
Open vSwitch

Release 2.11.90

Feb 12, 2019

1	Open vSwitch Documentation	1
1.1	How the Documentation is Organised	1
1.2	First Steps	1
1.3	Deeper Dive	2
1.4	The Open vSwitch Project	2
1.5	Getting Help	2
2	Getting Started	3
2.1	What Is Open vSwitch?	4
2.1.1	Overview	4
2.1.2	What's here?	5
2.2	Why Open vSwitch?	5
2.2.1	The mobility of state	6
2.2.2	Responding to network dynamics	6
2.2.3	Maintenance of logical tags	6
2.2.4	Hardware integration	6
2.2.5	Summary	7
2.3	Installing Open vSwitch	7
2.3.1	Installation from Source	7
2.3.2	Installation from Packages	41
2.3.3	Upgrades	49
2.3.4	Others	50
3	Tutorials	53
3.1	OVS Faucet Tutorial	53
3.1.1	Setting Up OVS	53
3.1.2	Setting up Faucet	54
3.1.3	Overview	55
3.1.4	Switching	56
3.1.5	Routing	66
3.1.6	ACLs	73
3.1.7	Finishing Up	75
3.1.8	Further Directions	75
3.2	OVS IPsec Tutorial	76
3.2.1	Requirements	76
3.2.2	Installing OVS and IPsec Packages	76
3.2.3	Configuring IPsec tunnel	77

3.2.4	Troubleshooting	80
3.2.5	Bug Reporting	81
3.3	Open vSwitch Advanced Features	81
3.3.1	Getting Started	81
3.3.2	Using GDB	82
3.3.3	Motivation	83
3.3.4	Scenario	83
3.3.5	Setup	84
3.3.6	Implementing Table 0: Admission control	85
3.3.7	Testing Table 0	85
3.3.8	Implementing Table 1: VLAN Input Processing	86
3.3.9	Testing Table 1	87
3.3.10	Implementing Table 2: MAC+VLAN Learning for Ingress Port	88
3.3.11	Testing Table 2	89
3.3.12	Implementing Table 3: Look Up Destination Port	90
3.3.13	Testing Table 3	91
3.3.14	Implementing Table 4: Output Processing	93
3.3.15	Testing Table 4	94
3.4	OVN Sandbox	96
3.4.1	Getting Started	96
3.4.2	Using GDB	96
3.4.3	Creating OVN Resources	96
3.4.4	Using ovn-trace	97
3.5	OVN OpenStack Tutorial	98
3.5.1	Setting Up DevStack	98
3.5.2	DevStack preliminaries	101
3.5.3	Shortening UUIDs	101
3.5.4	Overview	102
3.5.5	Switching	103
3.5.6	Routing	114
3.5.7	Adding a Gateway	118
3.5.8	IPv6	122
3.5.9	ACLs	125
3.5.10	DHCP	127
3.5.11	Further Directions	129
3.6	OVN Role-Based Access Control (RBAC) Tutorial	129
3.6.1	Generating Certificates and Keys	129
3.6.2	Configuring RBAC	130
3.7	OVN IPsec Tutorial	131
3.7.1	Generating Certificates and Keys	131
3.7.2	Configuring OVN IPsec	131
3.7.3	Enabling OVN IPsec	131
3.7.4	Troubleshooting	132
3.7.5	Bug Reporting	132
3.8	OVS Conntrack Tutorial	133
3.8.1	Definitions	133
3.8.2	Conntrack Related Fields	133
3.8.3	Sample Topology	134
3.8.4	Tool used to generate TCP segments	136
3.8.5	Matching TCP packets	136
3.8.6	Summary	140
4	Deep Dive	141
4.1	OVS	141

4.1.1	Design Decisions In Open vSwitch	141
4.1.2	Open vSwitch Datapath Development Guide	156
4.1.3	Integration Guide for Centralized Control	160
4.1.4	Porting Open vSwitch to New Software or Hardware	163
4.1.5	OpenFlow Support in Open vSwitch	167
4.1.6	Bonding	171
4.1.7	Open vSwitch Networking Namespaces on Linux	174
4.1.8	OVSDB Replication Implementation	175
4.1.9	The DPDK Datapath	177
4.1.10	OVS-on-Hyper-V Design	200
4.1.11	Language Bindings	207
4.1.12	Testing	208
4.1.13	Tracing packets inside Open vSwitch	215
4.1.14	C IDL Compound Indexes	216
4.2	OVN	220
4.2.1	OVN Gateway High Availability Plan	220
4.2.2	Role Based Access Control	226
4.2.3	What's New with OVS and OVN 2.8	227
5	How-to Guides	233
5.1	OVS	233
5.1.1	Open vSwitch with KVM	233
5.1.2	Encrypt Open vSwitch Tunnels with IPsec	234
5.1.3	Open vSwitch with SELinux	237
5.1.4	Open vSwitch with Libvirt	239
5.1.5	Open vSwitch with SSL	240
5.1.6	Using LISP tunneling	245
5.1.7	Connecting VMs Using Tunnels	246
5.1.8	Connecting VMs Using Tunnels (Userspace)	248
5.1.9	Isolating VM Traffic Using VLANs	252
5.1.10	Quality of Service (QoS) Rate Limiting	254
5.1.11	How to Use the VTEP Emulator	257
5.1.12	Monitoring VM Traffic Using sFlow	260
5.1.13	Using Open vSwitch with DPDK	263
5.2	OVN	268
5.2.1	Open Virtual Networking With Docker	268
5.2.2	Integration of Containers with OVN and OpenStack	273
5.2.3	Open Virtual Network With firewallD	274
6	Reference Guide	277
6.1	Man Pages	277
6.1.1	ovs-sim	277
6.1.2	ovs-test	280
6.1.3	ovs-vlan-test	282
6.1.4	ovsdb-server	284
6.1.5	ovsdb	290
6.1.6	ovsdb	294
7	Open vSwitch FAQ	303
7.1	Basic Configuration	303
7.2	Development	306
7.3	Implementation Details	307
7.4	General	309
7.5	Common Configuration Issues	310

7.6	Using OpenFlow	316
7.7	Quality of Service (QoS)	323
7.8	Releases	325
7.9	Terminology	329
7.10	VLANs	330
7.11	VXLANs	333
7.12	OVN	334
8	Open vSwitch Internals	335
8.1	Contributing to Open vSwitch	335
8.1.1	Submitting Patches	335
8.1.2	Backporting patches	341
8.1.3	Open vSwitch Coding Style	344
8.1.4	Open vSwitch Windows Datapath Coding Style	352
8.1.5	Open vSwitch Documentation Style	355
8.1.6	Open vSwitch Library ABI Updates	361
8.2	Mailing Lists	362
8.2.1	ovs-announce	362
8.2.2	ovs-discuss	363
8.2.3	ovs-dev	363
8.2.4	ovs-git	363
8.2.5	ovs-build	363
8.2.6	bugs	363
8.2.7	security	363
8.3	Patchwork	363
8.3.1	git-pw	363
8.3.2	pwclient	364
8.4	Open vSwitch Release Process	364
8.4.1	Release Strategy	364
8.4.2	Release Numbering	365
8.4.3	Release Scheduling	365
8.4.4	Contact	365
8.5	Reporting Bugs in Open vSwitch	365
8.6	Open vSwitch's Security Process	366
8.6.1	What is a vulnerability?	366
8.6.2	Step 1: Reception	367
8.6.3	Step 2: Assessment	367
8.6.4	Step 3a: Document	367
8.6.5	Step 3b: Fix	369
8.6.6	Step 4: Embargoed Disclosure	369
8.6.7	Step 5: Public Disclosure	370
8.7	The Linux Foundation Open vSwitch Project Charter	370
8.8	Emeritus Status for OVS Committers	372
8.9	Expectations for Developers with Open vSwitch Repo Access	373
8.9.1	Pre-requisites	373
8.9.2	Review	373
8.9.3	Git conventions	373
8.10	OVS Committer Grant/Revocation Policy	374
8.10.1	Granting Commit Access	374
8.10.2	Revoking Commit Access	375
8.10.3	Changing the Policy	376
8.10.4	Nomination to Grant Commit Access	376
8.10.5	Vote to Grant Commit Access	376
8.10.6	Vote Results for Grant of Commit Access	376

8.10.7	Invitation to Accepted Committer	377
8.10.8	Proposal to Revoke Commit Access for Detrimental Behavior	377
8.10.9	Vote to Revoke Commit Access	377
8.10.10	Vote Results for Revocation of Commit Access	377
8.10.11	Notification of Commit Revocation for Detrimental Behavior	378
8.11	Authors	378
8.12	Committers	391
8.13	How Open vSwitch's Documentation Works	391
8.13.1	reStructuredText and Sphinx	391
8.13.2	ovs-sphinx-theme	392
8.13.3	Read the Docs	392
8.13.4	openvswitch.org	392

1.1 How the Documentation is Organised

The Open vSwitch documentation is organised into multiple sections:

- *Installation guides* guide you through installing Open vSwitch (OVS) and Open Virtual Network (OVN) on a variety of different platforms
- *Tutorials* take you through a series of steps to configure OVS and OVN in sandboxed environments
- *Topic guides* provide a high level overview of OVS and OVN internals and operation
- *How-to guides* are recipes or use-cases for OVS and OVN. They are more advanced than the tutorials.
- *Frequently Asked Questions* provide general insight into a variety of topics related to configuration and operation of OVS and OVN.

1.2 First Steps

Getting started with Open vSwitch (OVS) or Open Virtual Network (OVN) for Open vSwitch? Start here.

- **Overview:** *What Is Open vSwitch? | Why Open vSwitch?*
- **Install:** *Open vSwitch on Linux, FreeBSD and NetBSD | Open vSwitch without Kernel Support | Open vSwitch on NetBSD | Open vSwitch on Windows | Open vSwitch on Citrix XenServer | Open vSwitch with DPDK | Installation FAQs*
- **Tutorials:** *OVS Faucet Tutorial | Open vSwitch Advanced Features | OVN Sandbox | OVN OpenStack Tutorial | OVS Conntrack Tutorial | OVS IPsec Tutorial | OVN IPsec Tutorial | OVN Role-Based Access Control (RBAC) Tutorial*

1.3 Deeper Dive

- **Architecture** *Design Decisions In Open vSwitch | OpenFlow Support in Open vSwitch | Integration Guide for Centralized Control | Porting Open vSwitch to New Software or Hardware*
- **DPDK** *Using Open vSwitch with DPDK | DPDK vHost User Ports*
- **Windows** *OVS-on-Hyper-V Design*
- **Integrations:** *Language Bindings*
- **Reference Guides:** *Reference Guide*
- **Testing** *Testing*
- **Packaging:** *Debian Packaging for Open vSwitch | RHEL 5.6, 6.x Packaging for Open vSwitch | Fedora, RHEL 7.x Packaging for Open vSwitch*

1.4 The Open vSwitch Project

Learn more about the Open vSwitch project and about how you can contribute:

- **Community:** *Open vSwitch Release Process | Authors | Mailing Lists | Patchwork | Reporting Bugs in Open vSwitch | Open vSwitch's Security Process*
- **Contributing:** *Submitting Patches | Backporting patches | Open vSwitch Coding Style | Open vSwitch Windows Datapath Coding Style*
- **Maintaining:** *The Linux Foundation Open vSwitch Project Charter | Committers | Expectations for Developers with Open vSwitch Repo Access | OVS Committer Grant/Revocation Policy | Emeritus Status for OVS Committers*
- **Documentation:** *Open vSwitch Documentation Style | Building Open vSwitch Documentation | How Open vSwitch's Documentation Works*

1.5 Getting Help

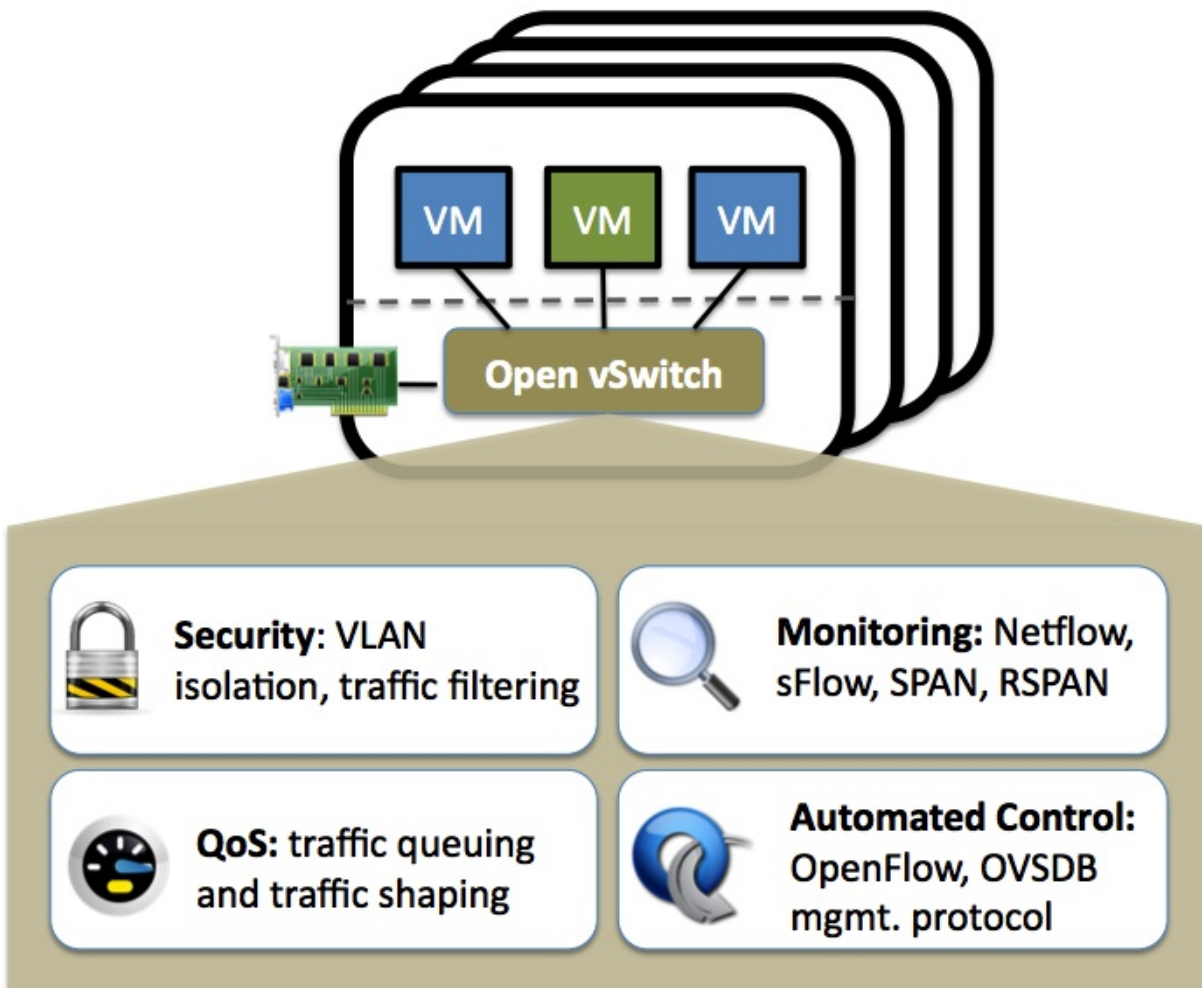
- Seeing an issue of potential bug? Report problems to bugs@openvswitch.org
- Looking for specific information? Try the [genindex](#), [modindex](#) or the *detailed table of contents*.

CHAPTER 2

Getting Started

How to get started with Open vSwitch.

2.1 What Is Open vSwitch?



2.1.1 Overview

Open vSwitch is a multilayer software switch licensed under the open source Apache 2 license. Our goal is to implement a production quality switch platform that supports standard management interfaces and opens the forwarding functions to programmatic extension and control.

Open vSwitch is well suited to function as a virtual switch in VM environments. In addition to exposing standard control and visibility interfaces to the virtual networking layer, it was designed to support distribution across multiple physical servers. Open vSwitch supports multiple Linux-based virtualization technologies including Xen/XenServer, KVM, and VirtualBox.

The bulk of the code is written in platform-independent C and is easily ported to other environments. The current release of Open vSwitch supports the following features:

- Standard 802.1Q VLAN model with trunk and access ports
- NIC bonding with or without LACP on upstream switch

- NetFlow, sFlow(R), and mirroring for increased visibility
- QoS (Quality of Service) configuration, plus policing
- Geneve, GRE, VXLAN, STT, and LISP tunneling
- 802.1ag connectivity fault management
- OpenFlow 1.0 plus numerous extensions
- Transactional configuration database with C and Python bindings
- High-performance forwarding using a Linux kernel module

The included Linux kernel module supports Linux 3.10 and up.

Open vSwitch can also operate entirely in userspace without assistance from a kernel module. This userspace implementation should be easier to port than the kernel-based switch. OVS in userspace can access Linux or DPDK devices. Note Open vSwitch with userspace datapath and non DPDK devices is considered experimental and comes with a cost in performance.

2.1.2 What's here?

The main components of this distribution are:

- `ovs-vswitchd`, a daemon that implements the switch, along with a companion Linux kernel module for flow-based switching.
- `ovsdb-server`, a lightweight database server that `ovs-vswitchd` queries to obtain its configuration.
- `ovs-dpctl`, a tool for configuring the switch kernel module.
- Scripts and specs for building RPMs for Citrix XenServer and Red Hat Enterprise Linux. The XenServer RPMs allow Open vSwitch to be installed on a Citrix XenServer host as a drop-in replacement for its switch, with additional functionality.
- `ovs-vsctl`, a utility for querying and updating the configuration of `ovs-vswitchd`.
- `ovs-appctl`, a utility that sends commands to running Open vSwitch daemons.

Open vSwitch also provides some tools:

- `ovs-ofctl`, a utility for querying and controlling OpenFlow switches and controllers.
- `ovs-pki`, a utility for creating and managing the public-key infrastructure

2.2 Why Open vSwitch?

Hypervisors need the ability to bridge traffic between VMs and with the outside world. On Linux-based hypervisors, this used to mean using the built-in L2 switch (the Linux bridge), which is fast and reliable. So, it is reasonable to ask why Open vSwitch is used.

The answer is that Open vSwitch is targeted at multi-server virtualization deployments, a landscape for which the previous stack is not well suited. These environments are often characterized by highly dynamic end-points, the maintenance of logical abstractions, and (sometimes) integration with or offloading to special purpose switching hardware.

The following characteristics and design considerations help Open vSwitch cope with the above requirements.

2.2.1 The mobility of state

All network state associated with a network entity (say a virtual machine) should be easily identifiable and migratable between different hosts. This may include traditional “soft state” (such as an entry in an L2 learning table), L3 forwarding state, policy routing state, ACLs, QoS policy, monitoring configuration (e.g. NetFlow, IPFIX, sFlow), etc.

Open vSwitch has support for both configuring and migrating both slow (configuration) and fast network state between instances. For example, if a VM migrates between end-hosts, it is possible to not only migrate associated configuration (SPAN rules, ACLs, QoS) but any live network state (including, for example, existing state which may be difficult to reconstruct). Further, Open vSwitch state is typed and backed by a real data-model allowing for the development of structured automation systems.

2.2.2 Responding to network dynamics

Virtual environments are often characterized by high-rates of change. VMs coming and going, VMs moving backwards and forwards in time, changes to the logical network environments, and so forth.

Open vSwitch supports a number of features that allow a network control system to respond and adapt as the environment changes. This includes simple accounting and visibility support such as NetFlow, IPFIX, and sFlow. But perhaps more useful, Open vSwitch supports a network state database (OVSDDB) that supports remote triggers. Therefore, a piece of orchestration software can “watch” various aspects of the network and respond if/when they change. This is used heavily today, for example, to respond to and track VM migrations.

Open vSwitch also supports OpenFlow as a method of exporting remote access to control traffic. There are a number of uses for this including global network discovery through inspection of discovery or link-state traffic (e.g. LLDP, CDP, OSPF, etc.).

2.2.3 Maintenance of logical tags

Distributed virtual switches (such as VMware vDS and Cisco’s Nexus 1000V) often maintain logical context within the network through appending or manipulating tags in network packets. This can be used to uniquely identify a VM (in a manner resistant to hardware spoofing), or to hold some other context that is only relevant in the logical domain. Much of the problem of building a distributed virtual switch is to efficiently and correctly manage these tags.

Open vSwitch includes multiple methods for specifying and maintaining tagging rules, all of which are accessible to a remote process for orchestration. Further, in many cases these tagging rules are stored in an optimized form so they don’t have to be coupled with a heavyweight network device. This allows, for example, thousands of tagging or address remapping rules to be configured, changed, and migrated.

In a similar vein, Open vSwitch supports a GRE implementation that can handle thousands of simultaneous GRE tunnels and supports remote configuration for tunnel creation, configuration, and tear-down. This, for example, can be used to connect private VM networks in different data centers.

2.2.4 Hardware integration

Open vSwitch’s forwarding path (the in-kernel datapath) is designed to be amenable to “offloading” packet processing to hardware chipsets, whether housed in a classic hardware switch chassis or in an end-host NIC. This allows for the Open vSwitch control path to be able to both control a pure software implementation or a hardware switch.

There are many ongoing efforts to port Open vSwitch to hardware chipsets. These include multiple merchant silicon chipsets (Broadcom and Marvell), as well as a number of vendor-specific platforms. The “Porting” section in the documentation discusses how one would go about making such a port.

The advantage of hardware integration is not only performance within virtualized environments. If physical switches also expose the Open vSwitch control abstractions, both bare-metal and virtualized hosting environments can be managed using the same mechanism for automated network control.

2.2.5 Summary

In many ways, Open vSwitch targets a different point in the design space than previous hypervisor networking stacks, focusing on the need for automated and dynamic network control in large-scale Linux-based virtualization environments.

The goal with Open vSwitch is to keep the in-kernel code as small as possible (as is necessary for performance) and to re-use existing subsystems when applicable (for example Open vSwitch uses the existing QoS stack). As of Linux 3.3, Open vSwitch is included as a part of the kernel and packaging for the userspace utilities are available on most popular distributions.

2.3 Installing Open vSwitch

A collection of guides detailing how to install Open vSwitch in a variety of different environments and using different configurations.

2.3.1 Installation from Source

Open vSwitch on Linux, FreeBSD and NetBSD

This document describes how to build and install Open vSwitch on a generic Linux, FreeBSD, or NetBSD host. For specifics around installation on a specific platform, refer to one of the other installation guides listed in *Installing Open vSwitch*.

Obtaining Open vSwitch Sources

The canonical location for Open vSwitch source code is its Git repository, which you can clone into a directory named “ovs” with:

```
$ git clone https://github.com/openvswitch/ovs.git
```

Cloning the repository leaves the “master” branch initially checked out. This is the right branch for general development. If, on the other hand, if you want to build a particular released version, you can check it out by running a command such as the following from the “ovs” directory:

```
$ git checkout v2.7.0
```

The repository also has a branch for each release series. For example, to obtain the latest fixes in the Open vSwitch 2.7.x release series, which might include bug fixes that have not yet been in any released version, you can check it out from the “ovs” directory with:

```
$ git checkout origin/branch-2.7
```

If you do not want to use Git, you can also obtain tarballs for Open vSwitch release versions via <http://openvswitch.org/download/>, or download a ZIP file for any snapshot from the web interface at <https://github.com/openvswitch/ovs>.

Build Requirements

To compile the userspace programs in the Open vSwitch distribution, you will need the following software:

- GNU make
- A C compiler, such as:
 - GCC 4.6 or later.
 - Clang 3.4 or later.
 - MSVC 2013. Refer to *Open vSwitch on Windows* for additional Windows build instructions.

While OVS may be compatible with other compilers, optimal support for atomic operations may be missing, making OVS very slow (see `lib/ovs-atomic.h`).

- libssl, from OpenSSL, is optional but recommended if you plan to connect the Open vSwitch to an OpenFlow controller. libssl is required to establish confidentiality and authenticity in the connections from an Open vSwitch to an OpenFlow controller. If libssl is installed, then Open vSwitch will automatically build with support for it.
- libcap-ng, written by Steve Grubb, is optional but recommended. It is required to run OVS daemons as a non-root user with dropped root privileges. If libcap-ng is installed, then Open vSwitch will automatically build with support for it.
- Python 2.7. You must also have the Python `six` library version 1.4.0 or later.
- Unbound library, from <http://www.unbound.net>, is optional but recommended if you want to enable `ovs-vswitchd` and other utilities to use DNS names when specifying OpenFlow and OVSDDB remotes. If unbound library is already installed, then Open vSwitch will automatically build with support for it. The environment variable `OVS_RESOLV_CONF` can be used to specify DNS server configuration file (the default file on Linux is `/etc/resolv.conf`).

On Linux, you may choose to compile the kernel module that comes with the Open vSwitch distribution or to use the kernel module built into the Linux kernel (version 3.3 or later). See the *Open vSwitch FAQ* question “What features are not available in the Open vSwitch kernel datapath that ships as part of the upstream Linux kernel?” for more information on this trade-off. You may also use the userspace-only implementation, at some cost in features and performance. Refer to *Open vSwitch without Kernel Support* for details.

To compile the kernel module on Linux, you must also install the following:

- A supported Linux kernel version.

For optional support of ingress policing, you must enable kernel configuration options `NET_CLS_BASIC`, `NET_SCH_INGRESS`, and `NET_ACT_POLICE`, either built-in or as modules. `NET_CLS_POLICE` is obsolete and not needed.)

On kernels before 3.11, the `ip_gre` module, for GRE tunnels over IP (`NET_IPGRE`), must not be loaded or compiled in.

To configure HTB or HFSC quality of service with Open vSwitch, you must enable the respective configuration options.

To use Open vSwitch support for TAP devices, you must enable `CONFIG_TUN`.

- To build a kernel module, you need the same version of GCC that was used to build that kernel.
- A kernel build directory corresponding to the Linux kernel image the module is to run on. Under Debian and Ubuntu, for example, each `linux-image` package containing a kernel binary has a corresponding `linux-headers` package with the required build infrastructure.

If you are working from a Git tree or snapshot (instead of from a distribution tarball), or if you modify the Open vSwitch build system or the database schema, you will also need the following software:

- Autoconf version 2.63 or later.
- Automake version 1.10 or later.
- libtool version 2.4 or later. (Older versions might work too.)

The datapath tests for userspace and Linux datapaths also rely upon:

- pyftplib. Version 1.2.0 is known to work. Earlier versions should also work.
- GNU wget. Version 1.16 is known to work. Earlier versions should also work.
- netcat. Several common implementations are known to work.
- curl. Version 7.47.0 is known to work. Earlier versions should also work.
- tftpy. Version 0.6.2 is known to work. Earlier versions should also work.
- netstat. Available from various distro specific packages

The `ovs-vswitchd.conf.db(5)` manpage will include an E-R diagram, in formats other than plain text, only if you have the following:

- dot from graphviz (<http://www.graphviz.org/>).

If you are going to extensively modify Open vSwitch, consider installing the following to obtain better warnings:

- “sparse” version 0.5.1 or later (<https://git.kernel.org/pub/scm/devel/sparse/sparse.git/>).
- GNU make.
- clang, version 3.4 or later
- flake8 along with the hacking flake8 plugin (for Python code). The automatic flake8 check that runs against Python code has some warnings enabled that come from the “hacking” flake8 plugin. If it’s not installed, the warnings just won’t occur until it’s run on a system with “hacking” installed.

You may find the `ovs-dev` script found in `utilities/ovs-dev.py` useful.

Installation Requirements

The machine you build Open vSwitch on may not be the one you run it on. To simply install and run Open vSwitch you require the following software:

- Shared libraries compatible with those used for the build.
- On Linux, if you want to use the kernel-based datapath (which is the most common use case), then a kernel with a compatible kernel module. This can be a kernel module built with Open vSwitch (e.g. in the previous step), or the kernel module that accompanies Linux 3.3 and later. Open vSwitch features and performance can vary based on the module and the kernel. Refer to [Releases](#) for more information.
- For optional support of ingress policing on Linux, the “tc” program from iproute2 (part of all major distributions and available at <https://wiki.linuxfoundation.org/networking/iproute2>).
- Python 2.7. You must also have the Python six library version 1.4.0 or later.

On Linux you should ensure that `/dev/urandom` exists. To support TAP devices, you must also ensure that `/dev/net/tun` exists.

Bootstrapping

This step is not needed if you have downloaded a released tarball. If you pulled the sources directly from an Open vSwitch Git tree or got a Git tree snapshot, then run `boot.sh` in the top source directory to build the “configure” script:

```
$ ./boot.sh
```

Configuring

Configure the package by running the configure script. You can usually invoke configure without any arguments. For example:

```
$ ./configure
```

By default all files are installed under `/usr/local`. Open vSwitch also expects to find its database in `/usr/local/etc/openvswitch` by default. If you want to install all files into, e.g., `/usr` and `/var` instead of `/usr/local` and `/usr/local/var` and expect to use `/etc/openvswitch` as the default database directory, add options as shown here:

```
$ ./configure --prefix=/usr --localstatedir=/var --sysconfdir=/etc
```

Note: Open vSwitch installed with packages like `.rpm` (e.g. via `yum install` or `rpm -ivh`) and `.deb` (e.g. via `apt-get install` or `dpkg -i`) use the above configure options.

By default, static libraries are built and linked against. If you want to use shared libraries instead:

```
$ ./configure --enable-shared
```

To use a specific C compiler for compiling Open vSwitch user programs, also specify it on the configure command line, like so:

```
$ ./configure CC=gcc-4.2
```

To use ‘clang’ compiler:

```
$ ./configure CC=clang
```

To supply special flags to the C compiler, specify them as `CFLAGS` on the configure command line. If you want the default `CFLAGS`, which include `-g` to build debug symbols and `-O2` to enable optimizations, you must include them yourself. For example, to build with the default `CFLAGS` plus `-mssse3`, you might run configure as follows:

```
$ ./configure CFLAGS="-g -O2 -mssse3"
```

For efficient hash computation special flags can be passed to leverage built-in intrinsics. For example on X86_64 with SSE4.2 instruction set support, CRC32 intrinsics can be used by passing `-msse4.2`:

```
$ ./configure CFLAGS="-g -O2 -msse4.2"
```

Also builtin `popcnt` instruction can be used to speedup the counting of the bits set in an integer. For example on X86_64 with `POPCNT` support, it can be enabled by passing `-mpopcnt`:

```
$ ./configure CFLAGS="-g -O2 -mpopcnt"
```

If you are on a different processor and don’t know what flags to choose, it is recommended to use `-march=native` settings:

```
$ ./configure CFLAGS="-g -O2 -march=native"
```

With this, GCC will detect the processor and automatically set appropriate flags for it. This should not be used if you are compiling OVS outside the target machine.

Note: CFLAGS are not applied when building the Linux kernel module. Custom CFLAGS for the kernel module are supplied using the EXTRA_CFLAGS variable when running make. For example:

```
$ make EXTRA_CFLAGS="-Wno-error=date-time"
```

If you are a developer and want to enable Address Sanitizer for debugging purposes, at about a 2x runtime cost, you can add `-fsanitize=address -fno-omit-frame-pointer -fno-common` to CFLAGS. For example:

```
$ ./configure CFLAGS="-g -O2 -fsanitize=address -fno-omit-frame-pointer -fno-common"
```

To build the Linux kernel module, so that you can run the kernel-based switch, pass the location of the kernel build directory on `--with-linux`. For example, to build for a running instance of Linux:

```
$ ./configure --with-linux=/lib/modules/$(uname -r)/build
```

Note: If `--with-linux` requests building for an unsupported version of Linux, then configure will fail with an error message. Refer to the [Open vSwitch FAQ](#) for advice in that case.

If you wish to build the kernel module for an architecture other than the architecture of the machine used for the build, you may specify the kernel architecture string using the KARCH variable when invoking the configure script. For example, to build for MIPS with Linux:

```
$ ./configure --with-linux=/path/to/linux KARCH=mips
```

If you plan to do much Open vSwitch development, you might want to add `--enable-Werror`, which adds the `-Werror` option to the compiler command line, turning warnings into errors. That makes it impossible to miss warnings generated by the build. For example:

```
$ ./configure --enable-Werror
```

If you're building with GCC, then, for improved warnings, install `sparse` (see "Prerequisites") and enable it for the build by adding `--enable-sparse`. Use this with `--enable-Werror` to avoid missing both compiler and `sparse` warnings, e.g.:

```
$ ./configure --enable-Werror --enable-sparse
```

To build with `gcov` code coverage support, add `--enable-coverage`:

```
$ ./configure --enable-coverage
```

The configure script accepts a number of other options and honors additional environment variables. For a full list, invoke configure with the `--help` option:

```
$ ./configure --help
```

You can also run configure from a separate build directory. This is helpful if you want to build Open vSwitch in more than one way from a single source directory, e.g. to try out both GCC and Clang builds, or to build kernel modules for more than one Linux version. For example:

```
$ mkdir _gcc && (cd _gcc && ./configure CC=gcc)
$ mkdir _clang && (cd _clang && ./configure CC=clang)
```

Under certain loads the ovsdb-server and other components perform better when using the jemalloc memory allocator, instead of the glibc memory allocator. If you wish to link with jemalloc add it to LIBS:

```
$ ./configure LIBS=-ljemalloc
```

Building

1. Run GNU make in the build directory, e.g.:

```
$ make
```

or if GNU make is installed as “gmake”:

```
$ gmake
```

If you used a separate build directory, run make or gmake from that directory, e.g.:

```
$ make -C _gcc
$ make -C _clang
```

Note: Some versions of Clang and ccache are not completely compatible. If you see unusual warnings when you use both together, consider disabling ccache.

2. Consider running the testsuite. Refer to [Testing](#) for instructions.
3. Run `make install` to install the executables and manpages into the running system, by default under `/usr/local`:

```
$ make install
```

5. If you built kernel modules, you may install them, e.g.:

```
$ make modules_install
```

It is possible that you already had a Open vSwitch kernel module installed on your machine that came from upstream Linux (in a different directory). To make sure that you load the Open vSwitch kernel module you built from this repository, you should create a `depmod.d` file that prefers your newly installed kernel modules over the kernel modules from upstream Linux. The following snippet of code achieves the same:

```
$ config_file="/etc/depmod.d/openvswitch.conf"
$ for module in datapath/linux/*.ko; do
  modname="$(basename ${module})"
  echo "override ${modname%.ko} * extra" >> "$config_file"
  echo "override ${modname%.ko} * weak-updates" >> "$config_file"
done
$ depmod -a
```

Finally, load the kernel modules that you need. e.g.:

```
$ /sbin/modprobe openvswitch
```

To verify that the modules have been loaded, run `/sbin/lsmmod` and check that `openvswitch` is listed:

```
$ /sbin/lsmmod | grep openvswitch
```

Note: If the `modprobe` operation fails, look at the last few kernel log messages (e.g. with `dmesg | tail`). Generally, issues like this occur when Open vSwitch is built for a kernel different from the one into which you are trying to load it. Run `modinfo` on `openvswitch.ko` and on a module built for the running kernel, e.g.:

```
$ /sbin/modinfo openvswitch.ko
$ /sbin/modinfo /lib/modules/$(uname -r)/kernel/net/bridge/bridge.ko
```

Compare the “vermagic” lines output by the two commands. If they differ, then Open vSwitch was built for the wrong kernel.

If you decide to report a bug or ask a question related to module loading, include the output from the `dmesg` and `modinfo` commands mentioned above.

Starting

On Unix-alike systems, such as BSDs and Linux, starting the Open vSwitch suite of daemons is a simple process. Open vSwitch includes a shell script, and helpers, called `ovs-ctl` which automates much of the tasks for starting and stopping `ovsdb-server`, and `ovs-vswitchd`. After installation, the daemons can be started by using the `ovs-ctl` utility. This will take care to setup initial conditions, and start the daemons in the correct order. The `ovs-ctl` utility is located in `$(pkgdatadir)/scripts`, and defaults to `/usr/local/share/openvswitch/scripts`. An example after install might be:

```
$ export PATH=$PATH:/usr/local/share/openvswitch/scripts
$ ovs-ctl start
```

Additionally, the `ovs-ctl` script allows starting / stopping the daemons individually using specific options. To start just the `ovsdb-server`:

```
$ export PATH=$PATH:/usr/local/share/openvswitch/scripts
$ ovs-ctl --no-ovs-vswitchd start
```

Likewise, to start just the `ovs-vswitchd`:

```
$ export PATH=$PATH:/usr/local/share/openvswitch/scripts
$ ovs-ctl --no-ovsdb-server start
```

Refer to `ovs-ctl(8)` for more information on `ovs-ctl`.

In addition to using the automated script to start Open vSwitch, you may wish to manually start the various daemons. Before starting `ovs-vswitchd` itself, you need to start its configuration database, `ovsdb-server`. Each machine on which Open vSwitch is installed should run its own copy of `ovsdb-server`. Before `ovsdb-server` itself can be started, configure a database that it can use:

```
$ mkdir -p /usr/local/etc/openvswitch
$ ovsdb-tool create /usr/local/etc/openvswitch/conf.db \
  vswitchd/vswitch.ovsschema
```

Configure `ovsdb-server` to use database created above, to listen on a Unix domain socket, to connect to any managers specified in the database itself, and to use the SSL configuration in the database:

```
$ mkdir -p /usr/local/var/run/openvswitch
$ ovsdb-server --remote=punix:/usr/local/var/run/openvswitch/db.sock \
  --remote=db:Open_vSwitch,Open_vSwitch,manager_options \
  --private-key=db:Open_vSwitch,SSL,private_key \
  --certificate=db:Open_vSwitch,SSL,certificate \
  --bootstrap-ca-cert=db:Open_vSwitch,SSL,ca_cert \
  --pidfile --detach --log-file
```

Note: If you built Open vSwitch without SSL support, then omit `--private-key`, `--certificate`, and `--bootstrap-ca-cert`.)

Initialize the database using `ovs-vsctl`. This is only necessary the first time after you create the database with `ovsdb-tool`, though running it at any time is harmless:

```
$ ovs-vsctl --no-wait init
```

Start the main Open vSwitch daemon, telling it to connect to the same Unix domain socket:

```
$ ovs-vswitchd --pidfile --detach --log-file
```

Validating

At this point you can use `ovs-vsctl` to set up bridges and other Open vSwitch features. For example, to create a bridge named `br0` and add ports `eth0` and `vif1.0` to it:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 vif1.0
```

Refer to `ovs-vsctl(8)` for more details. You may also wish to refer to [Testing](#) for information on more generic testing of OVS.

Upgrading

When you upgrade Open vSwitch from one version to another you should also upgrade the database schema:

Note: The following manual steps may also be accomplished by using `ovs-ctl` to stop and start the daemons after upgrade. The `ovs-ctl` script will automatically upgrade the schema.

1. Stop the Open vSwitch daemons, e.g.:

```
$ kill `cd /usr/local/var/run/openvswitch && cat ovsdb-server.pid ovs-vswitchd.
↪pid`
```

2. Install the new Open vSwitch release by using the same configure options as was used for installing the previous version. If you do not use the same configure options, you can end up with two different versions of Open vSwitch executables installed in different locations.
3. Upgrade the database, in one of the following two ways:

- If there is no important data in your database, then you may delete the database file and recreate it with `ovsdb-tool`, following the instructions under “Building and Installing Open vSwitch for Linux, FreeBSD or NetBSD”.
- If you want to preserve the contents of your database, back it up first, then use `ovsdb-tool convert` to upgrade it, e.g.:

```
$ ovsdb-tool convert /usr/local/etc/openvswitch/conf.db \
    vswitchd/vswitch.ovsschema
```

4. Start the Open vSwitch daemons as described under *Starting* above.

Hot Upgrading

Upgrading Open vSwitch from one version to the next version with minimum disruption of traffic going through the system that is using that Open vSwitch needs some considerations:

1. If the upgrade only involves upgrading the userspace utilities and daemons of Open vSwitch, make sure that the new userspace version is compatible with the previously loaded kernel module.
2. An upgrade of userspace daemons means that they have to be restarted. Restarting the daemons means that the OpenFlow flows in the `ovs-vswitchd` daemon will be lost. One way to restore the flows is to let the controller re-populate it. Another way is to save the previous flows using a utility like `ovs-ofctl` and then re-add them after the restart. Restoring the old flows is accurate only if the new Open vSwitch interfaces retain the old ‘ofport’ values.
3. When the new userspace daemons get restarted, they automatically flush the old flows setup in the kernel. This can be expensive if there are hundreds of new flows that are entering the kernel but userspace daemons are busy setting up new userspace flows from either the controller or an utility like `ovs-ofctl`. Open vSwitch database provides an option to solve this problem through the `other_config:flow-restore-wait` column of the `Open_vSwitch` table. Refer to the `ovs-vswitchd.conf.db(5)` manpage for details.
4. If the upgrade also involves upgrading the kernel module, the old kernel module needs to be unloaded and the new kernel module should be loaded. This means that the kernel network devices belonging to Open vSwitch is recreated and the kernel flows are lost. The downtime of the traffic can be reduced if the userspace daemons are restarted immediately and the userspace flows are restored as soon as possible.

The `ovs-ctl` utility’s `restart` function only restarts the userspace daemons, makes sure that the ‘ofport’ values remain consistent across restarts, restores userspace flows using the `ovs-ofctl` utility and also uses the `other_config:flow-restore-wait` column to keep the traffic downtime to the minimum. The `ovs-ctl` utility’s `force-reload-kmod` function does all of the above, but also replaces the old kernel module with the new one. Open vSwitch startup scripts for Debian, XenServer and RHEL use `ovs-ctl`’s functions and it is recommended that these functions be used for other software platforms too.

Reporting Bugs

Report problems to bugs@openvswitch.org.

Open vSwitch on NetBSD

On NetBSD, you might want to install requirements from `pkgsrc`. In that case, you need at least the following packages.

- automake
- libtool-base

- gmake
- python27
- py27-six
- py27-xml

Some components have additional requirements. Refer to *Open vSwitch on Linux, FreeBSD and NetBSD* for more information.

Assuming you are running NetBSD/amd64 6.1.2, you can download and install pre-built binary packages as the following:

```
$ PKG_PATH=http://ftp.netbsd.org/pub/pkgsrc/packages/NetBSD/amd64/7.0.2/All/  
$ export PKG_PATH  
$ pkg_add automake libtool-base gmake python27 py27-six py27-xml \  
  pkg_alternatives
```

Note: You might get some warnings about minor version mismatch. These can be safely ignored.

NetBSD's `/usr/bin/make` is not GNU make. GNU make is installed as `/usr/pkg/bin/gmake` by the above mentioned gmake package.

As all executables installed with pkgsrc are placed in `/usr/pkg/bin/` directory, it might be a good idea to add it to your PATH. Or install OVS by `gmake` and `gmake install`.

Open vSwitch on NetBSD is currently “userspace switch” implementation in the sense described in *Open vSwitch without Kernel Support* and *Porting Open vSwitch to New Software or Hardware*.

Open vSwitch on Windows

Build Requirements

Open vSwitch on Linux uses autoconf and automake for generating Makefiles. It will be useful to maintain the same build system while compiling on Windows too. One approach is to compile Open vSwitch in a MinGW environment that contains autoconf and automake utilities and then use Visual C++ as a compiler and linker.

The following explains the steps in some detail.

- Mingw

Install Mingw on a Windows machine by following the instructions on mingw.org.

This should install mingw at `C:\Mingw` and msys at `C:\Mingw\msys`. Add `C:\MinGW\bin` and `C:\Mingw\msys\1.0\bin` to PATH environment variable of Windows.

You can either use the MinGW installer or the command line utility `mingw-get` to install both the base packages and additional packages like automake and autoconf(version 2.68).

Also make sure that `/mingw` mount point exists. If its not, please add/create the following entry in `/etc/fstab`:

```
'C:/MinGW /mingw'.
```

- Python

Install the latest Python 2.x from python.org and verify that its path is part of Windows' PATH environment variable. We require that you have Python six and pypiwin32 libraries installed. The libraries can be installed via pip command:

```
$ pip install six
$ pip install pypiwin32
```

- Visual Studio

You will need at least Visual Studio 2013 (update 4) to compile userspace binaries. In addition to that, if you want to compile the kernel module you will also need to install Windows Driver Kit (WDK) 8.1 Update.

It is important to get the Visual Studio related environment variables and to have the \$PATH inside the bash to point to the proper compiler and linker. One easy way to achieve this for VS2013 is to get into the “VS2013 x86 Native Tools Command Prompt” (in a default installation of Visual Studio 2013 this can be found under the following location: C:\Program Files (x86)\Microsoft Visual Studio 12.0\Common7\Tools\Shortcuts) and through it enter into the bash shell available from msys by typing `bash --login`.

There is support for generating 64 bit binaries too. To compile under x64, open the “VS2013 x64 Native Tools Command Prompt” (if your current running OS is 64 bit) or “VS2013 x64 Cross Tools Command Prompt” (if your current running OS is not 64 bit) instead of opening its x86 variant. This will point the compiler and the linker to their 64 bit equivalent.

If after the above step, a `which link` inside MSYS's bash says, `/bin/link.exe`, rename `/bin/link.exe` to something else so that the Visual studio's linker is used. You should also see a ‘which sort’ report `/bin/sort.exe`.

- pthreads-win32

For pthread support, install the library, dll and includes of pthreads-win32 project from [sourceware](http://sourceware.org) to a directory (e.g.: C:\pthread). You should add the pthread-win32's dll path (e.g.: C:\pthread\dll\x86) to the Windows' PATH environment variable.

- OpenSSL

To get SSL support for Open vSwitch on Windows, you will need to install [OpenSSL for Windows](#)

Note down the directory where OpenSSL is installed (e.g.: C:/OpenSSL-Win32) for later use.

Note: Commands prefixed by \$ must be run in the Bash shell provided by MinGW. Open vSwitch commands, such as `ovs-dpctl` are shown running under the DOS shell (`cmd.exe`), as indicated by the > prefix, but will also run under Bash. The remainder, prefixed by >, are PowerShell commands and must be run in PowerShell.

Install Requirements

- Share network adaptors

We require that you don't disable the “Allow management operating system to share this network adapter” under ‘Virtual Switch Properties’ > ‘Connection type: External network’, in the Hyper-V virtual network switch configuration.

- Checksum Offloads

While there is some support for checksum/segmentation offloads in software, this is still a work in progress. Till the support is complete we recommend disabling TX/RX offloads for both the VM's as well as the Hyper-V.

Bootstrapping

This step is not needed if you have downloaded a released tarball. If you pulled the sources directly from an Open vSwitch Git tree or got a Git tree snapshot, then run `boot.sh` in the top source directory to build the “configure” script:

```
$ ./boot.sh
```

Configuring

Configure the package by running the configure script. You should provide some configure options to choose the right compiler, linker, libraries, Open vSwitch component installation directories, etc. For example:

```
$ ./configure CC=./build-aux/cccl LD="$(which link)" \
  LIBS="-lws2_32 -lShlwapi -liphlpapi -lbemuuid -lole32 -loleaut32" \
  --prefix="C:/openvswitch/usr" \
  --localstatedir="C:/openvswitch/var" \
  --sysconfdir="C:/openvswitch/etc" \
  --with-pthread="C:/pthread"
```

Note: By default, the above enables compiler optimization for fast code. For default compiler optimization, pass the `--with-debug` configure option.

To configure with SSL support, add the requisite additional options:

```
$ ./configure CC=./build-aux/cccl LD="`which link`" \
  LIBS="-lws2_32 -lShlwapi -liphlpapi -lbemuuid -lole32 -loleaut32" \
  --prefix="C:/openvswitch/usr" \
  --localstatedir="C:/openvswitch/var" \
  --sysconfdir="C:/openvswitch/etc" \
  --with-pthread="C:/pthread" \
  --enable-ssl --with-openssl="C:/OpenSSL-Win32"
```

Finally, to the kernel module also:

```
$ ./configure CC=./build-aux/cccl LD="`which link`" \
  LIBS="-lws2_32 -lShlwapi -liphlpapi -lbemuuid -lole32 -loleaut32" \
  --prefix="C:/openvswitch/usr" \
  --localstatedir="C:/openvswitch/var" \
  --sysconfdir="C:/openvswitch/etc" \
  --with-pthread="C:/pthread" \
  --enable-ssl --with-openssl="C:/OpenSSL-Win32" \
  --with-vstudiotarget="<target type>" \
  --with-vstudiotargetver="<target versions>"
```

Possible values for `<target type>` are: Debug and Release Possible values for `<target versions>` is a comma separated list of target versions to compile among: Win8, Win8.1, Win10

Note: You can directly use the Visual Studio 2013 IDE to compile the kernel datapath. Open the `ovsext.sln` file in the IDE and build the solution.

Refer to *Open vSwitch on Linux, FreeBSD and NetBSD* for information on additional configuration options.

Building

Once correctly configured, building Open vSwitch on Windows is similar to building on Linux, FreeBSD, or NetBSD.

1. Run make for the ported executables in the top source directory, e.g.:

```
$ make
```

For faster compilation, you can pass the `-j` argument to make. For example, to run 4 jobs simultaneously, run `make -j4`.

Note: MSYS 1.0.18 has a bug that causes parallel make to hang. You can overcome this by downgrading to MSYS 1.0.17. A simple way to downgrade is to exit all MinGW sessions and then run the below command from MSVC developers command prompt.:

```
> mingw-get upgrade msys-core-bin=1.0.17-1
```

2. To run all the unit tests in Open vSwitch, one at a time:

```
$ make check
```

To run all the unit tests in Open vSwitch, up to 8 in parallel:

```
$ make check TESTSUITEFLAGS="-j8"
```

3. To install all the compiled executables on the local machine, run:

```
$ make install
```

Note: This will install the Open vSwitch executables in `C:/openvswitch`. You can add `C:\openvswitch\usr\bin` and `C:\openvswitch\usr\sbin` to Windows' PATH environment variable for easy access.

The Kernel Module

If you are building the kernel module, you will need to copy the below files to the target Hyper-V machine.

- `./datapath-windows/x64/Win8.1Debug/package/ovsext.inf`
- `./datapath-windows/x64/Win8.1Debug/package/OVSExt.sys`
- `./datapath-windows/x64/Win8.1Debug/package/ovsext.cat`
- `./datapath-windows/misc/install.cmd`
- `./datapath-windows/misc/uninstall.cmd`

Note: The above path assumes that the kernel module has been built using Windows DDK 8.1 in Debug mode. Change the path appropriately, if a different WDK has been used.

Now run `./uninstall.cmd` to remove the old extension. Once complete, run `./install.cmd` to insert the new one. For this to work you will have to turn on `TESTSIGNING` boot option or 'Disable Driver Signature Enforcement' during boot. The following commands can be used:

```
> bcdedit /set LOADOPTIONS DISABLE_INTEGRITY_CHECKS
> bcdedit /set TESTSIGNING ON
> bcdedit /set nointegritychecks ON
```

Note: You may have to restart the machine for the settings to take effect.

In the Virtual Switch Manager configuration you can enable the Open vSwitch Extension on an existing switch or create a new switch. If you are using an existing switch, make sure to enable the “Allow Management OS” option for VXLAN to work (covered later).

The command to create a new switch named ‘OVS-Extended-Switch’ using a physical NIC named ‘Ethernet 1’ is:

```
PS > New-VMSwitch "OVS-Extended-Switch" -NetAdapterName "Ethernet 1"
```

Note: You can obtain the list of physical NICs on the host using ‘Get-NetAdapter’ command.

In the properties of any switch, you should now see “Open vSwitch Extension” under ‘Extensions’. Click the check box to enable the extension. An alternative way to do the same is to run the following command:

```
PS > Enable-VMSwitchExtension "Open vSwitch Extension" OVS-Extended-Switch
```

Note: If you enabled the extension using the command line, a delay of a few seconds has been observed for the change to be reflected in the UI. This is not a bug in Open vSwitch.

Starting

Important: The following steps assume that you have installed the Open vSwitch utilities in the local machine via ‘make install’.

Before starting ovs-vswitchd itself, you need to start its configuration database, ovsdb-server. Each machine on which Open vSwitch is installed should run its own copy of ovsdb-server. Before ovsdb-server itself can be started, configure a database that it can use:

```
> ovsdb-tool create C:\openvswitch\etc\openvswitch\conf.db \
  C:\openvswitch\usr\share\openvswitch\vswitch.ovsschema
```

Configure ovsdb-server to use database created above and to listen on a Unix domain socket:

```
> ovsdb-server -vfile:info --remote=punix:db.sock --log-file \
  --pidfile --detach
```

Note: The logfile is created at C:/openvswitch/var/log/openvswitch/

Initialize the database using ovs-vsctl. This is only necessary the first time after you create the database with ovsdb-tool, though running it at any time is harmless:

```
> ovs-vsctl --no-wait init
```

Tip: If you would later like to terminate the started ovsdb-server, run:

```
> ovs-appctl -t ovsdb-server exit
```

Start the main Open vSwitch daemon, telling it to connect to the same Unix domain socket:

```
> ovs-vswitchd -vfile:info --log-file --pidfile --detach
```

Tip: If you would like to terminate the started ovs-vswitchd, run:

```
> ovs-appctl exit
```

Note: The logfile is created at C:/openvswitch/var/log/openvswitch/

Validating

At this point you can use ovs-vsctl to set up bridges and other Open vSwitch features.

Add bridges

Let's start by creating an integration bridge, br-int and a PIF bridge, br-pif:

```
> ovs-vsctl add-br br-int
> ovs-vsctl add-br br-pif
```

Note: There's a known bug that running the ovs-vsctl command does not terminate. This is generally solved by having ovs-vswitchd running. If you face the issue despite that, hit Ctrl-C to terminate ovs-vsctl and check the output to see if your command succeeded.

Validate that ports are added by dumping from both ovs-dpctl and ovs-vsctl:

```
> ovs-dpctl show
system@ovs-system:
  lookups: hit:0 missed:0 lost:0
  flows: 0
  port 2: br-pif (internal)    <<< internal port on 'br-pif' bridge
  port 1: br-int (internal)    <<< internal port on 'br-int' bridge

> ovs-vsctl show
a56ec7b5-5b1f-49ec-a795-79f6eb63228b
  Bridge br-pif
    Port br-pif
      Interface br-pif
        type: internal
```

(continues on next page)

(continued from previous page)

```

Bridge br-int
  Port br-int
    Interface br-int
      type: internal

```

Note: There's a known bug that the ports added to OVSDB via `ovs-vsctl` don't get to the kernel datapath immediately, ie. they don't show up in the output of `ovs-dpctl show` even though they show up in output of `ovs-vsctl show`. In order to workaround this issue, restart `ovs-vswitchd`. (You can terminate `ovs-vswitchd` by running `ovs-appctl exit`.)

Add physical NICs (PIF)

Now, let's add the physical NIC and the internal port to `br-pif`. In OVS for Hyper-V, we use the name of the adapter on top of which the Hyper-V virtual switch was created, as a special name to refer to the physical NICs connected to the Hyper-V switch, e.g. if we created the Hyper-V virtual switch on top of the adapter named `Ethernet0`, then in OVS we use that name (`Ethernet0`) as a special name to refer to that adapter.

Note: We assume that the OVS extension is enabled Hyper-V switch.

Internal ports are the virtual adapters created on the Hyper-V switch using the `ovs-vsctl add-br <bridge>` command. By default they are created under the following rule "<name of bridge>" and the adapters are disabled. One needs to enable them and set the corresponding values to it to make them IP-able.

As a whole example, if we issue the following in a powershell console:

```

PS > Get-NetAdapter | select Name,InterfaceDescription
Name                               InterfaceDescription
----                               -
Ethernet1                          Intel(R) PRO/1000 MT Network Connection
br-pif                             Hyper-V Virtual Ethernet Adapter #2
Ethernet0                          Intel(R) PRO/1000 MT Network Connection #2
br-int                             Hyper-V Virtual Ethernet Adapter #3

PS > Get-VMSwitch
Name      SwitchType NetAdapterInterfaceDescription
-----
external External   Intel(R) PRO/1000 MT Network Connection #2

```

We can see that we have a switch(`external`) created upon adapter name 'Ethernet0' with the internal ports under name 'br-pif' and 'br-int'. Thus resulting into the following `ovs-vsctl` commands:

```
> ovs-vsctl add-port br-pif Ethernet0
```

Dumping the ports should show the additional ports that were just added:

```

> ovs-dpctl show
system@ovs-system:
  lookups: hit:0 missed:0 lost:0
  flows: 0
  port 2: br-pif (internal)          <<< internal port
                                      adapter on

```

(continues on next page)

(continued from previous page)

```

port 1: br-int (internal)          Hyper-V switch
                                   <<< internal port
                                   adapter on
                                   Hyper-V switch
port 3: Ethernet0                 <<< Physical NIC

> ovs-vsctl show
a56ec7b5-5b1f-49ec-a795-79f6eb63228b
    Bridge br-pif
        Port br-pif
            Interface br-pif
                type: internal
        Port "Ethernet0"
            Interface "Ethernet0"
    Bridge br-int
        Port br-int
            Interface br-int
                type: internal

```

Add virtual interfaces (VIFs)

Adding VIFs to Open vSwitch is a two step procedure. The first step is to assign a ‘OVS port name’ which is a unique name across all VIFs on this Hyper-V. The next step is to add the VIF to the ovssdb using its ‘OVS port name’ as key.

First, assign a unique ‘OVS port name’ to the VIF. The VIF needs to have been disconnected from the Hyper-V switch before assigning a ‘OVS port name’ to it. In the example below, we assign a ‘OVS port name’ called `ovs-port-a` to a VIF on a VM VM1. By using index 0 for `$vnic`, the first VIF of the VM is being addressed. After assigning the name `ovs-port-a`, the VIF is connected back to the Hyper-V switch with name `OVS-HV-Switch`, which is assumed to be the Hyper-V switch with OVS extension enabled.:

```

PS > import-module .\datapath-windows\misc\OVS.psml
PS > $vnic = Get-VMNetworkAdapter <Name of the VM>
PS > Disconnect-VMNetworkAdapter -VMNetworkAdapter $vnic[0]
PS > $vnic[0] | Set-VMNetworkAdapterOVSPort -OVSPortName ovs-port-a
PS > Connect-VMNetworkAdapter -VMNetworkAdapter $vnic[0] \
    -SwitchName OVS-Extended-Switch

```

Next, add the VIFs to `br-int`:

```
> ovs-vsctl add-port br-int ovs-port-a
```

Dumping the ports should show the additional ports that were just added:

```

> ovs-dpctl show
system@ovs-system:
    lookups: hit:0 missed:0 lost:0
    flows: 0
    port 4: ovs-port-a
    port 2: br-pif (internal)
    port 1: br-int (internal)
    port 3: Ethernet0

> ovs-vsctl show
4cd86499-74df-48bd-a64d-8d115b12a9f2
    Bridge br-pif

```

(continues on next page)

(continued from previous page)

```

Port "vEthernet (external)"
    Interface "vEthernet (external)"
Port "Ethernet0"
    Interface "Ethernet0"
Port br-pif
    Interface br-pif
        type: internal
Bridge br-int
    Port br-int
        Interface br-int
            type: internal
    Port "ovs-port-a"
        Interface "ovs-port-a"

```

Add multiple NICs to be managed by OVS

To leverage support of multiple NICs into OVS, we will be using the MSFT cmdlets for forwarding team extension. More documentation about them can be found at [technet](#).

For example, to set up a switch team combined from Ethernet0 2 and Ethernet1 2 named external:

```

PS > Get-NetAdapter
Name                               InterfaceDescription
----
br-int                             Hyper-V Virtual Ethernet Adapter #3
br-pif                             Hyper-V Virtual Ethernet Adapter #2
Ethernet3 2                       Intel(R) 82574L Gigabit Network Co...#3
Ethernet2 2                       Intel(R) 82574L Gigabit Network Co...#4
Ethernet1 2                       Intel(R) 82574L Gigabit Network Co...#2
Ethernet0 2                       Intel(R) 82574L Gigabit Network Conn...

PS > New-NetSwitchTeam -Name external -TeamMembers "Ethernet0 2","Ethernet1 2"

PS > Get-NetSwitchTeam
Name      : external
Members  : {Ethernet1 2, Ethernet0 2}

```

This will result in a new adapter bound to the host called external:

```

PS > Get-NetAdapter
Name                               InterfaceDescription
----
br-test                            Hyper-V Virtual Ethernet Adapter #4
br-pif                             Hyper-V Virtual Ethernet Adapter #2
external                          Microsoft Network Adapter Multiplexo...
Ethernet3 2                       Intel(R) 82574L Gigabit Network Co...#3
Ethernet2 2                       Intel(R) 82574L Gigabit Network Co...#4
Ethernet1 2                       Intel(R) 82574L Gigabit Network Co...#2
Ethernet0 2                       Intel(R) 82574L Gigabit Network Conn...

```

Next we will set up the Hyper-V VMSwitch on the new adapter external:

```

PS > New-VMSwitch -Name external -NetAdapterName external \
    -AllowManagementOS $false

```


Under OVS the adapters under the team external, Ethernet0 2 and Ethernet1 2, can be added either under a bond device or separately.

The following example shows how the bridges look with the NICs being separated:

```
> ovs-vsctl show
6cd9481b-c249-4ee3-8692-97b399dd29d8
    Bridge br-test
        Port br-test
            Interface br-test
                type: internal
        Port "Ethernet1 2"
            Interface "Ethernet1 2"
    Bridge br-pif
        Port "Ethernet0 2"
            Interface "Ethernet0 2"
        Port br-pif
            Interface br-pif
                type: internal
```

Add patch ports and configure VLAN tagging

The Windows Open vSwitch implementation support VLAN tagging in the switch. Switch VLAN tagging along with patch ports between br-int and br-pif is used to configure VLAN tagging functionality between two VMs on different Hyper-Vs. To start, add a patch port from br-int to br-pif:

```
> ovs-vsctl add-port br-int patch-to-pif
> ovs-vsctl set interface patch-to-pif type=patch \
    options:peer=patch-to-int
```

Add a patch port from br-pif to br-int:

```
> ovs-vsctl add-port br-pif patch-to-int
> ovs-vsctl set interface patch-to-int type=patch \
    options:peer=patch-to-pif
```

Re-Add the VIF ports with the VLAN tag:

```
> ovs-vsctl add-port br-int ovs-port-a tag=900
> ovs-vsctl add-port br-int ovs-port-b tag=900
```

Add tunnels

The Windows Open vSwitch implementation support VXLAN and STT tunnels. To add tunnels. For example, first add the tunnel port between 172.168.201.101 <-> 172.168.201.102:

```
> ovs-vsctl add-port br-int tun-1
> ovs-vsctl set Interface tun-1 type=<port-type>
> ovs-vsctl set Interface tun-1 options:local_ip=172.168.201.101
> ovs-vsctl set Interface tun-1 options:remote_ip=172.168.201.102
> ovs-vsctl set Interface tun-1 options:in_key=flow
> ovs-vsctl set Interface tun-1 options:out_key=flow
```

...and the tunnel port between 172.168.201.101 <-> 172.168.201.105:

```
> ovs-vsctl add-port br-int tun-2
> ovs-vsctl set Interface tun-2 type=<port-type>
> ovs-vsctl set Interface tun-2 options:local_ip=172.168.201.102
> ovs-vsctl set Interface tun-2 options:remote_ip=172.168.201.105
> ovs-vsctl set Interface tun-2 options:in_key=flow
> ovs-vsctl set Interface tun-2 options:out_key=flow
```

Where <port-type> is one of: stt or vxlan

Note: Any patch ports created between br-int and br-pif **MUST** be deleted prior to adding tunnels.

Windows Services

Open vSwitch daemons come with support to run as a Windows service. The instructions here assume that you have installed the Open vSwitch utilities and daemons via `make install`.

To start, create the database:

```
> ovssdb-tool create C:/openvswitch/etc/openvswitch/conf.db \
  "C:/openvswitch/usr/share/openvswitch/vswitch.ovsschema"
```

Create the ovssdb-server service and start it:

```
> sc create ovssdb-server \
  binpath="C:/openvswitch/usr/sbin/ovssdb-server.exe \
  C:/openvswitch/etc/openvswitch/conf.db \
  -vfile:info --log-file --pidfile \
  --remote=punix:db.sock --service --service-monitor"
> sc start ovssdb-server
```

Tip: One of the common issues with creating a Windows service is with mungled paths. You can make sure that the correct path has been registered with the Windows services manager by running:

```
> sc qc ovssdb-server
```

Check that the service is healthy by running:

```
> sc query ovssdb-server
```

Initialize the database:

```
> ovs-vsctl --no-wait init
```

Create the ovs-vswitchd service and start it:

```
> sc create ovs-vswitchd \
  binpath="C:/openvswitch/usr/sbin/ovs-vswitchd.exe \
  --pidfile -vfile:info --log-file --service --service-monitor"
> sc start ovs-vswitchd
```

Check that the service is healthy by running:

```
> sc query ovs-vswitchd
```

To stop and delete the services, run:

```
> sc stop ovs-vswitchd
> sc stop ovssdb-server
> sc delete ovs-vswitchd
> sc delete ovssdb-server
```

Windows CI Service

[AppVeyor](#) provides a free Windows autobuild service for open source projects. Open vSwitch has integration with AppVeyor for continuous build. A developer can build test his changes for Windows by logging into appveyor.com using a github account, creating a new project by linking it to his development repository in github and triggering a new build.

TODO

- Investigate the working of sFlow on Windows and re-enable the unit tests.
- Investigate and add the feature to provide QoS.
- Sign the driver & create an MSI for installing the different Open vSwitch components on Windows.

Open vSwitch on Citrix XenServer

This document describes how to build and install Open vSwitch on a Citrix XenServer host. If you want to install Open vSwitch on a generic Linux or BSD host, refer to [Open vSwitch on Linux, FreeBSD and NetBSD](#) instead.

Open vSwitch should work with XenServer 5.6.100 and later. However, Open vSwitch requires Python 2.7 or later, so using Open vSwitch with XenServer 6.5 or earlier requires installing Python 2.7.

Building

You may build from an Open vSwitch distribution tarball or from an Open vSwitch Git tree. The recommended build environment to build RPMs for Citrix XenServer is the DDK VM available from Citrix.

1. If you are building from an Open vSwitch Git tree, then you will need to first create a distribution tarball by running:

```
$ ./boot.sh
$ ./configure
$ make dist
```

You cannot run this in the DDK VM, because it lacks tools that are necessary to bootstrap the Open vSwitch distribution. Instead, you must run this on a machine that has the tools listed in [Installation Requirements](#) as prerequisites for building from a Git tree.

2. Copy the distribution tarball into `/usr/src/redhat/SOURCES` inside the DDK VM.
3. In the DDK VM, unpack the distribution tarball into a temporary directory and “cd” into the root of the distribution tarball.
4. To build Open vSwitch userspace, run:

```
$ rpmbuild -bb xenserver/openvswitch-xen.spec
```

This produces three RPMs in `/usr/src/redhat/RPMS/i386`:

- `openvswitch`
- `openvswitch-modules-xen`
- `openvswitch-debuginfo`

The above command automatically runs the Open vSwitch unit tests. To disable the unit tests, run:

```
$ rpmbuild -bb --without check xenserver/openvswitch-xen.spec
```

Build Parameters

`openvswitch-xen.spec` needs to know a number of pieces of information about the XenServer kernel. Usually, it can figure these out for itself, but if it does not do it correctly then you can specify them yourself as parameters to the build. Thus, the final `rpmbuild` step above can be elaborated as:

```
$ VERSION=<Open vSwitch version>
$ KERNEL_NAME=<Xen Kernel name>
$ KERNEL_VERSION=<Xen Kernel version>
$ KERNEL_FLAVOR=<Xen Kernel flavor(suffix)>
$ rpmbuild \
  -D "openvswitch_version $VERSION" \
  -D "kernel_name $KERNEL_NAME" \
  -D "kernel_version $KERNEL_VERSION" \
  -D "kernel_flavor $KERNEL_FLAVOR" \
  -bb xenserver/openvswitch-xen.spec
```

where:

<openvswitch version> is the version number that appears in the name of the Open vSwitch tarball, e.g. 0.90.0.

<Xen Kernel name> is the name of the XenServer kernel package, e.g. `kernel-xen` or `kernel-NAME-xen`, without the `kernel-` prefix.

<Xen Kernel version> is the output of:

```
$ rpm -q --queryformat "%{Version}-%{Release}" <kernel-devel-package>,
```

e.g. `2.6.32.12-0.7.1.xs5.6.100.323.170596`, where `<kernel-devel-package>` is the name of the `-devel` package corresponding to `<Xen Kernel name>`.

<Xen Kernel flavor (suffix)> is either `xen` or `kdump`, where `xen` flavor is the main running kernel flavor and the `kdump` flavor is the crashdump kernel flavor. Commonly, one would specify `xen` here.

For XenServer 6.5 or above, the kernel version naming no longer contains `KERNEL_FLAVOR`. In fact, only providing the `uname -r` output is enough. So, the final `rpmbuild` step changes to:

```
$ KERNEL_UNAME=`uname -r` output>
$ rpmbuild \
  -D "kernel_uname $KERNEL_UNAME" \
  -bb xenserver/openvswitch-xen.spec
```

Installing Open vSwitch for XenServer

To install Open vSwitch on a XenServer host, or to upgrade to a newer version, copy the `openvswitch` and `openvswitch-modules-xen` RPMs to that host with `scp`, then install them with `rpm -U`, e.g.:

```
$ scp openvswitch-$VERSION-1.i386.rpm \
    openvswitch-modules-xen-$XEN_KERNEL_VERSION-$VERSION-1.i386.rpm \
    root@<host>:
# Enter <host>'s root password.
$ ssh root@<host>
# Enter <host>'s root password again.
$ rpm -U openvswitch-$VERSION-1.i386.rpm \
    openvswitch-modules-xen-$XEN_KERNEL_VERSION-$VERSION-1.i386.rpm
```

To uninstall Open vSwitch from a XenServer host, remove the packages:

```
$ ssh root@<host>
# Enter <host>'s root password again.
$ rpm -e openvswitch openvswitch-modules-xen-$XEN_KERNEL_VERSION
```

After installing or uninstalling Open vSwitch, the XenServer should be rebooted as soon as possible.

Open vSwitch Boot Sequence on XenServer

When Open vSwitch is installed on XenServer, its startup script `/etc/init.d/openvswitch` runs early in boot. It does roughly the following:

- Loads the OVS kernel module, `openvswitch`.
- Starts `ovsdb-server`, the OVS configuration database.
- XenServer expects there to be no bridges configured at startup, but the OVS configuration database likely still has bridges configured from before reboot. To match XenServer expectations, the startup script deletes all configured bridges from the database.
- Starts `ovs-vswitchd`, the OVS switching daemon.

At this point in the boot process, then, there are no Open vSwitch bridges, even though all of the Open vSwitch daemons are running. Later on in boot, `/etc/init.d/management-interface` (part of XenServer, not Open vSwitch) creates the bridge for the XAPI management interface by invoking `/opt/xensource/libexec/interface-reconfigure`. Normally this program consults XAPI's database to obtain information about how to configure the bridge, but XAPI is not running yet(*) so it instead consults `/var/xapi/network.dbcache`, which is a cached copy of the most recent network configuration.

(*) Even if XAPI were running, if this XenServer node is a pool slave then the query would have to consult the master, which requires network access, which begs the question of how to configure the management interface.

XAPI starts later on in the boot process. XAPI can then create other bridges on demand using `/opt/xensource/libexec/interface-reconfigure`. Now that XAPI is running, that program consults XAPI directly instead of reading the cache.

As part of its own startup, XAPI invokes the Open vSwitch XAPI plugin script `/etc/xapi.d/openvswitch-cfg-update` passing the `update` command. The plugin script does roughly the following:

- Calls `/opt/xensource/libexec/interface-reconfigure` with the `rewrite` command, to ensure that the network cache is up-to-date.

- Queries the Open vSwitch manager setting (named `vswitch_controller`) from the XAPI database for the XenServer pool.
- If XAPI and OVS are configured for different managers, or if OVS is configured for a manager but XAPI is not, runs `ovs-vsctl emer-reset` to bring the Open vSwitch configuration to a known state. One effect of `emer-reset` is to deconfigure any manager from the OVS database.
- If XAPI is configured for a manager, configures the OVS manager to match with `ovs-vsctl set-manager`.

Notes

- The Open vSwitch boot sequence only configures an OVS configuration database manager. There is no way to directly configure an OpenFlow controller on XenServer and, as a consequence of the step above that deletes all of the bridges at boot time, controller configuration only persists until XenServer reboot. The configuration database manager can, however, configure controllers for bridges. See the BUGS section of `ovs-testcontroller(8)` for more information on this topic.
- The Open vSwitch startup script automatically adds a firewall rule to allow GRE traffic. This rule is needed for the XenServer feature called “Cross-Host Internal Networks” (CHIN) that uses GRE. If a user configures tunnels other than GRE (ex: Geneve, VXLAN, LISP), they will have to either manually add a iptables firewall rule to allow the tunnel traffic or add it through a startup script (Please refer to the “enable-protocol” command in the `ovs-ctl(8)` manpage).

Reporting Bugs

Please report problems to bugs@openvswitch.org.

Open vSwitch without Kernel Support

Open vSwitch can operate, at a cost in performance, entirely in userspace, without assistance from a kernel module. This file explains how to install Open vSwitch in such a mode.

This version of Open vSwitch should be built manually with `configure` and `make`. Debian packaging for Open vSwitch is also included, but it has not been recently tested, and so Debian packages are not a recommended way to use this version of Open vSwitch.

Warning: The userspace-only mode of Open vSwitch without DPDK is considered experimental. It has not been thoroughly tested.

Building and Installing

The requirements and procedure for building, installing, and configuring Open vSwitch are the same as those given in *Open vSwitch on Linux, FreeBSD and NetBSD*. You may omit configuring, building, and installing the kernel module, and the related requirements.

On Linux, the userspace switch additionally requires the kernel TUN/TAP driver to be available, either built into the kernel or loaded as a module. If you are not sure, check for a directory named `/sys/class/misc/tun`. If it does not exist, then attempt to load the module with `modprobe tun`.

The tun device must also exist as `/dev/net/tun`. If it does not exist, then create `/dev/net` (if necessary) with `mkdir /dev/net`, then create `/dev/net/tun` with `mknod /dev/net/tun c 10 200`.

On FreeBSD and NetBSD, the userspace switch additionally requires the kernel tap(4) driver to be available, either built into the kernel or loaded as a module.

Using the Userspace Datapath with ovs-vswitchd

To use ovs-vswitchd in userspace mode, create a bridge with datapath_type=netdev in the configuration database. For example:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl set bridge br0 datapath_type=netdev
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 eth1
$ ovs-vsctl add-port br0 eth2
```

ovs-vswitchd will create a TAP device as the bridge's local interface, named the same as the bridge, as well as for each configured internal interface.

Currently, on FreeBSD, the functionality required for in-band control support is not implemented. To avoid related errors, you can disable the in-band support with the following command:

```
$ ovs-vsctl set bridge br0 other_config:disable-in-band=true
```

Firewall Rules

On Linux, when a physical interface is in use by the userspace datapath, packets received on the interface still also pass into the kernel TCP/IP stack. This can cause surprising and incorrect behavior. You can use “iptables” to avoid this behavior, by using it to drop received packets. For example, to drop packets received on eth0:

```
$ iptables -A INPUT -i eth0 -j DROP
$ iptables -A FORWARD -i eth0 -j DROP
```

Other Settings

On NetBSD, depending on your network topology and applications, the following configuration might help. See sysctl(7):

```
$ sysctl -w net.inet.ip.checkinterface=1
```

Reporting Bugs

Report problems to bugs@openvswitch.org.

Open vSwitch with DPDK

This document describes how to build and install Open vSwitch using a DPDK datapath. Open vSwitch can use the DPDK library to operate entirely in userspace.

Important: The [releases FAQ](#) lists support for the required versions of DPDK for each version of Open vSwitch. If building OVS and DPDK outside of the master build tree users should consult this list first.

Build requirements

In addition to the requirements described in *Open vSwitch on Linux, FreeBSD and NetBSD*, building Open vSwitch with DPDK will require the following:

- DPDK 18.11
- A [DPDK supported NIC](#)
Only required when physical ports are in use
- A suitable kernel

On Linux Distros running kernel version ≥ 3.0 , only *IOMMU* needs to be enabled via the grub cmdline, assuming you are using **VFIO**. For older kernels, ensure the kernel is built with **UIO**, **HUGETLBFS**, **PROC_PAGE_MONITOR**, **HPET**, **HPET_MMAP** support. If these are not present, it will be necessary to upgrade your kernel or build a custom kernel with these flags enabled.

Detailed system requirements can be found at [DPDK requirements](#).

Installing

Install DPDK

1. Download the [DPDK sources](#), extract the file and set `DPDK_DIR`:

```
$ cd /usr/src/  
$ wget http://fast.dpdk.org/rel/dpdk-18.11.tar.xz  
$ tar xf dpdk-18.11.tar.xz  
$ export DPDK_DIR=/usr/src/dpdk-18.11  
$ cd $DPDK_DIR
```

2. (Optional) Configure DPDK as a shared library

DPDK can be built as either a static library or a shared library. By default, it is configured for the former. If you wish to use the latter, set `CONFIG_RTE_BUILD_SHARED_LIB=y` in `$DPDK_DIR/config/common_base`.

Note: Minor performance loss is expected when using OVS with a shared DPDK library compared to a static DPDK library.

3. Configure and install DPDK

Build and install the DPDK library:

```
$ export DPDK_TARGET=x86_64-native-linuxapp-gcc  
$ export DPDK_BUILD=$DPDK_DIR/$DPDK_TARGET  
$ make install T=$DPDK_TARGET DESTDIR=install
```

4. (Optional) Export the DPDK shared library location

If DPDK was built as a shared library, export the path to this library for use when building OVS:

```
$ export LD_LIBRARY_PATH=$DPDK_DIR/x86_64-native-linuxapp-gcc/lib
```


Install OVS

OVS can be installed using different methods. For OVS to use DPDK datapath, it has to be configured with DPDK support (`--with-dpdk`).

Note: This section focuses on generic recipe that suits most cases. For distribution specific instructions, refer to one of the more relevant guides.

1. Ensure the standard OVS requirements, described in [Build Requirements](#), are installed
2. Bootstrap, if required, as described in [Bootstrapping](#)
3. Configure the package using the `--with-dpdk` flag:

```
$ ./configure --with-dpdk=$DPDK_BUILD
```

where `DPDK_BUILD` is the path to the built DPDK library. This can be skipped if DPDK library is installed in its default location.

If no path is provided to `--with-dpdk`, but a `pkg-config` configuration for `libdpdk` is available the include paths will be generated via an equivalent `pkg-config --cflags libdpdk`.

Note: While `--with-dpdk` is required, you can pass any other configuration option described in [Configuring](#).

4. Build and install OVS, as described in [Building](#)

Additional information can be found in [Open vSwitch on Linux, FreeBSD and NetBSD](#).

Note: If you are running using the Fedora or Red Hat package, the Open vSwitch daemon will run as a non-root user. This implies that you must have a working IOMMU. Visit the [RHEL README](#) for additional information.

Setup

Setup Hugepages

Allocate a number of 2M Huge pages:

- For persistent allocation of huge pages, write to `hugepages.conf` file in `/etc/sysctl.d`:

```
$ echo 'vm.nr_hugepages=2048' > /etc/sysctl.d/hugepages.conf
```

- For run-time allocation of huge pages, use the `sysctl` utility:

```
$ sysctl -w vm.nr_hugepages=N # where N = No. of 2M huge pages
```

To verify hugepage configuration:

```
$ grep HugePages_ /proc/meminfo
```

Mount the hugepages, if not already mounted by default:

```
$ mount -t hugetlbfs none /dev/hugepages`
```

Note: The amount of hugepage memory required can be affected by various aspects of the datapath and device configuration. Refer to *DPDK Device Memory Models* for more details.

Setup DPDK devices using VFIO

VFIO is preferred to the UIO driver when using recent versions of DPDK. VFIO support requires support from both the kernel and BIOS. For the former, kernel version > 3.6 must be used. For the latter, you must enable VT-d in the BIOS and ensure this is configured via grub. To ensure VT-d is enabled via the BIOS, run:

```
$ dmesg | grep -e DMAR -e IOMMU
```

If VT-d is not enabled in the BIOS, enable it now.

To ensure VT-d is enabled in the kernel, run:

```
$ cat /proc/cmdline | grep iommu=pt
$ cat /proc/cmdline | grep intel_iommu=on
```

If VT-d is not enabled in the kernel, enable it now.

Once VT-d is correctly configured, load the required modules and bind the NIC to the VFIO driver:

```
$ modprobe vfio-pci
$ /usr/bin/chmod a+x /dev/vfio
$ /usr/bin/chmod 0666 /dev/vfio/*
$ $DPDK_DIR/usertools/dpdk-devbind.py --bind=vfio-pci eth1
$ $DPDK_DIR/usertools/dpdk-devbind.py --status
```

Setup OVS

Open vSwitch should be started as described in *Open vSwitch on Linux, FreeBSD and NetBSD* with the exception of `ovs-vswitchd`, which requires some special configuration to enable DPDK functionality. DPDK configuration arguments can be passed to `ovs-vswitchd` via the `other_config` column of the `Open_vSwitch` table. At a minimum, the `dpdk-init` option must be set to either `true` or `try`. For example:

```
$ export PATH=$PATH:/usr/local/share/openvswitch/scripts
$ export DB_SOCKET=/usr/local/var/run/openvswitch/db.sock
$ ovs-vsctl --no-wait set Open_vSwitch . other_config:dpdk-init=true
$ ovs-ctl --no-ovsdb-server --db-socket="$DB_SOCKET" start
```

There are many other configuration options, the most important of which are listed below. Defaults will be provided for all values not explicitly set.

dpdk-init Specifies whether OVS should initialize and support DPDK ports. This field can either be `true` or `try`. A value of `true` will cause the `ovs-vswitchd` process to abort on initialization failure. A value of `try` will imply that the `ovs-vswitchd` process should continue running even if the EAL initialization fails.

dpdk-lcore-mask Specifies the CPU cores on which dpdk lcore threads should be spawned and expects hex string (eg `'0x123'`).

dpdk-socket-mem Comma separated list of memory to pre-allocate from hugepages on specific sockets. If not specified, 1024 MB will be set for each numa node by default.

dpdk-hugepage-dir Directory where `hugetlbfs` is mounted

vhost-sock-dir Option to set the path to the vhost-user unix socket files.

If allocating more than one GB hugepage, you can configure the amount of memory used from any given NUMA nodes. For example, to use 1GB from NUMA node 0 and 0GB for all other NUMA nodes, run:

```
$ ovs-vsctl --no-wait set Open_vSwitch . \
    other_config:dpdk-socket-mem="1024,0"
```

or:

```
$ ovs-vsctl --no-wait set Open_vSwitch . \
    other_config:dpdk-socket-mem="1024"
```

Note: Changing any of these options requires restarting the ovs-vswitchd application

See the section [Performance Tuning](#) for important DPDK customizations.

Validating

DPDK support can be confirmed by validating the `dpdk_initialized` boolean value from the `ovsdb`. A value of `true` means that the DPDK EAL initialization succeeded:

```
$ ovs-vsctl get Open_vSwitch . dpdk_initialized
true
```

Additionally, the library version linked to `ovs-vswitchd` can be confirmed with either the `ovs-vswitchd` logs, or by running either of the commands:

```
$ ovs-vswitchd --version
ovs-vswitchd (Open vSwitch) 2.9.0
DPDK 17.11.0
$ ovs-vsctl get Open_vSwitch . dpdk_version
"DPDK 17.11.0"
```

At this point you can use `ovs-vsctl` to set up bridges and other Open vSwitch features. Seeing as we've configured the DPDK datapath, we will use DPDK-type ports. For example, to create a userspace bridge named `br0` and add two `dpdk` ports to it, run:

```
$ ovs-vsctl add-br br0 -- set bridge br0 datapath_type=netdev
$ ovs-vsctl add-port br0 myportnameone -- set Interface myportnameone \
    type=dpdk options:dpdk-devargs=0000:06:00.0
$ ovs-vsctl add-port br0 myportnametwo -- set Interface myportnametwo \
    type=dpdk options:dpdk-devargs=0000:06:00.1
```

DPDK devices will not be available for use until a valid `dpdk-devargs` is specified.

Refer to `ovs-vsctl(8)` and [Using Open vSwitch with DPDK](#) for more details.

Performance Tuning

To achieve optimal OVS performance, the system can be configured and that includes BIOS tweaks, Grub cmdline additions, better understanding of NUMA nodes and apt selection of PCIe slots for NIC placement.

Note: This section is optional. Once installed as described above, OVS with DPDK will work out of the box.

Recommended BIOS Settings

Table 1: Recommended BIOS Settings

Setting	Value
C3 Power State	Disabled
C6 Power State	Disabled
MLC Streamer	Enabled
MLC Spacial Prefetcher	Enabled
DCU Data Prefetcher	Enabled
DCA	Enabled
CPU Power and Performance	Performance
Memory RAS and Performance Config -> NUMA optimized	Enabled

PCIe Slot Selection

The fastpath performance can be affected by factors related to the placement of the NIC, such as channel speeds between PCIe slot and CPU or the proximity of PCIe slot to the CPU cores running the DPDK application. Listed below are the steps to identify right PCIe slot.

1. Retrieve host details using `dmidecode`. For example:

```
$ dmidecode -t baseboard | grep "Product Name"
```

2. Download the technical specification for product listed, e.g: S2600WT2
3. Check the Product Architecture Overview on the Riser slot placement, CPU sharing info and also PCIe channel speeds

For example: On S2600WT, CPU1 and CPU2 share Riser Slot 1 with Channel speed between CPU1 and Riser Slot1 at 32GB/s, CPU2 and Riser Slot1 at 16GB/s. Running DPDK app on CPU1 cores and NIC inserted in to Riser card Slots will optimize OVS performance in this case.

4. Check the Riser Card #1 - Root Port mapping information, on the available slots and individual bus speeds. In S2600WT slot 1, slot 2 has high bus speeds and are potential slots for NIC placement.

Advanced Hugepage Setup

Allocate and mount 1 GB hugepages.

- For persistent allocation of huge pages, add the following options to the kernel bootline:

```
default_hugepagesz=1GB hugepagesz=1G hugepages=N
```

For platforms supporting multiple huge page sizes, add multiple options:

```
default_hugepagesz=<size> hugepagesz=<size> hugepages=N
```

where:

N number of huge pages requested

size huge page size with an optional suffix [kKmMgG]

- For run-time allocation of huge pages:

```
$ echo N > /sys/devices/system/node/nodeX/hugepages/hugepages-1048576kB/nr_
↪ hugepages
```

where:

N number of huge pages requested

x NUMA Node

Note: For run-time allocation of 1G huge pages, Contiguous Memory Allocator (CONFIG_CMA) has to be supported by kernel, check your Linux distro.

Now mount the huge pages, if not already done so:

```
$ mount -t hugetlbfs -o pagesize=1G none /dev/hugepages
```

Isolate Cores

The `isolcpus` option can be used to isolate cores from the Linux scheduler. The isolated cores can then be used to dedicatedly run HPC applications or threads. This helps in better application performance due to zero context switching and minimal cache thrashing. To run platform logic on core 0 and isolate cores between 1 and 19 from scheduler, add `isolcpus=1-19` to GRUB cmdline.

Note: It has been verified that core isolation has minimal advantage due to mature Linux scheduler in some circumstances.

Compiler Optimizations

The default compiler optimization level is `-O2`. Changing this to more aggressive compiler optimization such as `-O3 -march=native` with gcc (verified on 5.3.1) can produce performance gains though not significant. `-march=native` will produce optimized code on local machine and should be used when software compilation is done on Testbed.

Multiple Poll-Mode Driver Threads

With pmd multi-threading support, OVS creates one pmd thread for each NUMA node by default, if there is at least one DPDK interface from that NUMA node added to OVS. However, in cases where there are multiple ports/rxq's producing traffic, performance can be improved by creating multiple pmd threads running on separate cores. These pmd threads can share the workload by each being responsible for different ports/rxq's. Assignment of ports/rxq's to pmd threads is done automatically.

A set bit in the mask means a pmd thread is created and pinned to the corresponding CPU core. For example, to run pmd threads on core 1 and 2:

```
$ ovs-vsctl set Open_vSwitch . other_config:pmc-cpu-mask=0x6
```

When using dpdk and dpdkvhostuser ports in a bi-directional VM loopback as shown below, spreading the workload over 2 or 4 pmc threads shows significant improvements as there will be more total CPU occupancy available:

```
NIC port0 <-> OVS <-> VM <-> OVS <-> NIC port 1
```

Refer to `ovs-vswitchd.conf.db(5)` for additional information on configuration options.

Affinity

For superior performance, DPDK pmc threads and Qemu vCPU threads needs to be affinityized accordingly.

- PMc thread Affinity

A poll mode driver (pmc) thread handles the I/O of all DPDK interfaces assigned to it. A pmc thread shall poll the ports for incoming packets, switch the packets and send to tx port. A pmc thread is CPU bound, and needs to be affinityized to isolated cores for optimum performance. Even though a PMc thread may exist, the thread only starts consuming CPU cycles if there is at least one receive queue assigned to the pmc.

Note: On NUMA systems, PCI devices are also local to a NUMA node. Unbound rx queues for a PCI device will be assigned to a pmc on it's local NUMA node if a non-isolated PMc exists on that NUMA node. If not, the queue will be assigned to a non-isolated pmc on a remote NUMA node. This will result in reduced maximum throughput on that device and possibly on other devices assigned to that pmc thread. If such a queue assignment is made a warning message will be logged: "There's no available (non-isolated) pmc thread on numa node N. Queue Q on port P will be assigned to the pmc on core C (numa node N'). Expect reduced performance."

Binding PMc threads to cores is described in the above section `Multiple Poll-Mode Driver Threads`.

- QEMU vCPU thread Affinity

A VM performing simple packet forwarding or running complex packet pipelines has to ensure that the vCPU threads performing the work has as much CPU occupancy as possible.

For example, on a multicore VM, multiple QEMU vCPU threads shall be spawned. When the DPDK `testpmc` application that does packet forwarding is invoked, the `taskset` command should be used to affinityize the vCPU threads to the dedicated isolated cores on the host system.

Enable HyperThreading

With HyperThreading, or SMT, enabled, a physical core appears as two logical cores. SMT can be utilized to spawn worker threads on logical cores of the same physical core thereby saving additional cores.

With DPDK, when pinning pmc threads to logical cores, care must be taken to set the correct bits of the `pmc-cpu-mask` to ensure that the pmc threads are pinned to SMT siblings.

Take a sample system configuration, with 2 sockets, 2 * 10 core processors, HT enabled. This gives us a total of 40 logical cores. To identify the physical core shared by two logical cores, run:

```
$ cat /sys/devices/system/cpu/cpuN/topology/thread_siblings_list
```

where N is the logical core number.

In this example, it would show that cores 1 and 21 share the same physical core. Logical cores can be specified in `pmd-cpu-masks` similarly to physical cores, as described in `Multiple Poll-Mode Driver Threads`.

NUMA/Cluster-on-Die

Ideally inter-NUMA datapaths should be avoided where possible as packets will go across QPI and there may be a slight performance penalty when compared with intra NUMA datapaths. On Intel Xeon Processor E5 v3, Cluster On Die is introduced on models that have 10 cores or more. This makes it possible to logically split a socket into two NUMA regions and again it is preferred where possible to keep critical datapaths within the one cluster.

It is good practice to ensure that threads that are in the datapath are pinned to cores in the same NUMA area. e.g. `pmd threads` and QEMU vCPUs responsible for forwarding. If DPDK is built with `CONFIG_RTE_LIBRTE_VHOST_NUMA=y`, vHost User ports automatically detect the NUMA socket of the QEMU vCPUs and will be serviced by a PMD from the same node provided a core on this node is enabled in the `pmd-cpu-mask`. `libnuma` packages are required for this feature.

