipv6-prefix>
          ipv4-prefix <ipv4-prefix>
          ea-length <ea-bits-length>
          psid-offset <psid-offset-len>
        :
        mapping-rule <rule-name>
          [no] shutdown
          rule-prefix <ipv6-prefix>
          ipv4-prefix <ipv4-prefix>
          ea-length <ea-bits-length>
          psid-offset <psid-offset-len>
        :
        up to 256 rules per domain
```

MAP-T domain instantiation:

```
configure
  service vprn <id> | router
    nat
      map-t
        map-domain <domain-name>
```

MAP domain example template:

The following example shows the MAP domain template for the BMRs defined [BMR rules implementation example](#).

```
configure
```

```

service
  nat
    map-domain domain_1 create
      no shutdown
      dmr-prefix 2001:db8:0100::/64
      mapping-rule rule_1
        no shutdown
        rule-prefix 2001:db8:0000::/48
        ipv4-prefix 10.11.11.0/24
        ea-length 12
        psid-offset 6
      mapping-rule rule_2
        no shutdown
        rule-prefix 2001:db8:0001::/48
        ipv4-prefix 10.12.12.0/24
        ea-length 12
        psid-offset 6
      mapping-rule rule_3
        no shutdown
        rule-prefix 2001:db8:0002::/48
        ipv4-prefix 10.13.13.0/24
        ea-length 12
        psid-offset 6

```

MAP-T domain instantiation example:

The following example shows the MAP-T domain instantiation for the BMRs defined in [BMR rules implementation example](#).

```

configure
  service
    vprn 10
      nat
        map-t
          map-domain domain_1

```

7.25.12.1 Modifying MAP-T parameters when the MAP-T domain is active

You can add new rules to an existing MAP-T domain while the MAP-T domain is instantiated and forwarding traffic. However, each rule must be in the **shutdown** state before any of its parameters are modified.

A MAP-T domain must be in the **shutdown** state to modify the *dmr-prefix* parameter. The remaining parameters (**tcp-mss-adjust**, **mtu**, **ip-fragmentation**) can be modified while the domain is active.

A MAP domain does not have to be in a **shutdown** state when rule modification is in progress.

7.25.13 Inter-chassis redundancy

MAP natively provides multi-chassis redundancy through the use of the anycast BR prefix that is advertised from multiple nodes.

As there is no state maintenance in the MAP-T BR, any BR node can process traffic for the same domain at all times. The only traffic interruption during the switch-over is for the fragmented traffic in the downstream direction being handled at the time of switchover (the flow record cache is not synchronized between the nodes).

7.26 Configuring NAT

This section provides information to configure NAT using the command line interface.

7.26.1 ISA redundancy

The 7750 SR supports ISA redundancy to provide reliable NAT even when an MDA fails. The **active-mda-limit** command allows an operator to specify how many MDAs are active in a NAT group. Any number of MDAs configured above the active-mda-limit are spare MDAs; they take over the NAT function if one of the current active MDAs fail.

A sample configuration is as follows:

```
configure
  isa
    nat-group 1 create
      active-mda-limit 1
      mda 1/2
      mda 2/2
      no shutdown
    exit
  exit
exit
```

Show commands are available to display the actual state of a nat-group and its corresponding MDAs:

```
show isa nat-group 1
=====
ISA NAT Group 1
=====
Admin state       : inService           Operational state : inService
Active MDA limit  : 1                   Reserved sessions : 0
High Watermark (%) : (Not Specified)    Low Watermark (%) : (Not Specified)
Last Mgmt Change  : 01/11/2010 15:05:36
=====
ISA NAT Group 1 members
=====
Group Member      State      Mda  Addresses  Blocks  Se-% Hi Se-Prio
-----
1      1      active    1/2  0         0         0   N  0
-----
No. of members: 1
=====
```

A nat-group cannot become active (no shutdown) if the number of configured MDAs is lower than the active-mda-limit.

An MDA can be configured in several nat-groups but it can only be active in a single nat-group at any moment in time. Spare MDAs can be shared in several nat-groups, but a spare can only become active in one nat-group at a time. Changing the active-mda-limit, adding or removing MDAs can only be done when the nat-group is shutdown.

Nat-groups that share spare MDAs must be configured with the same list of MDAs. It is possible to remove/add spare MDAs to a nat-group while the nat-group is admin enabled.

```
Configure
  isa
    nat-group 1 create
      active-mda-limit 1
      mda 1/2
      mda 2/2
      mda 3/1
      no shutdown
    exit
    nat-group 2 create
      active-mda-limit 1
      mda 1/2
      mda 2/2
      mda 3/1
      no shutdown
    exit
  exit
exit
```

Through show commands, it is possible to display an overview of all the nat-groups and MDAs.

```
show isa nat-group
=====
ISA NAT Group Summary
=====
Mda  Group 1          Group 2
-----
1/1  active            busy
2/2  busy               active
3/1  standby            standby
=====
```

If an MDA fails, the spare (if available) takes over. All active sessions are lost, but new incoming sessions makes use of the spare MDA.

In case of an MDA failure in a nat-group without any spare MDA, all traffic toward that MDA is black-holed.

For L2-aware NAT, the operator has the possibility to clear all the subscribers on the affected MDA (clear nat isa), terminating all the subscriber leases. New incoming subscribers make use of the MDAs that are still available in the nat-group.

7.26.2 NAT Layer 2-Aware configurations

The following sections provide NAT Layer 2-Aware configurations.

```
#-----
echo "Card Configuration"
#-----
  card 1
    card-type iom4-e-b
    mda 1
      mda-type m48-1gb-xp-tx
    exit
    mda 2
      mda-type isa-bb
    exit
```

```

exit
card 2
  card-type iom4-e-b
  mda 1
    mda-type m48-lgb-xp-tx
  exit
  mda 2
    mda-type isa-bb
  exit
exit

#-----
echo "ISA Configuration"
#-----
  isa
    nat-group 1 create
      description "1 active + 1 spare"
      active-mda-limit 1
      mda 1/2
      mda 2/2
      no shutdown
    exit
  exit

#-----
echo "Router (Network Side) Configuration"
#-----
  router
    ...

#-----
echo "NAT (Network Side) Configuration"
#-----
  nat
    outside
      pool "pool1" nat-group 1 type l2-aware create
      address-range 10.81.0.0 10.81.0.200 create
      exit
      no shutdown
    exit
  exit
exit

#-----
echo "Service Configuration"
#-----
  service
    customer 1 create
      description "Default customer"
    exit
    ...
    vprn 100 customer 1 create
      ...
      nat
        outside
          pool "pool2" nat-group 1 type l2-aware create
          address-range 10.0.0.0 10.0.0.200 create
          exit
          no shutdown
        exit
      exit
    exit
  exit

  vprn 101 customer 1 create
    ...
    nat

```

```

        inside
            l2-aware
                # Hosts in this service with IP addresses in these ranges
                # will be subject to l2-aware NAT.
                address 10.0.0.1/29
                address 10.1.0.1/29
            exit
        exit
    exit
exit
...
nat
    nat-policy "l2-aware-nat-policy1" create
        pool "pool1" router Base
    exit
    nat-policy "l2-aware-nat-policy2" create
        pool "pool2" router 100
    exit
exit
...
exit
#-----
echo "Subscriber-mgmt Configuration"
#-----
subscriber-mgmt
    # Subscribers using these sub-profiles will be subject to l2-aware NAT.
    # The configured nat-policies will determine which IP pool will be used.
    sub-profile "l2-aware-profile1" create
        nat-policy "l2-aware-nat-policy1"
    exit
    sub-profile "l2-aware-profile2" create
        nat-policy "l2-aware-nat-policy2"
    exit
    ...
exit

```

7.26.3 Large scale NAT configuration

The following sections provide Large Scale NAT configuration examples.

```

configure
#-----
echo "Card Configuration"
#-----
card 3
    card-type imm-2pac-fp3
    mda 1
        mda-type isa-bb
    exit
    mda 2
        mda-type isa-bb
    exit
exit
#-----
echo "ISA Configuration"
#-----
isa
    nat-group 1 create
        active-mds-limit 2
        mda 3/1
        mda 3/2

```



```

        no shutdown
    exit
exit
exit
#-----
echo "Filter Configuration"
#-----
filter
    ip-filter 123 create
        entry 10 create
            match
                src-ip 10.0.0.1/8
            exit
        action nat
    exit
exit
exit
exit
#-----
echo "NAT (Declarations) Configuration"
#-----
service
    nat
        nat-policy "ls-outPolicy" create
    exit
exit
exit
#-----
echo "Service Configuration"
#-----
service
    customer 1 create
        description "Default customer"
    exit
    vprn 500 customer 1 create
        interface "ip-192.168.113.1" create
        exit
        nat
            outside
                pool "nat1-pool" nat-group 1 type large-scale create
                port-reservation ports 200
                address-range 10.81.0.0 10.81.6.0 create
                exit
                no shutdown
            exit
        exit
    exit
    vprn 550 customer 1 create
        interface "ip-192.168.13.1" create
        exit
    exit
    nat
        nat-policy "ls-outPolicy" create
        pool "nat1-pool" router 500
        timeouts
            udp hrs 5
            udp-initial min 4
        exit
    exit
    vprn 500 customer 1 create
        router-id 10.21.1.2
        route-distinguisher 500:10
        vrf-target export target:500:1 import target:500:1
        interface "ip-192.168.113.1" create

```

```

        address 192.168.113.1/24
        static-arp 192.168.113.5 2001:db8:01:01:00:01
        sap 1/1/1:200 create
        exit
    exit
    no shutdown
exit
vprn 550 customer 1 create
    router-id 10.21.1.2
    route-distinguisher 550:10
    vrf-target export target:550:1 import target:550:1
    interface "ip-192.168.13.1" create
        address 192.168.13.1/8
        sap 1/2/1:900 create
            ingress
            filter ip 123
        exit
    exit
exit
nat
    inside
        nat-policy "ls-outPolicy"
    exit
exit
no shutdown
exit
exit
exit all

```

7.26.4 NAT configuration examples

The following output displays example configurations.

VPRN service example:

```

configure service vprn 100 nat
    inside
        nat-policy "priv-nat-policy"
        destination-prefix 0.0.0.0/0
        dual-stack-lite
            subscriber-prefix-length 128
            address 2001:db8:470:1F00:FFFF:190
            tunnel-mtu 1500
        exit
        no shutdown
    exit
    redundancy
        no peer
        no steering-route
    exit
    subscriber-identification
        shutdown
        no attribute
        no description
        no radius-proxy-server
    exit
    l2-aware
    exit
exit
outside
    no mtu

```

```
exit
```

Router NAT example:

```
configure router nat
  outside
    no mtu
    pool "privpool" nat-group 3 type large-scale create
      no description
      port-reservation blocks 128
      port-forwarding-range 1023
      redundancy
        no export
        no monitor
      exit
      subscriber-limit 65535
      no watermarks
      mode auto
      address-range 10.0.0.5 10.0.0.6 create
        no description
        no drain
      exit
      no shutdown
    exit
    pool "pubpool" nat-group 1 type large-scale create
      no description
      port-reservation blocks 1
      port-forwarding-range 1023
      redundancy
        no export
        no monitor
      exit
      subscriber-limit 65535
      no watermarks
      mode auto
      address-range 192.168.8.241 192.168.8.247 create
        no description
        no drain
      exit
      no shutdown
    exit
  exit
```

Service NAT example:

```
configure service nat
  nat-policy "priv-nat-policy" create
    alg
      ftp
      rtsp
      sip
    exit
    block-limit 4
    no destination-nat
    no description
    filtering endpoint-independent
    pool "privpool" router Base
    no ipfix-export-policy
    port-limits
      forwarding 64
      no reserved
      no watermarks
    exit
```

```
priority-sessions
exit
session-limits
  max 65535
  no reserved
  no watermarks
exit
timeouts
  icmp-query min 1
  sip min 2
  no subscriber-retention
  tcp-established hrs 2 min 4
  tcp-syn sec 15
  no tcp-time-wait
  tcp-transitory min 4
  udp min 5
  udp-initial sec 15
  udp-dns sec 15
exit
no tcp-mss-adjust
no udp-inbound-refresh
exit
nat-policy "pub-nat-policy" create
alg
  ftp
  no rtsp
  no sip
exit
block-limit 1
no destination-nat
no description
filtering endpoint-independent
pool "pubpool" router Base
no ipfix-export-policy
port-limits
  no forwarding
  no reserved
  no watermarks
exit
priority-sessions
exit
session-limits
  max 65535
  no reserved
  no watermarks
exit
timeouts
  icmp-query min 1
  sip min 2
  no subscriber-retention
  tcp-established hrs 2 min 4
  tcp-syn sec 15
  no tcp-time-wait
  tcp-transitory min 4
  udp min 5
  udp-initial sec 15
  udp-dns sec 15
exit
no tcp-mss-adjust
no udp-inbound-refresh
exit
```

7.27 Configuring VSR-NAT

This section provides information about the VSR-NAT functionality, including licensing requirements, statistics collection, and examples of **show** command output.

7.27.1 VSR-NAT licensing

Appropriate licensing is required to enable the VSR-NAT functionality in the system. However, no further licensing enforcement is performed based on resource utilization, such as the consumed bandwidth or the number of NAT bindings.


The following NAT-related functionality is enabled through licensing:


- LSN (LSN44, DS-Lite, and NAT64)
- L2-Aware NAT
- UPnP
- geo-redundancy

You can use the CLI or MIB on VSR-NAT to get more information about the number of LSN bindings and LSN bandwidth.

[Table 58: NAT licenses required to unlock NAT functionality](#) describes the licenses required to unlock the VSR-NAT functionality.

Table 58: NAT licenses required to unlock NAT functionality

NAT license title	Functionality enabled	License purchased
LSN	LSN Pool <ul style="list-style-type: none"> • configure router nat outside pool <i>name type large-scale</i> • configure service nat outside pool <i>name type large-scale</i> 	The following two scaling licenses are required: <ul style="list-style-type: none"> • license for the number of LSN bindings • license for consumed bandwidth You must purchase both licenses to enable the LSN functionality.
L2AWARE	L2Aware Pool <ul style="list-style-type: none"> • configure router nat outside pool <i>name type l2-aware</i> • configure service nat outside pool <i>name type l2-aware</i> 	Purchase the L2-Aware license to enable the functionality. The LSN scaling license is not required.  Note: The L2-Aware NAT functionality can only be used with the VBNG.
UPnP	UPnP commands: configure subscriber-mgmt sub-profile <i>sub-prof-name upnp-policy upnp-pol-name</i>	Purchase the UPnP license to enable the functionality.

NAT license title	Functionality enabled	License purchased
		 Note: The UPnP functionality can only be used with the L2-Aware NAT.
GEO REDUNDANCY	Geo-redundancy Pool <ul style="list-style-type: none"> • configure router nat outside pool redundancy • configure service nat outside pool redundancy 	Purchase the Geo Redundancy license to enable the functionality.

7.27.2 Statistics collection For LSN bindings

A NAT subscriber is an internal entity whose true identity is hidden outside the network. The NAT subscriber is represented by a binding that is a set of stateful mappings between the internal and external representations of the subscriber. From the licensing perspective, the terms "NAT bindings" and "NAT subscribers" can be used interchangeably.

VSR-NAT collects the number of LSN subscribers for licensing purposes; the L2-Aware NAT subscribers are excluded from this count. An LSN subscriber is defined as follows:

- **Large Scale NAT44 (or CGN)**

The subscriber is an internal IPv4 address.

- **DS-Lite**

The subscriber is identified by the CPE IPv6 address (B4 element) or an IPv6 prefix. The selection of the address or prefix as the representation of a DS-Lite subscriber is configuration-dependent.

- **NAT64**

The subscriber is an IPv6 address.

The number of LSN subscribers (LSN44, DS-Lite, and NAT64) in VSR-NAT is sampled every hour on the hour (for example, at 00:00 am, 01:00 am, 02:00 am, and so on). Each sample is a snapshot of the number of subscribers at the time that the statistics are collected.

The CLI can be used to view the following information:

- 24 samples (one per hour) in the current day
- Maximum value for each of the last 7 days
- Average value for each of the last 7 days
- Maximum value since the system booted

For the list of CLI commands available for use, see [VSR-NAT show command examples](#).

7.27.3 Statistics collection for LSN bandwidth

The measurement of LSN bandwidth includes translated packets and octets in the upstream and downstream direction. Packets that are rejected for any reason and traffic carrying logging information are both excluded from the statistics.

LSN bandwidth statistics for VSR-NAT are collected every 10 minutes. The bandwidth is derived as a difference in octet count between the two consecutive collection intervals, divided by a 10 minute interval. There is no bandwidth differentiation per LSN type (LSN44, DS-Lite, and NAT64) or per direction. Aggregate bandwidth values per node are maintained in kb/s units. L2-Aware NAT and WLAN gateway statistics are not included in the statistics collection.

The CLI can be used to view the following LSN bandwidth information:

- 144 bandwidth values for the current day (bandwidth statistics are collected every 10 minutes)
- Maximum bandwidth value for each of the last 7 days
- Average bandwidth value for each of the last 7 days
- Maximum bandwidth value since the system booted

For the list of CLI commands available for use, see [VSR-NAT show command examples](#).

7.27.4 VSR-NAT show command examples

The following CLI commands are available for use:

- **show system license-statistics 24-hours application nat**
- **show system license-statistics week application nat**
- **show system license-statistics peak application nat**

The following output shows examples of NAT statistics.

Weekly display example:

```
*A:Dut-A>show system license-statistics week application nat
=====
week license statistics for nat
=====
index      time                average             peak
-----
LSN subscribers
1          2016/02/01 00:00:00 370                456
2          2016/01/31 00:00:00 375                512
3          2016/01/30 00:00:00 374                510
4          2016/01/29 00:00:00 373                478
5          2016/01/28 00:00:00 360                450
6          2016/01/27 00:00:00 370                496
7          2016/01/26 00:00:00 373                503
LSN bandwidth
1          2016/02/01 00:00:00 12472623          12472623
2          2016/01/31 00:00:00 12472623          12472623
3          2016/01/30 00:00:00 12472623          12472623
4          2016/01/29 00:00:00 12472623          12472623
5          2016/01/28 00:00:00 12472623          12472623
6          2016/01/27 00:00:00 12472623          12472623
7          2016/01/26 00:00:00 12472623          12472623
-----
No. of license statistics entries: 14
=====
```

24-hour display example:

```
*A:Dut-A# show system license-statistics 24-hours application nat
=====
```

```

24 hours license statistics for nat
=====
index      time                value
-----
LSN subscribers
1          2016/06/29 19:00:00 512
2          2016/06/29 20:00:00 512
LSN bandwidth
1          2016/06/29 18:10:00 0
2          2016/06/29 18:20:00 0
3          2016/06/29 18:30:00 0
4          2016/06/29 18:40:00 2996286
5          2016/06/29 18:50:00 12472524
6          2016/06/29 19:00:00 12472424
7          2016/06/29 19:10:00 12471020
8          2016/06/29 19:20:00 12471980
9          2016/06/29 19:30:00 12471566
10         2016/06/29 19:40:00 12471881
11         2016/06/29 19:50:00 12472116
12         2016/06/29 20:00:00 12472623
-----
No. of license statistics entries: 14
=====
    
```

Peak display example:

```

=====
*A:Dut-A# show system license-statistics peak application nat
=====
peak license statistics for nat
=====
                                time                peak
-----
LSN subscribers                2016/06/29 19:00:00 512
LSN bandwidth                   2016/06/29 20:00:00 12472623
-----
No. of license statistics entries: 2
=====
    
```

Table 59: NAT statistics output fields describes the NAT statistics output fields.

Table 59: NAT statistics output fields

Label	Description
Index	The entry number of the displayed value. A weekly display contains 7 entries, one for each of the last 7 days. A 24-hour display can contain up to 24 values for NAT subscribers (statistics are collected hourly) and 144 values for NAT bandwidth (statistics are collected every 10 minutes).
Time	The timestamp of the statistics collection. The bandwidth is averaged in 10 minute intervals. Consequently, bandwidth value at a specific time represents the average bandwidth for the last 10 minute period.

Label	Description
Value	The value for the number of NAT subscribers at a specific time, or the average bandwidth in kb/s for the last 10 minute period.
Average	In the weekly display, the average daily value for the number of NAT subscribers or the NAT bandwidth.
Peak	In the weekly display, the daily peak value for the number of NAT subscribers or the NAT bandwidth.

7.28 VSR scaling profiles on BB-ISA

7.28.1 Scaling profiles on the VSR

To meet flexible scaling requirements in common compute platforms, operators can use the CLI to select VSR-NAT scaling profiles that correspond to the amount of memory allocated in the VM.

The following scaling profiles have predefined upper scaling limits and are available for VSR-NAT and IPv6 FW:

- **profile1** is a lower scaling profile.
- **profile2** is a higher scaling profile.

The default scaling profile is **profile1**.

Contact your Nokia representative for more information about NAT scaling figures in each profile.

For the number of required CPU control cores on a VSR-I in relation to profiles, see to the Sysinfo section in the *Virtualized Service Router Installation and Setup Guide*.

For the amount of required memory on VSR-I in relation to profiles, see Software Release Notes, section VM Memory Requirements by Function Mix.

Scaling profiles are applied under the following CLI hierarchy:

Classic CLI

```
*A:cses-V22# configure isa
*A:cses-V22>config>isa# nat-group 1 create
*A:cses-V22>config>isa>nat-group# scaling-profile
  - scaling-profile <scaling-profile-id>
  <scaling-profile-id> : {profile1|profile2}
```

MD-CLI

```
*[pr:/configure isa nat-group 1]
A:admin@cses-V27# scaling-profile ?
scaling-profile <keyword>
<keyword> - {profile1|profile2}
Default   - profile1
Scaling profile for the NAT group
Warning: Modifying this element toggles
'configure isa nat-group 1 admin-state' automatically for the new value to
take effect.
Warning: Modifying this element clears ISA state, such as flow state, for
the new value to take effect.
*[pr:/configure isa nat-group 1]
```

7.28.2 Scale profile modification

A scaling profile can be changed only when the NAT group is in a shut down state. After the scaling profile is changed and the NAT group is activated (**no shutdown**), the system tries to allocate necessary memory. If successful, the vISA transitions inservice; if unsuccessful (for example, if there is not enough memory in the system), the NAT group remains in the shut down state.

Without sufficient resources to accommodate the required scaling profile, the vMDA where the vISA resides transitions into a failed state, followed by logs describing the failure:

```
33 2018/05/28 09:36:30.484 UTC MAJOR: CHASSIS #2001 Base Mda 1/1
"Class MDA Module : failed, reason: Insufficient memory to boot"

36 2018/09/18 14:13:02.426 UTC MAJOR: CHASSIS #2001 Base Mda 1/1
"Class MDA Module : failed, reason: Insufficient mgmt cores to boot"
```

7.29 NAT scaling profiles on ESA

7.29.1 Scaling profiles for NAT on ESA

NAT on ESA offers the following scaling profiles, each of which are adapted to the amount of memory allocated to the VM:

- **profile1** – a lower scaling profile that requires 8 CPU cores and 32 GB of DRAM memory per ESA-VM
- **profile2** – a medium scaling profile that requires 11 CPU cores and 96 GB of DRAM memory per ESA-VM
- **profile3** – a high scaling profile that requires 15 CPU cores and 115 GB of DRAM memory per ESA-VM

The default scaling profile is **profile1**.

Contact your Nokia representative for more specific information about the NAT scaling figures in each profile.

Use the following command to configure scaling profiles.

