

MidoNet Operations Guide

5.4-SNAPSHOT (2017-03-22 12:04 UTC)

DRAFT



mido**net**

docs.midonet.org

MidoNet Operations Guide

5.4-SNAPSHOT (2017-03-22 12:04 UTC)

Copyright © 2017 Midokura SARL All rights reserved.

MidoNet is a network virtualization software for Infrastructure-as-a-Service (IaaS) clouds.

It decouples your IaaS cloud from your network hardware, creating an intelligent software abstraction layer between your end hosts and your physical network.

This guide includes instructions on creating routers, bridges, and ports. It also describes rule chains and several MidoNet features, including L4 load balancing, resource protection, NAT configuration, handling IP packet fragments, and L2 address matching.



Caution

This document is a DRAFT. It may be MISSING relevant information or contain UNTESTED information. Use it at your own risk.



Note

Please consult the [MidoNet Mailing Lists or Chat](#) if you need assistance.

Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

Table of Contents

Preface	ix
Conventions	ix
1. Configuring uplinks	1
Edge Router Setup	1
BGP Setup	2
Static Setup	6
2. Authentication and authorization	9
Available authentication services in MidoNet	9
Using the Keystone authentication service	11
3. Admitting resources to MidoNet	13
What are tunnel zones?	13
Working with hosts	14
4. Device abstractions	17
Creating a router	17
Adding a port to a router	17
Adding a bridge	18
Adding a port to a bridge	18
Binding an exterior port to a host	18
Stateful port groups	19
5. Connecting devices	21
Connecting a bridge to a router	21
Connecting two routers	22
6. Routing	23
Routing process overview	23
Viewing routes	25
Adding routes	26
Deleting routes	27
ECMP Limitations	28
7. Rule chains	29
A packet's flow within a router	29
A packet's flow within a rule chain	30
Rule types	31
Rule order	32
Rule conditions	33
MidoNet rule chain examples	37
8. Network Address Translation	40
Static NAT	40
Viewing NAT rule chain information	40
Configuring SNAT, DNAT, and REV_DNAT	42
DNAT/REV_DNAT example	42
SNAT REV_SNAT example	43
9. Layer 4 Load Balancing	48
Load balancer configuration	48
Sticky Source IP	50
Health monitor	51
10. Load Balancing as a Service (LBaaS)	54
Neutron LBaaS Support	54
11. L2 address masking	56
L2 address mask rule chain example	56
12. Handling fragmented packets	57
Definitions and allowed values	57
Fragmented packets rule chain creation example	58

Non-Fragmented and Fragmented Packets	59
Fragmented and Non-Fragmented Packets with Different Destinations	60
13. MidoNet resource protection	62
Introduction	62
Expected Behavior	62
Configuration	63
Disabling Resource Protection	63
14. MidoNet monitoring	64
Metering	64
Monitoring Network State Database	65
Monitoring Midolman Agents	71
Monitoring events	74
Packet Tracing	83
Port mirroring	84
15. VXLAN configuration	88
VXLAN Gateway	88
VXLAN Coordinator	90
VXLAN Flooding Proxy	90
Connecting to the VTEP	90
Setting up a connection between a VTEP and a Neutron network	92
Enabling connection between VTEP and MidoNet hosts	93
VXLAN Gateway high availability (VTEP side, active-passive mode)	94
Troubleshooting VTEP/VXGW configuration	95
CLI commands used for working with the VXGW	100
16. Setting up an L2 gateway	106
Configuring an L2 gateway	107
Fail-over/Fail-back	108
17. Service Containers	109
Configuration	111
Management	112
Scheduling	113
Troubleshooting	118
18. Service Insertion and Chaining	121
Service Function Model	121
L2Insertion Object	122
L2Insertion in MidoNet CLI	124
Alternatives	124
19. Router Peering	126
Diagrams	128
Creating tenant networks and (peer) routers	130
Creating a VTEP Router Gateway Device	130
Creating a Multi-site Network and L2 Gateway Connection	131
Peering a Tenant Router	132
Adding a Remote MAC Entry (Endpoint) to the Multi-site Network	134
Configure Edge Bindings	134
Deleting a Remote MAC Entry	134
Deleting an L2 Gateway Connection	135
Deleting an L2 Gateway	135
Deleting a Gateway Device	135
Neutron CLI Gateway Device Reference	135
20. FWaaS Logging	141
Creating a router and firewall	141
Creating a logging resource and firewall log	142
Updating a logging resource and firewall log	143
Deleting a logging resource and firewall log	143

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

List of Tables

2.1. Authentication Provider Classes	9
2.2. Keystone Version	11
2.3. Keystone Protocol	11
7.1. CLI Rule Chain Attributes	33
7.2. CLI Rule Chain Attributes That Match Packets	34
14.1. Configuration Files/Locations	75
14.2. Event Message Files/Locations	75
29.1. Recommended Configuration Values	189
29.2. Admin Roles	192
29.3. System Properties for the HTTPS Key Store	193

Preface

Conventions

The MidoNet documentation uses several typesetting conventions.

Notices

Notices take these forms:



Note

A handy tip or reminder.



Important

Something you must be aware of before proceeding.



Warning

Critical information about the risk of data loss or security issues.

Command prompts

\$ prompt

Any user, including the root user, can run commands that are prefixed with the \$ prompt.

prompt

The root user must run commands that are prefixed with the # prompt. You can also prefix these commands with the **sudo** command, if available, to run them.

Create a standard neutron router:

```
neutron router-create <EDGE ROUTER NAME>
```

Attach the edge router to an external network:

```
neutron router-interface-add <EDGE_ROUTER_ID> <EXT_SUBNET_ID>
```

Create a special network called `uplink` network, representing the physical network outside of the cloud:

```
neutron net-create <UPLINK_NET_NAME> --tenant_id admin --
provider:network_type uplink
```

Create a subnet for the uplink network matching the CIDR used in the uplink network (could just be /30 if linked directly to another router):

```
neutron subnet-create --tenant_id admin --disable-dhcp --name
    <UPLINK SUBNET NAME> <UPLINK NET NAME> <CIDR>
```

Create a port on the uplink network with a specific IP that you want to use and the binding details so that this virtual port gets bound to a specific NIC on the gateway host:

```
neutron port-create <UPLINK_NET_ID> --binding:host_id <HOST_NAME> --
binding:profile type=dict interface_name=<INTERFACE_NAME> --fixed-ip
ip address=<IP_ADDR>
```

Attach the uplink port to the Edge Router:

```
neutron router-interface-add <EDGE_ROUTER_ID> port=<UPLINK_PORT_ID>
```

BGP Setup

You set up a BGP link to connect MidoNet with an external Autonomous System (AS). This creates an up-link to an external network.

Typically, you connect a MidoNet network to the Internet through two independent up-link routers. In the simple case of connecting MidoNet to the Internet via two BGP-enabled routers, it's best to create two ports on the virtual router and bind them to network interfaces in two different hosts (two Gateway Nodes). This will distribute load between both Gateway Nodes hosts, as well as eliminate a single point of failure. The two ports must be configured as a stateful port pair, for details, refer to [the section called "Stateful port groups" \[19\]](#).

MidoNet uses quagga's bgpd to terminate BGP sessions on behalf of a virtual router. The routes that bgpd learns from its peers are added to the virtual router in MidoNet's topology. The quagga package is provided in the MidoNet release package repositories. Any system running Midolman should already have everything that it needs to set up the BGP. You have to keep in mind that bgpd processes run in the host where a particular virtual router port is bound.



Important

Make sure you have the following information before setting up the BGP:
Local and peer Autonomous System (AS) Numbers for the BGP session. The
IP address of the BGP peer.

BGP Uplink Configuration

MidoNet utilizes the Border Gateway Protocol (BGP) for external connectivity.


```
midonet> router router0 add bgp-peer asn 64514 address 203.0.113.1
router0:peer1

midonet> router router0 list bgp-peer
peer peer0 asn 64513 address 198.51.100.1
peer peer1 asn 64514 address 203.0.113.1
```

4. If needed, configure MD5 authentication:

```
midonet> router router0 bgp-peer peer0 set password BGP_PASSWORD
midonet> router router0 bgp-peer peer1 set password BGP_PASSWORD
```

5. If needed, configure custom timers that will take precedence over the default ones defined in the MidoNet configuration:

```
midonet> router router0 bgp-peer peer0 set connect-retry 10
midonet> router router0 bgp-peer peer0 set hold-time 5
midonet> router router0 bgp-peer peer0 set keep-alive 5
midonet> router router0 bgp-peer peer1 set connect-retry 10
midonet> router router0 bgp-peer peer1 set hold-time 5
midonet> router router0 bgp-peer peer1 set keep-alive 5
midonet> router router0 list bgp-peer
peer peer0 asn 64513 address 198.51.100.1 keep-alive 5 hold-time 5
  connect-retry 10
peer peer1 asn 64514 address 203.0.113.1 keep-alive 5 hold-time 5
  connect-retry 10
```

6. Add routes to the remote BGP peers

In order to be able to establish connections to the remote BGP peers, corresponding routes have to be added.

```
midonet> router router0 route add src 0.0.0.0/0 dst 198.51.100.0/30
  port router0:port0 type normal
router0:route0

midonet> router router0 route add src 0.0.0.0/0 dst 203.0.113.0/30
  port router0:port1 type normal
router0:route1
```

7. Advertise BGP routes

In order to provide external connectivity for hosted virtual machines, the floating IP network has to be advertised to the BGP peers.

```
midonet> router router0 add bgp-network net 192.0.2.0/24
router0:net0

midonet> router router0 list bgp-network
net net0 net 192.0.2.0/24
```

Adding a second session on the same router port

It may be useful or a good idea to add a second BGP session to this router port if there is a second uplink router available. That has two upsides as the host that owns the port binding for this router port will be able to load balance among both upstream routers and it will not be disconnected if only one of them fails.

To add a second peer to the same router port, you use the same command as for the first peer, adjusting its AS number and IP address. The router port on which MidoNet establishes the BGP session is chosen automatically based on the peer's IP address.


```

midonet> host list
host host0 name controller alive true
midonet> host host0 add binding port router router0 port port0 interface
<Binding_Interface>
host host0 interface <Binding_Interface> port router0:port

```

Now we can reach VMs from the underlay host using their floating IPs, and VMs can reach external networks as well (as long as the host has external connectivity).

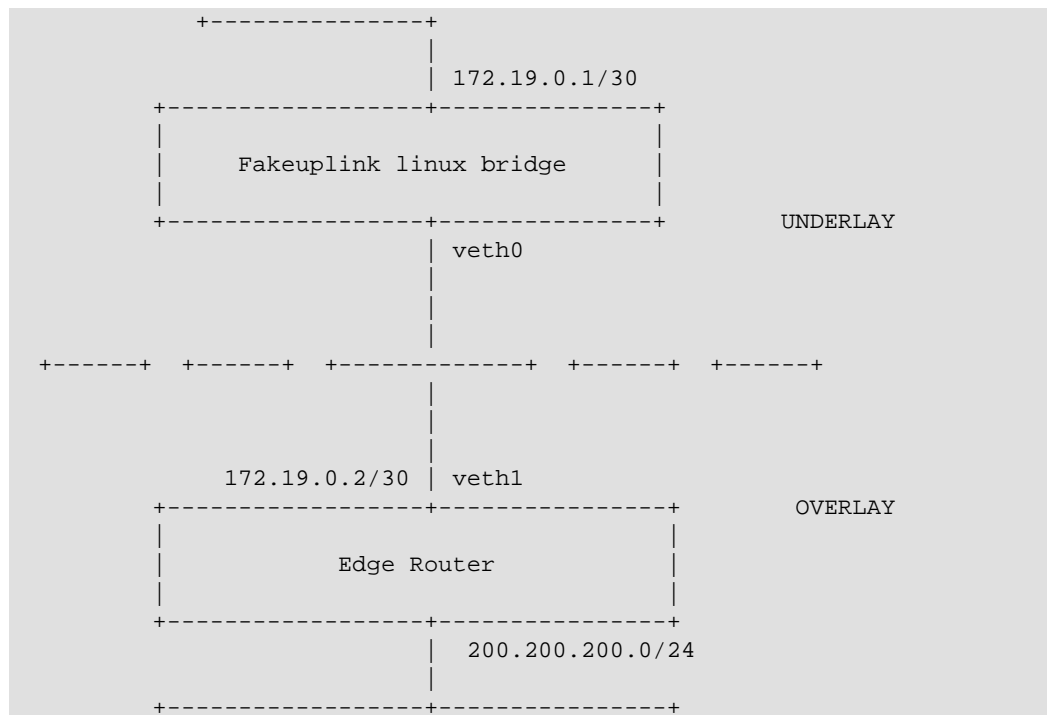
Fake Uplink Setup

If you are not connecting through a BGP link, or you just want to use static routing follow this section.

This creates a static up-link to connect VMs to the external network.

1. Create fake uplink

We are going to create the following topology to allow the VMs reach external networks:



2. Create a veth pair

```

# ip link add type veth
# ip link set dev veth0 up
# ip link set dev veth1 up

```

3. Create a bridge, set an IP address and attach veth0

```

# brctl addbr uplinkbridge
# brctl addif uplinkbridge veth0
# ip addr add 172.19.0.1/30 dev uplinkbridge
# ip link set dev uplinkbridge up

```

4. Enable IP forwarding

```

# sysctl -w net.ipv4.ip_forward=1

```


5. Route packets to 'external' network to the bridge

```
# ip route add 200.200.200.0/24 via 172.19.0.2
```

6. Create a port on the Edge Router and bind it to the veth:

```
$ midonet-cli
midonet> router list
router router0 name Edge Router state up
midonet> router router0 add port addresses 172.19.0.2/30
router0:port0
midonet> router router0 add route src 0.0.0.0/0 dst 0.0.0.0/0 type
normal port router router0 port port0 gw 172.19.0.1
midonet> host list
host host0 name controller alive true
midonet> host host0 add binding port router router0 port port0 interface
veth1
host host0 interface veth1 port router0:port
```

7. Add masquerading to your external interface so connections coming from the overlay with addresses that belong to the "fake" external network are NATed. Also make sure these packets can be forwarded:

```
# iptables -t nat -I POSTROUTING -o eth0 -s 200.200.200.0/24 -j
MASQUERADE
# iptables -I FORWARD -s 200.200.200.0/24 -j ACCEPT
```

Now we can reach VMs from the underlay host using their floating IPs, and VMs can reach external networks as well (as long as the host has external connectivity).



Note

MidoNet version 5.0.2 introduced support for Keystone version 3, and a new authentication provider class `org.midonet.cluster.auth.keystone.KeystoneService`, which replaces the previous `org.midonet.cluster.auth.keystone.v2_0.KeystoneService`. However, the legacy provider is kept for compatibility reasons and provides the same capabilities as the new Keystone authentication provider, supporting both Keystone version 2 and 3.



Important

Since MidoNet version 5.0.2, the default Keystone authentication is Keystone version 3. If you are upgrading from a previous MidoNet version, you must modify the Keystone configuration in order to continue using Keystone version 2. See [the section called “Enabling Keystone authentication” \[11\]](#).



Important

You must restart all MidoNet Cluster instances after changing the authentication provider.

To enable Keystone authentication, you must configure the MidoNet cluster with the credentials of an administrative user. The cluster uses this credentials to authenticate third-party requests to the MidoNet API. MidoNet supports two Keystone authentication methods:

- *Password authentication*, where you configure the user name and password credentials of an administrative user.
- *Token authentication*, where you configure the administrative token used by Keystone.



Important

The token authentication takes precedence over the password authentication.

For additional Keystone configuration options, see [the section called “Enabling Keystone authentication” \[11\]](#)

Mock authentication

Mock authentication disables the authentication system by returning fake tokens to authenticating clients, and ignoring the sent tokens during authorization. To enable the mock authentication configure the `MockAuthService` provider class.

```
echo "cluster.auth.provider_class : org.midonet.cluster.auth.
MockAuthService" | mn-conf set -t default
```



Warning

Mock authentication is the default authentication provider for the MidoNet Cluster. However, this mode is used for testing purposes but should not be used in production.

Using the Keystone authentication service

This section explains how to use the Keystone authentication service with MidoNet.

Enabling Keystone authentication

In order to use the OpenStack Keystone authentication service with MidoNet, you must configure the following keys in the MidoNet configuration.

`cluster.auth.provider_class`

The fully qualified path of the Java class that provides the Keystone authentication service.

```
org.midonet.cluster.auth.keystone.KeystoneService
```



Note

Since MidoNet version 5.0.2, the legacy authentication provider `org.midonet.cluster.auth.keystone.v2_0.KeystoneService` is available for compatibility reasons, but it provides the same capabilities as the new provider, supporting both Keystone version 2 and 3. The default Keystone version is 3.

`cluster.auth.admin_role`

Identifies the name of the admin role in MidoNet. The admin has read and write access to all the resources. We recommend re-using the OpenStack *admin* role. Optionally you can create a separate admin role for MidoNet.

`cluster.auth.keystone.version`

The version of the Keystone API. The following values are allowed:

Table 2.2. Keystone Version

Version	Description
2	Uses the Keystone version 2 API.
3	Uses the Keystone version 3 API (default).

`cluster.auth.keystone.protocol`

Identifies the protocol used for the Keystone service. The following values are allowed:

Table 2.3. Keystone Protocol

Class	Description
http	Uses plain text HTTP to communicate with the Keystone service (default).
https	Uses encrypted HTTPS to communicate with the Keystone service (recommended).

`cluster.auth.keystone.host`

Identifies the host of the Keystone service (default is *localhost*).

Deleting tunnel zones

Use this procedure to delete a tunnel zone.

1. Enter the `list tunnel-zone` command to list the tunnel zones. For example:

```
midonet> list tunnel-zone
tzone tzone0 name new-tz type vxlan
tzone tzone1 name vxlan type vxlan
```

2. Enter the `delete tunnel-zone tz-alias` command to delete the desired tunnel zone. For example:

```
midonet> delete tunnel-zone tzone0
```

Specify the dynamically assigned number of the alias for the tunnel zone to delete; in the above example, the assigned number is 0 (tzone0).

3. (Optional) Enter the command to list the tunnel zones to confirm the deletion:

```
midonet> list tunnel-zone
tzone tzone1 name vxlan type vxlan
```

Viewing tunnel zone information

Use this procedure to view tunnel zone information.

```
midonet> tunnel-zone tzone0 list member
zone tzone0 host host0 address 192.168.0.3
zone tzone0 host host1 address 192.168.0.5
zone tzone0 host host2 address 192.168.0.4
zone tzone0 host host3 address 192.168.0.6
```

The above output shows the:

- Aliases for the hosts in the tunnel zone (host0, host1, and so on)
- IP addresses assigned to the hosts

Working with hosts

This section shows how to view host information and admit new hosts to a tunnel zone.

Viewing host information

Use this procedure to view information about hosts.

- To list the hosts enter the command:

```
midonet> list host
host host0 name controller alive true
host host2 name compute1 alive true
host host3 name compute3 alive false
host host1 name compute2 alive false
```

- To list the interfaces on a certain host enter the command:

```
midonet> host host0 list interface
iface midonet host_id host0 status 0 addresses [] mac 12:6e:b7:d0:4f:f1
mtu 1500 type Virtual endpoint DATAPATH
```


In the above command example:

- tzone0 = the tunnel zone you want to add the member (host) to
- host0 = the alias of the host you want to add
- 10.1.2.200 = the IP address of the host you want to add

Removing a host from a tunnel zone

Use this procedure to remove a host from a tunnel zone.

1. Enter the `list tunnel-zone` command to list the tunnel zone. For example:

```
midonet> list tunnel-zone
tzone tzone0 name default_tz type gre
```

2. Enter the `tunnel-zone tunnel-zone list member` command to list the tunnel zone members (hosts). For example:

```
midonet> tunnel-zone tzone0 list member
zone tzone0 host host0 address 172.19.0.2
```

3. Enter the `tunnel-zone tunnel-zone member host host show` command to show information about a specific host. For example:

```
midonet> tunnel-zone tzone0 member host host0 show
tunnel-zone-host zone tzone0 host host0 address 172.19.0.2
```

4. Enter the `tunnel-zone tunnel-zone member host host delete` command to delete the desired host (identified by the host's alias). For example:

```
midonet> tunnel-zone tzone0 member host host0 delete
```

5. (Optional) You can add the host back to the tunnel zone, using the `tunnel-zone tunnel-zone member add host host address ip-address` command as shown below:

```
midonet> tunnel-zone tzone0 member add host host0 address 172.19.0.2
zone tzone0 host host0 address 172.19.0.2
```

Removing a host

Use this procedure to remove inactive hosts.

1. Enter the command to list the hosts:

```
midonet> list host
host host0 name precise64 alive true
```

2. Enter the command to delete the desired host (identified by its alias):

```
midonet> host host0 delete
```


AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

2. Enter the command to list the bridges on the current tenant:

```
midonet> list bridge
bridge bridge0 name External state up
bridge bridge1 name Management state up
bridge bridge2 name Internal state up
```

3. Enter the command to list the ports on the desired bridge:

```
midonet> bridge bridge0 list port
port port0 device bridge0 state up
port port1 device bridge0 state up
port port2 device bridge0 state up
```

4. Enter the command to list the interfaces on a certain host:

```
midonet> host host0 list interface
iface lo host_id host0 status 3 addresses [u'127.0.0.1',
u'0:0:0:0:0:0:1'] mac 00:00:00:00:00:00 mtu 65536 type Virtual
endpoint LOCALHOST
iface midonet host_id host0 status 0 addresses [] mac 8e:4d:60:c1:70:d7
mtu 1500 type Virtual endpoint DATAPATH
iface eth1 host_id host0 status 3 addresses
[u'fe80:0:0:0:250:56ff:fe93:7c35'] mac 00:50:56:93:7c:35 mtu 1500 type
Physical endpoint PHYSICAL
iface eth0 host_id host0 status 3 addresses [u'10.1.2.200',
u'fe80:0:0:0:250:56ff:fe93:c9a4'] mac 00:50:56:93:c9:a4 mtu 1500 type
Physical endpoint PHYSICAL
```

5. Enter the command to bind a certain host to virtual port:

```
midonet> host host0 add binding
host interface port
```

6. Enter the command to bind a virtual port on the bridge with a physical interface on the host:

```
midonet> host host0 add binding port bridge0:port0 interface eth1
host host0 interface eth1 port bridge0:port0
```

Stateful port groups

MidoNet features stateful port groups, which are groups of virtual ports (typically two) that are logically associated, usually to perform load balancing or for link redundancy.

For such ports MidoNet keeps state local to the two endpoints of a connection. In most cases, connections that traverse MidoNet do so between a single pair of ports. Typical cases include a router with two uplink BGP ports, or an L2GW with two ports connected to a physical L2 network. In both cases, the pair of ports becomes a set of ports because packets may return through different paths. Those port pairs will share state.

You configure stateful port groups in the MidoNet CLI, using the port-group command.

Creating stateful port groups

Follow the steps of this procedure to create a stateful group of ports, using the MidoNet CLI.

Before you launch the MidoNet CLI you need to find out the OpenStack UUID of the tenant on which you want to create your port group. To this end, you can use keystone. Issue the following commands in the terminal on the MidoNet host:

```
# keystone tenant-list
```

id	name	enabled
7a4937fa604a425e867f085427cc351e	admin	True
037b382a5706483a822d0f7b3b2a9555	alt_demo	True
0a1bf57198074c779894776a9d002146	demo	True
28c40ac757e746f08747cddb32a83c40b	services	True

The output of the command shows the full list of tenants. For this procedure we will use the 'admin' tenant, 7a4937fa604a425e867f085427cc351e.

1. In the MidoNet CLI determine the list of available routers.

```
midonet> list router
router router0 name Edge Router state up
router router1 name TenantRouter state up
```

Let's assume that the router whose ports you are going to add to the port group is Edge Router, router0.

2. Now list the ports on router0.

```
midonet> router router0 list port
port port0 device router0 state up mac 02:c2:0f:b0:f2:68 addresses 100.100.100.1/30
port port1 device router0 state up mac 02:cb:3d:85:89:2a addresses 172.168.0.1/16
port port2 device router0 state up mac 02:46:87:89:49:41 addresses 200.200.200.1/24 peer bridge0:port0
port port3 device router0 state up mac 02:6b:9f:0d:c4:a8 addresses 169.254.255.1/30
```

You want to add port0 and port1 on the router to load balance the BGP traffic on the Edge Router.

3. Load your tenant using the 'set' command.

```
midonet-cli> sett 7a4937fa604a425e867f085427cc351e
tenant_id: 7a4937fa604a425e867f085427cc351e
```

4. Create a stateful port group using the 'port-group create' command.

```
midonet-cli> port-group create name SPG stateful true
pgroup0
```

5. Add the two ports on the Edge Router that you want to participate in load balancing, to the port group you just created.

```
midonet> port-group pgroup0 add member port router0:port0
port-group pgroup0 port router0:port0
midonet> port-group pgroup0 add member port router0:port1
port-group pgroup0 port router0:port1
```

You have successfully added both router ports to the stateful port group, which you can verify by issuing the following command:

```
midonet> port-group pgroup0 list member
port-group pgroup0 port router0:port1
port-group pgroup0 port router0:port0
```


Connecting two routers

You can easily connect two virtual routers via virtual ports on each router.

Make sure you create the router ports on the two routers and assign the ports to the same subnet. See [Chapter 4, “Device abstractions” \[17\]](#) for information about creating routers and adding router ports.

To connect two routers:

1. Enter the command to list the routers on the current tenant:

```
midonet> list router
router router3 name test-router2 state up
router router1 name test-router state up
```

2. Enter the command to list the ports on one of the routers you want to connect:

```
midonet> router router1 list port
port port0 device router1 state up mac 02:a6:81:08:ab:5d addresses 10.
100.1.1/24 peer bridge1:port0
port port1 device router1 state up mac 02:fa:5f:87:bb:d2 addresses 10.
100.1.2/24
```

3. Enter the command to list the ports on the router you want to connect it to:

```
midonet> router router3 list port
port port0 device router3 state up mac 02:df:24:5b:19:9b addresses 10.
100.1.128/24
```

4. Enter the command to bind the port on one router (for example, port 1 on router1) to the port on another router (for example, port0 on router3):

```
midonet> router router1 port port1 set peer router3:port0
```

5. Enter the command to list the ports on one of the routers:

```
midonet> router router1 list port
port port0 device router1 state up mac 02:a6:81:08:ab:5d addresses 10.
100.1.1/24 peer bridge1:port0
port port1 device router1 state up mac 02:fa:5f:87:bb:d2 addresses 10.
100.1.2/24 peer router3:port0
```

The above output shows that port1 on router1 is connected to port0 on router3.

Source

Type

There are three types: Normal, Blackhole, and Reject.

Destination

Next Hop Gateway

The route has three options regarding what to do with a packet:



Note

Prior to adding the above route, a port was added to router2 using the following command:

```
midonet> router router2 add port addresses 169.254.255.3/30
router2:port2
```

2. Enter the command to list the routes on router2 to confirm the added route(s):

```
midonet> router router2 list route
route route0 type normal src 0.0.0.0/0 dst 10.100.1.1 port router2:port0
weight 0
route route1 type normal src 0.0.0.0/0 dst 10.100.1.2 port router2:port1
weight 0
route route2 type normal src 0.0.0.0/0 dst 169.254.255.3 port
router2:port2 weight 0
route route3 type normal src 0.0.0.0/0 dst 169.254.255.0/30 port
router2:port2 weight 0
```

Deleting routes

If you are using MidoNet as a standalone SDN controller, there are many situations where you might want to delete routes; all related to managing your physical network devices.