Binding PMD threads is described in the above section `Multiple Poll-Mode Driver Threads`.

DPDK Physical Port Rx Queues

```
$ ovs-vsctl set Interface <DPDK interface> options:n_rxq=<integer>
```

The above command sets the number of rx queues for DPDK physical interface. The rx queues are assigned to `pmd` threads on the same NUMA node in a round-robin fashion.

DPDK Physical Port Queue Sizes

```
$ ovs-vsctl set Interface dpdk0 options:n_rxq_desc=<integer>
$ ovs-vsctl set Interface dpdk0 options:n_txq_desc=<integer>
```

The above command sets the number of rx/tx descriptors that the NIC associated with `dpdk0` will be initialised with.

Different `n_rxq_desc` and `n_txq_desc` configurations yield different benefits in terms of throughput and latency for different scenarios. Generally, smaller queue sizes can have a positive impact for latency at the expense of throughput. The opposite is often true for larger queue sizes. Note: increasing the number of rx descriptors eg. to 4096 may have a negative impact on performance due to the fact that non-vectorised DPDK rx functions may be used. This is dependent on the driver in use, but is true for the commonly used `i40e` and `ixgbe` DPDK drivers.

Exact Match Cache

Each `pmd` thread contains one Exact Match Cache (EMC). After initial flow setup in the datapath, the EMC contains a single table and provides the lowest level (fastest) switching for DPDK ports. If there is a miss in the EMC then the next level where switching will occur is the datapath classifier. Missing in the EMC and looking up in the datapath classifier incurs a significant performance penalty. If lookup misses occur in the EMC because it is too small to handle the number of flows, its size can be increased. The EMC size can be modified by editing the define `EM_FLOW_HASH_SHIFT` in `lib/dpif-netdev.c`.

As mentioned above, an EMC is per `pmd` thread. An alternative way of increasing the aggregate amount of possible flow entries in EMC and avoiding datapath classifier lookups is to have multiple `pmd` threads running.

Rx Mergeable Buffers

Rx mergeable buffers is a virtio feature that allows chaining of multiple virtio descriptors to handle large packet sizes. Large packets are handled by reserving and chaining multiple free descriptors together. Mergeable buffer support is negotiated between the virtio driver and virtio device and is supported by the DPDK vhost library. This behavior is supported and enabled by default, however in the case where the user knows that rx mergeable buffers are not needed i.e. jumbo frames are not needed, it can be forced off by adding `mrq_rxbuf=off` to the QEMU command line options. By not reserving multiple chains of descriptors it will make more individual virtio descriptors available for rx to the guest using dpdkvhost ports and this can improve performance.

Output Packet Batching

To make advantage of batched transmit functions, OVS collects packets in intermediate queues before sending when processing a batch of received packets. Even if packets are matched by different flows, OVS uses a single send operation for all packets destined to the same output port.

Furthermore, OVS is able to buffer packets in these intermediate queues for a configurable amount of time to reduce the frequency of send bursts at medium load levels when the packet receive rate is high, but the receive batch size still very small. This is particularly beneficial for packets transmitted to VMs using an interrupt-driven virtio driver, where the interrupt overhead is significant for the OVS PMD, the host operating system and the guest driver.

The `tx-flush-interval` parameter can be used to specify the time in microseconds OVS should wait between two send bursts to a given port (default is 0). When the intermediate queue fills up before that time is over, the buffered packet batch is sent immediately:

```
$ ovs-vsctl set Open_vSwitch . other_config:tx-flush-interval=50
```

This parameter influences both throughput and latency, depending on the traffic load on the port. In general lower values decrease latency while higher values may be useful to achieve higher throughput.

Low traffic ($\text{packet rate} < 1 / \text{tx-flush-interval}$) should not experience any significant latency or throughput increase as packets are forwarded immediately.

At intermediate load levels ($1 / \text{tx-flush-interval} < \text{packet rate} < 32 / \text{tx-flush-interval}$) traffic should experience an average latency increase of up to $1 / 2 * \text{tx-flush-interval}$ and a possible throughput improvement.

Very high traffic ($\text{packet rate} \gg 32 / \text{tx-flush-interval}$) should experience the average latency increase equal to $32 / (2 * \text{packet rate})$. Most send batches in this case will contain the maximum number of packets (32).

A `tx-burst-interval` value of 50 microseconds has shown to provide a good performance increase in a PHY-VM-PHY scenario on x86 system for interrupt-driven guests while keeping the latency increase at a reasonable level:

<https://mail.openvswitch.org/pipermail/ovs-dev/2017-December/341628.html>

Note: Throughput impact of this option significantly depends on the scenario and the traffic patterns. For example: `tx-burst-interval` value of 50 microseconds shows performance degradation in PHY-VM-PHY with bonded PHY scenario while testing with 256 – 1024 packet flows:

<https://mail.openvswitch.org/pipermail/ovs-dev/2017-December/341700.html>

The average number of packets per output batch can be checked in PMD stats:


```
$ ovs-appctl dpif-netdev/pmd-stats-show
```

Limitations

- **Network Interface Firmware requirements:** Each release of DPDK is validated against a specific firmware version for a supported Network Interface. New firmware versions introduce bug fixes, performance improvements and new functionality that DPDK leverages. The validated firmware versions are available as part of the release notes for DPDK. It is recommended that users update Network Interface firmware to match what has been validated for the DPDK release.

The latest list of validated firmware versions can be found in the [DPDK release notes](#).

- **Upper bound MTU:** DPDK device drivers differ in how the L2 frame for a given MTU value is calculated e.g. i40e driver includes 2 x vlan headers in MTU overhead, em driver includes 1 x vlan header, ixgbe driver does not include a vlan header in overhead. Currently it is not possible for OVS DPDK to know what upper bound MTU value is supported for a given device. As such OVS DPDK must provision for the case where the L2 frame for a given MTU includes 2 x vlan headers. This reduces the upper bound MTU value for devices that do not include vlan headers in their L2 frames by 8 bytes e.g. ixgbe devices upper bound MTU is reduced from 9710 to 9702. This work around is temporary and is expected to be removed once a method is provided by DPDK to query the upper bound MTU value for a given device.

Reporting Bugs

Report problems to bugs@openvswitch.org.

2.3.2 Installation from Packages

Open vSwitch is packaged on a variety of distributions. The tooling required to build these packages is included in the Open vSwitch tree. The instructions are provided below.

Distributions packaging Open vSwitch

This document lists various popular distributions packaging Open vSwitch. Open vSwitch is packaged by various distributions for multiple platforms and architectures.

Note: The packaged version available with distributions may not be latest Open vSwitch release.

Debian

You can use `apt-get` or `aptitude` to install the `.deb` packages and must be superuser.

1. Debian has `openvswitch-switch` and `openvswitch-common` `.deb` packages that includes the core userspace components of the switch.
2. For kernel datapath, `openvswitch-datapath-dkms` can be installed to automatically build and install Open vSwitch kernel module for your running kernel.
3. For DPDK datapath, Open vSwitch with DPDK support is bundled in the package `openvswitch-switch-dpdk`.

Fedora

Fedora provides `openvswitch`, `openvswitch-devel`, `openvswitch-test` and `openvswitch-debuginfo` rpm packages. You can install `openvswitch` package in minimum installation. Use `yum` or `dnf` to install the rpm packages and must be superuser.

Red Hat

RHEL distributes `openvswitch` rpm package that supports kernel datapath. DPDK accelerated Open vSwitch can be installed using `openvswitch-dpdk` package.

OpenSuSE

OpenSUSE provides `openvswitch`, `openvswitch-switch` rpm packages. Also `openvswitch-dpdk` and `openvswitch-dpdk-switch` can be installed for Open vSwitch using DPDK accelerated datapath.

Debian Packaging for Open vSwitch

This document describes how to build Debian packages for Open vSwitch. To install Open vSwitch on Debian without building Debian packages, refer to *Open vSwitch on Linux, FreeBSD and NetBSD* instead.

Note: These instructions should also work on Ubuntu and other Debian derivative distributions.

Before You Begin

Before you begin, consider whether you really need to build packages yourself. Debian “wheezy” and “sid”, as well as recent versions of Ubuntu, contain pre-built Debian packages for Open vSwitch. It is easier to install these than to build your own. To use packages from your distribution, skip ahead to “Installing .deb Packages”, below.

Building Open vSwitch Debian packages

You may build from an Open vSwitch distribution tarball or from an Open vSwitch Git tree with these instructions.

You do not need to be the superuser to build the Debian packages.

1. Install the “build-essential” and “fakeroot” packages. For example:

```
$ apt-get install build-essential fakeroot
```

2. Obtain and unpack an Open vSwitch source distribution and `cd` into its top level directory.
3. Install the build dependencies listed under “Build-Depends:” near the top of `debian/control`. You can install these any way you like, e.g. with `apt-get install`.

Check your work by running `dpkg-checkbuilddeps` in the top level of your OVS directory. If you’ve installed all the dependencies properly, `dpkg-checkbuilddeps` will exit without printing anything. If you forgot to install some dependencies, it will tell you which ones.

4. Build the package:

```
$ fakeroot debian/rules binary
```

This will do a serial build that runs the unit tests. This will take approximately 8 to 10 minutes. If you prefer, you can run a faster parallel build:

```
$ DEB_BUILD_OPTIONS='parallel=8' fakeroot debian/rules binary
```

If you are in a big hurry, you can even skip the unit tests:

```
$ DEB_BUILD_OPTIONS='parallel=8 nocheck' fakeroot debian/rules binary
```

Note: There are a few pitfalls in the Debian packaging building system so that, occasionally, you may find that in a tree that you have using for a while, the build command above exits immediately without actually building anything. To fix the problem, run:

```
$ fakeroot debian/rules clean
```

or start over from a fresh copy of the source tree.

5. The generated .deb files will be in the parent directory of the Open vSwitch source distribution.

Installing .deb Packages

These instructions apply to installing from Debian packages that you built yourself, as described in the previous section. In this case, use a command such as `dpkg -i` to install the .deb files that you build. You will have to manually install any missing dependencies.

You can also use these instruction to install from packages provided by Debian or a Debian derivative distribution such as Ubuntu. In this case, use a program such as `apt-get` or `aptitude` to download and install the provided packages. These programs will also automatically download and install any missing dependencies.

Important: You must be superuser to install Debian packages.

1. Start by installing an Open vSwitch kernel module. See `debian/openvswitch-switch.README`. Debian for the available options.
2. Install the `openvswitch-switch` and `openvswitch-common` packages. These packages include the core userspace components of the switch.

Open vSwitch .deb packages not mentioned above are rarely useful. Refer to their individual package descriptions to find out whether any of them are useful to you.

Reporting Bugs

Report problems to bugs@openvswitch.org.

Fedora, RHEL 7.x Packaging for Open vSwitch

This document provides instructions for building and installing Open vSwitch RPM packages on a Fedora Linux host. Instructions for the installation of Open vSwitch on a Fedora Linux host without using RPM packages can be found in the *Open vSwitch on Linux, FreeBSD and NetBSD*.

These instructions have been tested with Fedora 23, and are also applicable for RHEL 7.x and its derivatives, including CentOS 7.x and Scientific Linux 7.x.

Build Requirements

You will need to install all required packages to build the RPMs. Newer distributions use `dnf` but if it's not available, then use `yum` instructions.

The command below will install RPM tools and generic build dependencies. And (optionally) include these packages: `libcap-ng libcap-ng-devel dpdk-devel`.

DNF:

```
$ dnf install @'Development Tools' rpm-build dnf-plugins-core
```

YUM:

```
$ yum install @'Development Tools' rpm-build yum-utils
```

Then it is necessary to install Open vSwitch specific build dependencies. The dependencies are listed in the SPEC file, but first it is necessary to replace the `VERSION` tag to be a valid SPEC.

The command below will create a temporary SPEC file:

```
$ sed -e 's/@VERSION@/0.0.1/' rhel/openvswitch-fedora.spec.in \  
> /tmp/ovs.spec
```

And to install specific dependencies, use the corresponding tool below. For some of the dependencies on RHEL you may need to add two additional repositories to help `yum-builddep`, e.g.:

```
$ subscription-manager repos --enable=rhel-7-server-extras-rpms  
$ subscription-manager repos --enable=rhel-7-server-optional-rpms
```

DNF:

```
$ dnf builddep /tmp/ovs.spec
```

YUM:

```
$ yum-builddep /tmp/ovs.spec
```

Once that is completed, remove the file `/tmp/ovs.spec`.

Bootstrapping

Refer to *Bootstrapping*.

Configuring

Refer to *Configuring*.

Building

User Space RPMs

To build Open vSwitch user-space RPMs, execute the following from the directory in which *./configure* was executed:

```
$ make rpm-fedora
```

This will create the RPMs *openvswitch*, *python-openvswitch*, *openvswitch-test*, *openvswitch-devel* and *openvswitch-debuginfo*.

To enable DPDK support in the *openvswitch* package, the `--with dpdk` option can be added:

```
$ make rpm-fedora RPMBUILD_OPT="--with dpdk --without check"
```

You can also have the above commands automatically run the Open vSwitch unit tests. This can take several minutes.

```
$ make rpm-fedora RPMBUILD_OPT="--with check"
```

To build OVN RPMs, execute the following from the directory in which *./configure* was executed:

```
$ make rpm-fedora-ovn
```

This will create the RPMs *ovn*, *ovn-common*, *ovn-central*, *ovn-host*, *ovn-docker* and *ovn-vtep*.

Kernel OVS Tree Datapath RPM

To build the Open vSwitch kernel module for the currently running kernel version, run:

```
$ make rpm-fedora-kmod
```

To build the Open vSwitch kernel module for another kernel version, the desired kernel version can be specified via the *kversion* macro. For example:

```
$ make rpm-fedora-kmod \
  RPMBUILD_OPT='-D "kversion 4.3.4-300.fc23.x86_64"'
```

Installing

RPM packages can be installed by using the command `rpm -i`. Package installation requires superuser privileges.

The *openvswitch-kmod* RPM should be installed first if the Linux OVS tree datapath module is to be used. The *openvswitch-kmod* RPM should not be installed if only the in-tree Linux datapath or user-space datapath is needed. Refer to the [Open vSwitch FAQ](#) for more information about the various Open vSwitch datapath options.

In most cases only the *openvswitch* RPM will need to be installed. The *python-openvswitch*, *openvswitch-test*, *openvswitch-devel*, and *openvswitch-debuginfo* RPMs are optional unless required for a specific purpose.

The *ovn-** packages are only needed when using OVN.

Refer to the [RHEL README](#) for additional usage and configuration information.

Reporting Bugs

Report problems to bugs@openvswitch.org.

RHEL 5.6, 6.x Packaging for Open vSwitch

This document describes how to build and install Open vSwitch on a Red Hat Enterprise Linux (RHEL) host. If you want to install Open vSwitch on a generic Linux host, refer to [Open vSwitch on Linux, FreeBSD and NetBSD](#) instead.

We have tested these instructions with RHEL 5.6 and RHEL 6.0.

For RHEL 7.x (or derivatives, such as CentOS 7.x), you should follow the instructions in the [Fedora, RHEL 7.x Packaging for Open vSwitch](#). The Fedora spec files are used for RHEL 7.x.

Prerequisites

You may build from an Open vSwitch distribution tarball or from an Open vSwitch Git tree.

The default RPM build directory, `_topdir`, has five directories in the top-level.

BUILD/ where the software is unpacked and built

RPMS/ where the newly created binary package files are written

SOURCES/ contains the original sources, patches, and icon files

SPECS/ contains the spec files for each package to be built

SRPMS/ where the newly created source package files are written

Before you begin, note the RPM sources directory on your version of RHEL. The command `rpmbuild --showrc` will show the configuration for each of those directories. Alternatively, the command `rpm --eval '%{_topdir}'` shows the current configuration for the top level directory and the command `rpm --eval '%{_sourcedir}'` does the same for the sources directory. On RHEL 5, the default RPM `_topdir` is `/usr/src/redhat` and the default RPM sources directory is `/usr/src/redhat/SOURCES`. On RHEL 6, the default `_topdir` is `$HOME/rpmbuild` and the default RPM sources directory is `$HOME/rpmbuild/SOURCES`.

Build Requirements

You will need to install all required packages to build the RPMs. The command below will install RPM tools and generic build dependencies:

```
$ yum install @'Development Tools' rpm-build yum-utils
```

Then it is necessary to install Open vSwitch specific build dependencies. The dependencies are listed in the SPEC file, but first it is necessary to replace the VERSION tag to be a valid SPEC.

The command below will create a temporary SPEC file:

```
$ sed -e 's/@VERSION@/0.0.1/' rhel/openvswitch.spec.in > /tmp/ovs.spec
```

And to install specific dependencies, use `yum-builddep` tool:

```
$ yum-builddep /tmp/ovs.spec
```

Once that is completed, remove the file `/tmp/ovs.spec`.

If `python-sphinx` package is not available in your version of RHEL, you can install it via pip with `'pip install sphinx'`.

Open vSwitch requires python 2.7 or newer which is not available in older distributions. In the case of RHEL 6.x and its derivatives, one option is to install `python34` and `python34-six` from [EPEL](#).

Bootstrapping and Configuring

If you are building from a distribution tarball, skip to [Building](#). If not, you must be building from an Open vSwitch Git tree. Determine what version of Autoconf is installed (e.g. run `autoconf --version`). If it is not at least version 2.63, then you must upgrade or use another machine to build the packages.

Assuming all requirements have been met, build the tarball by running:

```
$ ./boot.sh
$ ./configure
$ make dist
```

You must run this on a machine that has the tools listed in [Build Requirements](#) as prerequisites for building from a Git tree. Afterward, proceed with the rest of the instructions using the distribution tarball.

Now you have a distribution tarball, named something like `openvswitch-x.y.z.tar.gz`. Copy this file into the RPM sources directory, e.g.:

```
$ cp openvswitch-x.y.z.tar.gz $HOME/rpmbuild/SOURCES
```

Broken build symlink

Some versions of the RHEL 6 kernel-devel package contain a broken `build` symlink. If you are using such a version, you must fix the problem before continuing.

To find out whether you are affected, run:

```
$ cd /lib/modules/<version>
$ ls -l build/
```

where `<version>` is the version number of the RHEL 6 kernel.

Note: The trailing slash in the final command is important. Be sure to include it.

If the `ls` command produces a directory listing, your kernel-devel package is OK. If it produces a `No such file or directory` error, your kernel-devel package is buggy.

If your kernel-devel package is buggy, then you can fix it with:

```
$ cd /lib/modules/<version>
$ rm build
$ ln -s /usr/src/kernels/<target> build
```

where `<target>` is the name of an existing directory under `/usr/src/kernels`, whose name should be similar to `<version>` but may contain some extra parts. Once you have done this, verify the fix with the same procedure you used above to check for the problem.

Building

You should have a distribution tarball named something like `openvswitch-x.y.z.tar.gz`. Copy this file into the RPM sources directory:

```
$ cp openvswitch-x.y.z.tar.gz $HOME/rpmbuild/SOURCES
```

Make another copy of the distribution tarball in a temporary directory. Then unpack the tarball and `cd` into its root:

```
$ tar xzf openvswitch-x.y.z.tar.gz
$ cd openvswitch-x.y.z
```

Userspace

To build Open vSwitch userspace, run:

```
$ rpmbuild -bb rhel/openvswitch.spec
```

This produces two RPMs: “openvswitch” and “openvswitch-debuginfo”.

The above command automatically runs the Open vSwitch unit tests. To disable the unit tests, run:

```
$ rpmbuild -bb --without check rhel/openvswitch.spec
```

Note: If the build fails with `configure: error: source dir /lib/modules/2.6.32-279.el6.x86_64/build doesn't exist` or similar, then the kernel-devel package is missing or buggy.

Kernel Module

On RHEL 6, to build the Open vSwitch kernel module run:

```
$ rpmbuild -bb rhel/kmod-openvswitch-rhel6.spec
```

You might have to specify a kernel version and/or variants, e.g.:

```
$ rpmbuild -bb -D “kversion 2.6.32-131.6.1.el6.x86_64” -D “kflavors default debug kdump”
rhel/kmod-openvswitch-rhel6.spec
```

This produces an “kmod-openvswitch” RPM for each kernel variant, in this example: “kmod-openvswitch”, “kmod-openvswitch-debug”, and “kmod-openvswitch-kdump”.

Red Hat Network Scripts Integration

A RHEL host has default firewall rules that prevent any Open vSwitch tunnel traffic from passing through. If a user configures Open vSwitch tunnels like Geneve, GRE, VXLAN, LISP etc., they will either have to manually add iptables firewall rules to allow the tunnel traffic or add it through a startup script. Refer to the “enable-protocol” command in the `ovs-ctl(8)` manpage for more information.

In addition, simple integration with Red Hat network scripts has been implemented. Refer to [README.RHEL.rst](#) in the source tree or `/usr/share/doc/openvswitch/README.RHEL.rst` in the installed openvswitch package for details.

Reporting Bugs

Report problems to bugs@openvswitch.org.

2.3.3 Upgrades

OVN Upgrades

Since OVN is a distributed system, special consideration must be given to the process used to upgrade OVN across a deployment. This document discusses the recommended upgrade process.

Release Notes

You should always check the OVS and OVN release notes (NEWS file) for any release specific notes on upgrades.

OVS

OVN depends on and is included with OVS. It's expected that OVS and OVN are upgraded together, partly for convenience. OVN is included in OVS releases so it's easiest to upgrade them together. OVN may also make use of new features of OVS only available in that release.

Upgrade ovn-controller

You should start by upgrading ovn-controller on each host it's running on. First, you upgrade the OVS and OVN packages. Then, restart the ovn-controller service. You can restart with ovn-ctl:

```
$ sudo /usr/share/openvswitch/scripts/ovn-ctl restart_controller
```

or with systemd:

```
$ sudo systemctl restart ovn-controller
```

Upgrade OVN Databases and ovn-northd

The OVN databases and ovn-northd should be upgraded next. Since ovn-controller has already been upgraded, it will be ready to operate on any new functionality specified by the database or logical flows created by ovn-northd.

Upgrading the OVN packages installs everything needed for an upgrade. The only step required after upgrading the packages is to restart ovn-northd, which automatically restarts the databases and upgrades the database schema, as well.

You may perform this restart using the ovn-ctl script:

```
$ sudo /usr/share/openvswitch/scripts/ovn-ctl restart_northd
```

or if you're using a Linux distribution with systemd:

```
$ sudo systemctl restart ovn-northd
```

Schema Change

During database upgrading, if there is schema change, the DB file will be converted to the new schema automatically, if the schema change is backward compatible. OVN tries the best to keep the DB schemas backward compatible.

However, there can be situations that an incompatible change is reasonable. An example of such case is to add constraints in the table to ensure correctness. If there were already data that violates the new constraints got added somehow, it will result in DB upgrade failures. In this case, user should manually correct data using `ovn-nbctl` (for north-bound DB) or `ovn-sbctl` (for south-bound DB), and then upgrade again following previous steps. Below is a list of known impactable schema changes and how to fix when error encountered.

1. Release 2.11: index [type, ip] added for Encap table of south-bound DB to prevent duplicated IPs being used for same tunnel type. If there are duplicated data added already (e.g. due to improper chassis management), a convenient way to fix is to find the chassis that is using the IP with command:

```
$ ovn-sbctl show
```

Then delete the chassis with command:

```
$ ovn-sbctl chassis-del <chassis>
```

Upgrading OVN Integration

Lastly, you may also want to upgrade integration with OVN that you may be using. For example, this could be the OpenStack Neutron driver or `ovn-kubernetes`.

OVN's northbound database schema is a backwards compatible interface, so you should be able to safely complete an OVN upgrade before upgrading any integration in use.

2.3.4 Others

Bash command-line completion scripts

There are two completion scripts available: `ovs-appctl-bashcomp.bash` and `ovs-vsctl-bashcomp.bash`.

ovs-appctl-bashcomp

`ovs-appctl-bashcomp.bash` adds bash command-line completion support for `ovs-appctl`, `ovs-dpctl`, `ovs-ofctl` and `ovsdb-tool` commands.

Features

- Display available completion or complete on unfinished user input (long option, subcommand, and argument).
- Subcommand hints
- Convert between keywords like `bridge`, `port`, `interface`, or `dp` and the available record in `ovsdb`.

Limitations

- Only supports a small set of important keywords (`dp`, `datapath`, `bridge`, `switch`, `port`, `interface`, `iface`).
- Does not support parsing of nested options. For example:

```
$ ovsdb-tool create [db [schema]]
```

- Does not support expansion on repeated argument. For example:

```
$ ovs-dpctl show [dp...]).
```

- Only supports matching on long options, and only in the format `--option [arg]`. Do not use `--option=[arg]`.

ovs-vsctl-bashcomp

`ovs-vsctl-bashcomp.bash` adds Bash command-line completion support for `ovs-vsctl` command.

Features

- Display available completion and complete on user input for global/local options, command, and argument.
- Query database and expand keywords like `table`, `record`, `column`, or `key`, to available completions.
- Deal with argument relations like ‘one and more’, ‘zero or one’.
- Complete multiple `ovs-vsctl` commands cascaded via `--`.

Limitations

Completion of very long `ovs-vsctl` commands can take up to several seconds.

Usage

The `bashcomp` scripts should be placed at `/etc/bash_completion.d/` to be available for all bash sessions. Running `make install` will place the scripts to `$(sysconfdir)/bash_completion.d/`, thus, the user should specify `--sysconfdir=/etc` at configuration. If OVS is installed from packages, the scripts will automatically be placed inside `/etc/bash_completion.d/`.

If you just want to run the scripts in one bash, you can remove them from `/etc/bash_completion.d/` and run the scripts via `. ovs-appctl-bashcomp.bash` or `. ovs-vsctl-bashcomp.bash`.

Tests

Unit tests are added in `tests/completion.at` and integrated into autotest framework. To run the tests, just run `make check`.

Open vSwitch Documentation

This document describes how to build the OVS documentation for use offline. A continuously updated, online version can be found at docs.openvswitch.org.

Note: These instructions provide information on building the documentation locally. For information on writing documentation, refer to *Open vSwitch Documentation Style*

Build Requirements

As described in the *Open vSwitch Documentation Style*, the Open vSwitch documentation is written in reStructured-Text and built with Sphinx. A detailed guide on installing Sphinx in many environments is available on the [Sphinx website](#) but, for most Linux distributions, you can install with your package manager. For example, on Debian/Ubuntu run:

```
$ sudo apt-get install python-sphinx
```

Similarly, on RHEL/Fedora run:

```
$ sudo dnf install python-sphinx
```

A `requirements.txt` is also provided in the `/Documentation`, should you wish to install using `pip`:

```
$ virtualenv .venv
$ source .venv/bin/activate
$ pip install -r Documentation/requirements.txt
```

Configuring

It's unlikely that you'll need to customize any aspect of the configuration. However, the `Documentation/conf.py` is the go-to place for all configuration. This file is well documented and further information is available on the [Sphinx website](#).

Building

Once Sphinx installed, the documentation can be built using the provided Makefile targets:

```
$ make docs-check
```

Important: The `docs-check` target will fail if there are any syntax errors. However, it won't catch more succinct issues such as style or grammar issues. As a result, you should always inspect changes visually to ensure the result is as intended.

Once built, documentation is available in the `/Documentation/_build` folder. Open the root `index.html` to browse the documentation.

Getting started with Open vSwitch (OVS) and Open Virtual Network (OVN) for Open vSwitch.

3.1 OVS Faucet Tutorial

This tutorial demonstrates how Open vSwitch works with a general-purpose OpenFlow controller, using the Faucet controller as a simple way to get started. It was tested with the “master” branch of Open vSwitch and version 1.6.15 of Faucet. It does not use advanced or recently added features in OVS or Faucet, so other versions of both pieces of software are likely to work equally well.

The goal of the tutorial is to demonstrate Open vSwitch and Faucet in an end-to-end way, that is, to show how it works from the Faucet controller configuration at the top, through the OpenFlow flow table, to the datapath processing. Along the way, in addition to helping to understand the architecture at each level, we discuss performance and troubleshooting issues. We hope that this demonstration makes it easier for users and potential users to understand how Open vSwitch works and how to debug and troubleshoot it.

We provide enough details in the tutorial that you should be able to fully follow along by following the instructions.

3.1.1 Setting Up OVS

This section explains how to set up Open vSwitch for the purpose of using it with Faucet for the tutorial.

You might already have Open vSwitch installed on one or more computers or VMs, perhaps set up to control a set of VMs or a physical network. This is admirable, but we will be using Open vSwitch in a different way to set up a simulation environment called the OVS “sandbox”. The sandbox does not use virtual machines or containers, which makes it more limited, but on the other hand it is (in this writer’s opinion) easier to set up.

There are two ways to start a sandbox: one that uses the Open vSwitch that is already installed on a system, and another that uses a copy of Open vSwitch that has been built but not yet installed. The latter is more often used and thus better tested, but both should work. The instructions below explain both approaches:

1. Get a copy of the Open vSwitch source repository using Git, then `cd` into the new directory:

```
$ git clone https://github.com/openvswitch/ovs.git
$ cd ovs
```

The default checkout is the master branch. You can check out a tag (such as v2.8.0) or a branch (such as origin/branch-2.8), if you prefer.

2. If you do not already have an installed copy of Open vSwitch on your system, or if you do not want to use it for the sandbox (the sandbox will not disturb the functionality of any existing switches), then proceed to step 3. If you do have an installed copy and you want to use it for the sandbox, try to start the sandbox by running:

```
$ tutorial/ovs-sandbox
```

If it is successful, you will find yourself in a subshell environment, which is the sandbox (you can exit with `exit` or Control+D). If so, you're finished and do not need to complete the rest of the steps. If it fails, you can proceed to step 3 to build Open vSwitch anyway.

3. Before you build, you might want to check that your system meets the build requirements. Read [Open vSwitch on Linux, FreeBSD and NetBSD](#) to find out. For this tutorial, there is no need to compile the Linux kernel module, or to use any of the optional libraries such as OpenSSL, DPDK, or libcap-ng.
4. Configure and build Open vSwitch:

```
$ ./boot.sh
$ ./configure
$ make -j4
```

5. Try out the sandbox by running:

```
$ make sandbox
```

You can exit the sandbox with `exit` or Control+D.

3.1.2 Setting up Faucet

This section explains how to get a copy of Faucet and set it up appropriately for the tutorial. There are many other ways to install Faucet, but this simple approach worked well for me. It has the advantage that it does not require modifying any system-level files or directories on your machine. It does, on the other hand, require Docker, so make sure you have it installed and working.

It will be a little easier to go through the rest of the tutorial if you run these instructions in a separate terminal from the one that you're using for Open vSwitch, because it's often necessary to switch between one and the other.

1. Get a copy of the Faucet source repository using Git, then `cd` into the new directory:

```
$ git clone https://github.com/faucetsdn/faucet.git
$ cd faucet
```

At this point I checked out the latest tag:

```
$ latest_tag=$(git describe --tags $(git rev-list --tags --max-count=1))
$ git checkout $latest_tag
```

2. Build a docker container image:

```
$ docker build -t faucet/faucet .
```

This will take a few minutes.

3. Create an installation directory under the `faucet` directory for the docker image to use:

```
$ mkdir inst
```

The Faucet configuration will go in `inst/faucet.yaml` and its main log will appear in `inst/faucet.log`. (The official Faucet installation instructions call to put these in `/etc/ryu/faucet` and `/var/log/ryu/faucet`, respectively, but we avoid modifying these system directories.)

4. Create a container and start Faucet:

```
$ docker run -d --name faucet --restart=always -v $(pwd)/inst:/etc/faucet/ -v
↪$(pwd)/inst:/var/log/faucet/ -p 6653:6653 -p 9302:9302 faucet/faucet
```

5. Look in `inst/faucet.log` to verify that Faucet started. It will probably start with an exception and trace-back because we have not yet created `inst/faucet.yaml`.
6. Later on, to make a new or updated Faucet configuration take effect quickly, you can run:

```
$ docker exec faucet pkill -HUP -f faucet.faucet
```

Another way is to stop and start the Faucet container:

```
$ docker restart faucet
```

You can also stop and delete the container; after this, to start it again, you need to rerun the `docker run` command:

```
$ docker stop faucet
$ docker rm faucet
```

3.1.3 Overview

Now that Open vSwitch and Faucet are ready, here's an overview of what we're going to do for the remainder of the tutorial:

1. Switching: Set up an L2 network with Faucet.
2. Routing: Route between multiple L3 networks with Faucet.
3. ACLs: Add and modify access control rules.

At each step, we will take a look at how the features in question work from Faucet at the top to the data plane layer at the bottom. From the highest to lowest level, these layers and the software components that connect them are:

Faucet. As the top level in the system, this is the authoritative source of the network configuration.

Faucet connects to a variety of monitoring and performance tools, but we won't use them in this tutorial. Our main insights into the system will be through `faucet.yaml` for configuration and `faucet.log` to observe state, such as MAC learning and ARP resolution, and to tell when we've screwed up configuration syntax or semantics.

The OpenFlow subsystem in Open vSwitch. OpenFlow is the protocol, standardized by the Open Networking Foundation, that controllers like Faucet use to control how Open vSwitch and other switches treat packets in the network.

We will use `ovs-ofctl`, a utility that comes with Open vSwitch, to observe and occasionally modify Open vSwitch's OpenFlow behavior. We will also use `ovs-appctl`, a utility for communicating with `ovs-vswitchd` and other Open vSwitch daemons, to ask "what-if?" type questions.

In addition, the OVS sandbox by default raises the Open vSwitch logging level for OpenFlow high enough that we can learn a great deal about OpenFlow behavior simply by reading its log file.

Open vSwitch datapath. This is essentially a cache designed to accelerate packet processing. Open vSwitch includes a few different datapaths, such as one based on the Linux kernel and a userspace-only datapath (sometimes called the “DPDK” datapath). The OVS sandbox uses the latter, but the principles behind it apply equally well to other datapaths.

At each step, we discuss how the design of each layer influences performance. We demonstrate how Open vSwitch features can be used to debug, troubleshoot, and understand the system as a whole.

3.1.4 Switching

Layer-2 (L2) switching is the basis of modern networking. It’s also very simple and a good place to start, so let’s set up a switch with some VLANs in Faucet and see how it works at each layer. Begin by putting the following into `inst/faucet.yaml`:

```
dps:
  switch-1:
    dp_id: 0x1
    timeout: 3600
    arp_neighbor_timeout: 3600
    interfaces:
      1:
        native_vlan: 100
      2:
        native_vlan: 100
      3:
        native_vlan: 100
      4:
        native_vlan: 200
      5:
        native_vlan: 200
vlangs:
  100:
  200:
```

This configuration file defines a single switch (“datapath” or “dp”) named `switch-1`. The switch has five ports, numbered 1 through 5. Ports 1, 2, and 3 are in VLAN 100, and ports 4 and 5 are in VLAN 2. Faucet can identify the switch from its datapath ID, which is defined to be `0x1`.

Note: This also sets high MAC learning and ARP timeouts. The defaults are 5 minutes and about 8 minutes, which are fine in production but sometimes too fast for manual experimentation. (Don’t use a timeout bigger than about 65000 seconds because it will crash Faucet.)

Now restart Faucet so that the configuration takes effect, e.g.:

```
$ docker restart faucet
```

Assuming that the configuration update is successful, you should now see a new line at the end of `inst/faucet.log`:

```
Jan 06 15:14:35 faucet INFO      Add new datapath DPID 1 (0x1)
```

Faucet is now waiting for a switch with datapath ID `0x1` to connect to it over OpenFlow, so our next step is to create a switch with OVS and make it connect to Faucet. To do that, switch to the terminal where you checked out OVS and

start a sandbox with `make sandbox` or `tutorial/ovs-sandbox` (as explained earlier under *Setting Up OVS*). You should see something like this toward the end of the output:

```
-----
You are running in a dummy Open vSwitch environment.  You can use
ovs-vsctl, ovs-ofctl, ovs-appctl, and other tools to work with the
dummy switch.
```

```
Log files, pidfiles, and the configuration database are in the
"sandbox" subdirectory.
```

```
Exit the shell to kill the running daemons.
blp@sigabrt:~/nicira/ovs/tutorial(0)$
```

Inside the sandbox, create a switch (“bridge”) named `br0`, set its datapath ID to `0x1`, add simulated ports to it named `p1` through `p5`, and tell it to connect to the Faucet controller. To make it easier to understand, we request for port `p1` to be assigned OpenFlow port 1, `p2` port 2, and so on. As a final touch, configure the controller to be “out-of-band” (this is mainly to avoid some annoying messages in the `ovs-vswitchd` logs; for more information, run `man ovs-vswitchd.conf.db` and search for `connection_mode`):

```
$ ovs-vsctl add-br br0 \
    -- set bridge br0 other-config:datapath-id=0000000000000001 \
    -- add-port br0 p1 -- set interface p1 ofport_request=1 \
    -- add-port br0 p2 -- set interface p2 ofport_request=2 \
    -- add-port br0 p3 -- set interface p3 ofport_request=3 \
    -- add-port br0 p4 -- set interface p4 ofport_request=4 \
    -- add-port br0 p5 -- set interface p5 ofport_request=5 \
    -- set-controller br0 tcp:127.0.0.1:6653 \
    -- set controller br0 connection-mode=out-of-band
```

Note: You don’t have to run all of these as a single `ovs-vsctl` invocation. It is a little more efficient, though, and since it updates the OVS configuration in a single database transaction it means that, for example, there is never a time when the controller is set but it has not yet been configured as out-of-band.

Now, if you look at `inst/faucet.log` again, you should see that Faucet recognized and configured the new switch and its ports:

```
Jan 06 15:17:10 faucet      INFO      DPID 1 (0x1) connected
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Cold start configuring DP
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Configuring VLAN 100 vid:100
↪ports:Port 1,Port 2,Port 3
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Configuring VLAN 200 vid:200
↪ports:Port 4,Port 5
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Port 1 up, configuring
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Port 2 up, configuring
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Port 3 up, configuring
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Port 4 up, configuring
Jan 06 15:17:10 faucet.valve INFO      DPID 1 (0x1) Port 5 up, configuring
```

Over on the Open vSwitch side, you can see a lot of related activity if you take a look in `sandbox/ovs-vswitchd.log`. For example, here is the basic OpenFlow session setup and Faucet’s probe of the switch’s ports and capabilities:

```
rconn|INFO|br0<->tcp:127.0.0.1:6653: connecting...
vconn|DBG|tcp:127.0.0.1:6653: sent (Success): OFPT_HELLO (OF1.4) (xid=0x1):
  version bitmap: 0x01, 0x02, 0x03, 0x04, 0x05
```

(continues on next page)

(continued from previous page)

```

vconn|DBG|tcp:127.0.0.1:6653: received: OFPT_HELLO (OF1.3) (xid=0x2f24810a):
  version bitmap: 0x01, 0x02, 0x03, 0x04
vconn|DBG|tcp:127.0.0.1:6653: negotiated OpenFlow version 0x04 (we support version_
↳0x05 and earlier, peer supports version 0x04 and earlier)
rconn|INFO|br0<->tcp:127.0.0.1:6653: connected
vconn|DBG|tcp:127.0.0.1:6653: received: OFPT_ECHO_REQUEST (OF1.3) (xid=0x2f24810b): 0_
↳bytes of payload
vconn|DBG|tcp:127.0.0.1:6653: sent (Success): OFPT_ECHO_REPLY (OF1.3)_
↳(xid=0x2f24810b): 0 bytes of payload
vconn|DBG|tcp:127.0.0.1:6653: received: OFPT_FEATURES_REQUEST (OF1.3)_
↳(xid=0x2f24810c):
vconn|DBG|tcp:127.0.0.1:6653: sent (Success): OFPT_FEATURES_REPLY (OF1.3)_
↳(xid=0x2f24810c): dpid:0000000000000001
  n_tables:254, n_buffers:0
  capabilities: FLOW_STATS TABLE_STATS PORT_STATS GROUP_STATS QUEUE_STATS
vconn|DBG|tcp:127.0.0.1:6653: received: OFPST_PORT_DESC request (OF1.3)_
↳(xid=0x2f24810d): port=ANY
vconn|DBG|tcp:127.0.0.1:6653: sent (Success): OFPST_PORT_DESC reply (OF1.3)_
↳(xid=0x2f24810d):
  1(p1): addr:aa:55:aa:55:00:14
    config: PORT_DOWN
    state: LINK_DOWN
    speed: 0 Mbps now, 0 Mbps max
  2(p2): addr:aa:55:aa:55:00:15
    config: PORT_DOWN
    state: LINK_DOWN
    speed: 0 Mbps now, 0 Mbps max
  3(p3): addr:aa:55:aa:55:00:16
    config: PORT_DOWN
    state: LINK_DOWN
    speed: 0 Mbps now, 0 Mbps max
  4(p4): addr:aa:55:aa:55:00:17
    config: PORT_DOWN
    state: LINK_DOWN
    speed: 0 Mbps now, 0 Mbps max
  5(p5): addr:aa:55:aa:55:00:18
    config: PORT_DOWN
    state: LINK_DOWN
    speed: 0 Mbps now, 0 Mbps max
LOCAL(br0): addr:c6:64:ff:59:48:41
  config: PORT_DOWN
  state: LINK_DOWN
  speed: 0 Mbps now, 0 Mbps max

```

After that, you can see Faucet delete all existing flows and then start adding new ones:

```

vconn|DBG|tcp:127.0.0.1:6653: received: OFPT_FLOW_MOD (OF1.3) (xid=0x2f24810e): DEL_
↳table:255 priority=0 actions=drop
vconn|DBG|tcp:127.0.0.1:6653: received: OFPT_BARRIER_REQUEST (OF1.3) (xid=0x2f24810f):
vconn|DBG|tcp:127.0.0.1:6653: sent (Success): OFPT_BARRIER_REPLY (OF1.3)_
↳(xid=0x2f24810f):
vconn|DBG|tcp:127.0.0.1:6653: received: OFPT_FLOW_MOD (OF1.3) (xid=0x2f248110): ADD_
↳priority=0 cookie:0x5adc15c0 out_port:0 actions=drop
vconn|DBG|tcp:127.0.0.1:6653: received: OFPT_FLOW_MOD (OF1.3) (xid=0x2f248111): ADD_
↳table:1 priority=0 cookie:0x5adc15c0 out_port:0 actions=drop
...

```

OpenFlow Layer

Let's take a look at the OpenFlow tables that Faucet set up. Before we do that, it's helpful to take a look at `docs/architecture.rst` in the Faucet documentation to learn how Faucet structures its flow tables. In summary, this document says:

Table 0 Port-based ACLs

Table 1 Ingress VLAN processing

Table 2 VLAN-based ACLs

Table 3 Ingress L2 processing, MAC learning

Table 4 L3 forwarding for IPv4

Table 5 L3 forwarding for IPv6

Table 6 Virtual IP processing, e.g. for router IP addresses implemented by Faucet

Table 7 Egress L2 processing

Table 8 Flooding

With that in mind, let's dump the flow tables. The simplest way is to just run plain `ovs-ofctl dump-flows`:

```
$ ovs-ofctl dump-flows br0
```

If you run that bare command, it produces a lot of extra junk that makes the output harder to read, like statistics and "cookie" values that are all the same. In addition, for historical reasons `ovs-ofctl` always defaults to using OpenFlow 1.0 even though Faucet and most modern controllers use OpenFlow 1.3, so it's best to force it to use OpenFlow 1.3. We could throw in a lot of options to fix these, but we'll want to do this more than once, so let's start by defining a shell function for ourselves:

```
$ dump-flows () {
  ovs-ofctl -OOpenFlow13 --names --no-stat dump-flows "$@" \
    | sed 's/cookie=0x5adc15c0, //'
}
```

Let's also define `save-flows` and `diff-flows` functions for later use:

```
$ save-flows () {
  ovs-ofctl -OOpenFlow13 --no-names --sort dump-flows "$@"
}
$ diff-flows () {
  ovs-ofctl -OOpenFlow13 diff-flows "$@" | sed 's/cookie=0x5adc15c0 //'
}
```

Now let's take a look at the flows we've got and what they mean, like this:

```
$ dump-flows br0
```

First, table 0 has a flow that just jumps to table 1 for each configured port, and drops other unrecognized packets. Presumably it will do more if we configured port-based ACLs:

```
priority=9099,in_port=p1 actions=goto_table:1
priority=9099,in_port=p2 actions=goto_table:1
priority=9099,in_port=p3 actions=goto_table:1
priority=9099,in_port=p4 actions=goto_table:1
priority=9099,in_port=p5 actions=goto_table:1
priority=0 actions=drop
```

Table 1, for ingress VLAN processing, has a bunch of flows that drop inappropriate packets, such as LLDP and STP:

```
table=1, priority=9099,dl_dst=01:80:c2:00:00:00 actions=drop
table=1, priority=9099,dl_dst=01:00:0c:cc:cc:cd actions=drop
table=1, priority=9099,dl_type=0x88cc actions=drop
```

Table 1 also has some more interesting flows that recognize packets without a VLAN header on each of our ports (vlan_tci=0x0000/0x1fff), push on the VLAN configured for the port, and proceed to table 3. Presumably these skip table 2 because we did not configure any VLAN-based ACLs. There is also a fallback flow to drop other packets, which in practice means that if any received packet already has a VLAN header then it will be dropped:

```
table=1, priority=9000,in_port=p1,vlan_tci=0x0000/0x1fff actions=push_vlan:0x8100,set_
↪field:4196->vlan_vid,goto_table:3
table=1, priority=9000,in_port=p2,vlan_tci=0x0000/0x1fff actions=push_vlan:0x8100,set_
↪field:4196->vlan_vid,goto_table:3
table=1, priority=9000,in_port=p3,vlan_tci=0x0000/0x1fff actions=push_vlan:0x8100,set_
↪field:4196->vlan_vid,goto_table:3
table=1, priority=9000,in_port=p4,vlan_tci=0x0000/0x1fff actions=push_vlan:0x8100,set_
↪field:4296->vlan_vid,goto_table:3
table=1, priority=9000,in_port=p5,vlan_tci=0x0000/0x1fff actions=push_vlan:0x8100,set_
↪field:4296->vlan_vid,goto_table:3
table=1, priority=0 actions=drop
```

Note: The syntax `set_field:4196->vlan_vid` is curious and somewhat misleading. OpenFlow 1.3 defines the `vlan_vid` field as a 13-bit field where bit 12 is set to 1 if the VLAN header is present. Thus, since 4196 is 0x1064, this action sets VLAN value 0x64, which in decimal is 100.