```
configure isa nat-group scaling-profile
```

7.29.2 Scale profile modification

A scaling profile can be changed using CLI only when all ESA-VMs in a NAT group are removed from the configuration.

For example, in the following case, transitioning from **profile2** to **profile1** is not possible until **esa-vm 1/1** is removed from the CLI:

```
config
  isa
    nat-group 1
      esa-vm 1/1
        scaling-profile profile2
config>isa>nat-group# scaling-profile profile1
MINOR: BBGRPMGR #1052 Cannot change scale-profile with MDAs or ESA-VMs provisioned
```

7.30 Expanding a NAT group

Adding or removing an MDA from a NAT group affects all currently active subscribers and may invalidate existing static port forwards and mappings configured in deterministic NAT.

Store configurations offline before removing the configuration as part of the NAT group modification process that is described in the following information. You can restore the configuration to the node after the change is complete.

The procedure to add or remove an MDA from a NAT group is described in the following information.

Adding and removing an MDA from a NAT group in the MD-CLI:

1. Administratively disable deterministic prefix policies and delete their mappings. Perform this for every deterministic prefix and their mapping used in a NAT group in which the size is modified. Store the deterministic mapping configuration offline before removing it and then reapply after the change. When the NAT group size is modified and the deterministic mappings reapplied, the commit may fail. If the commit fails, you must create a new mapping. Use the following command to create a new mapping:

```
tools perform nat deterministic calculate-maps
```

Static port forwards configurations created with the **tools** command are automatically deleted during the commitment of the modified NAT group.

2. Commit the changes.
3. Change the active and failed MDA limit.
4. Commit the changes.
5. Re-apply deterministic mappings and static port forwards.

Adding and removing an MDA from a NAT group in the classic CLI:

1. Shut down the NAT group.
2. Remove all statically configured large-scale subscribers (such as deterministic, LI, debug, and subscriber-aware) in a NAT group that is being modified.
3. A static port forward configuration created via the **tools** commands is automatically deleted.
4. Manually delete any static port forward configurations created via classic CLI.
5. Shut down and remove the deterministic policies.
6. Delete NAT policy references in all inside routing contexts associated with the NAT group that is being modified.
7. Reconfigure the **active** and **failed-mds-limit** options.
8. Use the following command to enable the NAT group.

```
no shutdown
```

9. Restore previously removed NAT group references in all of the inside routing contexts associated with the modified NAT group.
10. Reapply the subscriber-aware and deterministic subscribers (prefixes and maps), static port-forwards, LI, and debug.

8 Residential firewall

8.1 Residential firewall overview

The residential firewall protects a home by tracking all flows to or from the home. Only inbound traffic that matches flows that originated inside of the home is allowed to pass through the firewall. By blocking other flows, an attacker cannot initiate a connection to a vulnerable service within the home. The residential firewall also provides protection against fingerprinting, port scanning, and DoS attacks. The dynamic flow tracking functionality provides a better user experience compared to static firewall rules because it does not limit any connection that has been set up within the home.

The residential firewall is based entirely on the tracking of Layer 3 and Layer 4 flows. Minimal application layer gateway (ALG) support is provided to allow protocols that use multiple flows, but application layer protection is not supported. The firewall only supports IPv6 flows. It is recommended to use Layer 2-aware NAT to provide similar protection for IPv4 flows within the same residential subscriber.

8.1.1 Supported protocols and extension headers

The residential firewall distinguishes between known and unknown protocols or known and unknown extension headers.

Unknown protocols create or match flows based only on Layer 3 information. For a known protocol, the firewall inspects Layer 4 information to create or match flows more precisely. The following known protocols are supported:

- TCP
- UDP
- ICMPv6

Known extension headers are allowed by the firewall and processing continues on the remainder of the packet. The following extension headers are treated as known:

- Hop-by-hop (0)
- Fragment Header (44)
- Authentication Header (51)
- Destination Options (60)
- Shim Header (140)

8.1.1.1 Unknown protocols

Unknown protocols are created and matched by a 3-tuple identifier that has the format <source IP, destination IP, protocol>. No Layer 4 data is used to differentiate between possible sub-flows. Because the firewall is unaware of unknown protocol states, removal of flows with unknown protocols is only governed by a single configurable timeout.

8.1.1.2 TCP and UDP

TCP and UDP flows are created and matched by a 5-tuple identifier that has the format <source IP, destination IP, protocol, source port, destination port>. Multiple configurable timeouts can apply depending on the exact flow state.

8.1.1.3 ICMPv6

ICMPv6 error messages (codes up to 127) are handled based on the encapsulated invoking packet. Layer 3 and Layer 4 information is re-extracted from the packet and is used to perform a flow lookup. If an existing flow is found, then the error message is forwarded; otherwise, it is dropped.

ICMPv6 echo flows are created and matched by a 4-tuple identifier that has the format <source IP, destination IP, protocol, identifier>. Echo replies must always match an existing flow. A single configurable timeout applies to these flows.

Other informational or non-transit ICMPv6 messages are dropped by the firewall.

8.1.2 Application Layer Gateway

Application layer gateways (ALGs) are used to track protocols where one flow triggers the creation of several associated flows. For example, a single session initiation protocol (SIP) session can trigger several additional media connections. These flows are not always triggered from inside the home, but traffic should still be allowed to pass. To support this, the residential firewall creates additional flows when a supported ALG connection is recognized and enabled.

8.1.3 Additional filtering control

The residential firewall has two filtering modes that control which action to take when an inbound packet does not match an existing flow.

In address and port-dependent filtering mode, security is considered most important and packets that do not match an existing flow are dropped. This could interfere with the operation of some applications that rely on multiple connections using the same host port.

In endpoint independent filtering mode, application transparency is considered most important. When a packet matches any flow that has the correct protocol and destination IP address, the packet is allowed to pass, and the IP address and port of the foreign endpoint are ignored. The assumption is that the application that triggers the original session may require additional remotely-triggered sessions for correct operation. This can be a security concern when an application with known vulnerabilities is used, as all firewall functionality for that application ceases as soon as the application itself opens one flow. Additionally, this exposes the host to fingerprinting attacks.

In addition to filtering, it is possible to limit the number of sessions, or flows, per subscriber. Sessions can be split into priority and non-priority categories based on their mapped forwarding class. Separate limits apply to each category to avoid starvation of priority sessions by non-priority sessions. This granularity of control helps to protect the firewall and the host against DoS attacks and resource starvation.

8.1.4 TCP MSS adjustment

TCP maximum segment size (MSS) adjustment can be used to clamp the MSS value that is sent during a TCP handshake. If the MSS option is not present, or is bigger than the configured value, then the firewall changes it to the configured value.

This is useful when a low-MTU link is used, such as during tunneling. If the MSS is changed to match the low MTU, IP layer packet fragmentation can be avoided, improving the performance of both the firewall and the end hosts.

8.1.5 Static port forwards and DMZ

The residential firewall supports static port forwards and DMZ to selectively allow inbound network-initiated traffic flows. Static port forwards allow operators to open up a specific subset of traffic. An exact IP address and a protocol must be provided. For TCP and UDP traffic, the system also requires at least one port. A foreign prefix or port may also be provided to limit the pinhole to a specific connection.

DMZ is enabled on a per-host basis and disables the firewall for that specific host. Before traffic can be forwarded on SLAAC hosts, the exact /128 address must be learned, either by DAD snooping, or initial upstream traffic. For security reasons, the system does not send any ND for a completely unknown /128 address for network-initiated flows.

Static port forwards are configured under the AAA Context. See the 7450 ESS, 7750 SR, and VSR RADIUS Attributes Reference Guide for more information.

8.2 Residential firewall provisioning

Residential firewalls are provisioned in three steps.

1. A firewall domain is created in the router or VPRN where the firewall is connected to an unsafe network, such as the Internet. In this domain, a list of prefixes specify which prefixes are subject to firewall rules.
2. A firewall policy is created that specifies operational rules for the firewall and which domain should be used.
3. The firewall policy is linked to an ESM subscriber using the subscriber profile.

```

Node# /configure service vprn 4 firewall
Node>config>service>vprn>firewall# info
-----
        domain "domain_4" nat-group 1 create
            prefix 2001:DB8::/32 create
            exit
            no shutdown
        exit
-----
Node# /configure service nat
Node>config>service>nat# info
-----
        firewall-policy "firewall_4" create
            description "IPv6 Firewall policy for VPRN 4"
            domain router 4 name "domain_4"
            filtering address-and-port-dependent
        exit
-----
Node# /configure subscriber-mgmt

```

```
Node>config>subscr-mgmt# info
-----
sub-profile "profile_1" create
  firewall-policy "firewall_4"
exit
-----
```

8.2.1 Domains and addressing

A firewall domain specifies both the network (router or VPRN) to which a firewall is connected and which IP prefixes in that network are protected by the firewall. Hosts of a firewall-enabled subscriber are automatically protected if they are assigned an IP address from a domain prefix. It is possible to mix protected and unprotected hosts within one subscriber, but unprotected hosts must receive an IP address that is outside of the firewall domain.

The router or VPRN where the firewall domain is configured must not be the same as the router or VPRN where the subscriber is terminated. This function replaces classic ESM wholesale/retail for firewall hosts.

9 TCP MSS adjustment

9.1 Overview

This feature adds support for adjustment of MSS of TCP packets with SYN flag according to access/aggregation network to prevent fragmentation of upstream and downstream TCP packets using ISA-BB.

There are two modes of adjustment operations supported: TCP MSS Adjustment for ESM Hosts, and TCP MSS Adjustment for NAT Services.

For TCP MSS adjust using ISA2-AA, see section [AQP](#) for the AQP rules.

9.2 TCP MSS adjustment for ESM hosts

This feature adds support for adjustment of the MSS size of TCP packets with SYN flag according to the access/aggregation network to prevent fragmentation of upstream and downstream TCP packets using ISA-BB diverted by IP/IPv6 filter actions.

The following ESM host types are supported:

- IPv4/IPv6 IPoE hosts
- locally terminated PPPoE hosts (without L2TP LAC)
- L2TP LNS hosts

The configuration steps are as follows:

1. Create a NAT group used for an MSS adjustment.

```
config>isa
nat-group 1
active-mda-limit 2
mda 1/1
mda 1/2
```

2. Associate the NAT group with a routing instance and configure the MSS value.

```
config>router
config>service>vprn
mss-adjust-group 1 segment-size 1452
```



Note: Unless there are dedicated ISAs or ESAs for MSS adjustment, an existing NAT group or WLAN-GW group can be referenced. If multiple NAT or WLAN-GW groups reference the same ISA or ESA, the NAT or WLAN-GW groups become inactive. MSS adjustment does not function correctly if it references an inactive NAT or WLAN-GW group.

3. Create an IPv4/IPv6 filter to perform an MSS adjust.

```
config>filter>ip-filter>entry
```

```
egress-pbr default-load-balancing
match tcp-syn
action tcp-mss-adjust
config>filter>ipv6-filter>entry
match tcp-syn
action tcp-mss-adjust
```

4. Apply an IPv4/IPv6 filter to the SLA profile.

9.3 TCP MSS adjustment for NAT services

This feature provides MSS adjustment for TCP packets to be translated by NAT services.

The configuration steps are as follows:

1. Create a NAT-group used for NAT services with MSS adjustment.

```
config>isa
nat-group 1
  active-mda-limit 2
  mda 1/1
  mda 1/2
```

2. Create a NAT-policy that also adjusts MSS.

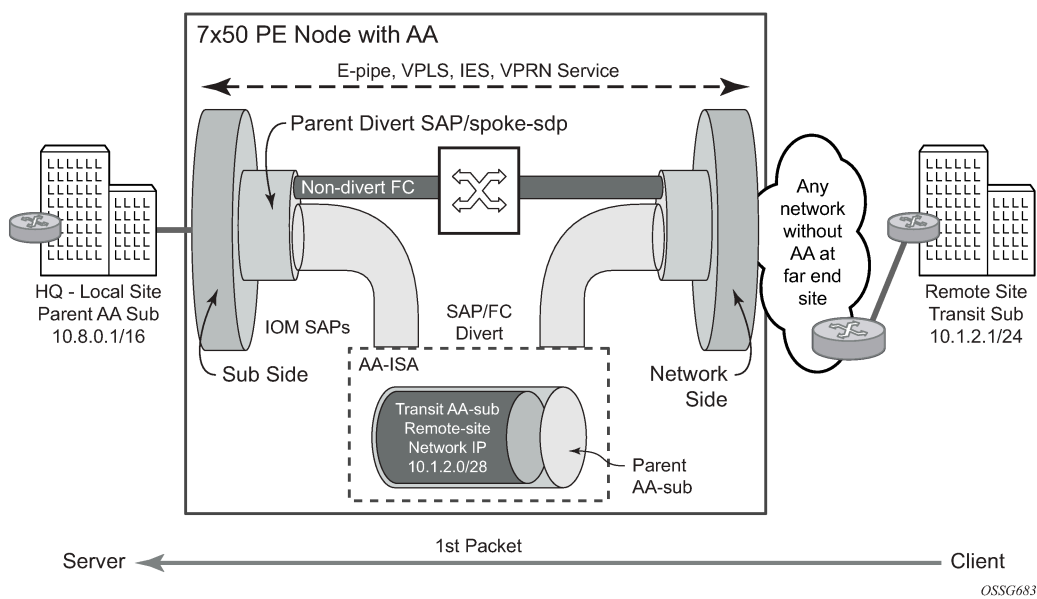
```
config>service>nat
nat-policy "policy-for-mss-adjust" crate
  tcp-mss-adjust 1452
```

10 L2TP network server

10.1 Subscriber agg-rate-limit on LNS

In non-LNS ESM environment, the existing **agg-rate-limit** command is applied to the subscriber within the subscriber profile (sub-profile). However, the **agg-rate-limit** cannot be the highest level in subscriber's HQoS hierarchy. The **agg-rate-limit** is only effective if it is applied to a subscriber that is tied to a port-scheduler. In other words, the port-scheduler in subscriber's HQoS hierarchy is a prerequisite for successful operation of **agg-rate-limit**. On regular MDAs, the port-scheduler is directly applied to a physical port. The port between the carrier IOM and the ISA is an internal port that is not exposed in the CLI. This is shown in [Figure 112: QoS hierarchy on LNS](#).

Figure 112: QoS hierarchy on LNS



OSSG683

The port-scheduler is applied to the internal *Ins-esm* port in the egress direction. The *Ins-esm* egress port is a port between the carrier IOM and the ISA that is passing traffic from all VRFs that have subscriber L2TP sessions terminated in the corresponding ISA.

Use the following CLI syntax to apply the port-scheduler to each *Ins-esm* port:

```
configure
  port-policy <port-policy-name>
    egress-scheduler-policy <port-scheduler-policy-name>
```

Port-policy at the root CLI level creates a port policy manager that can apply various policies (port scheduler) to hidden, dynamically created ports for WLAN GW/LNS/NAT.

CLI syntax:

```
configure
  isa
    lns-group <grp-id>
      mda <card>/<slot>
      mda <card>/<slot>
      :
      port-policy <port-policy-name>
```

The port policy itself is applied to internal LNS port under the lns-group CLI hierarchy. The port scheduler is automatically applied to egress lns-esm ports on carrier IOMs toward every LNS ISA in the lns-group. The port schedulers have the same configuration on every lns-esm port in the lns group but operate independently on each port. Additional consideration:

- An ISA can be assigned to a single lns-group. In other words, two or more LNS-groups cannot contain the same ISA. However, an ISA can belong simultaneously to an LNS-group and a NAT group. The port scheduler affects only LNS traffic.
- The port scheduler rates are wire rates that are based on the encapsulation between the carrier IOM and the ISA which is Ethernet QinQ. However, the queue rates, the billing stats and the agg-rate-limit rates can be optionally based on the last mile encapsulation in the same way as they have been supported in non-LNS environment with **queue-frame-based-accounting** and **encap-offset** commands.

The ability to calculate queue rates or the agg-rate-limit based on the last mile encapsulation is referred to as Last Mile Aware Shaping.

For example, the **encap-offset** command causes the queue rates, the billing stats and the agg-rate-limit to be based on the wire encapsulation in the last mile. For ATM in the last mile, the wire overhead is calculated per each packet (including ATM cellification overhead and padding). For Ethernet in the first mile, a fixed last mile encapsulation (defined with the **encap-offset** command or the RFC 5515, *Layer 2 Tunneling Protocol (L2TP) Access Line Information Attribute Value Pair (AVP) Extensions*) wire overhead is considered in rate calculation. In essence, the length of the PPPoE Ethernet QinQ header that is used on the link between the carrier IOM and the ISA is artificially modified so that it matches the length of the header used in the last mile. The net effect is rate shaping on LNS based on the virtual packet length that is present in the last mile.

The last mile encapsulation information that is used in Last Mile Aware Shaping can be obtained either statically through the explicit value in the **encap-offset** command or dynamically by the RFC 5515 method (AVP 144 in ICRQ). The latter is the case if the **encap-offset** command does not have any explicitly configured value.

In the absence of the **encap-offset** command, the queue rates, the billing stats and the agg-rate-limit rates are based on the Ethernet QinQ encapsulation between the carrier IOM and the ISA. Depending on the queue-frame-based-accounting configuration option, those rates can be wire based or data based (Layer 2 encapsulation only).

- The agg-rate-limit is not applicable to ingress direction (LNS or non LNS ESM).
- V-Port is not applicable in LNS configuration.

10.2 LNS reassembly

10.2.1 Overview

In some cases, PPPoE clients do not honor the negotiated MRU during the LCP phase and consequently, they send packets larger than the negotiated MRU. This applies to control and data packets.

In this case, the LAC fragments IPv4 packets which then have to be reassembled in LNS.

In general, reassembly processing applies only to the end nodes that are receiving fragments. In tunneled environment a fragmented packet must be reassembled before it is de-encapsulated.

10.2.2 Reassembly function

LNS reassembly is implemented through a generic IPv4 reassembly function that can be shared across multiple ISAs in a nat-group. The same ISA can be independently part of an lns-group and a nat-group.

Traffic that needs to be reassembled is steered to the nat-group via filters. After the fragmented traffic is in the nat-group, it is reassembled and injected back within the same routing context to the lns-group for further L2TP processing.

The configuration steps, along with corresponding CLI syntax examples, are as follows:

1. Configure two isa groups, a nat-group providing generic reassembly function and a lns-group providing the L2TP services. The ISAs can be shared amongst the groups, or they can be separated per each group:

```
configure
  isa
    nat-group 1
      active-mda-limit 2
      mda 1/1
      mda 1/1
    lns-group 1
      mda 1/1
      mda 1/2
```

2. Configure redirection of the L2TP traffic to the nat-group performing reassembly:

```
configure
  filter
    ip-filter 10
      entry 5
        match
          dst-ip 10.10.10.10 - traffic classification criteria ; in this case LNS
        tunnel endpoint.
          action reassemble
          default-action forward
```

3. Apply 'reassembly' filter on the incoming L2TP traffic:

```
configure
  router
    interface from-lac
      address 10.0.0.1/24
      port 2/2/2
      ingress
        filter ip 10
```

4. Associate the reassembly context with the same service where LNS is configured:

```

configure
  service
    vprn 10
      reassembly-group 1
        l2tp
          group "lns-vrf-10" create
          ppp
          authentication-policy "lns"
          proxy-authentication
          proxy-lcp
          tunnel "lns-test-tunnel" create
          lns-group 1
          no shutdown

          subscriber-interface "int1" create
          address 10.20.20.254/24
          group-interface "lns-grp-10" lns create
          sap-parameters
          sub-sla-mgmt
          sub-ident-policy "sub-ident"
        dhcp
          server 192.168.1.1
          trusted
          client-applications ppp
          gi-address 10.20.20.1

```

10.2.3 Load sharing between the ISAs

All traffic matching the criteria associated with the filter action reassemble is forwarded to the reassembly function, regardless of whether the traffic is fragmented or not.

In case that there are multiple ISAs in the NAT-group, traffic is load shared between them based on the source IP address and the incoming service ID (routing context).

10.2.4 Inter-chassis ISA redundancy

In case that an active ISA fails in a nat-group, the standby ISA takes over the reassembly function. However, the switchover is not stateful and consequently traffic destined for the failed ISA is lost until it is restarted.

10.3 MLPPPoE, MLPPP(oE)oA with LFI on LNS

MLPPPoX is generally used to address bandwidth constraints in the last mile. The following are other uses for MLPPPoX:

- To increase bandwidth in the access network by bundling multiple links/VCs together. For example it is less expensive for a customer with an E1 access to add another E1 link to increase the access b/w, instead of upgrading to the next circuit speed (E3).
- LFI on a single link to prioritize small packet size traffic over traffic with large size packets. This is needed in the upstream and downstream direction.

PPPoE and PPPoEoA/PPPoA v4/v6 host types are supported.

10.3.1 Terminology

The term MLPPPoX is used to reference MLPPP sessions over ATM transport (oA), Ethernet over ATM transport (oEoA) or Ethernet transport (oE). Although MLPPP in subscriber management context is not supported natively over PPP/HDLC links, the terms MLPPP and MLPPPoX terms can be used interchangeably. The reason for this is that link bundling, MLPPP encapsulation, fragmentation and interleaving can be in a broader scope observed independently of the transport in the first mile.

Terms speed and rate are interchangeably used throughout this section. Usually speed refers to the speed of the link in general context (high or low) while rate usually quantitatively describes the link speed and associates it with the specific value in b/s.

10.3.2 LNS MLPPPoX

This functionality is supported through LNS on BB-ISA. LNS MLPPPoX can be used then as a workaround for PTA deployments, whereby LAC and LNS can be run back-to-back in the same system (connected via an external loop or a VSM2 module), and therefore locally terminate PPP sessions.

MLPPPoX can:

- Increase bandwidth in the last mile by bundling multiple links together.
- LFI/reassembly over a single MLPPPoX capable link (plain PPP does not support LFI).

10.3.3 MLPPP encapsulation

After the MLPPP bundle is created in the 7750 SR, traffic can be transmitted by using MLPPP encapsulation. However, MLPPP encapsulation is not mandatory over an MLPPP bundle.

MLPPP header is primarily required for sequencing the fragments. But in case that a packet is not fragmented, it can be transmitted over the MLPPP bundle using either plain PPP encapsulation or MLPPP encapsulation.

10.3.4 MLPPPoX negotiation

MLPPPoX is negotiated during the LCP session negotiation phase by the presence of the Max-Received-Reconstructed Unit (MRRU) field in the LCP ConfReq. MRRU option is a mandatory field required in MLPPPoX negotiation. It represents the maximum number of octets in the Information field of a reassembled packet. The MRRU value negotiated in the LCP phase must be the same on all member links and it can be greater or lesser than the PPP negotiated MRU value of each member link. This means that the reassembled payload of the PPP packet can be greater than the transmission size limit imposed by individual member links within the MLPPPoX bundle. Packets are always be fragmented so that the fragments are within the MRU size of each member link.

Another field that could be optionally present in an MLPPPoX LCP Conf Req is an Endpoint Discriminator (ED). Along with the authentication information, this field can be used to associate the link with the bundle.

The last MLPPPoX negotiated option is the Short Sequence Number Header Format Option which allows the sequence numbers in MLPPPoX encapsulated frames/fragments to be 12-bit long (instead 24-bit long, by default).

After the multilink capability is successfully negotiated via LCP, PPP sessions can be bundled together over MLPPPoX capable links.

The basic operational principles are:

- LCP session is negotiated on each physical link with MLPPPoX capabilities between the two nodes.
- Based on the ED or the authentication outcome, a bundle is created. A subsequent IPCP negotiation is conveyed over this bundle. User traffic is sent over the bundle.
- If a new link tries to join the bundle by sending a new MLPPPoX LCP Conf Request, the LCP session is negotiated, authentication performed and the link is placed under the bundle containing the links with the same ED or authentication outcome.
- IPCP/IPv6CP is, in the whole process, negotiated only once over the bundle. This negotiation occurs at the beginning, when the first link is established and MLPPPoX bundle created. IPCP and IPv6CP messages are transmitted from the 7750 SR LNS without MLPPPoX encapsulation, while they can be received as MLPPPoX encapsulated or non-MLPPPoX encapsulated.

10.3.5 Enabling MLPPPoX

The lowest granularity at which MLPPPoX can be enabled is an L2TP tunnel. An MLPPPoX enabled tunnel is not limited to carrying only MLPPPoX sessions but can carry normal PPP(oE) sessions as well.

In addition to enabling MLPPPoX on the session terminating LNS node, MLPPPoX can also be enabled on the LAC via PPP policy. The purpose of enabling MLPPPoX on the LAC is to negotiate MLPPPoX LCP parameters with the client. When the LAC receives the MRRU option from the client in the initial LCP ConfReq, it changes its tunnel selection algorithm so that all sessions of an MLPPPoX bundle are mapped into the same tunnel.

The LAC negotiates MLPPPoX LCP parameters regardless of the transport technology connected to it (ATM or Ethernet). LCP negotiated parameters are passed by the LAC to the LNS via Proxy LCP in ICCN message. In this fashion the LNS has an option to accept the LCP parameters negotiated by the LAC or to reject them and restart the negotiation directly with the client.

The LAC transparently passes session traffic handed to it by the LNS in the downstream direction and the MLPPPoX client in the upstream direction. The LNS and the MLPPPoX client performs all data processing functions related to MLPPPoX such as fragmentation and interleaving.

When the LCP negotiation is completed and the LCP transition into an open state (configuration ACKs are sent and received), the Authentication phase on the LAC begins. During the Authentication phase, the L2TP parameters become known (l2tp group, tunnel, and so on), and the session is extended by the LAC to the LNS via L2TP. In case that the Authentication phase does not return L2TP parameters, the session is terminated because the 7750 SR does not support directly terminated MLPPPoX sessions.

In the case that MLPPPoX is not enabled on the LAC, the LAC negotiates plain PPP session with the client. In case that the client accepts plain PPP instead of MLPPPoX as offered by the LAC, when the session is extended to the LNS, the LNS re-negotiates MLPPPoX LCP with the client on a MLPPPoX enabled tunnel. The LNS learns about the MLPPPoX capability of the client via Proxy LCP message in ICCN (first Conf Req received from the client is also send in Proxy LCP). If there is no indication of the MLPPPoX capability of the client, the LNS establishes a plain PPP(oE) session with the client.

10.3.6 Link Fragmentation and Interleaving (LFI)

The purpose of LFI is to ensure that short high priority packets are not delayed by the transmission delay of large low priority packets on slow links.

For example it takes ~150ms to transmit a 5000B packet over a 256 kb/s link, while the same packet is transmitted in only 40us over a 1G link (~4000 times faster transmission). To avoid the delay of a high priority packet by waiting in the queue while the large packet is being transmitted, the large packet can be segmented into smaller chunks. The high priority packet can be then interleaved with the smaller fragments. This approach can significantly reduce the delay of high priority packets.

The interleaving functionality is only supported on MLPPPoX bundles with a single link. If more than one link is added into a interleaving capable MLPPPoX bundle, then interleaving is internally disabled and the `tmnxMlpppBundleIndicatorsChange` trap is generated.

With interleaving enabled on an MLPPPoX enabled tunnel, the following session types are supported:

- **multiple LCP sessions tied into a single MLPPPoX bundle**

This scenario assumes multiple physical links on the client side. Theoretically it would be possible to have multiple sessions running over the same physical link in the last mile. For example, two PPPoE sessions going over the same Ethernet link in the last mile. Whichever the case may be, the LAC/LNS is unaware of the physical topology in the last mile (single or multiple physical links). Interleaving functionality is internally disabled on such MLPPPoX bundles.

- **a single LCP session (including dual stack) over the MLPPPoX bundle**

This scenario assumes a single physical link on the client side. Interleaving is supported on such single session MLPPPoX bundle as long as the conditions for interleaving are met. Those conditions are governed by `max-fragment-delay` parameter and calculation of the fragment size as described in subsequent sections.

- **an LCP session (including dual stack) over a plain PPP/PPPoE session**

This type of session is a regular PPP(oE) session outside of any MLPPPoX bundle and therefore its traffic is not MLPPPoX encapsulated.

Packets on an MLPPPoX bundle are MLPPPoX encapsulated unless they are classified as high priority packets when interleaving is enabled.

10.3.6.1 MLPPPoX fragmentation, MRRU and MRU considerations

A packet of the size greater than the internally calculated fragment length cannot be natively transmitted over an MLPPPoX bundle. Such packets are MLPPPoX encapsulated and consequently fragmented. This is irrespective of whether the fragmentation is enabled or disabled. The size of the internally calculated fragment length depends on:

- the needed transmission delay in the last mile
- the fragment "payload to encapsulation overhead" efficiency ratio
- various MTU sizes in the 7750 SR dictated mainly by received MRU, received MRRU and configured PPP MTU under the following hierarchy:
 - `configure service vprn l2tp group ppp mtu`
 - `configure service vprn l2tp group tunnel ppp mtu`
 - `configure router l2tp group ppp mtu`

- configure router l2tp group tunnel ppp mtu

In cases where MLPPPoX fragmentation is disabled with the **no max-fragment-delay** command, it is expected that packets are not MLPPPoX fragmented but rather only MLPPPoX encapsulated to be load balanced over multiple physical links in the last mile. However, even if MLPPPoX fragmentation is disabled, it is possible that fragmentation occurs under specific circumstances. This behavior is related to the calculation of the MTU values on an MLPPPoX bundle.

MLPPPoX in the 7750 SR is concerned with two MTUs:

- **bundle-mtu** determines the maximum length of the original IP packet that can be transmitted over the entire bundle (collection of links) before any MLPPPoX processing takes place on the transmitting side. This is also the maximum size of the IP packet that the receiving node can accept after it de-encapsulates and assembles received MLPPPoX fragments of the same packet. Bundle-mtu is relevant in the context of the collection of links.
- **link-mtu** determines the maximum length of the payload before it is PPP encapsulated and transmitted over an individual link within the bundle. Link-mtu is relevant in the context of the single link within the bundle.

Assuming that the CPE advertised MRRU and MRU values are smaller than any configurable mtu on MLPPPoX processing modules in the 7750 SR (carrier IOM and BB-ISA), the bundle-mtu and the link-mtu are based on the received MRRU and MRU values, respectively. For example, the bundle-mtu is set to the received MRRU value while link-bundle is set to the MRU value minus the MLPPPoX encapsulation overhead (4 or 6 bytes).

Consider an example where received MRRU value sent by CPE is 1500B while received MRU is 1492B. In this case, our bundle-mtu is set to 1500B and our link-mtu is set to 1488B (or 1486B) to allow for the additional 4/6B of MLPPPoX encapsulation overhead. Consequently, IP payload of 1500B can be transmitted over the bundle but only 1488B can be transmitted over any individual link. In case that an IP packet with the size between 1489B and 1500B needs to be transmitted from the 7750 SR toward the CPE, this packet would be MLPPPoX fragmented in the 7750 SR as dictated by the link-mtu. This is irrespective of whether MLPPPoX fragmentation is enabled or disabled (as set by no max-fragment-delay flag).

To entirely avoid MLPPPoX fragmentation in this case, the received MRRU sent by CPE should be lower than the received MRU for the length of the MLPPPoX header (4 or 6 bytes). In this case, for IP packets larger than 1488B, IP fragmentation would occur (assuming that DF flag in the IP header allows it) and MLPPPoX fragmentation would be avoided.

On the 7750 SR side, it is not possible to set different advertised MRRU and MRU values with the **ppp-mtu** command. Both MRRU and MRU advertised values adhere to the same configured ppp mtu value.

10.3.7 LFI functionality implemented in LNS

As mentioned in the previous section, LFI on LNS is implemented only on MLPPPoX bundles with a single LCP session.

There are two major tasks associated with LFI on the LNS:

- executing subscriber QoS in the carrier IOM based on the last mile conditions. The subscriber QoS rates are the last mile on-the-wire rates. After traffic is QoS conditioned, it is sent to the BB-ISA for further processing.
- fragmentation and artificial delay (queuing) of the fragments so that high priority packets can be injected in-between low priority fragments (interleaved). This operation is performed by the BB-ISA.

Most of this is also applicable to non-lfi case. The only difference between lfi and non-lfi is that there is no artificial delay performed in non-lfi case.

Examine an example to further clarify functionality of LFI. The parameters, conditions and requirements that are used in our example to describe the wanted behavior are the following:

- High priority packets must not be delayed for more than 50ms in the last mile because of the transmission delay of the large low priority packets. Considering that tolerated end-to-end VoIP delay must be under 150ms, limiting the transmission delay to 50ms on the last mile link is a reasonable choosing.
- The link between the LNS and LAC is 1Gb/s Ethernet.
- The last mile link rate is 256 kb/s.
- Three packets arrive back-to-back on the network side of the LNS (in the downstream direction). The large 5000B low priority packet P1 arrives first, followed by two smaller high priority packets P2 and P3, each 100B in length. Packets P1, P2 and P3 can be originated by independent sources (PCs, servers, and so on.) and therefore can theoretically arrive in the LNS from the network side back-to-back at the full network link rate (10Gb/s or 100Gb/s).
- The transmission time on the internal 10G link between the BB-ISA and the carrier IOM for the large packet (5000B) is 4us while the transmission time for the small packet (100B) is 80ns.
- The transmission time on the 1G link (LNS->LAC) for the large packet (5000B) is 40us while the transmission time for the small packet (100B) is 0.8us.
- The transmission time in the last mile (256 kb/s) for the large packet is ~150ms while the transmission time for the small packet on the same link is ~3ms.
- Last mile transport is ATM.

To satisfy the delay requirement for the high priority packets, the large packets are fragmented into three smaller fragments. The fragments are carefully sized so that their individual transmission time in the last mile does not exceed 50ms. After the first 50ms interval, there is window of opportunity to interleave the two smaller high priority packets.

This entire process is further clarified by the five points (1-5) in the packet route from the LNS to the Residential Gateway (RG).

The five points are described in subsequent sections.

10.3.7.1 Last mile QoS awareness in the LNS

By implementing MLPPPoX in LNS, we are effectively transferring the traffic treatment functions (QoS/LFI) of the last mile to the node (LNS) that is multiple hops away.

The success of this operation depends on the accuracy at which we can simulate the last mile conditions in the LNS. The assumption is that the LNS is aware of the two most important parameters of the last mile:

- **the last mile encapsulation**

This is needed for the accurate calculation of the overhead associated of the transport medium in the last mile for traffic shaping and interleaving.

- **the last mile link rate**

This is crucial for the creation of artificial congestion and packet delay in the LNS.

The subscriber QoS in the LNS is implemented in the carrier IOM and is performed on a per packets basis before the packet is handed over to the BB-ISA. Per packet, instead of per fragment QoS processing

ensures a more efficient utilization of network resources in the downstream direction. Discarding fragments in the LNS would have detrimental effects in the RG as the RG would be unable to reconstruct a packet without all of its fragments.

High priority traffic within the bundle is classified into the high priority queue. This type of traffic is not MLPPPoX encapsulated unless its packet size exceeds the link MTU as described in [MLPPPoX fragmentation, MRRU and MRU considerations](#). Low priority traffic is classified into a low priority queue and is always MLPPPoX encapsulated. In case that the high priority traffic becomes MLPPPoX encapsulated/fragmented, the MLPPPoX processing module (BB-ISA) considers it as low-priority. The assumption is that the high priority traffic is small in size and consequently MLPPPoX encapsulation/fragmentation and degradation in priority can be avoided. The aggregate rate of the MLPPPoX bundle is on-the-wire rate of the last mile as shown in Figure 3.

ATM on the wire overhead for non MLPPPoX encapsulated high priority traffic includes:

- ATM encapsulation (VC MUX, LLC/NLPID, LCC/SNAP)
- AAL5 trailer (8B)
- AAL5 padding to 48B cell boundary (this makes the overhead dependent on the packet size)
- multiplication by 53/48 to account for the ATM cell headers

For low priority traffic which is always MLPPPoX encapsulated, an additional overhead related to MLPPPoX encapsulation and possibly fragmentation must be added (blue arrow in Figure 3). In other words, each fragment carries ATM+MLPPPoX overhead.

The 48B boundary padding can be avoided for all fragments except the last one. This can be done by choosing the fragment length so that it is aligned on the 48B boundary (rounded down if based on max-fragment-delay or rounded up if based on the encapsulation/utilization).

For Ethernet in the last mile, our implementation always assures that the fragment size plus the encapsulation overhead is always larger or equal to the minimum Ethernet packet length (64B).

10.3.7.2 BB-ISA processing

MLPPPoX encapsulation, fragmentation and interleaving are performed by the LNS in BB-ISA. If we refer to our example, a large low priority packet (P1) is received by the BB-ISA, immediately followed by the two small high priority packets (P2 and P3). Because our requirement stipulates that there is no more than 50ms of transmission delay in the last mile (including on-the-wire overhead), the large packet must be fragmented into three smaller fragments each of which does not cause more than 50ms of transmission delay.

The BB-ISA would normally send packets/fragments to the carrier IOM at the rate of 10Gb/s. In other words, by default the three fragments of the low priority packet would be sent out of the BB-ISA back-to-back at the very high rate before the high priority packets even arrive in the BB-ISA. To interleave, the BB-ISA must simulate the last mile conditions by delaying the transmission of the fragments. The fragments are paced out of the BB-ISA (and out of the box) at the rate of the last mile. High priority packets get the opportunity to be injected in front of the fragments while the fragments are being delayed.

As shown in [Figure 112: QoS hierarchy on LNS](#) (point 2) the first fragment F1 is sent out immediately (transmission delay at 10G is in the 1us range). The transmission of the next fragment F2 is delayed by 50ms. While the transmission of the second fragment F2 is being delayed, the two high priority packets (P1 and P2 in red) are received by the BB-ISA and are immediately transmitted ahead of fragments F2 and F3. This approach relies on the imperfection of the IOM shaper which is releasing traffic in bursts (P2 and P3 right after P1). The burst size is dependent on the depth of the rate token bucket associated with the IOM shaper.

By the time the second fragment F2 is transmitted, the first fragment F1 has traveled a long way (50ms) on high rate links toward the Access Node (assuming that there is no queuing delay along the way), and its transmission on the last mile link has already begun (if not already completed).

This is not applicable for this discussion, but nonetheless worth noticing is that the LNS BB-ISA also adds the L2TP encapsulation to each packet/fragment. The L2TP encapsulation is removed in the LAC before the packet/fragment is transmitted toward the AN.

10.3.7.3 LNS-LAC link

This is the high rate link (1Gb/s) on which the first fragment F1 and the two consecutive high priority packets, P2 and P3, are sent back-to-back by the BB-ISA

(BB-ISA->carrier IOM->egress IOM-> out-of-the-LNS).

The remaining fragments (F2 and F3) are still waiting in the BB-ISA to be transmitted. They are artificially delayed by 50ms each.

Additional QoS based on the L2TP header can be performed on the egress port in the LNS toward the LAC. This QoS is based on the classification fields inside of the packet/fragment headers (DSCP, dot1.p, EXP).

The LAC-AN link is not really relevant for the operation of LFI on the LNS. This link can be either Ethernet (in case of PPPoE) or ATM (PPPoE or PPP). The rate of the link between the LAC and the AN is still considered a high speed link compared to the slow last mile link.

10.3.7.4 AN-RG link

Finally, this is the slow link of the last mile, the reason why LFI is performed in the first place. Assuming that LFI played its role in the network as designed, by the time the transmission of one fragment on this link is completed, the next fragment arrives just in time for unblocked transmission. In between the two fragments, we can have one or more small high priority packets waiting in the queue for the transmission to complete.

On the AN-RG link in [Figure 112: QoS hierarchy on LNS](#) that packets P2 and P3 are ahead of fragments F2 and F3. Therefore the delay incurred on this link by the low priority packets is never greater than the transmission delay of the first fragment (50ms). The remaining two fragments, F2 and F3, can be queued and further delayed by the transmission time of packets P2 and P3 (which is normally small, in our example 3ms for each).

If many low priority packets are waiting in the queue, then they would have caused delay and would have further delayed the fragments that are in transit from the LNS to the LAC. This condition is normally caused by bursts and it should clear itself out over time.

10.3.7.5 Home link

High priority packets P2 and P3 are transmitted by the RG into the home network ahead of the packet P1 although the fragment F1 has arrived in the RG first. The reason for this is that the RG must wait for the fragments F2 and F3 before it can re-assemble packet P1.

10.3.7.6 Optimum fragment size calculation by LNS

Fragmentation in LFI is based on the optimal fragment size. LNS implementation calculates the two optimal fragment sizes, based on two different criteria:

- optimal fragment size based on the payload efficiency of the fragment given the fragmentation/transportation header overhead associated with the fragment encapsulation based fragment size
- optimal fragment size based on the maximum transmission delay of the fragment set by configuration delay-based fragment size

At the end, only one optimal fragment size is selected. The actual fragments length are of the optimal fragment size.

The parameters required to calculate the optimal fragment sizes are known to the LNS either via configuration or via signaling. These, in-advance known parameters are:

- last mile maximum transmission delay (max-fragment-delay obtained via CLI)
- last mile ATM Encapsulation (in our example the last mile is ATM but in general it can be Ethernet for MLPPPoE)
- MLPPP encapsulation length (depending on the fragment sequence number format)
- the last mile on-the-wire rate for the MLPPPoX bundle

Examine closer each of the two optimal fragment sizes.

10.3.7.6.1 Encapsulation based fragment size

One needs to be mindful of the fact that fragmentation may cause low link utilization. In other words, during fragmentation a node may end up transporting mainly overhead bytes in the fragment as opposed to payload bytes. This would only intensify the problem that fragmentation is intended to solve, especially on an ATM access link that tend to carry larger encapsulation overhead.

To reduce the overhead associated with fragmentation, the following is enforced in the 7750 SR:

The minimum fragment payload size is at least 10 times greater than the overhead (MLPPP header, ATM Encapsulation and AAL5 trailer) associated with the fragment.

The optimal fragment length (including the MLPPP header, the ATM Encapsulation and the AAL5 trailer) is a multiple of 48B. Otherwise, the AAL5 layer would add an additional 48B boundary padding to each fragment, which would unnecessarily expand the overhead associated with fragmentation. By aligning all-but-last fragments to a 48B boundary, only the last fragment potentially contains the AAL5 48B boundary padding which is no different from a non-fragmented packet. For future reference, we will refer to all fragments except for the last fragment as non-padded fragments. The last fragment will obviously be padded if it is not already natively aligned to a 48B boundary.

As an example, calculate the optimal fragment size based on the encapsulation criteria with the maximum fragment overhead of 22B. To achieve >10x transmission efficiency the fragment payload size must be 220B (10*22B). To avoid the AAL5 padding, the entire fragment (overhead + payload) is rounded UP on a 48B boundary. The final fragment size is 288B [22B + 22B*10 + 48B_alignment].

In conclusion, an optimal fragment size was selected that carries the payload with at least 90% efficiency. The last fragment of the packet cannot be artificially aligned on a 48B boundary (it is a natural reminder), so it is padded by the AAL5 layer. Therefore the efficiency of the last fragment is probably be less than 90% in our example. In the extreme case, the efficiency of this last fragment may be only 2%.

The fragment size chosen in this manner is purely chosen based on the overhead length. The maximum transmission delay did not play any role in the calculations.

For Ethernet based last mile, the CPM always makes sure that the fragment size plus encapsulation overhead is larger or equal to the minimum Ethernet packet length of 64B.

10.3.7.6.2 Fragment size based on the max transmission delay

The first criterion in selecting the optimal fragment size based on the maximum transmission delay mandates that the transmission time for the fragment, including all overheads (MLPPP header, ATM encapsulation header, AAL5 overhead and ATM cell overhead) must be less than the configured max-fragment-delay time.

The second criterion mandates that each fragment, including the MLPPP header, the ATM Encapsulation header, the AAL5 trailer and the ATM cellification overhead be a multiple of 48B. The fragment size is rounded down to the nearest 48B boundary during the calculations to minimize the transmission delay. Aligning the fragment on the 48B boundary eliminates the AAL5 padding and therefore reduces the overhead associated with the fragment. The overhead reduction does not only improve the transmission time, but also increases the efficiency of the fragment.

Given these two criteria along with the configuration parameters (ATM Encapsulation, MLPPP header length, max-fragment-delay time, rate in the last mile), the implementation calculates the optimal non-padded fragment length as well as the transmission time for this optimal fragment length.

10.3.7.6.3 Selection of the optimum fragment length

So far the implementation has calculated the two optimum fragment lengths, one based on the length of the MLPPP/transport encapsulation overhead of the fragment, the other one based on the maximum transmission delay of the fragment. Both of them are aligned on a 48B boundary. The larger of the two is chosen and the BB-ISA performs LFI based on this selected optimal fragment length.

10.3.8 Upstream traffic considerations

Fragmentation and interleaving is implemented on the originating end of the traffic. In other words, in the upstream direction the CPE (or RG) is fragmenting and interleaving traffic. There is no interleaving or fragmentation processing in the upstream direction in the 7750 SR. The 7750 SR is on the receiving end and is only concerned with the reassembly of the fragments arriving from the CPE. Fragments are buffered until the packet can be reconstructed. If all fragments of a packet are not received within a preconfigured timeframe, the received fragments of the partial packet are discarded (a packet cannot be reconstructed without all of its fragments). This time-out and discard is necessary to prevent buffer starvation in the BB-ISA. Two values for the time-out can be configured: 100ms and 1s.

10.3.9 Multiple links MLPPPoX with no interleaving

Interleaving over MLPPPoX bundles with multiple links are not supported. However, fragmentation is supported.

To preserve packet order, all packets on an MLPPPoX bundle with multiple links are MLPPPoX encapsulated (monotonically increased sequence numbers).

We do not support multiclass MLPPP (RFC 2686, *The Multi-Class Extension to Multi-Link PPP*). Multiclass MLPPP would require another level of intelligent queuing in the BB-ISA which we do not have.

10.3.10 MLPPPoX session support

MLPPPoE is the only session type in the last mile that is supported:

MLPPPoE can be a single physical link or multilink. The last mile encapsulation is Ethernet over copper (This could be Ethernet over VDSL or HSDSL). The access rates (especially upstream) are still limited by the xDSL distance limitation and therefore, interleaving is required on a slow speed single link in the last mile. It is possible that the last mile encapsulation is Ethernet over fiber (FTTH) but in this case, users would not be concerned with the link speed to the point where interleaving and link aggregation is required.

Finally, this is the slow link of the last mile, the reason why LFI is performed in the first place. Assuming that LFI played its role in the network as designed, by the time the transmission of one fragment on this link is completed, the next fragment arrives just in time for unblocked transmission. In between the two fragments, we can have one or more small high priority packets waiting in the queue for the transmission to complete.

We can see on the AN-RG link in Figure 2 that packets P2 and P3 are ahead of fragments F2 and F3. Therefore the delay incurred on this link by the low priority packets is never greater than the transmission delay of the first fragment (50ms). The remaining two fragments, F2 and F3, can be queued and further delayed by the transmission time of packets P2 and P3 (which is normally small, in our example 3ms for each).

If many low priority packets were waiting in the queue, then they would have caused delay for each other and would have further delayed the fragments in transit from the LNS to the LAC. This condition is normally caused by bursts and it should clear itself out over time.

MLPPP(oEo)A can be a single physical link or multilink. The last mile encapsulation is ATM over xDSL.

Some other combinations are also possible (ATM in the last mile, Ethernet in the aggregation) but they all come down to one of the above models that are characterized by:

- Ethernet or ATM in the last mile.
- Ethernet or ATM access on the LAC.
- MLPPP/PPPoE termination on the LNS

10.3.11 Session load balancing across multiple BB-ISAs

PPP/PPPoE sessions are by default load balanced across multiple BB-ISAs (max 6) in the same group. The load balancing algorithm considers the number of active session on each BB-ISA in the same group.

The load balancing algorithm does not take into account the number of queues consumed on the carrier IOM. Therefore a session can be refused if queues are depleted on the carrier IOM even though the BB-ISA may be lightly loaded in terms of the number of sessions that is hosting.

With MLPPPoX, it is important that multiple sessions per bundle be terminated on the same LNS BB-ISA. This can be achieved by per tunnel load balancing mode where all sessions of a tunnel are terminated in the same BB-ISA. Per tunnel load balancing mode is mandatory on LNS BB-ISAs that are in the group that supports MLPPPoX.

On the LAC side, all sessions in an MLPPPoX bundle are automatically assigned to the same tunnel. In other words an MLPPPoX bundle is assigned to the tunnel. There can be multiple tunnels created between the same pair of LAC/LNS nodes.

10.3.12 BB-ISA hashing considerations

All downstream traffic on an MLPPPoX bundle with multiple links is always MLPPPoX encapsulated. Some traffic is fragmented and served in a octet oriented round robin fashion over multiple member links. However, fragments are never delayed in case that the bundle contains multiple links.

In a per fragment/packet load sharing algorithm, there is always the possibility that there is uneven load utilization between the member links. A single link overload most likely goes unnoticed in the network all the way to the Access Node. The access node is the only node in the network that actually has multiple physical links connected to it. All other session-aware nodes (LAC and LNS) only see MLPPPoX as a bundle with multiple sessions without any mechanism to shape traffic per physical link. Other nodes in this case being 7750 SRs. Other vendors may have the ability to condition (shape) traffic per session.

If one of the member sessions is perpetually overloaded by the LNS, traffic is dropped in the last mile because the corresponding physical link cannot absorb traffic beyond its physical capabilities. This would have detrimental effects on the whole operation of the MLPPPoX bundle. To prevent this perpetual overloading of the member links that can be caused by per packet/fragment load balancing scheme, the load balancing scheme that takes into account the number of octets transmitted over each member link. The octet counter of a new link is initialized to the lowest value of any existing link counter. Otherwise the load balancing mechanism would show significant bias toward the new link until the byte counter catches up with the rest of the links.

10.3.13 Last mile rate and encapsulation parameters

The last mile rate information along with the encapsulation information is used for fragmentation (to determine the maximum fragment length) and interleaving (delaying fragments in the BB-ISA). In addition, the aggregate subscriber rate (aggregate-rate-limit) on the LNS is automatically adjusted based on the last mile link rate and the number of links in the MLPPPoX bundle.

- **downstream data rate in the last mile**

The subscriber aggregate rates (agg-rate-limit) used in (H)QoS on the carrier IOM and in the BB-ISA (for interleaving) must be wire based in the last mile. This rule applies equally to both, the LAC and LNS.

The last mile on-the-wire rates of the subscriber can be submitted to the LAC and the LNS via various means. The following bullets describe how the last mile wire rates are passed to each entity:

- **LAC**

The last mile link rate is taken via the following methods in the order of listed priority:

1. LUDB (**rate-down** command under the host hierarchy in LUDB)
2. RADIUS Alc-Access-Loop-Rate-Down VSA. Although this VSA is stored in the state of plain PPP(oE) sessions (MLPPPoX bundled or not), it is applicable only to MLPPPoX bundles.
3. PPPoE tags; Vendor Specific Tags (RFC 2516, *A Method for Transmitting PPP Over Ethernet (PPPoE)*); tag type 0x0105; tag value is Enterprise Number 3561 followed by the TLV sub-options as specified in TR-101 -> Actual Data Rate Downstream 0x82)

As long as the link rate information is available in the LAC, it is always passed to the LNS in the ICRQ message using the standard L2TP encoding. This cannot be disabled.

In addition, an option is available to control the source of the rate information can be conveyed to the LNS via TX Connect Speed AVP in the ICCN message. This can be used for compatibility reasons with other vendors that can only use TX Connect Speed to pass the link rate information to the LNS. By default, the maximum port speed (or the sum of the maximum speeds of all member ports in the LAG) is reported in TX Connect Speed. Unlike the rate conveyed in ICRQ message, The TX Connect Speed content is configurable via the following command:

```
config>subscr-mgmt
  sla-profile <name>
    egress
      report-rate agg-rate-limit | scheduler <scheduler-name> | pppoe-actual-rate
    | rfc5515-actual-rate
```

The report-rate configuration option dictates which rate is reported in the TX Connect Speed as follows:

- agg-rate-limit => statically configured agg-rate-limit value or RADIUS QoS override is reported.
- scheduler <scheduler-name> => virtual schedulers are not supported in MLPPPoX.
- pppoe-actual-rate => rate taken from PPPoE Tags are reported. Rate reported via RFC5515 can still be different if the source for both methods is not the same.
- rfc5515-actual-speed => the rate is taken from RFC5515.

The RFC 5515 relies on the same encoding as PPPoE tags (vendor ID is ADSL Forum and the type for Actual Data Rate Downstream is 0x82). The two methods of passing the line rate to the LNS are using different message types (ICRQ and ICCN).

The LAC on the 7750 SR is not aware of MLPPPoX bundles. As such, the aggregate subscriber bandwidth on the LAC is configured statically via usual means (sub-profile, scheduler-policy) or dynamically modified via RADIUS. The aggregate subscriber (or MLPPPoX bundle) bandwidth on the LAC is not automatically adjusted according to the rates of the individual links in the bundle and the number of the links in the bundle. As such, an operator must ensure that the statically provided rate value for aggregate-rate-limit is the sum of the bandwidth of each member link in the MLPPPoX bundle. The number of member links and their bandwidth must be therefore known in advance. The alternative is to have the aggregate rate of the MLPPPoX bundle set to a high value and rely on the QoS treatment performed on the LNS.

• LNS

The sources of information for the last mile link rate on the LNS are taken in the following order:

1. LUDB (during user authentication phase, same as in LAC)
2. RADIUS (same as in LAC)
3. ICRQ message, Actual Data Downstream Rate (RFC 5515)
4. ICCN message, TX Connect Speed

There is no configuration option to determine the priority of the source of information for the last mile link rate. TX Connect Speed in ICCN message is only taken into consideration as a last resort in absence of any other source of last mile rate information.

After the last mile rate information is obtained, the subscriber aggregate rate (aggregate-rate-limit) is automatically adjusted to the minimum value of:

- the smallest link speed in the MLPPPoX bundle multiplied by the number of links in the bundle
- statically configured aggregate-rate-limit

The link speed of each link in the bundle must be the same, that is, different link speeds within the bundle are not supported. In the case that we receive different link speed values for last mile links within the bundle, we adopt the minimum received speed and apply it to all links.

In case that the obtained rate information from the last mile for a session within the MLPPP bundle is out of bounds (1 kb/s to 100 Mb/s), the session within the bundle is terminated.

- **encapsulation**

Wire-rates are dependent on the encapsulation of the link to which they apply. The last mile encapsulation information can be extracted via various means.

- **LAC**

- static configuration via LUDB
- RADIUS (Alc-Access_Loop-Encap-Offset VSA)
- PPPoE tags; Vendor Specific Tags (RFC 2516; tag type 0x0105; tag value is Enterprise Number 3561 followed by the TLV sub-options as specified in TR-101 -> Actual Data Rate Downstream 0x82).

The LAC passes the line encapsulation information to the LNS via ICRQ message using the encoding defined in the RFC 5515.

- **LNS**

The LNS extracts the encapsulation information in the following order:

- static configuration via LUDB
- RADIUS (Alc-Access-Loop-Encap-Offset VSA)
- ICRQ message (RFC 5515)

In case that the encapsulation information is not provided by any of the existing means (LUDB, RADIUS, AVP signaling, PPPoE Tags), then by default pppoa-null encapsulation is in effect. This applies to LAC and LNS.

10.3.14 Link failure detection

The link failure in the last mile is detected via the expiration of session keepalives (LCP). The LNS tears down the session over the failed link and notify the LAC via a CDN message.

10.3.15 CoA support

CoA request for the subscriber aggregate-rate-limit change is honored on the LAC and the LNS.

CoA for the rate change of an individual link within the bundle is supported through the same VSA that can be used to initially assign the rate parameter to each member link. This is supported only on LNS. The rate override via CoA is applied to all active link members within the bundle.

Change of the access link parameters via CoA is supported in the following fashion:

- Change of access loop encap: refused (NAK)

- Change of access loop rate down:
- On L2TP LAC session: refused (NAK). On LAC the access loop rate down is not locally used for any rate limiting function but instead it is just passed to the LNS at the beginning when the session is first established. Mid-session changes on LAC via CoA are not propagated to the LNS.
- On L2TP LNS session:
 - Plain session: ignored. The rate is stored in the MIB table but no rate limiting action is taken. In other words, this parameter is internally excluded from rate calculations and advertisements. However, it is shown in the output of the relevant show commands.
- Bundle session: applied on all link sessions. The aggregate rate limit of the bundle is set to the minimum of the:
 - CoA obtained local loop down rate multiplied by the number of links in the bundle
 - The aggregate rate limit configured statically or obtained via CoA.
- Fragment length is affected by this change. In case that interleaving is enabled on a single link bundle, the interleave interval is affected.
- Non-L2TP: ignored. The rate is stored in the MIB table but no rate limiting action is taken. In other words, this parameter is internally excluded from rate calculations and advertisements. However, it is shown in the output of the relevant show commands.

Similar behavior is exhibited if at mid-session, the parameters are changed via LUDB with the exception of the rate-down parameter in LAC. If this parameter is changed on the LAC, all sessions are disconnected.

10.3.16 Accounting

Accounting counters on the LNS include all packet overhead (wire overhead from the last mile). There is only one accounting session per bundle.

On the LAC, there is one accounting session per pppoe session (link).

In tunnel-accounting mode there is one accounting session per link.

On LNS only the stop-link of the last link of the bundle carries all accounting data for the bundle.

10.3.17 Filters and mirroring

Filters and mirrors (LI) are not supported on an MLPPPoX bundle on LAC. However, filters and ip-only mirror type are supported on the LNS.

10.3.18 PTA considerations

Locally terminated MLPPPoX (PTA) solution is offered based on the LAC and the LNS hosted in the same system. An external loop (or VSM2) is used to connect the LAC to the LNS within the same box. The subscribers are terminated on the LNS.

10.3.19 QoS considerations

10.3.19.1 Dual-pass

HQoS and LFI are performed in two stages that involve double traversal (dual-pass) of traffic through the carrier IOM and the BB-ISA. The following are the functions performed in each pass:

- In the first pass through the carrier IOM, traffic is marked (dot1p bits) as high or low priority. This plays a crucial role in the execution of LFI in the BB-ISA.
- In the first pass through the BB-ISA this prioritization from the 1st step, is an indication (along with the internally calculated fragment size) of whether the traffic is interleaved (non MLPPP encapsulated) or not (MLPPP encapsulated). Consequently the BB-ISA adds the necessary padding related to last mile wire overhead to each packet. This padding is used in the second pass on the carrier IOM to perform last mile wire based QoS functions.
- In the second pass through the carrier IOM, the last mile wire based HQoS is performed based on the padding added in the first pass through the BB-ISA.
- In the second pass through the BB-ISA, previously added overhead is stripped off and LFI/MLPPP encapsulation functions are performed.

10.3.19.2 Traffic prioritization in LFI

The delivery of high priority traffic within predefined delay bounds on a slow speed last mile link is ensured by correct QoS classification and prioritization. High priority traffic is interleaved with low priority fragments on a single link MLPPPoX bundle with LFI enabled. The classification of traffic into the correct (high or low priority) forwarding class is performed on the downstream ingress interface. However, traffic can be re-classified (re-mapped into another forwarding class) on the egress access interface of the carrier IOM, just before packets are transmitted to the BB-ISA for MLPPPoX processing. This can be achieved via QoS sap-egress policy referenced in the LNS sla-profile.

The priority of the forwarding class in regular QoS (on IOM) is determined by the properties (Expedited, non-expedited queue type, CIR and PIR rates) of the queue to which the forwarding class is mapped. In contracts, traffic prioritization in LFI domain (in BB-ISA) is determined by the outer dot1p bits that are set by the carrier IOM while transmitting packets toward the BB-ISA. The outer dot1p bits are marked based on the forwarding class information determined by classification/re-classification on ingress/carrier IOM. This marking of outer dot1p bits in the Ethernet header between the carrier IOM and the BB-ISA is fixed and defined in the default sap-egress LNS ESM policy 65537. The marking definition is as follows:

```
FC be -> dot1p 0
FC l2 -> dot1p 1
FC af -> dot1p 2
FC l1 -> dot1p 3
FC h2 -> dot1p 4
FC ef -> dot1p 5
FC h1 -> dot1p 6
FC nc -> dot1p 7
```

In LFI (on BB-ISA), dot1p bits [0,1,2 and 3] are considered low priority while dot1p bits (4,5,6 and 7) are considered high priority. Consequently, forwarding classes BE, L2, AF and L1 are considered low priority while forwarding classes H2, EF, H1 and NC are considered high priority. High priority traffic (assuming that the packet size does not exceed maximum fragment size) is interleaved with low priority traffic.

The following describes the reference points in traffic prioritization for the purpose of LFI in the 7750 SR:

- **classification on downstream ingress interface (entrance point into the 7750 SR)**

Packets can be classified into one of the following eight forwarding classes: be, l2, af, l1, h2, ef, h1 and nc. Depending on the type of the ingress interface (access or network), traffic can be classified based on dot1p, exp, DSCP, ToS bits or ip-match criteria (dscp, dst-ip, dst-port, fragment, src-ip, src-port and protocol-id).

- **re-classification on downstream access egress interface between the carrier IOM and the BB-ISA**

In the carrier IOM, downstream traffic can be re-classified into another forwarding class, just before it is forwarded to the BB-ISA. Re-classification on access egress is based on the same fields as on ingress except for the dot1p and exp bits because Ethernet or MPLS headers from ingress are not carried from ingress to egress.

- **marking on downstream access egress interface between the carrier IOM and the BB-ISA**

When the forwarding class is available on the carrier IOM in the egress direction (toward BB-ISA), it is used to mark outer dot1p bits in the new Ethernet header that is used to transport the frame from the carrier IOM to the BB-ISA. The marking of the dot1p bits on the egress SAP between the carrier IOM and the BB-ISA cannot be changed for MLPPPoX even if the **no qos-marking-from-sap** command is configured under the sla-profile on egress.

10.3.19.3 Shaping based on the last mile wire rates

Accurate QoS, amongst other things, require that the subscriber rates in the first mile on an MLPPPoX bundle be properly represented in the LNS. In other words, the rate limiting functions in the LNS must account for the last mile on-the-wire encapsulation overhead. The last mile encapsulation can be Ethernet or ATM.

For ATM in the last mile, the LNS accounts for the following per fragment overhead:

- PID
- MLPPP encapsulation header
- ATM Fixed overhead (ATM encap + fixed AAL5 trailer)
- 48B boundary padding as part of AAL5 trailer
- 5B per each 48B of data in ATM cell

In case of Ethernet encapsulation in the last mile, the overhead is:

- PID
- MLPPP header per fragment
- Ethernet Header + FCS per fragment
- preamble + IPG overhead per fragment

The **encap-offset** command under the sub-profile egress CLI node is ignored in case of MLPPPoX. MLPPPoX rate calculation is, by default, always based on the last mile wire overhead.

The HQoS rates (port-scheduler, aggregate-rate-limit and scheduler) on LNS are based on the wire overhead of the entity to which the HQoS is applied. For example, if the port-scheduler is managing bandwidth on the link between the BB-ISA and the carrier IOM, then the rate of such scheduler accounts for the q-in-q Ethernet encapsulation on that link along with the preamble and inter packet gap (20B).

While virtual schedulers (attached via sub-profile) are supported on LNS for plain PPPoX sessions, they are not supported for MLPPPoX bundles. Only aggregate- rate-limit along with the port-scheduler can be used in MLPPPoX deployments.

10.3.19.4 Downstream bandwidth management on egress port

Bandwidth management on the egress physical ports (Physical Port 1 and Physical Port 2 in Figure 8) is performed at the egress port itself on the egress IOM instead on the carrier IOM. By default, the forwarding class (FC) information is preserved from network ingress to network egress. However, this can be changed via QoS configuration applied to the egress SAP of the carrier IOM toward the BB-ISA.

L2TP traffic originated locally in LNS can be marked via the router/service vprn->sgt-qos hierarchy.

10.3.20 Sub/sla-profile considerations

- **sub-profile**

In the MLPPPoX case on LNS, multiple sessions are tied into the same subscriber aggregate-rate-limit via a sub-profile. The consequence is that the aggregate rate of the subscriber can be adjusted dynamically depending on the advertised link speed in the last mile and the number of links in the bundle. Shaping in the LNS is performed per the entire MLPPPoX bundle (subscriber) instead of per individual member links within the bundle. The exception is obviously a MLPPPoX bundle with the single member link (interleaving case) where the relationship between the session and the MLPPPoX bundle is 1:1.

In the LAC, the subscriber aggregate rate cannot be dynamically changed based on the number of links in the bundle and their rate. The LAC has no notion of MLPPPoX bundles. However, multiple sessions that in reality belong to an MLPPPoX bundle under the subscriber are shaped as an aggregate (agg-rate-limit under the sub-profile). This in essence yields the same shaping behavior as on LNS.

- **sla-profile**

Sessions within the MLPPPoX bundle in LNS share a single sla-profile instances (queues).

In the LAC, as long as the sessions within the subscriber6 are on the same SAP, they can also share the same sla-profile. This is the case in MLPPPoX.

The manner in which sub/sla-profile are applied to MLPPPoX bundles and the individual sessions within results in aggregate shaping per MLPPPoX bundle as well as allocation of unique set of queues per MLPPPoX bundle. This is valid irrespective of the location where shaping is executed (LAC or LNS). Other vendors may have implemented shaping per session within the bundle and this is something that needs to be taken into consideration during the migration process.

10.3.21 Example of MLPPPoX session setup flow

- **LAC behavior**

1. A new PPP(oEoA) session request arrives on the LAC (PADI or LCP Conf Req).
2. The LAC negotiates PADx session if applicable.
3. The LAC may negotiate MLPPPoX LCP phase with its own endpoint discriminator, or it may reject MLPPPoX specific options in LCP if MLPPPoX on the LAC is disabled (that is, no accept-mrru in the LAC's ppp-policy). If MLPPPoX options (seq num header format, ED, MRRU) are rejected, the assumption is that the client renegotiates plain PPP(oEoA) session with the LAC.
4. When LCP (MLPPPoX capable or not) is negotiated, the session is authenticated (PAP/CHAP).
5. Upon successful authentication, an L2TP tunnel is identified to which the session belongs.

6. If the session is a non-L2TP session (PTA MLPPPoX capable session for which the tunnel cannot be determined), the session is terminated.
7. Otherwise, the QoS constructs are created for the subscriber hosts: the session is assigned to a sub/sla-profiles.
8. The session LCP parameters are sent to the LNS via call management messages.
9. If another LCP session is requested on the same bundle, the LAC creates a new LCP session and join this session to the existing subscriber as another host. In other words, the LAC is bundle agnostic and the two sessions appear as two hosts under the same subscriber.

- **LNS behavior**

The following assumes that MLPPPoX is configured on the LNS under the L2TP group or the tunnel hierarchy.

- The LNS has the option to accept the LCP parameters or to reject them and start renegotiating LCP parameters directly with the client.
- If the LNS choose to renegotiate LCP parameters with the client directly, this renegotiation is completely transparent to the LAC by the means of a T-bit (control vs. data) in the L2TP header. LCP is renegotiated on the LNS with all the options necessary to support MLPPPoX. Endpoint Discriminator is not mandatory in the MLPPPoX negotiation. If the client rejects it, the LNS must still be able to negotiate MLPPPoX capable session (same is valid for the LAC). If the client's endpoint discriminator is invalid (bad format, invalid class, and so on), the 7750 SR does not negotiate MLPPPoX and instead a plain PPP session is created.
- If the LNS is configured to accept the LCP Proxy parameters, the LNS determines the capability of the client.

If there is no indication of MLPPPoX capability in the Proxy LCP (not even in the original ConfReq), the LNS may accept plain (non MLPPPoX capable) LCP session or renegotiate from scratch the non MLPPPoX capable session.

If there is an indication of MLPPPoX capability in the Proxy LCP (either completely negotiated on the LAC or at least attempted from the client), the LNS tries to either accept the MLPPPoX negotiated session by the LAC or renegotiate the MLPPPoX capable session directly with the client.

If the LCP Proxy parameters with MLPPPoX capability are accepted by the LNS, then the endpoint as negotiated on the LAC is also accepted.

- After the MLPPPoX capable LCP session is negotiated or accepted, authentication can be performed on the LNS. Authentication on the LNS can be restarted (CHAP challenge/response with the client), or accepted (chap challenge/response accepted and verified by the LNS via RADIUS).
- If the authentication is successful, depending on the evaluation of the parameters negotiated up to this point, a new MLPPPoX bundle is created or an existing MLPPPoX bundle is joined. In case that a new bundle is established, the QoS constructs for the subscriber(-host) are created (sub/sla-profile). Session negotiation advances to IPCP phase.
- The decision whether a new session should join an existing MLPPPoX bundle, or trigger creation of a new one is governed by RFC 1990, *The PPP Multilink Protocol (MP)*, section 5.1.3, page 16, cases 1,2,3, and 4.
- Interleaving is supported only on MLPPPoX bundles with single session in them.

10.3.22 Other considerations

- IPv6 is supported.
- AA is supported at LNS where full IP packets can be redirected via AA policies.
- Intra-chassis redundancy is supported:
 - CPM (stateful failover)
 - BB-ISA (non-stateful failover)

10.4 LNS support on ESA

The recommendations for the LNS support on ESA are:

- The entire ESA should be dedicated to the LNS application.
- Dedicating all cores and memory for a single VM can allow higher throughput per L2TP session while dividing the cores and memory among VMs can allow a higher number of L2TP sessions. Nokia recommends dividing the ESA into a maximum of 2 VMs with an equal number of cores and memory to allow a higher number of L2TP sessions. Contact your local Nokia representatives for more information.

The limitations of the LNS support on ESA are:

- This feature is supported on FP3-based line cards and later.
- ISA and ESA cannot be used in the same LNS group.

10.5 Configuration notes

MLPPP in subscriber management context is supported only over ATM, Ethernet over ATM or plain Ethernet transport (MLPPPoX). Native MLPPP over PPP/HDLC links is supported outside of the subscriber management context on the ASAP MDA.

MLPPPoX is supported only on LNS.

Interleaving is supported only on MLPPPoX bundles with a single member link. If more than one link is present in an MLPPPoX bundle, the interleaving is automatically disabled and a SNMP trap is generated. The MIB for this even is defined as `tmnxMlpppBundleIndicatorsChange`.

If MLPPPoX is enabled on LNS, the load balancing mode between the BB-ISAs within the group should be set to per tunnel. This ensures that all sessions of the same MLPPPoX bundle are terminated on the same BB-ISA. On the LAC, sessions of the same bundle are setup in the same tunnel.

Virtual schedulers are not supported on MLPPPoX tunnels on LNS. However, aggregate-rate-limit is supported.

The aggregate-rate-limit on LNS is automatically adjusted to the minimum value of:

- configured aggregate-rate-limit
- minimum last mile rate (obtained via LUDB, RADIUS or PPPoE tags) multiplied by the number of links in the bundle

The aggregate-rate-limit on the LAC is not adjusted automatically. Therefore, if configured it should be set to a high value and therefore, the traffic treatment should rely on QoS performed on the LNS.

The rate (rate-down information) of the member links within the bundle must be the same. Otherwise the lowest rate is selected and applied to all member links.

A single CoA for a rate change (Alc-Access-Loop-Rate-Down) of an individual link in an MLPPPoX bundle modifies rates of all links in the bundle. This is applicable on LNS only.

The range of supported last mile rate (rate-down information) for the member links on an MLPPPoX session is 1 kb/s — 100 Mb/s. On the LNS the last mile rate can be obtained:

- from the LAC via Tx-Connect-Speed AVP or by standard L2TP encoding as described in the RFC 5515, *Layer 2 Tunneling Protocol (L2TP) Access Line Information Attribute Value Pair (AVP) Extensions*.
- from the LAC via LUDB or RADIUS
- directly on the LNS via LUDB or RADIUS.

The session fails to come up if the obtained rate-down information is outside of the allowable range (1 kb/s — 100 Mb/s).

A session within the MLPPPoX bundle is terminated if the rate-down information for the session is out of bounds (1kb/s — 100 Mb/s).

If a member link in the last mile fails, traffic is blackholed until the LNS is notified of this failure. The failure detection in the LNS relies on PPP keepalives.

Shaping is performed per MLPPPoX bundle and not individually per member links.

If encapsulation overhead associated with fragmentation is too large in comparison to payload, the fragments are sized based on the encapsulation overhead (to increase link efficiency) instead of on maximum transmission delay.

There can be only a single MLPPPoX bundle per subscriber.

MLPPPoX bundles and non-MLPPPoX (plain L2TP PPPoE) sessions cannot coexist under the same subscriber.

Filters and mirrors (LI) are not supported on MLPPPoX bundles on LAC.

ip-only type mirrors are supported on MLPPPoX bundles.

In MLPPP scenario, downstream traffic is traversing Carrier IOM and BB-ISA twice. This is referred to as dual-pass and effectively cuts the throughput for MLPPP in half (for example, 5Gb/s of MLPPP traffic on a 10Gb/s capable BB-ISA).

11 Standards and protocol support

**Note:**

The information provided in this chapter is subject to change without notice and may not apply to all platforms.

Nokia assumes no responsibility for inaccuracies.

11.1 Access Node Control Protocol (ANCP)

draft-ietf-ancp-protocol-02, *Protocol for Access Node Control Mechanism in Broadband Networks*

RFC 5851, *Framework and Requirements for an Access Node Control Mechanism in Broadband Multi-Service Networks*

11.2 Bidirectional Forwarding Detection (BFD)

draft-ietf-lsr-ospf-bfd-strict-mode-10, *OSPF BFD Strict-Mode*

RFC 5880, *Bidirectional Forwarding Detection (BFD)*

RFC 5881, *Bidirectional Forwarding Detection (BFD) IPv4 and IPv6 (Single Hop)*

RFC 5882, *Generic Application of Bidirectional Forwarding Detection (BFD)*

RFC 5883, *Bidirectional Forwarding Detection (BFD) for Multihop Paths*

RFC 7130, *Bidirectional Forwarding Detection (BFD) on Link Aggregation Group (LAG) Interfaces*

RFC 7880, *Seamless Bidirectional Forwarding Detection (S-BFD)*

RFC 7881, *Seamless Bidirectional Forwarding Detection (S-BFD) for IPv4, IPv6, and MPLS*

RFC 7883, *Advertising Seamless Bidirectional Forwarding Detection (S-BFD) Discriminators in IS-IS*

RFC 7884, *OSPF Extensions to Advertise Seamless Bidirectional Forwarding Detection (S-BFD) Target Discriminators*

RFC 9247, *BGP - Link State (BGP-LS) Extensions for Seamless Bidirectional Forwarding Detection (S-BFD)*

11.3 Border Gateway Protocol (BGP)

draft-gredler-idr-bgplu-epe-14, *Egress Peer Engineering using BGP-LU*

draft-hares-idr-update-attrib-low-bits-fix-01, *Update Attribute Flag Low Bits Clarification*

draft-ietf-idr-add-paths-guidelines-08, *Best Practices for Advertisement of Multiple Paths in IBGP*

draft-ietf-idr-best-external-03, *Advertisement of the best external route in BGP*

draft-ietf-idr-bgp-flowspec-oid-03, *Revised Validation Procedure for BGP Flow Specifications*
draft-ietf-idr-bgp-gr-notification-01, *Notification Message support for BGP Graceful Restart*
draft-ietf-idr-bgp-ls-app-specific-attr-16, *Application-Specific Attributes Advertisement with BGP Link-State*
draft-ietf-idr-bgp-ls-flex-algo-06, *Flexible Algorithm Definition Advertisement with BGP Link-State*
draft-ietf-idr-bgp-optimal-route-reflection-10, *BGP Optimal Route Reflection (BGP-ORR)*
draft-ietf-idr-error-handling-03, *Revised Error Handling for BGP UPDATE Messages*
draft-ietf-idr-flowspec-interfaceset-03, *Applying BGP flowspec rules on a specific interface set*
draft-ietf-idr-flowspec-path-redirect-05, *Flowspec Indirection-id Redirect – localised ID*
draft-ietf-idr-flowspec-redirect-ip-02, *BGP Flow-Spec Redirect to IP Action*
draft-ietf-idr-link-bandwidth-03, *BGP Link Bandwidth Extended Community*
draft-ietf-idr-long-lived-gr-00, *Support for Long-lived BGP Graceful Restart*
RFC 1772, *Application of the Border Gateway Protocol in the Internet*
RFC 1997, *BGP Communities Attribute*
RFC 2385, *Protection of BGP Sessions via the TCP MD5 Signature Option*
RFC 2439, *BGP Route Flap Damping*
RFC 2545, *Use of BGP-4 Multiprotocol Extensions for IPv6 Inter-Domain Routing*
RFC 2858, *Multiprotocol Extensions for BGP-4*
RFC 2918, *Route Refresh Capability for BGP-4*
RFC 4271, *A Border Gateway Protocol 4 (BGP-4)*
RFC 4360, *BGP Extended Communities Attribute*
RFC 4364, *BGP/MPLS IP Virtual Private Networks (VPNs)*
RFC 4456, *BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)*
RFC 4486, *Subcodes for BGP Cease Notification Message*
RFC 4659, *BGP-MPLS IP Virtual Private Network (VPN) Extension for IPv6 VPN*
RFC 4684, *Constrained Route Distribution for Border Gateway Protocol/MultiProtocol Label Switching (BGP/MPLS) Internet Protocol (IP) Virtual Private Networks (VPNs)*