For example, if you want to reverse something you did that required manually adding routes, you can delete the routes.



Warning

It is not recommended to delete routes that were added automatically as a result of OpenStack Neutron operations.

To delete a route:

1. Enter the command to list the routes on a certain router:

```
midonet> router router2 list route
route route0 type normal src 0.0.0.0/0 dst 10.100.1.1 port router2:port0
weight 0
route route1 type normal src 0.0.0.0/0 dst 10.100.1.2 port router2:port1
weight 0
route route2 type normal src 0.0.0.0/0 dst 169.254.255.3 port
router2:port2 weight 0
route route3 type normal src 0.0.0.0/0 dst 169.254.255.0/30 port
router2:port2 weight 0
```

The above command lists the routes on router2.

2. Enter the command(s) to delete the desired route(s) from the desired router:

```
midonet> router router2 delete route route2
midonet> router router2 delete route route3
```

The above commands delete route2 and route3 from router2.

3. Enter the command to list the routes on the router to confirm the deletions:

```
midonet> router router2 list route
```

```
route route0 type normal src 0.0.0.0/0 dst 10.100.1.1 port router2:port0
weight 0
route route1 type normal src 0.0.0.0/0 dst 10.100.1.2 port router2:port1
weight 0
```

ECMP Limitations

Limitations of MidoNet ECMP for load-balancing to a set of stateful services:

- ECMP route selection is "flow-hash mod N"
- MidoNet does not test that the routes are ordered the same in every Agent so that "hash mod N" will return the same result on different L3 Gateways as long as N does not change.
- The hash algorithm does not satisfy "consistent hashing" properties, in particular, it does not minimise the impact of changes to the load-balanced set. Therefore, when the load-balanced set changes, any active flow that is recomputed (because it lasts more than 30 seconds or because upstream routing causes it to ingress a new gateway) will be balanced to a different back-end, usually leading to a broken connection.
 - Note that MidoNet's DNAT rules as well as L4 LBaaS implementation do generate per-flow state so those features are more resilient to backend set changes.
- The hash algorithm is not symmetric, so return packets are unlikely to traverse the same VNF instance.
 - This matters for stateful services that need to see traffic in both directions. Some stateful services can live without. And some stateful service instances forward traffic with their own address and therefore receive the return traffic without ECMP configuration.
- MidoNet's L3 Gateway makes use of ECMP for choosing an egress gateway for packets heading out of the logical network (to the rest of the data center or to the Internet). That use-case works fine because it is load-balancing to a set of stateless upstream routers. There is no harm to suddenly sending packets in a flow out of a different gateway node (apart from the possibility of some packets arriving out of order at the destination endpoint).

Despite the limitations discussed above, MidoNet's ECMP works pretty well for a number of use-cases including load-balancing a VIP to an Active-Active HA set of L7 load-balancers.

Remember that you can program MidoNet's virtual router's tables using either MidoNet's native API or via the Neutron ExtraRoute API extension.

Rule types

DNAT, SNAT

Every rule has a single Condition object that a packet must match in order for the rule to be applied.

Taking a jump rule as an example, if a packet matches the jump's Condition object, then rule processing for that packet will continue in the jump's target chain; if the packet doesn't match, then processing continues with the rule following the jump in the jump's own rule-chain.

Condition objects specify a set or combination of attributes. Attributes are simple statements about the contents of a packet's headers. Examples of attributes are:

- 'the packet's TCP/UDP port number is between 500 and 1000'
- 'the packet's source IP address is in 10.0.0.0/16'



Conditions are checked against the packet in the state the packet is in when it reaches a rule. In other words, for example, if a previous rule modified the packet's port number, then the current rule's condition will be checked against the modified, not the original, port number.

In order to form a Condition, you specify a number of attributes (optionally, you can invert most attributes using the CLI). Enter an exclamation point (!) or "bang" symbol to invert it, as shown in the "CLI Rule Chain Attributes That Match Packets" table. For example, if you invert the src attribute, this matches packets whose source does not match the specified IP address or network.

Below is the list of Condition attributes (attributes, invert-flags, and arguments) and their descriptions.



The ports identified in rules are virtual ports on virtual routers or bridges. A virtual port may be bound to a specific Ethernet interface (like a tap) on a physical host OR it may be peered with another virtual port (in which case it connects two virtual devices). In either case, think of the virtual port as virtual because the rules only exist in the virtual topology AND nothing is known during rule evaluation about possible bindings of the virtual port to physical Ethernet interfaces.

Table 7.1. CLI Rule Chain Attributes

Attributes	Description
pos <INTEGER>:	The rule's position in the chain
type <TYPE>:	The rule <TYPE>; this is mostly used to distinguish between regular filtering rules and different types of NAT rules. The recognized <TYPE> values are: accept, continue, drop, jump, reject, return, dnat, snat, rev_dnat, rev_snat.
action accept	continue
return:	The rule action, meaningful for NAT rules only.
jump-to <CHAIN>:	The chain to jump to (if this is a jump rule).
target <IP_ADDRESS[-IP_ADDRESS][:INTEGER[-INTEGER]]>:	The NAT target, if this is a dnat or snat rule. At least one IP address must be given, optionally the NAT target may

Attributes	Description
	also contain a second address to form an address range and L4 port number or range of ports.

Table 7.2. CLI Rule Chain Attributes That Match Packets

Attributes That Match Packets	Description
hw-src [!]<MAC_ADDRESS>:	The source hardware address
hw-dst [!]<MAC_ADDRESS>:	The destination hardware address
ethertype [!]<STRING>:	Sets the data link layer (EtherType) of packets matched by this rule.
in-ports [!]<PORT[,PORT...]>:	Matches the virtual port through which the packet ingresses the virtual device that is currently processing the packet.
out-ports [!]<PORT[,PORT...]>:	Matches the port through which the packet egresses the virtual device that is currently processing the packet.
tos [!]<INTEGER>:	The value of the packet's Type of Service (TOS) field to match. Use this field to match the differentiated services value. See TOS for information.
proto [!]<INTEGER>:	The IP protocol number to match. See Protocol Numbers for information. Examples: ICMP = 1, IGMP = 2, TCP = 6, UDP = 17
src [!]<CIDR>:	The source IP address or CIDR block
dst [!]<CIDR>:	The destination IP address or CIDR block
src-port [!]<INTEGER[-INTEGER]>:	The TCP or UDP source port or port range
dst-port [!]<INTEGER[-INTEGER]>:	The TCP or UDP destination port or port range
flow <fwd-flow return-flow>:	Matches the connection-tracking status of the packet. If the packet is starting a new connection, fwd-flow will match. Alternatively, if the packet belongs to a connection already known to MidoNet, return-flow will match.
port-group [!]<PORT_GROUP>:	Matches a port group. Port groups allow the grouping of virtual ports to ease the creation of chain rules. See the CLI commands help for information.
ip-address-group-src [!]<IP_ADDRESS_GROUP>:	Matches a source IP address group. IP address groups allow the grouping of IP addresses to ease the creation of chain rules. See the CLI commands help for information.
ip-address-group-dst [!]<IP_ADDRESS_GROUP>:	Matches a destination IP address group. IP address groups allow the grouping of IP addresses to ease the creation of chain rules. See the CLI commands help for information.
hw-src-mask	<p>Source MAC address mask - A 48-bit bitmask in the form xxxx.xxxx.xxxx, where x is any hexadecimal digit. Specifies which bits are to be considered when applying the rule's hw-src test.</p> <p>Default value = ffff.ffff.ffff: All bits are considered when applying the hw-src test, so a packet's source MAC address must match hw-src exactly.</p> <p>ffff.0000.0000: Only the first sixteen bits are considered when applying the hw-src test, the first sixteen bits of a packet's source MAC address must match the first sixteen bits of hw-src.</p> <p>0000.0000.0000: No bits are considered when applying the hw-src test, so any packet will match.</p>
hw-dst-mask	<p>Destination MAC address mask - A 48-bit bitmask in the form xxxx.xxxx.xxxx, where x is any hexadecimal digit. Specifies which bits are to be considered when applying the rule's hw-dst test.</p> <p>Default value = ffff.ffff.ffff: All bits are considered when applying the hw-dst test, so a packet's destination MAC address must match hw-dst exactly.</p>

Attributes That Match Packets	Description
	ffff.0000.0000: Only the first sixteen bits are considered when applying the hw-dst test, the first sixteen bits of a packet's destination MAC address must match the first sixteen bits of hw-dst. 0000.0000.0000: No bits are considered when applying the hw-dst test, so any packet will match.
fragment-policy header nonheader any unfragmented	fragment-policy - Specifies the fragment type to match. ANY: Matches any packet. HEADER: Matches any packet that has a full header, that is, a header fragment or unfragmented packet. NONHEADER: Matches only nonheader fragments. UNFRAGMENTED: Matches only unfragmented packets. In general, ANY is the default policy. However, if a rule has a value for the src or dst field, the NONHEADER and ANY policies are disallowed and the default is HEADER. Furthermore, if the rule's type is dnat or snat and its target is not a single IP address with no ports specified, then the policy will default to UNFRAGMENTED, which is the only policy permitted for such rules. Unlike other rule properties, fragment-policy may not be inverted.

Example condition 1

Only packets whose network source is in 10.0.0.0/16 are allowed through to network 10.0.5.0/24. You can accomplish this a few different ways.

One way is to construct a DROP or REJECT rule that has a Condition and an ACCEPT rule with these attributes specify:

1. DROP when (ethertype equal 2048) AND (src NOT equal (10.0.0.0, 16))
2. ACCEPT when (dst equal (10.0.5.0, 24))
3. DROP



Note

The unconditional drop is needed to make rule 2 meaningful.

To create a rule chain with the above attributes:

1. If necessary, use the sett command or some other means to access the desired tenant.

```
midonet> sett 10a83af63f9342118433c3a43a329528
tenant_id: 10a83af63f9342118433c3a43a329528
```

2. Enter the command to create a new rule chain and assign it a name:

```
midonet> chain create name "drop_not_src_mynetwork_INBOUND"
chain5
```

3. Enter the command to drop IPv4 traffic that does not have the source 10.0.0.0/16:

```
midonet> chain chain5 add rule ethertype 2048 src !10.0.0.0/16 type drop
chain5:rule0
```

4. Enter the command to accept IPv4 traffic with the destination 10.0.5.0/24:

```
midonet> chain chain5 add rule ethertype 2048 dst 10.0.5.0/24 pos 2 type
accept
chain5:rule2
```

5. Enter the command to list the rules added to the new rule chain:

```
midonet> chain chain5 list rule
rule rule3 ethertype 2048 src !10.0.0.0/16 proto 0 tos 0 pos 1 type drop
rule rule2 ethertype 2048 dst 10.0.5.0/24 proto 0 tos 0 pos 2 type
accept
```

Example condition 2

Same as Example Condition 1, except here assume that you're structuring your rules differently. You want to have one DROP rule at the end of the chain that matches all packets; earlier in the chain you place ACCEPT rules that match packets/flows that are specifically allowed through.

The ACCEPT rule for the traffic allowed by Example Condition 1 would have a Condition with these attributes:

In the rule language, the chain would have:

ACCEPT when src=(10.0.0.0, 16) OR dst=(10.0.5.0, 24)

Rule at the end:

DROP all other packets

To create a rule chain with the above attributes:

1. If necessary, use the sett command or some other means to access the desired tenant.

```
midonet> sett 10a83af63f9342118433c3a43a329528
tenant_id: 10a83af63f9342118433c3a43a329528
```

2. Enter the command to create a new rule chain and assign it a name:

```
midonet> chain create name "accept_src_dst_mynetwork_INBOUND"
chain11
```

3. Enter the command to accept IPv4 traffic from the source 10.0.0.0/16:

```
midonet> chain chain11 add rule ethertype 2048 src 10.0.0.0/16 type
accept
chain11:rule0
```

4. Enter the command to accept IPv4 traffic with the destination 10.0.5.0/24:

```
midonet> chain chain11 add rule ethertype 2048 dst 10.0.5.0/24 type
accept
chain11:rule1
```

5. Enter the command to drop all IPv4 traffic (that didn't match the attributes in the preceding rules):

```
midonet> chain chain11 add rule ethertype 2048 pos 3 type drop
chain11:rule2
```

6. Enter the command to list the rules added to the new rule chain:



Note

Ports with infilter (pre-routing) and outfilter (post-routing) chains are connected to VMs. port1 is connected to a VM.

Listing the Rules for the Inbound Chain on a Port

To list the pre-routing rule chain for port1, enter the command:

```
midonet> chain chain2 list rule
rule rule0 ethertype 2048 src !172.16.3.3 proto 0 tos 0 pos 1 type drop
rule rule1 hw-src !fa:16:3e:fb:19:07 proto 0 tos 0 pos 2 type drop
rule rule2 proto 0 tos 0 flow return-flow pos 3 type accept
rule rule3 proto 0 tos 0 pos 4 type jump jump-to chain4
rule rule4 ethertype !2054 proto 0 tos 0 pos 5 type drop
```

The pre-routing rule chain contains the following instructions:

- rule0 says: for packets that match the ethertype 2048 (IPv4) that do not match the source IP address 172.16.3.3 (the private IP address of the VM), drop these packets. This prevents the port's VM from sending packets with a forged IP address.
- rule1 says: for packets with a hardware source that does not match the listed source MAC address (the VM's MAC address), drop these packets. This prevents the VM from sending packets with a forged MAC address.
- rule2 says: for packets that match a return flow (that is, a packet that belongs to a connection already known to MidoNet), accept these packets.
- rule3 says: for packets that were not dropped or accepted as a result of matching previous rules, allow these packets to jump to the indicated chain (chain4).
- rule4 says: for packets that do not match the ethertype 2054 (ARP packets), drop these packets.

Listing the OpenStack Security Group Rule Chain

You can list all the rule chains and then look at the rule chain for the OpenStack security group.

To list all the rule chains and examine specific rule chains, enter the command:

```
midonet> list chain
chain chain5 name OS_SG_050593ed-56ad-44ef-8489-4052d02d99ff_INGRESS
chain chain0 name OS_PRE_ROUTING_5a151b0b-dea7-4918-bd17-876c1f7f5c64
chain chain1 name OS_POST_ROUTING_5a151b0b-dea7-4918-bd17-876c1f7f5c64
chain chain6 name OS_SG_01fcelb8-c277-4a37-a8cc-86732eea186d_INGRESS
chain chain4 name OS_SG_050593ed-56ad-44ef-8489-4052d02d99ff_EGRESS
chain chain7 name OS_SG_01fcelb8-c277-4a37-a8cc-86732eea186d_EGRESS
chain chain2 name OS_PORT_6f72342b-4947-432f-8d01-0cf4e4b8d049_INBOUND
chain chain3 name OS_PORT_6f72342b-4947-432f-8d01-0cf4e4b8d049_OUTBOUND
```

Note chain5, identified as a chain for an OpenStack security group (OS_SG) for INGRESS traffic.

To look at rule chain5, enter the command:

```
midonet> chain chain5 list rule
rule rule0 ethertype 2048 src 0.0.0.0/0 proto 6 tos 0 dst-port 5900 pos 1
type accept
```

```
rule rule1 ethertype 2048 src 0.0.0.0/0 proto 1 tos 0 pos 2 type accept
rule rule2 ethertype 2048 src 0.0.0.0/0 proto 6 tos 0 dst-port 22 pos 3
  type accept
rule rule3 ethertype 2048 src 0.0.0.0/0 proto 6 tos 0 dst-port 80 pos 4
  type accept
```

The above output shows the rule chain used to implement the security group you configured in OpenStack. These rules contain the following instructions:

- All the rules match ethertype 2048 (IPv4) packets.
- All the rules match traffic from any source network (0.0.0.0/0).
- All the rules, except rule1, match packets of IP protocol 6 (TCP) and accept them. rule1 matches packets of the ICMP type and accepts them.
- In addition to the other matches already mentioned, the rules match and accept the packets according to the security group rules you defined in OpenStack, specifically, packets with the destinations:
 - TCP port 5900 (VNC)
 - TCP port 22 (SSH)
 - TCP port 80 (HTTP)

Setting and clearing Rule Chains

To set or change a rule chain on a device, use the `set` command:

```
midonet> bridge bridge0 port port1 set infiltrer chain8
```

To remove a rule chain from a device, use the `clear` command:

```
midonet> bridge bridge0 port port1 clear infiltrer
```


- Display information about the infilter (pre-routing) and outfilter (post-routing) chains configured on a tenant router.
- List the rules for the rule chains.

Assumptions

For the example below, assume the following network conditions exist:

- A tenant router named "tenant-router"
- A private network with a 172.16.3.0/24 network address
- A public network with a 198.51.100.0/24 network address
- A VM with a 172.16.3.3 private IP address and a 198.51.100.3 public (floating) IP address

Viewing a Pre-Routing rule

To list routers on the current tenant and the router rule-chain information on the router(s), enter the command:

```
midonet> list router
router router0 name tenant-router state up infilter chain0 outfilter chain1
```

As shown in the above output, chain0 is the router's pre-routing (infilter) rule chain and chain1 is the post-routing (outfilter) rule chain.

To list information about the router's pre-routing rule chain, enter the command:

```
midonet> chain chain0 list rule
rule rule0 dst 198.51.100.3 proto 0 tos 0 in-ports router0:port0 pos 1 type
  dn timer action accept target 172.16.3.3
rule rule1 dst 198.51.100.2 proto 0 tos 0 in-ports router0:port0 pos 2 type
  rev_sn timer action accept
```

rule0 of the pre-routing rule chain on the tenant router contains the following instructions:

- For packets with the destination of 198.51.100.3 (the floating IP address associated with the VM):
 - Perform a destination NAT (DNAT) translation to change the destination IP address from the VM's floating IP address (198.51.100.3) to the VM's private IP address (172.16.3.3).

Viewing a Post-Routing rule

To list the post-routing rule on the tenant router, enter the command:

```
midonet> chain chain1 list rule
rule rule0 src 172.16.3.3 proto 0 tos 0 out-ports router1:port0 pos 1 type
  snat action accept target 198.51.100.3
rule rule1 proto 0 tos 0 out-ports router1:port0 pos 2 type snat action
  accept target 198.51.100.2:1--1
```

rule0 of the post-routing rule on tenant-router contains the following instructions:

- You want to translate traffic with the VM's public IP address to the VM's private IP address.
- You want to perform a reverse destination address translation.

To create the rule chain for the above DNAT configuration:

1. Create a new rule chain:

```
midonet> chain create name "dnat-test"
chain10
```

2. Create a rule that accepts traffic on a router port with the destination 198.51.100.4 and translates the destination to 10.100.1.150:

```
midonet> chain chain10 add rule dst 198.51.100.4 in-ports router1:port0
pos 1 type dnat action accept target 10.100.1.150
chain10:rule2
```

3. Create a rule that accepts traffic with the destination of the router's gateway to the public network and performs a reverse address translation from the public network address to the private network address.

```
midonet> chain chain10 add rule dst 198.51.100.2 in-ports router0:port0
pos 2 type rev_snat action accept
chain10:rule3
```

4. List the rules to check them:

```
midonet> chain chain10 list rule
rule rule2 dst 198.51.100.4 proto 0 tos 0 in-ports router1:port0 pos 1
type dnat action accept target 10.100.1.150
rule rule3 dst 198.51.100.2 proto 0 tos 0 in-ports router0:port0 pos 2
type rev_snat action accept
```

SNAT REV_SNAT example

This example shows how to use MidoNet CLI commands to configure source NAT (known as SNAT) and reverse source nat (we call that REV_SNAT) on a router using static NAT (one private ip maps to one public ip). Please be aware that reverse source nat and DNAT are not the same thing (it depends on the connection and the associated flow state).

You can configure SNAT for traffic of any network device (like a VM or a container) that enters the virtual topology. This means you can also use a MidoNet gateway to do static source NAT for physical servers sitting in a physical network connected to a MidoNet gateway acting as a L2 gateway or L3 router.

For demonstration purposes and to keep things simple we will use a veth pair and a network namespace in a simple installation with one midonet gateway acting as a static L3 router.

Note that the router does not have to have its virtual ports configured with ip addresses that we are using for NAT. Also you do not need to modify the routing table in the router that is actually responsible for doing the NAT.

This also holds true for L4 load balancing (which is also implemented in MidoNet using NAT): the ip of the load balancer does not have to be an ip on a port of a router.

For this example to work, your physical network must be configured to route traffic for 200.200.200.0/24 to the ip address of your edge router.

The following assumptions are made for this example code: * the host for setting up the veth pair is called 'os004' * the edge router is simply called 'edge' * for brevity we have only one gw with static routing (you have to configure this yourself) in this example. * the underlay is 192.168.7.0/24 * the overlay is 192.168.11.0/24 * we want to make 192.168.11.0/24 appear as SNAT traffic coming from 200.200.200.0/24 * the reverse source NAT path should convert traffic going to 200.200.200.0/24 back to 192.168.11.0/24

On a machine with the MidoNet agent running, we create a small lab environment with a veth-pair and a separate network namespace:

```
ip link | grep veth0-snat || ip link add veth0-snat type veth peer name veth1-snat
ip link set dev veth0-snat up
ip netns show | grep snat-example || ip netns add snat-example
ip link set veth1-snat netns snat-example
ip netns exec snat-example ip link set dev veth1-snat up
ip netns exec snat-example ip addr add 192.168.11.11/24 dev veth1-snat
ip netns exec snat-example ip route add default via 192.168.11.1
```

We also create a new bridge, create ports on the bridge and set up the edge router port.

```
midonet-cli -e 'bridge list name snat-example-bridge' || midonet-cli -e 'bridge create name snat-example-bridge'
export BRIDGE="$(midonet-cli -e 'bridge list name snat-example-bridge' | awk '{print $2;}')"

midonet-cli -e "bridge $BRIDGE port list" | wc -l | grep '2' || midonet-cli -e "bridge $BRIDGE port create"
midonet-cli -e "bridge $BRIDGE port list" | wc -l | grep '2' || midonet-cli -e "bridge $BRIDGE port create"

export VETH_PORT="$(midonet-cli -e "bridge $BRIDGE port list" | head -n1 | awk '{print $2;}')"
export ROUTER_PORT="$(midonet-cli -e "bridge $BRIDGE port list" | tail -n1 | awk '{print $2;}')"

export HOST="$(midonet-cli -e "host list name os004" | awk '{print $2;}')"
midonet-cli -e "host $HOST add binding port $VETH_PORT interface veth0-snat"

export ROUTER="$(midonet-cli -e "router list name edge" | awk '{print $2;}')"
midonet-cli -e "router $ROUTER port list" | grep 'address 192.168.11.1' | grep 'net 192.168.11.0/24' || \
midonet-cli -e "router $ROUTER port create address 192.168.11.1 net 192.168.11.0/24"

export RPORT="$(midonet-cli -e "router $ROUTER port list" | grep 'address 192.168.11.1' | grep 'net 192.168.11.0/24' | awk '{print $2;}')"
midonet-cli -e "router $ROUTER route list" | grep 'dst 192.168.11.0/24' || \
midonet-cli -e "router $ROUTER route add type normal src 0.0.0.0/0 dst 192.168.11.0/24 port $RPORT"

midonet-cli -e "port $RPORT set peer $ROUTER_PORT"
```

Now you should be able to ping the router from inside the network namespace.

```
root@os004:~# ip netns exec snat-example ping -c1 192.168.11.1
```

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

```
midonet-cli -e "router $ROUTER set infilter $CHAIN"
```

By setting the infilter on the router, this SNAT setup automatically becomes a firewall, this means that only this machine can now use the router and its traffic will always be SNATed when going over the router.

Here you can see the final results, tcpdump from the gateway (arp traffic omitted):

```
root@os002:~# clear; tcpdump -i enp0s8 -l -nnn -vvv -X -e ether host
ac:ab:ac:ab:ac:ab and icmp
tcpdump: listening on enp0s8, link-type EN10MB (Ethernet), capture size
262144 bytes
14:36:41.512662 ac:ab:ac:ab:ac:ab > 80:ee:73:83:93:56, ethertype IPv4
(0x0800), length 98: (tos 0x0, ttl 63, id 13058, offset 0, flags [DF],
proto ICMP (1), length 84)
    200.200.200.11 > 192.168.7.1: ICMP echo request, id 4927, seq 1, length
64
    0x0000:  4500 0054 3302 4000 3f01 b029 c8c8 c80b  E..T3.@.?..)...
    0x0010:  c0a8 0701 0800 138c 133f 0001 69b7 be58  .....?..i..X
    0x0020:  0000 0000 e350 0700 0000 0000 1011 1213  ....P.....
    0x0030:  1415 1617 1819 1a1b 1c1d 1e1f 2021 2223  .....!"#
    0x0040:  2425 2627 2829 2a2b 2c2d 2e2f 3031 3233  $%&'()*+,-./0123
    0x0050:  3435 3637                                     4567
14:36:41.513747 80:ee:73:83:93:56 > ac:ab:ac:ab:ac:ab, ethertype IPv4
(0x0800), length 98: (tos 0x0, ttl 64, id 0, offset 0, flags [DF], proto
ICMP (1), length 84)
    192.168.7.1 > 200.200.200.11: ICMP echo reply, id 4927, seq 1, length
64
    0x0000:  4500 0054 0000 4000 4001 e22b c0a8 0701  E..T...@...+....
    0x0010:  c8c8 c80b 0000 1b8c 133f 0001 69b7 be58  .....?..i..X
    0x0020:  0000 0000 e350 0700 0000 0000 1011 1213  ....P.....
    0x0030:  1415 1617 1819 1a1b 1c1d 1e1f 2021 2223  .....!"#
    0x0040:  2425 2627 2829 2a2b 2c2d 2e2f 3031 3233  $%&'()*+,-./0123
    0x0050:  3435 3637                                     4567
```

This is the traffic in the VXLAN tunnels going between the node with the veth pair and the MidoNet gateway:

```
root@os002:~# clear; tcpdump -i enp0s3 -l -nnn -vvv -X -e -Tvxlan port 6677
tcpdump: listening on enp0s3, link-type EN10MB (Ethernet), capture size
262144 bytes
14:37:26.709018 08:00:27:97:e4:37 > 08:00:27:c5:6b:60, ethertype IPv4
(0x0800), length 148: (tos 0x0, ttl 255, id 24228, offset 0, flags [none],
proto UDP (17), length 134)
    192.168.7.189.49719 > 192.168.7.190.6677: VXLAN, flags [I] (0x08), vni
8072578
ac:ab:ac:ab:ac:ab > 80:ee:73:83:93:56, ethertype IPv4 (0x0800), length 98:
(tos 0x0, ttl 63, id 17418, offset 0, flags [DF], proto ICMP (1), length
84)
    200.200.200.11 > 192.168.7.1: ICMP echo request, id 4930, seq 1, length
64
    0x0000:  4500 0086 5ea4 0000 ff11 cbf6 c0a8 07bd  E...^.....
    0x0010:  c0a8 07be c237 1a15 0072 0000 0800 0000  ....7...r.....
    0x0020:  7b2d 8200 80ee 7383 9356 acab acab acab  {-....s..V.....
    0x0030:  0800 4500 0054 440a 4000 3f01 9f21 c8c8  ..E..TD.@.?..!..
    0x0040:  c80b c0a8 0701 0800 a70e 1342 0001 96b7  .....B....
    0x0050:  be58 0000 0000 1fcb 0a00 0000 0000 1011  .X.....
    0x0060:  1213 1415 1617 1819 1a1b 1c1d 1e1f 2021  .....!
    0x0070:  2223 2425 2627 2829 2a2b 2c2d 2e2f 3031  "##$%&'()*+,-./01
    0x0080:  3233 3435 3637                                     234567
14:37:26.714107 08:00:27:c5:6b:60 > 08:00:27:97:e4:37, ethertype IPv4
(0x0800), length 148: (tos 0x0, ttl 255, id 45283, offset 0, flags [none],
proto UDP (17), length 134)
```

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -


```
midonet> list router
router router0 name Edge Router state up
router router1 name Tenant Router state up infiltrer chain0 outfilter chain1
```

As you can see, your Tenant Router where you are going to create a load balancer is router1.