Table 2 isn't used because there are no VLAN-based ACLs. It just has a drop flow:

```
table=2, priority=0 actions=drop
```

Table 3 is used for MAC learning but the controller hasn't learned any MAC yet. It also drops some inappropriate packets such as those that claim to be from a broadcast source address (why not from all multicast source addresses, though?). We'll come back here later:

```
table=3, priority=9099,dl_src=ff:ff:ff:ff:ff:ff actions=drop
table=3, priority=9001,dl_src=0e:00:00:00:00:01 actions=drop
table=3, priority=0 actions=drop
table=3, priority=9000 actions=CONTROLLER:96,goto_table:7
```

Tables 4, 5, and 6 aren't used because we haven't configured any routing:

```
table=4, priority=0 actions=drop
table=5, priority=0 actions=drop
table=6, priority=0 actions=drop
```

Table 7 is used to direct packets to learned MACs but Faucet hasn't learned any MACs yet, so it just sends all the packets along to table 8:

```
table=7, priority=0 actions=drop
table=7, priority=9000 actions=goto_table:8
```

Table 8 implements flooding, broadcast, and multicast. The flows for broadcast and flood are easy to understand: if the packet came in on a given port and needs to be flooded or broadcast, output it to all the other ports in the same VLAN:

```

table=8, priority=9008,in_port=p1,d1_vlan=100,d1_dst=ff:ff:ff:ff:ff:ff actions=pop_
↪vlan,output:p2,output:p3
table=8, priority=9008,in_port=p2,d1_vlan=100,d1_dst=ff:ff:ff:ff:ff:ff actions=pop_
↪vlan,output:p1,output:p3
table=8, priority=9008,in_port=p3,d1_vlan=100,d1_dst=ff:ff:ff:ff:ff:ff actions=pop_
↪vlan,output:p1,output:p2
table=8, priority=9008,in_port=p4,d1_vlan=200,d1_dst=ff:ff:ff:ff:ff:ff actions=pop_
↪vlan,output:p5
table=8, priority=9008,in_port=p5,d1_vlan=200,d1_dst=ff:ff:ff:ff:ff:ff actions=pop_
↪vlan,output:p4
table=8, priority=9000,in_port=p1,d1_vlan=100 actions=pop_vlan,output:p2,output:p3
table=8, priority=9000,in_port=p2,d1_vlan=100 actions=pop_vlan,output:p1,output:p3
table=8, priority=9000,in_port=p3,d1_vlan=100 actions=pop_vlan,output:p1,output:p2
table=8, priority=9000,in_port=p4,d1_vlan=200 actions=pop_vlan,output:p5
table=8, priority=9000,in_port=p5,d1_vlan=200 actions=pop_vlan,output:p4

```

Note: These flows could apparently be simpler because OpenFlow says that `output : <port>` is ignored if `<port>` is the input port. That means that the first three flows above could apparently be collapsed into just:

```

table=8, priority=9008,d1_vlan=100,d1_dst=ff:ff:ff:ff:ff:ff actions=pop_vlan,
↪output:p1,output:p2,output:p3

```

There might be some reason why this won't work or isn't practical, but that isn't obvious from looking at the flow table.

There are also some flows for handling some standard forms of multicast, and a fallback drop flow:

```

table=8, priority=9006,in_port=p1,d1_vlan=100,d1_dst=33:33:00:00:00:00/
↪ff:ff:00:00:00:00 actions=pop_vlan,output:p2,output:p3
table=8, priority=9006,in_port=p2,d1_vlan=100,d1_dst=33:33:00:00:00:00/
↪ff:ff:00:00:00:00 actions=pop_vlan,output:p1,output:p3
table=8, priority=9006,in_port=p3,d1_vlan=100,d1_dst=33:33:00:00:00:00/
↪ff:ff:00:00:00:00 actions=pop_vlan,output:p1,output:p2
table=8, priority=9006,in_port=p4,d1_vlan=200,d1_dst=33:33:00:00:00:00/
↪ff:ff:00:00:00:00 actions=pop_vlan,output:p5
table=8, priority=9006,in_port=p5,d1_vlan=200,d1_dst=33:33:00:00:00:00/
↪ff:ff:00:00:00:00 actions=pop_vlan,output:p4
table=8, priority=9002,in_port=p1,d1_vlan=100,d1_dst=01:80:c2:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p2,output:p3
table=8, priority=9002,in_port=p2,d1_vlan=100,d1_dst=01:80:c2:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p1,output:p3
table=8, priority=9002,in_port=p3,d1_vlan=100,d1_dst=01:80:c2:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p1,output:p2
table=8, priority=9004,in_port=p1,d1_vlan=100,d1_dst=01:00:5e:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p2,output:p3
table=8, priority=9004,in_port=p2,d1_vlan=100,d1_dst=01:00:5e:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p1,output:p3
table=8, priority=9004,in_port=p3,d1_vlan=100,d1_dst=01:00:5e:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p1,output:p2
table=8, priority=9002,in_port=p4,d1_vlan=200,d1_dst=01:80:c2:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p5
table=8, priority=9002,in_port=p5,d1_vlan=200,d1_dst=01:80:c2:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p4
table=8, priority=9004,in_port=p4,d1_vlan=200,d1_dst=01:00:5e:00:00:00/
↪ff:ff:ff:00:00:00 actions=pop_vlan,output:p5

```

(continues on next page)

(continued from previous page)

```
table=8, priority=9004, in_port=p5, dl_vlan=200, dl_dst=01:00:5e:00:00:00/
→ff:ff:ff:00:00:00 actions=pop_vlan, output:p4
table=8, priority=0 actions=drop
```

Tracing

Let's go a level deeper. So far, everything we've done has been fairly general. We can also look at something more specific: the path that a particular packet would take through Open vSwitch. We can use OVN `ofproto/trace` command to play "what-if?" games. This command is one that we send directly to `ovs-vswitchd`, using the `ovs-appctl` utility.

Note: `ovs-appctl` is actually a very simple-minded JSON-RPC client, so you could also use some other utility that speaks JSON-RPC, or access it from a program as an API.

The `ovs-vswitchd(8)` manpage has a lot of detail on how to use `ofproto/trace`, but let's just start by building up from a simple example. You can start with a command that just specifies the datapath (e.g. `br0`), an input port, and nothing else; unspecified fields default to all-zeros. Let's look at the full output for this trivial example:

```
$ ovs-appctl ofproto/trace br0 in_port=p1
Flow: in_port=1, vlan_tci=0x0000, dl_src=00:00:00:00:00:00, dl_dst=00:00:00:00:00:00, dl_
→type=0x0000

bridge("br0")
-----
0. in_port=1, priority 9099, cookie 0x5adc15c0
   goto_table:1
1. in_port=1, vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4196->vlan_vid
   goto_table:3
3. priority 9000, cookie 0x5adc15c0
   CONTROLLER:96
   goto_table:7
7. priority 9000, cookie 0x5adc15c0
   goto_table:8
8. in_port=1, dl_vlan=100, priority 9000, cookie 0x5adc15c0
   pop_vlan
   output:2
   output:3

Final flow: unchanged
MegafLOW: recirc_id=0, eth, in_port=1, vlan_tci=0x0000, dl_src=00:00:00:00:00:00, dl_
→dst=00:00:00:00:00:00, dl_type=0x0000
Datapath actions: push_vlan(vid=100, pcp=0), userspace(pid=0, controller(reason=1,
→flags=1, recirc_id=1, rule_cookie=0x5adc15c0, controller_id=0, max_len=96)), pop_vlan, 2, 3
```

The first line of output, beginning with `Flow:`, just repeats our request in a more verbose form, including the L2 fields that were zeroed.

Each of the numbered items under `bridge("br0")` shows what would happen to our hypothetical packet in the table with the given number. For example, we see in table 1 that the packet matches a flow that push on a VLAN header, set the VLAN ID to 100, and goes on to further processing in table 3. In table 3, the packet gets sent to the controller to allow MAC learning to take place, and then table 8 floods the packet to the other ports in the same VLAN.

Summary information follows the numbered tables. The packet hasn't been changed (overall, even though a VLAN was pushed and then popped back off) since ingress, hence Final flow: unchanged. We'll look at the MegafLOW information later. The Datapath actions summarize what would actually happen to such a packet.

Triggering MAC Learning

We just saw how a packet gets sent to the controller to trigger MAC learning. Let's actually send the packet and see what happens. But before we do that, let's save a copy of the current flow tables for later comparison:

```
$ save-flows br0 > flows1
```

Now use `ofproto/trace`, as before, with a few new twists: we specify the source and destination Ethernet addresses and append the `-generate` option so that side effects like sending a packet to the controller actually happen:

```
$ ovs-appctl ofproto/trace br0 in_port=p1,dl_src=00:11:11:00:00:00,dl_
↪dst=00:22:22:00:00:00 -generate
```

The output is almost identical to that before, so it is not repeated here. But, take a look at `inst/faucet.log` now. It should now include a line at the end that says that it learned about our MAC 00:11:11:00:00:00, like this:

```
Jan 06 15:56:02 faucet.valve INFO DPID 1 (0x1) L2 learned 00:11:11:00:00:00 (L2_
↪type 0x0000, L3 src None) on Port 1 on VLAN 100 (1 hosts total)
```

Now compare the flow tables that we saved to the current ones:

```
diff-flows flows1 br0
```

The result should look like this, showing new flows for the learned MACs:

```
+table=3 priority=9098,in_port=1,dl_vlan=100,dl_src=00:11:11:00:00:00 hard_
↪timeout=3601 actions=goto_table:7
+table=7 priority=9099,dl_vlan=100,dl_dst=00:11:11:00:00:00 idle_timeout=3601_
↪actions=pop_vlan,output:1
```

To demonstrate the usefulness of the learned MAC, try tracing (with side effects) a packet arriving on p2 (or p3) and destined to the address learned on p1, like this:

```
$ ovs-appctl ofproto/trace br0 in_port=p2,dl_src=00:22:22:00:00:00,dl_
↪dst=00:11:11:00:00:00 -generate
```

The first time you run this command, you will notice that it sends the packet to the controller, to learn p2's 00:22:22:00:00:00 source address:

```
bridge("br0")
-----
0. in_port=2, priority 9099, cookie 0x5adc15c0
   goto_table:1
1. in_port=2,vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4196->vlan_vid
   goto_table:3
3. priority 9000, cookie 0x5adc15c0
   CONTROLLER:96
   goto_table:7
7. dl_vlan=100,dl_dst=00:11:11:00:00:00, priority 9099, cookie 0x5adc15c0
   pop_vlan
   output:1
```

If you check `inst/faucet.log`, you can see that p2's MAC has been learned too:

```
Jan 06 15:58:09 faucet.valve INFO DPID 1 (0x1) L2 learned 00:22:22:00:00:00 (L2_
↳type 0x0000, L3 src None) on Port 2 on VLAN 100 (2 hosts total)
```

Similarly for `diff-flows`:

```
$ diff-flows flows1 br0
+table=3 priority=9098,in_port=1,d1_vlan=100,d1_src=00:11:11:00:00:00 hard_
↳timeout=3601 actions=goto_table:7
+table=3 priority=9098,in_port=2,d1_vlan=100,d1_src=00:22:22:00:00:00 hard_
↳timeout=3604 actions=goto_table:7
+table=7 priority=9099,d1_vlan=100,d1_dst=00:11:11:00:00:00 idle_timeout=3601_
↳actions=pop_vlan,output:1
+table=7 priority=9099,d1_vlan=100,d1_dst=00:22:22:00:00:00 idle_timeout=3604_
↳actions=pop_vlan,output:2
```

Then, if you re-run either of the `ofproto/trace` commands (with or without `-generate`), you can see that the packets go back and forth without any further MAC learning, e.g.:

```
$ ovs-appctl ofproto/trace br0 in_port=p2,d1_src=00:22:22:00:00:00,d1_
↳dst=00:11:11:00:00:00 -generate
Flow: in_port=2,vlan_tci=0x0000,d1_src=00:22:22:00:00:00,d1_dst=00:11:11:00:00:00,d1_
↳type=0x0000

bridge("br0")
-----
0. in_port=2, priority 9099, cookie 0x5adc15c0
   goto_table:1
1. in_port=2,vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4196->vlan_vid
   goto_table:3
3. in_port=2,d1_vlan=100,d1_src=00:22:22:00:00:00, priority 9098, cookie 0x5adc15c0
   goto_table:7
7. d1_vlan=100,d1_dst=00:11:11:00:00:00, priority 9099, cookie 0x5adc15c0
   pop_vlan
   output:1

Final flow: unchanged
MegafLOW: recirc_id=0,eth,in_port=2,vlan_tci=0x0000/0x1fff,d1_src=00:22:22:00:00:00,
↳d1_dst=00:11:11:00:00:00,d1_type=0x0000
Datapath actions: 1
```

Performance

Open vSwitch has a concept of a “fast path” and a “slow path”; ideally all packets stay in the fast path. This distinction between slow path and fast path is the key to making sure that Open vSwitch performs as fast as possible.

Some factors can force a flow or a packet to take the slow path. As one example, all CFM, BFD, LACP, STP, and LLDP processing takes place in the slow path, in the cases where Open vSwitch processes these protocols itself instead of delegating to controller-written flows. As a second example, any flow that modifies ARP fields is processed in the slow path. These are corner cases that are unlikely to cause performance problems in practice because these protocols send packets at a relatively slow rate, and users and controller authors do not normally need to be concerned about them.

To understand what cases users and controller authors should consider, we need to talk about how Open vSwitch optimizes for performance. The Open vSwitch code is divided into two major components which, as already mentioned, are called the “slow path” and “fast path” (aka “datapath”). The slow path is embedded in the `ovs-vswitchd` userspace program. It is the part of the Open vSwitch packet processing logic that understands OpenFlow. Its job is to take a packet and run it through the OpenFlow tables to determine what should happen to it. It outputs a list of actions in a form similar to OpenFlow actions but simpler, called “ODP actions” or “datapath actions”. It then passes the ODP actions to the datapath, which applies them to the packet.

Note: Open vSwitch contains a single slow path and multiple fast paths. The difference between using Open vSwitch with the Linux kernel versus with DPDK is the datapath.

If every packet passed through the slow path and the fast path in this way, performance would be terrible. The key to getting high performance from this architecture is caching. Open vSwitch includes a multi-level cache. It works like this:

1. A packet initially arrives at the datapath. Some datapaths (such as DPDK and the in-tree version of the OVS kernel module) have a first-level cache called the “microflow cache”. The microflow cache is the key to performance for relatively long-lived, high packet rate flows. If the datapath has a microflow cache, then it consults it and, if there is a cache hit, the datapath executes the associated actions. Otherwise, it proceeds to step 2.
2. The datapath consults its second-level cache, called the “megaflow cache”. The megaflow cache is the key to performance for shorter or low packet rate flows. If there is a megaflow cache hit, the datapath executes the associated actions. Otherwise, it proceeds to step 3.
3. The datapath passes the packet to the slow path, which runs it through the OpenFlow table to yield ODP actions, a process that is often called “flow translation”. It then passes the packet back to the datapath to execute the actions and to, if possible, install a megaflow cache entry so that subsequent similar packets can be handled directly by the fast path. (We already described above most of the cases where a cache entry cannot be installed.)

The megaflow cache is the key cache to consider for performance tuning. Open vSwitch provides tools for understanding and optimizing its behavior. The `ofproto/trace` command that we have already been using is the most common tool for this use. Let’s take another look at the most recent `ofproto/trace` output:

```
$ ovs-appctl ofproto/trace br0 in_port=p2,d1_src=00:22:22:00:00:00,d1_
↳dst=00:11:11:00:00:00 -generate
Flow: in_port=2,vlan_tci=0x0000,d1_src=00:22:22:00:00:00,d1_dst=00:11:11:00:00:00,d1_
↳type=0x0000

bridge("br0")
-----
0. in_port=2, priority 9099, cookie 0x5adc15c0
   goto_table:1
1. in_port=2,vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4196->vlan_vid
   goto_table:3
3. in_port=2,d1_vlan=100,d1_src=00:22:22:00:00:00, priority 9098, cookie 0x5adc15c0
   goto_table:7
7. d1_vlan=100,d1_dst=00:11:11:00:00:00, priority 9099, cookie 0x5adc15c0
   pop_vlan
   output:1

Final flow: unchanged
Megaflow: recirc_id=0,eth,in_port=2,vlan_tci=0x0000/0x1fff,d1_src=00:22:22:00:00:00,
↳d1_dst=00:11:11:00:00:00,d1_type=0x0000
Datapath actions: 1
```

This time, it's the last line that we're interested in. This line shows the entry that Open vSwitch would insert into the megaflow cache given the particular packet with the current flow tables. The megaflow entry includes:

- `recirc_id`. This is an implementation detail that users don't normally need to understand.
- `eth`. This just indicates that the cache entry matches only Ethernet packets; Open vSwitch also supports other types of packets, such as IP packets not encapsulated in Ethernet.
- All of the fields matched by any of the flows that the packet visited:
 - `in_port` In tables 0, 1, and 3.
 - `vlan_tci` In tables 1, 3, and 7 (`vlan_tci` includes the VLAN ID and PCP fields and "`dl_vlan`" is just the VLAN ID).
 - `dl_src` In table 3
 - `dl_dst` In table 7.
- All of the fields matched by flows that had to be ruled out to ensure that the ones that actually matched were the highest priority matching rules.

The last one is important. Notice how the megaflow matches on `dl_type=0x0000`, even though none of the tables matched on `dl_type` (the Ethernet type). One reason is because of this flow in OpenFlow table 1 (which shows up in `dump-flows` output):

```
table=1, priority=9099,dl_type=0x88cc actions=drop
```

This flow has higher priority than the flow in table 1 that actually matched. This means that, to put it in the megaflow cache, `ovs-vswitchd` has to add a match on `dl_type` to ensure that the cache entry doesn't match LLDP packets (with Ethertype 0x88cc).

Note: In fact, in some cases `ovs-vswitchd` matches on fields that aren't strictly required according to this description. `dl_type` is actually one of those, so deleting the LLDP flow probably would not have any effect on the megaflow. But the principle here is sound.

So why does any of this matter? It's because, the more specific a megaflow is, that is, the more fields or bits within fields that a megaflow matches, the less valuable it is from a caching viewpoint. A very specific megaflow might match on L2 and L3 addresses and L4 port numbers. When that happens, only packets in one (half-)connection match the megaflow. If that connection has only a few packets, as many connections do, then the high cost of the slow path translation is amortized over only a few packets, so the average cost of forwarding those packets is high. On the other hand, if a megaflow only matches a relatively small number of L2 and L3 packets, then the cache entry can potentially be used by many individual connections, and the average cost is low.

For more information on how Open vSwitch constructs megaflows, including about ways that it can make megaflow entries less specific than one would infer from the discussion here, please refer to the 2015 NSDI paper, "The Design and Implementation of Open vSwitch", which focuses on this algorithm.

3.1.5 Routing

We've looked at how Faucet implements switching in OpenFlow, and how Open vSwitch implements OpenFlow through its datapath architecture. Now let's start over, adding L3 routing into the picture.

It's remarkably easy to enable routing. We just change our `vlangs` section in `inst/faucet.yaml` to specify a router IP address for each VLAN and define a router between them. The `dps` section is unchanged:

```
dps:
  switch-1:
    dp_id: 0x1
    timeout: 3600
    arp_neighbor_timeout: 3600
    interfaces:
      1:
        native_vlan: 100
      2:
        native_vlan: 100
      3:
        native_vlan: 100
      4:
        native_vlan: 200
      5:
        native_vlan: 200
vlangs:
  100:
    faucet_vips: ["10.100.0.254/24"]
  200:
    faucet_vips: ["10.200.0.254/24"]
routers:
  router-1:
    vlans: [100, 200]
```

Then we restart Faucet:

```
$ docker restart faucet
```

Note: One should be able to tell Faucet to re-read its configuration file without restarting it. I sometimes saw anomalous behavior when I did this, although I didn't characterize it well enough to make a quality bug report. I found restarting the container to be reliable.

OpenFlow Layer

Back in the OVS sandbox, let's see how the flow table has changed, with:

```
$ diff-flows flows1 br0
```

First, table 3 has new flows to direct ARP packets to table 6 (the virtual IP processing table), presumably to handle ARP for the router IPs. New flows also send IP packets destined to a particular Ethernet address to table 4 (the L3 forwarding table); we can make the educated guess that the Ethernet address is the one used by the Faucet router:

```
+table=3 priority=9131,arp,dl_vlan=100 actions=goto_table:6
+table=3 priority=9131,arp,dl_vlan=200 actions=goto_table:6
+table=3 priority=9099,ip,dl_vlan=100,dl_dst=0e:00:00:00:00:01 actions=goto_table:4
+table=3 priority=9099,ip,dl_vlan=200,dl_dst=0e:00:00:00:00:01 actions=goto_table:4
```

The new flows in table 4 appear to be verifying that the packets are indeed addressed to a network or IP address that Faucet knows how to route:

```
+table=4 priority=9131,ip,dl_vlan=100,nw_dst=10.100.0.254 actions=goto_table:6
+table=4 priority=9131,ip,dl_vlan=200,nw_dst=10.200.0.254 actions=goto_table:6
```

(continues on next page)

(continued from previous page)

```
+table=4 priority=9123,ip,d1_vlan=100,nw_dst=10.100.0.0/24 actions=goto_table:6
+table=4 priority=9123,ip,d1_vlan=200,nw_dst=10.100.0.0/24 actions=goto_table:6
+table=4 priority=9123,ip,d1_vlan=100,nw_dst=10.200.0.0/24 actions=goto_table:6
+table=4 priority=9123,ip,d1_vlan=200,nw_dst=10.200.0.0/24 actions=goto_table:6
```

Table 6 has a few different things going on. It sends ARP requests for the router IPs to the controller; presumably the controller will generate replies and send them back to the requester. It switches other ARP packets, either broadcasting them if they have a broadcast destination or attempting to unicast them otherwise. It sends all other IP packets to the controller:

```
+table=6 priority=9133,arp,arp_tpa=10.100.0.254 actions=CONTROLLER:128
+table=6 priority=9133,arp,arp_tpa=10.200.0.254 actions=CONTROLLER:128
+table=6 priority=9132,arp,d1_dst=ff:ff:ff:ff:ff:ff actions=goto_table:8
+table=6 priority=9131,arp actions=goto_table:7
+table=6 priority=9130,ip actions=CONTROLLER:128
```

Performance is clearly going to be poor if every packet that needs to be routed has to go to the controller, but it's unlikely that's the full story. In the next section, we'll take a closer look.

Tracing

As in our switching example, we can play some “what-if?” games to figure out how this works. Let's suppose that a machine with IP 10.100.0.1, on port p1, wants to send a IP packet to a machine with IP 10.200.0.1 on port p4. Assuming that these hosts have not been in communication recently, the steps to accomplish this are normally the following:

1. Host 10.100.0.1 sends an ARP request to router 10.100.0.254.
2. The router sends an ARP reply to the host.
3. Host 10.100.0.1 sends an IP packet to 10.200.0.1, via the router's Ethernet address.
4. The router broadcasts an ARP request to p4 and p5, the ports that carry the 10.200.0.<x> network.
5. Host 10.200.0.1 sends an ARP reply to the router.
6. Either the router sends the IP packet (which it buffered) to 10.200.0.1, or eventually 10.100.0.1 times out and resends it.

Let's use `ofproto/trace` to see whether Faucet and OVS follow this procedure.

Before we start, save a new snapshot of the flow tables for later comparison:

```
$ save-flows br0 > flows2
```

Step 1: Host ARP for Router

Let's simulate the ARP from 10.100.0.1 to its gateway router 10.100.0.254. This requires more detail than any of the packets we've simulated previously:

```
$ ovs-appctl ofproto/trace br0 in_port=p1,d1_src=00:01:02:03:04:05,d1_
→dst=ff:ff:ff:ff:ff:ff,d1_type=0x806,arp_spa=10.100.0.1,arp_tpa=10.100.0.254,arp_
→sha=00:01:02:03:04:05,arp_tha=ff:ff:ff:ff:ff:ff,arp_op=1 -generate
```

The important part of the output is where it shows that the packet was recognized as an ARP request destined to the router gateway and therefore sent to the controller:

```
6. arp,arp_tpa=10.100.0.254, priority 9133, cookie 0x5adc15c0
   CONTROLLER:128
```

The Faucet log shows that Faucet learned the host's MAC address, its MAC-to-IP mapping, and responded to the ARP request:

```
Jan 06 16:12:23 faucet.valve INFO      DPID 1 (0x1) Adding new route 10.100.0.1/32 via
↳10.100.0.1 (00:01:02:03:04:05) on VLAN 100
Jan 06 16:12:23 faucet.valve INFO      DPID 1 (0x1) Responded to ARP request for 10.
↳100.0.254 from 10.100.0.1 (00:01:02:03:04:05) on VLAN 100
Jan 06 16:12:23 faucet.valve INFO      DPID 1 (0x1) L2 learned 00:01:02:03:04:05 (L2
↳type 0x0806, L3 src 10.100.0.1) on Port 1 on VLAN 100 (1 hosts total)
```

We can also look at the changes to the flow tables:

```
$ diff-flows flows2 br0
+table=3 priority=9098,in_port=1,dl_vlan=100,dl_src=00:01:02:03:04:05 hard_
↳timeout=3600 actions=goto_table:7
+table=4 priority=9131,ip,dl_vlan=100,nw_dst=10.100.0.1 actions=set_field:4196->vlan_
↳vid,set_field:0e:00:00:00:00:01->eth_src,set_field:00:01:02:03:04:05->eth_dst,dec_
↳ttl,goto_table:7
+table=4 priority=9131,ip,dl_vlan=200,nw_dst=10.100.0.1 actions=set_field:4196->vlan_
↳vid,set_field:0e:00:00:00:00:01->eth_src,set_field:00:01:02:03:04:05->eth_dst,dec_
↳ttl,goto_table:7
+table=7 priority=9099,dl_vlan=100,dl_dst=00:01:02:03:04:05 idle_timeout=3600
↳actions=pop_vlan,output:1
```

The new flows include one in table 3 and one in table 7 for the learned MAC, which have the same forms we saw before. The new flows in table 4 are different. They matches packets directed to 10.100.0.1 (in two VLANs) and forward them to the host by updating the Ethernet source and destination addresses appropriately, decrementing the TTL, and skipping ahead to unicast output in table 7. This means that packets sent to 10.100.0.1 should now get to their destination.

Step 2: Router Sends ARP Reply

inst/faucet.log said that the router sent an ARP reply. How can we see it? Simulated packets just get dropped by default. One way is to configure the dummy ports to write the packets they receive to a file. Let's try that. First configure the port:

```
$ ovs-vsctl set interface p1 options:pcap=p1.pcap
```

Then re-run the “trace” command:

```
$ ovs-appctl ofproto/trace br0 in_port=p1,dl_src=00:01:02:03:04:05,dl_
↳dst=ff:ff:ff:ff:ff:ff,dl_type=0x806,arp_spa=10.100.0.1,arp_tpa=10.100.0.254,arp_
↳sha=00:01:02:03:04:05,arp_tha=ff:ff:ff:ff:ff:ff,arp_op=1 -generate
```

And dump the reply packet:

```
$ /usr/sbin/tcpdump -evvvr sandbox/p1.pcap
reading from file sandbox/p1.pcap, link-type EN10MB (Ethernet)
16:14:47.670727 0e:00:00:00:00:01 (oui Unknown) > 00:01:02:03:04:05 (oui Unknown),
↳ethertype ARP (0x0806), length 60: Ethernet (len 6), IPv4 (len 4), Reply 10.100.0.
↳254 is-at 0e:00:00:00:00:01 (oui Unknown), length 46
```

We clearly see the ARP reply, which tells us that the Faucet router's Ethernet address is 0e:00:00:00:00:01 (as we guessed before from the flow table).

Let's configure the rest of our ports to log their packets, too:

```
$ for i in 2 3 4 5; do ovs-vsctl set interface p$i options:pcap=p$i.pcap; done
```

Step 3: Host Sends IP Packet

Now that host 10.100.0.1 has the MAC address for its router, it can send an IP packet to 10.200.0.1 via the router's MAC address, like this:

```
$ ovs-appctl ofproto/trace br0 in_port=p1,dl_src=00:01:02:03:04:05,dl_
↪dst=0e:00:00:00:00:01,udp,nw_src=10.100.0.1,nw_dst=10.200.0.1,nw_ttl=64 -generate
Flow: udp,in_port=1,vlan_tci=0x0000,dl_src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,
↪nw_src=10.100.0.1,nw_dst=10.200.0.1,nw_tos=0,nw_ecn=0,nw_ttl=64,tp_src=0,tp_dst=0

bridge("br0")
-----
0. in_port=1, priority 9099, cookie 0x5adc15c0
   goto_table:1
1. in_port=1,vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4196->vlan_vid
   goto_table:3
3. ip,dl_vlan=100,dl_dst=0e:00:00:00:00:01, priority 9099, cookie 0x5adc15c0
   goto_table:4
4. ip,dl_vlan=100,nw_dst=10.200.0.0/24, priority 9123, cookie 0x5adc15c0
   goto_table:6
6. ip, priority 9130, cookie 0x5adc15c0
   CONTROLLER:128

Final flow: udp,in_port=1,dl_vlan=100,dl_vlan_pcp=0,vlan_tci1=0x0000,dl_
↪src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,nw_src=10.100.0.1,nw_dst=10.200.0.1,
↪nw_tos=0,nw_ecn=0,nw_ttl=64,tp_src=0,tp_dst=0
Megaflow: recirc_id=0,eth,ip,in_port=1,vlan_tci=0x0000/0x1fff,dl_
↪src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,nw_dst=10.200.0.0/25,nw_frag=no
Datapath actions: push_vlan(vid=100,pcp=0),userspace(pid=0,controller(reason=1,
↪flags=0,recirc_id=6,rule_cookie=0x5adc15c0,controller_id=0,max_len=128))
```

Observe that the packet gets recognized as destined to the router, in table 3, and then as properly destined to the 10.200.0.0/24 network, in table 4. In table 6, however, it gets sent to the controller. Presumably, this is because Faucet has not yet resolved an Ethernet address for the destination host 10.200.0.1. It probably sent out an ARP request. Let's take a look in the next step.

Step 4: Router Broadcasts ARP Request

The router needs to know the Ethernet address of 10.200.0.1. It knows that, if this machine exists, it's on port p4 or p5, since we configured those ports as VLAN 200.

Let's make sure:

```
$ /usr/sbin/tcpdump -evvvv sandbox/p4.pcap
reading from file sandbox/p4.pcap, link-type EN10MB (Ethernet)
16:17:43.174006 0e:00:00:00:00:01 (oui Unknown) > Broadcast, ethertype ARP (0x0806),
↪length 60: Ethernet (len 6), IPv4 (len 4), Request who-has 10.200.0.1 tell 10.200.0.1
↪254, length 46
```

(continues on next page)

(continued from previous page)

and:

```
$ /usr/sbin/tcpdump -evvvr sandbox/p5.pcap
reading from file sandbox/p5.pcap, link-type EN10MB (Ethernet)
16:17:43.174268 0e:00:00:00:00:01 (oui Unknown) > Broadcast, ethertype ARP (0x0806),
↳ length 60: Ethernet (len 6), IPv4 (len 4), Request who-has 10.200.0.1 tell 10.200.0.
↳ 254, length 46
```

For good measure, let's make sure that it wasn't sent to p3:

```
$ /usr/sbin/tcpdump -evvvr sandbox/p3.pcap
reading from file sandbox/p3.pcap, link-type EN10MB (Ethernet)
```

Step 5: Host 2 Sends ARP Reply

The Faucet controller sent an ARP request, so we can send an ARP reply:

```
$ ovs-appctl ofproto/trace br0 in_port=p4,dl_src=00:10:20:30:40:50,dl_
↳ dst=0e:00:00:00:00:01,dl_type=0x806,arp_spa=10.200.0.1,arp_tpa=10.200.0.254,arp_
↳ sha=00:10:20:30:40:50,arp_tha=0e:00:00:00:00:01,arp_op=2 -generate
Flow: arp,in_port=4,vlan_tci=0x0000,dl_src=00:10:20:30:40:50,dl_dst=0e:00:00:00:00:01,
↳ arp_spa=10.200.0.1,arp_tpa=10.200.0.254,arp_op=2,arp_sha=00:10:20:30:40:50,arp_
↳ tha=0e:00:00:00:00:01

bridge("br0")
-----
0. in_port=4, priority 9099, cookie 0x5adc15c0
   goto_table:1
1. in_port=4,vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4296->vlan_vid
   goto_table:3
3. arp,dl_vlan=200, priority 9131, cookie 0x5adc15c0
   goto_table:6
6. arp,arp_tpa=10.200.0.254, priority 9133, cookie 0x5adc15c0
   CONTROLLER:128

Final flow: arp,in_port=4,dl_vlan=200,dl_vlan_pcp=0,vlan_tci1=0x0000,dl_
↳ src=00:10:20:30:40:50,dl_dst=0e:00:00:00:00:01,arp_spa=10.200.0.1,arp_tpa=10.200.0.
↳ 254,arp_op=2,arp_sha=00:10:20:30:40:50,arp_tha=0e:00:00:00:00:01
MegafLOW: recirc_id=0,eth,arp,in_port=4,vlan_tci=0x0000/0x1fff,dl_
↳ dst=0e:00:00:00:00:01,arp_tpa=10.200.0.254
Datapath actions: push_vlan(vid=200,pcp=0),userspace(pid=0,controller(reason=1,
↳ flags=0,recirc_id=7,rule_cookie=0x5adc15c0,controller_id=0,max_len=128))
```

It shows up in inst/faucet.log:

```
Jan 06 03:20:11 faucet.valve INFO      DPID 1 (0x1) Adding new route 10.200.0.1/32 via
↳ 10.200.0.1 (00:10:20:30:40:50) on VLAN 200
Jan 06 03:20:11 faucet.valve INFO      DPID 1 (0x1) ARP response 10.200.0.1
↳ (00:10:20:30:40:50) on VLAN 200
Jan 06 03:20:11 faucet.valve INFO      DPID 1 (0x1) L2 learned 00:10:20:30:40:50 (L2
↳ type 0x0806, L3 src 10.200.0.1) on Port 4 on VLAN 200 (1 hosts total)
```

and in the OVS flow tables:

```
$ diff-flows flows2 br0
+table=3 priority=9098,in_port=4,dl_vlan=200,dl_src=00:10:20:30:40:50 hard_
↳timeout=3601 actions=goto_table:7
...
+table=4 priority=9131,ip,dl_vlan=200,nw_dst=10.200.0.1 actions=set_field:4296->vlan_
↳vid,set_field:0e:00:00:00:00:01->eth_src,set_field:00:10:20:30:40:50->eth_dst,dec_
↳ttl,goto_table:7
+table=4 priority=9131,ip,dl_vlan=100,nw_dst=10.200.0.1 actions=set_field:4296->vlan_
↳vid,set_field:0e:00:00:00:00:01->eth_src,set_field:00:10:20:30:40:50->eth_dst,dec_
↳ttl,goto_table:7
...
+table=4 priority=9123,ip,dl_vlan=100,nw_dst=10.200.0.0/24 actions=goto_table:6
+table=7 priority=9099,dl_vlan=200,dl_dst=00:10:20:30:40:50 idle_timeout=3601,
↳actions=pop_vlan,output:4
```

Step 6: IP Packet Delivery

Now both the host and the router have everything they need to deliver the packet. There are two ways it might happen. If Faucet's router is smart enough to buffer the packet that trigger ARP resolution, then it might have delivered it already. If so, then it should show up in `p4.pcap`. Let's take a look:

```
$ /usr/sbin/tcpdump -evvvr sandbox/p4.pcap ip
reading from file sandbox/p4.pcap, link-type EN10MB (Ethernet)
```

Nope. That leaves the other possibility, which is that Faucet waits for the original sending host to re-send the packet. We can do that by re-running the trace:

```
$ ovs-appctl ofproto/trace br0 in_port=1,dl_src=00:01:02:03:04:05,dl_
↳dst=0e:00:00:00:00:01,udp,nw_src=10.100.0.1,nw_dst=10.200.0.1,nw_ttl=64 -generate
Flow: udp,in_port=1,vlan_tci=0x0000,dl_src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,
↳nw_src=10.100.0.1,nw_dst=10.200.0.1,nw_tos=0,nw_ecn=0,nw_ttl=64,tp_src=0,tp_dst=0

bridge("br0")
-----
0. in_port=1, priority 9099, cookie 0x5adc15c0
   goto_table:1
1. in_port=1,vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4196->vlan_vid
   goto_table:3
3. ip,dl_vlan=100,dl_dst=0e:00:00:00:00:01, priority 9099, cookie 0x5adc15c0
   goto_table:4
4. ip,dl_vlan=100,nw_dst=10.200.0.1, priority 9131, cookie 0x5adc15c0
   set_field:4296->vlan_vid
   set_field:0e:00:00:00:00:01->eth_src
   set_field:00:10:20:30:40:50->eth_dst
   dec_ttl
   goto_table:7
7. dl_vlan=200,dl_dst=00:10:20:30:40:50, priority 9099, cookie 0x5adc15c0
   pop_vlan
   output:4

Final flow: udp,in_port=1,vlan_tci=0x0000,dl_src=0e:00:00:00:00:01,dl_
↳dst=00:10:20:30:40:50,nw_src=10.100.0.1,nw_dst=10.200.0.1,nw_tos=0,nw_ecn=0,nw_
↳ttl=63,tp_src=0,tp_dst=0
```

(continues on next page)

(continued from previous page)

```

MegafLOW: recirc_id=0,eth,ip,in_port=1,vlan_tci=0x0000/0x1fff,dl_
↪src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,nw_dst=10.200.0.1,nw_ttl=64,nw_
↪frag=no
Datapath actions: set (eth(src=0e:00:00:00:00:01,dst=00:10:20:30:40:50)),
↪set (ipv4(dst=10.200.0.1,ttl=63)),4

```

Finally, we have working IP packet forwarding!

Performance

Take another look at the megafLOW line above:

```

MegafLOW: recirc_id=0,eth,ip,in_port=1,vlan_tci=0x0000/0x1fff,dl_
↪src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,nw_dst=10.200.0.1,nw_ttl=64,nw_
↪frag=no

```

This means that (almost) any packet between these Ethernet source and destination hosts, destined to the given IP host, will be handled by this single megafLOW cache entry. So regardless of the number of UDP packets or TCP connections that these hosts exchange, Open vSwitch packet processing won't need to fall back to the slow path. It is quite efficient.

Note: The exceptions are packets with a TTL other than 64, and fragmented packets. Most hosts use a constant TTL for outgoing packets, and fragments are rare. If either of those did change, then that would simply result in a new megafLOW cache entry.

The datapath actions might also be worth a look:

```

Datapath actions: set (eth(src=0e:00:00:00:00:01,dst=00:10:20:30:40:50)),
↪set (ipv4(dst=10.200.0.1,ttl=63)),4

```

This just means that, to process these packets, the datapath changes the Ethernet source and destination addresses and the IP TTL, and then transmits the packet to port p4 (also numbered 4). Notice in particular that, despite the OpenFlow actions that pushed, modified, and popped back off a VLAN, there is nothing in the datapath actions about VLANs. This is because the OVS flow translation code “optimizes out” redundant or unneeded actions, which saves time when the cache entry is executed later.

Note: It's not clear why the actions also re-set the IP destination address to its original value. Perhaps this is a minor performance bug.

3.1.6 ACLs

Let's try out some ACLs, since they do a good job illustrating some of the ways that OVS tries to optimize megafLOWS. Update `inst/faucet.yaml` to the following:

```

dps:
  switch-1:
    dp_id: 0x1
    timeout: 3600
    arp_neighbor_timeout: 3600
    interfaces:
      1:

```

(continues on next page)

(continued from previous page)

```

        native_vlan: 100
        acl_in: 1
    2:
        native_vlan: 100
    3:
        native_vlan: 100
    4:
        native_vlan: 200
    5:
        native_vlan: 200
vlangs:
    100:
        faucet_vips: ["10.100.0.254/24"]
    200:
        faucet_vips: ["10.200.0.254/24"]
routers:
    router-1:
        vlans: [100, 200]
acls:
    1:
        - rule:
            dl_type: 0x800
            nw_proto: 6
            tcp_dst: 8080
            actions:
                allow: 0
        - rule:
            actions:
                allow: 1

```

Then restart Faucet:

```
$ docker restart faucet
```

On port 1, this new configuration blocks all traffic to TCP port 8080 and allows all other traffic. The resulting change in the flow table shows this clearly too:

```

$ diff-flows flows2 br0
-priority=9099,in_port=1 actions=goto_table:1
+priority=9098,in_port=1 actions=goto_table:1
+priority=9099,tcp,in_port=1,tp_dst=8080 actions=drop

```

The most interesting question here is performance. If you recall the earlier discussion, when a packet through the flow table encounters a match on a given field, the resulting megaflow has to match on that field, even if the flow didn't actually match. This is expensive.

In particular, here you can see that any TCP packet is going to encounter the ACL flow, even if it is directed to a port other than 8080. If that means that every megaflow for a TCP packet is going to have to match on the TCP destination, that's going to be bad for caching performance because there will be a need for a separate megaflow for every TCP destination port that actually appears in traffic, which means a lot more megaflows than otherwise. (Really, in practice, if such a simple ACL blew up performance, OVS wouldn't be a very good switch!)

Let's see what happens, by sending a packet to port 80 (instead of 8080):

```

$ ovs-appctl ofproto/trace br0 in_port=p1,d1_src=00:01:02:03:04:05,d1_
→dst=0e:00:00:00:00:01,tcp,nw_src=10.100.0.1,nw_dst=10.200.0.1,nw_ttl=64,tp_dst=80 -
→generate

```

(continues on next page)

(continued from previous page)

```

Flow: tcp,in_port=1,vlan_tci=0x0000,dl_src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,
↪nw_src=10.100.0.1,nw_dst=10.200.0.1,nw_tos=0,nw_ecn=0,nw_ttl=64,tp_src=0,tp_dst=80,
↪tcp_flags=0

bridge("br0")
-----
0. in_port=1, priority 9098, cookie 0x5adc15c0
   goto_table:1
1. in_port=1,vlan_tci=0x0000/0x1fff, priority 9000, cookie 0x5adc15c0
   push_vlan:0x8100
   set_field:4196->vlan_vid
   goto_table:3
3. ip,dl_vlan=100,dl_dst=0e:00:00:00:00:01, priority 9099, cookie 0x5adc15c0
   goto_table:4
4. ip,dl_vlan=100,nw_dst=10.200.0.0/24, priority 9123, cookie 0x5adc15c0
   goto_table:6
6. ip, priority 9130, cookie 0x5adc15c0
   CONTROLLER:128

Final flow: tcp,in_port=1,dl_vlan=100,dl_vlan_pcp=0,vlan_tci1=0x0000,dl_
↪src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,nw_src=10.100.0.1,nw_dst=10.200.0.1,
↪nw_tos=0,nw_ecn=0,nw_ttl=64,tp_src=0,tp_dst=80,tcp_flags=0
Megaflow: recirc_id=0,eth,tcp,in_port=1,vlan_tci=0x0000/0x1fff,dl_
↪src=00:01:02:03:04:05,dl_dst=0e:00:00:00:00:01,nw_dst=10.200.0.1,nw_frag=no,tp_
↪dst=0x0/0xf000
Datapath actions: push_vlan(vid=100,pcp=0)

```

Take a look at the Megaflow line and in particular the match on `tp_dst`, which says `tp_dst=0x0/0xf000`. What this means is that the megaflow matches on only the top 4 bits of the TCP destination port. That works because:

```

80 (base 10) == 0000,0000,0101,0000 (base 2)
8080 (base 10) == 0001,1111,1001,0000 (base 2)

```

and so by matching on only the top 4 bits, rather than all 16, the OVS fast path can distinguish port 80 from port 8080. This allows this megaflow to match one-sixteenth of the TCP destination port address space, rather than just 1/65536th of it.

Note: The algorithm OVS uses for this purpose isn't perfect. In this case, a single-bit match would work (e.g. `tp_dst=0x0/0x1000`), and would be superior since it would only match half the port address space instead of one-sixteenth.

For details of this algorithm, please refer to `lib/classifier.c` in the Open vSwitch source tree, or our 2015 NSDI paper "The Design and Implementation of Open vSwitch".

3.1.7 Finishing Up

When you're done, you probably want to exit the sandbox session, with `Control+D` or `exit`, and stop the Faucet controller with `docker stop faucet; docker rm faucet`.

3.1.8 Further Directions

We've looked a fair bit at how Faucet interacts with Open vSwitch. If you still have some interest, you might want to explore some of these directions:

- Adding more than one switch. Faucet can control multiple switches but we've only been simulating one of them. It's easy enough to make a single OVS instance act as multiple switches (just `ovs-vsctl add-br` another bridge), or you could use genuinely separate OVS instances.
- Additional features. Faucet has more features than we've demonstrated, such as IPv6 routing and port mirroring. These should also interact gracefully with Open vSwitch.
- Real performance testing. We've looked at how flows and traces **should** demonstrate good performance, but of course there's no proof until it actually works in practice. We've also only tested with trivial configurations. Open vSwitch can scale to millions of OpenFlow flows, but the scaling in practice depends on the particular flow tables and traffic patterns, so it's valuable to test with large configurations, either in the way we've done it or with real traffic.