RFC 4724, *Graceful Restart Mechanism for BGP – helper mode*
RFC 4760, *Multiprotocol Extensions for BGP-4*
RFC 4798, *Connecting IPv6 Islands over IPv4 MPLS Using IPv6 Provider Edge Routers (6PE)*
RFC 5004, *Avoid BGP Best Path Transitions from One External to Another*
RFC 5065, *Autonomous System Confederations for BGP*
RFC 5291, *Outbound Route Filtering Capability for BGP-4*
RFC 5396, *Textual Representation of Autonomous System (AS) Numbers – asplain*
RFC 5492, *Capabilities Advertisement with BGP-4*
RFC 5668, *4-Octet AS Specific BGP Extended Community*
RFC 6286, *Autonomous-System-Wide Unique BGP Identifier for BGP-4*

RFC 6368, *Internal BGP as the Provider/Customer Edge Protocol for BGP/MPLS IP Virtual Private Networks (VPNs)*

RFC 6793, *BGP Support for Four-Octet Autonomous System (AS) Number Space*

RFC 6810, *The Resource Public Key Infrastructure (RPKI) to Router Protocol*

RFC 6811, *Prefix Origin Validation*

RFC 6996, *Autonomous System (AS) Reservation for Private Use*

RFC 7311, *The Accumulated IGP Metric Attribute for BGP*

RFC 7606, *Revised Error Handling for BGP UPDATE Messages*

RFC 7607, *Codification of AS 0 Processing*

RFC 7674, *Clarification of the Flowspec Redirect Extended Community*

RFC 7752, *North-Bound Distribution of Link-State and Traffic Engineering (TE) Information Using BGP*

RFC 7854, *BGP Monitoring Protocol (BMP)*

RFC 7911, *Advertisement of Multiple Paths in BGP*

RFC 7999, *BLACKHOLE Community*

RFC 8092, *BGP Large Communities Attribute*

RFC 8097, *BGP Prefix Origin Validation State Extended Community*

RFC 8212, *Default External BGP (EBGP) Route Propagation Behavior without Policies*

RFC 8277, *Using BGP to Bind MPLS Labels to Address Prefixes*

RFC 8571, *BGP - Link State (BGP-LS) Advertisement of IGP Traffic Engineering Performance Metric Extensions*

RFC 8950, *Advertising IPv4 Network Layer Reachability Information (NLRI) with an IPv6 Next Hop*

RFC 8955, *Dissemination of Flow Specification Rules*

RFC 8956, *Dissemination of Flow Specification Rules for IPv6*

RFC 9086, *Border Gateway Protocol - Link State (BGP-LS) Extensions for Segment Routing BGP Egress Peer Engineering*

11.4 Bridging and management

IEEE 802.1AB, *Station and Media Access Control Connectivity Discovery*

IEEE 802.1ad, *Provider Bridges*

IEEE 802.1ag, *Connectivity Fault Management*

IEEE 802.1ah, *Provider Backbone Bridges*

IEEE 802.1ak, *Multiple Registration Protocol*

IEEE 802.1aq, *Shortest Path Bridging*

IEEE 802.1AX, *Link Aggregation*

IEEE 802.1D, *MAC Bridges*

IEEE 802.1p, *Traffic Class Expediting*

IEEE 802.1Q, *Virtual LANs*
IEEE 802.1s, *Multiple Spanning Trees*
IEEE 802.1w, *Rapid Reconfiguration of Spanning Tree*
IEEE 802.1X, *Port Based Network Access Control*

11.5 Broadband Network Gateway (BNG) Control and User Plane Separation (CUPS)

3GPP TS 23.003, *Numbering, addressing and identification*
3GPP TS 23.007, *Restoration procedures*
3GPP TS 23.501, *System architecture for the 5G System (5GS)*
3GPP TS 23.502, *Procedures for the 5G System (5GS)*
3GPP TS 23.503, *Policy and charging control framework for the 5G System (5GS)*
3GPP TS 24.501, *Non-Access-Stratum (NAS) protocol for 5G System (5GS)*
3GPP TS 29.244, *Interface between the Control Plane and the User Plane nodes*
3GPP TS 29.281, *General Packet Radio System (GPRS) Tunnelling Protocol User Plane (GTPv1-U)*
3GPP TS 29.500, *Technical Realization of Service Based Architecture*
3GPP TS 29.501, *Principles and Guidelines for Services Definition*
3GPP TS 29.502, *Session Management Services*
3GPP TS 29.503, *Unified Data Management Services*
3GPP TS 29.512, *Session Management Policy Control Service*
3GPP TS 29.518, *Access and Mobility Management Services*
3GPP TS 32.255, *5G data connectivity domain charging*
3GPP TS 32.290, *Services, operations and procedures of charging using Service Based Interface (SBI)*
3GPP TS 32.291, *5G system, charging service*
BBF TR-459, *Control and User Plane Separation for a Disaggregated BNG*
BBF TR-459.2, *Multi-Service Disaggregated BNG with CUPS: Integrated Carrier Grade NAT function*
RFC 8300, *Network Service Header (NSH)*

11.6 Certificate management

RFC 4210, *Internet X.509 Public Key Infrastructure Certificate Management Protocol (CMP)*
RFC 4211, *Internet X.509 Public Key Infrastructure Certificate Request Message Format (CRMF)*
RFC 5280, *Internet X.509 Public Key Infrastructure Certificate and Certificate Revocation List (CRL) Profile*
RFC 6712, *Internet X.509 Public Key Infrastructure -- HTTP Transfer for the Certificate Management Protocol (CMP)*

RFC 7030, *Enrollment over Secure Transport*

RFC 7468, *Textual Encodings of PKIX, PKCS, and CMS Structures*

11.7 Circuit emulation

RFC 4553, *Structure-Agnostic Time Division Multiplexing (TDM) over Packet (SAToP)*

RFC 5086, *Structure-Aware Time Division Multiplexed (TDM) Circuit Emulation Service over Packet Switched Network (CESoPSN)*

RFC 5287, *Control Protocol Extensions for the Setup of Time-Division Multiplexing (TDM) Pseudowires in MPLS Networks*

11.8 Ethernet

IEEE 802.3ah, *Media Access Control Parameters, Physical Layers, and Management Parameters for Subscriber Access Networks*

IEEE 802.3x, *Ethernet Flow Control*

ITU-T G.8031/Y.1342, *Ethernet Linear Protection Switching*

ITU-T G.8032/Y.1344, *Ethernet Ring Protection Switching*

ITU-T Y.1731, *OAM functions and mechanisms for Ethernet based networks*

11.9 Ethernet VPN (EVPN)

draft-ietf-bess-evpn-ipvpn-interworking-06, *EVPN Interworking with IPVPN*

draft-ietf-bess-evpn-irb-mcast-04, *EVPN Optimized Inter-Subnet Multicast (OISM) Forwarding – ingress replication*

draft-ietf-bess-evpn-pref-df-06, *Preference-based EVPN DF Election*

draft-ietf-bess-evpn-unequal-lb-16, *Weighted Multi-Path Procedures for EVPN Multi-Homing – section 9*

draft-ietf-bess-evpn-virtual-eth-segment-06, *EVPN Virtual Ethernet Segment*

draft-ietf-bess-pbb-evpn-isid-cmacflush-00, *PBB-EVPN ISID-based CMAC-Flush*

draft-sajassi-bess-evpn-ip-aliasing-05, *EVPN Support for L3 Fast Convergence and Aliasing/Backup Path – IP Prefix routes*

draft-sr-bess-evpn-vpws-gateway-03, *Ethernet VPN Virtual Private Wire Services Gateway Solution*

draft-trr-bess-bgp-srv6-args-02, *SRv6 Argument Signaling for BGP Services*

RFC 7432, *BGP MPLS-Based Ethernet VPN*

RFC 7623, *Provider Backbone Bridging Combined with Ethernet VPN (PBB-EVPN)*

RFC 8214, *Virtual Private Wire Service Support in Ethernet VPN*

RFC 8317, *Ethernet-Tree (E-Tree) Support in Ethernet VPN (EVPN) an Provider Backbone Bridging EVPN (PBB-EVPN)*

RFC 8365, *A Network Virtualization Overlay Solution Using Ethernet VPN (EVPN)*
RFC 8560, *Seamless Integration of Ethernet VPN (EVPN) with Virtual Private LAN Service (VPLS) and Their Provider Backbone Bridge (PBB) Equivalents*
RFC 8584, *DF Election and AC-influenced DF Election*
RFC 9047, *Propagation of ARP/ND Flags in an Ethernet Virtual Private Network (EVPN)*
RFC 9135, *Integrated Routing and Bridging in Ethernet VPN (EVPN) – Asymmetric IRB Procedures and Mobility Procedure*
RFC 9136, *IP Prefix Advertisement in Ethernet VPN (EVPN)*
RFC 9161, *Operational Aspects of Proxy ARP/ND in Ethernet Virtual Private Networks*
RFC 9251, *Internet Group Management Protocol (IGMP) and Multicast Listener Discovery (MLD) Proxies for Ethernet VPN (EVPN)*

11.10 gRPC Remote Procedure Calls (gRPC)

cert.proto version 0.1.0, *gRPC Network Operations Interface (gNOI) Certificate Management Service*
file.proto version 0.1.0, *gRPC Network Operations Interface (gNOI) File Service*
gnmi.proto version 0.8.0, *gRPC Network Management Interface (gNMI) Service Specification*
PROTOCOL-HTTP2, *gRPC over HTTP2*
system.proto Version 1.0.0, *gRPC Network Operations Interface (gNOI) System Service*

11.11 Intermediate System to Intermediate System (IS-IS)

draft-ietf-isis-mi-02, *IS-IS Multi-Instance*
draft-kaplan-isis-ext-eth-02, *Extended Ethernet Frame Size Support*
ISO/IEC 10589:2002 Second Edition, *Intermediate system to Intermediate system intra-domain routing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode Network Service (ISO 8473)*
RFC 1195, *Use of OSI IS-IS for Routing in TCP/IP and Dual Environments*
RFC 2973, *IS-IS Mesh Groups*
RFC 3359, *Reserved Type, Length and Value (TLV) Codepoints in Intermediate System to Intermediate System*
RFC 3719, *Recommendations for Interoperable Networks using Intermediate System to Intermediate System (IS-IS)*
RFC 3787, *Recommendations for Interoperable IP Networks using Intermediate System to Intermediate System (IS-IS)*
RFC 5120, *M-ISIS: Multi Topology (MT) Routing in IS-IS*
RFC 5130, *A Policy Control Mechanism in IS-IS Using Administrative Tags*
RFC 5301, *Dynamic Hostname Exchange Mechanism for IS-IS*
RFC 5302, *Domain-wide Prefix Distribution with Two-Level IS-IS*

RFC 5303, *Three-Way Handshake for IS-IS Point-to-Point Adjacencies*
RFC 5304, *IS-IS Cryptographic Authentication*
RFC 5305, *IS-IS Extensions for Traffic Engineering TE*
RFC 5306, *Restart Signaling for IS-IS – helper mode*
RFC 5308, *Routing IPv6 with IS-IS*
RFC 5309, *Point-to-Point Operation over LAN in Link State Routing Protocols*
RFC 5310, *IS-IS Generic Cryptographic Authentication*
RFC 6119, *IPv6 Traffic Engineering in IS-IS*
RFC 6213, *IS-IS BFD-Enabled TLV*
RFC 6232, *Purge Originator Identification TLV for IS-IS*
RFC 6233, *IS-IS Registry Extension for Purges*
RFC 6329, *IS-IS Extensions Supporting IEEE 802.1aq Shortest Path Bridging*
RFC 7775, *IS-IS Route Preference for Extended IP and IPv6 Reachability*
RFC 7794, *IS-IS Prefix Attributes for Extended IPv4 and IPv6 Reachability – sections 2.1 and 2.3*
RFC 7981, *IS-IS Extensions for Advertising Router Information*
RFC 7987, *IS-IS Minimum Remaining Lifetime*
RFC 8202, *IS-IS Multi-Instance – single topology*
RFC 8570, *IS-IS Traffic Engineering (TE) Metric Extensions – Min/Max Unidirectional Link Delay metric for flex-algo, RSVP, SR-TE*
RFC 8919, *IS-IS Application-Specific Link Attributes*

11.12 Internet Protocol (IP) Fast Reroute (FRR)

draft-ietf-rtgwg-lfa-manageability-08, *Operational management of Loop Free Alternates*
RFC 5286, *Basic Specification for IP Fast Reroute: Loop-Free Alternates*
RFC 7431, *Multicast-Only Fast Reroute*
RFC 7490, *Remote Loop-Free Alternate (LFA) Fast Reroute (FRR)*
RFC 8518, *Selection of Loop-Free Alternates for Multi-Homed Prefixes*

11.13 Internet Protocol (IP) general

draft-grant-tacacs-02, *The TACACS+ Protocol*
RFC 768, *User Datagram Protocol*
RFC 793, *Transmission Control Protocol*
RFC 854, *Telnet Protocol Specifications*
RFC 1350, *The TFTP Protocol (revision 2)*

RFC 2347, *TFTP Option Extension*
RFC 2348, *TFTP Blocksize Option*
RFC 2349, *TFTP Timeout Interval and Transfer Size Options*
RFC 2428, *FTP Extensions for IPv6 and NATs*
RFC 2617, *HTTP Authentication: Basic and Digest Access Authentication*
RFC 2784, *Generic Routing Encapsulation (GRE)*
RFC 2818, *HTTP Over TLS*
RFC 2890, *Key and Sequence Number Extensions to GRE*
RFC 3164, *The BSD syslog Protocol*
RFC 4250, *The Secure Shell (SSH) Protocol Assigned Numbers*
RFC 4251, *The Secure Shell (SSH) Protocol Architecture*
RFC 4252, *The Secure Shell (SSH) Authentication Protocol – publickey, password*
RFC 4253, *The Secure Shell (SSH) Transport Layer Protocol*
RFC 4254, *The Secure Shell (SSH) Connection Protocol*
RFC 4511, *Lightweight Directory Access Protocol (LDAP): The Protocol*
RFC 4513, *Lightweight Directory Access Protocol (LDAP): Authentication Methods and Security Mechanisms – TLS*
RFC 4632, *Classless Inter-domain Routing (CIDR): The Internet Address Assignment and Aggregation Plan*
RFC 5082, *The Generalized TTL Security Mechanism (GTSM)*
RFC 5246, *The Transport Layer Security (TLS) Protocol Version 1.2 – TLS client, RSA public key*
RFC 5425, *Transport Layer Security (TLS) Transport Mapping for Syslog – RFC 3164 with TLS*
RFC 5656, *Elliptic Curve Algorithm Integration in the Secure Shell Transport Layer – ECDSA*
RFC 5925, *The TCP Authentication Option*
RFC 5926, *Cryptographic Algorithms for the TCP Authentication Option (TCP-AO)*
RFC 6398, *IP Router Alert Considerations and Usage – MLD*
RFC 6528, *Defending against Sequence Number Attacks*
RFC 7011, *Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of Flow Information*
RFC 7012, *Information Model for IP Flow Information Export*
RFC 7230, *Hypertext Transfer Protocol (HTTP/1.1): Message Syntax and Routing*
RFC 7231, *Hypertext Transfer Protocol (HTTP/1.1): Semantics and Content*
RFC 7232, *Hypertext Transfer Protocol (HTTP/1.1): Conditional Requests*
RFC 7301, *Transport Layer Security (TLS) Application Layer Protocol Negotiation Extension*
RFC 7616, *HTTP Digest Access Authentication*
RFC 8446, *The Transport Layer Security (TLS) Protocol Version 1.3*

11.14 Internet Protocol (IP) multicast

cisco-ipmulticast/pim-autorp-spec01, *Auto-RP: Automatic discovery of Group-to-RP mappings for IP multicast* – version 1

draft-ietf-bier-pim-signaling-08, *PIM Signaling Through BIER Core*

draft-ietf-idmr-traceroute-ipm-07, *A "traceroute" facility for IP Multicast*

draft-ietf-l2vpn-vpls-pim-snooping-07, *Protocol Independent Multicast (PIM) over Virtual Private LAN Service (VPLS)*

RFC 1112, *Host Extensions for IP Multicasting*

RFC 2236, *Internet Group Management Protocol, Version 2*

RFC 2365, *Administratively Scoped IP Multicast*

RFC 2375, *IPv6 Multicast Address Assignments*

RFC 2710, *Multicast Listener Discovery (MLD) for IPv6*

RFC 3306, *Unicast-Prefix-based IPv6 Multicast Addresses*

RFC 3376, *Internet Group Management Protocol, Version 3*

RFC 3446, *Anycast Rendezvous Point (RP) mechanism using Protocol Independent Multicast (PIM) and Multicast Source Discovery Protocol (MSDP)*

RFC 3590, *Source Address Selection for the Multicast Listener Discovery (MLD) Protocol*

RFC 3618, *Multicast Source Discovery Protocol (MSDP)*

RFC 3810, *Multicast Listener Discovery Version 2 (MLDv2) for IPv6*

RFC 3956, *Embedding the Rendezvous Point (RP) Address in an IPv6 Multicast Address*

RFC 3973, *Protocol Independent Multicast - Dense Mode (PIM-DM): Protocol Specification (Revised) – auto-RP groups*

RFC 4541, *Considerations for Internet Group Management Protocol (IGMP) and Multicast Listener Discovery (MLD) Snooping Switches*

RFC 4604, *Using Internet Group Management Protocol Version 3 (IGMPv3) and Multicast Listener Discovery Protocol Version 2 (MLDv2) for Source-Specific Multicast*

RFC 4607, *Source-Specific Multicast for IP*

RFC 4608, *Source-Specific Protocol Independent Multicast in 232/8*

RFC 4610, *Anycast-RP Using Protocol Independent Multicast (PIM)*

RFC 4611, *Multicast Source Discovery Protocol (MSDP) Deployment Scenarios*

RFC 5059, *Bootstrap Router (BSR) Mechanism for Protocol Independent Multicast (PIM)*

RFC 5186, *Internet Group Management Protocol Version 3 (IGMPv3) / Multicast Listener Discovery Version 2 (MLDv2) and Multicast Routing Protocol Interaction*

RFC 5384, *The Protocol Independent Multicast (PIM) Join Attribute Format*

RFC 5496, *The Reverse Path Forwarding (RPF) Vector TLV*

RFC 6037, *Cisco Systems' Solution for Multicast in MPLS/BGP IP VPNs*

RFC 6512, *Using Multipoint LDP When the Backbone Has No Route to the Root*

RFC 6513, *Multicast in MPLS/BGP IP VPNs*
RFC 6514, *BGP Encodings and Procedures for Multicast in MPLS/IP VPNs*
RFC 6515, *IPv4 and IPv6 Infrastructure Addresses in BGP Updates for Multicast VPNs*
RFC 6516, *IPv6 Multicast VPN (MVPN) Support Using PIM Control Plane and Selective Provider Multicast Service Interface (S-PMSI) Join Messages*
RFC 6625, *Wildcards in Multicast VPN Auto-Discover Routes*
RFC 6826, *Multipoint LDP In-Band Signaling for Point-to-Multipoint and Multipoint-to-Multipoint Label Switched Path*
RFC 7246, *Multipoint Label Distribution Protocol In-Band Signaling in a Virtual Routing and Forwarding (VRF) Table Context*
RFC 7385, *IANA Registry for P-Multicast Service Interface (PMSI) Tunnel Type Code Points*
RFC 7716, *Global Table Multicast with BGP Multicast VPN (BGP-MVPN) Procedures*
RFC 7761, *Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol Specification (Revised)*
RFC 8279, *Multicast Using Bit Index Explicit Replication (BIER)*
RFC 8296, *Encapsulation for Bit Index Explicit Replication (BIER) in MPLS and Non-MPLS Networks – MPLS encapsulation*
RFC 8401, *Bit Index Explicit Replication (BIER) Support via IS-IS*
RFC 8444, *OSPFv2 Extensions for Bit Index Explicit Replication (BIER)*
RFC 8487, *Mtrace Version 2: Traceroute Facility for IP Multicast*
RFC 8534, *Explicit Tracking with Wildcard Routes in Multicast VPN – (C-*,C-*) wildcard*
RFC 8556, *Multicast VPN Using Bit Index Explicit Replication (BIER)*

11.15 Internet Protocol (IP) version 4

RFC 791, *Internet Protocol*
RFC 792, *Internet Control Message Protocol*
RFC 826, *An Ethernet Address Resolution Protocol*
RFC 951, *Bootstrap Protocol (BOOTP) – relay*
RFC 1034, *Domain Names - Concepts and Facilities*
RFC 1035, *Domain Names - Implementation and Specification*
RFC 1191, *Path MTU Discovery – router specification*
RFC 1519, *Classless Inter-Domain Routing (CIDR): an Address Assignment and Aggregation Strategy*
RFC 1534, *Interoperation between DHCP and BOOTP*
RFC 1542, *Clarifications and Extensions for the Bootstrap Protocol*
RFC 1812, *Requirements for IPv4 Routers*
RFC 1918, *Address Allocation for Private Internets*
RFC 2003, *IP Encapsulation within IP*

RFC 2131, *Dynamic Host Configuration Protocol*
RFC 2132, *DHCP Options and BOOTP Vendor Extensions*
RFC 2401, *Security Architecture for Internet Protocol*
RFC 3021, *Using 31-Bit Prefixes on IPv4 Point-to-Point Links*
RFC 3046, *DHCP Relay Agent Information Option (Option 82)*
RFC 3768, *Virtual Router Redundancy Protocol (VRRP)*
RFC 4884, *Extended ICMP to Support Multi-Part Messages – ICMPv4 and ICMPv6 Time Exceeded*

11.16 Internet Protocol (IP) version 6

RFC 2464, *Transmission of IPv6 Packets over Ethernet Networks*