Important

In MidoNet, routers have inbound and outbound filters. If the load balancer on a router balances traffic, these filters will be skipped. When using MidoNet with OpenStack, these filters usually only contain NAT rules that are irrelevant to load-balanced traffic, but this is worth taking into account if you are adding custom rules to the router's filters.

1. Create a load balancer and assign it to the Tenant Router.

```
midonet> load-balancer create
lb0
midonet> router router1 set load-balancer lb0
```

The load balancer assigned to the router will act on traffic flowing through that router.

2. Create a pool to which target back-end servers will be assigned.

```
midonet> load-balancer lb0 create pool lb-method ROUND_ROBIN
lb0:pool0
midonet> load-balancer lb0 pool pool0 show
pool pool0 load-balancer lb0 lb-method ROUND_ROBIN state up
```

3. Next, add target back-end servers to the pool you just created.

```
midonet> load-balancer lb0 pool pool0 create member address 192.168.100.
100 protocol-port 80
lb0:pool0:pm0
midonet> load-balancer lb0 pool pool0 create member address 192.168.100.
101 protocol-port 80
lb0:pool0:pm1
midonet> load-balancer lb0 pool pool0 member pm0 show
pm pm0 address 192.168.100.1 protocol-port 80 weight 0 state up
```

For each back-end server you must add its IP address and port to the pool.

4. Create a virtual IP address (VIP) and port, then assign it to the pool against which load balancing will be performed (lb0:pool0). Typically, a VIP is an IP address from the public IP space.

```
midonet> load-balancer lb0 pool pool0 list vip
midonet> load-balancer lb0 pool pool0 create vip address 203.0.113.2
persistence SOURCE_IP protocol-port 8080
lb0:pool0:vip0
midonet> load-balancer lb0 pool pool0 vip vip0 show
vip vip0 load-balancer lb0 address 203.0.113.2 protocol-port 8080
persistence SOURCE_IP state up
```



Note

Above example uses sticky source IP address persistence, read more about it in [the section called "Sticky Source IP" \[50\]](#).



Health monitor

HAProxy configuration

When using a Layer 4 load balancer, you can configure a health monitor to perform checks on the back-end servers.

Run the following command to check if a host is enabled for health monitoring:

- Create a health monitor object using the CLI or API server, and set the relevant delay, timeout, and max_retries values (see "Health Monitor" for information).
- Attach the health monitor object to the pool that you want to be monitored. A single health monitor can be attached to any number of pools, but pools may only have a single health monitor.
- Set the admin_state_up on the health monitor object to true.

CLI Example

The example below shows how to use the MidoNet CLI to configure health monitoring.

```
midonet> health-monitor list
midonet> health-monitor create type TCP delay 100 max-retries 50 timeout
500
hm0
midonet> load-balancer lb0 pool pool0 set health-monitor hm0
midonet> load-balancer lb0 pool pool0 health-monitor show
hm hm0 delay 100 timeout 500 max-retries 50 state down
midonet> health-monitor hm0 pool list
pool pool0 load-balancer lb0 health-monitor hm0 lb-method ROUND_ROBIN state
up
```

Disabling Health Monitoring

To disable health monitoring on a pool you can do perform one of the following:

- Set the admin_state_up on the health monitor to false. Note that all pools that are using this health monitor will have health monitoring disabled.
- Set the health_monitor_id on the pool to null.
- Delete the health monitor object.

10. Load Balancing as a Service (LBaaS)

Table of Contents

Neutron LBaaS Support	54
-----------------------------	----

Load Balancers are used to distribute traffic to a set of members based on a configured algorithm. It can be used for fault tolerance, high availability, performance improvements, as well as in many other situations. Load balancers can be configured in a cloud environment (Load Balancer as a Service, or LBaaS).

Neutron LBaaS Support

Neutron and MidoNet has supported LBaaS for some time, but the first implementation of LBaaS was suboptimal, especially for scaling and performance purposes. Thus, Neutron released a second version of the LBaaS standard, and MidoNet now also supports the newer LBaaS V2 model as well as the older, deprecated LBaaS V1 model.

There have been many changes from LBaaS V1 to V2, so this guide will both cover LBaaS V2 from a fresh-starter perspective, as well as offering guides on using V2 for those familiar with those familiar with V1.

Regarding LBaaS configuration, in order for load balancing to work in terms of a virtually defined network, there must be a route from the network on which the load balancer resides to each individual pool member. This also means no network address translation should be performed in between the load balancer and the pool members. Note that this is a general configuration requirement for load balancing, and not only limited to MidoNet.

Neutron LBaaS V1 Specification Support

MidoNet supports the Neutron LBaaS V1 specification (as documented in https://wiki.openstack.org/wiki/Neutron/LBaaS/API_1.0). To set up and use load balancers, users can set up and manipulate load balancers, pools, members, health monitors, etc. via the standard Neutron commands.

MidoNet load balancer limitations

Although MidoNet supports the Neutron LBaaS V1 API, not all Neutron LBaaS features are supported in MidoNet:

- L7 load balancing is not supported.
- There are no pool statistics.
- Only round robin is supported for the load balancer method.
- The Neutron provider model is not supported (for example, assigning a specific provider for each pool is not supported).
- Only one health monitor per pool (the first one in the list).
- Only TCP health check (no UDP, HTTP, also no ICMP/pings).

- Only Source IP session persistence (no cookie or URL).
- You cannot associate a floating IP to a virtual IP address (VIP).
- A VIP must not be on the same subnet as the pool member IF health monitoring is enabled.
- There is no connection limit.
- Only TCP load balancing is supported.
- When using sticky-source IP, directly connecting to the selected back-end host on the same port that is being used by the load balancer is not possible.
- To use an external network for VIPs, it has to be "shared".

Neutron LBaaS V2 Specification Support

MidoNet also supports the Neutron LBaaS V2 specification (as documented in https://wiki.openstack.org/wiki/Neutron/LBaaS/API_2.0). To set up and use load balancers, users can set up and manipulate load balancers, pools, members, health monitors, etc. via the standard Neutron commands as outlined in <http://docs.openstack.org/mitaka/networking-guide/config-lbaas.html>.

MidoNet load balancer limitations

Although MidoNet supports the Neutron LBaaS V2 API, not all Neutron LBaaS features are supported in MidoNet:

- L7 load balancing is not supported, which means several parameters for Layer 7 options, such as `sni-container-refs` and `default-tls-container-ref`, are not supported. Also, the protocol field for listeners, pools, and health monitor type can only be set as `TCP`.
- Only round robin is supported for the load balancer method.
- Only Source IP session persistence is supported.



A single L4 flow may generate up to two different flows: one to handle non-header fragments, another to handle all other packets.

Fragmented packets rule chain creation example

Create a chain (this creates a rule chain with an alias, "chain0" in the example, pointing to the chain created):

```
create chain name chain0
```

Add a rule to the chain that drops header fragments:

```
chain chain0 add rule fragment-policy header pos 2 header type drop
```

Example 1 Firewall, Does Not Account for Fragmented Packets

The example below only handles non-fragmented packets. These are the firewall rules you start with before you decide to handle fragmented packets.

Initially, you design your firewall to:

- Only allow incoming TCP port 80 (HTTP) traffic
- Drop all other packets

Without addressing fragmented packets, you create a rule chain with the following two rules:

- Rule at position 1
 - By default, this rule matches only non-fragmented packets and header fragments.
 - ACCEPTs packets with protocol=TCP and destination=80.

```
midonet> chain chain0 add rule ethertype 2048 src 0.0.0.0/0 proto 6 in-ports router2:port0 dst-port 80 pos 1 type accept
```

- Rule at position 2
 - DROPS all packets.

```
midonet> chain chain0 add rule ethertype 2048 src 0.0.0.0/0 dst 0.0.0.0/0 in-ports router2:port0 pos 2 type drop
```

```
midonet> chain chain0 list rule
rule rule0 ethertype 2048 src 0.0.0.0/0 proto 6 tos 0 dst-port 80 in-ports router2:port0 pos 1 type accept
rule rule1 ethertype 2048 src 0.0.0.0/0 dst 0.0.0.0/0 proto 0 tos 0 in-ports router2:port0 pos 2 type drop
```

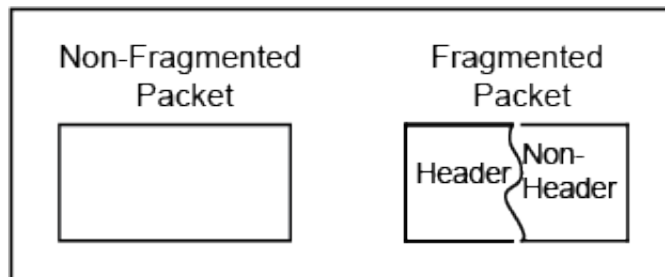
With the above rule chain, MidoNet handles fragmented packet with the destination TCP port 80 as follows:

- The first half of the packet, which contains the TCP header, reaches the rule at position 1, and is accepted.

- However, the second half of the fragmented, which does not have the destination port, reaches the rule at position 1, does not match the rule's condition, and is dropped. This means the fragmented packets do not reach the Web server.

Example 2 Firewall, Addresses Fragmented Packets

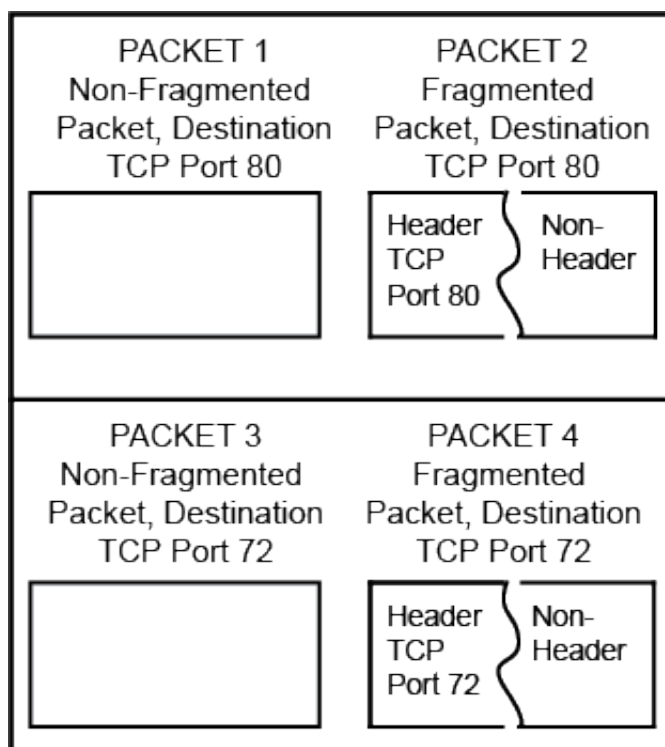
To address this problem, MidoNet provides a mechanism to handle the fragmented packets. This mechanism allows the fragmented packets to reach their destination, as shown in the following example. The drawing below simply depicts a whole, non-fragmented packet and a fragmented packet that consists of two parts, a header and non-header.



Non-Fragmented and Fragmented Packets

For this example, consider the following packets:

- Non-fragmented packet with the destination, TCP port 80
- Fragmented packet with the destination, TCP port 80
- Non-fragmented packet with the destination, TCP port 72
- Fragmented packet with the destination, TCP port 72



AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

Configuration

You can specify resource protection configuration parameters using `mn-conf(1)`, in the `agent.datapath` section. The available parameters are:

- `global_incoming_burst_capacity` - this sets the size of the top-level bucket and also defines the total number of tokens in the system (corresponding to in-flight packets) that will be divided among the different levels in the HTB; the rate at which tokens are placed back in the bucket is a function of the rate at which they are processed.
- `tunnel_incoming_burst_capacity` - this sets the capacity of the bucket associated with tunnel traffic, enforcing the rate at which a MidoNet Agent can communicate with the other Agents.
- `vm_incoming_burst_capacity` - this sets the capacity of each VM leaf bucket, which is below the shared VM bucket. This parameter enforces the rate at which individual VMs can send traffic.
- `vtep_incoming_burst_capacity` - this sets the capacity of the bucket associated with the VxLAN VTEP functionality, which enforces the rate at which the MidoNet Agent can communicate with the VxLAN domain.

See the `mn-conf(1)` schemas for more information about the above parameters.

Recommended values for these properties depend on the role of the MidoNet host (Gateway vs. Compute Node) and interaction with other resource-related properties, like JVM-memory and flow-table size. Midolman RPM and Debian packages include versions of each configuration tuned for Compute/Gateway hosts respectively. You can find these configurations in `/etc/midolman`, alongside the default configuration files. See "Recommended Values" for a table of recommended values.

Disabling Resource Protection

You can disable the resource protection feature.

To disable resource protection:

1. Specify a size value of "0" for all the parameters described in the Configuration section, except for the `global_incoming_burst_capacity` parameter.

This will cause all tokens to accumulate in the global bucket and all the traffic will be distributed from this single bucket.

14. MidoNet monitoring

Table of Contents

Metering	64
Monitoring Network State Database	65
Monitoring Midolman Agents	71
Monitoring events	74
Packet Tracing	83
Port mirroring	84

MidoNet is composed of various services; each service exposes a variety of metrics that can be fetched from typical monitoring services.

This chapter describes the main available metrics for each service.

Metering

Note: This feature is in **experimental** status.

Overview

The goal of metering is to provide packets and bytes traffic counters for arbitrary slices of the traffic that travels through MidoNet.

A meter is a counter of bytes and packets, associated with a name. In order to be incremented, the meter needs to have flows associated with it. MidoNet agents will automatically associate flows with certain meters, and users can create their own custom meters setting the `meterName` attribute in chain rules.

For example, all traffic going through bridge with uuid `FOO` in MidoNet be counted under meter `meters:device:FOO`. All traffic egressing port `BAR` will also be reflected in meter `meters:port:tx:BAR`.

MidoNet agents offer these counters for their partial view of overlay traffic. In other words, each agent will provide meters that only account for the traffic that agent has simulated. For a given meter, the MidoNet-wide real value is the sum of the value of this meter across all agents.

Note: Metering data is meant to be polled and stored by a monitoring layer onto a time series database. Thus agents don't persist the metering data they gather, and meter values will reset to zero when an agent reboots. **Any metering data collection layer should account for this effect and detect counter resets.**

Querying meters

Agents publish meters over JMX and a command line tool, `mm-meter`, uses their JMX interface to list, fetch and monitor meter values.

For example code on the JMX interface, the best source is [the code of mm-meter itself](#).

Querying meters with `mm-meter` is very simple:

```
$ mm-meter --help
```

```

-h, --host <arg>    Host (default = localhost)
-p, --port <arg>    JMX port (default = 7200)
--help              Show help message

Subcommand: list - list all active meters
--help             Show help message
Subcommand: get
-n, --meter-name <arg> name of the meter
--help             Show help message

trailing arguments:
delay (not required)  delay between updates, in seconds. If no delay is
                      specified, only one report is printed. (default = 0)
count (not required) number of updates, defaults to infinity
                      (default = 2147483647)

```

The `list` command will print a list of all meters known to this agent:

```

$ mm-meter list
meters:user:port0-on-the-bridge
meters:port:rx:cf453c9d-94c4-4c27-ba32-529b7cbacf1d
meters:device:845a54bf-b702-4dc2-8958-bbe7156bc4ef
meters:port:tx:cf453c9d-94c4-4c27-ba32-529b7cbacf1d
meters:port:tx:9182485b-8f86-462d-a8be-62586060eeb9
meters:device:f0d1f093-2de7-49a1-a5ec-898f94769e34
meters:device:9182485b-8f86-462d-a8be-62586060eeb9
meters:port:rx:9182485b-8f86-462d-a8be-62586060eeb9
meters:device:cf453c9d-94c4-4c27-ba32-529b7cbacf1d

```

And the `get` command will print the *current*, *local* counters for a meter. It takes a delay, in which case it will poll the meter and print deltas periodically:

```

$ mm-meter get -n meters:port:rx:cf453c9d-94c4-4c27-ba32-529b7cbacf1d 10
      packets      bytes
      568935      4215888475
           0           0
           0           0
          23         5834
           0           0

```

Creating a custom meter

Operators may want to meter a custom slice of their virtual network traffic. This is possible by matching on that slice using one or several chain rules in the virtual topology. The `meterName` property in chain rules will assign matching flows to the meter referred to by its value, in addition to the meters that flow would naturally feed.

Besides using the REST API, operators can use `midonet-cli` to set up such rules. The following rule will assign to meter `my-meter` all traffic that hits the rule after having traversed device `9182485b-8f86-462d-a8be-62586060eeb9`:

```

midonet> chain chain0 list rule
rule rule0 proto 0 tos 0 traversed-device 9182485b-8f86-462d-
a8be-62586060eeb9 fragment-policy any pos 1 type accept meter my-meter

```

Note that, when inspecting meters `my-meter` will turn into `meters:user:my-meter`, to avoid naming conflicts with built-in meters.

Monitoring Network State Database

The Network State Database is deployed with ZooKeeper and optionally, Cassandra, instances. Both offer JMX bindings.

The configuration provided with MidoNet uses only a subset of the most relevant for our use cases. Details in the sections below provide additional information about the metrics configured by MidoNet's deployment scripts, as well as an explanation about what to watch.

JXM Metrics

In addition to monitoring ZooKeeper and Cassandra directly, both the MidoNet Agent and the MidoNet Cluster expose a set of internal metrics that you can use to monitor the performance and health of the connection to the ZooKeeper Network State Database, as well as gaining visibility into various statistics such as the number of reads and writes.

These metrics are available using JMX in the `metrics` domain. The name of each JMX MBean is prefixed by `org.midonet.cluster.monitoring.metrics`. To access a particular metric, prepend the above prefix to the metric names listed below. For example, the full name of a JMX metric is:

```
metrics:name=org.midonet.cluster.monitoring.metrics.StorageGauge.  
connectionState
```

General Metrics

These metrics are gauges for the current ZooKeeper connection state and open observables to the topology objects.

Name	Description
<code>StorageGauge.connectionState</code>	The state of the connection to ZooKeeper used for normal reading and writing.
<code>StorageGauge.failFastConnectionState</code>	The state of the fail fast connection to ZooKeeper, which is used by some services to detect whether a certain node has lost connectivity to the NSDB with a low latency.
<code>StorageGauge.objectObservableCount</code>	The number of open observables for topology objects.
<code>StorageGauge.classObservableCount</code>	The number of open observables for topology classes.
<code>StorageGauge.objectObservableCountByClass</code>	The number of open observables for topology objects of a particular class.

Error Metrics

These metrics expose counters for errors that occurred while writing or reading the topology or state data in ZooKeeper.

Name	Description
<code>StorageCounter.concurrentModificationExceptions</code>	The number of exceptions thrown because a topology object was modified concurrently by two or more clients.
<code>StorageCounter.conflictExceptions</code>	The number of exceptions thrown because an attempt was made to modify a topology object that conflicted with the current references to/from that object.
<code>StorageCounter.objectReferenceExceptions</code>	The number of exceptions thrown because an attempt was made to delete a topology object that was still referenced by another object.
<code>StorageCounter.objectExistsExceptions</code>	The number of exceptions thrown because an attempt was made to create a topology object with an identifier that already exists.
<code>StorageCounter.objectNotFoundExceptions</code>	The number of exceptions thrown because an attempt was made to access (update or delete) a topology object that does not exist.
<code>StorageCounter.nodeExistsExceptions</code>	The number of exceptions thrown because an attempt was made to create a ZooKeeper node that already exists.
<code>StorageCounter.nodeNotFoundExceptions</code>	The number of exceptions thrown because an attempt was made to access (update or delete) a ZooKeeper node that does not exist.

Name	Description
<code>StorageCounter.objectObserved</code>	The number of occurrences where an attempt was made to read a topology object via an observable stream that was already closed. This measures the degree of concurrent access to the same objects and it is automatically recovered by opening a new observable.
<code>StorageCounter.classObserved</code>	The number of occurrences where an attempt was made to read all the objects in a topology class via an observable stream that was already closed. This measures the degree of concurrent access to the same objects and it is automatically recovered by opening a new observable.
<code>StorageCounter.stateObserved</code>	The number of occurrences where an attempt was made to read the state of a topology object via an observable stream that was already closed. This measures the degree of concurrent access to the same object state and it is automatically recovered by opening a new observable.
<code>StorageCounter.objectObservedErrors</code>	The number of unrecoverable errors that were emitted via the observable stream of a topology object.
<code>StorageCounter.classObservedErrors</code>	The number of unrecoverable errors that were emitted via the observable stream of a topology class.
<code>StorageCounter.stateObservedErrors</code>	The number of unrecoverable errors that were emitted via the observable stream of a topology object state.

Performance Metrics

These metrics expose information about the ZooKeeper connection changes and data operations.

Name	Description
<code>StorageMeter.connectionsCreated</code>	A meter for the number of created ZooKeeper connections.
<code>StorageMeter.connectionsLost</code>	A meter for the number of lost ZooKeeper connections.
<code>StorageTimer.read</code>	A timer for the number of ZooKeeper node reads.
<code>StorageTimer.readChildren</code>	A timer for the number of ZooKeeper node children reads.
<code>StorageTimer.write</code>	A timer for the number of ZooKeeper node writes.
<code>StorageTimer.multi</code>	A timer for the number of ZooKeeper transaction writes.
<code>StorageHistogram.stateTableLatency</code>	A histogram for the latency of reading the entire list of entries of a state table from ZooKeeper.
<code>StorageHistogram.stateTableLatencyDirect</code>	A histogram for the latency of adding a new entry to a state table by writing directly to ZooKeeper.
<code>StorageHistogram.stateTableLatencyProxy</code>	A histogram for the latency between adding a new entry to a state table and receiving a notification the entry was accepted by the NSDB. When the <i>state proxy</i> service is disabled, this latency measures the round-trip refresh of a state table via ZooKeeper. When the <i>state proxy</i> service is enabled, this latency measures the delay of writing the entry to ZooKeeper and reading the updated table via the state table proxy.

Session Metrics

These metrics expose statistics for the ZooKeeper connection.

Name	Description
<code>StorageHistogram.timeConnected</code>	A histogram for the duration of the ZooKeeper connection.
<code>StorageHistogram.timeDisconnected</code>	A histogram for the duration of the ZooKeeper connection loss.

Watcher Metrics

These metrics expose statistics for the watcher triggered by changes to the topology or state data from ZooKeeper.

Name	Description
<code>StorageCounter.nodeTriggered</code>	The number of triggered watchers for changes to ZooKeeper node data.

Name	Description
StorageCounter.childrenTr	The number of triggered watchers for changes to ZooKeeper node children.

Cassandra

By default, Cassandra uses port 7199 for JMX service from all its nodes and you can connect using jconsole for a comprehensive view.

Additionally, Cassandra's own nodetool utility offers commands like cfstats and tpstats that allow access to valuable stats into keyspace, tables, column families, and so on on a given node.

For a rich reference into Cassandra monitoring, visit the official documentation (go to <http://www.datastax.com/>, and search for "monitoring a Cassandra cluster").

Cache Reqs vs. Hits

This is self-descriptive, ideally you want the cache hits to be as close to the requests as possible. Note that by default MidoNet Cassandra nodes only enable the Partition Key Cache, but not Row Cache, so it's normal that these stay at 0. For MidoNet the Partition Key Cache should effectively be very similar to the Row Key Cache because our column families (CF) have only one column and therefore rows are not spread across several SSTables.

Compactions

This indicates the number of bytes being compacted. Typical workloads will present regular small spikes when the minor compaction jobs are run, and infrequent large spikes when major compactions are run. A large number of compactions indicates the need to add capacity to the cluster.

Internal Tasks

These are internal Cassandra tasks. The most important are:

- Gossip: MidoNet's Cassandra nodes are expected to spend a fair amount of their time busy in Gossip (wherein state information transfers among peers).
- MemTable Post Flusher: memtable flushes that are waiting to be written to the commit log. These should be as low as possible, and definitely not sustained.
- Hinted Handoff tasks: the appearance of these tasks indicates cases where replicas are detected as unavailable, so non-replica nodes need to temporarily store data until the replicas become available. Frequent Hinted Handoff spikes may hint at nodes being partitioned from the cluster.
- Anti-Entropy spikes: indicate data inconsistencies detected and being resolved.
- Stream activity: involves transferring or requesting data from other nodes. Ideally these should be infrequent and short-spaced.

Messaging Service Tasks

These are tasks received and responded to each of the peer nodes. Expect an even distribution with all peers.

NAT Column Family Latency

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

collection. This will correlate with high latency accessing Midolman's column families, which will propagate to Midolman. The Midolman agent will increase simulation latency and also degrade the utilization of CPU resources (experiencing more idle time while waiting for responses from Cassandra).

ZooKeeper

You can find ZooKeeper stats in the "zookeeper" category of the MidoStorage group. ZooKeeper classifies metrics in separate MBeans for leader/follower, so each node will report some values twice, one in the "Follower" role, another in the "Leader" role. Take into account that a given node may change roles (for example, if the leader shuts down, a follower node may be promoted to leader). You can easily spot these events. For example, the line in the "Connection Count as Follower" will suddenly blank out, and another will appear in the "Connection Count as Leader" graph.

Connection Count (as Follower/Leader)

These two graphs display the number of live connections to this node in its role at a given point in time.

In Memory Data Tree (as Follower/Leader)

Exposes the size of the in-memory znode database, both data nodes and watch count.

Latency (as Follower/Leader)

Exposes the average and maximum latency experienced in connections.

Packet Count (as Follower/Leader)

Exposes the count of packets sent/received by the node in its role at a given point in time.

Quorum Size

Exposes each node's view of the number of nodes agreeing on the leader's election.

ZooKeeper also exposes some information about each specific connection, this may be useful to watch for troubleshooting. Using jconsole (go to <http://www.oracle.com/technetwork/java/index.html> and search on "jconsole" for information), you can:

1. Connect to any ZooKeeper node at port 9199.
2. Navigate to org.apache.ZooKeeperService, ReplicatedServer_idX.
3. Choose the desired replica.
4. Go into Leader or Follower, Connections to see a list of IP addresses of the connected clients. Information shown here includes:
 - latency
 - packet sent/received count
 - session ID, etc. for that specific client.