3.2 OVS IPsec Tutorial

This document provides a step-by-step guide for running IPsec tunnel in Open vSwitch. A more detailed description on OVS IPsec tunnel and its configuration modes can be found in [Encrypt Open vSwitch Tunnels with IPsec](#).

3.2.1 Requirements

OVS IPsec tunnel requires Linux kernel (\geq v3.10.0) and OVS out-of-tree kernel module. The compatible IKE daemons are LibreSwan (\geq v3.23) and StrongSwan (\geq v5.3.5).

3.2.2 Installing OVS and IPsec Packages

OVS IPsec has .deb and .rpm packages. You should use the right package based on your Linux distribution. This tutorial uses Ubuntu 16.04 and Fedora 27 as examples.

Ubuntu

1. Follow [Debian Packaging for Open vSwitch](#) to build debian packages.

Note: If you have already installed OVS, then you only need to install `openvswitch-pki_*.deb` and `openvswitch-ipsec_*.deb` in the following step. If your kernel version is below v4.13.0, update your kernel to v4.13.0 or above.

2. Install the related packages:

```
$ apt-get install dkms strongswan
$ dpkg -i libopenvswitch_*.deb openvswitch-common_*.deb \
    openvswitch-switch_*.deb openvswitch-datapath-dkms_*.deb \
    python-openvswitch_*.deb openvswitch-pki_*.deb \
    openvswitch-ipsec_*.deb
```

If the installation is successful, you should be able to see the `ovs-monitor-ipsec` daemon is running in your system.

Fedora

1. Follow *Fedora, RHEL 7.x Packaging for Open vSwitch* to build RPM packages.
2. Install the related packages:

```
$ dnf install python2-openvswitch libreswan \
    "kernel-devel-uname-r == $(uname -r) "
$ rpm -i openvswitch-*.rpm openvswitch-kmod-*.rpm \
    openvswitch-openvswitch-ipsec-*.rpm
```

3. Install firewall rules to allow ESP and IKE traffic:

```
$ iptables -A IN_FedoraServer_allow -p esp -j ACCEPT
$ iptables -A IN_FedoraServer_allow -p udp --dport 500 -j ACCEPT
```

4. Run the openvswitch-ipsec service:

```
$ systemctl start openvswitch-ipsec.service
```

Note: The SELinux policies might prevent openvswitch-ipsec.service to access certain resources. You can configure SELinux to remove such restrictions.

3.2.3 Configuring IPsec tunnel

Suppose you want to build IPsec tunnel between two hosts. Assume *host_1*'s external IP is 1.1.1.1, and *host_2*'s external IP is 2.2.2.2. Make sure *host_1* and *host_2* can ping each other via these external IPs.

0. Set up some variables to make life easier. On both hosts, set *ip_1* and *ip_2* variables, e.g.:

```
$ ip_1=1.1.1.1
$ ip_2=2.2.2.2
```

1. Set up OVS bridges in both hosts.

In *host_1*:

```
$ ovs-vsctl add-br br-ipsec
$ ip addr add 192.0.0.1/24 dev br-ipsec
$ ip link set br-ipsec up
```

In *host_2*:

```
$ ovs-vsctl add-br br-ipsec
$ ip addr add 192.0.0.2/24 dev br-ipsec
$ ip link set br-ipsec up
```

2. Set up IPsec tunnel.

There are three authentication methods. You can choose one to set up your IPsec tunnel.

- (a) Using pre-shared key:

In *host_1*:

```
$ ovs-vsctl add-port br-ipsec tun -- \
    set interface tun type=gre \
        options:remote_ip=$ip_2 \
        options:psk=swordfish
```

In *host_2*:

```
$ ovs-vsctl add-port br-ipsec tun -- \
    set interface tun type=gre \
        options:remote_ip=$ip_1 \
        options:psk=swordfish
```

Note: Pre-shared key (PSK) based authentication is easy to set up but less secure compared with other authentication methods. You should use it cautiously in production systems.

(b) Using self-signed certificate:

Generate self-signed certificate in both *host_1* and *host_2*. Then copy the certificate of *host_1* to *host_2* and the certificate of *host_2* to *host_1*.

In *host_1*:

```
$ ovs-pki req -u host_1
$ ovs-pki self-sign host_1
$ scp host_1-cert.pem $ip_2:/etc/keys/host_1-cert.pem
```

In *host_2*:

```
$ ovs-pki req -u host_2
$ ovs-pki self-sign host_2
$ scp host_2-cert.pem $ip_1:/etc/keys/host_2-cert.pem
```

Note: If you use StrongSwan as IKE daemon, please move the host certificates to */etc/ipsec.d/certs/* and private key to */etc/ipsec.d/private/* so that StrongSwan has permission to access those files.

Configure IPsec tunnel to use self-signed certificates.

In *host_1*:

```
$ ovs-vsctl set Open_vSwitch . \
    other_config:certificate=/etc/keys/host_1-cert.pem \
    other_config:private_key=/etc/keys/host_1-privkey.pem
$ ovs-vsctl add-port br-ipsec tun -- \
    set interface tun type=gre \
        options:remote_ip=$ip_2 \
        options:remote_cert=/etc/keys/host_2-cert.pem
```

In *host_2*:

```
$ ovs-vsctl set Open_vSwitch . \
    other_config:certificate=/etc/keys/host_2-cert.pem \
    other_config:private_key=/etc/keys/host_2-privkey.pem
$ ovs-vsctl add-port br-ipsec tun -- \
    set interface tun type=gre \
```

(continues on next page)

(continued from previous page)

```
options:remote_ip=$ip_1 \
options:remote_cert=/etc/keys/host_1-cert.pem
```

Note: The confidentiality of the private key is very critical. Don't copy it to places where it might be compromised. (The certificate need not be kept confidential.)

(c) Using CA-signed certificate:

First you need to establish a public key infrastructure (PKI). Suppose you choose *host_1* to host PKI.

In *host_1*:

```
$ ovs-pki init
```

Generate certificate requests and copy the certificate request of *host_2* to *host_1*.

In *host_1*:

```
$ ovs-pki req -u host_1
```

In *host_2*:

```
$ ovs-pki req -u host_2
$ scp host_2-req.pem $ip_1:/etc/keys/host_2-req.pem
```

Sign the certificate requests with the CA key. Copy *host_2*'s signed certificate and the CA certificate to *host_2*.

In *host_1*:

```
$ ovs-pki sign host_1 switch
$ ovs-pki sign host_2 switch
$ scp host_2-cert.pem $ip_2:/etc/keys/host_2-cert.pem
$ scp /var/lib/openvswitch/pki/switchca/cacert.pem \
    $ip_2:/etc/keys/cacert.pem
```

Note: If you use StrongSwan as IKE daemon, please move the host certificates to */etc/ipsec.d/certs/*, CA certificate to */etc/ipsec.d/cacerts/*, and private key to */etc/ipsec.d/private/* so that StrongSwan has permission to access those files.

Configure IPsec tunnel to use CA-signed certificate.

In *host_1*:

```
$ ovs-vsctl set Open_vSwitch . \
    other_config:certificate=/etc/keys/host_1-cert.pem \
    other_config:private_key=/etc/keys/host_1-privkey.pem \
    other_config:ca_cert=/etc/keys/cacert.pem
$ ovs-vsctl add-port br-ipsec tun -- \
    set interface tun type=gre \
    options:remote_ip=$ip_2 \
    options:remote_name=host_2
```

In *host_2*:

```
$ ovs-vsctl set Open_vSwitch . \
    other_config:certificate=/etc/keys/host_2-cert.pem \
    other_config:private_key=/etc/keys/host_2-privkey.pem \
    other_config:ca_cert=/etc/keys/cacert.pem
$ ovs-vsctl add-port br-ipsec tun -- \
    set interface tun type=gre \
    options:remote_ip=$ip_1 \
    options:remote_name=host_1
```

Note: remote_name is the common name (CN) of the signed-certificate. It must match the name given as the argument to the `ovs-pki sign` command. It ensures that only certificate with the expected CN can be authenticated; otherwise, any certificate signed by the CA would be accepted.

3. Test IPsec tunnel.

Now you should have an IPsec GRE tunnel running between two hosts. To verify it, in *host_1*:

```
$ ping 192.0.0.2 &
$ tcpdump -ni any net $ip_2
```

You should be able to see that ESP packets are being sent from *host_1* to *host_2*.

3.2.4 Troubleshooting

The `ovs-monitor-ipsec` daemon manages and monitors the IPsec tunnel state. Use the following `ovs-appctl` command to view `ovs-monitor-ipsec` internal representation of tunnel configuration:

```
$ ovs-appctl -t ovs-monitor-ipsec tunnels/show
```

If there is misconfiguration, then `ovs-appctl` should indicate why. For example:

```
Interface name: gre0 v5 (CONFIGURED) <--- Should be set to CONFIGURED.
                                                Otherwise, error message will
                                                be provided

Tunnel Type:      gre
Remote IP:        2.2.2.2
SKB mark:         None
Local cert:       None
Local name:       None
Local key:        None
Remote cert:      None
Remote name:      None
CA cert:          None
PSK:              swordfish
Ofport:           1      <--- Whether ovs-vswitchd has assigned Ofport
                        number to this Tunnel Port
CFM state:        Up      <--- Whether CFM declared this tunnel healthy
Kernel policies installed:
...               <--- IPsec policies for this OVS tunnel in
                        Linux Kernel installed by strongSwan
Kernel security associations installed:
...               <--- IPsec security associations for this OVS
                        tunnel in Linux Kernel installed by
                        strongswan
```

(continues on next page)

(continued from previous page)

```
IPsec connections that are active:
...                <--- IPsec "connections" for this OVS
                    tunnel
```

If you don't see any active connections, try to run the following command to refresh the `ovs-monitor-ipsec` daemon:

```
$ ovs-appctl -t ovs-monitor-ipsec refresh
```

You can also check the logs of the `ovs-monitor-ipsec` daemon and the IKE daemon to locate issues. `ovs-monitor-ipsec` outputs log messages to `/var/log/openvswitch/ovs-monitor-ipsec.log`.

3.2.5 Bug Reporting

If you think you may have found a bug with security implications, like

1. IPsec protected tunnel accepted packets that came unencrypted; OR
2. IPsec protected tunnel allowed packets to leave unencrypted;

Then report such bugs according to *Open vSwitch's Security Process*.

If bug does not have security implications, then report it according to instructions in *Reporting Bugs in Open vSwitch*.

If you have suggestions to improve this tutorial, please send a email to ovs-discuss@openvswitch.org.

3.3 Open vSwitch Advanced Features

Many tutorials cover the basics of OpenFlow. This is not such a tutorial. Rather, a knowledge of the basics of OpenFlow is a prerequisite. If you do not already understand how an OpenFlow flow table works, please go read a basic tutorial and then continue reading here afterward.

It is also important to understand the basics of Open vSwitch before you begin. If you have never used `ovs-vsctl` or `ovs-ofctl` before, you should learn a little about them before proceeding.

Most of the features covered in this tutorial are Open vSwitch extensions to OpenFlow. Also, most of the features in this tutorial are specific to the software Open vSwitch implementation. If you are using an Open vSwitch port to an ASIC-based hardware switch, this tutorial will not help you.

This tutorial does not cover every aspect of the features that it mentions. You can find the details elsewhere in the Open vSwitch documentation, especially `ovs-ofctl(8)` and the comments in the `include/openflow/nicira-ext.h` and `include/openvswitch/meta-flow.h` header files.

3.3.1 Getting Started

This is a hands-on tutorial. To get the most out of it, you will need Open vSwitch binaries. You do not, on the other hand, need any physical networking hardware or even supervisor privilege on your system. Instead, we will use a script called `ovs-sandbox`, which accompanies the tutorial, that constructs a software simulated network environment based on Open vSwitch.

You can use `ovs-sandbox` three ways:

- If you have already installed Open vSwitch on your system, then you should be able to just run `ovs-sandbox` from this directory without any options.

- If you have not installed Open vSwitch (and you do not want to install it), then you can build Open vSwitch according to the instructions in *Open vSwitch on Linux, FreeBSD and NetBSD*, without installing it. Then run `./ovs-sandbox -b DIRECTORY` from this directory, substituting the Open vSwitch build directory for `DIRECTORY`.
- As a slight variant on the latter, you can run `make sandbox` from an Open vSwitch build directory.

When you run `ovs-sandbox`, it does the following:

1. **CAUTION:** Deletes any subdirectory of the current directory named “sandbox” and any files in that directory.
2. Creates a new directory “sandbox” in the current directory.
3. Sets up special environment variables that ensure that Open vSwitch programs will look inside the “sandbox” directory instead of in the Open vSwitch installation directory.
4. If you are using a built but not installed Open vSwitch, installs the Open vSwitch manpages in a subdirectory of “sandbox” and adjusts the `MANPATH` environment variable to point to this directory. This means that you can use, for example, `man ovs-vsctl` to see a manpage for the `ovs-vsctl` program that you built.
5. Creates an empty Open vSwitch configuration database under “sandbox”.
6. Starts `ovsdb-server` running under “sandbox”.
7. Starts `ovs-vswitchd` running under “sandbox”, passing special options that enable a special “dummy” mode for testing.
8. Starts a nested interactive shell inside “sandbox”.

At this point, you can run all the usual Open vSwitch utilities from the nested shell environment. You can, for example, use `ovs-vsctl` to create a bridge:

```
$ ovs-vsctl add-br br0
```

From Open vSwitch’s perspective, the bridge that you create this way is as real as any other. You can, for example, connect it to an OpenFlow controller or use `ovs-ofctl` to examine and modify it and its OpenFlow flow table. On the other hand, the bridge is not visible to the operating system’s network stack, so `ip` cannot see it or affect it, which means that utilities like `ping` and `tcpdump` will not work either. (That has its good side, too: you can’t screw up your computer’s network stack by manipulating a sandboxed OVS.)

When you’re done using OVS from the sandbox, exit the nested shell (by entering the “exit” shell command or pressing Control+D). This will kill the daemons that `ovs-sandbox` started, but it leaves the “sandbox” directory and its contents in place.

The sandbox directory contains log files for the Open vSwitch daemons. You can examine them while you’re running in the sandboxed environment or after you exit.

3.3.2 Using GDB

GDB support is not required to go through the tutorial. It is added in case user wants to explore the internals of OVS programs.

GDB can already be used to debug any running process, with the usual `gdb <program> <process-id>` command.

`ovs-sandbox` also has a `-g` option for launching `ovs-vswitchd` under GDB. This option can be handy for setting break points before `ovs-vswitchd` runs, or for catching early segfaults. Similarly, a `-d` option can be used to run `ovsdb-server` under GDB. Both options can be specified at the same time.

In addition, a `-e` option also launches `ovs-vswitchd` under GDB. However, instead of displaying a `gdb>` prompt and waiting for user input, `ovs-vswitchd` will start to execute immediately. `-r` option is the corresponding option for running `ovsdb-server` under `gdb` with immediate execution.

To avoid GDB mangling with the sandbox sub shell terminal, `ovs-sandbox` starts a new xterm to run each GDB session. For systems that do not support X windows, GDB support is effectively disabled.

When launching sandbox through the build tree's make file, the `-g` option can be passed via the `SANDBOXFLAGS` environment variable. `make sandbox SANDBOXFLAGS=-g` will start the sandbox with `ovs-vswitchd` running under GDB in its own xterm if X is available.

In addition, a set of GDB macros are available in `utilities/gdb/ovs_gdb.py`. Which are able to dump various internal data structures. See the header of the file itself for some more details and an example.

3.3.3 Motivation

The goal of this tutorial is to demonstrate the power of Open vSwitch flow tables. The tutorial works through the implementation of a MAC-learning switch with VLAN trunk and access ports. Outside of the Open vSwitch features that we will discuss, OpenFlow provides at least two ways to implement such a switch:

1. An OpenFlow controller to implement MAC learning in a “reactive” fashion. Whenever a new MAC appears on the switch, or a MAC moves from one switch port to another, the controller adjusts the OpenFlow flow table to match.
2. The “normal” action. OpenFlow defines this action to submit a packet to “the traditional non-OpenFlow pipeline of the switch”. That is, if a flow uses this action, then the packets in the flow go through the switch in the same way that they would if OpenFlow was not configured on the switch.

Each of these approaches has unfortunate pitfalls. In the first approach, using an OpenFlow controller to implement MAC learning, has a significant cost in terms of network bandwidth and latency. It also makes the controller more difficult to scale to large numbers of switches, which is especially important in environments with thousands of hypervisors (each of which contains a virtual OpenFlow switch). MAC learning at an OpenFlow controller also behaves poorly if the OpenFlow controller fails, slows down, or becomes unavailable due to network problems.

The second approach, using the “normal” action, has different problems. First, little about the “normal” action is standardized, so it behaves differently on switches from different vendors, and the available features and how those features are configured (usually not through OpenFlow) varies widely. Second, “normal” does not work well with other OpenFlow actions. It is “all-or-nothing”, with little potential to adjust its behavior slightly or to compose it with other features.

3.3.4 Scenario

We will construct Open vSwitch flow tables for a VLAN-capable, MAC-learning switch that has four ports:

- p1** a trunk port that carries all VLANs, on OpenFlow port 1.
- p2** an access port for VLAN 20, on OpenFlow port 2.
- p3, p4** both access ports for VLAN 30, on OpenFlow ports 3 and 4, respectively.

Note: The ports' names are not significant. You could call them `eth1` through `eth4`, or any other names you like.

Note: An OpenFlow switch always has a “local” port as well. This scenario won't use the local port.

Our switch design will consist of five main flow tables, each of which implements one stage in the switch pipeline:

Table 0 Admission control.

Table 1 VLAN input processing.

Table 2 Learn source MAC and VLAN for ingress port.

Table 3 Look up learned port for destination MAC and VLAN.

Table 4 Output processing.

The section below describes how to set up the scenario, followed by a section for each OpenFlow table.

You can cut and paste the `ovs-vsctl` and `ovs-ofctl` commands in each of the sections below into your `ovs-sandbox` shell. They are also available as shell scripts in this directory, named `t-setup`, `t-stage0`, `t-stage1`, ..., `t-stage4`. The `ovs-appctl` test commands are intended for cutting and pasting and are not supplied separately.

3.3.5 Setup

To get started, start `ovs-sandbox`. Inside the interactive shell that it starts, run this command:

```
$ ovs-vsctl add-br br0 -- set Bridge br0 fail-mode=secure
```

This command creates a new bridge “br0” and puts “br0” into so-called “fail-secure” mode. For our purpose, this just means that the OpenFlow flow table starts out empty.

Note: If we did not do this, then the flow table would start out with a single flow that executes the “normal” action. We could use that feature to yield a switch that behaves the same as the switch we are currently building, but with the caveats described under “Motivation” above.)

The new bridge has only one port on it so far, the “local port” br0. We need to add p1, p2, p3, and p4. A shell for loop is one way to do it:

```
for i in 1 2 3 4; do
    ovs-vsctl add-port br0 p$i -- set Interface p$i ofport_request=$i
    ovs-ofctl mod-port br0 p$i up
done
```

In addition to adding a port, the `ovs-vsctl` command above sets its `ofport_request` column to ensure that port p1 is assigned OpenFlow port 1, p2 is assigned OpenFlow port 2, and so on.

Note: We could omit setting the `ofport_request` and let Open vSwitch choose port numbers for us, but it’s convenient for the purposes of this tutorial because we can talk about OpenFlow port 1 and know that it corresponds to p1.

The `ovs-ofctl` command above brings up the simulated interfaces, which are down initially, using an OpenFlow request. The effect is similar to `ip link up`, but the sandbox’s interfaces are not visible to the operating system and therefore `ip` would not affect them.

We have not configured anything related to VLANs or MAC learning. That’s because we’re going to implement those features in the flow table.

To see what we’ve done so far to set up the scenario, you can run a command like `ovs-vsctl show` or `ovs-ofctl show br0`.

3.3.6 Implementing Table 0: Admission control

Table 0 is where packets enter the switch. We use this stage to discard packets that for one reason or another are invalid. For example, packets with a multicast source address are not valid, so we can add a flow to drop them at ingress to the switch with:

```
$ ovs-ofctl add-flow br0 \
    "table=0, dl_src=01:00:00:00:00:00/01:00:00:00:00:00, actions=drop"
```

A switch should also not forward IEEE 802.1D Spanning Tree Protocol (STP) packets, so we can also add a flow to drop those and other packets with reserved multicast protocols:

```
$ ovs-ofctl add-flow br0 \
    "table=0, dl_dst=01:80:c2:00:00:00/ff:ff:ff:ff:ff:f0, actions=drop"
```

We could add flows to drop other protocols, but these demonstrate the pattern.

We need one more flow, with a priority lower than the default, so that flows that don't match either of the "drop" flows we added above go on to pipeline stage 1 in OpenFlow table 1:

```
$ ovs-ofctl add-flow br0 "table=0, priority=0, actions=resubmit(,1)"
```

Note: The "resubmit" action is an Open vSwitch extension to OpenFlow.

3.3.7 Testing Table 0

If we were using Open vSwitch to set up a physical or a virtual switch, then we would naturally test it by sending packets through it one way or another, perhaps with common network testing tools like `ping` and `tcpdump` or more specialized tools like Scapy. That's difficult with our simulated switch, since it's not visible to the operating system.

But our simulated switch has a few specialized testing tools. The most powerful of these tools is `ofproto/trace`. Given a switch and the specification of a flow, `ofproto/trace` shows, step-by-step, how such a flow would be treated as it goes through the switch.

Example 1

Try this command:

```
$ ovs-appctl ofproto/trace br0 in_port=1,dl_dst=01:80:c2:00:00:05
```

The output should look something like this:

```
Flow: in_port=1,vlan_tci=0x0000,dl_src=00:00:00:00:00:00,dl_dst=01:80:c2:00:00:05,dl_
↳type=0x0000

bridge("br0")
-----
 0. dl_dst=01:80:c2:00:00:00/ff:ff:ff:ff:ff:f0, priority 32768
    drop

Final flow: unchanged
Megaflow: recirc_id=0,in_port=1,dl_src=00:00:00:00:00:00/01:00:00:00:00:00,dl_
↳dst=01:80:c2:00:00:00/ff:ff:ff:ff:ff:f0,dl_type=0x0000
Datapath actions: drop
```

The first line shows the flow being traced, in slightly greater detail than specified on the command line. It is mostly zeros because unspecified fields default to zeros.

The second group of lines shows the packet's trip through bridge br0. We see, in table 0, the OpenFlow flow that the fields matched, along with its priority, followed by its actions, one per line. In this case, we see that this packet that has a reserved multicast destination address matches the flow that drops those packets.

The final block of lines summarizes the results, which are not very interesting here.

Example 2

Try another command:

```
$ ovs-appctl ofproto/trace br0 in_port=1,dl_dst=01:80:c2:00:00:10
```

The output should be:

```
Flow: in_port=1,vlan_tci=0x0000,dl_src=00:00:00:00:00:00,dl_dst=01:80:c2:00:00:10,dl_
↪type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. No match.
   drop

Final flow: unchanged
Megaflow: recirc_id=0,in_port=1,dl_src=00:00:00:00:00:00/01:00:00:00:00:00,dl_
↪dst=01:80:c2:00:00:10/ff:ff:ff:ff:ff:f0,dl_type=0x0000
Datapath actions: drop
```

This time the flow we handed to `ofproto/trace` doesn't match any of our "drop" flows in table 0, so it falls through to the low-priority "resubmit" flow. The "resubmit" causes a second lookup in OpenFlow table 1, described by the block of text that starts with "1." We haven't yet added any flows to OpenFlow table 1, so no flow actually matches in the second lookup. Therefore, the packet is still actually dropped, which means that the externally observable results would be identical to our first example.

3.3.8 Implementing Table 1: VLAN Input Processing

A packet that enters table 1 has already passed basic validation in table 0. The purpose of table 1 is validate the packet's VLAN, based on the VLAN configuration of the switch port through which the packet entered the switch. We will also use it to attach a VLAN header to packets that arrive on an access port, which allows later processing stages to rely on the packet's VLAN always being part of the VLAN header, reducing special cases.

Let's start by adding a low-priority flow that drops all packets, before we add flows that pass through acceptable packets. You can think of this as a "default drop" flow:

```
$ ovs-ofctl add-flow br0 "table=1, priority=0, actions=drop"
```

Our trunk port `p1`, on OpenFlow port 1, is an easy case. `p1` accepts any packet regardless of whether it has a VLAN header or what the VLAN was, so we can add a flow that resubmits everything on input port 1 to the next table:

```
$ ovs-ofctl add-flow br0 \
    "table=1, priority=99, in_port=1, actions=resubmit(,2) "
```

On the access ports, we want to accept any packet that has no VLAN header, tag it with the access port's VLAN number, and then pass it along to the next stage:

```
$ ovs-ofctl add-flows br0 - <<'EOF'
table=1, priority=99, in_port=2, vlan_tci=0, actions=mod_vlan_vid:20, resubmit(,2)
table=1, priority=99, in_port=3, vlan_tci=0, actions=mod_vlan_vid:30, resubmit(,2)
table=1, priority=99, in_port=4, vlan_tci=0, actions=mod_vlan_vid:30, resubmit(,2)
EOF
```

We don't write any flows that match packets with 802.1Q that enter this stage on any of the access ports, so the "default drop" flow we added earlier causes them to be dropped, which is ordinarily what we want for access ports.

Note: Another variation of access ports allows ingress of packets tagged with VLAN 0 (aka 802.1p priority tagged packets). To allow such packets, replace `vlan_tci=0` by `vlan_tci=0/0xffff` above.

3.3.9 Testing Table 1

`ofproto/trace` allows us to test the ingress VLAN flows that we added above.

Example 1: Packet on Trunk Port

Here's a test of a packet coming in on the trunk port:

```
$ ovs-appctl ofproto/trace br0 in_port=1,vlan_tci=5
```

The output shows the lookup in table 0, the resubmit to table 1, and the resubmit to table 2 (which does nothing because we haven't put anything there yet):

```
Flow: in_port=1,vlan_tci=0x0005,dl_src=00:00:00:00:00:00,dl_dst=00:00:00:00:00:00,dl_
→type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. in_port=1, priority 99
   resubmit(,2)
2. No match.
   drop

Final flow: unchanged
Megaflow: recirc_id=0,in_port=1,dl_src=00:00:00:00:00:00/01:00:00:00:00:00,dl_
→dst=00:00:00:00:00:00/ff:ff:ff:ff:ff:f0,dl_type=0x0000
Datapath actions: drop
```

Example 2: Valid Packet on Access Port

Here's a test of a valid packet (a packet without an 802.1Q header) coming in on access port p2:

```
$ ovs-appctl ofproto/trace br0 in_port=2
```

The output is similar to that for the previous case, except that it additionally tags the packet with p2's VLAN 20 before it passes it along to table 2:

```
Flow: in_port=2,vlan_tci=0x0000,dl_src=00:00:00:00:00:00,dl_dst=00:00:00:00:00:00,dl_
↳type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. in_port=2,vlan_tci=0x0000, priority 99
   mod_vlan_vid:20
   resubmit(,2)
2. No match.
   drop

Final flow: in_port=2,dl_vlan=20,dl_vlan_pcp=0,dl_src=00:00:00:00:00:00,dl_
↳dst=00:00:00:00:00:00,dl_type=0x0000
Megaflow: recirc_id=0,in_port=2,vlan_tci=0x0000,dl_src=00:00:00:00:00:00/
↳01:00:00:00:00:00,dl_dst=00:00:00:00:00:00/ff:ff:ff:ff:ff:f0,dl_type=0x0000
Datapath actions: drop
```

Example 3: Invalid Packet on Access Port

This tests an invalid packet (one that includes an 802.1Q header) coming in on access port p2:

```
$ ovs-appctl ofproto/trace br0 in_port=2,vlan_tci=5
```

The output shows the packet matching the default drop flow:

```
Flow: in_port=2,vlan_tci=0x0005,dl_src=00:00:00:00:00:00,dl_dst=00:00:00:00:00:00,dl_
↳type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. priority 0
   drop

Final flow: unchanged
Megaflow: recirc_id=0,in_port=2,vlan_tci=0x0005,dl_src=00:00:00:00:00:00/
↳01:00:00:00:00:00,dl_dst=00:00:00:00:00:00/ff:ff:ff:ff:ff:f0,dl_type=0x0000
Datapath actions: drop
```

3.3.10 Implementing Table 2: MAC+VLAN Learning for Ingress Port

This table allows the switch we're implementing to learn that the packet's source MAC is located on the packet's ingress port in the packet's VLAN.

Note: This table is a good example why table 1 added a VLAN tag to packets that entered the switch through an access port. We want to associate a MAC+VLAN with a port regardless of whether the VLAN in question was originally part of the packet or whether it was an assumed VLAN associated with an access port.

It only takes a single flow to do this. The following command adds it:


```
$ ovs-ofctl add-flow br0 \
    "table=2 actions=learn(table=10, NXM_OF_VLAN_TCI[0..11], \
        NXM_OF_ETH_DST[]=NXM_OF_ETH_SRC[], \
        load:NXM_OF_IN_PORT[]->NXM_NX_REG0[0..15]), \
    resubmit(,3) "
```

The “learn” action (an Open vSwitch extension to OpenFlow) modifies a flow table based on the content of the flow currently being processed. Here’s how you can interpret each part of the “learn” action above:

table=10 Modify flow table 10. This will be the MAC learning table.

NXM_OF_VLAN_TCI[0..11] Make the flow that we add to flow table 10 match the same VLAN ID that the packet we’re currently processing contains. This effectively scopes the MAC learning entry to a single VLAN, which is the ordinary behavior for a VLAN-aware switch.

NXM_OF_ETH_DST[]=NXM_OF_ETH_SRC[] Make the flow that we add to flow table 10 match, as Ethernet destination, the Ethernet source address of the packet we’re currently processing.

load:NXM_OF_IN_PORT[]->NXM_NX_REG0[0..15] Whereas the preceding parts specify fields for the new flow to match, this specifies an action for the flow to take when it matches. The action is for the flow to load the ingress port number of the current packet into register 0 (a special field that is an Open vSwitch extension to OpenFlow).

Note: A real use of “learn” for MAC learning would probably involve two additional elements. First, the “learn” action would specify a `hard_timeout` for the new flow, to enable a learned MAC to eventually expire if no new packets were seen from a given source within a reasonable interval. Second, one would usually want to limit resource consumption by using the `Flow_Table` table in the Open vSwitch configuration database to specify a maximum number of flows in table 10.

This definitely calls for examples.

3.3.11 Testing Table 2

Example 1

Try the following test command:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,vlan_tci=20,d1_src=50:00:00:00:00:01 -generate
```

The output shows that “learn” was executed in table 2 and the particular flow that was added:

```
Flow: in_port=1,vlan_tci=0x0014,d1_src=50:00:00:00:00:01,d1_dst=00:00:00:00:00:00,d1_
↳type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. in_port=1, priority 99
   resubmit(,2)
2. priority 32768
   learn(table=10,NXM_OF_VLAN_TCI[0..11],NXM_OF_ETH_DST[]=NXM_OF_ETH_SRC[],load:NXM_
↳OF_IN_PORT[]->NXM_NX_REG0[0..15])
   -> table=10 vlan_tci=0x0014/0x0fff,d1_dst=50:00:00:00:00:01 priority=32768_
↳actions=load:0x1->NXM_NX_REG0[0..15]
```

(continues on next page)

(continued from previous page)

```

    resubmit(,3)
3. No match.
    drop

Final flow: unchanged
Megaflow: recirc_id=0,in_port=1,vlan_tci=0x0014/0x1fff,dl_src=50:00:00:00:00:01,dl_
↳dst=00:00:00:00:00:00/ff:ff:ff:ff:ff:f0,dl_type=0x0000
Datapath actions: drop

```

The `-generate` keyword is new. Ordinarily, `ofproto/trace` has no side effects: “output” actions do not actually output packets, “learn” actions do not actually modify the flow table, and so on. With `-generate`, though, `ofproto/trace` does execute “learn” actions. That’s important now, because we want to see the effect of the “learn” action on table 10. You can see that by running:

```
$ ovs-ofctl dump-flows br0 table=10
```

which (omitting the duration and `idle_age` fields, which will vary based on how soon you ran this command after the previous one, as well as some other uninteresting fields) prints something like:

```

NXST_FLOW reply (xid=0x4):
  table=10, vlan_tci=0x0014/0x0fff,dl_dst=50:00:00:00:00:01 actions=load:0x1->NXM_NX_
↳REG0[0..15]

```

You can see that the packet coming in on VLAN 20 with source MAC 50:00:00:00:00:01 became a flow that matches VLAN 20 (written in hexadecimal) and destination MAC 50:00:00:00:00:01. The flow loads port number 1, the input port for the flow we tested, into register 0.

Example 2

Here’s a second test command:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=2,dl_src=50:00:00:00:00:01 -generate
```

The flow that this command tests has the same source MAC and VLAN as example 1, although the VLAN comes from an access port VLAN rather than an 802.1Q header. If we again dump the flows for table 10 with:

```
$ ovs-ofctl dump-flows br0 table=10
```

then we see that the flow we saw previously has changed to indicate that the learned port is port 2, as we would expect:

```

NXST_FLOW reply (xid=0x4):
  table=10, vlan_tci=0x0014/0x0fff,dl_dst=50:00:00:00:00:01 actions=load:0x2->NXM_NX_
↳REG0[0..15]

```

3.3.12 Implementing Table 3: Look Up Destination Port

This table figures out what port we should send the packet to based on the destination MAC and VLAN. That is, if we’ve learned the location of the destination (from table 2 processing some previous packet with that destination as its source), then we want to send the packet there.

We need only one flow to do the lookup:

```
$ ovs-ofctl add-flow br0 \
    "table=3 priority=50 actions=resubmit(,10), resubmit(,4) "
```

The flow's first action resubmits to table 10, the table that the "learn" action modifies. As you saw previously, the learned flows in this table write the learned port into register 0. If the destination for our packet hasn't been learned, then there will be no matching flow, and so the "resubmit" turns into a no-op. Because registers are initialized to 0, we can use a register 0 value of 0 in our next pipeline stage as a signal to flood the packet.

The second action resubmits to table 4, continuing to the next pipeline stage.

We can add another flow to skip the learning table lookup for multicast and broadcast packets, since those should always be flooded:

```
$ ovs-ofctl add-flow br0 \
    "table=3 priority=99 dl_dst=01:00:00:00:00:00/01:00:00:00:00:00 \
    actions=resubmit(,4) "
```

Note: We don't strictly need to add this flow, because multicast addresses will never show up in our learning table. (In turn, that's because we put a flow into table 0 to drop packets that have a multicast source address.)

3.3.13 Testing Table 3

Example

Here's a command that should cause OVS to learn that `f0:00:00:00:00:01` is on `p1` in VLAN 20:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_vlan=20,dl_src=f0:00:00:00:00:01,dl_dst=90:00:00:00:00:01 \
    -generate
```

The output shows (from the "no match" looking up the resubmit to table 10) that the flow's destination was unknown:

```
Flow: in_port=1,dl_vlan=20,dl_vlan_pcp=0,dl_src=f0:00:00:00:00:01,dl_
↳dst=90:00:00:00:00:01,dl_type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. in_port=1, priority 99
   resubmit(,2)
2. priority 32768
   learn(table=10,NXM_OF_VLAN_TCI[0..11],NXM_OF_ETH_DST[]=NXM_OF_ETH_SRC[],load:NXM_
↳OF_IN_PORT[]->NXM_NX_REG0[0..15])
   -> table=10 vlan_tci=0x0014/0x0fff,dl_dst=f0:00:00:00:00:01 priority=32768_
↳actions=load:0x1->NXM_NX_REG0[0..15]
   resubmit(,3)
3. priority 50
   resubmit(,10)
10. No match.
    drop
    resubmit(,4)
4. No match.
    drop
```

(continues on next page)

(continued from previous page)

```
Final flow: unchanged
Megaflow: recirc_id=0,in_port=1,dl_vlan=20,dl_src=f0:00:00:00:00:01,dl_
↳dst=90:00:00:00:00:01,dl_type=0x0000
Datapath actions: drop
```

There are two ways that you can verify that the packet's source was learned. The most direct way is to dump the learning table with:

```
$ ovs-ofctl dump-flows br0 table=10
```

which ought to show roughly the following, with extraneous details removed:

```
table=10, vlan_tci=0x0014/0x0fff,dl_dst=f0:00:00:00:00:01 actions=load:0x1->NXM_NX_
↳REG0[0..15]
```

Note: If you tried the examples for the previous step, or if you did some of your own experiments, then you might see additional flows there. These additional flows are harmless. If they bother you, then you can remove them with *ovs-ofctl del-flows br0 table=10*.

The other way is to inject a packet to take advantage of the learning entry. For example, we can inject a packet on p2 whose destination is the MAC address that we just learned on p1:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=2,dl_src=90:00:00:00:00:01,dl_dst=f0:00:00:00:00:01 -generate
```

Here is this command's output. Take a look at the lines that trace the `resubmit(,10)`, showing that the packet matched the learned flow for the first MAC we used, loading the OpenFlow port number for the learned port p1 into register 0:

```
Flow: in_port=2,vlan_tci=0x0000,dl_src=90:00:00:00:00:01,dl_dst=f0:00:00:00:00:01,dl_
↳type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. in_port=2,vlan_tci=0x0000, priority 99
   mod_vlan_vid:20
   resubmit(,2)
2. priority 32768
   learn(table=10,NXM_OF_VLAN_TCI[0..11],NXM_OF_ETH_DST[]=NXM_OF_ETH_SRC[],load:NXM_
↳OF_IN_PORT[]->NXM_NX_REG0[0..15])
   -> table=10 vlan_tci=0x0014/0x0fff,dl_dst=90:00:00:00:00:01 priority=32768_
↳actions=load:0x2->NXM_NX_REG0[0..15]
   resubmit(,3)
3. priority 50
   resubmit(,10)
   10. vlan_tci=0x0014/0x0fff,dl_dst=f0:00:00:00:00:01, priority 32768
      load:0x1->NXM_NX_REG0[0..15]
   resubmit(,4)
4. No match.
   drop
```

(continues on next page)

(continued from previous page)

```
Final flow: reg0=0x1,in_port=2,d_l_vlan=20,d_l_vlan_pcp=0,d_l_src=90:00:00:00:00:01,d_l_
→dst=f0:00:00:00:00:01,d_l_type=0x0000
Megaflow: recirc_id=0,in_port=2,vlan_tci=0x0000,d_l_src=90:00:00:00:00:01,d_l_
→dst=f0:00:00:00:00:01,d_l_type=0x0000
Datapath actions: drop
```

If you read the commands above carefully, then you might have noticed that they simply have the Ethernet source and destination addresses exchanged. That means that if we now rerun the first `ovs-appctl` command above, e.g.:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,d_l_vlan=20,d_l_src=f0:00:00:00:00:01,d_l_dst=90:00:00:00:00:01 \
    -generate
```

then we see in the output, looking at the indented “load” action executed in table 10, that the destination has now been learned:

```
Flow: in_port=1,d_l_vlan=20,d_l_vlan_pcp=0,d_l_src=f0:00:00:00:00:01,d_l_
→dst=90:00:00:00:00:01,d_l_type=0x0000

bridge("br0")
-----
0. priority 0
   resubmit(,1)
1. in_port=1, priority 99
   resubmit(,2)
2. priority 32768
   learn(table=10,NXM_OF_VLAN_TCI[0..11],NXM_OF_ETH_DST[]=NXM_OF_ETH_SRC[],load:NXM_
→OF_IN_PORT[]->NXM_NX_REG0[0..15])
   → table=10 vlan_tci=0x0014/0x0fff,d_l_dst=f0:00:00:00:00:01 priority=32768,
→actions=load:0x1->NXM_NX_REG0[0..15]
   resubmit(,3)
3. priority 50
   resubmit(,10)
   10. vlan_tci=0x0014/0x0fff,d_l_dst=90:00:00:00:00:01, priority 32768
      load:0x2->NXM_NX_REG0[0..15]
   resubmit(,4)
4. No match.
   drop
```

3.3.14 Implementing Table 4: Output Processing

At entry to stage 4, we know that register 0 contains either the desired output port or is zero if the packet should be flooded. We also know that the packet’s VLAN is in its 802.1Q header, even if the VLAN was implicit because the packet came in on an access port.

The job of the final pipeline stage is to actually output packets. The job is trivial for output to our trunk port `p1`:

```
$ ovs-ofctl add-flow br0 "table=4 reg0=1 actions=1"
```

For output to the access ports, we just have to strip the VLAN header before outputting the packet:

```
$ ovs-ofctl add-flows br0 - <<'EOF'
table=4 reg0=2 actions=strip_vlan,2
table=4 reg0=3 actions=strip_vlan,3
```

(continues on next page)

(continued from previous page)

```
table=4 reg0=4 actions=strip_vlan,4
EOF
```

The only slightly tricky part is flooding multicast and broadcast packets and unicast packets with unlearned destinations. For those, we need to make sure that we only output the packets to the ports that carry our packet's VLAN, and that we include the 802.1Q header in the copy output to the trunk port but not in copies output to access ports:

```
$ ovs-ofctl add-flows br0 - <<'EOF'
table=4 reg0=0 priority=99 dl_vlan=20 actions=1,strip_vlan,2
table=4 reg0=0 priority=99 dl_vlan=30 actions=1,strip_vlan,3,4
table=4 reg0=0 priority=50 actions=1
EOF
```

Note: Our flows rely on the standard OpenFlow behavior that an output action will not forward a packet back out the port it came in on. That is, if a packet comes in on p1, and we've learned that the packet's destination MAC is also on p1, so that we end up with `actions=1` as our actions, the switch will not forward the packet back out its input port. The multicast/broadcast/unknown destination cases above also rely on this behavior.

3.3.15 Testing Table 4

Example 1: Broadcast, Multicast, and Unknown Destination

Try tracing a broadcast packet arriving on p1 in VLAN 30:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_dst=ff:ff:ff:ff:ff:ff,dl_vlan=30
```

The interesting part of the output is the final line, which shows that the switch would remove the 802.1Q header and then output the packet to p3 and p4, which are access ports for VLAN 30:

```
Datapath actions: pop_vlan,3,4
```

Similarly, if we trace a broadcast packet arriving on p3:

```
$ ovs-appctl ofproto/trace br0 in_port=3,dl_dst=ff:ff:ff:ff:ff:ff
```

then we see that it is output to p1 with an 802.1Q tag and then to p4 without one:

```
Datapath actions: push_vlan(vid=30,pcp=0),1,pop_vlan,4
```

Note: Open vSwitch could simplify the datapath actions here to just `4,push_vlan(vid=30,pcp=0),1` but it is not smart enough to do so.

The following are also broadcasts, but the result is to drop the packets because the VLAN only belongs to the input port:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_dst=ff:ff:ff:ff:ff:ff
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_dst=ff:ff:ff:ff:ff:ff,dl_vlan=55
```

Try some other broadcast cases on your own:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_dst=ff:ff:ff:ff:ff:ff,dl_vlan=20
$ ovs-appctl ofproto/trace br0 \
    in_port=2,dl_dst=ff:ff:ff:ff:ff:ff
$ ovs-appctl ofproto/trace br0 \
    in_port=4,dl_dst=ff:ff:ff:ff:ff:ff
```

You can see the same behavior with multicast packets and with unicast packets whose destination has not been learned, e.g.:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=4,dl_dst=01:00:00:00:00:00
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_dst=90:12:34:56:78:90,dl_vlan=20
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_dst=90:12:34:56:78:90,dl_vlan=30
```

Example 2: MAC Learning

Let's follow the same pattern as we did for table 3. First learn a MAC on port p1 in VLAN 30:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_vlan=30,dl_src=10:00:00:00:00:01,dl_dst=20:00:00:00:00:01 \
    -generate
```

You can see from the last line of output that the packet's destination is unknown, so it gets flooded to both p3 and p4, the other ports in VLAN 30:

```
Datapath actions: pop_vlan,3,4
```

Then reverse the MACs and learn the first flow's destination on port p4:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=4,dl_src=20:00:00:00:00:01,dl_dst=10:00:00:00:00:01 -generate
```

The last line of output shows that the this packet's destination is known to be p1, as learned from our previous command:

```
Datapath actions: push_vlan(vid=30,pcp=0),1
```

Now, if we rerun our first command:

```
$ ovs-appctl ofproto/trace br0 \
    in_port=1,dl_vlan=30,dl_src=10:00:00:00:00:01,dl_dst=20:00:00:00:00:01 \
    -generate
```

... we can see that the result is no longer a flood but to the specified learned destination port p4:

```
Datapath actions: pop_vlan,4
```

Contact

bugs@openvswitch.org <http://openvswitch.org/>

3.4 OVN Sandbox

This tutorial shows you how to explore features using `ovs-sandbox` as a simulated test environment. It's assumed that you have an understanding of OVS before going through this tutorial. Detail about OVN is covered in [ovn-architecture](#), but this tutorial lets you quickly see it in action.

3.4.1 Getting Started

For some general information about `ovs-sandbox`, see the “Getting Started” section of *Open vSwitch Advanced Features*.

`ovs-sandbox` does not include OVN support by default. To enable OVN, you must pass the `--ovn` flag. For example, if running it straight from the OVS git tree you would run:

```
$ make sandbox SANDBOXFLAGS="--ovn"
```

Running the sandbox with OVN enabled does the following additional steps to the environment:

1. Creates the `OVN_Northbound` and `OVN_Southbound` databases as described in [ovn-nb\(5\)](#) and [ovn-sb\(5\)](#).
2. Creates a backup server for `OVN_Southbound` database. Sandbox launch screen provides the instructions on accessing the backup database. However access to the backup server is not required to go through the tutorial.
3. Creates the `hardware_vtep` database as described in [vtep\(5\)](#).
4. Runs the `ovn-northd(8)`, `ovn-controller(8)`, and `ovn-controller-vtep(8)` daemons.
5. Makes OVN and VTEP utilities available for use in the environment, including `vtep-ctl(8)`, `ovn-nbctl(8)`, and `ovn-sbctl(8)`.

3.4.2 Using GDB

GDB support is not required to go through the tutorial. See the “Using GDB” section of *Open vSwitch Advanced Features* for more info. Additional flags exist for launching the debugger for the OVN programs:

```
--gdb-ovn-northd
--gdb-ovn-controller
--gdb-ovn-controller-vtep
```

3.4.3 Creating OVN Resources

Once you have `ovs-sandbox` running with OVN enabled, you can start using OVN utilities to create resources in OVN. As an example, we will create an environment that has two logical switches connected by a logical router.

Create the first logical switch with one port:

```
$ ovn-nbctl ls-add sw0
$ ovn-nbctl lsp-add sw0 sw0-port1
$ ovn-nbctl lsp-set-addresses sw0-port1 "50:54:00:00:00:01 192.168.0.2"
```

Create the second logical switch with one port:

```
$ ovn-nbctl ls-add sw1
$ ovn-nbctl lsp-add sw1 sw1-port1
$ ovn-nbctl lsp-set-addresses sw1-port1 "50:54:00:00:00:03 11.0.0.2"
```


Create the logical router and attach both logical switches:

```
$ ovn-nbctl lr-add lr0
$ ovn-nbctl lsp-add lr0 lrp0 00:00:00:00:ff:01 192.168.0.1/24
$ ovn-nbctl lsp-add sw0 lrp0-attachment
$ ovn-nbctl lsp-set-type lrp0-attachment router
$ ovn-nbctl lsp-set-addresses lrp0-attachment 00:00:00:00:ff:01
$ ovn-nbctl lsp-set-options lrp0-attachment router-port=lrp0
$ ovn-nbctl lsp-add lr0 lrp1 00:00:00:00:ff:02 11.0.0.1/24
$ ovn-nbctl lsp-add sw1 lrp1-attachment
$ ovn-nbctl lsp-set-type lrp1-attachment router
$ ovn-nbctl lsp-set-addresses lrp1-attachment 00:00:00:00:ff:02
$ ovn-nbctl lsp-set-options lrp1-attachment router-port=lrp1
```

View a summary of OVN's current logical configuration:

```
$ ovn-nbctl show
  switch 1396cf55-d176-4082-9a55-1c06cef626e4 (sw1)
    port lrp1-attachment
      addresses: ["00:00:00:00:ff:02"]
    port sw1-port1
      addresses: ["50:54:00:00:00:03 11.0.0.2"]
  switch 2c9d6d03-09fc-4e32-8da6-305f129b0d53 (sw0)
    port lrp0-attachment
      addresses: ["00:00:00:00:ff:01"]
    port sw0-port1
      addresses: ["50:54:00:00:00:01 192.168.0.2"]
  router f8377e8c-f75e-4fc8-8751-f3ea03c6dd98 (lr0)
    port lrp0
      mac: "00:00:00:00:ff:01"
      networks: ["192.168.0.1/24"]
    port lrp1
      mac: "00:00:00:00:ff:02"
      networks: ["11.0.0.1/24"]
```

The tutorial directory of the OVS source tree includes a script that runs all of the commands for you:

```
$ ./ovn-setup.sh
```

3.4.4 Using ovn-trace

Once you have configured resources in OVN, try using `ovn-trace` to see how OVN would process a sample packet through its logical pipeline.

For example, we can trace an IP packet from `sw0-port1` to `sw1-port1`. The `--minimal` output shows each visible action performed on the packet, which includes:

1. The logical router will decrement the IP TTL field.
2. The logical router will change the source and destination MAC addresses to reflect the next hop.
3. The packet will be output to `sw1-port1`.