RFC 2529, *Transmission of IPv6 over IPv4 Domains without Explicit Tunnels*
RFC 3122, *Extensions to IPv6 Neighbor Discovery for Inverse Discovery Specification*
RFC 3315, *Dynamic Host Configuration Protocol for IPv6 (DHCPv6)*
RFC 3587, *IPv6 Global Unicast Address Format*
RFC 3596, *DNS Extensions to Support IP version 6*
RFC 3633, *IPv6 Prefix Options for Dynamic Host Configuration Protocol (DHCP) version 6*
RFC 3646, *DNS Configuration options for Dynamic Host Configuration Protocol for IPv6 (DHCPv6)*
RFC 3736, *Stateless Dynamic Host Configuration Protocol (DHCP) Service for IPv6*
RFC 3971, *SEcure Neighbor Discovery (SEND)*
RFC 3972, *Cryptographically Generated Addresses (CGA)*
RFC 4007, *IPv6 Scoped Address Architecture*
RFC 4191, *Default Router Preferences and More-Specific Routes – Default Router Preference*
RFC 4193, *Unique Local IPv6 Unicast Addresses*
RFC 4291, *Internet Protocol Version 6 (IPv6) Addressing Architecture*
RFC 4443, *Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6) Specification*
RFC 4861, *Neighbor Discovery for IP version 6 (IPv6)*
RFC 4862, *IPv6 Stateless Address Autoconfiguration – router functions*
RFC 4890, *Recommendations for Filtering ICMPv6 Messages in Firewalls*
RFC 4941, *Privacy Extensions for Stateless Address Autoconfiguration in IPv6*
RFC 5007, *DHCPv6 Leasequery*
RFC 5095, *Deprecation of Type 0 Routing Headers in IPv6*
RFC 5722, *Handling of Overlapping IPv6 Fragments*
RFC 5798, *Virtual Router Redundancy Protocol (VRRP) Version 3 for IPv4 and IPv6 – IPv6*
RFC 5952, *A Recommendation for IPv6 Address Text Representation*

RFC 6092, *Recommended Simple Security Capabilities in Customer Premises Equipment (CPE) for Providing Residential IPv6 Internet Service – Internet Control and Management, Upper-Layer Transport Protocols, UDP Filters, IPsec and Internet Key Exchange (IKE), TCP Filters*

RFC 6106, *IPv6 Router Advertisement Options for DNS Configuration*

RFC 6164, *Using 127-Bit IPv6 Prefixes on Inter-Router Links*

RFC 6221, *Lightweight DHCPv6 Relay Agent*

RFC 6437, *IPv6 Flow Label Specification*

RFC 6603, *Prefix Exclude Option for DHCPv6-based Prefix Delegation*

RFC 8021, *Generation of IPv6 Atomic Fragments Considered Harmful*

RFC 8200, *Internet Protocol, Version 6 (IPv6) Specification*

RFC 8201, *Path MTU Discovery for IP version 6*

11.17 Internet Protocol Security (IPsec)

draft-ietf-ipsec-isakmp-mode-cfg-05, *The ISAKMP Configuration Method*

draft-ietf-ipsec-isakmp-xauth-06, *Extended Authentication within ISAKMP/Oakley (XAUTH)*

RFC 2401, *Security Architecture for the Internet Protocol*

RFC 2403, *The Use of HMAC-MD5-96 within ESP and AH*

RFC 2404, *The Use of HMAC-SHA-1-96 within ESP and AH*

RFC 2405, *The ESP DES-CBC Cipher Algorithm With Explicit IV*

RFC 2406, *IP Encapsulating Security Payload (ESP)*

RFC 2407, *IPsec Domain of Interpretation for ISAKMP (IPsec DoI)*

RFC 2408, *Internet Security Association and Key Management Protocol (ISAKMP)*

RFC 2409, *The Internet Key Exchange (IKE)*

RFC 2410, *The NULL Encryption Algorithm and Its Use With IPsec*

RFC 2560, *X.509 Internet Public Key Infrastructure Online Certificate Status Protocol - OCSP*

RFC 3526, *More Modular Exponential (MODP) Diffie-Hellman group for Internet Key Exchange (IKE)*

RFC 3566, *The AES-XCBC-MAC-96 Algorithm and Its Use With IPsec*

RFC 3602, *The AES-CBC Cipher Algorithm and Its Use with IPsec*

RFC 3706, *A Traffic-Based Method of Detecting Dead Internet Key Exchange (IKE) Peers*

RFC 3947, *Negotiation of NAT-Traversal in the IKE*

RFC 3948, *UDP Encapsulation of IPsec ESP Packets*

RFC 4106, *The Use of Galois/Counter Mode (GCM) in IPsec ESP*

RFC 4109, *Algorithms for Internet Key Exchange version 1 (IKEv1)*

RFC 4301, *Security Architecture for the Internet Protocol*

RFC 4303, *IP Encapsulating Security Payload*

RFC 4307, *Cryptographic Algorithms for Use in the Internet Key Exchange Version 2 (IKEv2)*

RFC 4308, *Cryptographic Suites for IPsec*
RFC 4434, *The AES-XCBC-PRF-128 Algorithm for the Internet Key Exchange Protocol (IKE)*
RFC 4543, *The Use of Galois Message Authentication Code (GMAC) in IPsec ESP and AH*
RFC 4754, *IKE and IKEv2 Authentication Using the Elliptic Curve Digital Signature Algorithm (ECDSA)*
RFC 4835, *Cryptographic Algorithm Implementation Requirements for Encapsulating Security Payload (ESP) and Authentication Header (AH)*
RFC 4868, *Using HMAC-SHA-256, HMAC-SHA-384, and HMAC-SHA-512 with IPsec*
RFC 4945, *The Internet IP Security PKI Profile of IKEv1/ISAKMP, IKEv2 and PKIX*
RFC 5019, *The Lightweight Online Certificate Status Protocol (OCSP) Profile for High-Volume Environments*
RFC 5282, *Using Authenticated Encryption Algorithms with the Encrypted Payload of the IKEv2 Protocol*
RFC 5903, *ECP Groups for IKE and IKEv2*
RFC 5996, *Internet Key Exchange Protocol Version 2 (IKEv2)*
RFC 5998, *An Extension for EAP-Only Authentication in IKEv2*
RFC 6379, *Suite B Cryptographic Suites for IPsec*
RFC 6380, *Suite B Profile for Internet Protocol Security (IPsec)*
RFC 6960, *X.509 Internet Public Key Infrastructure Online Certificate Status Protocol - OCSP*
RFC 7296, *Internet Key Exchange Protocol Version 2 (IKEv2)*
RFC 7321, *Cryptographic Algorithm Implementation Requirements and Usage Guidance for Encapsulating Security Payload (ESP) and Authentication Header (AH)*
RFC 7383, *Internet Key Exchange Protocol Version 2 (IKEv2) Message Fragmentation*
RFC 7427, *Signature Authentication in the Internet Key Exchange Version 2 (IKEv2)*

11.18 Label Distribution Protocol (LDP)

draft-pdutta-mpls-ldp-adj-capability-00, *LDP Adjacency Capabilities*
draft-pdutta-mpls-ldp-v2-00, *LDP Version 2*
draft-pdutta-mpls-mlldp-up-redundancy-00, *Upstream LSR Redundancy for Multi-point LDP Tunnels*
draft-pdutta-mpls-multi-ldp-instance-00, *Multiple LDP Instances*
draft-pdutta-mpls-tldp-hello-reduce-04, *Targeted LDP Hello Reduction*
RFC 3037, *LDP Applicability*
RFC 3478, *Graceful Restart Mechanism for Label Distribution Protocol – helper mode*
RFC 5036, *LDP Specification*
RFC 5283, *LDP Extension for Inter-Area Label Switched Paths (LSPs)*
RFC 5443, *LDP IGP Synchronization*
RFC 5561, *LDP Capabilities*
RFC 5919, *Signaling LDP Label Advertisement Completion*

RFC 6388, *Label Distribution Protocol Extensions for Point-to-Multipoint and Multipoint-to-Multipoint Label Switched Paths*

RFC 6512, *Using Multipoint LDP When the Backbone Has No Route to the Root*

RFC 6826, *Multipoint LDP in-band signaling for Point-to-Multipoint and Multipoint-to-Multipoint Label Switched Paths*

RFC 7032, *LDP Downstream-on-Demand in Seamless MPLS*

RFC 7473, *Controlling State Advertisements of Non-negotiated LDP Applications*

RFC 7552, *Updates to LDP for IPv6*

11.19 Layer Two Tunneling Protocol (L2TP) Network Server (LNS)

draft-mammoliti-l2tp-accessline-avp-04, *Layer 2 Tunneling Protocol (L2TP) Access Line Information Attribute Value Pair (AVP) Extensions*

RFC 2661, *Layer Two Tunneling Protocol "L2TP"*

RFC 2809, *Implementation of L2TP Compulsory Tunneling via RADIUS*

RFC 3438, *Layer Two Tunneling Protocol (L2TP) Internet Assigned Numbers: Internet Assigned Numbers Authority (IANA) Considerations Update*

RFC 3931, *Layer Two Tunneling Protocol - Version 3 (L2TPv3)*

RFC 4719, *Transport of Ethernet Frames over Layer 2 Tunneling Protocol Version 3 (L2TPv3)*

RFC 4951, *Fail Over Extensions for Layer 2 Tunneling Protocol (L2TP) "failover"*

11.20 Multiprotocol Label Switching (MPLS)

draft-ietf-mpls-lsp-ping-ospfv3-codepoint-02, *OSPFv3 CodePoint for MPLS LSP Ping*

RFC 3031, *Multiprotocol Label Switching Architecture*

RFC 3032, *MPLS Label Stack Encoding*

RFC 3270, *Multi-Protocol Label Switching (MPLS) Support of Differentiated Services – E-LSP*

RFC 3443, *Time To Live (TTL) Processing in Multi-Protocol Label Switching (MPLS) Networks*

RFC 4023, *Encapsulating MPLS in IP or Generic Routing Encapsulation (GRE)*

RFC 4182, *Removing a Restriction on the use of MPLS Explicit NULL*

RFC 4950, *ICMP Extensions for Multiprotocol Label Switching*

RFC 5332, *MPLS Multicast Encapsulations*

RFC 5884, *Bidirectional Forwarding Detection (BFD) for MPLS Label Switched Paths (LSPs)*

RFC 6374, *Packet Loss and Delay Measurement for MPLS Networks – Delay Measurement, Channel Type 0x000C*

RFC 6424, *Mechanism for Performing Label Switched Path Ping (LSP Ping) over MPLS Tunnels*

RFC 6425, *Detecting Data Plane Failures in Point-to-Multipoint Multiprotocol Label Switching (MPLS) - Extensions to LSP Ping*

RFC 6790, *The Use of Entropy Labels in MPLS Forwarding*
RFC 7308, *Extended Administrative Groups in MPLS Traffic Engineering (MPLS-TE)*
RFC 7510, *Encapsulating MPLS in UDP*
RFC 7746, *Label Switched Path (LSP) Self-Ping*
RFC 7876, *UDP Return Path for Packet Loss and Delay Measurement for MPLS Networks – Delay Measurement*
RFC 8029, *Detecting Multiprotocol Label Switched (MPLS) Data-Plane Failures*

11.21 Multiprotocol Label Switching - Transport Profile (MPLS-TP)

RFC 5586, *MPLS Generic Associated Channel*
RFC 5921, *A Framework for MPLS in Transport Networks*
RFC 5960, *MPLS Transport Profile Data Plane Architecture*
RFC 6370, *MPLS Transport Profile (MPLS-TP) Identifiers*
RFC 6378, *MPLS Transport Profile (MPLS-TP) Linear Protection*
RFC 6426, *MPLS On-Demand Connectivity and Route Tracing*
RFC 6427, *MPLS Fault Management Operations, Administration, and Maintenance (OAM)*
RFC 6428, *Proactive Connectivity Verification, Continuity Check and Remote Defect indication for MPLS Transport Profile*
RFC 6478, *Pseudowire Status for Static Pseudowires*
RFC 7213, *MPLS Transport Profile (MPLS-TP) Next-Hop Ethernet Addressing*

11.22 Network Address Translation (NAT)

draft-ietf-behave-address-format-10, *IPv6 Addressing of IPv4/IPv6 Translators*
draft-ietf-behave-v6v4-xlate-23, *IP/ICMP Translation Algorithm*
draft-miles-behave-l2nat-00, *Layer2-Aware NAT*
draft-nishitani-cgn-02, *Common Functions of Large Scale NAT (LSN)*
RFC 4787, *Network Address Translation (NAT) Behavioral Requirements for Unicast UDP*
RFC 5382, *NAT Behavioral Requirements for TCP*
RFC 5508, *NAT Behavioral Requirements for ICMP*
RFC 6146, *Stateful NAT64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers*
RFC 6333, *Dual-Stack Lite Broadband Deployments Following IPv4 Exhaustion*
RFC 6334, *Dynamic Host Configuration Protocol for IPv6 (DHCPv6) Option for Dual-Stack Lite*
RFC 6887, *Port Control Protocol (PCP)*
RFC 6888, *Common Requirements For Carrier-Grade NATs (CGNs)*
RFC 7753, *Port Control Protocol (PCP) Extension for Port-Set Allocation*

RFC 7915, *IP/ICMP Translation Algorithm*

11.23 Network Configuration Protocol (NETCONF)

RFC 5277, *NETCONF Event Notifications*

RFC 6020, *YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)*

RFC 6022, *YANG Module for NETCONF Monitoring*

RFC 6241, *Network Configuration Protocol (NETCONF)*

RFC 6242, *Using the NETCONF Protocol over Secure Shell (SSH)*

RFC 6243, *With-defaults Capability for NETCONF*

RFC 8342, *Network Management Datastore Architecture (NMDA) – Startup, Candidate, Running and Intended datastores*

RFC 8525, *YANG Library*

RFC 8526, *NETCONF Extensions to Support the Network Management Datastore Architecture – <get-data> operation*

11.24 Open Shortest Path First (OSPF)

RFC 1765, *OSPF Database Overflow*

RFC 2328, *OSPF Version 2*

RFC 3101, *The OSPF Not-So-Stubby Area (NSSA) Option*

RFC 3509, *Alternative Implementations of OSPF Area Border Routers*

RFC 3623, *Graceful OSPF Restart Graceful OSPF Restart – helper mode*

RFC 3630, *Traffic Engineering (TE) Extensions to OSPF Version 2*

RFC 4222, *Prioritized Treatment of Specific OSPF Version 2 Packets and Congestion Avoidance*

RFC 4552, *Authentication/Confidentiality for OSPFv3*

RFC 4576, *Using a Link State Advertisement (LSA) Options Bit to Prevent Looping in BGP/MPLS IP Virtual Private Networks (VPNs)*

RFC 4577, *OSPF as the Provider/Customer Edge Protocol for BGP/MPLS IP Virtual Private Networks (VPNs)*

RFC 5185, *OSPF Multi-Area Adjacency*

RFC 5187, *OSPFv3 Graceful Restart – helper mode*

RFC 5243, *OSPF Database Exchange Summary List Optimization*

RFC 5250, *The OSPF Opaque LSA Option*

RFC 5309, *Point-to-Point Operation over LAN in Link State Routing Protocols*

RFC 5340, *OSPF for IPv6*

RFC 5642, *Dynamic Hostname Exchange Mechanism for OSPF*

RFC 5709, *OSPFv2 HMAC-SHA Cryptographic Authentication*
RFC 5838, *Support of Address Families in OSPFv3*
RFC 6549, *OSPFv2 Multi-Instance Extensions*
RFC 6987, *OSPF Stub Router Advertisement*
RFC 7471, *OSPF Traffic Engineering (TE) Metric Extensions – Min/Max Unidirectional Link Delay metric for flex-algo, RSVP, SR-TE*
RFC 7684, *OSPFv2 Prefix/Link Attribute Advertisement*
RFC 7770, *Extensions to OSPF for Advertising Optional Router Capabilities*
RFC 8362, *OSPFv3 Link State Advertisement (LSA) Extensibility*
RFC 8920, *OSPF Application-Specific Link Attributes*

11.25 OpenFlow

TS-007 Version 1.3.1, *OpenFlow Switch Specification* – OpenFlow-hybrid switches

11.26 Path Computation Element Protocol (PCEP)

draft-ietf-pce-binding-label-sid-15, *Carrying Binding Label/Segment Identifier (SID) in PCE-based Networks*. – MPLS binding SIDs
draft-alvarez-pce-path-profiles-04, *PCE Path Profiles*
draft-dhs-spring-pce-sr-p2mp-policy-00, *PCEP extensions for p2mp sr policy*
RFC 5440, *Path Computation Element (PCE) Communication Protocol (PCEP)*
RFC 8231, *Path Computation Element Communication Protocol (PCEP) Extensions for Stateful PCE*
RFC 8253, *PCEPS: Usage of TLS to Provide a Secure Transport for the Path Computation Element Communication Protocol (PCEP)*
RFC 8281, *PCEP Extensions for PCE-initiated LSP Setup in a Stateful PCE Model*
RFC 8408, *Conveying Path Setup Type in PCE Communication Protocol (PCEP) Messages*
RFC 8664, *Path Computation Element Communication Protocol (PCEP) Extensions for Segment Routing*

11.27 Point-to-Point Protocol (PPP)

RFC 1332, *The PPP Internet Protocol Control Protocol (IPCP)*
RFC 1990, *The PPP Multilink Protocol (MP)*
RFC 1994, *PPP Challenge Handshake Authentication Protocol (CHAP)*
RFC 2516, *A Method for Transmitting PPP Over Ethernet (PPPoE)*
RFC 4638, *Accommodating a Maximum Transit Unit/Maximum Receive Unit (MTU/MRU) Greater Than 1492 in the Point-to-Point Protocol over Ethernet (PPPoE)*

RFC 5072, *IP Version 6 over PPP*

11.28 Policy management and credit control

3GPP TS 29.212 Release 11, *Policy and Charging Control (PCC); Reference points – Gx support as it applies to wireline environment (BNG)*

RFC 4006, *Diameter Credit-Control Application*

RFC 6733, *Diameter Base Protocol*

11.29 Pseudowire (PW)

draft-ietf-l2vpn-vpws-iw-oam-04, *OAM Procedures for VPWS Interworking*

MFA Forum 12.0.0, *Multiservice Interworking - Ethernet over MPLS*

MFA Forum 13.0.0, *Fault Management for Multiservice Interworking v1.0*

MFA Forum 16.0.0, *Multiservice Interworking - IP over MPLS*

RFC 3916, *Requirements for Pseudo-Wire Emulation Edge-to-Edge (PWE3)*

RFC 3985, *Pseudo Wire Emulation Edge-to-Edge (PWE3)*

RFC 4385, *Pseudo Wire Emulation Edge-to-Edge (PWE3) Control Word for Use over an MPLS PSN*

RFC 4446, *IANA Allocations for Pseudowire Edge to Edge Emulation (PWE3)*

RFC 4447, *Pseudowire Setup and Maintenance Using the Label Distribution Protocol (LDP)*

RFC 4448, *Encapsulation Methods for Transport of Ethernet over MPLS Networks*

RFC 5085, *Pseudowire Virtual Circuit Connectivity Verification (VCCV): A Control Channel for Pseudowires*

RFC 5659, *An Architecture for Multi-Segment Pseudowire Emulation Edge-to-Edge*

RFC 5885, *Bidirectional Forwarding Detection (BFD) for the Pseudowire Virtual Circuit Connectivity Verification (VCCV)*

RFC 6073, *Segmented Pseudowire*

RFC 6310, *Pseudowire (PW) Operations, Administration, and Maintenance (OAM) Message Mapping*

RFC 6391, *Flow-Aware Transport of Pseudowires over an MPLS Packet Switched Network*

RFC 6575, *Address Resolution Protocol (ARP) Mediation for IP Interworking of Layer 2 VPNs*

RFC 6718, *Pseudowire Redundancy*

RFC 6829, *Label Switched Path (LSP) Ping for Pseudowire Forwarding Equivalence Classes (FECs) Advertised over IPv6*

RFC 6870, *Pseudowire Preferential Forwarding Status bit*

RFC 7023, *MPLS and Ethernet Operations, Administration, and Maintenance (OAM) Interworking*

RFC 7267, *Dynamic Placement of Multi-Segment Pseudowires*

RFC 7392, *Explicit Path Routing for Dynamic Multi-Segment Pseudowires – ER-TLV and ER-HOP IPv4 Prefix*

RFC 8395, *Extensions to BGP-Signaled Pseudowires to Support Flow-Aware Transport Labels*

11.30 Quality of Service (QoS)

RFC 2430, *A Provider Architecture for Differentiated Services and Traffic Engineering (PASTE)*

RFC 2474, *Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers*

RFC 2597, *Assured Forwarding PHB Group*

RFC 3140, *Per Hop Behavior Identification Codes*

RFC 3246, *An Expedited Forwarding PHB (Per-Hop Behavior)*

11.31 Remote Authentication Dial In User Service (RADIUS)

draft-oscca-cfrg-sm3-02, *The SM3 Cryptographic Hash Function*

RFC 2865, *Remote Authentication Dial In User Service (RADIUS)*

RFC 2866, *RADIUS Accounting*

RFC 2867, *RADIUS Accounting Modifications for Tunnel Protocol Support*

RFC 2868, *RADIUS Attributes for Tunnel Protocol Support*

RFC 2869, *RADIUS Extensions*

RFC 3162, *RADIUS and IPv6*

RFC 4818, *RADIUS Delegated-IPv6-Prefix Attribute*

RFC 5176, *Dynamic Authorization Extensions to RADIUS*

RFC 6613, *RADIUS over TCP – with TLS*

RFC 6614, *Transport Layer Security (TLS) Encryption for RADIUS*

RFC 6929, *Remote Authentication Dial-In User Service (RADIUS) Protocol Extensions*

RFC 6911, *RADIUS attributes for IPv6 Access Networks*

11.32 Resource Reservation Protocol - Traffic Engineering (RSVP-TE)

draft-newton-mpls-te-dynamic-overbooking-00, *A Diffserv-TE Implementation Model to dynamically change booking factors during failure events*

RFC 2702, *Requirements for Traffic Engineering over MPLS*

RFC 2747, *RSVP Cryptographic Authentication*

RFC 2961, *RSVP Refresh Overhead Reduction Extensions*