Some of the MBeans whose values are exposed in our graphs also contain (computationally intensive) operations that also offer interesting information. Using jconsole,

expand `org.apache.ZooKeeperService`, and then the appropriate replica, and either Leader or Follower, according to its role:

- `InMemoryDataTree.approximateDataSize`: tells the size of the in-memory data store.
- `InMemoryDataTree.countEphemerals`: tells the count of ephemeral nodes.

The installation script also provides graphs to monitor the state of ZooKeeper's JVM:

- JVM GC times
- JVM Heap Summary
- JVM Non-Heap Summary

Descriptions of these graphs are beyond the scope of this document, but high JVM GC times are the best indication that ZooKeeper may be the source of problems in Midolman, which will be exhibited in high latency and under-utilization of CPU resources. ZooKeeper uses the Parallel collector, and the JVM GC times track the time of the last collection in all spaces from the `java.lang:type=GarbageCollector,name=PS Scavenge` MBean, which deals with all generations.

Monitoring Midolman Agents

MidoNet Agents expose a set of internal metrics that you can use to monitor the performance and health of agent nodes.

These metrics are available using JMX in the `metrics` domain. The name of each JMX MBean is prefixed by `org.midonet.midolman.monitoring.metrics`. To access a particular metric, prepend the above prefix to the metric names listed below. For example, the full name of a JMX metric is:

```
metrics:name=org.midonet.midolman.monitoring.metrics.DatapathMeter.flows.created
```

Additionally, some graphs are provided to monitor the state of the JVM running the MidoNet Agent.

Datapath Metrics

These metrics include meters for the created and deleted datapath flows, and errors that occur while changing these flows.

Name	Description
<code>DatapathMeter.flows.created</code>	A meter for the created datapath flows.
<code>DatapathMeter.flows.createdErrors</code>	A meter for the errors occurring during the creation of a datapath flow.
<code>DatapathMeter.flows.createdDuplicateErrors</code>	A meter for the errors occurring during the creation of a duplicate datapath flow.
<code>DatapathMeter.flows.deleted</code>	A meter for the deleted datapath flows.
<code>DatapathMeter.flows.deletedErrors</code>	A meter for the errors occurring during the deletion of a datapath flow.

Netlink Metrics

These metrics include meters for Netlink messages.

Name	Description
<code>NetlinkMeter.notifications</code>	A meter for the notifications received from the OVS datapath via Netlink.

Name	Description
NetlinkMeter.errors	A meter for the errors received from the OVS datapath via Netlink.
NetlinkMeter.htbDrops	A meter for the packets dropped because of the MidoNet Agent rate limiting mechanism.

Packet Pipeline Metrics

These metrics expose counters and meters for the simulation packet pipeline. The MidoNet Agent exposes certain metrics for each packet worker, where the worker index is included in the metric name. For example, when configuring MidoNet to use two (2) packet workers, the following distinct metrics are available: `PacketPipelineHistogram.worker-0.packetsProcessed` and `PacketPipelineHistogram.worker-1.packetsProcessed`.

Global Metrics

Name	Description
PacketPipelineCounter.contextsAllocated	The number of allocated packet context. A packet context is an internal variable that tracks the simulation progress for a packet that is being processed. Since each packet context requires a certain amount of memory, the counter reflects the memory consumption of the MidoNet agent when processing the user-space packets.
PacketPipelineCounter.contextsPooled	The number of packet contexts that have been pooled to be reused on subsequent packets.
PacketPipelineCounter.contextsSimulated	The number of packet contexts corresponding to packets that are currently simulated in the pipeline.

Packet Worker Metrics

Name	Description
PacketPipelineCounter.worker-{id}.packetsOnHold	The number of packets on hold, that are waiting for the agent to read the virtual topology data from the Network State Database.
PacketPipelineMeter.worker-{id}.packetsPostponed	A meter for the packets that were postponed because their corresponding virtual topology devices was not yet available at the agent.
PacketPipelineMeter.worker-{id}.packetsDropped	A meter for the packets dropped during simulation. These packets are dropped because the current virtual topology does not allow them to reach their intended destination.
PacketPipelineHistogram.worker-{id}.packetsProcessed	A histogram for the simulation latency of the processed packets. This value is fundamental for network latency.
PacketPipelineHistogram.worker-{id}.packetsExecuted	A histogram for the packet execution latency, which includes sending back a processed packet to the OVS datapath.
PacketPipelineMeter.worker-{id}.statePacketsProcessed	A meter for the processed packets that carried flow state data.
PacketPipelineMeter.worker-{id}.packetQueue.overflow	A meter for the packet queue overflows.
FlowTablesGauge.worker-{id}.currentDatapathFlows	A gauge for the current datapath flows.
FlowTablesMeter.worker-{id}.datapathFlowsCreated	A meter for created datapath flows.
FlowTablesMeter.worker-{id}.datapathFlowsRemoved	A meter for removed datapath flows.

Topology Metrics

The metrics measure the virtual topology cached by a given MidoNet Agent. The virtual topology consists of virtual devices, such as bridges, routers, ports, chains, rules, etc. When a device is requested in order to process a packet that must traverse or otherwise requires that device, the MidoNet Agent loads the device from the Network State Data-

base, and caches it for later usage. In addition, the agent keeps a notification stream open to the NSDB for that devices, such that future changes of the device are notified to the agent, updating its internal cache and invalidating any flows that have previously been established using that device's older configuration.

Some virtual topology metrics are global, whereas others are reported for each device class, which can be one of the following: Bridge, Chain, Host, IPAddrGroup, LoadBalancer, Mirror, Pool, PoolHealthMonitorMap, Port, PortGroup, Router, RuleLogger and TunnelZone.

Global Metrics

Name	Description
VirtualTopologyGauge.device	A gauge for the number of devices stored in the virtual topology cache.
VirtualTopologyGauge.observe	A gauge for the number of open observable streams that report changes for the devices stored in the virtual topology cache.
VirtualTopologyGauge.cache	A gauge for the cache hits when a device was requested by the packet pipeline when simulating a packet. A cache hit means that the device was available immediately, and the packet simulation was not interrupted.
VirtualTopologyGauge.cache	A gauge for the cache misses when a device was requested by the packet pipeline when simulating a packet. A cache miss means that the device needed to be loaded from the NSDB, requiring the packet simulation to be postponed until the device became available. Cache misses lead to higher packet processing latencies.
VirtualTopologyCounter.device	The number of updates the virtual topology caches receives for all cached devices.
VirtualTopologyCounter.device	The number of errors the virtual topology caches receives for all cached devices.
VirtualTopologyCounter.device	A counter for the cached devices that were deleted.
VirtualTopologyMeter.device	A meter for the device updates received by the virtual topology for the cached devices.
VirtualTopologyMeter.device	A meter for the device errors received by the virtual topology for the cached devices.
VirtualTopologyMeter.device	A meter for the devices deleted from the virtual topology.
VirtualTopologyHistogram.device	A histogram with the latency of loading a device from the NSDB.
VirtualTopologyHistogram.device	A histogram with the lifetime of a device in the virtual topology cache.

JVM Non-Heap Summary

Shows off-heap memory usage, which consists of mainly buffer pools used for messages to/from the Netlink layer.

JVM Heap Summary

Shows per-generation stats. The MidoNet Agent has very specific memory-usage constraints because it aims for a low memory and CPU footprint. At the same time, simulations generate a significant amount of short-lived garbage.

- The Eden is configured as the largest generation trying to hold as much garbage as possible. However, it is likely that it fills up frequently under high traffic, which may imply that some short-lived objects get promoted to the old generation and garbage is collected soon afterward.
- The Old Generation is expected to contain a baseline of long-lived objects that get reused during simulations. An amount of short-lived objects may also be pushed from the young generation, eventually also filling the old generation and triggering a GC event that will collect them. This will show as a see-saw pattern in the "Old used". The

Table 14.1. Configuration Files/Locations

Type of Node	Location of the Configuration File
MidoNet Network Agent	/etc/midolman/logback.xml
MidoNet Cluster server	/etc/midonet-cluster/logback.xml

Below are the behaviors with the default configuration shipped with the MidoNet release, but you can configure the behaviors as you like. See <http://logback.qos.ch/manual/index.html> for instructions on how to configure the logback.xml file.

Event log files locations

Event messages are stored locally on the filesystem in a separate file, in addition to the ordinary log file.

Table 14.2. Event Message Files/Locations

Type of Node	Location
MidoNet Network Agent	/var/log/midolman/midolman.event.log
MidoNet API server	/var/log/tomcat6/midonet-api.event.log (on Red Hat)
	/var/log/tomcat7/midonet-api.event.log (on Ubuntu)



Tip

In addition to midolman.event.log, /var/log/midolman/midolman.log contains additional debug information. You do not normally need to use it, but it may contain useful troubleshooting information.

Message format

By default, event messages have the following format.

```
<pattern>%d{yyyy.MM.dd HH:mm:ss.SSS} ${HOSTNAME} %-5level %logger - %m%n
%rEx </pattern>
```

See <http://logback.qos.ch/manual/layouts.html> for details about the above placeholders.

List of event messages

This section lists the event messages.

The event messages are organized in the following major categories:

- Virtual topology events
- API server events
- MidoNet Agent events

Virtual topology events

This section describes the messages associated with virtual topology events.

Router

Logger	org.midonet.event.topology.Router.CREATE
--------	--

Message	CREATE routerId={0}, data={1}.
Level	INFO
Explanation	Router with routerId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.Router.UPDATE
Message	UPDATE routerId={0}, data={1}.
Level	INFO
Explanation	Router with routerId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.Router.DELETE
Message	DELETE routerId={0}.
Level	INFO
Explanation	Router with routerId={0} was deleted.
Corrective Action	N/A

Logger	org.midonet.event.topology.Router.ROUTE_CREATE
Message	ROUTE_CREATE routerId={0}, data={1}.
Level	INFO
Explanation	Route={1} was created in routerId={0}.
Corrective Action	N/A

Logger	org.midonet.event.topology.Router.ROUTE_DELETE
Message	ROUTE_DELETE routerId={0}, routeId={1}.
Level	INFO
Explanation	routeId={1} was deleted in routerId={0}.
Corrective Action	N/A

Bridge

Logger	org.midonet.event.topology.Bridge.CREATE
Message	CREATE bridgeId={0}, data={1}.
Level	INFO
Explanation	Bridge with bridgeId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.Bridge.UPDATE
Message	UPDATE bridgeId={0}, data={1}.
Level	INFO
Explanation	Bridge with bridgeId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.Bridge.DELETE
Message	DELETE bridgeId={0}.
Level	INFO
Explanation	Bridge with bridgeId={0} was deleted.
Corrective Action	N/A

Port

Logger	org.midonet.event.topology.Port.CREATE
--------	--

Message	CREATE portId={0}, data={1}.
Level	INFO
Explanation	Port with portId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.Port.UPDATE
Message	UPDATE portId={0}, data={1}.
Level	INFO
Explanation	Port with portId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.Port.DELETE
Message	DELETE portId={0}.
Level	INFO
Explanation	Port with portId={0} was deleted.
Corrective Action	N/A

Logger	org.midonet.event.topology.Port.LINK
Message	LINK portId={0}, data={1}.
Level	INFO
Explanation	Port with portId={0} was linked.
Corrective Action	N/A

Logger	org.midonet.event.topology.Port.UNLINK
Message	UNLINK portId={0}, data={1}.
Level	INFO
Explanation	Port with portId={0} was unlinked.
Corrective Action	N/A

Logger	org.midonet.event.topology.Port.BIND
Message	BIND portId={0}, data={1}.
Level	INFO
Explanation	Port with portId={0} was bound.
Corrective Action	N/A

Logger	org.midonet.event.topology.Port.UNBIND
Message	UNBIND portId={0}.
Level	INFO
Explanation	Port with portId={0} was unbound.
Corrective Action	N/A

Chain

Logger	org.midonet.event.topology.Chain.CREATE
Message	CREATE chainId={0}, data={1}.
Level	INFO
Explanation	Chain with chainId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.Chain.DELETE
--------	---

Message	DELETE chainId={0}.
Level	INFO
Explanation	Chain with chainId={0} was deleted.
Corrective Action	N/A

Rule

Logger	org.midonet.event.topology.Rule.CREATE
Message	CREATE ruleId={0}, data={1}.
Level	INFO
Explanation	Rule with ruleId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.Rule.DELETE
Message	DELETE ruleId={0}.
Level	INFO
Explanation	Rule with ruleId={0} was deleted.
Corrective Action	N/A

Tunnel Zone

Logger	org.midonet.event.topology.TunnelZone.CREATE
Message	CREATE tunnelZoneId={0}, data={1}.
Level	INFO
Explanation	TunnelZone with tunnelZoneId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.TunnelZone.UPDATE
Message	UPDATE tunnelZoneId={0}, data={1}.
Level	INFO
Explanation	TunnelZone with tunnelZoneId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.TunnelZone.DELETE
Message	DELETE tunnelZoneId={0}.
Level	INFO
Explanation	TunnelZone with tunnelZoneId={0} was deleted.
Corrective Action	N/A

Logger	org.midonet.event.topology.TunnelZone.MEMBER_CREATE
Message	MEMBER_CREATE tunnelZoneId={0}, data={1}.
Level	INFO
Explanation	TunnelZone member={1} was added to tunnel-ZoneId={0}.
Corrective Action	N/A

Logger	org.midonet.event.topology.TunnelZone.MEMBER_DELETE
Message	MEMBER_DELETE tunnelZoneId={0}, data={1}.
Level	INFO
Explanation	TunnelZone member={1} was deleted from tunnel-ZoneId={0}.

BGP

Corrective Action	N/A
-------------------	-----

Logger	org.midonet.event.topology.Bgp.CREATE
Message	CREATE bgpId={0}, data={1}.
Level	INFO
Explanation	Bgp with bgpId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.Bgp.UPDATE
Message	UPDATE bgpId={0}, data={1}.
Level	INFO
Explanation	Bgp with bgpId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.Bgp.DELETE
Message	DELETE bgpId={0}.
Level	INFO
Explanation	Bgp with bgpId={0} was deleted.
Corrective Action	N/A

Logger	org.midonet.event.topology.Bgp.ROUTE_CREATE
Message	ROUTE_CREATE bgpId={0}, data={1}.
Level	INFO
Explanation	Route={1} was added to bgpId={0}.
Corrective Action	N/A

Logger	org.midonet.event.topology.Bgp.ROUTE_DELETE
Message	ROUTE_DELETE bgpId={0}, data={1}.
Level	INFO
Explanation	Route={1} was deleted from bgpId={0}.
Corrective Action	N/A

LoadBalancer

Logger	org.midonet.event.topology.LoadBalancer.CREATE
Message	CREATE loadBalancerId={0}, data={1}.
Level	INFO
Explanation	LoadBalancer with loadBalancerId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.LoadBalancer.UPDATE
Message	UPDATE loadBalancerId={0}, data={1}.
Level	INFO
Explanation	LoadBalancer with loadBalancerId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.LoadBalancer.DELETE
Message	DELETE loadBalancerId={0}.

VIP

Level	INFO
Explanation	LoadBalancer with loadBalancerId={0} was deleted.
Corrective Action	N/A

Logger	org.midonet.event.topology.VIP.CREATE
Message	CREATE vipId={0}, data={1}.
Level	INFO
Explanation	VIP with vipId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.VIP.UPDATE
Message	UPDATE vipId={0}, data={1}.
Level	INFO
Explanation	VIP with vipId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.VIP.DELETE
Message	DELETE vipId={0}.
Level	INFO
Explanation	VIP with vipId={0} was deleted.
Corrective Action	N/A

Pool

Logger	org.midonet.event.topology.Pool.CREATE
Message	CREATE poolId={0}, data={1}.
Level	INFO
Explanation	Pool with poolId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.Pool.UPDATE
Message	UPDATE poolId={0}, data={1}.
Level	INFO
Explanation	Pool with poolId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.Pool.DELETE
Message	DELETE poolId={0}.
Level	INFO
Explanation	Pool with poolId={0} was deleted.
Corrective Action	N/A

PoolMember

Logger	org.midonet.event.topology.PoolMember.CREATE
Message	CREATE poolMemberId={0}, data={1}.
Level	INFO
Explanation	PoolMember with poolMemberId={0} was created.

Corrective Action	N/A
-------------------	-----

Logger	org.midonet.event.topology.PoolMember.UPDATE
Message	UPDATE poolMemberId={0}, data={1}.
Level	INFO
Explanation	PoolMember with poolMemberId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.PoolMember.DELETE
Message	DELETE poolMemberId={0}.
Level	INFO
Explanation	PoolMember with poolMemberId={0} was deleted.
Corrective Action	N/A

HealthMonitor

Logger	org.midonet.event.topology.HealthMonitor.CREATE
Message	CREATE healthMonitorId={0}, data={1}.
Level	INFO
Explanation	HealthMonitor with healthMonitorId={0} was created.
Corrective Action	N/A

Logger	org.midonet.event.topology.HealthMonitor.UPDATE
Message	UPDATE healthMonitorId={0}, data={1}.
Level	INFO
Explanation	HealthMonitor with healthMonitorId={0} was updated to {1}.
Corrective Action	N/A

Logger	org.midonet.event.topology.HealthMonitor.DELETE
Message	DELETE healthMonitorId={0}.
Level	INFO
Explanation	HealthMonitor with healthMonitorId={0} was deleted.
Corrective Action	N/A

API server events

This section describes the messages associated with API server events.

NSDB (Network State Database)

Logger	org.midonet.event.api.Nsdb.CONNECT
Message	CONNECT Connected to the NSDB cluster.
Level	INFO
Explanation	API server was connected to the NSDB cluster.
Corrective Action	N/A

Logger	org.midonet.event.api.Nsdb.DISCONNECT
Message	DISCONNECT Disconnected from the NSDB cluster.
Level	WARNING
Explanation	API server was disconnected from the NSDB cluster.

Corrective Action	If the connection is restored after this event, no corrective action is required. If this event continues, check the network connection between the API server and the NSDB cluster.
Logger	org.midonet.event.api.Nsdb.CONN_EXPIRE
Message	CONN_EXPIRE Connection to the NSDB cluster expired.
Level	ERROR
Explanation	The connection from the API server to the NSDB cluster expired.
Corrective Action	Check the network connection between the API server and the NSDB cluster and restart the MidoNet API server so it reconnects to the NSDB cluster.

MidoNet Agent events

This section describes the messages associated with MidoNet Agent events.

NSDB

Logger	org.midonet.event.agent.Nsdb.CONNECT
Message	CONNECT Connected to the NSDB cluster.
Level	INFO
Explanation	MidoNet Agent was connected to the NSDB cluster.
Corrective Action	N/A

Logger	org.midonet.event.agent.Nsdb.DISCONNECT
Message	DISCONNECT Disconnected from the NSDB cluster.
Level	WARNING
Explanation	MidoNet Agent was disconnected from the NSDB cluster.
Corrective Action	If the connection is restored after this event, no corrective action is required. If this event continues, check the network connection between the MidoNet Agent and the NSDB cluster.

Logger	org.midonet.event.agent.Nsdb.CONN_EXPIRE
Message	CONN_EXPIRE Connection to the NSDB cluster expired. Shutting down the MidoNet Agent.
Level	ERROR
Explanation	The connection from the MidoNet Agent to the NSDB cluster expired. Shutting down the MidoNet Agent.
Corrective Action	Check the network connection between the MidoNet Agent node and the NSDB cluster and restart the MidoNet Agent service on the node so it reconnects to the NSDB cluster.

Interface

Logger	org.midonet.event.agent.Interface.DETECT
Message	NEW interface={0}
Level	INFO
Explanation	MidoNet Agent detected a new interface={0}.
Corrective Action	N/A

Logger	org.midonet.event.agent.Interface.UPDATE
--------	--

Message	UPDATE interface={0} was updated.
Level	INFO
Explanation	MidoNet Agent detected an update in interface={0}.
Corrective Action	N/A

Logger	org.midonet.event.agent.Interface.DELETE
Message	DELETE interface={0} was deleted.
Level	INFO
Explanation	MidoNet Agent detected that interface={0} was deleted.
Corrective Action	N/A

Service

Logger	org.midonet.event.agent.Service.START
Message	START Service started.
Level	INFO
Explanation	Service started.
Corrective Action	N/A

Logger	org.midonet.event.agent.Service.EXIT
Message	EXIT Service exited.
Level	WARNING
Explanation	Service exited.
Corrective Action	Restart the MidoNet Agent service if this event happened unintentionally. If this event recurs, file a ticket in the bug tracker for further investigation by developers.

Packet Tracing

To configure packet tracing (via logging) in a MidoNet Agent (Midolman), the 'mm-trace' command can be used.

A MidoNet Agent can hold a set of filters that, when matching on an incoming packet, will cause it to log everything about its simulation to the agent's log file, regardless of the configured log level.

All trace messages have a "cookie:" prefix to identify its packet, and that can be used as a grep expression to filter out any non-tracing messages.



Important

The filters are not persistent, they are lost every time the agent is rebooted.

However, mm-trace prints the filters in exactly the same syntax that it will accept to re-add them again, allowing operators to easily replay the commands.

Usage

All available options can be displayed with the '-help' option:

```
$ mm-trace --help
-h, --host <arg>  Host (default = localhost)
```

```

-p, --port <arg>    JMX port (default = 7200)
--help              Show help message

Subcommand: add - add a packet tracing match
-d, --debug          logs at debug level
--dst-port <arg>     match on TCP/UDP destination port
--ethertype <arg>    match on ethertype
--ip-dst <arg>       match on ip destination address
--ip-protocol <arg>  match on ip protocol field
--ip-src <arg>       match on ip source address
-l, --limit <arg>    number of packets to match before disabling
this trace
--mac-dst <arg>      match on destination MAC address
--mac-src <arg>      match on source MAC address
--src-port <arg>     match on TCP/UDP source port
-t, --trace          logs at trace level
--help              Show help message

Subcommand: remove - remove a packet tracing match
-d, --debug          logs at debug level
--dst-port <arg>     match on TCP/UDP destination port
--ethertype <arg>    match on ethertype
--ip-dst <arg>       match on ip destination address
--ip-protocol <arg>  match on ip protocol field
--ip-src <arg>       match on ip source address
-l, --limit <arg>    number of packets to match before disabling
this trace
--mac-dst <arg>      match on destination MAC address
--mac-src <arg>      match on source MAC address
--src-port <arg>     match on TCP/UDP source port
-t, --trace          logs at trace level
--help              Show help message

Subcommand: flush - clear the list of tracing matches
-D, --dead-only      flush expired tracers only
--help              Show help message

Subcommand: list - list all active tracing matches
-L, --live-only      list active tracers only
--help              Show help message

```

Example

```

$ mm-trace list
$ mm-trace add --debug --ip-dst 192.0.2.1
$ mm-trace add --trace --ip-src 192.0.2.1 --dst-port 80
$ mm-trace list
tracer: --debug --ip-dst 192.0.2.1
tracer: --trace --ip-src 192.0.2.1 --dst-port 80
$ mm-trace remove --trace --ip-src 192.0.2.1 --dst-port 80
Removed 1 tracer(s)
$ mm-trace flush
Removed 1 tracer(s)

```

Port mirroring

Port mirroring lets operators monitor arbitrary subsets of traffic in the overlay in specified vports. This can be useful for passive monitoring or for active troubleshooting.

MidoNet v5.0 introduces port mirroring based on these concepts:

1. A new type of virtual device: **mirror**.
2. Each mirror is associated with a destination virtual port, through its **to-port** attribute. This is where mirror traffic will be copied to.

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

- AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

Now the operator can see all local traffic in that bridge by tcpdump'ing on the monitoring port:

```
hypervisor01$ sudo ip netns exec mon tcpdump -nei monns
tcpdump: verbose output suppressed, use -v or -vv for full protocol decode
listening on monns, link-type EN10MB (Ethernet), capture size 65535 bytes
```

By the same means, the operator could mirror any other slice of traffic and do so from any point in the virtual overlay. If a mirror is applied to the upstream facing port of the router, the mirror will see the MAC and IP addresses as that port sees them.

Each mirror can be applied at any number of devices, and can hold several match conditions to capture different slices of traffic. Similarly, each mirroring hook in a device, can have several mirrors applied. Thus the operator has total freedom in selecting which traffic to monitor in his monitoring port, or, by creating different network interfaces and adding more vports to the monitoring bridge, he could also send different kinds of traffic to different monitoring ports.

Removing port mirrors. To remove port mirrors from a bridge, use the `clear` command:

```
midonet> bridge bridge0 clear in-mirrors
midonet> bridge bridge0 show
bridge bridge0 name a-tenant state up
```


15. VXLAN configuration

Table of Contents

VXLAN Gateway	88
VXLAN Coordinator	90
VXLAN Flooding Proxy	90
Connecting to the VTEP	90
Setting up a connection between a VTEP and a Neutron network	92
Enabling connection between VTEP and MidoNet hosts	93
VXLAN Gateway high availability (VTEP side, active-passive mode)	94
Troubleshooting VTEP/VXGW configuration	95
CLI commands used for working with the VXGW	100

MidoNet supports the Virtual Extensible LAN (VXLAN) technology.

What is VXLAN?

VXLAN is a network virtualization technology that uses a VLAN-like encapsulation technique to encapsulate MAC-based OSI layer 2 Ethernet frames within layer 3 UDP packets.

This type of encapsulation (Ethernet-in-IP) is much better suited to Software Defined Networks than either VLANs (802.1q) or even stacked VLANs (Q-in-Q).

Another important advantage of VXLAN over traditional VLAN is its 24-bit VXLAN ID thanks to which VXLAN can scale up to over 16 million logical networks. By comparison - the maximum number of VLANs is 4096.

How is VXLAN supported in MidoNet?

MidoNet provides VXLAN implementation through:

- VXLAN Gateway, to bridge the overlay with physical L3 hosts in the underlay.
- VXLAN tunneling between MidoNet hosts.

VXLAN Gateway

The VXLAN Gateway (VXGW) allows a virtual bridge to be extended to a physical L2 segment that is reachable via an L3 network and a VXLAN-capable physical switch.

A VXLAN-capable physical switch is also referred to as a *hardware VTEP* (VXLAN Tunnel End Point). The VXGW allows creating one or many VXLAN-based Logical Switches that span any number of hardware VTEPs and a single MidoNet-ODP cloud.

The VXGW has the following advantages:

- Provides L2 connectivity between VMs in an overlay and servers in a physical L2 segment.

The MidoNet controller will handle the exchange of learned MACs among all VTEPs and MidoNet's Network State Database (NSDB) automatically.

VXLAN Coordinator

The Coordinator is the component of the MidoNet architecture responsible for VXLAN support.

The Coordinator has the following responsibilities:

- Exposing VTEP state through the MidoNet REST API.
- Configuring the VTEP switch in order to implement the bindings configured through the MidoNet REST API.
- Acting as an L2 control plane for traffic flowing between MN and the VTEP.

VXLAN Flooding Proxy

The VXLAN Gateway controller running in the MidoNet Cluster nodes will try to populate the MAC Remote tables in VTEPs so that the switch can tunnel traffic directly to the exact hypervisor that hosts the destination VM.

Depending on the virtual topology, it may not always be possible to instruct the VTEP to tunnel to a specific physical location. This will typically happen on BUM (Broadcast, Unknown and Multicast) traffic. In these cases MidoNet will instruct the VTEP to tunnel the packet to a service node in order for it to be simulated and delivered to the right destination. This node is called the "Flooding Proxy", and it has the same properties:

- The Flooding Proxy (FP) is one single node elected among all the MidoNet hosts that belong to the same tunnel zone as the VTEP.
- The FP will be in charge of simulating BUM traffic, and tunnelling the packet to their destination (typically a hypervisor).
- Upon failure of the currently elected Flooding Proxy, the MidoNet cluster will use a weighted algorithm to elect a new "Flooding Proxy" role, and instruct the VTEP to tunnel all BUM traffic to it for simulation.

The weight assigned to a MidoNet Agent defaults to 1, and can be altered issuing the following command on the MidoNet CLI:

```
host <host-alias> set flooding-proxy-weight <new-weight>
```

Higher weights will imply a higher probability of the host being chosen as Flooding Proxy.

To exclude an Agent from the candidate set for Flooding Proxy, assign a weight of 0.

Note that the Flooding Proxy may potentially process a large volume of traffic. In these circumstances it is recommended to assign a much higher weight to a dedicated host.

Connecting to the VTEP

Use this procedure to connect MidoNet to a hardware VTEP. This step is required before any Neutron networks can be bound to port-VLAN pairs on that VTEP.

- MidoNet will expect that the `Physical_Switch` table on the VTEP contains a record with the management IP, management port and tunnel IP of this VTEP. Keep these details at hand as they will be needed to configure the VTEP and any bindings to Neutron networks. Use the following command to dump the contents of this table:

Make sure that your VTEP also registers all the physical ports. You can verify this by examining the `Physical_Ports` table in the VTEP. Only ports present in this table will be available for binding to a Neutron network. Use the following command would display all physical ports, replacing `<vtep-name>` with the name assigned to the `Physical_Switch` (you can check this using the previous `vtep-ctl list Physical_Switch` command).

2. After setting up the VTEP you might need to test connectivity to both the tunnel and management interfaces. Both should be *up*.

```
$ telnet <management-ip> <management-port>
Trying <management-ip>...
Connected to <management-ip>
Escape character is '^['.
```

At this point, paste this on the console:

And you should see this reply:

This just verified that the OVSDB server that holds the configuration of this VTEP is active and handling connections. If you did not get a similar reply, please review the configuration of the VTEP.

3. The VXLAN service in Midonet is enabled by default. However, if the service has been previously disabled, you can enable it by typing the following command and restarting at least one cluster instance in order to become the VXLAN coordinator:

4. Having ensured that the VTEP has been properly configured, now you're ready to add the VTEP to the MidoNet configuration.

For information, see [the section called “Adding a VTEP” \[101\]](#).



Important

Apart from the information on the VLAN-port assignment, and VTEP management interface IP and port, you will also need the identifier of a TunnelZone of type *vtep*. All the hosts running MidoNet Agent daemons that you want to create VXLAN tunnels with the VTEP should be members of this tunnel zone, using the local IP that each host uses as a VXLAN tunnel endpoint.

After you successfully add a VTEP to the MidoNet configuration, the API Server connects to its (VTEP's) management interface and collects all the required information for creating a Logical Bridge. For more information, see the "Logical Bridge" section.

Setting up a connection between a VTEP and a Neutron network

Use this procedure to set up a connection between a VTEP and a Neutron network in MidoNet.

For this procedure you will need to know the VTEP's management IP and port, the physical port on the VTEP and the VLAN ID at this port to which you are connecting, the UUID of the Neutron network which you want to connect to the VTEP, and the IP addresses of all the hosts on the Neutron network that you want to communicate with the VTEP.

1. Create a tunnel zone of type *vtep*.

All hosts that want to communicate with the VTEP using VXLAN tunnels are required to belong to a tunnel zone of type VTEP, using the IP that each of them will use as VXLAN tunnel endpoint.

To create a tunnel zone issue the following command in the MidoNet CLI:

```
midonet> tunnel-zone create name vtep_zone type vtep
tzone0
```

As you can see you just created a *vtep* type tunnel zone, *tzone1*.

2. Add a VTEP to MidoNet and assign it to the *vtep* tunnel zone that you created, using the local IP that this host it meant to use in VXLAN tunnels to the VTEP. Note that this IP may be the same that the host uses to communicate with other MidoNet hosts.

```
midonet> vtep add management-ip 192.168.2.10 management-port 6632
tunnel-zone tzone0
vtep0
```

You can determine the state of the VTEP connection by listing the current VTEPs. If your VTEP had been added successfully you should see a similar message, saying `connection-state connected`.

```
vtep vtep0 name VTEP-NAME description OVS VTEP Emulator management-ip
192.168.2.10 management-port 6632 tunnel-zone tzone0 connection-state
connected
```

3. Create a binding between the VTEP and a Neutron network behind a MidoNet bridge. For that, you will need the UUID of the Neutron network behind that bridge. To find out the UUID use these commands:

```
midonet> list bridge
bridge bridge0 name public state up
midonet> show bridge bridge0 id
765cf657-3cf4-4d79-9621-7d71af38a298
```

The Neutron network you are binding the VTEP to is behind `bridge0`, and it has the UUID of `765cf657-3cf4-4d79-9621-7d71af38a298` as you can see in the output of the command.

- a. To find out their addresses, use these commands:

b. Add the host's IP address to the same tunnel zone as the VTEP:

Repeat these steps for every host in the Neutron network that you want to communicate with the VTEP.



5. Create a binding between the VTEP's VLAN 10 and the Neutron network behind the bridge0.

```
midonet> vtep vtep0 add binding network-id
765cf657-3cf4-4d79-9621-7d71af38a298 physical-port swp1 vlan 10
```

Congratulations, you have just created a binding between the network behind the VTEP's VLAN 10 physical port `swp1` and the Neutron network with the UUID `765cf657-3cf4-4d79-9621-7d71af38a298` in MidoNet.



To be able to test the connection between the VTEP and MidoNet (i.e. to *ping* MidoNet from a host on the VTEP) you have to modify the default ingress security rule, by adding to it the IP address of the host (pinging the host from MidoNet should work without any additional configuration). For more information, see [the section called “Enabling connection between VTEP and MidoNet hosts” \[93\]](#).

Enabling connection between VTEP and MidoNet hosts

By default Neutron includes a security rule on all networks that restricts forwarding only to traffic addressed to IP/MAC of VMs on that network.

By binding a network to physical ports in a VTEP, we're effectively adding hosts to the L2 segment of this Neutron network that Neutron itself does not know about, and thus traffic addressed to these physical hosts will be dropped.

1. In the MidoNet CLI find out what the ingress default security rule is by issuing this command:

Locate the ingress security rule that is assigned to the neutron network. In this case, we'll use chain0 (OS_SG_64d9f3df-9875-4896-ad0c-ffc5bba84c5e_INGRESS) rule chain, the ingress chain.

```
midonet> chain chain0 list rule
rule rule0 ethertype 2048 proto 0 tos 0 ip-address-group-src ip-address-
group0 fragment-policy unfragmented pos 1 type accept
rule rule1 ethertype -31011 proto 0 tos 0 ip-address-group-src ip-
address-group0 fragment-policy unfragmented pos 2 type accept
```

3. Now, go ahead and add the IP address of the host on the VTEP to the security group `ip-address-group0`.

```
midonet> ip-address-group ip-address-group0 add ip address 172.16.0.3
address 172.16.0.3
```

You should now be able to ping a host in MidoNet from host 172.16.0.3 on the VTEP (providing they are in the same tunnel zone).

VXLAN Gateway high availability (VTEP side, active-passive mode)

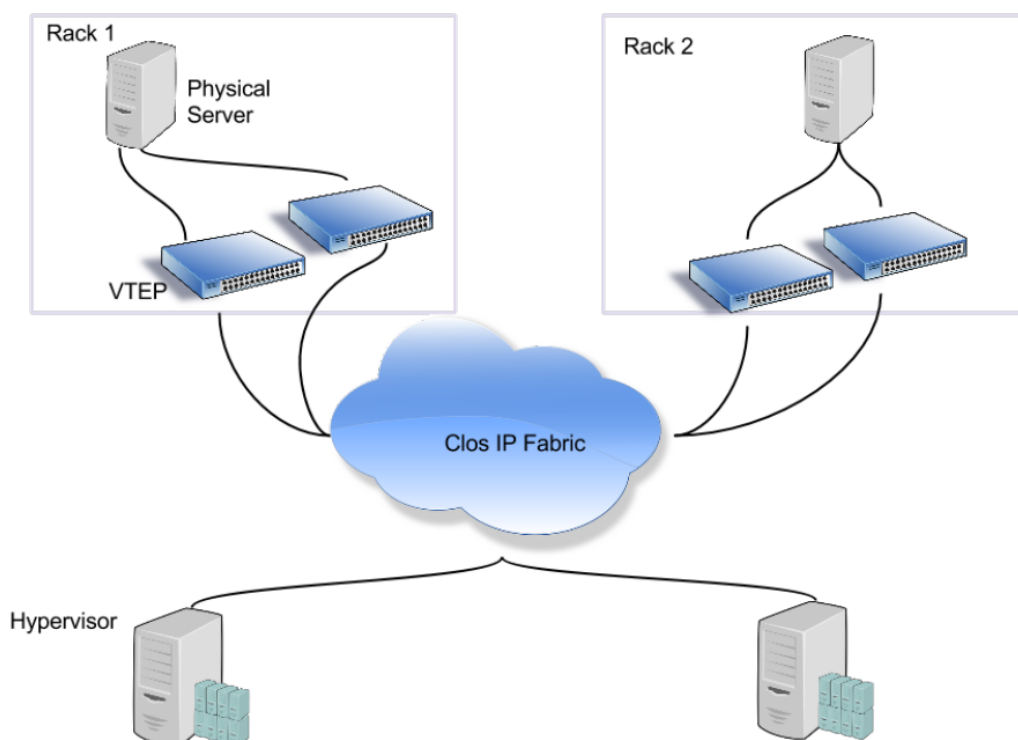
MidoNet supports VTEP HA by using two ToR VTEPs, both connected to each host in the rack with an active-passive configuration.

There are two possible failure scenarios:

- The active link from a host is lost. The passive link will take over, so traffic from the host will now flow through the other VTEP and its MACs will be learned there.
- A VTEP dies. Any host whose active link connects to this VTEP will now move to the passive link, with the same result: MACs and IPs migrate to a different VTEP.

MAC migrations will be applied from the VxLAN Gateway Service as soon as the OVSDB notifies about them, the new MacLocations will be propagated to VTEPs and NSDB, and traffic from agents and other VTEPs will now flow through the new VTEP.

Figure 15.1. VTEP HA



Troubleshooting VTEP/VXGW configuration

VTEP deployments have a relatively large number of moving pieces and potential failure points. This guide will focus on troubleshooting MidoNet and the integration with the VTEP. For specifics on the configuration of the logical switch please refer to your vendor's documentation.

Is the MidoNet API able to connect to the VTEP

After following the procedure to add a VTEP as described in [the section called "Adding a VTEP" \[101\]](#), the expected output should be as follows:

```
midonet> vtep add management-ip 192.168.2.10 management-port 6633 tunnel-
zone tzone0
vtep0
midonet> vtep list
vtep vtep0 name VTEP-NAME management-ip 192.168.2.10 management-port 6633
tunnel-zone tzone0 connection-state connected
```

The same output should appear for VTEPs already added to MidoNet.

Note that the state is `connected`. An `error` state will indicate that the VTEP's management IP is unreachable from the MidoNet API.

Is the VTEP well configured?

A typical reason for the VTEP being on `error` state is a misconfiguration of the VTEP OVDSB instance. You can verify this by executing the following command on the console:

```
ovsdb-client dump hardware_vtep
```

Scroll down to the `Physical_Switch` table, which will look like this:

Obtaining information about a VTEP

Use this command to obtain information about a selected VTEP.

Syntax

```
vtep <vtep-alias> show <property>
```

where *property* is one of the following VTEP's attributes:

- name
- description
- management-ip
- management-port
- connection-state
- tunnel-zone

Result

The command returns the following information about the VTEP:

- name
- description
- management IP address (the same as the IP used with the command)
- mgmt_port (the same as the port values used by the command)
- tunnel IP addresses
- connection state (one of: connected, disconnected, error. The state is error if the end-point is not a VXLAN End Point)
- the tunnel-zone to which this VTEP belongs.

Example

Successful command:

```
midonet> vtep vtep0 show id  
ba2739df-87cf-458f-9ad2-39885cab217d
```

```
midonet> vtep vtep0 show management-ip  
192.168.2.10
```

Adding a VTEP binding

Use this command to bridge the port-VLAN pair on a VTEP to a specified Neutron network.

Syntax

```
vtep <vtep-alias> add binding physical-port <port-name> vlan <vlan-id>  
network-id <neutron-network-id>
```


An example of a unsuccessful command:

```
midonet> vtep vtep0 delete binding
Syntax error at: ...binding
```

Deleting a VTEP

Use this command to delete a VTEP.

Syntax

```
vtep <vtep-alias> delete
```

Result

Issuing this command completely deletes a VTEP from MidoNet's list of known VTEPs.

The command will fail if any of the VTEP's port-VLAN pairs are bound to any Neutron networks. For more information, see [the section called "Removing a VTEP binding" \[104\]](#)

Example

```
midonet> vtep vtep0 delete
```



```
midonet> host host1 binding add interface eth1 port bridge0:port3
```

Fail-over/Fail-back

In combination with the Spanning Tree Protocol (STP) enabled on the physical bridges, MidoNet VABs are able to provide fail-over capabilities by forwarding Bridge Protocol Data Unit (BPDU) frames across their trunk ports.

Assuming that both physical switches belong to the same bridged network, as a result of the STP, both devices detect a loop through MidoNet's VAB and one switch chooses to block its trunk. For example, let's assume the left switch blocks. The VAB only sees ingress traffic from the right trunk, and thus associates all source MAC addresses seen in those frames to the right trunk.

A variety of events, including failures in the network, may result in the switches deciding to invert the state of the trunks. An example could be MidoNet losing connection to the left switch, and thus stop forwarding BPDUs to/from the right bridge and undoing the loop.

In such a fail-over scenario, traffic would start flowing from the other switch. With this change, MidoNet now detects ingress traffic on a new port, and thus updates its internal MAC-port associations. If the former state of the topology is restored (that is, MidoNet recovers connectivity to the left switch), MidoNet will again react and update its MAC-port associations.

The fail-over/fail-back times depend on the STP configuration on the switches, mainly the "forward delay," and the nature of the traffic. With standard values, and continuous traffic ingressing from the trunks, fail-over and fail-back cycles should be completed in 50 seconds, plus MAC learning time.

17. Service Containers

Table of Contents

Configuration	111
Management	112
Scheduling	113
Troubleshooting	118

MidoNet supports the network functions virtualization, such as L4LB, VPN or BGP, using third party software that executes in namespaces connected to the MidoNet virtual topology. To make the integration of these components easier and to control the namespace scheduling, MidoNet 5.1 introduces the concept of service containers. Currently, they are used only for VPNaaS.

Service containers are Linux namespaces with a port connected to a MidoNet virtual device, and they run a set of programs that are configured according to the type of the service container.

Service containers are launched and destroyed automatically by the MidoNet Cluster with the topology object to which they correspond. For example, an IPSec container, used to encrypt traffic for VPN connections and handle the security associations, launches when creating the first Neutron `VpnService` on a router, and it is destroyed when deleting the last.

Container scheduling on the available physical computes is also managed by the MidoNet Cluster, and it takes into account the computes' availability, one of the several configurable scheduling policies, and a per-compute host container weight parameter allowing the cloud operator to differentiate the container load between different computes.

Architecture

A service container is an object that exists in the MidoNet topology, like a virtual device or port. When created, it is linked to an existing virtual port, which becomes the *container port* through which the container namespace exchanges traffic with the other devices inside the MidoNet virtual topology.

A *service type* describes the container namespace and the processes that are executed inside it, while a *configuration identifier* keeps a reference to an object in the MidoNet topology that indicates how the namespaces and the container processes are configured.

The container *scheduling policy* is specified using a *service container group* object. Therefore, several containers can share the same scheduling policy by referencing the same container group. The MidoNet Cluster uses the policy as defined by the container group to determine the physical compute where a container is launched. This scheduling is controlled via the binding of the container port on a given host. For more information see [the section called "Scheduling Policies" \[116\]](#).

In addition, a service container group specifies an *allocation policy*, which specifies how existing compute hosts are allocated to new containers given the host pool determined from the *scheduling policy*. Whereas the scheduling policy considers configurable prop-

erties of the hosts, such as membership to a particular host or port group, the allocation policy uses live attributes of the hosts, such as their *container weight* or number of launched containers. For more information see [the section called “Allocation Policies” \[117\]](#).

When a MidoNet Agent launches a scheduled container, it runs a set of routines that are customized for a particular container type, and which configures and starts the container processes. These can be loaded from external libraries, allowing the cloud operator to deploy additional containerized services as needed, without necessarily installing a new version of the agent.

MidoNet agents that run service containers use the same container library to monitor and report the container status to the NSDB. This allows the MidoNet Cluster to report it via the REST API, and to drive the container scheduling upon a container failure.

General Fields

The following table lists the fields that describe a service container in the MidoNet topology.

Field Name	Type	Description
id	UUID	The service container identifier.
serviceType	String	A unique name that identifies the type of the service container.
groupId	UUID	The identifier of the service container group that describes the scheduling policy.
portId	UUID	The identifier of the container virtual port. This is an exterior port connected to a device from the virtual topology. The port's binding information describes the compute host where the container is launched.
configurationId	UUID	The identifier of an object from the MidoNet topology that provides the container configuration. The type of this object is specific for a particular service type, and the agent library handling that service type is responsible for giving it the proper interpretation.

Status Fields

The following table lists the fields that describe the status of a service container.

Field Name	Type	Description
statusCode	String	Indicates the running state of a container, it can be one of the following: <i>starting</i> , <i>running</i> , <i>stopping</i> , <i>stopped</i> and <i>error</i> .
statusMessage	String	A custom status message, usually containing information about the processes running inside the container.
hostId	UUID	The identifier of the agent host where the container is running. This field is present only for starting, running or stopping, containers.
namespaceName	String	The name of the Linux namespace hosting the container.
interfaceName	String	The name of the virtual Ethernet interface connected to the container.

Container Group Fields

The following table lists the fields that describe the service container group. For more information on the service container scheduling and scheduling policies see [the section called “Scheduling” \[113\]](#).

Field Name	Type	Description
id	UUID	The service container group identifier.

Field Name	Type	Description
hostGroupId	UUID	The identifier of the host group when using host group scheduling policy.
portGroupId	UUID	The identifier of the port group when using port group scheduling policy.
policy	String	The hosts allocation policy for the containers belonging to this group.

Configuration

The containers are configured using `mn-conf`. Both the MidoNet Agent and Cluster run a Container Service, which is responsible for launching the containers at the agent, and scheduling the containers across different computes at the cluster.

Agent Configuration

The container configuration at the agent is found under the `agent.containers` namespace. The following table lists the configuration keys, and modifying the value of a configuration key requires restarting the agent in order to take effect.

Field Name	Type	Description
enabled	Boolean	Whether an agent instance runs the containers service. This allows the operator to disable containers across all compute hosts or at specific compute hosts. The default is <code>true</code> .
timeout	Duration	The timeout for a container operation to complete. The default is 30 seconds.
shutdown_grace_time	Duration	The timeout for all containers to stop gracefully when the containers service is stopped. The default is 30 seconds.
thread_count	Integer	The number of threads that manage the containers executing at this agent. This determines the number of containers that can be changed in parallel. Once created a container is bound to the same thread on which it was created.
log_directory	String	The name of the log directory where the agent writes a record of all container operations. The log is used to cleanup containers following an agent failure. The default value is <code>containers</code> . The directory is created in the same directory as the agent log, which is configurable using the <code>midolman.log.dir</code> JVM system property in the agent startup script, and which defaults to <code>/var/log/midolman</code> .

Specific container types may include additional configuration keys, as follows.

IPSec Containers Configuration

The following table lists the configuration keys for IPSec containers, which are used to implement VPNaaS. All keys are found under the `agent.containers.ipsec` namespace.

Field Name	Type	Description
logging_enabled	Boolean	Indicates whether the IPSec container writes to the agent log the log messages reported by the IPSec processes running inside the container. The default is <code>true</code> . When enabled, these messages are logged using the logger name <code>org.midonet.containers.ipsec.ipsec-pluto</code> .
logging_poll_interval	Duration	The polling interval for the IPSec container log. The default is 250 milliseconds.
logging_timeout	Duration	The timeout when closing the log file. The default is 3 seconds.

Field Name	Type	Description
status_update_interval	Duration	The interval for updating the container status and health. The default is 5 seconds.

Cluster Configuration

The containers configuration at the cluster is found under the `cluster.containers` namespace. The following table lists the configuration keys, and modifying the value of a configuration key requires restarting the cluster in order to take effect.

Field Name	Type	Description
enabled	Boolean	Whether a cluster node runs the containers service. This allows the operator to disable containers across all cluster nodes or at specific cluster nodes. The default is <code>true</code> .
scheduler_timeout	Duration	The timeout for a container to reach <code>running</code> state when scheduled at a compute host. The default is 10 seconds. For more information on container scheduling see the section called "Scheduling" [113] .
scheduler_retry	Duration	Indicates the time interval after which the scheduler attempts to reschedule a down container.
scheduler_max_retries	Integer	The maximum number of attempts the scheduler will retry to schedule a down container. Once this number is reached, the scheduler will schedule a new container only as a response to an external event, such as the available hosts have changed or the operator triggered a manual scheduling.
scheduler_bad_host_lifetime	Duration	When an agent fails to start a container within the expected timeout interval, that host is considered bad and not used for subsequent scheduling operations of the same container. This key controls for how long a host is considered bad after failing to launch a container. The default is 5 minutes. Set this value to zero (0) to disable bad hosts for scheduling. For more information on container scheduling see the section called "Scheduling" [113] .



Note

If the containers service is enabled at multiple cluster nodes, only one cluster node determined by via distributed leader election will actively perform the container scheduling at a time. If the leader node fails, another node will take over the scheduling role automatically. The failover latency equals the configured ZooKeeper session timeout at `zookeeper.session_timeout`.

Management

Normally, the service containers are launched and destroyed automatically by the MidoNet Cluster. Therefore, the operator would typically use the MidoNet CLI only to inspect the containers configuration and running state.

Listing the Service Containers

To list information about all service containers, enter the command:

```
midonet> list container
container container0 container-group cgroup0 configuration 63c6fa43-
abcd-43b0-870b-c72ab6f2c4dc port port0 type IPSEC status running host host0
namespace vpn-e9a728b2 interface vpn-e9a728b2
container container1 container-group cgroup1 configuration
401994b8-4fec-428d-9133-40c51988578a port port1 type IPSEC status running
host host0 namespace vpn-a55d15fe interface vpn-a55d15fe
```

The example above shows two service containers of IPSEC type that have been scheduled at host `host0` and are currently in the `running` state. If the containers are not running, then the output of the command would be as follows:

```
midonet> list container
container container0 container-group cgroup0 configuration 63c6fa43-
abcd-43b0-870b-c72ab6f2c4dc port port0 type IPSEC status stopped
container container1 container-group cgroup1 configuration
401994b8-4fec-428d-9133-40c51988578a port port1 type IPSEC status stopped
```

Listing the Service Container Group

To list information about all service container groups, enter the command:

```
midonet> list container-group
cgroup cgroup0 policy least
cgroup cgroup1 policy least
cgroup cgroup2 policy least
```

Viewing the Container Status Message

The container status message is usually a text with multiple lines, and therefore it is not included in the container listing. Therefore, to view the message, enter the command:

```
midonet> container container0 show status-message
000 using kernel interface: netkey
000 interface lo/lo ::1@500
000 interface lo/lo 127.0.0.1@4500
000 interface lo/lo 127.0.0.1@500
000 interface vpn-e9a728b2-ns/vpn-e9a728b2-ns 1.0.0.2@4500
000 interface vpn-e9a728b2-ns/vpn-e9a728b2-ns 1.0.0.2@500
000 interface vpn-e9a728b2-ns/vpn-e9a728b2-ns 169.254.0.2@4500
000 interface vpn-e9a728b2-ns/vpn-e9a728b2-ns 169.254.0.2@500
...
```



Note

You can also use the `show` command for any other fields of a service container or service container group object.



Warning

MidoNet also provides `create` and `delete` commands for the service containers. However, these should not be used under normal circumstances and should be reserved only for troubleshooting purposes.

Scheduling

The MidoNet Cluster uses the Containers Service to schedule existing service containers at different compute hosts. Container scheduling considers the following input data:

- The MidoNet agent must be running at the compute host, and it must have the Containers Service enabled. See [the section called “Agent Configuration” \[111\]](#) for how to enable or disable the service at the agent.
- The configuration of the agent host in the NSDB should have a positive value for the *service container weight*. See [the section called “Host Container Weight” \[115\]](#).
- The configuration of the agent host in the NSDB should have a positive value for the *service container limit*. See [the section called “Host Container Limit” \[115\]](#).

- The container service scheduling policy as defined by the service container group. See [the section called “Scheduling Policies” \[116\]](#).
- The container service host allocation policy as defined by the service container group. See [the section called “Allocation Policies” \[117\]](#).
- The container status reported by the last host where the container has been scheduled.

Scheduling is performed at the container level, and it considers the following possible states for a container:

State	Description
DOWN	The container is not scheduled at a host.
SCHED- ULED	The container is scheduled at a host, but that host has not yet reported the container status.
UP	The container is scheduled at a host, and that host has reported the container status as <code>running</code> .

Only `DOWN` containers are being scheduled, and scheduling is performed when the container is created and subsequently whenever the available hosts and their *state* changes. By host state, we understand whether the agent at that host is running the container service and its corresponding *container weight*.

A host is selected to run a container if it is *eligible*, that is it meets the first three criteria from above about running the container service, having a positive *container weight*, does not currently run more containers than its *container limit*, and if the host has not been previously marked as *bad* for the same container. A *bad* host is a host that has previously failed to launch a container. This may happen when the host is overloaded, or it does not have installed the necessary software that is needed to launch the container. Marking the hosts as *bad* is a temporary action, and it prevents the rescheduling logic for retrying indefinitely with the same agent. The interval during which a host is marked as *bad* after failing to launch a container is configurable in the cluster containers configuration. See [the section called “Cluster Configuration” \[112\]](#).



Important

The values of the `scheduler_timeout` and `scheduler_bad_host_lifetime` configuration keys are paramount to ensure that containers are scheduled properly, and their values should consider both the acceptable downtime during a container or agent failure and the time required to start the container at the host. Some containers, such as the `IPSEC` containers, requires a non-negligible time interval to start, and therefore setting the container timeout to a very low value coupled with a long bad host lifetime will lead to containers not being scheduled. This situation may occur during the bulk creation of containers with a small number of compute hosts. In such cases, we recommend increasing the `scheduler_timeout` or setting the `scheduler_bad_host_lifetime` to zero.