RFC 3097, *RSVP Cryptographic Authentication -- Updated Message Type Value*

RFC 3209, *RSVP-TE: Extensions to RSVP for LSP Tunnels*

RFC 3477, *Signalling Unnumbered Links in Resource ReSerVation Protocol - Traffic Engineering (RSVP-TE)*

RFC 3564, *Requirements for Support of Differentiated Services-aware MPLS Traffic Engineering*
RFC 3906, *Calculating Interior Gateway Protocol (IGP) Routes Over Traffic Engineering Tunnels*
RFC 4090, *Fast Reroute Extensions to RSVP-TE for LSP Tunnels*
RFC 4124, *Protocol Extensions for Support of Diffserv-aware MPLS Traffic Engineering*
RFC 4125, *Maximum Allocation Bandwidth Constraints Model for Diffserv-aware MPLS Traffic Engineering*
RFC 4127, *Russian Dolls Bandwidth Constraints Model for Diffserv-aware MPLS Traffic Engineering*
RFC 4561, *Definition of a Record Route Object (RRO) Node-Id Sub-Object*
RFC 4875, *Extensions to Resource Reservation Protocol - Traffic Engineering (RSVP-TE) for Point-to-Multipoint TE Label Switched Paths (LSPs)*
RFC 5712, *MPLS Traffic Engineering Soft Preemption*
RFC 5817, *Graceful Shutdown in MPLS and Generalized MPLS Traffic Engineering Networks*

11.33 Routing Information Protocol (RIP)

RFC 1058, *Routing Information Protocol*
RFC 2080, *RIPng for IPv6*
RFC 2082, *RIP-2 MD5 Authentication*
RFC 2453, *RIP Version 2*

11.34 Segment Routing (SR)

draft-ietf-bess-mvpn-evpn-sr-p2mp-07, *Multicast and Ethernet VPN with Segment Routing P2MP and Ingress Replication – MVPN*
draft-bashandy-rtgwg-segment-routing-uloop-15, *Loop avoidance using Segment Routing*
draft-filsfils-spring-net-pgm-extension-srv6-usid-15, *Network Programming extension: SRv6 uSID instruction*
draft-filsfils-spring-srv6-net-pgm-insertion-08, *SRv6 NET-PGM extension: Insertion*
draft-ietf-idr-bgppls-srv6-ext-14, *BGP Link State Extensions for SRv6*
draft-ietf-idr-segment-routing-te-policy-23, *Advertising Segment Routing Policies in BGP*
draft-ietf-idr-ts-flowspec-srv6-policy-03, *Traffic Steering using BGP FlowSpec with SR Policy*
draft-ietf-pim-p2mp-policy-ping-03, *P2MP Policy Ping*
draft-ietf-pim-sr-p2mp-policy-06, *Segment Routing Point-to-Multipoint Policy – MPLS*
draft-ietf-rtgwg-segment-routing-ti-lfa-11, *Topology Independent Fast Reroute using Segment Routing*
draft-ietf-spring-conflict-resolution-05, *Segment Routing MPLS Conflict Resolution*
draft-ietf-spring-sr-replication-segment-16, *SR Replication segment for Multi-point Service Delivery – MPLS*
draft-ietf-spring-srv6-srh-compression-xx, *Compressed SRv6 Segment List Encoding in SRH*

draft-voyer-6man-extension-header-insertion-10, *Deployments With Insertion of IPv6 Segment Routing Headers*

RFC 8287, *Label Switched Path (LSP) Ping/Traceroute for Segment Routing (SR) IGP-Prefix and IGP-Adjacency Segment Identifiers (SIDs) with MPLS Data Planes*

RFC 8426, *Recommendations for RSVP-TE and Segment Routing (SR) Label Switched Path (LSP) Coexistence*

RFC 8476, *Signaling Maximum SID Depth (MSD) Using OSPF – node MSD*

RFC 8491, *Signaling Maximum SID Depth (MSD) Using IS-IS – node MSD*

RFC 8660, *Segment Routing with the MPLS Data Plane*

RFC 8661, *Segment Routing MPLS Interworking with LDP*

RFC 8663, *MPLS Segment Routing over IP – BGP SR with SR-MPLS-over-UDP/IP*

RFC 8665, *OSPF Extensions for Segment Routing*

RFC 8666, *OSPFv3 Extensions for Segment Routing*

RFC 8667, *IS-IS Extensions for Segment Routing*

RFC 8669, *Segment Routing Prefix Segment Identifier Extensions for BGP*

RFC 8754, *IPv6 Segment Routing Header (SRH)*

RFC 8814, *Signaling Maximum SID Depth (MSD) Using the Border Gateway Protocol - Link State*

RFC 8986, *Segment Routing over IPv6 (SRv6) Network Programming*

RFC 9085, *Border Gateway Protocol - Link State (BGP-LS) Extensions for Segment Routing*

RFC 9088, *Signaling Entropy Label Capability and Entropy Readable Label Depth Using IS-IS – advertising ELC*

RFC 9089, *Signaling Entropy Label Capability and Entropy Readable Label Depth Using OSPF – advertising ELC*

RFC 9252, *BGP Overlay Services Based on Segment Routing over IPv6 (SRv6)*

RFC 9256, *Segment Routing Policy Architecture*

RFC 9259, *Operations, Administration, and Maintenance (OAM) in Segment Routing over IPv6 (SRv6)*

RFC 9350, *IGP Flexible Algorithm*

RFC 9352, *IS-IS Extensions to Support Segment Routing over the IPv6 Data Plane*

11.35 Simple Network Management Protocol (SNMP)

draft-blumenthal-aes-usm-04, *The AES Cipher Algorithm in the SNMP's User-based Security Model – CFB128-AES-192 and CFB128-AES-256*

draft-ietf-isis-wg-mib-06, *Management Information Base for Intermediate System to Intermediate System (IS-IS)*

draft-ietf-mboned-msdp-mib-01, *Multicast Source Discovery protocol MIB*

draft-ietf-mpls-ldp-mib-07, *Definitions of Managed Objects for the Multiprotocol Label Switching, Label Distribution Protocol (LDP)*

draft-ietf-mpls-lsr-mib-06, Multiprotocol Label Switching (MPLS) Label Switching Router (LSR) Management Information Base Using SMIv2

draft-ietf-mpls-te-mib-04, Multiprotocol Label Switching (MPLS) Traffic Engineering Management Information Base

draft-ietf-ospf-mib-update-08, OSPF Version 2 Management Information Base

draft-ietf-rrp-unified-mib-06, Definitions of Managed Objects for the VRRP over IPv4 and IPv6 – IPv6

ESO-CONSORTIUM-MIB revision 200406230000Z, esoConsortiumMIB

IANA-ADDRESS-FAMILY-NUMBERS-MIB revision 200203140000Z, ianaAddressFamilyNumbers

IANAifType-MIB revision 200505270000Z, ianaifType

IANA-RTPROTO-MIB revision 200009260000Z, ianaRtProtoMIB

IEEE8021-CFM-MIB revision 200706100000Z, ieee8021CfmMib

IEEE8021-PAE-MIB revision 200101160000Z, ieee8021paeMIB

IEEE8023-LAG-MIB revision 200006270000Z, lagMIB

LLDP-MIB revision 200505060000Z, lldpMIB

RFC 1157, A Simple Network Management Protocol (SNMP)

RFC 1212, Concise MIB Definitions

RFC 1215, A Convention for Defining Traps for use with the SNMP

RFC 1724, RIP Version 2 MIB Extension

RFC 1901, Introduction to Community-based SNMPv2

RFC 2021, Remote Network Monitoring Management Information Base Version 2 using SMIv2

RFC 2206, RSVP Management Information Base using SMIv2

RFC 2213, Integrated Services Management Information Base using SMIv2

RFC 2494, Definitions of Managed Objects for the DS0 and DS0 Bundle Interface Type

RFC 2578, Structure of Management Information Version 2 (SMIv2)

RFC 2579, Textual Conventions for SMIv2

RFC 2580, Conformance Statements for SMIv2

RFC 2787, Definitions of Managed Objects for the Virtual Router Redundancy Protocol

RFC 2819, Remote Network Monitoring Management Information Base

RFC 2856, Textual Conventions for Additional High Capacity Data Types

RFC 2863, The Interfaces Group MIB

RFC 2864, The Inverted Stack Table Extension to the Interfaces Group MIB

RFC 2933, Internet Group Management Protocol MIB

RFC 3014, Notification Log MIB

RFC 3165, Definitions of Managed Objects for the Delegation of Management Scripts

RFC 3231, Definitions of Managed Objects for Scheduling Management Operations

RFC 3273, Remote Network Monitoring Management Information Base for High Capacity Networks

RFC 3410, Introduction and Applicability Statements for Internet Standard Management Framework

RFC 3411, *An Architecture for Describing Simple Network Management Protocol (SNMP) Management Frameworks*

RFC 3412, *Message Processing and Dispatching for the Simple Network Management Protocol (SNMP)*

RFC 3413, *Simple Network Management Protocol (SNMP) Applications*

RFC 3414, *User-based Security Model (USM) for version 3 of the Simple Network Management Protocol (SNMPv3)*

RFC 3415, *View-based Access Control Model (VACM) for the Simple Network Management Protocol (SNMP)*

RFC 3416, *Version 2 of the Protocol Operations for the Simple Network Management Protocol (SNMP)*

RFC 3417, *Transport Mappings for the Simple Network Management Protocol (SNMP) – SNMP over UDP over IPv4*

RFC 3418, *Management Information Base (MIB) for the Simple Network Management Protocol (SNMP)*

RFC 3419, *Textual Conventions for Transport Addresses*

RFC 3498, *Definitions of Managed Objects for Synchronous Optical Network (SONET) Linear Automatic Protection Switching (APS) Architectures*

RFC 3584, *Coexistence between Version 1, Version 2, and Version 3 of the Internet-standard Network Management Framework*

RFC 3592, *Definitions of Managed Objects for the Synchronous Optical Network/Synchronous Digital Hierarchy (SONET/SDH) Interface Type*

RFC 3593, *Textual Conventions for MIB Modules Using Performance History Based on 15 Minute Intervals*

RFC 3635, *Definitions of Managed Objects for the Ethernet-like Interface Types*

RFC 3637, *Definitions of Managed Objects for the Ethernet WAN Interface Sublayer*

RFC 3826, *The Advanced Encryption Standard (AES) Cipher Algorithm in the SNMP User-based Security Model*

RFC 3877, *Alarm Management Information Base (MIB)*

RFC 3895, *Definitions of Managed Objects for the DS1, E1, DS2, and E2 Interface Types*

RFC 3896, *Definitions of Managed Objects for the DS3/E3 Interface Type*

RFC 4001, *Textual Conventions for Internet Network Addresses*

RFC 4022, *Management Information Base for the Transmission Control Protocol (TCP)*

RFC 4113, *Management Information Base for the User Datagram Protocol (UDP)*

RFC 4220, *Traffic Engineering Link Management Information Base*

RFC 4273, *Definitions of Managed Objects for BGP-4*

RFC 4292, *IP Forwarding Table MIB*

RFC 4293, *Management Information Base for the Internet Protocol (IP)*

RFC 4631, *Link Management Protocol (LMP) Management Information Base (MIB)*

RFC 4878, *Definitions and Managed Objects for Operations, Administration, and Maintenance (OAM) Functions on Ethernet-Like Interfaces*

RFC 7420, *Path Computation Element Communication Protocol (PCEP) Management Information Base (MIB) Module*

RFC 7630, *HMAC-SHA-2 Authentication Protocols in the User-based Security Model (USM) for SNMPv3*
SFLOW-MIB revision 200309240000Z, *sFlowMIB*

11.36 Timing

GR-1244-CORE Issue 3, *Clocks for the Synchronized Network: Common Generic Criteria*
GR-253-CORE Issue 3, *SONET Transport Systems: Common Generic Criteria*
IEEE 1588-2008, *IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems*
ITU-T G.781, *Synchronization layer functions*
ITU-T G.811, *Timing characteristics of primary reference clocks*
ITU-T G.813, *Timing characteristics of SDH equipment slave clocks (SEC)*
ITU-T G.8261, *Timing and synchronization aspects in packet networks*
ITU-T G.8262, *Timing characteristics of synchronous Ethernet equipment slave clock (EEC)*
ITU-T G.8262.1, *Timing characteristics of an enhanced synchronous Ethernet equipment slave clock (eEEC)*
ITU-T G.8264, *Distribution of timing information through packet networks*
ITU-T G.8265.1, *Precision time protocol telecom profile for frequency synchronization*
ITU-T G.8272, *Timing characteristics of primary reference time clocks – PRTC-A, PRTC-B*
ITU-T G.8275.1, *Precision time protocol telecom profile for phase/time synchronization with full timing support from the network*
ITU-T G.8275.2, *Precision time protocol telecom profile for phase/time synchronization with partial timing support from the network*
RFC 3339, *Date and Time on the Internet: Timestamps*
RFC 5905, *Network Time Protocol Version 4: Protocol and Algorithms Specification*

11.37 Two-Way Active Measurement Protocol (TWAMP)

RFC 5357, *A Two-Way Active Measurement Protocol (TWAMP) – server, unauthenticated mode*
RFC 5938, *Individual Session Control Feature for the Two-Way Active Measurement Protocol (TWAMP)*
RFC 6038, *Two-Way Active Measurement Protocol (TWAMP) Reflect Octets and Symmetrical Size Features*
RFC 8545, *Well-Known Port Assignments for the One-Way Active Measurement Protocol (OWAMP) and the Two-Way Active Measurement Protocol (TWAMP) – TWAMP*
RFC 8762, *Simple Two-Way Active Measurement Protocol – unauthenticated*
RFC 8972, *Simple Two-Way Active Measurement Protocol Optional Extensions – unauthenticated*

11.38 Virtual Private LAN Service (VPLS)

RFC 4761, *Virtual Private LAN Service (VPLS) Using BGP for Auto-Discovery and Signaling*
RFC 4762, *Virtual Private LAN Service (VPLS) Using Label Distribution Protocol (LDP) Signaling*
RFC 5501, *Requirements for Multicast Support in Virtual Private LAN Services*
RFC 6074, *Provisioning, Auto-Discovery, and Signaling in Layer 2 Virtual Private Networks (L2VPNs)*
RFC 7041, *Extensions to the Virtual Private LAN Service (VPLS) Provider Edge (PE) Model for Provider Backbone Bridging*
RFC 7117, *Multicast in Virtual Private LAN Service (VPLS)*

11.39 Voice and video

DVB BlueBook A86, *Transport of MPEG-2 TS Based DVB Services over IP Based Networks*
ETSI TS 101 329-5 Annex E, *QoS Measurement for VoIP - Method for determining an Equipment Impairment Factor using Passive Monitoring*
ITU-T G.1020 Appendix I, *Performance Parameter Definitions for Quality of Speech and other Voiceband Applications Utilizing IP Networks - Mean Absolute Packet Delay Variation & Markov Models*
ITU-T G.107, *The E Model - A computational model for use in planning*
ITU-T P.564, *Conformance testing for voice over IP transmission quality assessment models*
RFC 3550, *RTP: A Transport Protocol for Real-Time Applications – Appendix A.8*
RFC 4585, *Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)*
RFC 4588, *RTP Retransmission Payload Format*

11.40 Wireless Local Area Network (WLAN) gateway

3GPP TS 23.402, *Architecture enhancements for non-3GPP accesses – S2a roaming based on GPRS*

11.41 Yet Another Next Generation (YANG)

RFC 6991, *Common YANG Data Types*
RFC 7950, *The YANG 1.1 Data Modeling Language*
RFC 7951, *JSON Encoding of Data Modeled with YANG*

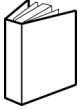
11.42 Yet Another Next Generation (YANG) OpenConfig Models

openconfig-aaa.yang version 0.4.0, *OpenConfig AAA Model*
openconfig-aaa-radius.yang version 0.3.0, *OpenConfig AAA RADIUS Model*

openconfig-aaa-tacacs.yang version 0.3.0, *OpenConfig AAA TACACS+ Model*
openconfig-acl.yang version 1.0.0, *OpenConfig ACL Model*
openconfig-alarms.yang version 0.3.2, *OpenConfig System Alarms Model*
openconfig-bfd.yang version 0.2.2, *OpenConfig BFD Model*
openconfig-bgp.yang version 6.1.0, *OpenConfig BGP Model*
openconfig-bgp-common.yang version 6.0.0, *OpenConfig BGP Common Model*
openconfig-bgp-common-multiprotocol.yang version 6.0.0, *OpenConfig BGP Common Multiprotocol Model*
openconfig-bgp-common-structure.yang version 6.0.0, *OpenConfig BGP Common Structure Model*
openconfig-bgp-global.yang version 6.0.0, *OpenConfig BGP Global Model*
openconfig-bgp-neighbor.yang version 6.1.0, *OpenConfig BGP Neighbor Model*
openconfig-bgp-peer-group.yang version 6.1.0, *OpenConfig BGP Peer Group Model*
openconfig-bgp-policy.yang version 4.0.1, *OpenConfig BGP Policy Model*
openconfig-if-aggregate.yang version 2.4.3, *OpenConfig Interfaces Aggregated Model*
openconfig-if-ethernet.yang version 2.12.1, *OpenConfig Interfaces Ethernet Model*
openconfig-if-ip.yang version 3.1.0, *OpenConfig Interfaces IP Model*
openconfig-if-ip-ext.yang version 2.3.1, *OpenConfig Interfaces IP Extensions Model*
openconfig-igmp.yang version 0.2.0, *OpenConfig IGMP Model*
openconfig-interfaces.yang version 3.0.0, *OpenConfig Interfaces Model*
openconfig-isis.yang version 1.1.0, *OpenConfig IS-IS Model*
openconfig-isis-policy.yang version 0.5.0, *OpenConfig IS-IS Policy Model*
openconfig-isis-routing.yang version 1.1.0, *OpenConfig IS-IS Routing Model*
openconfig-lacp.yang version 1.3.0, *OpenConfig LACP Model*
openconfig-lldp.yang version 0.1.0, *OpenConfig LLDP Model*
openconfig-local-routing.yang version 1.2.0, *OpenConfig Local Routing Model*
openconfig-mpls.yang version 2.3.0, *OpenConfig MPLS Model*
openconfig-mpls-ldp.yang version 3.0.2, *OpenConfig MPLS LDP Model*
openconfig-mpls-rsvp.yang version 2.3.0, *OpenConfig MPLS RSVP Model*
openconfig-mpls-te.yang version 2.3.0, *OpenConfig MPLS TE Model*
openconfig-network-instance.yang version 1.1.0, *OpenConfig Network Instance Model*
openconfig-network-instance-l3.yang version 0.11.1, *OpenConfig L3 Network Instance Model – static routes*
openconfig-ospfv2.yang version 0.4.0, *OpenConfig OSPFv2 Model*
openconfig-ospfv2-area.yang version 0.4.0, *OpenConfig OSPFv2 Area Model*
openconfig-ospfv2-area-interface.yang version 0.4.0, *OpenConfig OSPFv2 Area Interface Model*
openconfig-ospfv2-common.yang version 0.4.0, *OpenConfig OSPFv2 Common Model*
openconfig-ospfv2-global.yang version 0.4.0, *OpenConfig OSPFv2 Global Model*
openconfig-packet-match.yang version 1.0.0, *OpenConfig Packet Match Model*

openconfig-pim.yang version 0.2.0, *OpenConfig PIM Model*
openconfig-platform.yang version 0.15.0, *OpenConfig Platform Model*
openconfig-platform-fan.yang version 0.1.1, *OpenConfig Platform Fan Model*
openconfig-platform-linecard.yang version 0.1.2, *OpenConfig Platform Linecard Model*
openconfig-platform-port.yang version 0.4.2, *OpenConfig Port Model*
openconfig-platform-transceiver.yang version 0.9.0, *OpenConfig Transceiver Model*
openconfig-procmon.yang version 0.4.0, *OpenConfig Process Monitoring Model*
openconfig-relay-agent.yang version 0.1.0, *OpenConfig Relay Agent Model*
openconfig-routing-policy.yang version 3.0.0, *OpenConfig Routing Policy Model*
openconfig-rsvp-sr-ext.yang version 0.1.0, *OpenConfig RSVP-TE and SR Extensions Model*
openconfig-system.yang version 0.10.1, *OpenConfig System Model*
openconfig-system-grpc.yang version 1.0.0, *OpenConfig System gRPC Model*
openconfig-system-logging.yang version 0.3.1, *OpenConfig System Logging Model*
openconfig-system-terminal.yang version 0.3.0, *OpenConfig System Terminal Model*
openconfig-telemetry.yang version 0.5.0, *OpenConfig Telemetry Model*
openconfig-terminal-device.yang version 1.9.0, *OpenConfig Terminal Optics Device Model*
openconfig-vlan.yang version 2.0.0, *OpenConfig VLAN Model*

Customer document and product support



Customer documentation

[Customer documentation welcome page](#)



Technical support

[Product support portal](#)



Documentation feedback

[Customer documentation feedback](#)