A host is selected from the eligible hosts using the current allocation policy, specified by the service container group. See [the section called “Allocation Policies” \[117\]](#).



Note

Running containers are not affected by changes to a host's *container weight*, except when it is set to zero.



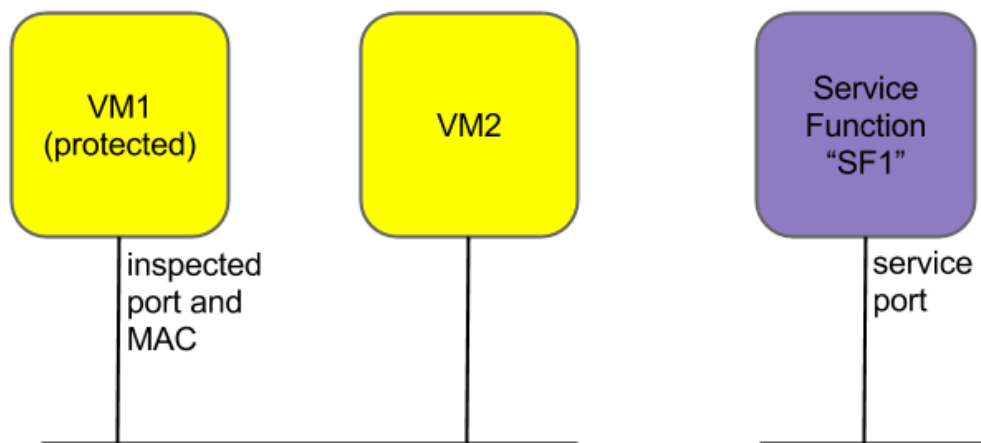
L2Insertion Object

The L2Insertion object has these fields (see the REST documentation for more detail):

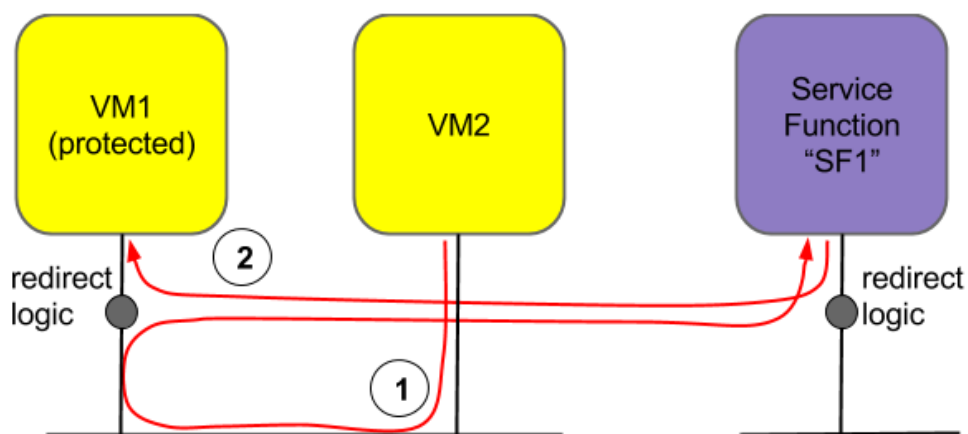
1. inspected VM port UUID
2. inspected VM MAC
3. service port UUID
4. VLAN tag
5. fail-open (true/false)
6. position

The first 3 fields are the most important to consider. Imagine an L2Insertion configured as described by the following diagram.

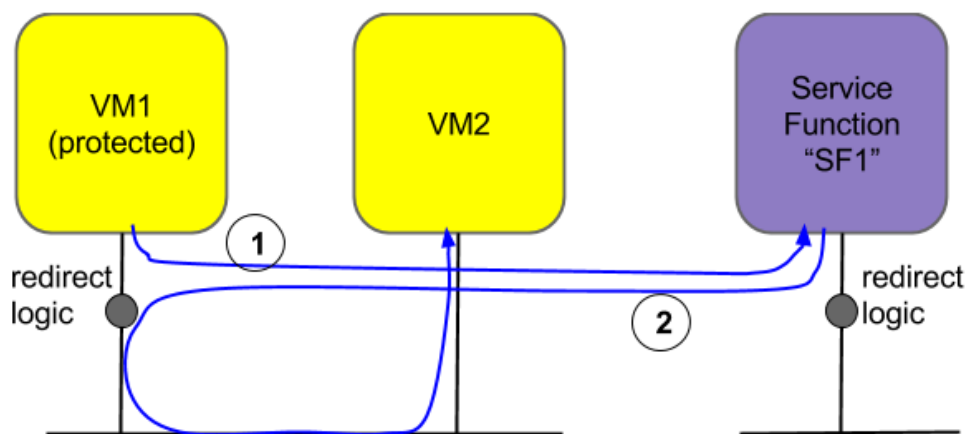
The position is an absolute position in the list of insertions. If two insertions have the same position the order is undefined. When we load the insertions we basically sort by the position field. In many ways, it acts more like a weight.



In such a scenario, MidoNet will set up redirect logic as shown in the diagram below. When VM2 sends a packet to VM1 (and specifically to VM1's MAC), at the point labelled "redirect logic" (after port-level firewall, not shown) the packet is redirected to the service function (red arrow 1). If the service function allows the packet and therefore re-emits it, the redirect logic at the service function redirects the packet so that it completes its journey to VM1 (red arrow 2).



The next diagram shows how packets from the protected port (and MAC) are redirected.



The diagrams do now show how field number 4, the VLAN tag, is used. If the VLAN tag is set to a value above zero, then MidoNet will push a VLAN tag with that value onto packets redirected into the service function. MidoNet will also expect the VLAN tag to be left on the packet if it is re-emitted from the service function. MidoNet also uses the PCP bits in the VLAN header for its own purposes, these should not be modified by the service function.

The use of field number 5, fail-open, is illustrated in the diagram below. If fail-open is false, then when the service function fails all redirect traffic is essentially dropped. If fail-open is true, there are cases when MidoNet can detect failure of the service function link (link status only, we do not yet have active monitoring neither in-band or out-of band) and the redirection logic will automatically allow the traffic to skip the failed service function.



19. Router Peering

Table of Contents

Diagrams	128
Creating tenant networks and (peer) routers	130
Creating a VTEP Router Gateway Device	130
Creating a Multi-site Network and L2 Gateway Connection	131
Peering a Tenant Router	132
Adding a Remote MAC Entry (Endpoint) to the Multi-site Network	134
Configure Edge Bindings	134
Deleting a Remote MAC Entry	134
Deleting an L2 Gateway Connection	135
Deleting an L2 Gateway	135
Deleting a Gateway Device	135
Neutron CLI Gateway Device Reference	135

The router peering feature of MidoNet provides overlay connectivity between multiple sites with VXLAN tunneling. The following section describes the Neutron CLI and REST API commands to set up router peering. All of the following operations are executed only by the cloud operator.

The following terminologies are used throughout this document:

- VTEP Router:

Admin owned routers that are VXLAN endpoints in both sites. They are the virtual gateway devices of each site that enables site-to-site connectivity using VXLAN technology.

- Gateway Device:

An abstraction of a gateway device in a cloud. A VTEP router is an example of a gateway device. This abstraction lets us define other types of gateway devices in the future, such as hardware VTEP.

- Tunnel IPs:

The IP addresses used by the gateway device when constructing VXLAN header (the outer header). These IP addresses are known only to the gateway devices. It is also important to note that these are not associated with any particular port, which allows VXLAN traffic to egress/ingress from any edge router port.

- Multi-Site Network:

Admin-owned network that stretches across multiple sites. Another terminology used to refer to this network is "Inter-AZ Network". This network contains a private CIDR allocated for the tenant that emulates a single L2 segment uplink network of their routers, which in actuality maps to networks across sites (with the same CIDR).

- Peer Router:

A regular tenant router that connects to the multi-site network. This router typically already has an external network attached for Internet. In the router peering scenario, there will be an additional router interface that connects to the multi-site network.

Gateway Device Tunnel IP: 200.200.0.1

Site B

Tenant Network CIDR: 10.0.1.0/24

Multi-Site Network (must match in Site A):

- CIDR: 192.168.0.0/24
- VNI: 100

Router Interface Port (Multi-Site Network < - > Peer Router):

- IP: 192.168.0.2
- MAC: 6F:E4:5A:FA:8E:09

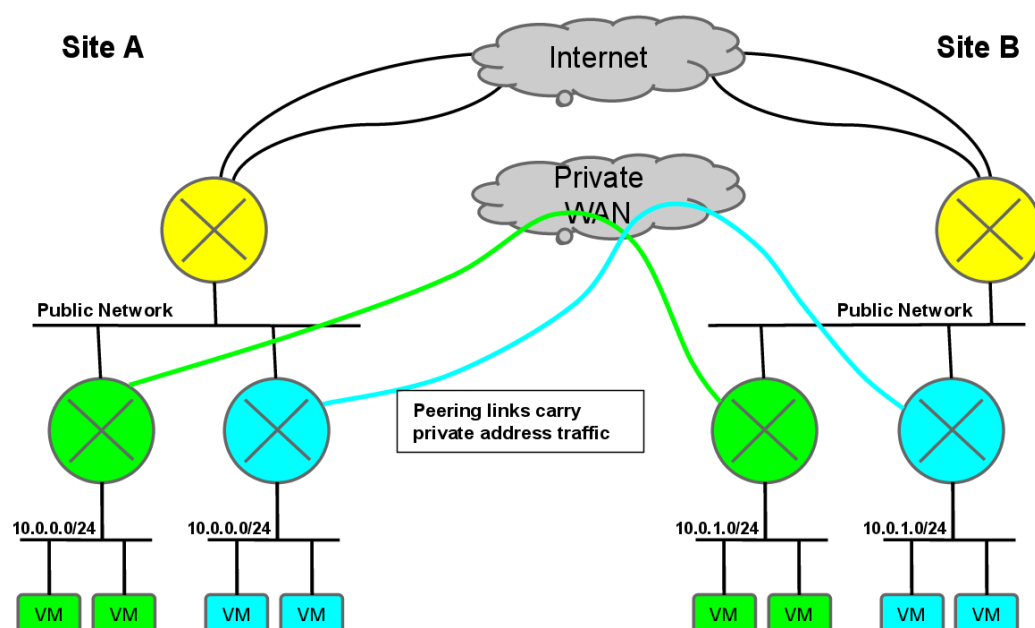
Gateway Device Tunnel IP: 200.200.0.2

Note that Gateway Device Tunnel IPs are in the same subnet in the example above, but in practice they do not need to be. They just need to be IP reachable.

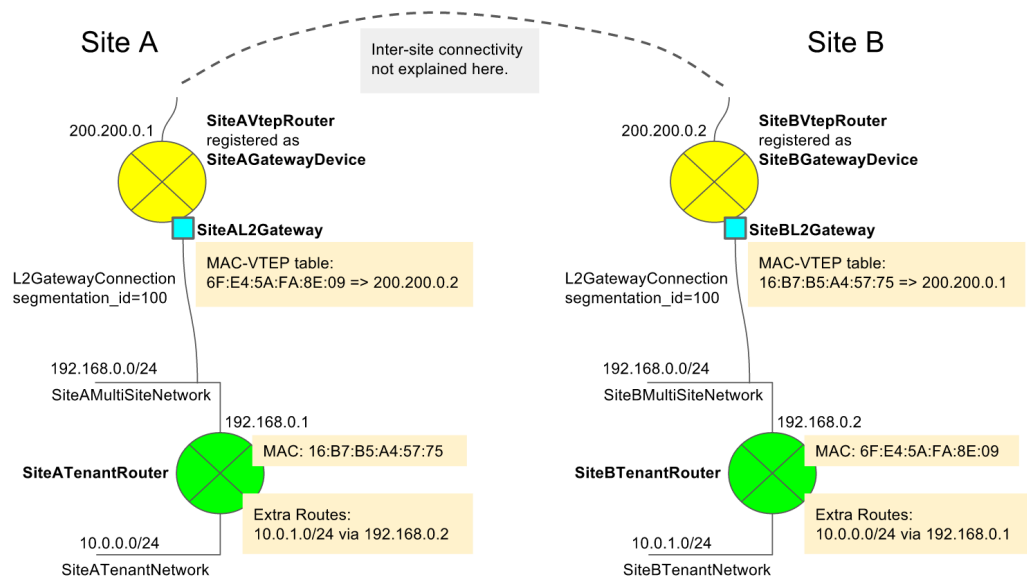
As you will see below, setting up multi-site peering requires a fairly large number of API operations and coordination among the sites that are mistake prone. Thus, it is recommended that these steps below are orchestrated and automated by an external service.

Diagrams

Overview



Details



When VM 10.0.0.3 in Site A sends an IP packet to VM 10.0.1.3 in Site B, the packet appears as follows when it egresses SiteAGatewayDevice:

- Outer Ethernet header: src is MAC of 200.200.0.1, dst is MAC of SiteAVtepRouter's peer
- Outer IP header: src is 200.200.0.1, dst is 200.200.0.2, proto is UDP
- UDP header: src is hash of inner packet, dst is standard VXLAN port
- VXLAN header: VNI is 100
- Inner Ethernet header:
 - src is 16:B7:B5:A4:57:75 (MAC of SiteATenantRouter on SiteAMultiSiteNetwork)
 - dst is 6F:E4:5A:FA:8E:09 (MAC of SiteBTenantRouter on SiteBMultiSiteNetwork)
- Inner IP header: src is 10.0.0.3, dst is 10.0.1.3

Note that SiteATenantRouter's and SiteBTenantRouter's IP addresses do not appear in the packet because they are just intermediate routers. However, traceroute from VM 10.0.0.3 to 10.0.1.3 should find these hops:

- 10.0.0.1
- 192.168.0.1
- 192.168.0.2
- 10.0.1.1
- 10.0.1.3

Traceroute will not detect SiteAVtepRouter or SiteBVtepRouter, which are part of the tunneling infrastructure.


```

gateway-device-update
    Update a given Gateway Device.

gateway-device-remote-mac-entry-create
    Create Gateway Device Remote Mac Entry information.

gateway-device-remote-mac-entry-delete
    Delete a given Gateway Device Remote Mac Entry.

gateway-device-remote-mac-entry-list
    List Gateway Device Remote Mac Entries.

gateway-device-remote-mac-entry-show
    Show information of a given Gateway Device Remote Mac Entry.

```

neutron gateway-device-create

```

usage: neutron gateway-device-create [-h]
                                     [-f {json,shell,table,value,yaml}]
                                     [-c COLUMN] [--max-width <integer>]
                                     [--noindent] [--prefix PREFIX]
                                     [--request-format {json}]
                                     [--tenant-id TENANT_ID]
                                     [--management-ip MANAGEMENT_IP]
                                     [--management-port MANAGEMENT_PORT]
                                     [--management-protocol
MANAGEMENT_PROTOCOL]
                                     [--type <hw_vtep | router_vtep>]
                                     [--resource-id RESOURCE_ID] [--name
NAME]
                                     [--tunnel-ip TUNNEL_IP]

```

Create Gateway Device information.

Optional arguments:

```

-h, --help            Show this help message and exit.
--request-format {json}
                        DEPRECATED! Only JSON request format is supported.
--tenant-id TENANT_ID
                        The owner tenant ID.
--management-ip MANAGEMENT_IP
                        Management IP to the device. Defaults to None.
--management-port MANAGEMENT_PORT
                        Management port to the device. Defaults to None.
--management-protocol MANAGEMENT_PROTOCOL
                        Management protocol to manage the device: ovssdb or
                        none. If management IP and port are specified,
                        defaults to ovssdb. Otherwise to none.
--type <hw_vtep | router_vtep>
                        Type of the device: hw_vtep or router_vtep. Defaults
                        to hw_vtep.
--resource-id RESOURCE_ID
                        Resource UUID or None (for type router_vtep will be
                        router UUID).
--name NAME           User defined device name.
--tunnel-ip TUNNEL_IP
                        IP address on which gateway device originates or
                        terminates tunnel.

```

neutron gateway-device-delete

```

usage: neutron gateway-device-delete [-h]

```

```

[--request-format {json}]
GATEWAY_DEVICE

```

Delete a given gateway-device.

Positional arguments:

```

GATEWAY_DEVICE      ID or name of gateway_device to delete.

```

Optional arguments:

```

-h, --help          Show this help message and exit.
--request-format {json}
                    DEPRECATED! Only JSON request format is supported.

```

neutron gateway-device-show

```

usage: neutron gateway-device-show [-h]
                                   [-f {json,shell,table,value,yaml}]
                                   [-c COLUMN] [--max-width <integer>]
                                   [--noindent] [--prefix PREFIX]
                                   [--request-format {json}] [-D] [-F
                                   FIELD]
                                   GATEWAY_DEVICE

```

Show information of a given gateway-device.

Positional arguments:

```

GATEWAY_DEVICE      ID or name of Gateway Device to look up.

```

Optional arguments:

```

-h, --help          Show this help message and exit.
--request-format {json}
                    DEPRECATED! Only JSON request format is supported.
-D, --show-details  Show detailed information.
-F FIELD, --field FIELD
                    Specify the field(s) to be returned by server. You
                    can repeat this option.

```

neutron gateway-device-update

```

usage: neutron gateway-device-update [-h]
                                     [--request-format {json}]
                                     [--name NAME] [--tunnel-ip TUNNEL_IP]
                                     GATEWAY_DEVICE

```

Update a given gateway-device.

Positional arguments:

```

GATEWAY_DEVICE      ID or name of Gateway Device to update.

```

Optional arguments:

```

-h, --help          Show this help message and exit.
--request-format {json}
                    DEPRECATED! Only JSON request format is supported.
--name NAME          User defined device name.
--tunnel-ip TUNNEL_IP
                    IP address on which gateway device originates or
                    terminates tunnel.

```

neutron gateway-device-remote-mac-entry-create

```
usage: neutron gateway-device-remote-mac-entry-create [-h]
                                                    [-f {json,shell,
table,value,yaml}]
                                                    [-c COLUMN]
                                                    [--max-width
<integer>]
                                                    [--noindent]
                                                    [--prefix PREFIX]
                                                    [--request-format
{json}]
                                                    --mac-address
MAC_ADDRESS
                                                    --vtep-address
VTEP_ADDRESS
                                                    --segmentation-id
SEGMENTATION_ID
                                                    GATEWAY_DEVICE
```

Create Gateway Device Remote Mac Entry information.

Positional arguments:

GATEWAY_DEVICE	ID or name of the Gateway Device.
----------------	-----------------------------------

Optional arguments:

```
-h, --help            Show this help message and exit.
--request-format {json}
                        DEPRECATED! Only JSON request format is supported.
--mac-address MAC_ADDRESS
                        Remote MAC address.
--vtep-address VTEP_ADDRESS
                        Remote VTEP Tunnel IP.
--segmentation-id SEGMENTATION_ID
                        VNI to be used.
```

neutron gateway-device-remote-mac-entry-delete

```
usage: neutron gateway-device-remote-mac-entry-delete [-h]
                                                    [--request-format
{json}]
                                                    REMOTE_MAC_ENTRY
                                                    GATEWAY_DEVICE
```

Delete a given Gateway Device Remote MAC Entry.

Positional arguments:

REMOTE_MAC_ENTRY	ID of Remote MAC Entry to delete.
GATEWAY_DEVICE	ID or name of the Gateway Device.

Optional arguments:

```
-h, --help            Show this help message and exit.
--request-format {json}
                        DEPRECATED! Only JSON request format is supported.
```

neutron gateway-device-remote-mac-entry-list

```
usage: neutron gateway-device-remote-mac-entry-list [-h]
```


Optional arguments:

```
-h, --help          Show this help message and exit.
--request-format {json}
                    DEPRECATED! Only JSON request format is supported.
-D, --show-details  Show detailed information.
-F FIELD, --field FIELD
                    Specify the field(s) to be returned by server. You
                    can repeat this option.
```

20. FWaaS Logging

Table of Contents

Creating a router and firewall	141
Creating a logging resource and firewall log	142
Updating a logging resource and firewall log	143
Deleting a logging resource and firewall log	143

Midonet's implementation of Neutron FWaaS Logging offers a single method of configuring firewalls to log events. FWaaS Logging can be configured by either cloud operator or tenants. The log files are generated on the physical hosts, so only the cloud operators can access them. It is expected that an external service collects and aggregates the generated logs. FWaaS Logging lets you audit and debug firewall rules that you manage.

There are two resources defined for FWaaS Logging.

- Logging Resource - Top level model that represents a collection of firewall logging objects.
- Firewall Logging - Object that represents logging of a particular firewall object.

The flow of operations is as follows:

1. Create a router that will be associated with the firewall for which events will be logged
2. Create a firewall for which events will be logged
3. Create a logging resource
4. Create a firewall log associated with the firewall created in step 2, and the logging resource in step 3.

Multiple firewall logs can be associated with a single logging resource. The logging resource can then be enabled or disabled, providing the ability to enable or disable event logging for multiple firewall logs in a single command.

The following is a concrete example of the topology creation.

Creating a router and firewall

```
$ neutron router-create example_router
Created a new router:
```

Field	Value
admin_state_up	True
description	
external_gateway_info	
id	d7acbd7b-8a1b-44ff-a92f-291fac96b4ca
name	example_router
routes	
status	ACTIVE
tenant_id	admin

```

+-----+
$ neutron firewall-policy-create example_policy
Created a new firewall_policy:
+-----+
| Field          | Value                                     |
+-----+-----+
| audited        | False                                   |
| description    |                                         |
| firewall_rules |                                         |
| id             | 52c3a819-4846-4a8b-8418-453abf965210 |
| name          | example_policy                         |
| shared         | False                                  |
| tenant_id     | admin                                  |
+-----+-----+

$ neutron firewall-create example_policy --router d7acbd7b-8alb-44ff-
a92f-291fac96b4ca
Created a new firewall:
+-----+
| Field          | Value                                     |
+-----+-----+
| admin_state_up | True                                    |
| description    |                                         |
| firewall_policy_id | 52c3a819-4846-4a8b-8418-453abf965210 |
| id            | 2b3cd306-df9b-47c8-adeb-ab8381e42bb2 |
| name          |                                         |
| router_ids    | d7acbd7b-8alb-44ff-a92f-291fac96b4ca |
| status        | PENDING_CREATE                         |
| tenant_id     | admin                                  |
+-----+-----+

$ neutron firewall-rule-create --protocol tcp --destination-port 80 --
action deny
Created a new firewall_rule:
+-----+
| Field          | Value                                     |
+-----+-----+
| action         | deny                                    |
| description    |                                         |
| destination_ip_address |                                         |
| destination_port | 80                                    |
| enabled        | True                                   |
| firewall_policy_id |                                         |
| id            | 386fa9fd-eb39-41c6-aadc-f1eaeef1714cf |
| ip_version     | 4                                      |
| name          |                                         |
| position       |                                         |
| protocol       | tcp                                    |
| shared         | False                                  |
| source_ip_address |                                         |
| source_port    |                                         |
| tenant_id     | admin                                  |
+-----+-----+

$ neutron firewall-policy-insert-rule example_policy 386fa9fd-eb39-41c6-
aadc-f1eaeef1714cf
Inserted firewall rule in firewall policy example_policy

```

Creating a logging resource and firewall log

```

$ neutron logging-create example_lr
Created a new logging_resource:
+-----+
| Field          | Value                                     |
+-----+-----+

```

description	
enabled	False
firewall_logs	
id	09e17388-479c-40ea-a124-a4fa53935322
name	example_lr
tenant_id	admin

```
$ neutron logging-firewall-create --firewall-id 2b3cd306-df9b-47c8-adeb-ab8381e42bb2 example_lr
Created a new firewall_log:
```

Field	Value
description	
firewall_id	2b3cd306-df9b-47c8-adeb-ab8381e42bb2
fw_event	ALL
id	36a5440b-dde5-47b6-9cf1-85ee325b4600
logging_resource_id	09e17388-479c-40ea-a124-a4fa53935322
tenant_id	admin

Updating a logging resource and firewall log

```
$ neutron logging-update --enabled True 09e17388-479c-40ea-a124-a4fa53935322
Updated logging_resource: 09e17388-479c-40ea-a124-a4fa53935322
```

```
$ logging-firewall-update 2b3cd306-df9b-47c8-adeb-ab8381e42bb2
09e17388-479c-40ea-a124-a4fa53935322 --fw-event DROP
Updated firewall_log: 2b3cd306-df9b-47c8-adeb-ab8381e42bb2
```

Deleting a logging resource and firewall log

```
$ logging-firewall-delete 09e17388-479c-40ea-a124-a4fa53935322
09e17388-479c-40ea-a124-a4fa53935322
Deleted firewall_log: 09e17388-479c-40ea-a124-a4fa53935322
```

```
$ neutron logging-delete 09e17388-479c-40ea-a124-a4fa53935322
Deleted logging_resource: 09e17388-479c-40ea-a124-a4fa53935322
```


AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

AFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT - DRAFT -

enable_dhcp	True
gateway_ip	10.1.0.1
host_routes	
id	538fd8eb-b36f-400e-8bba-995b8f674f21
ip_version	4
ipv6_address_mode	
ipv6_ra_mode	
name	LEFT_SUB
network_id	1331e08a-8f9c-4d75-b7ff-661a474290d0
subnetpool_id	
tenant_id	admin
updated_at	2016-07-08T08:15:20

```
# neutron router-create LEFT_ROUTER
```

Created a new router:

Field	Value
admin_state_up	True
description	
external_gateway_info	
id	384ed43a-3297-4fac-9f1a-a29e1749bffa
name	LEFT_ROUTER
routes	
status	ACTIVE
tenant_id	admin

```
# neutron router-interface-add \
  384ed43a-3297-4fac-9f1a-a29e1749bffa \
  538fd8eb-b36f-400e-8bba-995b8f674f21
```

Added interface b0c5fbcb-d944-4eaf-83b5-a30039449900
to router 384ed43a-3297-4fac-9f1a-a29e1749bffa.

```
# neutron router-gateway-set \
  384ed43a-3297-4fac-9f1a-a29e1749bffa \
  EXTERNAL
```

Set gateway for router 384ed43a-3297-4fac-9f1a-a29e1749bffa

Create the right side:

```
# neutron net-create RIGHT
```

Created a new network:

Field	Value
admin_state_up	True
created_at	2016-07-11T01:37:37
description	
id	75b9789e-6e13-4bfe-9668-c38dcdb7a67
name	RIGHT
port_security_enabled	True
provider:network_type	midonet
router:external	False
shared	False
status	ACTIVE
subnets	
tags	
tenant_id	admin
updated_at	2016-07-11T01:37:37

```
# neutron subnet-create \
  --name RIGHT_SUB \
  75b9789e-6e13-4bfe-9668-c38dcdb7a67 \
  10.2.0.0/24 \
  --gateway 10.2.0.1
```

Created a new subnet:

Field	Value
allocation_pools	{"start": "10.2.0.2", "end": "10.2.0.254"}
cidr	10.2.0.0/24
created_at	2016-07-11T01:38:15
description	
dns_nameservers	
enable_dhcp	True
gateway_ip	10.2.0.1
host_routes	
id	8058c633-4616-42ec-9838-2d2a8786441d
ip_version	4
ipv6_address_mode	
ipv6_ra_mode	
name	RIGHT_SUB
network_id	75b9789e-6e13-4bfe-9668-c38dcdb7a67
subnetpool_id	
tenant_id	admin
updated_at	2016-07-11T01:38:15

```
# neutron router-create RIGHT_ROUTER
```

Created a new router:

Field	Value
admin_state_up	True
description	
external_gateway_info	
id	24011587-0b8c-484b-9e35-f6779aa27b98
name	RIGHT_ROUTER
routes	
status	ACTIVE
tenant_id	admin

```
# neutron router-interface-add \
  24011587-0b8c-484b-9e35-f6779aa27b98 \
  8058c633-4616-42ec-9838-2d2a8786441d
```

Added interface 5dbf6189-a24b-415e-a934-bbc3f4761b8b
to router 24011587-0b8c-484b-9e35-f6779aa27b98.

```
# neutron router-gateway-set \
  24011587-0b8c-484b-9e35-f6779aa27b98 \
  EXTERNAL
```

Set gateway for router 24011587-0b8c-484b-9e35-f6779aa27b98

Create the VPN Policies

```
# neutron vpn-ikepolicy-create IKEPOLICY
```

Created a new ikepolicy:

Field	Value
-------	-------

auth_algorithm	sha1
description	
encryption_algorithm	aes-128
id	4d1bdf17-4c1a-48c9-a880-9d7a31356ab3
ike_version	v1
lifetime	{"units": "seconds", "value": 3600}
name	IKEPOLICY
pfs	group5
phase1_negotiation_mode	main
tenant_id	admin

```
# neutron vpn-ipsecpolicy-create IPSECPOLICY
```

Created a new ipsecpolicy:

Field	Value
auth_algorithm	sha1
description	
encapsulation_mode	tunnel
encryption_algorithm	aes-128
id	7607ac27-8708-451c-9df7-d913ec99c11a
lifetime	{"units": "seconds", "value": 3600}
name	IPSECPOLICY
pfs	group5
tenant_id	admin
transform_protocol	esp

Create the VPN Services and Connections

For the left side:

```
# neutron vpn-service-create \
  --name LEFT_CONN \
  384ed43a-3297-4fac-9f1a-a29e1749bffa \
  538fd8eb-b36f-400e-8bba-995b8f674f21
```

Created a new vpnservice:

Field	Value
admin_state_up	True
description	
external_v4_ip	200.200.200.2
external_v6_ip	
id	54b02ff5-698c-421e-807b-b1fd9ee69e45
name	LEFT_CONN
router_id	384ed43a-3297-4fac-9f1a-a29e1749bffa
status	PENDING_CREATE
subnet_id	538fd8eb-b36f-400e-8bba-995b8f674f21
tenant_id	admin

```
> neutron ipsec-site-connection-create \
  --name LEFT_SITE_CONN \
  --vpnservice-id 54b02ff5-698c-421e-807b-b1fd9ee69e45 \
  --ikepolicy-id 4d1bdf17-4c1a-48c9-a880-9d7a31356ab3 \
  --ipsecpolicy-id 7607ac27-8708-451c-9df7-d913ec99c11a \
  --peer-address 200.200.200.3 \
  --peer-id 200.200.200.3 \
  --peer-cidr 10.2.0.0/24 \
```



```
--psk secret

Created a new ipsec_site_connection:
```

Field	Value
admin_state_up	True
auth_mode	psk
description	
dpd	{"action": "hold", "interval": 30, "timeout": 120}
id	52684388-c74e-4c06-bf37-a16a045e6ecc
ikepolicy_id	4d1bdf17-4c1a-48c9-a880-9d7a31356ab3
initiator	bi-directional
ipsecpolicy_id	7607ac27-8708-451c-9df7-d913ec99c11a
local_ep_group_id	
mtu	1500
name	LEFT_SITE_CONN
peer_address	200.200.200.3
peer_cidrs	10.2.0.0/24
peer_ep_group_id	
peer_id	200.200.200.3
psk	secret
route_mode	static
status	PENDING_CREATE
tenant_id	admin
vpnservice_id	54b02ff5-698c-421e-807b-b1fd9ee69e45

For the right side:

```
# neutron vpn-service-create \
  --name RIGHT_CONN \
  24011587-0b8c-484b-9e35-f6779aa27b98 \
  8058c633-4616-42ec-9838-2d2a8786441d

Created a new vpnservice:
```

Field	Value
admin_state_up	True
description	
external_v4_ip	200.200.200.3
external_v6_ip	
id	0784e875-5ad0-4757-8005-e6b00aab9bd3
name	RIGHT_CONN
router_id	24011587-0b8c-484b-9e35-f6779aa27b98
status	PENDING_CREATE
subnet_id	8058c633-4616-42ec-9838-2d2a8786441d
tenant_id	admin

```
# neutron ipsec-site-connection-create \
  --name RIGHT_SITE_CONN \
  --vpnservice-id 0784e875-5ad0-4757-8005-e6b00aab9bd3 \
  --ikepolicy-id 4d1bdf17-4c1a-48c9-a880-9d7a31356ab3 \
  --ipsecpolicy-id 7607ac27-8708-451c-9df7-d913ec99c11a \
  --peer-address 200.200.200.2 \
  --peer-id 200.200.200.2 \
  --peer-cidr 10.1.0.0/24 \
  --psk secret
```

```
Created a new ipsec_site_connection:
```

Field	Value
-------	-------

admin_state_up	True
auth_mode	psk
description	
dpd	{"action": "hold", "interval": 30, "timeout": 120}
id	53368c20-9f1c-49c3-9f69-2566cd8656bd
ikepolicy_id	4dlbdf17-4c1a-48c9-a880-9d7a31356ab3
initiator	bi-directional
ipsecpolicy_id	7607ac27-8708-451c-9df7-d913ec99c11a
local_ep_group_id	
mtu	1500
name	RIGHT_SITE_CONN
peer_address	200.200.200.2
peer_cidrs	10.1.0.0/24
peer_ep_group_id	
peer_id	200.200.200.2
psk	secret
route_mode	static
status	PENDING_CREATE
tenant_id	admin
vpnservice_id	0784e875-5ad0-4757-8005-e6b00aab9bd3

riod of time (referred to as "*max_kbps*"), and maximum data burst allowed with a very short window (referred to as "*max_burst_kbps*", which unfortunately implies a data rate, when in practice, this is a maximum limit on data sent at a particular time).

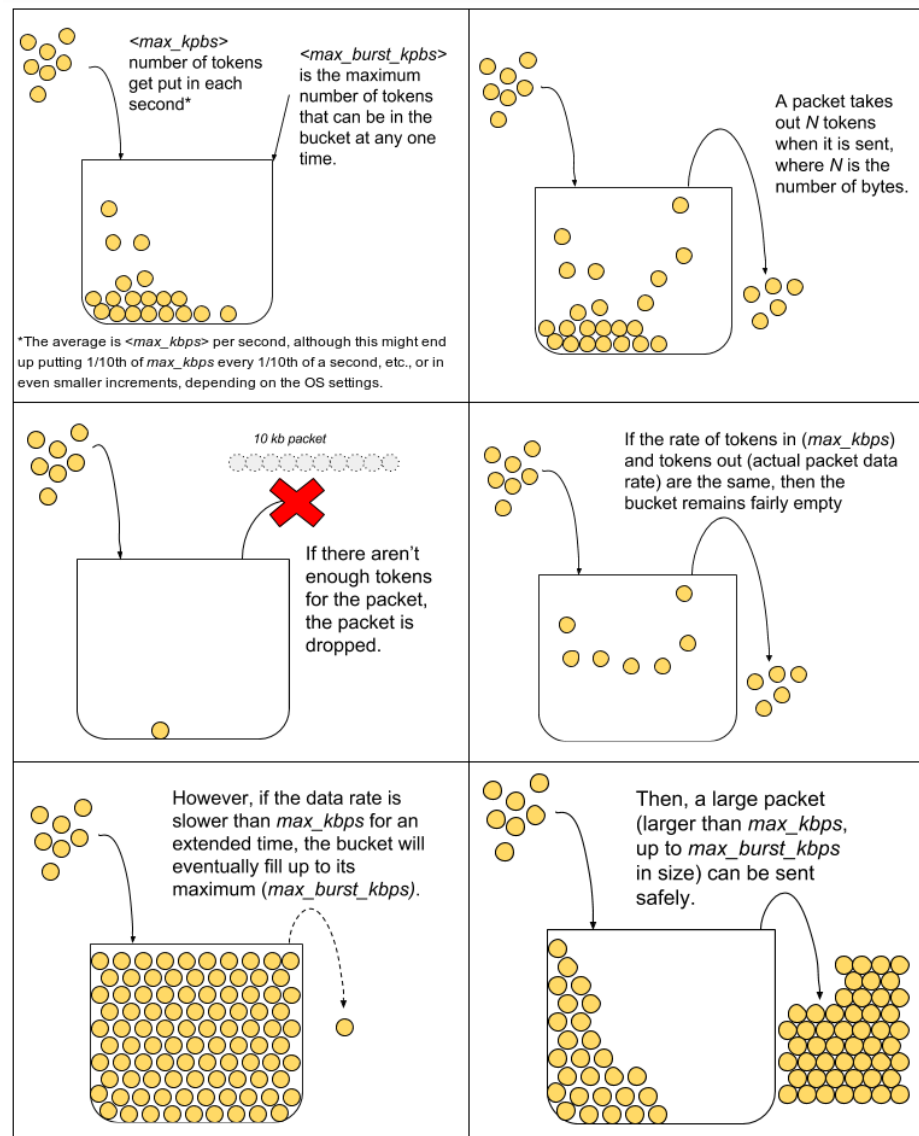


Note

How Linux Meters Traffic

The best way to contextualize this is to actually explain briefly how the mechanics of how Linux implements data rate limiting. Basically, picture a bucket with tokens. The *max_kbps* represents how many tokens put into the bucket each second (e.g. 1,000 kbps means 1,000 tokens are put into the bucket every second). The *max_burst_kbps* represents the maximum number of tokens allowed in the bucket. For every byte of data that passes through the port, a token is removed from the bucket. If there are not enough tokens in the bucket for the packet to take, the packet is dropped.

This means that the tokens can fill up over time up to a maximum of *max_burst_kbps*, which can then allow a large packet to pass through, however if the traffic is constant at *max_kbps*, this means tokens are being removed from the bucket as fast as they are being added, and larger packets will be dropped as there will never be enough tokens available.



To add a rule directly, use the command:

```
midonet> qos-policy <policy_id> bw-limit-rule create max-kbps <max_kbps>
max-burst-kbps <max_burst_kbps>
```

Where the parameters are the same as in the Neutron command, except that *policy ID* is the UUID or ID specifier of the policy to be updated.

An ID specifier is the short ID given for a policy in the `qos-policy list` command inside the CLI. It is usually of the form `qos-policyN` (where N is an increasing number from 0). Other objects also have ID specifiers when using the interactive CLI. *Note that the ID specifier is neither returned nor can it be used when using `midonet-cli -e op` to run a single command from the shell. In this case, use the UUID instead.*

DSCP Marking Rules

A QoS policy can have rules added to alter the IPv4 ToS field (now used for DSCP) of every packet that enters the cloud via a port. This will affect all IPv4 packets coming in to the cloud on the port with the QoS policy set. Packets exiting the cloud on the port will not be affected. At this time, the DSCP marking rules will affect all packets, mean-

Table of Contents

Starting with release v5.4, MidoNet supports mapping a floating IPv6 address to a host with a fixed IPv4 address. This capability allows IPv6 Internet clients to reach virtual instances in IPv4 tenant networks.

The feature implements stateful NAT64 as specified by [RFC 6146](#), which is a mechanism for IPv4-IPv6 transition and IPv4-IPv6 coexistence. It allows an IPv6 client from the Internet to initiate communication with an IPv4-only server. MidoNet does not support communication initiated by the IPv4-only hosts through statically configured bindings in the stateful NAT64.



Floating IPv6-IPv4 works alongside floating IPv4, and you can use this dual-stack configuration to support communication with both IPv4 and IPv6 hosts from a tenant network. Similarly, the physical uplink interfaces work simultaneously with IPv4 and IPv6.



Floating IPv6-IPv4 is only supported for virtual topologies where the IPv6 external network is directly connected to the edge router.

Floating IPv6-IPv4 is configured from Neutron. Use the following commands to create a virtual topology with floating IPv6-IPv4. You can skip certain steps if you already have an existing network topology.

1. Create an edge router.

2. Create an IPv6 uplink network and subnet. In principle, these will use a globally routable IPv6 prefix. In this example, we use `2001::/64` as the uplink IPv6 prefix.

```
neutron subnet-create \
  --tenant_id admin \
  --disable-dhcp \
```

```
--ip-version 6 \
--name uplink-subnet-6 \
uplink-network 2001::0/64
```

3. Create a port on the uplink network. This port will be bound to the physical uplink interface, <UPLINK>. The IPv6 address assigned to the port will be an address from the uplink IPv6 address, in this example we use 2001::2. We denote by <UPLINK-PORT-ID> the identifier of the uplink port, and by <HOST-ID> the name of the gateway host.

```
neutron port-create uplink-network \
--binding:host_id <HOST-ID> \
--binding:profile type=dict interface_name=<UPLINK> \
--fixed-ip ip_address=2001::2
```



Note

If you setup the uplink port for a dual-stack IPv4-IPv6 configuration, during this step you can specify multiple fixed IP addresses for the uplink port, one for each subnet created on the uplink network. For example, to setup a dual-stack configuration where the IPv4 uplink subnet is 200.0.0.0/30, create a second IPv4 subnet on the uplink network and then assign both an IPv4 and IPv6 address to the uplink port.

```
neutron subnet-create \
--tenant_id admin \
--disable-dhcp --ip-version 4 \
--name uplink-subnet-4 \
uplink-network 200.0.0.0/30
```

```
neutron port-create uplink-network \
--binding:host_id <HOST-ID> \
--binding:profile type=dict interface_name=<UPLINK> \
--fixed-ip ip_address=2001::2 \
--fixed-ip ip_address=200.0.0.2
```

4. Add a router interface from the edge router to the uplink network.

```
neutron router-interface-add edge-router \
port=<UPLINK-PORT-ID>
```

At the point you should be able to ping the IPv6 address of the uplink port from the uplink network.

5. Add to the edge router one or more static routes to allow the tenants to communicate with other IPv6 networks. For example, to add a default route for the IPv6 Internet, use:

```
neutron router-update \
--route destination=::/0,nextthop=2001::1 \
edge-router
```

To add a specific route for an IPv6 network, for instance 3001::/64, use:

```
neutron router-update \
--route destination=3001::/64,nextthop=2001::1 \
edge-router
```



Important

The MidoNet 5.4 BGP implementation does not support IPv6. Therefore, all IPv6 routes must be configured statically on the edge router.

6. Create an external network to interconnect the edge and tenant routers. For floating IPv6-IPv4, this network must have an IPv6 subnet that is also globally routable. In this example, we use 1000::/120 as the public IPv6 prefix.

```
neutron net-create public-network \
  --router:external true
```

```
neutron subnet-create \
  --name public-subnet-6 \
  --ip-version 6 \
  public-network 1000::/120
```



Note

If you setup a dual-stack IPv4 and IPv6 external network to allow both IPv4 and IPv6 communication for the tenant network, during this step you can also create an IPv4 subnet on the external network. Subsequently, this will allow you to create both IPv4 and IPv6 floating IPs. Alternatively, you can also create multiple external networks with one subnet, one for each IP version.

7. Add a router interface from the edge router to the external network.

```
neutron router-interface-add edge-router public-subnet-6
```

8. Create a tenant router.

```
neutron router-create tenant-router
```

9. Set the gateway for the tenant router to the external network.

```
neutron router-gateway-set tenant-router public-network
```



Note

If the external network is configured with dual-stack IPv4 and IPv6 subnets, this step will create a network port with multiple IP addresses, one for each subnet from the external network.

10. Create an IPv4 tenant network. In this example, we use 192.168.0.0/24 for the tenant network.

```
neutron net-create tenant-network
```

```
neutron subnet-create \
  --name tenant-subnet \
  tenant-network 192.168.0.0/24
```

11. Add a router interface from the tenant router to the tenant network.

```
neutron router-interface-add tenant-router tenant-subnet
```

12. Create a floating IP on the external network.

```
neutron floatingip-create public-network
```



Important

If the external network is configured with dual-stack IPv4 and IPv6 subnets, you must specify a subnet in the previous command that indicates the subnet for which the floating IP is created. Otherwise, Neutron will select the first subnet as the default one.

```
neutron floatingip-create --subnet public-subnet-6 public-network
```

- 13 Associate the floating IP with identifier <FIP-ID> created in the previous step to an instance port <VIF-PORT-ID>.

```
neutron floatingip-associate <FIP-ID> <VIF-PORT-ID>
```



Note

You may also create and associate a floating IP to a port during the same command.

```
neutron floatingip-create \
  --port-id <VIF-PORT-ID> \
  --subnet public-subnet-6 \
  public-network
```

Cleanup

Use the following steps to tear-down a virtual network topology that uses floating IPv6-IPv4.

1. Disassociate any floating IPs associated with virtual ports on the tenant network.

```
neutron floatingip-disassociate <FIP-ID>
```

2. Delete all ports from the tenant network (delete any corresponding compute instances as needed).

3. Delete the router interface from the tenant router to the tenant network.

```
neutron router-interface-delete tenant-router tenant-subnet
```

4. Delete the tenant network

```
neutron subnet-delete tenant-subnet
neutron net-delete tenant-network
```

5. Clear the tenant router gateway.

```
neutron router-gateway-clear tenant-router
```

6. Delete the tenant router.

```
neutron router-delete tenant-router
```

7. Delete the router interfaces from the edge router to the external network.

```
neutron router-interface-delete edge-router public-subnet-6
```



Note

If the external network has multiple subnets, delete the router interfaces for all subnets, as needed.

8. Delete the external network.

```
neutron subnet-delete public-subnet-6
neutron net-delete public-network
```



Note

If the external network has multiple subnets, delete all subnets from the external network, as needed.

9. Delete the static routes from the edge router.

```
neutron router-update --no-routes edge-router
```

- 10.Delete the router interface from the edge router to the uplink network.

```
neutron router-interface-delete edge-router \
    port=<UPLINK-PORT-ID>
```

- 11.Delete the uplink network.

```
neutron subnet-delete uplink-subnet-6
neutron net-delete uplink-network
```



Note

If the uplink network has multiple subnets, delete all subnet from the uplink network, as needed.

- 12.Delete the edge router.

```
neutron router-delete edge-router
```

Stateful NAT64

Floating IPv6 to fixed IPv4 uses NAT64 implemented between the edge and tenant routers, where inbound IPv6 traffic is translated to IPv4 according to IPv6-IPv4 translation mappings stored in an internal Binding Information Base (BIB) [RFC 6146]. These mappings are allocated on demand for new flows on a first-come-first-served basis, such that:

- The destination floating IPv6 address maps to the corresponding IPv4 address from the tenant network.
- The source IPv6 address maps to an allocated IPv4 address from a special IPv4 subnet called a NAT64 pool.

NAT64 Pool

There exists a separate NAT64 IPv4 address pool for each tenant router using the Shared Address Space for Carrier-Grade NAT with the IPv4 address range 100.64.0.0/10, as defined by [RFC 6598](#).



Important

Floating IPv6-IPv4 cannot be used with IPv4 tenant networks that use address ranges overlapping with the NAT64 IPv4 address pool.

The following diagram illustrates the NAT64 address translation for the addresses configured in the previous steps. The IP addresses are as follows:

- The floating IPv6 address associated to an instance from a tenant network is 1000::100.

riod of time (referred to as "*max_kbps*"), and maximum data burst allowed with a very short window (referred to as "*max_burst_kbps*", which unfortunately implies a data rate, when in practice, this is a maximum limit on data sent at a particular time).

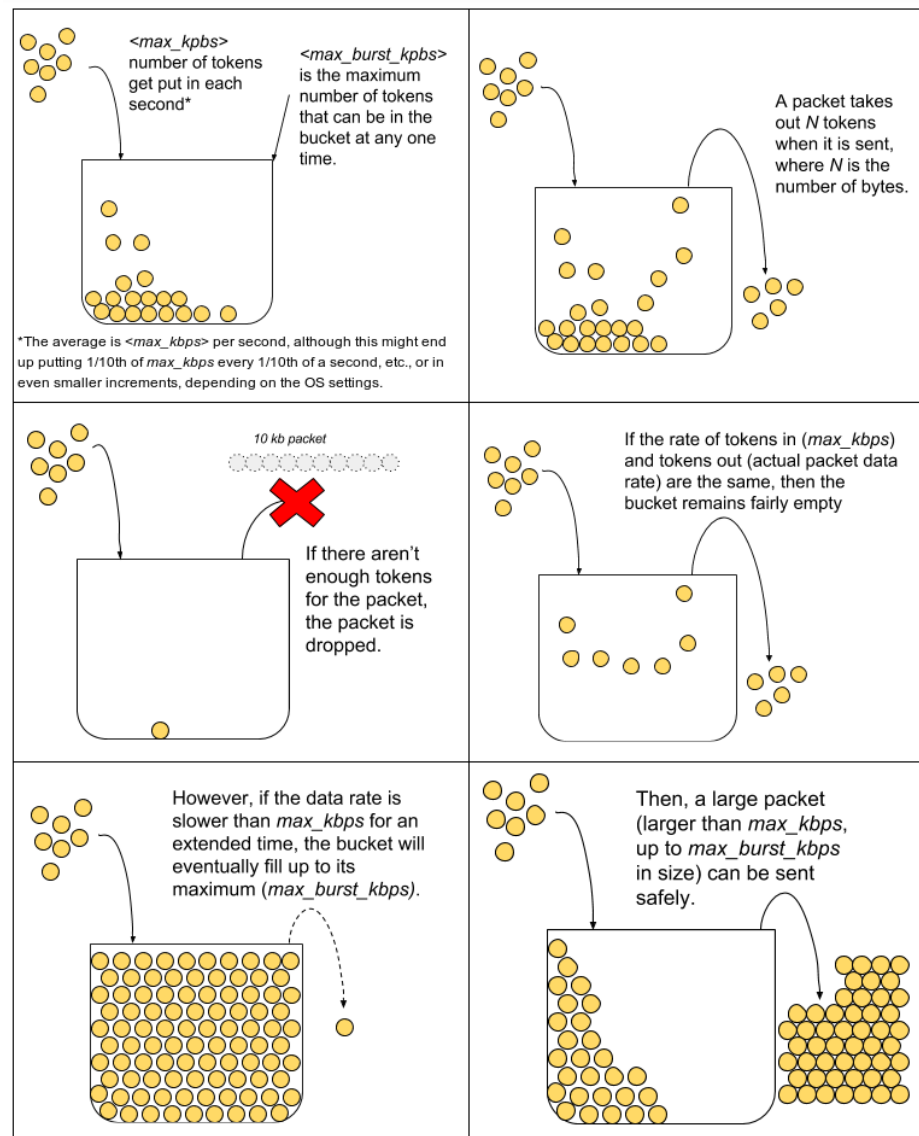


Note

How Linux Meters Traffic

The best way to contextualize this is to actually explain briefly how the mechanics of how Linux implements data rate limiting. Basically, picture a bucket with tokens. The *max_kbps* represents how many tokens put into the bucket each second (e.g. 1,000 kbps means 1,000 tokens are put into the bucket every second). The *max_burst_kbps* represents the maximum number of tokens allowed in the bucket. For every byte of data that passes through the port, a token is removed from the bucket. If there are not enough tokens in the bucket for the packet to take, the packet is dropped.

This means that the tokens can fill up over time up to a maximum of *max_burst_kbps*, which can then allow a large packet to pass through, however if the traffic is constant at *max_kbps*, this means tokens are being removed from the bucket as fast as they are being added, and larger packets will be dropped as there will never be enough tokens available.



Creating a Bandwidth Limit Rule

Bandwidth limit rules can be created via the Neutron command:

```
$ neutron qos-bandwidth-limit-rule-create --max-kbps <max_kbps> --max-burst-kbps <max_burst_kbps> <policy_id>
```

Where *policy ID* is the UUID of the policy to add the rule to, *max_kbps* is the maximum sustained data rate allowed (rate to fill the token bucket in the example above), and *max_burst_kbps* is the maximum packet size allowed (maximum tokens allowed in the token bucket in the example above).

DSCP Marking Rules

A QoS policy can have rules added to alter the IPv4 ToS field (now used for DSCP) of every packet that enters the cloud via a port. This will affect all IPv4 packets coming into the cloud on the port with the QoS policy set. Packets exiting the cloud on the port will not be affected. At this time, the DSCP marking rules will affect all packets, meaning the last DSCP marking rule will replace the ToS field with its marking value, overwriting any previous values set (even by other DSCP marking rules). This is to say that it really only makes sense to have a single DSCP marking rule on any QoS policy.

Applying a QoS Policy to All Ports on a Network

QoS Policies can also be applied to an entire network, in which case all VM/container ports on that network will have to rules applied. Again, if a particular port ever has a separate QoS policy assigned to it (via the above command, for example), that port's specific policy will take effect instead of the network's policy. The QoS policy that will be applied will be determined at any time by: a) If the port has a policy, use that, otherwise b) if the network has a policy, use that, otherwise c) use no QoS policy. Any updates to either the port's or network's QoS policy settings will take effect immediately.

A network can be updated to have a QoS policy via the Neutron command:

```
$ neutron net-update --qos-policy <policy_id> <port_id>
```

The *policy ID* is the UUID of the created QoS policy object, and the *network ID* is the UUID of the network to be updated.

Networks can also be created with the QoS policy set:

```
$ neutron net-create <...other params...> --qos-policy <policy_id> <name>
```

Where *policy ID* is the UUID of the policy object and name is the desired name of the network. Other parameters can be specified along with *qos-policy*.

Listing Active Policies and their Rules

QoS policies and their rules can be listed, displaying their parameters.

Listing All Policies

These commands will list all policy objects.

```
$ neutron qos-policy-list
```

This will list all policies along with their rules.

Listing Rules on a Policy

These commands will list all rules for a given policy.

```
$ neutron qos-bandwidth-limit-rule-list <policy_id>
```

Where *policy ID* is the UUID of the policy object to list the rules from.

Listing the Policy Active for a Port or Network

These commands show the configured QoS policy for a port or network.

```
$ neutron port-show <port_id>
```

Where *port ID* is the UUID of the port to show the configured QoS policy. This will display the port information. The UUID of the configured QoS policy will be set under the "qos_policy_id" field. If the field is absent or not set, then no QoS policy is set on the port.

For a network, the process is the same:

Changing Rules on a Policy

Changing which rules are active on a QoS policy basically boils down to either creating or deleting rules directly. You cannot affect the rules on a policy by manipulating the policy object itself, it must be done by targeting the individual rule on a specific policy. This also means that bulk changes of rules on a policy (like deleting many rules at a time, clearing all rules, creating many rules at once) are not supported. Please see [the section called “Creating QoS Policy Rules for a Policy” \[154\]](#) and [the section called “Deleting a Rule on a Policy” \[160\]](#).

Deleting QoS Policies

QoS policy objects or individual rules can be deleted. Rules must be deleted individually from a specific QoS policy object. If a QoS policy object is deleted, all rules that belong to that object are also deleted. If a deleted QoS policy is attached to a port or network, that link is cleared and the *qos-policy-id* field of the port or network object is set to empty.

Deleting a Rule on a Policy

A single rule on a QoS policy object can be deleted to remove its effects from the policy. Once a rule is deleted from a policy, traffic through any associated ports or networks will cease to be affected by the rule’s parameters immediately after it is deleted.

```
$ neutron qos-bandwidth-limit-rule-delete <rule_id> <policy_id>
```

Or:

```
$ neutron qos-dscp-marking-rule-delete <rule_id> <policy_id>
```

Where *rule ID* is the UUID of the rule to delete, *policy ID* is the UUID of the policy the rule belongs to.

Deleting an Entire Policy

An entire QoS policy, and any attached rules, can be deleted at a single time with the following commands. Once a policy is deleted, it will immediately cease to have any effect on any network traffic. All rules created on the deleted policy will be cleaned up and removed, and all networks and ports attached to this policy will have their *qos-policy-id* fields cleared.

```
$ neutron qos-policy-delete <policy_id>
```

Where *policy ID* is the UUID of the policy to delete.



You may also choose to enable or disable the feature for specific nodes, using the `-h` switch of the `mn-conf(1)` command.

The State Proxy is enabled at all NSDB clients by default (`true`)

Advanced Options

The client features an automatic reconnection mechanism with fallback to ZooKeeper when no MidoNet Cluster nodes running the State Proxy service can be reached after a specified time interval.

When there are no servers available, or when an existing connection to a server is lost, the service enters the *soft* reconnection state. During this phase, the subscriptions to the state tables are maintained while the client tries to contact other proxy servers. To avoid flooding a server with requests, there is a configurable interval between connection attempts.

After a fixed number of unsuccessful retries, the client enters the *hard* reconnection phase during which the State Proxy feature is disabled at the NSDB client and reading the state tables falls-back to ZooKeeper. The client will periodically attempt to reconnect and if successful, the feature is re-enabled.

You can use the following configuration options to tune the client connection behavior during the *soft* and *hard* phases.

Name	Default value	description
<code>state_proxy.soft_reconnect_delay</code>	1 second	The interval between reconnect attempts while in <i>soft</i> phase.
<code>state_proxy.max_soft_reconnect_attempts</code>	50	The maximum number of reconnect attempts before deeming the service unavailable.
<code>state_proxy.hard_reconnect_delay</code>	30 seconds	The delay between reconnection attempts when the service is unavailable.

The State Proxy client uses a configurable number of threads to handle updates received from the State Proxy service. You can change the number of allocated threads using the following configuration option:

Name	Default	description
<code>state_proxy.network_threads</code>	2	The number of threads used to handle network events

State Proxy Service Configuration

The State Proxy service runs on the MidoNet Cluster nodes.

Enabling and Disabling

Use the following command to enable or disable the State Proxy service on the cluster.

```
$ echo "cluster.state_proxy.enabled : { true | false }" | mn-conf set -t default
```

You can also use the `-h` switch of the `mn-conf(1)` command to enable or disable the service on specific cluster nodes.



Important

The State Proxy service at the cluster can be enabled or disabled independently of enabling or disabling the feature on the NSDB clients, mainly the MidoNet Agents. However, if you disable the State Proxy service while keeping the feature enabled at the client, the clients will automatically fallback to using ZooKeeper.

Advanced Options

Name	Default value	Description
state_proxy.cluster.cache_threads	4	Number of threads used to handle updates for the state tables. Increase this value on large deployments if big latencies are experienced.
state_proxy.cluster.notify_batch_size	64	This option specifies the maximum number of state table entry updates that can be sent to subscribed clients in a single notification message.
state_proxy.cluster.initial_subscriber_queue_size	16	Default size for the queue of messages for a subscriber.
state_proxy.cluster.server.address	0.0.0.0	The local-interface IP address.
state_proxy.cluster.server.port	2346	The listening TCP port number.
state_proxy.cluster.server.supervisor_threads	1	Number of threads used to accept inbound connections.
state_proxy.cluster.server.worker_threads	4	The number of threads used to handle client requests.
state_proxy.cluster.server.max_pending_connections	1024	The maximum number of half-opened incoming connections (backlog).
state_proxy.cluster.server.bind_retry_interval	60	The retry interval if the service fails to bind to a local TCP port because the port is already in use. Set to 0 to disable retries.
state_proxy.cluster.server.channel_timeout	15	The timeout for a channel operation such as closing a connection.
state_proxy.cluster.server.shutdown_quiet_period	0	The interval during server shutdown when the server may continue accepting new connections or requests. If new connections or requests are received during the quiet period, the server guarantees to process them and the quiet period starts over postponing the shutdown.
state_proxy.cluster.server.shutdown_timeout	15	The shutdown timeout interval, which begins after the expiration of the quiet period. The timeout interval allows tasks that are in progress to complete before the server shuts down.



```
$ echo "agent.minions.flow_state.local_push_state : false" | mn-conf set -t default
```

Advanced Configuration Options

A complete list of the configuration options (under the `agent.minions.flow_state mn-conf` key) available for the service is as follows:

Name	Default value	Description
legacy_push_state	false	Whether the FlowState minion will save incoming flow state messages Cassandra.
local_push_state	true	Whether the FlowState minion will save incoming flow state messages to local storage on the MidoNet Agent.
port	6688	The value of the UDP and TCP ports used to listen for incoming flow state messages from the parent Agent process through the loopback interface, and flow state transfer requests.
connection_timeout	2s	Timeout for TCP requests during flow state exchanges.
block_size	262144	The size in bytes of the compressed block for the flow state storage. Changing this has an impact on the amount of memory used as well as the compression ratio achieved by the snappy algorithm. Setting a higher value implies an increased memory usage and potentially a higher compression ratio. Defaults to 256 KB.
blocks_per_port	512	The number of allowed blocks of compressed flow state per virtual port. By default, it will use up to 128 MB (512 blocks of 256 KB) per virtual port. It does not pre-allocate blocks so disk space will not be used unless it's filled with data. With this space, we can hold around 500k messages at a rate of 4k flows per second for a given port (not considering compression).
expiration_time	120s	How long should we keep flow state stored. Flow state entries older than this period of time will be eligible for removal.
expiration_delay	30s	The delay between consecutive flow state invalidation tasks that remove and clear the blocks that are older than the expiration time. This task only marks the block headers as invalid if they are expired so the overhead is minimum.
clean_unused_files_delay	12h	The delay between consecutive runs of the flow state file cleaner task. This is a background task that looks into the current list of flow state files and removes those not being used (written to or reading from). This is a background housekeeping activity to prevent storage from being used needlessly.
log_directory	flowstate	The name of the flow state log directory. The MidoNet Agent will write to this directory the records of the

Name	Default value	Description
		current flow state associated to the ports bound to this agent. This directory will be created in /var/db/midolman by default. You can change the location of the base directory by specifying it in the MIDO_DB_DIR environment variable.

28. Working with the MidoNet CLI

Table of Contents

Using the MidoNet CLI 181

You can explore and edit all of MidoNet's virtual topology through the CLI, however, you should use write operations with caution, as they are likely to create inconsistencies between MidoNet's idea of the virtual network and OpenStack's view of it.



Note

When using MidoNet with OpenStack, please be careful not to introduce inconsistencies between OpenStack and MidoNet virtual topologies.

With that warning in mind, there are certain tasks for which the CLI can be particularly useful:

- Creating an Edge Router
- Setting up the cloud's up-link using the Border Gateway Protocol (BGP)
- Upgrading MidoNet
- Registering MidoNet Agents
- Setting up an L2 gateway
- Troubleshooting problems with the MidoNet API. Because the CLI uses the MidoNet API directly, it's the easiest way to make requests to it to verify that it works.
- Viewing and exploring MidoNet's virtual topology and the status and network addresses of all hosts running the MidoNet agent

Using the MidoNet CLI

To use the MidoNet CLI you need to connect to the MidoNet host running it.

1. Using SSH to connect to the host running the MidoNet CLI.

You must know the IP address of the machine you are connecting to as well as your login credentials, that is your username and password. Example of an SSH command:

```
$ ssh root@192.168.17.5
root@192.168.17.5's password:
```

You have already provided your username, 'root' as part of the command, and now the server is prompting you for a password. Type in your password to get in.

2. The CLI is documented in a set of man pages. To view the man pages, from the system command line, enter:

```
$ man midonet-cli
```

3. To start the midonet-cli, at the system prompt, enter:

```
$ midonet-cli  
midonet>
```

The midonet> prompt that gets displayed indicates system readiness to accept MidoNet commands. Type in help and hit Enter to get a list of all available commands. You can also infer proper usage and syntax from context-aware auto-complete and by using the describe command.





MidoNet behavior after ZooKeeper cluster failure

Nodes running the MidoNet Agent, Midolman, depend on a live ZooKeeper session to load pieces of a virtual network topology on-demand and watch for updates to those virtual devices.

When ZooKeeper becomes inaccessible, a MidoNet Agent instance will continue operating for as long as there's a chance to recover connectivity while keeping the same ZooKeeper session. The amount of operating time is thus dictated by the session timeout, which you can control by editing the `zookeeper session_gracetime` setting in `mn-conf(1)`.

Once the session expires, the MidoNet Agent will give up and shut itself down, prompting upstart to re-launch it. If the ZooKeeper connection and session are recovered within the `session_gracetime`, MidoNet Agent operation will resume uneventfully. The MidoNet Agent will learn of all the updates that happened to the virtual topology while it was disconnected and will update its internal state and flow tables accordingly.

While the MidoNet Agent is running disconnected from ZooKeeper, waiting for the session to come back, traffic will still be processed, but with reduced functionality, as follows:

- The MidoNet Agent will not see updates to the virtual topology, thus packets may be processed with a version of the network topology that's up to `session_gracetime` too old.
- The MidoNet Agent will be unable to load new pieces of the network topology. Packets that traverse devices that had never been loaded on a particular MidoNet Agent will error out.
- The MidoNet Agent will not be able to perform or see updates to Address Resolution Protocol (ARP) tables and Media Access Control (MAC) learning tables.

As time passes, a disconnected MidoNet Agent will become less and less useful. The trade-offs presented above are key to choosing a sensible `session_gracetime` value; the default is 30 seconds.

ZooKeeper connectivity is not an issue for the MidoNet API server. The API requests are stateless and will simply fail when there is no ZooKeeper connectivity.

ZooKeeper configuration

You may use the ZooKeeper configuration section in `mn-conf(1)` to adjust:

- the ZooKeeper session timeout value (in milliseconds). This value determines when the system considers the connection between ZooKeeper and the MidoNet Agent to be interrupted.
- the session grace timeout value (in milliseconds). This value determines the period of time during which the Agent can reconnect to ZooKeeper without causing node outage.
- the root path for MidoNet data.

```
zookeeper {  
    zookeeper_hosts = <comma separated IPs>
```


If you notice decreased performance because packet sizes exceed the maximum buffer size, you can increase the value for the `buf_size_kb` setting. This setting controls the buffer size (in KB). Be aware that the buffer size puts a limit on the packet size that the MidoNet Agent can send. In a network that jumbo frames traverse, adjust the size so one buffer will accommodate a whole frame, plus enough room for the flow's actions.

BGP failover configuration

The default BGP fail-over time is 2-3 minutes. However, you can reduce this time by changing some parameters on both ends of the session: in `mn-conf(1)` (the MidoNet side) and the remote end BGP peer configuration. The example below shows how to reduce the BGP fail-over time to one minute on the MidoNet side:

```
agent {
    midolman {
        bgp_connect_retry=1
        bgp_holdtime=3
        bgp_keepalive=1
    }
}
```

The settings in `mn-conf` must match those on the remote end BGP peer configuration. For more information about how to set them, see [the section called “BGP failover configuration on a BGP peer”](#) [5].

Advanced MidoNet REST API configuration options

This section describes the configuration options for advanced users. The values for the following keys are modified using the MidoNet configuration tool, `mn-conf`. For more information on the MidoNet configuration, see [the section called “MidoNet Configuration: mn-conf” \[183\]](#).

Configuration options

Table 29.2. Admin Roles

Configuration Key	Default Value	Description
<code>cluster.rest_api.enabled</code>	<code>true</code>	Specifies whether the REST API service is enabled at the MidoNet cluster instance where the configuration applies.
<code>cluster.rest_api.http_port</code>	<code>8181</code>	Specifies the listening port for the HTTP end-point.
<code>cluster.rest_api.https_port</code>	<code>8443</code>	Specifies the listening port for the HTTPS end-point.

Specifying the private and public keys for HTTPS

To enable the HTTPS end-point of the MidoNet Cluster REST API service, you must configure a JKS key store containing the private and public key X.509 certificate used for encrypting such connections.

The location of the key store file and the password for the private key are specified as the following Java system properties.

Table 29.3. System Properties for the HTTPS Key Store

Property Name	Default Value	Description
midonet.keystore_path	/etc/midonet-cluster/ssl/ midonet.jks	The name of the key store file.
midonet.keystore_password	none	The password for the private key entry. If not set, the HTTPS end-point of the REST API will be disabled (default).

To change the previous properties, and enable HTTPS, you can add the corresponding property values to the environmental MidoNet Cluster script file found at `/etc/midonet-cluster/midonet-cluster-env.sh`:

```
JVM_OPTS="$JVM_OPTS -Dmidonet.keystore_path=<key-store-file>"
JVM_OPTS="$JVM_OPTS -Dmidonet.keystore_password=<key-entry-password>"
```

Generating self-signed keys

To generate a self-signed key, you can use the following procedure. Note that you will be prompted for passwords during this process, and need to keep the keystore password for later use.

```
openssl genrsa -des3 -out midonet.key 2048
openssl rsa -in midonet.key -out midonet.key
openssl req -sha256 -new -key midonet.key -out midonet.csr -subj '/CN=localhost'
openssl x509 -req -days 365 -in midonet.csr -signkey midonet.key -out midonet.crt
```

Now we will combine the private key into the cert, because we generated them separately:

```
openssl pkcs12 -inkey midonet.key -in midonet.crt -export -out midonet.pkcs12
```

And load the certificate into the keystore:

```
keytool -importkeystore -srckeystore midonet.pkcs12 -srcstoretype PKCS12 -destkeystore midonet.jks
```

Now place the keystore in the default location:

```
mv midonet.jks /etc/midonet-cluster/ssl
```

For more advanced key management, including adding your own certificate to the keystore, please refer to the following documentation:

<https://www.eclipse.org/jetty/documentation/current/configuring-ssl.html>

Cassandra cache

This section describes details relative to the Cassandra-based implementation of the Cassandra cache used for connection tracking and NAT mappings.

Client library and normal operation

Cassandra operations are now asynchronous, so losing connectivity to Cassandra should not impact simulations.



zookeeper-session_timeout

Sets the timeout value (in milliseconds) after which ZooKeeper considers clients to be disconnected from the ZooKeeper server:

```
<context-param>  
  <param-name>zookeeper-session_timeout</param-name>  
  <param-value>30000</param-value>  
</context-param>
```

30. MidoNet and OpenStack TCP/UDP service ports

Table of Contents

Services on the Controller node	196
Services on the Network State Database nodes	197
Services on the Compute nodes	198
Services on the Gateway Nodes	198

This section lists the TCP/UDP ports used by services in MidoNet and OpenStack.

Services on the Controller node

This section lists the TCP/UDP ports used by the services on the Controller node.

Category	Service	Prot ocol	Port	Self	Controller	Compute	Mgmt. PC
OpenStack	glance-api	TCP	9292	x	x	x	x
OpenStack	httpd (Hori- zon)	TCP	80	x			x
MidoNet	midonet-api	TCP	8181	x	x		x
OpenStack	swift-ob- ject-server	TCP	6000	x	x	x	
OpenStack	swift-contain- er-server	TCP	6001	x	x	x	
OpenStack	swift-ac- count-server	TCP	6002	x	x	x	
OpenStack	keystone	TCP	35357	x	x	x	x
OpenStack	neutron-serv- er	TCP	9696	x	x	x	x
OpenStack	nova-novnc- proxy	TCP	6080	x	x		x
OpenStack	heat-api	TCP	8004	x	x		x
OpenStack	nova-api	TCP	8773	x	x		x
Tomcat	Tomcat shut- down control channel	TCP	8005	x	x		
OpenStack	nova-api	TCP	8774	x	x	x	x
OpenStack	nova-api	TCP	8775	x	x	x	x
OpenStack	glance-reg- istry	TCP	9191	x	x	x	
OpenStack	qpidd	TCP	5672	x	x	x	
OpenStack	keystone	TCP	5000	x	x	x	x
OpenStack	cinder-api	TCP	8776	x	x	x	x
Tomcat	Tomcat man- agement port (not used)	TCP	8009	x	x		
OpenStack	ceilome- ter-api	TCP	8777	x	x	x	x

Category	Service	Prot ocol	Port	Self	Controller	Compute	Mgmt. PC
OpenStack	mongod (ceilometer)	TCP	27017	x	x	x	
OpenStack	MySQL	TCP	3306	x	x	x	

Services on the Network State Database nodes

This section lists the TCP/UDP ports used by the services on the Network State Database nodes.

Category	Service	Prot ocol	Port	Self	Controller	NSDB	Compute	Comment
MidoNet	ZooKeeper communication	TCP	3888	x		x		
MidoNet	ZooKeeper leader	TCP	2888	x		x		
MidoNet	ZooKeeper/Cassandra	TCP	random	x				ZooKeeper/Cassandra "LISTEN" to a TCP high number port. The port number is randomly selected on each ZooKeeper/Cassandra host.
MidoNet	Cassandra Query Language (CQL) native transport port	TCP	9042					
MidoNet	Cassandra cluster communication	TCP	7000	x		x		
MidoNet	Cassandra cluster communication (Transport Layer Security (TLS) support)	TCP	7001	x		x		
MidoNet	Cassandra JMX	TCP	7199	x				JMX monitoring port. If you're using this port to monitor Cassandra health, enable communication from the monitoring server.
MidoNet	ZooKeeper client	TCP	2181	x	x	x	x	
MidoNet	Cassandra clients	TCP	9160	x	x	x	x	

Services on the Compute nodes

This section lists the TCP/UDP ports used by the services on the Compute nodes.

Category	Service	Protocol	Port	Self	Controller	Comment
OpenStack	qemu-kvm vnc	TCP	From 5900 to 5900 + # of VM		x	
MidoNet	midolman	TCP	random	x		midolman "LISTEN"s to a TCP high number port. The port number is randomly selected on each MN Agent host.
OpenStack	libvirtd	TCP	16509	x	x	
MidoNet	midolman	TCP	7200	x		JMX monitoring port If you're using this port to monitor health, enable communication from the monitoring server.
MidoNet	midolman	TCP	9697	x		If enabled, MidoNet Metadata Proxy listens on 169.254.169.254:9697 to accept metadata requests.

Services on the Gateway Nodes

This section lists the TCP/UDP ports used by the services on the Gateway Nodes.

Category	Service	Protocol	Port	Self	Misc.	Comment
MidoNet	midolman	TCP	random	x		midolman LISTEN"s to a TCP high number port. The port number is randomly selected on each MN Agent host.
MidoNet	midolman	TCP	7200	x	x	JMX monitoring port If you're using this port to monitor health, enable communication from the monitoring server.
MidoNet	quagga bgpd control	TCP	2606	x		Network-NameSpace mbgp[Peer Number]_ns
MidoNet	quagga bgpd bgp	TCP	179		BGP neighbor	Network-NameSpace mbgp[Peer Number]_ns




```
# cp /var/lib/zookeeper/data/myid /tmp/zk_backup/config/
```

RHEL:

```
# cp /etc/zookeeper/zoo.cfg /tmp/zk_backup/config/
# cp /var/lib/zookeeper/myid /tmp/zk_backup/config/
```

3. Backup ZooKeeper Data

On all NSDB nodes, make a copy of the ZooKeeper data:

Ubuntu:

```
# cp /var/lib/zookeeper/version-2/* /tmp/zk_backup/data/
```

RHEL:

```
# cp /var/lib/zookeeper/data/version-2/* /tmp/zk_backup/data/
```

4. Create MySQL Backup

Create the MySQL backup (as described in [the section called “MySQL / MariaDB” \[204\]](#)) to ensure that ZooKeeper and MySQL data are in sync.

5. Reset MidoNet API to READWRITE

Reset the MidoNet API to READWRITE:

```
# midonet-cli
midonet> set system-state availability READWRITE
```

Restore

1. Stop ZooKeeper

On all NSDB nodes, stop the ZooKeeper service:

```
# service zookeeper stop
```

```
# systemctl stop zookeeper.service
```

2. Restore ZooKeeper Configuration

On all NSDB nodes, restore the ZooKeeper configuration from your backup:

Ubuntu:

```
# cp /tmp/zk_backup/config/zoo.cfg /etc/zookeeper/conf/zoo.cfg
# cp /tmp/zk_backup/config/myid /var/lib/zookeeper/data/myid
```

RHEL:

```
# cp /tmp/zk_backup/config/zoo.cfg /etc/zookeeper/zoo.cfg
# cp /tmp/zk_backup/config/myid /var/lib/zookeeper/myid
```

3. Restore ZooKeeper Data

On all NSDB nodes, restore the ZooKeeper data from your backup:

Ubuntu:

```
# rm /var/lib/zookeeper/version-2/*
# cp /tmp/zk_backup/data/* /var/lib/zookeeper/version-2/
```

RHEL:

```
# rm /var/lib/zookeeper/data/version-2/*  
# cp /tmp/zk_backup/data/* /var/lib/zookeeper/data/version-2/
```

4. Restart ZooKeeper

On all NSDB nodes, start the ZooKeeper service:

Ubuntu:

```
# service zookeeper start
```

RHEL:

```
# systemctl start zookeeper.service
```

5. Verify ZooKeeper Operation

On all NSDB nodes, verify that ZooKeeper is operating properly.

A basic check can be done by executing the `ruok` (Are you ok?) command. This will reply with `imok` (I am ok.) if the server is running in a non-error state:

```
# echo ruok | nc 127.0.0.1 2181  
imok
```

More detailed information can be requested with the `stat` command, which lists statistics about performance and connected clients:

```
# echo stat | nc 127.0.0.1 2181  
Zookeeper version: 3.4.5--1, built on 06/10/2013 17:26 GMT  
Clients:  
 /127.0.0.1:34768[0](queued=0,recved=1,sent=0)  
 /192.0.2.1:49703[1](queued=0,recved=1053,sent=1053)  
  
Latency min/avg/max: 0/4/255  
Received: 1055  
Sent: 1054  
Connections: 2  
Outstanding: 0  
Zxid: 0x260000013d  
Mode: follower  
Node count: 3647
```

Cassandra

MidoNet exchanges connection tracking and NAT information directly between its Agents, but also stores this data in Cassandra as a backup.

Whilst normal packet processing does not require access to Cassandra, it is necessary to support port migrations and to keep connections working across Agent reboots. Agents pull this information from Cassandra anytime they bind a new local interface to the virtual topology.

If you do not rely on above, a full data backup of Cassandra is not required.

Backup

1. Stop Cassandra

RHEL:

```
# systemctl restart midolman.service
```

MidoNet Cluster

Backup

1. Backup Cluster Configuration

On the Cluster nodes, make a copy of the Cluster configuration:

```
# cp /etc/midonet/midonet.conf /tmp/mn-cluster_backup/config/
```

Restore

1. Restore Cluster Configuration

On the Cluster nodes, restore the Cluster configuration from your backup:

```
# cp /tmp/mn-cluster_backup/config/* /etc/midonet/
```

2. Restart Cluster

On the Cluster nodes, restart the Cluster service:

Ubuntu:

```
# service midonet-cluster restart
```

RHEL:

```
# systemctl restart midonet-cluster.service
```

MySQL / MariaDB

Technically MidoNet only relies on the Neutron and Keystone tables, but for completeness and ease of use a full backup of the database is recommended.

Backup



Important

To ensure that Neutron's MySQL database and MidoNet's ZooKeeper data are in sync, set the MidoNet API to read-only before starting the backup.

Backup both MySQL database and ZooKeeper data (as described in [the section called "ZooKeeper" \[199\]](#)) at the same time!

1. Set MidoNet API to READONLY

Set the MidoNet API to READONLY to avoid topology changes during the backup process:

```
# midonet-cli
midonet> set system-state availability READONLY
```

2. Create Database Dump

Use the `mysqldump` command to create a database backup:

```
# mysqldump \
  -uroot -p*_password*_ \
  --lock-tables --events --routines --triggers --all-databases | \
  gzip > /tmp/ost_backup/data/ost_database.sql.gz
```

3. Create ZooKeeper Backup

Create the ZooKeeper backup (as described in [the section called “ZooKeeper” \[199\]](#)) to ensure that MySQL and ZooKeeper data are in sync.

4. Reset MidoNet API to READWRITE

Reset the MidoNet API to READWRITE:

```
# midonet-cli
midonet> set system-state availability READWRITE
```

Restore

1. Stop Affected Services

Stop all services that access the database.

This includes, but is **not limited** to the following services (depending on your OpenStack installation):

Ubuntu:

```
# service apache2 stop
# service keystone stop
# service neutron-server stop
[...]
```

RHEL:

```
# systemctl stop openstack-keystone.service
# systemctl stop neutron-server.service
# systemctl stop httpd.service
[...]
```

2. Restore Database Dump

```
# gunzip -c /tmp/ost_backup/data/ost_database.sql.gz | mysql -uroot -
p*_password*_
```

3. Restart Services

Restart the services.

Ubuntu:

```
# service apache2 start
# service keystone start
# service neutron-server start
[...]
```

RHEL:

```
# systemctl start openstack-keystone.service
```

```
# systemctl start neutron-server.service
# systemctl start httpd.service
[...]
```

Neutron

Backup

1. Backup Neutron Configuration

Make a copy of the Neutron configuration:

```
# cp /etc/neutron/neutron.conf /tmp/neutron_backup/config/
# cp /etc/neutron/plugins/midonet/midonet.ini /tmp/neutron_backup/
config/
```

Restore

1. Restore Neutron Configuration

Restore the Neutron configuration from your backup:

```
# cp /tmp/neutron_backup/config/neutron.conf /etc/neutron/
# cp /tmp/neutron_backup/config/midonet.ini /etc/neutron/plugins/
midonet/
```

2. Restart Neutron

Restart the Neutron service.

Ubuntu:

```
# service neutron-server restart
```

RHEL:

```
# systemctl restart neutron-server.service
```

Libvirt

Backup

1. Backup Libvirt Configuration

On the Compute nodes, make a copy of the Libvirt configuration:

```
# cp /etc/libvirt/qemu.conf /tmp/libvirt_backup/
```

Restore

1. Restore Libvirt Configuration

Restore the Libvirt configuration from your backup:

```
# cp /tmp/libvirt_backup/qemu.conf /etc/libvirt/
```

2. Restart Libvirt

Restart the Libvirt service.

Ubuntu:

```
# service libvirt-bin restart
```

RHEL:

```
# systemctl restart libvirtd.service
```