```
$ ovn-trace --minimal sw0 'inport == "sw0-port1" \
> && eth.src == 50:54:00:00:00:01 && ip4.src == 192.168.0.2 \
> && eth.dst == 00:00:00:00:ff:01 && ip4.dst == 11.0.0.2 \
> && ip.ttl == 64'
```

(continues on next page)

(continued from previous page)

```
# ip,reg14=0x1,vlan_tci=0x0000,d1_src=50:54:00:00:00:01,d1_dst=00:00:00:00:ff:01,nw_  
→src=192.168.0.2,nw_dst=11.0.0.2,nw_proto=0,nw_tos=0,nw_ecn=0,nw_ttl=64  
ip.ttl--;  
eth.src = 00:00:00:00:ff:02;  
eth.dst = 50:54:00:00:00:03;  
output("sw1-port1");
```

The `ovn-trace` utility can also provide much more detail on how the packet would be processed through OVN's logical pipeline, as well as correlate that to OpenFlow flows programmed by `ovn-controller`. See the [ovn-trace\(8\)](#) man page for more detail.

3.5 OVN OpenStack Tutorial

This tutorial demonstrates how OVN works in an OpenStack “DevStack” environment. It was tested with the “master” branches of DevStack and Open vSwitch near the beginning of May 2017. Anyone using an earlier version is likely to encounter some differences. In particular, we noticed some shortcomings in OVN utilities while writing the tutorial and pushed out some improvements, so it's best to use recent Open vSwitch at least from that point of view.

The goal of this tutorial is to demonstrate OVN in an end-to-end way, that is, to show how it works from the cloud management system at the top (in this case, OpenStack and specifically its Neutron networking subsystem), through the OVN northbound and southbound databases, to the bottom at the OVN local controller and Open vSwitch data plane. We hope that this demonstration makes it easier for users and potential users to understand how OVN works and how to debug and troubleshoot it.

In addition to new material, this tutorial incorporates content from `testing.rst` in OpenStack networking-ovn, by Russell Bryant and others. Without that example, this tutorial could not have been written.

We provide enough details in the tutorial that you should be able to fully follow along, by creating a DevStack VM and cloning DevStack and so on. If you want to do this, start out from [Setting Up DevStack](#) below.

3.5.1 Setting Up DevStack

This section explains how to install DevStack, a kind of OpenStack packaging for developers, in a way that allows you to follow along with the tutorial in full.

Unless you have a spare computer laying about, it's easiest to install DevStack in a virtual machine. This tutorial was built using a VM implemented by KVM and managed by virt-manager. I recommend configuring the VM configured for the x86-64 architecture, 4 GB RAM, 2 VCPUs, and a 20 GB virtual disk.

Note: If you happen to run your Linux-based host with 32-bit userspace, then you will have some special issues, even if you use a 64-bit kernel:

- You may find that you can get 32-bit DevStack VMs to work to some extent, but I personally got tired of finding workarounds. I recommend running your VMs in 64-bit mode. To get this to work, I had to go to the CPUs tab for the VM configuration in virt-manager and change the CPU model from the one originally listed to “Hypervisor Default” (it is curious that this is not the default!).
- On a host with 32-bit userspace, KVM supports VMs with at most 2047 MB RAM. This is adequate, barely, to start DevStack, but it is not enough to run multiple (nested) VMs. To prevent out-of-memory failures, set up extra swap space in the guest. For example, to add 2 GB swap:

```
$ sudo dd if=/dev/zero of=/swapfile bs=1M count=2048
$ sudo mkswap /swapfile
$ sudo swapon /swapfile
```

and then add a line like this to `/etc/fstab` to add the new swap automatically upon reboot:

```
/swapfile swap swap defaults 0 0
```

Here are step-by-step instructions to get started:

1. Install a VM.

I tested these instructions with Centos 7.3. Download the “minimal install” ISO and booted it. The install is straightforward. Be sure to enable networking, and set a host name, such as “ovn-devstack-1”. Add a regular (non-root) user, and check the box “Make this user administrator”. Also, set your time zone.

2. You can SSH into the DevStack VM, instead of running from a console. I recommend it because it’s easier to cut and paste commands into a terminal than a VM console. You might also consider using a very wide terminal, perhaps 160 columns, to keep tables from wrapping.

To improve convenience further, you can make it easier to log in with the following steps, which are optional:

- (a) On your host, edit your `~/.ssh/config`, adding lines like the following:

```
Host ovn-devstack-1
    Hostname VMIP
    User VMUSER
```

where VMIP is the VM’s IP address and VMUSER is your username inside the VM. (You can omit the `User` line if your username is the same in the host and the VM.) After you do this, you can SSH to the VM by name, e.g. `ssh ovn-devstack-1`, and if command-line completion is set up in your host shell, you can shorten that to something like `ssh ovn` followed by hitting the Tab key.

- (b) If you have SSH public key authentication set up, with an SSH agent, run on your host:

```
$ ssh-copy-id ovn-devstack-1
```

and type your password once. Afterward, you can log in without typing your password again.

(If you don’t already use SSH public key authentication and an agent, consider looking into it—it will save you time in the long run.)

- (c) Optionally, inside the VM, append the following to your `~/.bash_profile`:

```
. $HOME/devstack/openrc admin
```

It will save you running it by hand each time you log in. But it also prints garbage to the console, which can screw up services like `ssh-copy-id`, so be careful.

2. Boot into the installed system and log in as the regular user, then install Git:

```
$ sudo yum install git
```

Note: If you installed a 32-bit i386 guest (against the advice above), install a non-PAE kernel and reboot into it at this point:

```
$ sudo yum install kernel-core kernel-devel
$ sudo reboot
```

Be sure to select the non-PAE kernel from the list at boot. Without this step, DevStack will fail to install properly later.

3. Get copies of DevStack and OVN and set them up:

```
$ git clone http://git.openstack.org/openstack-dev/devstack.git
$ git clone http://git.openstack.org/openstack/networking-ovn.git
$ cd devstack
$ cp ../networking-ovn/devstack/local.conf.sample local.conf
```

Note: If you installed a 32-bit i386 guest (against the advice above), at this point edit `local.conf` to add the following line:

```
CIRROS_ARCH=i386
```

4. Initialize DevStack:

```
$ ./stack.sh
```

This will spew many screenfuls of text, and the first time you run it, it will download lots of software from the Internet. The output should eventually end with something like this:

```
This is your host IP address: 172.16.189.6
This is your host IPv6 address: ::1
Horizon is now available at http://172.16.189.6/dashboard
Keystone is serving at http://172.16.189.6/identity/
The default users are: admin and demo
The password: password
2017-03-09 15:10:54.117 | stack.sh completed in 2110 seconds.
```

If there's some kind of failure, you can restart by running `./stack.sh` again. It won't restart exactly where it left off, but steps up to the one where it failed will skip the download steps. (Sometimes blindly restarting after a failure will allow it to succeed.) If you reboot your VM, you need to rerun this command. (If you run into trouble with `stack.sh` after rebooting your VM, try running `./unstack.sh`.)

At this point you can navigate a web browser on your host to the Horizon dashboard URL. Many OpenStack operations can be initiated from this UI. Feel free to explore, but this tutorial focuses on the alternative command-line interfaces because they are easier to explain and to cut and paste.

5. As of this writing, you need to run the following to fix a problem with using VM consoles from the OpenStack web instance:

```
$ (cd /opt/stack/noVNC && git checkout v0.6.0)
```

See <https://serenity-networks.com/how-to-fix-setkeycodes-00-and-unknown-key-pressed-console-errors-on-openstack/> for more details.

6. The firewall in the VM by default allows SSH access but not HTTP. You will probably want HTTP access to use the OpenStack web interface. The following command enables that. (It also enables every other kind of network access, so if you're concerned about security then you might want to find a more targeted approach.)

```
$ sudo iptables -F
```

(You need to re-run this if you reboot the VM.)

7. To use OpenStack command line utilities in the tutorial, run:

```
$ . ~/devstack/openrc admin
```

This needs to be re-run each time you log in (but see the following section).

3.5.2 DevStack preliminaries

Before we really jump in, let's set up a couple of things in DevStack. This is the first real test that DevStack is working, so if you get errors from any of these commands, it's a sign that `stack.sh` didn't finish properly, or perhaps that you didn't run the `openrc admin` command at the end of the previous instructions.

If you stop and restart DevStack via `unstack.sh` followed by `stack.sh`, you have to rerun these steps.

1. For SSH access to the VMs we're going to create, we'll need a SSH keypair. Later on, we'll get OpenStack to install this keypair into VMs. Create one with:

```
$ openstack keypair create demo > ~/id_rsa_demo
$ chmod 600 ~/id_rsa_demo
```

2. By default, DevStack security groups drop incoming traffic, but to test networking in a reasonable way we need to enable it. You only need to actually edit one particular security group, but DevStack creates multiple and it's somewhat difficult to figure out which one is important because all of them are named "default". So, the following adds rules to allow SSH and ICMP traffic into **every** security group:

```
$ for group in $(openstack security group list -f value -c ID); do \
openstack security group rule create --ingress --ethertype IPv4 --dst-port 22 --
→protocol tcp $group; \
openstack security group rule create --ingress --ethertype IPv4 --protocol ICMP
→$group; \
done
```

3. Later on, we're going to create some VMs and we'll need an operating system image to install. DevStack comes with a very simple image built-in, called "cirros", which works fine. We need to get the UUID for this image. Our later commands assume shell variable `IMAGE_ID` holds this UUID. You can set this by hand, e.g.:

```
$ openstack image list
+-----+-----+-----+
| ID                               | Name                               | Status |
+-----+-----+-----+
| 77f37d2c-3d6b-4e99-a01b-1fa5d78d1fa1 | cirros-0.3.5-x86_64-disk          | active |
+-----+-----+-----+
$ IMAGE_ID=73ca34f3-63c4-4c10-a62f-4540afc24eaa
```

or by parsing CLI output:

```
$ IMAGE_ID=$(openstack image list -f value -c ID)
```

Note: Your image ID will differ from the one above, as will every UUID in this tutorial. They will also change every time you run `stack.sh`. The UUIDs are generated randomly.

3.5.3 Shortening UUIDs

OpenStack, OVN, and Open vSwitch all really like UUIDs. These are great for uniqueness, but 36-character strings are terrible for readability. Statistically, just the first few characters are enough for uniqueness in small environments,

so let's define a helper to make things more readable:

```
$ abbrev() { a='[0-9a-fA-F]' b=$a$a c=$b$b; sed "s/$b-$c-$c-$c-$c$c$c//g"; }
```

You can use this as a filter to abbreviate UUIDs. For example, use it to abbreviate the above image list:

```
$ openstack image list -f yaml | abbrev
- ID: 77f37d
  Name: cirros-0.3.5-x86_64-disk
  Status: active
```

The command above also adds `-f yaml` to switch to YAML output format, because abbreviating UUIDs screws up the default table-based formatting and because YAML output doesn't produce wrap columns across lines and therefore is easier to cut and paste.

3.5.4 Overview

Now that DevStack is ready, with OVN set up as the networking back-end, here's an overview of what we're going to do in the remainder of the demo, all via OpenStack:

1. Switching: Create an OpenStack network `n1` and VMs `a` and `b` attached to it.
An OpenStack network is a virtual switch; it corresponds to an OVN logical switch.
2. Routing: Create a second OpenStack network `n2` and VM `c` attached to it, then connect it to network `n1` by creating an OpenStack router and attaching `n1` and `n2` to it.
3. Gateways: Make VMs `a` and `b` available via an external network.
4. IPv6: Add IPv6 addresses to our VMs to demonstrate OVN support for IPv6 routing.
5. ACLs: Add and modify OpenStack stateless and stateful rules in security groups.
6. DHCP: How it works in OVN.
7. Further directions: Adding more compute nodes.

At each step, we will take a look at how the features in question work from OpenStack's Neutron networking layer at the top to the data plane layer at the bottom. From the highest to lowest level, these layers and the software components that connect them are:

- OpenStack Neutron, which as the top level in the system is the authoritative source of the virtual network configuration.

We will use OpenStack's `openstack` utility to observe and modify Neutron and other OpenStack configuration.

- `networking-ovn`, the Neutron driver that interfaces with OVN and translates the internal Neutron representation of the virtual network into OVN's representation and pushes that representation down the OVN northbound database.

In this tutorial it's rarely worth distinguishing Neutron from `networking-ovn`, so we usually don't break out this layer separately.

- The OVN Northbound database, aka NB DB. This is an instance of OVSDB, a simple general-purpose database that is used for multiple purposes in Open vSwitch and OVN. The NB DB's schema is in terms of networking concepts such as switches and routers. The NB DB serves the purpose that in other systems might be filled by some kind of API; for example, in place of calling an API to create or delete a logical switch, `networking-ovn` performs these operations by inserting or deleting a row in the NB DB's `Logical_Switch` table.

We will use OVN's `ovn-nbctl` utility to observe the NB DB. (We won't directly modify data at this layer or below. Because configuration trickles down from Neutron through the stack, the right way to make changes is to use the `openstack` utility or another OpenStack interface and then wait for them to percolate through to lower layers.)

- The `ovn-northd` daemon, a program that runs centrally and translates the NB DB's network representation into the lower-level representation used by the OVN Southbound database in the next layer. The details of this daemon are usually not of interest, although without it OVN will not work, so this tutorial does not often mention it.
- The OVN Southbound database, aka SB DB, which is also an OVSDB database. Its schema is very different from the NB DB. Instead of familiar networking concepts, the SB DB defines the network in terms of collections of match-action rules called "logical flows", which while similar in concept to OpenFlow flows use logical concepts, such as virtual machine instances, in place of physical concepts like physical Ethernet ports.

We will use OVN's `ovn-sbctl` utility to observe the SB DB.

- The `ovn-controller` daemon. A copy of `ovn-controller` runs on each hypervisor. It reads logical flows from the SB DB, translates them into OpenFlow flows, and sends them to Open vSwitch's `ovs-vswitchd` daemon. Like `ovn-northd`, usually the details of what this daemon are not of interest, even though it's important to the operation of the system.
- `ovs-vswitchd`. This program runs on each hypervisor. It is the core of Open vSwitch, which processes packets according to the OpenFlow flows set up by `ovn-controller`.
- Open vSwitch datapath. This is essentially a cache designed to accelerate packet processing. Open vSwitch includes a few different datapaths but OVN installations typically use one based on the Open vSwitch Linux kernel module.

3.5.5 Switching

Switching is the basis of networking in the real world and in virtual networking as well. OpenStack calls its concept of a virtual switch a "network", and OVN calls its corresponding concept a "logical switch".

In this step, we'll create an OpenStack network `n1`, then create VMs `a` and `b` and attach them to `n1`.

Creating network `n1`

Let's start by creating the network:

```
$ openstack network create --project admin --provider-network-type geneve n1
```

OpenStack needs to know the subnets that a network serves. We inform it by creating subnet objects. To keep it simple, let's give our network a single subnet for the 10.1.1.0/24 network. We have to give it a name, in this case `n1subnet`:

```
$ openstack subnet create --subnet-range 10.1.1.0/24 --network n1 n1subnet
```

If you ask Neutron to show us the available networks, we see `n1` as well as the two networks that DevStack creates by default:

```
$ openstack network list -f yaml | abbrev
- ID: 5b6baf
  Name: n1
  Subnets: 5e67e7
- ID: c02c4d
```

(continues on next page)

(continued from previous page)

```

Name: private
Subnets: d88a34, fd87f9
- ID: dlac28
Name: public
Subnets: 0b1e79, c87dc1

```

Neutron pushes this network setup down to the OVN northbound database. We can use `ovn-nbctl show` to see an overview of what's in the NB DB:

```

$ ovn-nbctl show | abbrev
switch 5b3d5f (neutron-c02c4d) (aka private)
  port b256dd
    type: router
    router-port: lrp-b256dd
  port f264e7
    type: router
    router-port: lrp-f264e7
switch 2579f4 (neutron-dlac28) (aka public)
  port provnet-dlac28
    type: localnet
    addresses: ["unknown"]
  port ae9b52
    type: router
    router-port: lrp-ae9b52
switch 3eb263 (neutron-5b6baf) (aka n1)
router c59ad2 (neutron-9b057f) (aka router1)
  port lrp-ae9b52
    mac: "fa:16:3e:b2:d2:67"
    networks: ["172.24.4.9/24", "2001:db8::b/64"]
  port lrp-b256dd
    mac: "fa:16:3e:35:33:db"
    networks: ["fdb0:5860:4ba8::1/64"]
  port lrp-f264e7
    mac: "fa:16:3e:fc:c8:da"
    networks: ["10.0.0.1/26"]
  nat 80914c
    external ip: "172.24.4.9"
    logical ip: "10.0.0.0/26"
    type: "snat"

```

This output shows that OVN has three logical switches, each of which corresponds to a Neutron network, and a logical router that corresponds to the Neutron router that DevStack creates by default. The logical switch that corresponds to our new network `n1` has no ports yet, because we haven't added any. The `public` and `private` networks that DevStack creates by default have router ports that connect to the logical router.

Using `ovn-northd`, OVN translates the NB DB's high-level switch and router concepts into lower-level concepts of "logical datapaths" and logical flows. There's one logical datapath for each logical switch or router:

```

$ ovn-sbctl list datapath_binding | abbrev
__uuid          : 0ad69d
external_ids     : {logical-switch="5b3d5f", name="neutron-c02c4d", "name2
↪="private}
tunnel_key       : 1

__uuid          : a8a758
external_ids     : {logical-switch="3eb263", name="neutron-5b6baf", "name2"="n1"}

```

(continues on next page)

(continued from previous page)

```

tunnel_key      : 4

_uuid          : 191256
external_ids    : {logical-switch="2579f4", name="neutron-d1ac28", "name2"=public}
tunnel_key      : 3

_uuid          : b87bec
external_ids    : {logical-router="c59ad2", name="neutron-9b057f", "name2"=
↳ "router1"}
tunnel_key      : 2

```

This output lists the NB DB UUIDs in `external_ids:logical-switch` and Neutron UUIDs in `external_ids:uuid`. We can dive in deeper by viewing the OVN logical flows that implement a logical switch. Our new logical switch is a simple and almost pathological example given that it doesn't yet have any ports attached to it. We'll look at the details a bit later:

```

$ ovn-sbctl lflow-list n1 | abbrev
Datapath: "neutron-5b6baf" aka "n1" (a8a758) Pipeline: ingress
  table=0 (ls_in_port_sec_l2 ), priority=100 , match=(eth.src[40]), action=(drop;)
  table=0 (ls_in_port_sec_l2 ), priority=100 , match=(vlan.present), action=(drop;)
...
Datapath: "neutron-5b6baf" aka "n1" (a8a758) Pipeline: egress
  table=0 (ls_out_pre_lb      ), priority=0 , match=(1), action=(next;)
  table=1 (ls_out_pre_acl     ), priority=0 , match=(1), action=(next;)
...

```

We have one hypervisor (aka “compute node”, in OpenStack parlance), which is the one where we're running all these commands. On this hypervisor, `ovn-controller` is translating OVN logical flows into OpenFlow flows (“physical flows”). It makes sense to go deeper, to see the OpenFlow flows that get generated from this datapath. By adding `--ovs` to the `ovn-sbctl` command, we can see OpenFlow flows listed just below their logical flows. We also need to use `sudo` because connecting to Open vSwitch is privileged. Go ahead and try it:

```

$ sudo ovn-sbctl --ovs lflow-list n1 | abbrev
Datapath: "neutron-5b6baf" aka "n1" (a8a758) Pipeline: ingress
  table=0 (ls_in_port_sec_l2 ), priority=100 , match=(eth.src[40]), action=(drop;)
  table=0 (ls_in_port_sec_l2 ), priority=100 , match=(vlan.present), action=(drop;)
...
Datapath: "neutron-5b6baf" aka "n1" (a8a758) Pipeline: egress
  table=0 (ls_out_pre_lb      ), priority=0 , match=(1), action=(next;)
  table=1 (ls_out_pre_acl     ), priority=0 , match=(1), action=(next;)
...

```

You were probably disappointed: the output didn't change, and no OpenFlow flows were printed. That's because no OpenFlow flows are installed for this logical datapath, which in turn is because there are no VIFs for this logical datapath on the local hypervisor. For a better example, you can try `ovn-sbctl --ovs` on one of the other logical datapaths.

Attaching VMs

A switch without any ports is not very interesting. Let's create a couple of VMs and attach them to the switch. Run the following commands, which create VMs named `a` and `b` and attaches them to our network `n1` with IP addresses `10.1.1.5` and `10.1.1.6`, respectively. It is not actually necessary to manually assign IP address assignments, since OpenStack is perfectly happy to assign them itself from the subnet's IP address range, but predictable addresses are useful for our discussion:

```
$ openstack server create --nic net-id=n1,v4-fixed-ip=10.1.1.5 --flavor m1.nano --  
→image $IMAGE_ID --key-name demo a  
$ openstack server create --nic net-id=n1,v4-fixed-ip=10.1.1.6 --flavor m1.nano --  
→image $IMAGE_ID --key-name demo b
```

These commands return before the VMs are really finished being built. You can run `openstack server list` a few times until each of them is shown in the state `ACTIVE`, which means that they're not just built but already running on the local hypervisor.

These operations had the side effect of creating separate “port” objects, but without giving those ports any easy-to-read names. It'll be easier to deal with them later if we can refer to them by name, so let's name a's port `ap` and b's port `bp`:

```
$ openstack port set --name ap $(openstack port list --server a -f value -c ID)  
$ openstack port set --name bp $(openstack port list --server b -f value -c ID)
```

We'll need to refer to these ports' MAC addresses a few times, so let's put them in variables:

```
$ AP_MAC=$(openstack port show -f value -c mac_address ap)  
$ BP_MAC=$(openstack port show -f value -c mac_address bp)
```

At this point you can log into the consoles of the VMs if you like. You can do that from the OpenStack web interface or get a direct URL to paste into a web browser using a command like:

```
$ openstack console url show -f yaml a
```

(The option `-f yaml` keeps the URL in the output from being broken into noncontiguous pieces on a 80-column console.)

The VMs don't have many tools in them but `ping` and `ssh` from one to the other should work fine. The VMs do not have any external network access or DNS configuration.

Let's chase down what's changed in OVN. Start with the NB DB at the top of the system. It's clear that our logical switch now has the two logical ports attached to it:

```
$ ovn-nbctl show | abbrev  
...  
switch 3eb263 (neutron-5b6baf) (aka n1)  
  port c29d41 (aka bp)  
    addresses: ["fa:16:3e:99:7a:17 10.1.1.6"]  
  port 820c08 (aka ap)  
    addresses: ["fa:16:3e:a9:4c:c7 10.1.1.5"]  
...
```

We can get some more details on each of these by looking at their NB DB records in the `Logical_Switch_Port` table. Each port has addressing information, port security enabled, and a pointer to DHCP configuration (which we'll look at much later in *DHCP*):

```
$ ovn-nbctl list logical_switch_port ap bp | abbrev  
_uuid          : ef17e5  
addresses      : ["fa:16:3e:a9:4c:c7 10.1.1.5"]  
dhcpv4_options : 165974  
dhcpv6_options : []  
dynamic_addresses : []  
enabled        : true  
external_ids   : {"neutron:port_name"=ap}  
name           : "820c08"
```

(continues on next page)

(continued from previous page)

```

options          : {}
parent_name      : []
port_security    : ["fa:16:3e:a9:4c:c7 10.1.1.5"]
tag              : []
tag_request      : []
type             : ""
up               : true

_uuid            : e8af12
addresses        : ["fa:16:3e:99:7a:17 10.1.1.6"]
dhcpv4_options   : 165974
dhcpv6_options   : []
dynamic_addresses : []
enabled          : true
external_ids     : {"neutron:port_name"=bp}
name             : "c29d41"
options          : {}
parent_name      : []
port_security    : ["fa:16:3e:99:7a:17 10.1.1.6"]
tag              : []
tag_request      : []
type             : ""
up               : true

```

Now that the logical switch is less pathological, it's worth taking another look at the SB DB logical flow table. Try a command like this:

```
$ ovn-sbctl lflow-list nl | abbrev | less -S
```

and then glance through the flows. Packets that egress a VM into the logical switch travel through the flow table's ingress pipeline starting from table 0. At each table, the switch finds the highest-priority logical flow that matches and executes its actions, or if there's no matching flow then the packet is dropped. The `ovn-sb(5)` manpage gives all the details, but with a little thought it's possible to guess a lot without reading the manpage. For example, consider the flows in ingress pipeline table 0, which are the first flows encountered by a packet traversing the switch:

```

table=0 (ls_in_port_sec_l2 ), priority=100 , match=(eth.src[40]), action=(drop;)
table=0 (ls_in_port_sec_l2 ), priority=100 , match=(vlan.present), action=(drop;)
table=0 (ls_in_port_sec_l2 ), priority=50   , match=(inport == "820c08" && eth.src_
↳== {fa:16:3e:a9:4c:c7}), action=(next;)
table=0 (ls_in_port_sec_l2 ), priority=50   , match=(inport == "c29d41" && eth.src_
↳== {fa:16:3e:99:7a:17}), action=(next;)

```

The first two flows, with priority 100, immediately drop two kinds of invalid packets: those with a multicast or broadcast Ethernet source address (since multicast is only for packet destinations) and those with a VLAN tag (because OVN doesn't yet support VLAN tags inside logical networks). The next two flows implement L2 port security: they advance to the next table for packets with the correct Ethernet source addresses for their ingress ports. A packet that does not match any flow is implicitly dropped, so there's no need for flows to deal with mismatches.

The logical flow table includes many other flows, some of which we will look at later. For now, it's most worth looking at ingress table 13:

```

table=13(ls_in_l2_lkup      ), priority=100 , match=(eth.mcast), action=(output = "_
↳MC_flood"; output;)
table=13(ls_in_l2_lkup      ), priority=50   , match=(eth.dst == fa:16:3e:99:7a:17),
↳action=(output = "c29d41"; output;)
table=13(ls_in_l2_lkup      ), priority=50   , match=(eth.dst == fa:16:3e:a9:4c:c7),
↳action=(output = "820c08"; output;)

```

(continues on next page)

(continued from previous page)

The first flow in table 13 checks whether the packet is an Ethernet multicast or broadcast and, if so, outputs it to a special port that egresses to every logical port (other than the ingress port). Otherwise the packet is output to the port corresponding to its Ethernet destination address. Packets addressed to any other Ethernet destination are implicitly dropped.

(It's common for an OVN logical switch to know all the MAC addresses supported by its logical ports, like this one. That's why there's no logic here for MAC learning or flooding packets to unknown MAC addresses. OVN does support unknown MAC handling but that's not in play in our example.)

Note: If you're interested in the details for the multicast group, you can run a command like the following and then look at the row for the correct datapath:

```
$ ovn-sbctl find multicast_group name=_MC_flood | abbrev
```

Now if you want to look at the OpenFlow flows, you can actually see them. For example, here's the beginning of the output that lists the first four logical flows, which we already looked at above, and their corresponding OpenFlow flows. If you want to know more about the syntax, the `ovs-fields(7)` manpage explains OpenFlow matches and `ovs-ofctl(8)` explains OpenFlow actions:

```
$ sudo ovn-sbctl --ovs lflow-list n1 | abbrev
Datapath: "neutron-5b6baf" aka "n1" (a8a758) Pipeline: ingress
  table=0 (ls_in_port_sec_l2 ), priority=100 , match=(eth.src[40]), action=(drop;)
  table=8 metadata=0x4,d1_src=01:00:00:00:00:00/01:00:00:00:00:00 actions=drop
  table=0 (ls_in_port_sec_l2 ), priority=100 , match=(vlan.present), action=(drop;)
  table=8 metadata=0x4,vlan_tci=0x1000/0x1000 actions=drop
  table=0 (ls_in_port_sec_l2 ), priority=50 , match=(inport == "820c08" && eth.src_
↳ == {fa:16:3e:a9:4c:c7}), action=(next;)
  table=8 reg14=0x1,metadata=0x4,d1_src=fa:16:3e:a9:4c:c7 actions=resubmit(,9)
  table=0 (ls_in_port_sec_l2 ), priority=50 , match=(inport == "c29d41" && eth.src_
↳ == {fa:16:3e:99:7a:17}), action=(next;)
  table=8 reg14=0x2,metadata=0x4,d1_src=fa:16:3e:99:7a:17 actions=resubmit(,9)
...
```

Logical Tracing

Let's go a level deeper. So far, everything we've done has been fairly general. We can also look at something more specific: the path that a particular packet would take through OVN, logically, and Open vSwitch, physically.

Let's use OVN's `ovn-trace` utility to see what happens to packets from a logical point of view. The `ovn-trace(8)` manpage has a lot of detail on how to do that, but let's just start by building up from a simple example. You can start with a command that just specifies the logical datapath, an input port, and nothing else; unspecified fields default to all-zeros. This doesn't do much:

```
$ ovn-trace n1 'inport == "ap"'
...
ingress(dp="n1", inport="ap")
-----
0. ls_in_port_sec_l2: no match (implicit drop)
```

We see that the packet was dropped in logical table 0, "ls_in_port_sec_l2", the L2 port security stage (as we discussed earlier). That's because we didn't use the right Ethernet source address for a. Let's see what happens if we do:

```
$ ovn-trace nl 'inport == "ap" && eth.src == '$AP_MAC
...
ingress(dp="nl", inport="ap")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "ap" && eth.src ==
↪{fa:16:3e:a9:4c:c7}, priority 50, uuid 6dcc418a
  next;
13. ls_in_l2_lkup: no match (implicit drop)
```

Now the packet passes through L2 port security and skips through several other tables until it gets dropped in the L2 lookup stage (because the destination is unknown). Let's add the Ethernet destination for b:

```
$ ovn-trace nl 'inport == "ap" && eth.src == '$AP_MAC' && eth.dst == '$BP_MAC
...
ingress(dp="nl", inport="ap")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "ap" && eth.src ==
↪{fa:16:3e:a9:4c:c7}, priority 50, uuid 6dcc418a
  next;
13. ls_in_l2_lkup (ovn-northd.c:3529): eth.dst == fa:16:3e:99:7a:17, priority 50, ↪
↪uuid 57a4c46f
  output = "bp";
  output;

egress(dp="nl", inport="ap", output="bp")
-----
8. ls_out_port_sec_l2 (ovn-northd.c:3654): output == "bp" && eth.dst ==
↪{fa:16:3e:99:7a:17}, priority 50, uuid 8aa6426d
  output;
  /* output to "bp", type "" */
```

You can see that in this case the packet gets properly switched from a to b.

Physical Tracing for Hypothetical Packets

ovn-trace showed us how a hypothetical packet would travel through the system in a logical fashion, that is, without regard to how VMs are distributed across the physical network. This is a convenient representation for understanding how OVN is **supposed** to work abstractly, but sometimes we might want to know more about how it actually works in the real systems where it is running. For this, we can use the tracing tool that Open vSwitch provides, which traces a hypothetical packet through the OpenFlow tables.

We can actually get two levels of detail. Let's start with the version that's easier to interpret, by physically tracing a packet that looks like the one we logically traced before. One obstacle is that we need to know the OpenFlow port number of the input port. One way to do that is to look for a port whose "attached-mac" is the one we expect and print its ofport number:

```
$ AP_PORT=$(ovs-vsctl --bare --columns=ofport find interface external-ids:attached-
↪mac=\ "$AP_MAC\ ")
$ echo $AP_PORT
3
```

(You could also just do a plain `ovs-vsctl list interface` and then look through for the right row and pick its ofport value.)

Now we can feed this input port number into `ovs-appctl ofproto/trace` along with the correct Ethernet source and destination addresses and get a physical trace:

```
$ sudo ovs-appctl ofproto/trace br-int in_port=$AP_PORT,d1_src=$AP_MAC,d1_dst=$BP_MAC
Flow: in_port=3,vlan_tci=0x0000,d1_src=fa:16:3e:a9:4c:c7,d1_dst=fa:16:3e:99:7a:17,d1_
→type=0x0000
```

```
bridge("br-int")
```

```
-----
0. in_port=3, priority 100
   set_field:0x8->reg13
   set_field:0x9->reg11
   set_field:0xa->reg12
   set_field:0x4->metadata
   set_field:0x1->reg14
   resubmit(,8)
8. reg14=0x1,metadata=0x4,d1_src=fa:16:3e:a9:4c:c7, priority 50, cookie 0x6dcc418a
   resubmit(,9)
9. metadata=0x4, priority 0, cookie 0x8fe8689e
   resubmit(,10)
10. metadata=0x4, priority 0, cookie 0x719549d1
   resubmit(,11)
11. metadata=0x4, priority 0, cookie 0x39c99e6f
   resubmit(,12)
12. metadata=0x4, priority 0, cookie 0x838152a3
   resubmit(,13)
13. metadata=0x4, priority 0, cookie 0x918259e3
   resubmit(,14)
14. metadata=0x4, priority 0, cookie 0xcad14db2
   resubmit(,15)
15. metadata=0x4, priority 0, cookie 0x7834d912
   resubmit(,16)
16. metadata=0x4, priority 0, cookie 0x87745210
   resubmit(,17)
17. metadata=0x4, priority 0, cookie 0x34951929
   resubmit(,18)
18. metadata=0x4, priority 0, cookie 0xd7a8c9fb
   resubmit(,19)
19. metadata=0x4, priority 0, cookie 0xd02e9578
   resubmit(,20)
20. metadata=0x4, priority 0, cookie 0x42d35507
   resubmit(,21)
21. metadata=0x4,d1_dst=fa:16:3e:99:7a:17, priority 50, cookie 0x57a4c46f
   set_field:0x2->reg15
   resubmit(,32)
32. priority 0
   resubmit(,33)
33. reg15=0x2,metadata=0x4, priority 100
   set_field:0xb->reg13
   set_field:0x9->reg11
   set_field:0xa->reg12
   resubmit(,34)
34. priority 0
   set_field:0->reg0
   set_field:0->reg1
   set_field:0->reg2
   set_field:0->reg3
   set_field:0->reg4
   set_field:0->reg5
   set_field:0->reg6
```

(continues on next page)

(continued from previous page)

```

    set_field:0->reg7
    set_field:0->reg8
    set_field:0->reg9
    resubmit(,40)
40. metadata=0x4, priority 0, cookie 0xde9f3899
    resubmit(,41)
41. metadata=0x4, priority 0, cookie 0x74074eff
    resubmit(,42)
42. metadata=0x4, priority 0, cookie 0x7789c8b1
    resubmit(,43)
43. metadata=0x4, priority 0, cookie 0xa6b002c0
    resubmit(,44)
44. metadata=0x4, priority 0, cookie 0xaeab2b45
    resubmit(,45)
45. metadata=0x4, priority 0, cookie 0x290cc4d4
    resubmit(,46)
46. metadata=0x4, priority 0, cookie 0xa3223b88
    resubmit(,47)
47. metadata=0x4, priority 0, cookie 0x7ac2132e
    resubmit(,48)
48. reg15=0x2,metadata=0x4,d1_dst=fa:16:3e:99:7a:17, priority 50, cookie 0x8aa6426d
    resubmit(,64)
64. priority 0
    resubmit(,65)
65. reg15=0x2,metadata=0x4, priority 100
    output:4

```

```

Final flow: reg11=0x9,reg12=0xa,reg13=0xb,reg14=0x1,reg15=0x2,metadata=0x4,in_port=3,
↳vlan_tci=0x0000,d1_src=fa:16:3e:a9:4c:c7,d1_dst=fa:16:3e:99:7a:17,d1_type=0x0000
MegafLOW: recirc_id=0,ct_state=new-est-rel-rpl-inv-trk,ct_label=0/0x1,in_port=3,vlan_
↳tci=0x0000/0x1000,d1_src=fa:16:3e:a9:4c:c7,d1_dst=fa:16:3e:99:7a:17,d1_type=0x0000
Datapath actions: 4

```

There's a lot there, which you can read through if you like, but the important part is:

```

65. reg15=0x2,metadata=0x4, priority 100
    output:4

```

which means that the packet is ultimately being output to OpenFlow port 4. That's port b, which you can confirm with:

```

$ sudo ovs-vsctl find interface ofport=4
_uuid                : 840a5aca-ea8d-4c16-a11b-a94e0f408091
admin_state          : up
bfd                  : {}
bfd_status           : {}
cfm_fault            : []
cfm_fault_status     : []
cfm_flap_count       : []
cfm_health           : []
cfm_mpid             : []
cfm_remote_mpid      : []
cfm_remote_opstate   : []
duplex               : full
error                : []
external_ids         : {attached-mac="fa:16:3e:99:7a:17", iface-id="c29d4120-20a4-4c44-
↳bd83-8d91f5f447fd", iface-status=active, vm-id="2db969ca-ca2a-4d9a-b49e-f287d39c5645
↳"}

```

(continues on next page)

(continued from previous page)

```

ifindex          : 9
ingress_policing_burst: 0
ingress_policing_rate: 0
lacp_current     : []
link_resets      : 1
link_speed       : 10000000
link_state       : up
lldp             : {}
mac              : []
mac_in_use       : "fe:16:3e:99:7a:17"
mtu              : 1500
mtu_request      : []
name             : "tapc29d4120-20"
ofport           : 4
ofport_request   : []
options          : {}
other_config     : {}
statistics       : {collisions=0, rx_bytes=4254, rx_crc_err=0, rx_dropped=0, rx_
↳ errors=0, rx_frame_err=0, rx_over_err=0, rx_packets=39, tx_bytes=4188, tx_dropped=0,
↳ tx_errors=0, tx_packets=39}
status           : {driver_name=tun, driver_version="1.6", firmware_version=""}
type             : ""

```

or:

```

$ BP_PORT=$(ovs-vsctl --bare --columns=ofport find interface external-ids:attached-
↳ mac=\ "$BP_MAC\ ")
$ echo $BP_PORT
4

```

Physical Tracing for Real Packets

In the previous sections we traced a hypothetical L2 packet, one that's honestly not very realistic: we didn't even supply an Ethernet type, so it defaulted to zero, which isn't anything one would see on a real network. We could refine our packet so that it becomes a more realistic TCP or UDP or ICMP, etc. packet, but let's try a different approach: working from a real packet.

Pull up a console for VM a and start `ping 10.1.1.6`, then leave it running for the rest of our experiment.

Now go back to your DevStack session and run:

```
$ sudo watch ovs-dpctl dump-flows
```

We're working with a new program. `ovn-dpctl` is an interface to Open vSwitch datapaths, in this case to the Linux kernel datapath. Its `dump-flows` command displays the contents of the in-kernel flow cache, and by running it under the `watch` program we see a new snapshot of the flow table every 2 seconds.

Look through the output for a flow that begins with `recirc_id(0)` and matches the Ethernet source address for a. There is one flow per line, but the lines are very long, so it's easier to read if you make the window very wide. This flow's packet counter should be increasing at a rate of 1 packet per second. It looks something like this:

```

recirc_id(0),in_port(3),eth(src=fa:16:3e:f5:2a:90),eth_type(0x0800),ipv4(src=10.1.1.5,
↳ frag=no), packets:388, bytes:38024, used:0.977s, actions:ct(zone=8),recirc(0x18)

```

We can hand the first part of this (everything up to the first space) to `ofproto/trace`, and it will tell us what happens:


```
$ sudo ovs-appctl ofproto/trace 'recirc_id(0),in_port(3),eth(src=fa:16:3e:a9:4c:c7),
↪eth_type(0x0800),ipv4(src=10.1.1.5,dst=10.1.0.0/255.255.0.0,frag=no)'
Flow: ip,in_port=3,vlan_tci=0x0000,dl_src=fa:16:3e:a9:4c:c7,dl_dst=00:00:00:00:00:00,
↪nw_src=10.1.1.5,nw_dst=10.1.0.0,nw_proto=0,nw_tos=0,nw_ecn=0,nw_ttl=0

bridge("br-int")
-----
0. in_port=3, priority 100
   set_field:0x8->reg13
   set_field:0x9->reg11
   set_field:0xa->reg12
   set_field:0x4->metadata
   set_field:0x1->reg14
   resubmit(,8)
8. reg14=0x1,metadata=0x4,dl_src=fa:16:3e:a9:4c:c7, priority 50, cookie 0x6dcc418a
   resubmit(,9)
9. ip,reg14=0x1,metadata=0x4,dl_src=fa:16:3e:a9:4c:c7,nw_src=10.1.1.5, priority 90,
↪cookie 0x343af48c
   resubmit(,10)
10. metadata=0x4, priority 0, cookie 0x719549d1
    resubmit(,11)
11. ip,metadata=0x4, priority 100, cookie 0x46c089e6
    load:0x1->NXM_NX_XXREG0[96]
    resubmit(,12)
12. metadata=0x4, priority 0, cookie 0x838152a3
    resubmit(,13)
13. ip,reg0=0x1/0x1,metadata=0x4, priority 100, cookie 0xd1941634
    ct(table=22,zone=NXM_NX_REG13[0..15])
    drop

Final flow: ip,reg0=0x1,reg11=0x9,reg12=0xa,reg13=0x8,reg14=0x1,metadata=0x4,in_
↪port=3,vlan_tci=0x0000,dl_src=fa:16:3e:a9:4c:c7,dl_dst=00:00:00:00:00:00,nw_src=10.
↪1.1.5,nw_dst=10.1.0.0,nw_proto=0,nw_tos=0,nw_ecn=0,nw_ttl=0
Megaflow: recirc_id=0,ip,in_port=3,vlan_tci=0x0000/0x1000,dl_src=fa:16:3e:a9:4c:c7,nw_
↪src=10.1.1.5,nw_dst=10.1.0.0/16,nw_frag=no
Datapath actions: ct(zone=8),recirc(0xb)
```

Note: Be careful cutting and pasting `ovs-dpctl dump-flows` output into `ofproto/trace` because the latter has terrible error reporting. If you add an extra line break, etc., it will likely give you a useless error message.

There's no output action in the output, but there are `ct` and `recirc` actions (which you can see in the Datapath actions at the end). The `ct` action tells the kernel to pass the packet through the kernel connection tracking for fire-walling purposes and the `recirc` says to go back to the flow cache for another pass based on the firewall results. The `0xb` value inside the `recirc` gives us a hint to look at the kernel flows for a cached flow with `recirc_id(0xb)`. Indeed, there is one:

```
recirc_id(0xb),in_port(3),ct_state(-new+est-rel-rpl-inv+trk),ct_label(0/0x1),
↪eth(src=fa:16:3e:a9:4c:c7,dst=fa:16:3e:99:7a:17),eth_type(0x0800),ipv4(dst=10.1.1.4/
↪255.255.255.252,frag=no), packets:171, bytes:16758, used:0.271s,
↪actions:ct(zone=11),recirc(0xc)
```

We can then repeat our command with the match part of this kernel flow:

```
$ sudo ovs-appctl ofproto/trace 'recirc_id(0xb),in_port(3),ct_state(-new+est-rel-rpl-
↪inv+trk),ct_label(0/0x1),eth(src=fa:16:3e:a9:4c:c7,dst=fa:16:3e:99:7a:17),eth_
↪type(0x0800),ipv4(dst=10.1.1.4/255.255.255.252,frag=no)'
```

(continues on next page)

(continued from previous page)

```
...
Datapath actions: ct(zone=11),recirc(0xc)
```

In other words, the flow passes through the connection tracker a second time. The first time was for a's outgoing firewall; this second time is for b's incoming firewall. Again, we continue tracing with `recirc_id(0xc)`:

```
$ sudo ovs-appctl ofproto/trace 'recirc_id(0xc),in_port(3),ct_state(-new+est-rel-rpl-
→inv+trk),ct_label(0/0x1),eth(src=fa:16:3e:a9:4c:c7,dst=fa:16:3e:99:7a:17),eth_
→type(0x0800),ipv4(dst=10.1.1.6,proto=1,frag=no) '
...
Datapath actions: 4
```

It took multiple hops, but we finally came to the end of the line where the packet was output to b after passing through both firewalls. The port number here is a datapath port number, which is usually different from an OpenFlow port number. To check that it is b's port, we first list the datapath ports to get the name corresponding to the port number:

```
$ sudo ovs-dpctl show
system@ovs-system:
  lookups: hit:1994 missed:56 lost:0
  flows: 6
  masks: hit:2340 total:4 hit/pkt:1.14
  port 0: ovs-system (internal)
  port 1: br-int (internal)
  port 2: br-ex (internal)
  port 3: tap820c0888-13
  port 4: tapc29d4120-20
```

and then confirm that this is the port we think it is with a command like this:

```
$ ovs-vsctl --columns=external-ids list interface tapc29d4120-20
external_ids      : {attached-mac="fa:16:3e:99:7a:17", iface-id="c29d4120-20a4-4c44-
→bd83-8d91f5f447fd", iface-status=active, vm-id="2db969ca-ca2a-4d9a-b49e-f287d39c5645
→"} }
```

Finally, we can relate the OpenFlow flows from our traces back to OVN logical flows. For individual flows, cut and paste a “cookie” value from `ofproto/trace` output into `ovn-sbctl lflow-list`, e.g.:

```
$ ovn-sbctl lflow-list 0x6dcc418a|abbrev
Datapath: "neutron-5b6baf" aka "n1" (a8a758) Pipeline: ingress
  table=0 (ls_in_port_sec_l2 ), priority=50 , match=(inport == "820c08" && eth.src_
→== {fa:16:3e:a9:4c:c7}), action=(next;)
```

Or, you can pipe `ofproto/trace` output through `ovn-detrace` to annotate every flow:

```
$ sudo ovs-appctl ofproto/trace 'recirc_id(0xc),in_port(3),ct_state(-new+est-rel-rpl-
→inv+trk),ct_label(0/0x1),eth(src=fa:16:3e:a9:4c:c7,dst=fa:16:3e:99:7a:17),eth_
→type(0x0800),ipv4(dst=10.1.1.6,proto=1,frag=no) ' | ovn-detrace
...
```

3.5.6 Routing

Previously we set up a pair of VMs a and b on a network n1 and demonstrated how packets make their way between them. In this step, we'll set up a second network n2 with a new VM c, connect a router r to both networks, and demonstrate how routing works in OVN.

There's nothing really new for the network and the VM so let's just go ahead and create them:

```
$ openstack network create --project admin --provider-network-type geneve n2
$ openstack subnet create --subnet-range 10.1.2.0/24 --network n2 n2subnet
$ openstack server create --nic net-id=n2,v4-fixed-ip=10.1.2.7 --flavor m1.nano --
→ image $IMAGE_ID --key-name demo c
$ openstack port set --name cp $(openstack port list --server c -f value -c ID)
$ CP_MAC=$(openstack port show -f value -c mac_address cp)
```

The new network n2 is not yet connected to n1 in any way. You can try tracing a broadcast packet from a to see, for example, that it doesn't make it to c:

```
$ ovn-trace n1 'inport == "ap" && eth.src == '$AP_MAC' && eth.dst == '$CP_MAC'
...
```

Now create an OpenStack router and connect it to n1 and n2:

```
$ openstack router create r
$ openstack router add subnet r n1subnet
$ openstack router add subnet r n2subnet
```

Now a, b, and c should all be able to reach other. You can get some verification that routing is taking place by running you ping between c and one of the other VMs: the reported TTL should be one less than between a and b (63 instead of 64).

Observe via `ovn-nbctl` the new OVN logical switch and router and then ports that connect them together:

```
$ ovn-nbctl show|abbrev
...
switch f51234 (neutron-332346) (aka n2)
  port 82b983
    type: router
    router-port: lrp-82b983
  port 2e585f (aka cp)
    addresses: ["fa:16:3e:89:f2:36 10.1.2.7"]
switch 3eb263 (neutron-5b6baf) (aka n1)
  port c29d41 (aka bp)
    addresses: ["fa:16:3e:99:7a:17 10.1.1.6"]
  port 820c08 (aka ap)
    addresses: ["fa:16:3e:a9:4c:c7 10.1.1.5"]
  port 17d870
    type: router
    router-port: lrp-17d870
...
router dde06c (neutron-f88ebc) (aka r)
  port lrp-82b983
    mac: "fa:16:3e:19:9f:46"
    networks: ["10.1.2.1/24"]
  port lrp-17d870
    mac: "fa:16:3e:f6:e2:8f"
    networks: ["10.1.1.1/24"]
```

We have not yet looked at the logical flows for an OVN logical router. You might find it of interest to look at them on your own:

```
$ ovn-sbctl lflow-list r | abbrev | less -S
...
```

Let's grab the n1subnet router porter MAC address to simplify later commands:

```
$ N1SUBNET_MAC=$(ovn-nbctl --bare --columns=mac find logical_router_port networks=10.
↪1.1.1/24)
```

Let's see what happens at the logical flow level for an ICMP packet from `a` to `c`. This generates a long trace but an interesting one, so we'll look at it bit by bit. The first three stanzas in the output show the packet's ingress into `n1` and processing through the firewall on that side (via the "ct_next" connection-tracking action), and then the selection of the port that leads to router `r` as the output port:

```
$ ovn-trace n1 'inport == "ap" && eth.src == '$AP_MAC' && eth.dst == '$N1SUBNET_MAC' &
↪ & ip4.src == 10.1.1.5 && ip4.dst == 10.1.2.7 && ip.ttl == 64 && icmp4.type == 8'
...
ingress(dp="n1", inport="ap")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "ap" && eth.src ==
↪ {fa:16:3e:a9:4c:c7}, priority 50, uuid 6dcc418a
  next;
1. ls_in_port_sec_ip (ovn-northd.c:2364): inport == "ap" && eth.src ==
↪ fa:16:3e:a9:4c:c7 && ip4.src == {10.1.1.5}, priority 90, uuid 343af48c
  next;
3. ls_in_pre_acl (ovn-northd.c:2646): ip, priority 100, uuid 46c089e6
  reg0[0] = 1;
  next;
5. ls_in_pre_stateful (ovn-northd.c:2764): reg0[0] == 1, priority 100, uuid d1941634
  ct_next;

ct_next(ct_state=est|trk /* default (use --ct to customize) */)
-----
6. ls_in_acl (ovn-northd.c:2925): !ct.new && ct.est && !ct.rpl && ct_label.blocked_
↪ == 0 && (inport == "ap" && ip4), priority 2002, uuid a12b39f0
  next;
13. ls_in_l2_lkup (ovn-northd.c:3529): eth.dst == fa:16:3e:f6:e2:8f, priority 50,
↪ uuid c43ead31
  output = "17d870";
  output;

egress(dp="n1", inport="ap", output="17d870")
-----
1. ls_out_pre_acl (ovn-northd.c:2626): ip && output == "17d870", priority 110, uuid_
↪ 60395450
  next;
8. ls_out_port_sec_l2 (ovn-northd.c:3654): output == "17d870", priority 50, uuid_
↪ 91b5cab0
  output;
  /* output to "17d870", type "patch" */
```

The next two stanzas represent processing through logical router `r`. The processing in table 5 is the core of the routing implementation: it recognizes that the packet is destined for an attached subnet, decrements the TTL and updates the Ethernet source address. Table 6 then selects the Ethernet destination address based on the IP destination. The packet then passes to switch `n2` via an OVN "logical patch port":

```
ingress(dp="r", inport="lrp-17d870")
-----
0. lr_in_admission (ovn-northd.c:4071): eth.dst == fa:16:3e:f6:e2:8f && inport ==
↪ "lrp-17d870", priority 50, uuid fa5270b0
  next;
5. lr_in_ip_routing (ovn-northd.c:3782): ip4.dst == 10.1.2.0/24, priority 49, uuid_
↪ 5f9d469f
```

(continues on next page)

(continued from previous page)

```

ip.ttl--;
reg0 = ip4.dst;
reg1 = 10.1.2.1;
eth.src = fa:16:3e:19:9f:46;
output = "lrp-82b983";
flags.loopback = 1;
next;
6. lr_in_arp_resolve (ovn-northd.c:5088): output == "lrp-82b983" && reg0 == 10.1.2.
→7, priority 100, uuid 03d506d3
   eth.dst = fa:16:3e:89:f2:36;
   next;
8. lr_in_arp_request (ovn-northd.c:5260): 1, priority 0, uuid 6dacdd82
   output;

egress(dp="r", inport="lrp-17d870", output="lrp-82b983")
-----
3. lr_out_delivery (ovn-northd.c:5288): output == "lrp-82b983", priority 100, uuid
→00bea4f2
   output;
/* output to "lrp-82b983", type "patch" */

```

Finally the logical switch for n2 runs through the same logic as n1 and the packet is delivered to VM c:

```

ingress(dp="n2", inport="82b983")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "82b983", priority 50, uuid
→9a789e06
   next;
3. ls_in_pre_acl (ovn-northd.c:2624): ip && inport == "82b983", priority 110, uuid
→ab52f21a
   next;
13. ls_in_l2_lkup (ovn-northd.c:3529): eth.dst == fa:16:3e:89:f2:36, priority 50,
→uuid dcafb3e9
   output = "cp";
   output;

egress(dp="n2", inport="82b983", output="cp")
-----
1. ls_out_pre_acl (ovn-northd.c:2648): ip, priority 100, uuid cd9cfa74
   reg0[0] = 1;
   next;
2. ls_out_pre_stateful (ovn-northd.c:2766): reg0[0] == 1, priority 100, uuid 9e8e22c5
   ct_next;

ct_next(ct_state=est|trk /* default (use --ct to customize) */)
-----
4. ls_out_acl (ovn-northd.c:2925): !ct.new && ct.est && !ct.rpl && ct_label.blocked
→== 0 && (output == "cp" && ip4 && ip4.src == $as_ip4_0fc1b6cf_f925_49e6_8f00_
→6dd13beca9dc), priority 2002, uuid a746fa0d
   next;
7. ls_out_port_sec_ip (ovn-northd.c:2364): output == "cp" && eth.dst ==
→fa:16:3e:89:f2:36 && ip4.dst == {255.255.255.255, 224.0.0.0/4, 10.1.2.7}, priority
→90, uuid 4d9862b5
   next;
8. ls_out_port_sec_l2 (ovn-northd.c:3654): output == "cp" && eth.dst ==
→{fa:16:3e:89:f2:36}, priority 50, uuid 0242cdc3
   output;

```

(continues on next page)

(continued from previous page)

```
/* output to "cp", type " " */
```

Physical Tracing

It's possible to use `ofproto/trace`, just as before, to trace a packet through OpenFlow tables, either for a hypothetical packet or one that you get from a real test case using `ovs-dpctl`. The process is just the same as before and the output is almost the same, too. Using a router doesn't actually introduce any interesting new wrinkles, so we'll skip over this for this case and for the remainder of the tutorial, but you can follow the steps on your own if you like.

3.5.7 Adding a Gateway

The VMs that we've created can access each other but they are isolated from the physical world. In OpenStack, the dominant way to connect a VM to external networks is by creating what is called a "floating IP address", which uses network address translation to connect an external address to an internal one.

DevStack created a pair of networks named "private" and "public". To use a floating IP address from a VM, we first add a port to the VM with an IP address from the "private" network, then we create a floating IP address on the "public" network, then we associate the port with the floating IP address.

Let's add a new VM `d` with a floating IP:

```
$ openstack server create --nic net-id=private --flavor m1.nano --image $IMAGE_ID --
↪key-name demo d
$ openstack port set --name dp $(openstack port list --server d -f value -c ID)
$ DP_MAC=$(openstack port show -f value -c mac_address dp)
$ openstack floating ip create --floating-ip-address 172.24.4.8 public
$ openstack server add floating ip d 172.24.4.8
```

(We specified a particular floating IP address to make the examples easier to follow, but without that OpenStack will automatically allocate one.)

It's also necessary to configure the "public" network because DevStack does not do it automatically:

```
$ sudo ip link set br-ex up
$ sudo ip route add 172.24.4.0/24 dev br-ex
$ sudo ip addr add 172.24.4.1/24 dev br-ex
```

Now you should be able to "ping" VM `d` from the OpenStack host:

```
$ ping 172.24.4.8
PING 172.24.4.8 (172.24.4.8) 56(84) bytes of data.
64 bytes from 172.24.4.8: icmp_seq=1 ttl=63 time=56.0 ms
64 bytes from 172.24.4.8: icmp_seq=2 ttl=63 time=1.44 ms
64 bytes from 172.24.4.8: icmp_seq=3 ttl=63 time=1.04 ms
64 bytes from 172.24.4.8: icmp_seq=4 ttl=63 time=0.403 ms
^C
--- 172.24.4.8 ping statistics ---
4 packets transmitted, 4 received, 0% packet loss, time 3003ms
rtt min/avg/max/mdev = 0.403/14.731/56.028/23.845 ms
```

You can also SSH in with the key that we created during setup:

```
$ ssh -i ~/id_rsa_demo cirros@172.24.4.8
```

Let's dive in and see how this gets implemented in OVN. First, the relevant parts of the NB DB for the “public” and “private” networks and the router between them:

```
$ ovn-nbctl show | abbrev
switch 2579f4 (neutron-d1ac28) (aka public)
  port provnet-d1ac28
    type: localnet
    addresses: ["unknown"]
  port ae9b52
    type: router
    router-port: lrp-ae9b52
switch 5b3d5f (neutron-c02c4d) (aka private)
  port b256dd
    type: router
    router-port: lrp-b256dd
  port f264e7
    type: router
    router-port: lrp-f264e7
  port cae25b (aka dp)
    addresses: ["fa:16:3e:c1:f5:a2 10.0.0.6 fdb0:5860:4ba8:0:f816:3eff:fec1:f5a2"]
...
router c59ad2 (neutron-9b057f) (aka router1)
  port lrp-ae9b52
    mac: "fa:16:3e:b2:d2:67"
    networks: ["172.24.4.9/24", "2001:db8::b/64"]
  port lrp-b256dd
    mac: "fa:16:3e:35:33:db"
    networks: ["fdb0:5860:4ba8::1/64"]
  port lrp-f264e7
    mac: "fa:16:3e:fc:c8:da"
    networks: ["10.0.0.1/26"]
  nat 788c6d
    external ip: "172.24.4.8"
    logical ip: "10.0.0.6"
    type: "dnat_and_snat"
  nat 80914c
    external ip: "172.24.4.9"
    logical ip: "10.0.0.0/26"
    type: "snat"
...
```

What we see is:

- VM `d` is on the “private” switch under its private IP address 10.0.0.8. The “private” switch is connected to “router1” via two router ports (one for IPv4, one for IPv6).
- The “public” switch is connected to “router1” and to the physical network via a “localnet” port.
- “router1” is in the middle between “private” and “public”. In addition to the router ports that connect to these switches, it has “nat” entries that direct network address translation. The translation between floating IP address 172.24.4.8 and private address 10.0.0.8 makes perfect sense.

When the NB DB gets translated into logical flows at the southbound layer, the “nat” entries get translated into IP matches that then invoke “ct_snat” and “ct_dnat” actions. The details are intricate, but you can get some of the idea by just looking for relevant flows:

```
$ ovn-sbctl lflow-list router1 | abbrev | grep nat | grep -E '172.24.4.8|10.0.0.8'
table=3 (lr_in_unsnat      ), priority=100 , match=(ip && ip4.dst == 172.24.4.8 &&
→ inport == "lrp-ae9b52" && is_chassis_resident("cr-lrp-ae9b52")), action=(ct_snat;)
```

(continues on next page)

(continued from previous page)

```

    table=3 (lr_in_unsnat      ), priority=50    , match=(ip && ip4.dst == 172.24.4.8),
↪action=(reg9[0] = 1; next;)
    table=4 (lr_in_dnat       ), priority=100   , match=(ip && ip4.dst == 172.24.4.8 &&
↪inport == "lrp-ae9b52" && is_chassis_resident("cr-lrp-ae9b52")), action=(ct_
↪dnat(10.0.0.6);)
    table=4 (lr_in_dnat       ), priority=50    , match=(ip && ip4.dst == 172.24.4.8),
↪action=(reg9[0] = 1; next;)
    table=1 (lr_out_snat      ), priority=33    , match=(ip && ip4.src == 10.0.0.6 &&
↪outport == "lrp-ae9b52" && is_chassis_resident("cr-lrp-ae9b52")), action=(ct_
↪snat(172.24.4.8);)

```

Let's take a look at how a packet passes through this whole gauntlet. The first two stanzas just show the packet traveling through the “public” network and being forwarded to the “router1” network:

```

$ ovn-trace public 'inport == "provnet-dlac2896-18a7-4bca-8f46-b21e2370e5b1" && eth.
↪src == 00:01:02:03:04:05 && eth.dst == fa:16:3e:b2:d2:67 && ip4.src == 172.24.4.1 &&
↪ip4.dst == 172.24.4.8 && ip.ttl == 64 && icmp4.type==8'
...
ingress(dp="public", inport="provnet-dlac28")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "provnet-dlac28", priority 50,
↪uuid 8d86fb06
  next;
10. ls_in_arp_rsp (ovn-northd.c:3266): inport == "provnet-dlac28", priority 100, uuid
↪21313eff
  next;
13. ls_in_l2_lkup (ovn-northd.c:3571): eth.dst == fa:16:3e:b2:d2:67 && is_chassis_
↪resident("cr-lrp-ae9b52"), priority 50, uuid 7f28f51f
  outport = "ae9b52";
  output;

egress(dp="public", inport="provnet-dlac28", outport="ae9b52")
-----
8. ls_out_port_sec_l2 (ovn-northd.c:3654): outport == "ae9b52", priority 50, uuid
↪72fea396
  output;
  /* output to "ae9b52", type "patch" */

```

In “router1”, first the `ct_snat` action without an argument attempts to “un-SNAT” the packet. `ovn-trace` treats this as a no-op, because it doesn't have any state for tracking connections. As an alternative, it invokes `ct_dnat(10.0.0.8)` to NAT the destination IP:

```

ingress(dp="router1", inport="lrp-ae9b52")
-----
0. lr_in_admission (ovn-northd.c:4071): eth.dst == fa:16:3e:b2:d2:67 && inport ==
↪"lrp-ae9b52" && is_chassis_resident("cr-lrp-ae9b52"), priority 50, uuid 8c6945c2
  next;
3. lr_in_unsnat (ovn-northd.c:4591): ip && ip4.dst == 172.24.4.8 && inport == "lrp-
↪ae9b52" && is_chassis_resident("cr-lrp-ae9b52"), priority 100, uuid e922f541
  ct_snat;

ct_snat /* assuming no un-snat entry, so no change */
-----
4. lr_in_dnat (ovn-northd.c:4649): ip && ip4.dst == 172.24.4.8 && inport == "lrp-
↪ae9b52" && is_chassis_resident("cr-lrp-ae9b52"), priority 100, uuid 02f41b79
  ct_dnat(10.0.0.6);

```


Still in “router1”, the routing and output steps transmit the packet to the “private” network:

```
ct_dnat (ip4.dst=10.0.0.6)
-----
5. lr_in_ip_routing (ovn-northd.c:3782): ip4.dst == 10.0.0.0/26, priority 53, uuid_
↪86e005b0
   ip.ttl--;
   reg0 = ip4.dst;
   reg1 = 10.0.0.1;
   eth.src = fa:16:3e:fc:c8:da;
   output = "lrp-f264e7";
   flags.loopback = 1;
   next;
6. lr_in_arp_resolve (ovn-northd.c:5088): output == "lrp-f264e7" && reg0 == 10.0.0.
↪6, priority 100, uuid 2963d67c
   eth.dst = fa:16:3e:c1:f5:a2;
   next;
8. lr_in_arp_request (ovn-northd.c:5260): 1, priority 0, uuid eea419b7
   output;

egress(dp="router1", inport="lrp-ae9b52", outport="lrp-f264e7")
-----
3. lr_out_delivery (ovn-northd.c:5288): output == "lrp-f264e7", priority 100, uuid_
↪42dadcd23
   output;
   /* output to "lrp-f264e7", type "patch" */
```

In the “private” network, the packet passes through VM d’s firewall and is output to d:

```
ingress(dp="private", inport="f264e7")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "f264e7", priority 50, uuid_
↪5b721214
   next;
3. ls_in_pre_acl (ovn-northd.c:2624): ip && inport == "f264e7", priority 110, uuid_
↪5bdc3209
   next;
13. ls_in_l2_lkup (ovn-northd.c:3529): eth.dst == fa:16:3e:c1:f5:a2, priority 50, _
↪uuid 7957f80f
   output = "dp";
   output;

egress(dp="private", inport="f264e7", outport="dp")
-----
1. ls_out_pre_acl (ovn-northd.c:2648): ip, priority 100, uuid 4981c79d
   reg0[0] = 1;
   next;
2. ls_out_pre_stateful (ovn-northd.c:2766): reg0[0] == 1, priority 100, uuid 247e02eb
   ct_next;

ct_next(ct_state=est|trk /* default (use --ct to customize) */)
-----
4. ls_out_acl (ovn-northd.c:2925): !ct.new && ct.est && !ct.rpl && ct_label.blocked_
↪== 0 && (outport == "dp" && ip4 && ip4.src == 0.0.0.0/0 && icmp4), priority 2002, _
↪uuid b860fc9f
   next;
7. ls_out_port_sec_ip (ovn-northd.c:2364): outport == "dp" && eth.dst == _
↪fa:16:3e:c1:f5:a2 && ip4.dst == {255.255.255.255, 224.0.0.0/4, 10.0.0.6}, priority_
↪90, uuid 15655a98
```

(continues on next page)

(continued from previous page)

```

next;
8. ls_out_port_sec_l2 (ovn-northd.c:3654): outport == "dp" && eth.dst ==
↪{fa:16:3e:c1:f5:a2}, priority 50, uuid 5916f94b
output;
/* output to "dp", type " " */

```

3.5.8 IPv6

OVN supports IPv6 logical routing. Let's try it out.

The first step is to add an IPv6 subnet to networks `n1` and `n2`, then attach those subnets to our router `r`. As usual, though OpenStack can assign addresses itself, we use fixed ones to make the discussion easier:

```

$ openstack subnet create --ip-version 6 --subnet-range fc11::/64 --network n1_
↪n1subnet6
$ openstack subnet create --ip-version 6 --subnet-range fc22::/64 --network n2_
↪n2subnet6
$ openstack router add subnet r n1subnet6
$ openstack router add subnet r n2subnet6

```

Then we add an IPv6 address to each of our VMs:

```

$ A_PORT_ID=$(openstack port list --server a -f value -c ID)
$ openstack port set --fixed-ip subnet=n1subnet6,ip-address=fc11::5 $A_PORT_ID
$ B_PORT_ID=$(openstack port list --server b -f value -c ID)
$ openstack port set --fixed-ip subnet=n1subnet6,ip-address=fc11::6 $B_PORT_ID
$ C_PORT_ID=$(openstack port list --server c -f value -c ID)
$ openstack port set --fixed-ip subnet=n2subnet6,ip-address=fc22::7 $C_PORT_ID

```

At least for me, the new IPv6 addresses didn't automatically get propagated into the VMs. To do it by hand, pull up the console for `a` and run:

```

$ sudo ip addr add fc11::5/64 dev eth0
$ sudo ip route add via fc11::1

```

Then in `b`:

```

$ sudo ip addr add fc11::6/64 dev eth0
$ sudo ip route add via fc11::1

```

Finally in `c`:

```

$ sudo ip addr add fc22::7/64 dev eth0
$ sudo ip route add via fc22::1

```

Now you should have working IPv6 routing through router `r`. The relevant parts of the NB DB look like the following. The interesting parts are the new `fc11::` and `fc22::` addresses on the ports in `n1` and `n2` and the new IPv6 router ports in `r`:

```

$ ovn-nbctl show | abbrev
...
switch f51234 (neutron-332346) (aka n2)
  port 1a8162
    type: router
    router-port: lrp-1a8162

```

(continues on next page)

(continued from previous page)

```

port 82b983
    type: router
    router-port: lrp-82b983
port 2e585f (aka cp)
    addresses: ["fa:16:3e:89:f2:36 10.1.2.7 fc22::7"]
switch 3eb263 (neutron-5b6baf) (aka n1)
    port ad952e
        type: router
        router-port: lrp-ad952e
    port c29d41 (aka bp)
        addresses: ["fa:16:3e:99:7a:17 10.1.1.6 fc11::6"]
    port 820c08 (aka ap)
        addresses: ["fa:16:3e:a9:4c:c7 10.1.1.5 fc11::5"]
    port 17d870
        type: router
        router-port: lrp-17d870
...
router dde06c (neutron-f88ebc) (aka r)
    port lrp-1a8162
        mac: "fa:16:3e:06:de:ad"
        networks: ["fc22::1/64"]
    port lrp-82b983
        mac: "fa:16:3e:19:9f:46"
        networks: ["10.1.2.1/24"]
    port lrp-ad952e
        mac: "fa:16:3e:ef:2f:8b"
        networks: ["fc11::1/64"]
    port lrp-17d870
        mac: "fa:16:3e:f6:e2:8f"
        networks: ["10.1.1.1/24"]

```

Try tracing a packet from a to c. The results correspond closely to those for IPv4 which we already discussed back under [Routing](#):

```

$ N1SUBNET6_MAC=$(ovn-nbctl --bare --columns=mac find logical_router_port networks=\
↳ "fc11::1/64")
$ ovn-trace n1 'inport == "ap" && eth.src == '$AP_MAC' && eth.dst == '$N1SUBNET6_MAC'
↳ && ip6.src == fc11::5 && ip6.dst == fc22::7 && ip.ttl == 64 && icmp6.type == 8'
...
ingress(dp="n1", inport="ap")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "ap" && eth.src ==
↳ {fa:16:3e:a9:4c:c7}, priority 50, uuid 6dcc418a
    next;
1. ls_in_port_sec_ip (ovn-northd.c:2390): inport == "ap" && eth.src ==
↳ fa:16:3e:a9:4c:c7 && ip6.src == {fe80::f816:3eff:fea9:4cc7, fc11::5}, priority 90,
↳ uuid 604810ea
    next;
3. ls_in_pre_acl (ovn-northd.c:2646): ip, priority 100, uuid 46c089e6
    reg0[0] = 1;
    next;
5. ls_in_pre_stateful (ovn-northd.c:2764): reg0[0] == 1, priority 100, uuid d1941634
    ct_next;

ct_next(ct_state=est|trk /* default (use --ct to customize) */)
-----

```

(continues on next page)

(continued from previous page)

```

6. ls_in_acl (ovn-northd.c:2925): !ct.new && ct.est && !ct.rpl && ct_label.blocked_
↪== 0 && (inport == "ap" && ip6), priority 2002, uuid 7fdd607e
    next;
13. ls_in_l2_lkup (ovn-northd.c:3529): eth.dst == fa:16:3e:ef:2f:8b, priority 50,
↪uuid eld87fc5
    output = "ad952e";
    output;

egress(dp="n1", inport="ap", output="ad952e")
-----
1. ls_out_pre_acl (ovn-northd.c:2626): ip && output == "ad952e", priority 110, uuid_
↪88f68988
    next;
8. ls_out_port_sec_l2 (ovn-northd.c:3654): output == "ad952e", priority 50, uuid_
↪5935755e
    output;
    /* output to "ad952e", type "patch" */

ingress(dp="r", inport="lrp-ad952e")
-----
0. lr_in_admission (ovn-northd.c:4071): eth.dst == fa:16:3e:ef:2f:8b && inport ==
↪"lrp-ad952e", priority 50, uuid ddfef712
    next;
5. lr_in_ip_routing (ovn-northd.c:3782): ip6.dst == fc22::/64, priority 129, uuid_
↪cc2130ec
    ip.ttl--;
    xxreg0 = ip6.dst;
    xxreg1 = fc22::1;
    eth.src = fa:16:3e:06:de:ad;
    output = "lrp-1a8162";
    flags.loopback = 1;
    next;
6. lr_in_arp_resolve (ovn-northd.c:5122): output == "lrp-1a8162" && xxreg0 ==_
↪fc22::7, priority 100, uuid bcf75288
    eth.dst = fa:16:3e:89:f2:36;
    next;
8. lr_in_arp_request (ovn-northd.c:5260): 1, priority 0, uuid 6dacdd82
    output;

egress(dp="r", inport="lrp-ad952e", output="lrp-1a8162")
-----
3. lr_out_delivery (ovn-northd.c:5288): output == "lrp-1a8162", priority 100, uuid_
↪5260dfc5
    output;
    /* output to "lrp-1a8162", type "patch" */

ingress(dp="n2", inport="1a8162")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "1a8162", priority 50, uuid_
↪10957d1b
    next;
3. ls_in_pre_acl (ovn-northd.c:2624): ip && inport == "1a8162", priority 110, uuid_
↪a27ebd00
    next;
13. ls_in_l2_lkup (ovn-northd.c:3529): eth.dst == fa:16:3e:89:f2:36, priority 50,
↪uuid dcafb3e9
    output = "cp";

```

(continues on next page)

(continued from previous page)

```

    output;

egress(dp="n2", inport="1a8162", outport="cp")
-----
1. ls_out_pre_acl (ovn-northd.c:2648): ip, priority 100, uuid cd9cfa74
   reg0[0] = 1;
   next;
2. ls_out_pre_stateful (ovn-northd.c:2766): reg0[0] == 1, priority 100, uuid 9e8e22c5
   ct_next;

ct_next(ct_state=est|trk /* default (use --ct to customize) */)
-----
4. ls_out_acl (ovn-northd.c:2925): !ct.new && ct.est && !ct.rpl && ct_label.blocked_
↪ == 0 && (outport == "cp" && ip6 && ip6.src == $as_ip6_0fc1b6cf_f925_49e6_8f00_
↪ 6ddl3beca9dc), priority 2002, uuid 12fc96f9
   next;
7. ls_out_port_sec_ip (ovn-northd.c:2390): outport == "cp" && eth.dst ==
↪ fa:16:3e:89:f2:36 && ip6.dst == {fe80::f816:3eff:fe89:f236, ff00::/8, fc22::7},
↪ priority 90, uuid c622596a
   next;
8. ls_out_port_sec_l2 (ovn-northd.c:3654): outport == "cp" && eth.dst ==
↪ {fa:16:3e:89:f2:36}, priority 50, uuid 0242cdc3
   output;
   /* output to "cp", type "" */

```

3.5.9 ACLs

Let's explore how ACLs work in OpenStack and OVN. In OpenStack, ACL rules are part of “security groups”, which are “default deny”, that is, packets are not allowed by default and the rules added to security groups serve to allow different classes of packets. The default group (named “default”) that is assigned to each of our VMs so far allows all traffic from our other VMs, which isn't very interesting for testing. So, let's create a new security group, which we'll name “custom”, add rules to it that allow incoming SSH and ICMP traffic, and apply this security group to VM c:

```

$ openstack security group create custom
$ openstack security group rule create --dst-port 22 custom
$ openstack security group rule create --protocol icmp custom
$ openstack server remove security group c default
$ openstack server add security group c custom

```

Now we can do some experiments to test security groups. From the console on a or b, it should now be possible to “ping” c or to SSH to it, but attempts to initiate connections on other ports should be blocked. (You can try to connect on another port with `ssh -p PORT IP` or `nc PORT IP`.) Connection attempts should time out rather than receive the “connection refused” or “connection reset” error that you would see between a and b.

It's also possible to test ACLs via `ovn-trace`, with one new wrinkle. `ovn-trace` can't simulate connection tracking state in the network, so by default it assumes that every packet represents an established connection. That's good enough for what we've been doing so far, but for checking properties of security groups we want to look at more detail.

If you look back at the VM-to-VM traces we've done until now, you can see that they execute two `ct_next` actions:

- The first of these is for the packet passing outward through the source VM's firewall. We can tell `ovn-trace` to treat the packet as starting a new connection or adding to an established connection by adding a `--ct` option: `--ct new` or `--ct est`, respectively. The latter is the default and therefore what we've been using so far. We can also use `--ct est, rpl`, which in addition to `--ct est` means that the connection was initiated by the destination VM rather than by the VM sending this packet.

- The second is for the packet passing inward through the destination VM's firewall. For this one, it makes sense to tell `ovn-trace` that the packet is starting a new connection, with `--ct new`, or that it is a packet sent in reply to a connection established by the destination VM, with `--ct est,rpl`.

`ovn-trace` uses the `--ct` options in order, so if we want to override the second `ct_next` behavior we have to specify two options.

Another useful `ovn-trace` option for this testing is `--minimal`, which reduces the amount of output. In this case we're really just interested in finding out whether the packet reaches the destination VM, that is, whether there's an eventual output action to `c`, so `--minimal` works fine and the output is easier to read.

Try a few traces. For example:

- VM `a` initiates a new SSH connection to `c`:

```
$ ovn-trace --ct new --ct new --minimal nl 'inport == "ap" && eth.src == '$AP_MAC
↪' && eth.dst == '$N1SUBNET6_MAC' && ip4.src == 10.1.1.5 && ip4.dst == 10.1.2.7 &
↪& ip.ttl == 64 && tcp.dst == 22'
...
ct_next(ct_state=new|trk) {
    ip.ttl--;
    eth.src = fa:16:3e:19:9f:46;
    eth.dst = fa:16:3e:89:f2:36;
    ct_next(ct_state=new|trk) {
        output("cp");
    };
};
```

This succeeds, as you can see since there is an `output` action.

- VM `a` initiates a new Telnet connection to `c`:

```
$ ovn-trace --ct new --ct new --minimal nl 'inport == "ap" && eth.src == '$AP_MAC
↪' && eth.dst == '$N1SUBNET6_MAC' && ip4.src == 10.1.1.5 && ip4.dst == 10.1.2.7 &
↪& ip.ttl == 64 && tcp.dst == 23'
ct_next(ct_state=new|trk) {
    ip.ttl--;
    eth.src = fa:16:3e:19:9f:46;
    eth.dst = fa:16:3e:89:f2:36;
    ct_next(ct_state=new|trk);
};
```

This fails, as you can see from the lack of an `output` action.

- VM `a` replies to a packet that is part of a Telnet connection originally initiated by `c`:

```
$ ovn-trace --ct est,rpl --ct est,rpl --minimal nl 'inport == "ap" && eth.src == '
↪$AP_MAC' && eth.dst == '$N1SUBNET6_MAC' && ip4.src == 10.1.1.5 && ip4.dst == 10.
↪1.2.7 && ip.ttl == 64 && tcp.dst == 23'
...
ct_next(ct_state=est|rpl|trk) {
    ip.ttl--;
    eth.src = fa:16:3e:19:9f:46;
    eth.dst = fa:16:3e:89:f2:36;
    ct_next(ct_state=est|rpl|trk) {
        output("cp");
    };
};
```

This succeeds, as you can see from the `output` action, since traffic received in reply to an outgoing connection is always allowed.

3.5.10 DHCP

As a final demonstration of the OVN architecture, let's examine the DHCP implementation. Like switching, routing, and NAT, the OVN implementation of DHCP involves configuration in the NB DB and logical flows in the SB DB.

Let's look at the DHCP support for a's port `ap`. The port's `Logical_Switch_Port` record shows that `ap` has DHCPv4 options:

```
$ ovn-nbctl list logical_switch_port ap | abbrev
_uuid                : ef17e5
addresses            : ["fa:16:3e:a9:4c:c7 10.1.1.5 fc11::5"]
dhcpv4_options       : 165974
dhcpv6_options       : 26f7cd
dynamic_addresses    : []
enabled              : true
external_ids          : {"neutron:port_name"=ap}
name                 : "820c08"
options              : {}
parent_name          : []
port_security        : ["fa:16:3e:a9:4c:c7 10.1.1.5 fc11::5"]
tag                  : []
tag_request          : []
type                 : ""
up                   : true
```

We can then list them either by UUID or, more easily, by port name:

```
$ ovn-nbctl list dhcp_options ap | abbrev
_uuid                : 165974
cidr                 : "10.1.1.0/24"
external_ids         : {"subnet_id"=5e67e7}
options              : {"lease_time"=43200, "mtu"=1442, "router"=10.1.1.1, "server_id"=
↪ "10.1.1.1", "server_mac"=fa:16:3e:bb:94:72}
```

These options show the basic DHCP configuration for the subnet. They do not include the IP address itself, which comes from the `Logical_Switch_Port` record. This allows a whole Neutron subnet to share a single `DHCP_Options` record. You can see this sharing in action, if you like, by listing the record for port `bp`, which is on the same subnet as `ap`, and see that it is the same record as before:

```
$ ovn-nbctl list dhcp_options bp | abbrev
_uuid                : 165974
cidr                 : "10.1.1.0/24"
external_ids         : {"subnet_id"=5e67e7}
options              : {"lease_time"=43200, "mtu"=1442, "router"=10.1.1.1, "server_id"=
↪ "10.1.1.1", "server_mac"=fa:16:3e:bb:94:72}
```

You can take another look at the southbound flow table if you like, but the best demonstration is to trace a DHCP packet. The following is a trace of a DHCP request inbound from `ap`. The first part is just the usual travel through the firewall:

```
$ ovn-trace nl 'inport == "ap" && eth.src == '$AP_MAC' && eth.dst == ↪
↪ ff:ff:ff:ff:ff:ff && ip4.dst == 255.255.255.255 && udp.src == 68 && udp.dst == 67 &&
↪ ip.ttl == 1'
```

(continues on next page)

(continued from previous page)

```

...
ingress(dp="n1", inport="ap")
-----
0. ls_in_port_sec_l2 (ovn-northd.c:3234): inport == "ap" && eth.src ==
↪{fa:16:3e:a9:4c:c7}, priority 50, uuid 6dcc418a
  next;
1. ls_in_port_sec_ip (ovn-northd.c:2325): inport == "ap" && eth.src ==
↪fa:16:3e:a9:4c:c7 && ip4.src == 0.0.0.0 && ip4.dst == 255.255.255.255 && udp.src ==
↪68 && udp.dst == 67, priority 90, uuid e46bed6f
  next;
3. ls_in_pre_acl (ovn-northd.c:2646): ip, priority 100, uuid 46c089e6
  reg0[0] = 1;
  next;
5. ls_in_pre_stateful (ovn-northd.c:2764): reg0[0] == 1, priority 100, uuid d1941634
  ct_next;

```

The next part is the new part. First, an ACL in table 6 allows a DHCP request to pass through. In table 11, the special `put_dhcp_opts` action replaces a DHCPDISCOVER or DHCPREQUEST packet by a reply. Table 12 flips the packet's source and destination and sends it back the way it came in:

```

6. ls_in_acl (ovn-northd.c:2925): !ct.new && ct.est && !ct.rpl && ct_label.blocked
↪== 0 && (inport == "ap" && ip4 && ip4.dst == {255.255.255.255, 10.1.1.0/24} && udp &
↪& udp.src == 68 && udp.dst == 67), priority 2002, uuid 9c90245d
  next;
11. ls_in_dhcp_options (ovn-northd.c:3409): inport == "ap" && eth.src ==
↪fa:16:3e:a9:4c:c7 && ip4.src == 0.0.0.0 && ip4.dst == 255.255.255.255 && udp.src ==
↪68 && udp.dst == 67, priority 100, uuid 8d63f29c
  reg0[3] = put_dhcp_opts(offerip = 10.1.1.5, lease_time = 43200, mtu = 1442,
↪netmask = 255.255.255.0, router = 10.1.1.1, server_id = 10.1.1.1);
  /* We assume that this packet is DHCPDISCOVER or DHCPREQUEST. */
  next;
12. ls_in_dhcp_response (ovn-northd.c:3438): inport == "ap" && eth.src ==
↪fa:16:3e:a9:4c:c7 && ip4 && udp.src == 68 && udp.dst == 67 && reg0[3], priority 100,
↪uuid 995eeaa9
  eth.dst = eth.src;
  eth.src = fa:16:3e:bb:94:72;
  ip4.dst = 10.1.1.5;
  ip4.src = 10.1.1.1;
  udp.src = 67;
  udp.dst = 68;
  output = inport;
  flags.loopback = 1;
  output;

```

Then the last part is just traveling back through the firewall to VM a:

```

egress(dp="n1", inport="ap", outport="ap")
-----
1. ls_out_pre_acl (ovn-northd.c:2648): ip, priority 100, uuid 3752b746
  reg0[0] = 1;
  next;
2. ls_out_pre_stateful (ovn-northd.c:2766): reg0[0] == 1, priority 100, uuid 0c066ea1
  ct_next;

ct_next(ct_state=est|trk /* default (use --ct to customize) */)
-----

```

(continues on next page)

(continued from previous page)

```

4. ls_out_acl (ovn-northd.c:3008): output == "ap" && eth.src == fa:16:3e:bb:94:72 &&
↪ ip4.src == 10.1.1.1 && udp && udp.src == 67 && udp.dst == 68, priority 34000, uuid
↪ 0b383e77
    ct_commit;
    next;
7. ls_out_port_sec_ip (ovn-northd.c:2364): output == "ap" && eth.dst ==
↪ fa:16:3e:a9:4c:c7 && ip4.dst == {255.255.255.255, 224.0.0.0/4, 10.1.1.5}, priority
↪ 90, uuid 7b8cbcd5
    next;
8. ls_out_port_sec_l2 (ovn-northd.c:3654): output == "ap" && eth.dst ==
↪ {fa:16:3e:a9:4c:c7}, priority 50, uuid b874ece8
    output;
    /* output to "ap", type " " */

```

3.5.11 Further Directions

We’ve looked at a fair bit of how OVN works and how it interacts with OpenStack. If you still have some interest, then you might want to explore some of these directions:

- Adding more than one hypervisor (“compute node”, in OpenStack parlance). OVN connects compute nodes by tunneling packets with the STT or Geneve protocols. OVN scales to 1000 compute nodes or more, but two compute nodes demonstrate the principle. All of the tools and techniques we demonstrated also work with multiple compute nodes.
- Container support. OVN supports seamlessly connecting VMs to containers, whether the containers are hosted on “bare metal” or nested inside VMs. OpenStack support for containers, however, is still evolving, and too difficult to incorporate into the tutorial at this point.
- Other kinds of gateways. In addition to floating IPs with NAT, OVN supports directly attaching VMs to a physical network and connecting logical switches to VTEP hardware.

3.6 OVN Role-Based Access Control (RBAC) Tutorial

This document provides a step-by-step guide for setting up role-based access control (RBAC) in OVN. In OVN, hypervisors (chassis) read and write the southbound database to do configuration. Without restricting hypervisor’s access to the southbound database, a compromised hypervisor might disrupt the entire OVN deployment by corrupting the database. RBAC ensures that each hypervisor can only modify its own data and thus improves the security of OVN. More details about the RBAC design can be found in `ovn-architecture(7)` manpage.

This document assumes OVN is installed in your system and runs normally.

3.6.1 Generating Certificates and Keys

In the OVN RBAC deployment, `ovn-controller` connects to the southbound database via SSL connection. The southbound database uses CA-signed certificate to authenticate `ovn-controller`.

Suppose there are three machines in your deployment. *machine_1* runs *chassis_1* and has IP address *machine_1-ip*. *machine_2* runs *chassis_2* and has IP address *machine_2-ip*. *machine_3* hosts southbound database and has IP address *machine_3-ip*. *machine_3* also hosts public key infrastructure (PKI).

1. Initiate PKI.

In *machine_3*:

```
$ ovs-pki init
```

2. Generate southbound database's certificate request. Sign the certificate request with the CA key.

In *machine_3*:

```
$ ovs-pki req -u sbdb
$ ovs-pki sign sbdb switch
```

3. Generate chassis certificate requests. Copy the certificate requests to *machine_3*.

In *machine_1*:

```
$ ovs-pki req -u chassis_1
$ scp chassis_1-req.pem \
    machine_3@machine_3-ip:/path/to/chassis_1-req.pem
```

In *machine_2*:

```
$ ovs-pki req -u chassis_2
$ scp chassis_2-req.pem \
    machine_3@machine_3-ip:/path/to/chassis_2-req.pem
```

Note: *chassis_1* must be the same string as *external_ids:system-id* in the Open_vSwitch table (the chassis name) of *machine_1*. Same applies for *chassis_2*.

4. Sign the chassis certificate requests with the CA key. Copy *chassis_1*'s signed certificate and the CA certificate to *machine_1*. Copy *chassis_2*'s signed certificate and the CA certificate to *machine_2*.

In *machine_3*:

```
$ ovs-pki sign chassis_1 switch
$ ovs-pki sign chassis_2 switch
$ scp chassis_1-cert.pem \
    machine_1@machine_1-ip:/path/to/chassis_1-cert.pem
$ scp /var/lib/openvswitch/pki/switchca/cacert.pem \
    machine_1@machine_1-ip:/path/to/cacert.pem
$ scp chassis_2-cert.pem \
    machine_2@machine_2-ip:/path/to/chassis_2-cert.pem
$ scp /var/lib/openvswitch/pki/switchca/cacert.pem \
    machine_2@machine_2-ip:/path/to/cacert.pem
```

3.6.2 Configuring RBAC

1. Set certificate, private key, and CA certificate for the southbound database. Configure the southbound database to listen on SSL connection and enforce role-based access control.

In *machine_3*:

```
$ ovn-sbctl set-ssl /path/to/sbdb-privkey.pem \
    /path/to/sbdb-cert.pem /path/to/cacert.pem
$ ovn-sbctl set-connection role=ovn-controller pssl:6642
```

2. Set certificate, private key, and CA certificate for *chassis_1* and *chassis_2*. Configure *chassis_1* and *chassis_2* to connect southbound database via SSL.

In *machine_1*:

```
$ ovs-vsctl set-ssl /path/to/chassis_1-privkey.pem \
                    /path/to/chassis_1-cert.pem /path/to/cacert.pem
$ ovs-vsctl set open_vswitch . \
                    external_ids:ovn-remote=ssl:machine_3-ip:6642
```

In *machine_2*:

```
$ ovs-vsctl set-ssl /path/to/chassis_2-privkey.pem \
                    /path/to/chassis_2-cert.pem /path/to/cacert.pem
$ ovs-vsctl set open_vswitch . \
                    external_ids:ovn-remote=ssl:machine_3-ip:6642
```

3.7 OVN IPsec Tutorial

This document provides a step-by-step guide for encrypting tunnel traffic with IPsec in Open Virtual Network (OVN). OVN tunnel traffic is transported by physical routers and switches. These physical devices could be untrusted (devices in public network) or might be compromised. Enabling IPsec encryption for the tunnel traffic can prevent the traffic data from being monitored and manipulated. More details about the OVN IPsec design can be found in [ovn-architecture\(7\)](#) manpage.

This document assumes OVN is installed in your system and runs normally. Also, you need to install OVS IPsec packages in each chassis (refer to [Installing OVS and IPsec Packages](#)).

3.7.1 Generating Certificates and Keys

OVN chassis uses CA-signed certificate to authenticate peer chassis for building IPsec tunnel. If you have enabled Role-Based Access Control (RBAC) in OVN, you can use the RBAC SSL certificates and keys to set up OVN IPsec. Or you can generate separate certificates and keys with `ovs-pki` (refer to [Generating Certificates and Keys](#)).

Note: OVN IPsec requires x.509 version 3 certificate with the subjectAltName DNS field setting the same string as the common name (CN) field. CN should be set as the chassis name. `ovs-pki` in Open vSwitch 2.10.90 and later generates such certificates. Please generate compatible certificates if you use another PKI tool, or an older version of `ovs-pki`, to manage certificates.

3.7.2 Configuring OVN IPsec

You need to install the CA certificate, chassis certificate and private key in each chassis. Use the following command:

```
$ ovs-vsctl set Open_vSwitch . \
    other_config:certificate=/path/to/chassis-cert.pem \
    other_config:private_key=/path/to/chassis-privkey.pem \
    other_config:ca_cert=/path/to/cacert.pem
```

3.7.3 Enabling OVN IPsec

To enable OVN IPsec, set `ipsec` column in `NB_Global` table of the northbound database to true:

```
$ ovn-nbctl set nb_global . ipsec=true
```

With OVN IPsec enabled, all tunnel traffic in OVN will be encrypted with IPsec. To disable it, set `ipsec` column in `NB_Global` table of the northbound database to `false`:

```
$ ovn-nbctl set nb_global . ipsec=false
```

3.7.4 Troubleshooting

The `ovs-monitor-ipsec` daemon in each chassis manages and monitors the IPsec tunnel state. Use the following `ovs-appctl` command to view `ovs-monitor-ipsec` internal representation of tunnel configuration:

```
$ ovs-appctl -t ovs-monitor-ipsec tunnels/show
```

If there is a misconfiguration, then `ovs-appctl` should indicate why. For example:

```
Interface name: ovn-host_2-0 v1 (CONFIGURED) <--- Should be set
                                                    to CONFIGURED. Otherwise,
                                                    error message will be
                                                    provided

Tunnel Type:      geneve
Remote IP:        2.2.2.2
SKB mark:         None
Local cert:       /path/to/chassis-cert.pem
Local name:       host_1
Local key:        /path/to/chassis-privkey.pem
Remote cert:      None
Remote name:      host_2
CA cert:         /path/to/cacert.pem
PSK:              None
Ofport:          2      <--- Whether ovs-vswitchd has assigned Ofport
                        number to this Tunnel Port
CFM state:        Disabled <--- Whether CFM declared this tunnel healthy
Kernel policies installed:
...              <--- IPsec policies for this OVS tunnel in
                        Linux Kernel installed by strongSwan
Kernel security associations installed:
...              <--- IPsec security associations for this OVS
                        tunnel in Linux Kernel installed by
                        strongswan
IPsec connections that are active:
...              <--- IPsec "connections" for this OVS
                        tunnel
```

If you don't see any active connections, try to run the following command to refresh the `ovs-monitor-ipsec` daemon:

```
$ ovs-appctl -t ovs-monitor-ipsec refresh
```

You can also check the logs of the `ovs-monitor-ipsec` daemon and the IKE daemon to locate issues. `ovs-monitor-ipsec` outputs log messages to `/var/log/openvswitch/ovs-monitor-ipsec.log`.

3.7.5 Bug Reporting

If you think you may have found a bug with security implications, like

1. IPsec protected tunnel accepted packets that came unencrypted; OR
2. IPsec protected tunnel allowed packets to leave unencrypted;

Then report such bugs according to *Open vSwitch's Security Process*.

If bug does not have security implications, then report it according to instructions in *Reporting Bugs in Open vSwitch*.

If you have suggestions to improve this tutorial, please send a email to ovs-discuss@openvswitch.org.

3.8 OVS Conntrack Tutorial

OVS can be used with the Connection tracking system where OpenFlow flow can be used to match on the state of a TCP, UDP, ICMP, etc., connections. (Connection tracking system supports tracking of both stateful and stateless protocols)

This tutorial demonstrates how OVS can use the connection tracking system to match on the TCP segments from connection setup to connection tear down. It will use OVS with the Linux kernel module as the datapath for this tutorial. (The datapath that utilizes the openvswitch kernel module to do the packet processing in the Linux kernel) It was tested with the “master” branch of Open vSwitch.

3.8.1 Definitions

conntrack: is a connection tracking module for stateful packet inspection.

pipeline: is the packet processing pipeline which is the path taken by the packet when traversing through the tables where the packet matches the match fields of a flow in the table and performs the actions present in the matched flow.

network namespace: is a way to create virtual routing domains within a single instance of linux kernel. Each network namespace has it's own instance of network tables (arp, routing) and certain interfaces attached to it.

flow: used in this tutorial refers to the OpenFlow flow which can be programmed using an OpenFlow controller or OVS command line tools like `ovs-ofctl` which is used here. A flow will have match fields and actions.

3.8.2 Conntrack Related Fields

Match Fields

OVS supports following match fields related to conntrack:

1. **ct_state:** The state of a connection matching the packet. Possible values:

- *new*
- *est*
- *rel*
- *rpl*
- *inv*
- *trk*
- *snat*
- *dnat*

Each of these flags is preceded by either a “+” for a flag that must be set, or a “-” for a flag that must be unset. Multiple flags can also be specified e.g. `ct_state=+trk+new`. We will see the usage of some of these flags below. For a detailed description, please see the OVS fields documentation at: <http://openvswitch.org/support/dist-docs/ovs-fields.7.txt>

2. **ct_zone**: A zone is an independent connection tracking context which can be set by a `ct` action. A 16-bit `ct_zone` set by the most recent `ct` action (by an OpenFlow flow on a conntrack entry) can be used as a match field in another flow entry.
3. **ct_mark**: The 32-bit metadata committed, by an action within the `exec` parameter to the `ct` action, to the connection to which the current packet belongs.
4. **ct_label**: The 128-bit label committed by an action within the `exec` parameter to the `ct` action, to the connection to which the current packet belongs.
5. **ct_nw_src** / **ct_ipv6_src**: Matches IPv4/IPv6 conntrack original direction tuple source address.
6. **ct_nw_dst** / **ct_ipv6_dst**: Matches IPv4/IPv6 conntrack original direction tuple destination address.
7. **ct_nw_proto**: Matches conntrack original direction tuple IP protocol type.
8. **ct_tp_src**: Matches on the conntrack original direction tuple transport source port.
9. **ct_tp_dst**: Matches on the conntrack original direction tuple transport destination port.

Actions

OVS supports “`ct`” action related to conntrack.

ct([argument][,argument...])

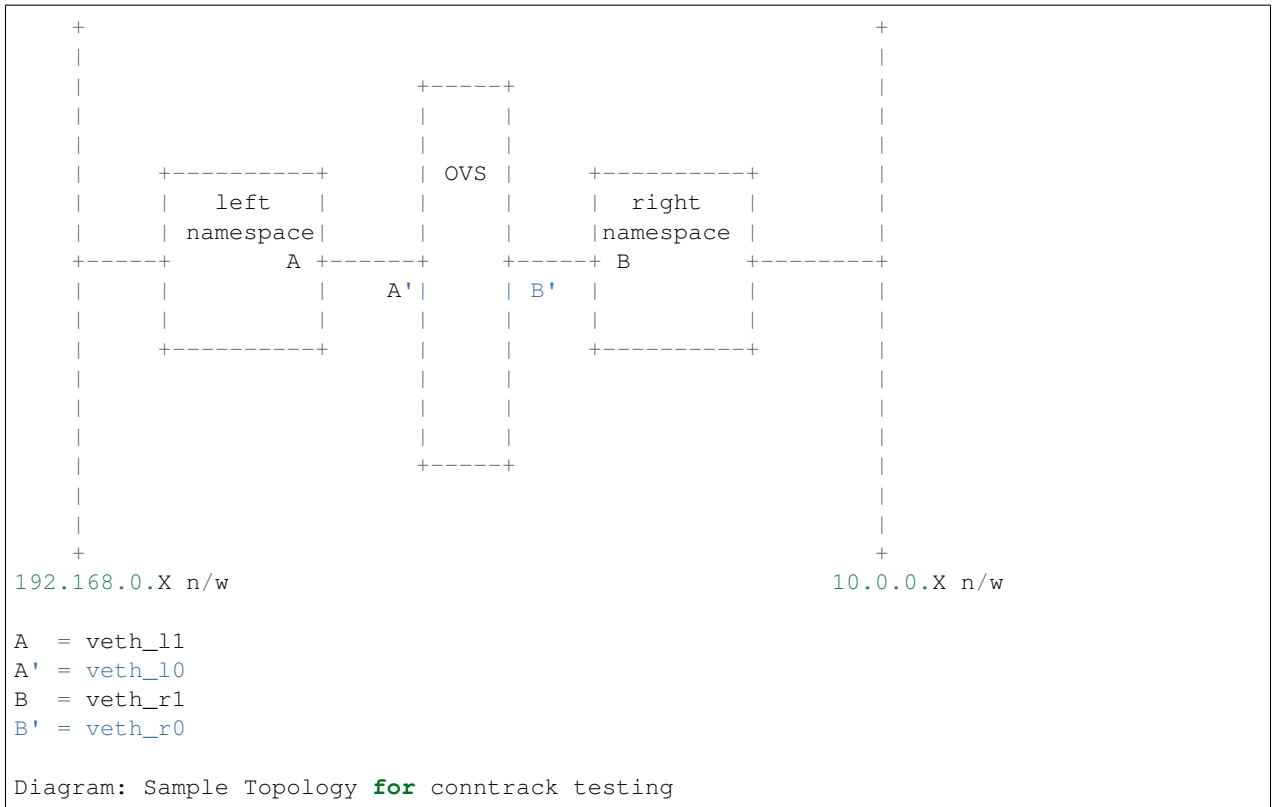
The **ct** action sends the packet through the connection tracker.

The following arguments are supported:

1. **commit**: Commit the connection to the connection tracking module which will be stored beyond the lifetime of packet in the pipeline.
2. **force**: The force flag may be used in addition to commit flag to effectively terminate the existing connection and start a new one in the current direction.
3. **table=number**: Fork pipeline processing in two. The original instance of the packet will continue processing the current actions list as an untracked packet. An additional instance of the packet will be sent to the connection tracker, which will be re-injected into the OpenFlow pipeline to resume processing in table number, with the `ct_state` and other `ct` match fields set.
4. **zone=value OR zone=src[start..end]**: A 16-bit context id that can be used to isolate connections into separate domains, allowing over-lapping network addresses in different zones. If a zone is not provided, then the default is to use zone zero.
5. **exec([action][,action...])**: Perform restricted set of actions within the context of connection tracking. Only actions which modify the `ct_mark` or `ct_label` fields are accepted within the `exec` action.
6. **alg=<ftp/tftp>**: Specify alg (application layer gateway) to track specific connection types.
7. **nat**: Specifies the address and port translation for the connection being tracked.

3.8.3 Sample Topology

This tutorial uses the following topology to carry out the tests.



The steps for creation of the setup are mentioned below.

Create “left” network namespace:

```
$ ip netns add left
```

Create “right” network namespace:

```
$ ip netns add right
```

Create first pair of veth interfaces:

```
$ ip link add veth_l0 type veth peer name veth_l1
```

Add veth_l1 to “left” network namespace:

```
$ ip link set veth_l1 netns left
```

Create second pair of veth interfaces:

```
$ ip link add veth_r0 type veth peer name veth_r1
```

Add veth_r1 to “right” network namespace:

```
$ ip link set veth_r1 netns right
```

Create a bridge br0:

```
$ ovs-vsctl add-br br0
```

Add veth_l0 and veth_r0 to br0:

```
$ ovs-vsctl add-port br0 veth_l0
$ ovs-vsctl add-port br0 veth_r0
```

Packets generated with src/dst IP set to 192.168.0.X / 10.0.0.X in the “left” and the inverse in the “right” namespaces will appear to OVS as hosts in two networks (192.168.0.X and 10.0.0.X) communicating with each other. This is basically a simulation of two networks / subnets with hosts communicating with each other with OVS in middle.

3.8.4 Tool used to generate TCP segments

You can use scapy to generate the TCP segments. We used scapy on Ubuntu 16.04 for the steps carried out in this testing. (Installation of scapy is not discussed and is out of scope of this document.)

You can keep two scapy sessions active on each of the namespaces:

```
$ sudo ip netns exec left sudo `which scapy`
$ sudo ip netns exec right sudo `which scapy`
```

Note: In case you encounter this error:

```
ifreq = ioctl(s, SIOCGIFADDR, struct.pack("16s16x", LOOPBACK_NAME))
IOError: [Errno 99] Cannot assign requested address
```

run the command:

```
$ sudo ip netns exec <namespace> sudo ip link set lo up
```

3.8.5 Matching TCP packets

TCP Connection setup

Two simple flows can be added in OVS which will forward packets from “left” to “right” and from “right” to “left”:

```
$ ovs-ofctl add-flow br0 \
    "table=0, priority=10, in_port=veth_l0, actions=veth_r0"
$ ovs-ofctl add-flow br0 \
    "table=0, priority=10, in_port=veth_r0, actions=veth_l0"
```

Instead of adding these two flows, we will add flows to match on the states of the TCP segments.

We will send the TCP connection setup segments namely: syn, syn-ack and ack between hosts 192.168.0.2 in the “left” namespace and 10.0.0.2 in the “right” namespace.

First, let’s add a flow to start “tracking” a packet received at OVS.

How do we start tracking a packet?

To start tracking a packet, it first needs to match a flow, which has action as “ct”. This action sends the packet through the connection tracker. To identify that a packet is an “untracked” packet, the ct_state in the flow match field must be set to “-trk”, which means it is not a tracked packet. Once the packet is sent to the connection tracker, then only we will know about its conntrack state. (i.e. whether this packet represents start of a new connection or the packet belongs to an existing connection or it is a malformed packet and so on.)

Let's add that flow:

```
(flow #1)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=-trk, tcp, in_port=veth_l0, actions=ct(table=0)"
```

A TCP syn packet sent from “left” namespace will match flow #1 because the packet is coming to OVS from veth_l0 port and it is not being tracked. (as the packet just entered OVS. All packets entering OVS for the first time are “untracked”) The flow will send the packet to the connection tracker due to the action “ct”. Also “table=0” in the “ct” action forks the pipeline processing in two. The original instance of packet will continue processing the current action list as untracked packet. (Since there are no actions after this, the original packet gets dropped.) The forked instance of the packet will be sent to the connection tracker, which will be re-injected into the OpenFlow pipeline to resume processing in table number, with the ct_state and other ct match fields set. In this case, the packet with the ct_state and other ct match fields comes back to table 0.

Next, we add a flow to match on the packet coming back from conntrack:

```
(flow #2)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=+trk,+new, tcp, in_port=veth_l0, \
    →actions=ct(commit),veth_r0"
```

Now that the packet is coming back from conntrack, the ct_state would have the “trk” set. Also, if this is the first packet of the TCP connection, the ct_state “new” would be set. (Which is the condition here as there does not exist any TCP connection between hosts 192.168.0.2 and 10.0.0.2) The ct argument “commit” will commit the connection to the connection tracking module. The significance of this action is that the information about the connection will now be stored beyond the lifetime of the packet in the pipeline.

Let's send the TCP syn segment using scapy (at the “left” scapy session) (flags=0x02 is syn):

```
$ >>> sendp(Ether()/IP(src="192.168.0.2", dst="10.0.0.2")/TCP(sport=1024, dport=2048, \
    →flags=0x02, seq=100), iface="veth_l1")
```

This packet will match flow #1 and flow #2.

The conntrack module will now have an entry for this connection:

```
$ ovs-appctl dpctl/dump-conntrack | grep "192.168.0.2"
tcp,orig=(src=192.168.0.2,dst=10.0.0.2,sport=1024,dport=2048),reply=(src=10.0.0.2,
    →dst=192.168.0.2,sport=2048,dport=1024),protoinfo=(state=SYN_SENT)
```

Note: At this stage, if the TCP syn packet is re-transmitted, it will again match flow #1 (since a new packet is untracked) and it will match flow #2. The reason it will match flow #2 is that although conntrack has information about the connection, but it is not in “ESTABLISHED” state, therefore it matches the “new” state again.

Next for the TCP syn-ack from the opposite/server direction, we need following flows at OVS:

```
(flow #3)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=-trk, tcp, in_port=veth_r0, actions=ct(table=0)"
(flow #4)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=+trk,+est, tcp, in_port=veth_r0, actions=veth_l0"
```

flow #3 matches untracked packets coming back from server (10.0.0.2) and sends this to conntrack. (Alternatively, we could have also combined flow #1 and flow #3 into one flow by not having the “in_port” match)

The syn-ack packet which has now gone through the conntrack has the ct_state of “est”.

Note: Conntrack puts the `ct_state` of the connection to “est” state when it sees bidirectional traffic, but till it does not get the third ack from client, it puts a short cleanup timer on the conntrack entry.

Sending TCP syn-ack segment using scapy (at the “right” scapy session) (flags=0x12 is ack and syn):

```
$ >>> sendp(Ether()/IP(src="10.0.0.2", dst="192.168.0.2")/TCP(sport=2048, dport=1024,
↳ flags=0x12, seq=200, ack=101), iface="veth_r1")
```

This packet will match flow #3 and flow #4.

conntrack entry:

```
$ ovs-appctl dpctl/dump-conntrack | grep "192.168.0.2"

tcp,orig=(src=192.168.0.2,dst=10.0.0.2,sport=1024,dport=2048),reply=(src=10.0.0.2,
↳ dst=192.168.0.2,sport=2048,dport=1024),protoinfo=(state=ESTABLISHED)
```

The conntrack state is “ESTABLISHED” on receiving just syn and syn-ack packets, but at this point if it does not receive the third ack (from client), the connection gets cleared up from conntrack quickly.

Next, for a TCP ack from client direction, we can add following flows to match on the packet:

```
(flow #5)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=+trk,+est, tcp, in_port=veth_l0, actions=veth_r0"
```

Send the third TCP ack segment using scapy (at the “left” scapy session) (flags=0x10 is ack):

```
$ >>> sendp(Ether()/IP(src="192.168.0.2", dst="10.0.0.2")/TCP(sport=1024, dport=2048,
↳ flags=0x10, seq=101, ack=201), iface="veth_l1")
```

This packet will match on flow #1 and flow #5.

conntrack entry:

```
$ ovs-appctl dpctl/dump-conntrack | grep "192.168.0.2"

tcp,orig=(src=192.168.0.2,dst=10.0.0.2,sport=1024,dport=2048), \
    reply=(src=10.0.0.2,dst=192.168.0.2,sport=2048,dport=1024), \
    protoinfo=(state=ESTABLISHED)
```

The conntrack state stays in “ESTABLISHED” state, but now since it has received the ack from client, it will stay in this state for a longer time even without receiving any data on this connection.

TCP Data

When a data segment, carrying one byte of TCP payload, is sent from 192.168.0.2 to 10.0.0.2, the packet carrying the segment would hit flow #1 and then flow #5.

Send a TCP segment with one byte data using scapy (at the “left” scapy session) (flags=0x10 is ack):

```
$ >>> sendp(Ether()/IP(src="192.168.0.2", dst="10.0.0.2")/TCP(sport=1024, dport=2048,
↳ flags=0x10, seq=101, ack=201)/"X", iface="veth_l1")
```

Send the TCP ack for the above segment using scapy (at the “right” scapy session) (flags=0x10 is ack):

```
$ >>> sendp(Ether()/IP(src="10.0.0.2", dst="192.168.0.2")/TCP(sport=2048, dport=1024,
↳ flags=0x10, seq=201, ack=102), iface="veth_r1")
```

The acknowledgement for the data would hit flow #3 and flow #4.

TCP Connection Teardown

There are different ways to tear down TCP connection. We will tear down the connection by sending “fin” from client, “fin-ack” from server followed by the last “ack” by client.

All the packets from client to server would hit flow #1 and flow #5. All the packets from server to client would hit flow #3 and flow #4. Interesting point to note is that even when the TCP connection is going down, all the packets (which are actually tearing down the connection) still hits “+est” state. A packet, for which the conntrack entry *is* or *was* in “ESTABLISHED” state, would continue to match “+est” ct_state in OVS.

Note: In fact, when the conntrack connection state is in “TIME_WAIT” state (after all the TCP fins and their acks are exchanged), a re-transmitted data packet (from 192.168.0.2 -> 10.0.0.2), still hits flows #1 and #5.

Sending TCP fin segment using scapy (at the “left” scapy session) (flags=0x11 is ack and fin):

```
$ >>> sendp(Ether()/IP(src="192.168.0.2", dst="10.0.0.2")/TCP(sport=1024, dport=2048,
↳ flags=0x11, seq=102, ack=201), iface="veth_l1")
```

This packet hits flow #1 and flow #5.

conntrack entry:

```
$ sudo ovs-appctl dpctl/dump-conntrack | grep "192.168.0.2"

tcp,orig=(src=192.168.0.2,dst=10.0.0.2,sport=1024,dport=2048),reply=(src=10.0.0.2,
↳ dst=192.168.0.2,sport=2048,dport=1024),protoinfo=(state=FIN_WAIT_1)
```

Sending TCP fin-ack segment using scapy (at the “right” scapy session) (flags=0x11 is ack and fin):

```
$ >>> sendp(Ether()/IP(src="10.0.0.2", dst="192.168.0.2")/TCP(sport=2048, dport=1024,
↳ flags=0x11, seq=201, ack=103), iface="veth_r1")
```

This packet hits flow #3 and flow #4.

conntrack entry:

```
$ sudo ovs-appctl dpctl/dump-conntrack | grep "192.168.0.2"

tcp,orig=(src=192.168.0.2,dst=10.0.0.2,sport=1024,dport=2048),reply=(src=10.0.0.2,
↳ dst=192.168.0.2,sport=2048,dport=1024),protoinfo=(state=LAST_ACK)
```

Sending TCP ack segment using scapy (at the “left” scapy session) (flags=0x10 is ack):

```
$ >>> sendp(Ether()/IP(src="192.168.0.2", dst="10.0.0.2")/TCP(sport=1024, dport=2048,
↳ flags=0x10, seq=103, ack=202), iface="veth_l1")
```

This packet hits flow #1 and flow #5.

conntrack entry:

```
$ sudo ovs-appctl dpctl/dump-conntrack | grep "192.168.0.2"

tcp,orig=(src=192.168.0.2,dst=10.0.0.2,sport=1024,dport=2048),reply=(src=10.0.0.2,
↳ dst=192.168.0.2,sport=2048,dport=1024),protoinfo=(state=TIME_WAIT)
```

3.8.6 Summary

Following table summarizes the TCP segments exchanged against the flow match fields

TCP Segment	ct_state(flow#)
Connection Setup	
192.168.0.2 → 10.0.0.2 [SYN] Seq=0	-trk(#1) then +trk+new(#2)
10.0.0.2 → 192.168.0.2 [SYN, ACK] Seq=0 Ack=1	-trk(#3) then +trk+est(#4)
192.168.0.2 → 10.0.0.2 [ACK] Seq=1 Ack=1	-trk(#1) then +trk+est(#5)
Data Transfer	
192.168.0.2 → 10.0.0.2 [ACK] Seq=1 Ack=1	-trk(#1) then +trk+est(#5)
10.0.0.2 → 192.168.0.2 [ACK] Seq=1 Ack=2	-trk(#3) then +trk+est(#4)
Connection Teardown	
192.168.0.2 → 10.0.0.2 [FIN, ACK] Seq=2 Ack=1	-trk(#1) then +trk+est(#5)
10.0.0.2 → 192.168.0.2 [FIN, ACK] Seq=1 Ack=3	-trk(#3) then +trk+est(#4)
192.168.0.2 → 10.0.0.2 [ACK] Seq=3 Ack=2	-trk(#1) then +trk+est(#5)

Note: Relative sequence number and acknowledgement numbers are shown as captured from tshark.

Flows

```
(flow #1)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=-trk, tcp, in_port=veth_l0, actions=ct(table=0)"

(flow #2)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=+trk,+new, tcp, in_port=veth_l0, \
    →actions=ct(commit),veth_r0"

(flow #3)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=-trk, tcp, in_port=veth_r0, actions=ct(table=0)"

(flow #4)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=+trk,+est, tcp, in_port=veth_r0, actions=veth_l0"

(flow #5)
$ ovs-ofctl add-flow br0 \
    "table=0, priority=50, ct_state=+trk,+est, tcp, in_port=veth_l0, actions=veth_r0"
```

How Open vSwitch and OVN are implemented and, where necessary, why it was implemented that way.

4.1 OVS

4.1.1 Design Decisions In Open vSwitch

This document describes design decisions that went into implementing Open vSwitch. While we believe these to be reasonable decisions, it is impossible to predict how Open vSwitch will be used in all environments. Understanding assumptions made by Open vSwitch is critical to a successful deployment. The end of this document contains contact information that can be used to let us know how we can make Open vSwitch more generally useful.

Asynchronous Messages

Over time, Open vSwitch has added many knobs that control whether a given controller receives OpenFlow asynchronous messages. This section describes how all of these features interact.

First, a service controller never receives any asynchronous messages unless it changes its `miss_send_len` from the service controller default of zero in one of the following ways:

- Sending an `OFPT_SET_CONFIG` message with nonzero `miss_send_len`.
- Sending any `NXT_SET_ASYNC_CONFIG` message: as a side effect, this message changes the `miss_send_len` to `OFP_DEFAULT_MISS_SEND_LEN` (128) for service controllers.

Second, `OFPT_FLOW_REMOVED` and `NXT_FLOW_REMOVED` messages are generated only if the flow that was removed had the `OFPFF_SEND_FLOW_REM` flag set.

Third, `OFPT_PACKET_IN` and `NXT_PACKET_IN` messages are sent only to OpenFlow controller connections that have the correct connection ID (see `struct nx_controller_id` and `struct nx_action_controller`):

- For packet-in messages generated by a `NXAST_CONTROLLER` action, the controller ID specified in the action.

- For other packet-in messages, controller ID zero. (This is the default ID when an OpenFlow controller does not configure one.)

Finally, Open vSwitch consults a per-connection table indexed by the message type, reason code, and current role. The following table shows how this table is initialized by default when an OpenFlow connection is made. An entry labeled *yes* means that the message is sent, an entry labeled *---* means that the message is suppressed.

Table 1: OFPT_PACKET_IN / NXT_PACKET_IN

message and reason code	other	slave
OFPR_NO_MATCH	yes	---
OFPR_ACTION	yes	---
OFPR_INVALID_TTL	---	---
OFPR_ACTION_SET (OF1.4+)	yes	---
OFPR_GROUP (OF1.4+)	yes	---
OFPR_PACKET_OUT (OF1.4+)	yes	---

Table 2: OFPT_FLOW_REMOVED / NXT_FLOW_REMOVED

message and reason code	other	slave
OFPRR_IDLE_TIMEOUT	yes	---
OFPRR_HARD_TIMEOUT	yes	---
OFPRR_DELETE	yes	---
OFPRR_GROUP_DELETE (OF1.3+)	yes	---
OFPRR_METER_DELETE (OF1.4+)	yes	---
OFPRR_EVICTION (OF1.4+)	yes	---

Table 3: OFPT_PORT_STATUS

message and reason code	other	slave
OFPPR_ADD	yes	yes
OFPPR_DELETE	yes	yes
OFPPR_MODIFY	yes	yes

Table 4: OFPT_ROLE_REQUEST / OFPT_ROLE_REPLY (OF1.4+)

message and reason code	other	slave
OFPCRR_MASTER_REQUEST	---	---
OFPCRR_CONFIG	---	---
OFPCRR_EXPERIMENTER	---	---

Table 5: OFPT_TABLE_STATUS (OF1.4+)

message and reason code	other	slave
OFPTR_VACANCY_DOWN	---	---
OFPTR_VACANCY_UP	---	---

Table 6: OFPT_REQUESTFORWARD (OF1.4+)

message and reason code	other	slave
OFPRFR_GROUP_MOD	---	---
OFPRFR_METER_MOD	---	---

The `NXT_SET_ASYNC_CONFIG` message directly sets all of the values in this table for the current connection. The `OFPC_INVALID_TTL_TO_CONTROLLER` bit in the `OFPT_SET_CONFIG` message controls the setting for `OFPR_INVALID_TTL` for the “master” role.

OFPAT_ENQUEUE

The OpenFlow 1.0 specification requires the output port of the `OFPAT_ENQUEUE` action to “refer to a valid physical port (i.e. `< OFPP_MAX`) or `OFPP_IN_PORT`”. Although `OFPP_LOCAL` is not less than `OFPP_MAX`, it is an ‘internal’ port which can have QoS applied to it in Linux. Since we allow the `OFPAT_ENQUEUE` to apply to ‘internal’ ports whose port numbers are less than `OFPP_MAX`, we interpret `OFPP_LOCAL` as a physical port and support `OFPAT_ENQUEUE` on it as well.

OFPT_FLOW_MOD

The OpenFlow specification for the behavior of `OFPT_FLOW_MOD` is confusing. The following tables summarize the Open vSwitch implementation of its behavior in the following categories:

- “**match on priority**” Whether the `flow_mod` acts only on flows whose priority matches that included in the `flow_mod` message.
- “**match on out_port**” Whether the `flow_mod` acts only on flows that output to the `out_port` included in the `flow_mod` message (if `out_port` is not `OFPP_NONE`). OpenFlow 1.1 and later have a similar feature (not listed separately here) for `out_group`.
- “**match on flow_cookie**”: Whether the `flow_mod` acts only on flows whose `flow_cookie` matches an optional controller-specified value and mask.
- “**updates flow_cookie**”: Whether the `flow_mod` changes the `flow_cookie` of the flow or flows that it matches to the `flow_cookie` included in the `flow_mod` message.
- “**updates OFPFF_flags**”: Whether the `flow_mod` changes the `OFPFF_SEND_FLOW_REM` flag of the flow or flows that it matches to the setting included in the flags of the `flow_mod` message.
- “**honors OFPFF_CHECK_OVERLAP**”: Whether the `OFPFF_CHECK_OVERLAP` flag in the `flow_mod` is significant.
- “**updates idle_timeout**” and “**updates hard_timeout**”: Whether the `idle_timeout` and `hard_timeout` in the `flow_mod`, respectively, have an effect on the flow or flows matched by the `flow_mod`.
- “**updates idle timer**”: Whether the `flow_mod` resets the per-flow timer that measures how long a flow has been idle.
- “**updates hard timer**”: Whether the `flow_mod` resets the per-flow timer that measures how long it has been since a flow was modified.
- “**zeros counters**”: Whether the `flow_mod` resets per-flow packet and byte counters to zero.
- “**may add a new flow**”: Whether the `flow_mod` may add a new flow to the flow table. (Obviously this is always true for “add” commands but in some OpenFlow versions “modify” and “modify-strict” can also add new flows.)
- “**sends flow_removed message**”: Whether the `flow_mod` generates a `flow_removed` message for the flow or flows that it affects.

An entry labeled `yes` means that the flow mod type does have the indicated behavior, `---` means that it does not, an empty cell means that the property is not applicable, and other values are explained below the table.

OpenFlow 1.0

RULE	ADD	MODIFY	STRICT	DELETE	STRICT
match on priority	yes	—	yes	—	yes
match on out_port	—	—	—	yes	yes
match on flow_cookie	—	—	—	—	—
match on table_id	—	—	—	—	—
controller chooses table_id	—	—	—		
updates flow_cookie	yes	yes	yes		
updates OFPFF_SEND_FLOW_REM	yes	•	•		
honors OFPFF_CHECK_OVERLAP	yes	•	•		
updates idle_timeout	yes	•	•		
updates hard_timeout	yes	•	•		
resets idle timer	yes	•	•		
resets hard timer	yes	yes	yes		
zeros counters	yes	•	•		
may add a new flow	yes	yes	yes		
sends flow_removed message	—	—	—	%	%

where:

- + “modify” and “modify-strict” only take these actions when they create a new flow, not when they update an existing flow.
- % “delete” and “delete-strict” generates a flow_removed message if the deleted flow or flows have the OFPFF_SEND_FLOW_REM flag set. (Each controller can separately control whether it wants to receive the generated messages.)

OpenFlow 1.1

OpenFlow 1.1 makes these changes:

- The controller now must specify the table_id of the flow match searched and into which a flow may be inserted. Behavior for a table_id of 255 is undefined.

- A `flow_mod`, except an “add”, can now match on the `flow_cookie`.
- When a `flow_mod` matches on the `flow_cookie`, “modify” and “modify-strict” never insert a new flow.

RULE	ADD	MODIFY	STRICT	DELETE	STRICT
match on priority	yes	—	yes	—	yes
match on out_port	—	—	—	yes	yes
match on flow_cookie	—	yes	yes	yes	yes
match on table_id	yes	yes	yes	yes	yes
controller chooses table_id	yes	yes	yes		
updates flow_cookie	yes	—	—		
updates OFPFF_SEND_FLOW_REM	yes	•	•		
honors OFPFF_CHECK_OVERLAP	yes	•	•		
updates idle_timeout	yes	•	•		
updates hard_timeout	yes	•	•		
resets idle timer	yes	•	•		
resets hard timer	yes	yes	yes		
zeros counters	yes	•	•		
may add a new flow	yes	#	#		
sends flow_removed message	—	—	—	%	%

where:

- + “modify” and “modify-strict” only take these actions when they create a new flow, not when they update an existing flow.
- % “delete” and “delete-strict” generates a `flow_removed` message if the deleted flow or flows have the `OFPFF_SEND_FLOW_REM` flag set. (Each controller can separately control whether it wants to receive the generated messages.)
- # “modify” and “modify-strict” only add a new flow if the `flow_mod` does not match on any bits of the flow cookie

OpenFlow 1.2

OpenFlow 1.2 makes these changes:

- Only “add” commands ever add flows, “modify” and “modify-strict” never do.
- A new flag `OFPPF_RESET_COUNTS` now controls whether “modify” and “modify-strict” reset counters, whereas previously they never reset counters (except when they inserted a new flow).

RULE	ADD	MODIFY	STRICT	DELETE	STRICT
match on priority	yes	—	yes	—	yes
match on out_port	—	—	—	yes	yes
match on flow_cookie	—	yes	yes	yes	yes
match on table_id	yes	yes	yes	yes	yes
controller chooses table_id	yes	yes	yes		
updates flow_cookie	yes	—	—		
updates <code>OFPPF_SEND_FLOW_REM</code>	yes	—	—		
honors <code>OFPPF_CHECK_OVERLAP</code>	yes	—	—		
updates idle_timeout	yes	—	—		
updates hard_timeout	yes	—	—		
resets idle timer	yes	—	—		
resets hard timer	yes	yes	yes		
zeros counters	yes	&	&		
may add a new flow	yes	—	—		
sends flow_removed message	—	—	—	%	%

% “delete” and “delete-strict” generates a flow_removed message if the deleted flow or flows have the `OFPPF_SEND_FLOW_REM` flag set. (Each controller can separately control whether it wants to receive the generated messages.)

& “modify” and “modify-strict” reset counters if the `OFPPF_RESET_COUNTS` flag is specified.

OpenFlow 1.3

OpenFlow 1.3 makes these changes:

- Behavior for a table_id of 255 is now defined, for “delete” and “delete-strict” commands, as meaning to delete from all tables. A table_id of 255 is now explicitly invalid for other commands.
- New flags `OFPPF_NO_PKT_COUNTS` and `OFPPF_NO_BYT_COUNTS` for “add” operations.

The table for 1.3 is the same as the one shown above for 1.2.

OpenFlow 1.4

OpenFlow 1.4 makes these changes:

- Adds the “importance” field to `flow_mods`, but it does not explicitly specify which kinds of `flow_mods` set the importance. For consistency, Open vSwitch uses the same rule for importance as for `idle_timeout` and `hard_timeout`, that is, only an “ADD” flow_mod sets the importance. (This issue has been filed with the ONF as EXT-496.)
- Eviction Mechanism to automatically delete entries of lower importance to make space for newer entries.

OpenFlow 1.4 Bundles

Open vSwitch makes all flow table modifications atomically, i.e., any datapath packet only sees flow table configurations either before or after any change made by any `flow_mod`. For example, if a controller removes all flows with a

single OpenFlow `flow_mod`, no packet sees an intermediate version of the OpenFlow pipeline where only some of the flows have been deleted.

It should be noted that Open vSwitch caches datapath flows, and that the cached flows are *NOT* flushed immediately when a flow table changes. Instead, the datapath flows are revalidated against the new flow table as soon as possible, and usually within one second of the modification. This design amortizes the cost of datapath cache flushing across multiple flow table changes, and has a significant performance effect during simultaneous heavy flow table churn and high traffic load. This means that different cached datapath flows may have been computed based on a different flow table configurations, but each of the datapath flows is guaranteed to have been computed over a coherent view of the flow tables, as described above.

With OpenFlow 1.4 bundles this atomicity can be extended across an arbitrary set of `flow_mod`. Bundles are supported for `flow_mod` and `port_mod` messages only. For `flow_mod`, both `atomic` and `ordered` bundle flags are trivially supported, as all bundled messages are executed in the order they were added and all flow table modifications are now atomic to the datapath. Port mods may not appear in atomic bundles, as port status modifications are not atomic.

To support bundles, `ovs-ofctl` has a `--bundle` option that makes the flow mod commands (`add-flow`, `add-flows`, `mod-flows`, `del-flows`, and `replace-flows`) use an OpenFlow 1.4 bundle to operate the modifications as a single atomic transaction. If any of the flow mods in a transaction fail, none of them are executed. All flow mods in a bundle appear to datapath lookups simultaneously.

Furthermore, `ovs-ofctl add-flow` and `add-flows` commands now accept arbitrary flow mods as an input by allowing the flow specification to start with an explicit `add`, `modify`, `modify_strict`, `delete`, or `delete_strict` keyword. A missing keyword is treated as `add`, so this is fully backwards compatible. With the new `--bundle` option all the flow mods are executed as a single atomic transaction using an OpenFlow 1.4 bundle. Without the `--bundle` option the flow mods are executed in order up to the first failing `flow_mod`, and in case of an error the earlier successful `flow_mod` calls are not rolled back.

OFPT_PACKET_IN

The OpenFlow 1.1 specification for `OFPT_PACKET_IN` is confusing. The definition in `OF1.1 openflow.h` is[*]:

```
/* Packet received on port (datapath -> controller). */
struct ofp_packet_in {
    struct ofp_header header;
    uint32_t buffer_id;      /* ID assigned by datapath. */
    uint32_t in_port;       /* Port on which frame was received. */
    uint32_t in_phy_port;    /* Physical Port on which frame was received. */
    uint16_t total_len;      /* Full length of frame. */
    uint8_t reason;         /* Reason packet is being sent (one of OFPR_*) */
    uint8_t table_id;       /* ID of the table that was looked up */
    uint8_t data[0];        /* Ethernet frame, halfway through 32-bit word,
                             so the IP header is 32-bit aligned. The
                             amount of data is inferred from the length
                             field in the header. Because of padding,
                             offsetof(struct ofp_packet_in, data) ==
                             sizeof(struct ofp_packet_in) - 2. */
};
OFP_ASSERT(sizeof(struct ofp_packet_in) == 24);
```

The confusing part is the comment on the `data[]` member. This comment is a leftover from `OF1.0 openflow.h`, in which the comment was correct: `sizeof(struct ofp_packet_in)` is 20 in `OF1.0` and `offsetof(struct ofp_packet_in, data)` is 18. When `OF1.1` was written, the structure members were changed but the comment was carelessly not updated, and the comment became wrong: `sizeof(struct ofp_packet_in)` and `offsetof(struct ofp_packet_in, data)` are both 24 in `OF1.1`.

That leaves the question of how to implement `ofp_packet_in` in OF1.1. The OpenFlow reference implementation for OF1.1 does not include any padding, that is, the first byte of the encapsulated frame immediately follows the `table_id` member without a gap. Open vSwitch therefore implements it the same way for compatibility.

For an earlier discussion, please see the thread archived at: <https://mailman.stanford.edu/pipermail/openflow-discuss/2011-August/002604.html>

[*] The quoted definition is directly from OF1.1. Definitions used inside OVS omit the 8-byte `ofp_header` members, so the sizes in this discussion are 8 bytes larger than those declared in OVS header files.

VLAN Matching

The 802.1Q VLAN header causes more trouble than any other 4 bytes in networking. More specifically, three versions of OpenFlow and Open vSwitch have among them four different ways to match the contents and presence of the VLAN header. The following table describes how each version works.

Match	NXM	OF1.0	OF1.1	OF1.2
[1]	0000/0000	????/1, ??/?	????/1, ??/?	0000/0000, --
[2]	0000/ffff	ffff/0, ??/?	ffff/0, ??/?	0000/ffff, --
[3]	1xxx/1fff	0xxx/0, ??/1	0xxx/0, ??/1	1xxx/ffff, --
[4]	z000/f000	????/1, 0y/0	fffe/0, 0y/0	1000/1000, 0y
[5]	zxxx/ffff	0xxx/0, 0y/0	0xxx/0, 0y/0	1xxx/ffff, 0y
[6]	0000/0fff	<none>	<none>	<none>
[7]	0000/f000	<none>	<none>	<none>
[8]	0000/efff	<none>	<none>	<none>
[9]	1001/1001	<none>	<none>	1001/1001, --
[10]	3000/3000	<none>	<none>	<none>
[11]	1000/1000	<none>	fffe/0, ??/1	1000/1000, --

where:

Match: See the list below.

NXM: `xxxx/yyyy` means `NXM_OF_VLAN_TCI_W` with value `xxxx` and mask `yyyy`. A mask of `0000` is equivalent to omitting `NXM_OF_VLAN_TCI (_W)`, a mask of `ffff` is equivalent to `NXM_OF_VLAN_TCI`.

OF1.0, OF1.1: `www/x, yy/z` means `dl_vlan www`, `OFFFW_DL_VLAN x`, `dl_vlan_pcp yy`, and `OFFFW_DL_VLAN_PCP z`. If `OFFFW_DL_VLAN` or `OFFFW_DL_VLAN_PCP` is 1, the corresponding field value is wildcarded, otherwise it is matched. ? means that the given bits are ignored (their conventional values are `0000/x, 00/0` in OF1.0, `0000/x, 00/1` in OF1.1; `x` is never ignored). `<none>` means that the given match is not supported.

OF1.2: `xxxx/yyyy, zz` means `OXM_OF_VLAN_VID_W` with value `xxxx` and mask `yyyy`, and `OXM_OF_VLAN_PCP` (which is not maskable) with value `zz`. A mask of `0000` is equivalent to omitting `OXM_OF_VLAN_VID (_W)`, a mask of `ffff` is equivalent to `OXM_OF_VLAN_VID`. `--` means that `OXM_OF_VLAN_PCP` is omitted. `<none>` means that the given match is not supported.

The matches are:

[1]: Matches any packet, that is, one without an 802.1Q header or with an 802.1Q header with any TCI value.

[2] Matches only packets without an 802.1Q header.

NXM: Any match with `vlan_tci == 0` and `(vlan_tci_mask & 0x1000) != 0` is equivalent to the one listed in the table.

OF1.0: The spec doesn't define behavior if `dl_vlan` is set to `0xffff` and `OFFFW_DL_VLAN_PCP` is not set.

OF1.1: The spec says explicitly to ignore `dl_vlan_pcp` when `dl_vlan` is set to `0xffff`.

OF1.2: The spec doesn't say what should happen if `vlan_vid == 0` and `(vlan_vid_mask & 0x1000) != 0` but `vlan_vid_mask != 0x1000`, but it would be straightforward to also interpret as [2].

[3] Matches only packets that have an 802.1Q header with VID `xxx` (and any PCP).

[4] Matches only packets that have an 802.1Q header with PCP `y` (and any VID).

NXM: `z` is `(y << 1) | 1`.

OF1.0: The spec isn't very clear, but OVS implements it this way.

OF1.2: Presumably other masks such that `(vlan_vid_mask & 0x1fff) == 0x1000` would also work, but the spec doesn't define their behavior.

[5] Matches only packets that have an 802.1Q header with VID `xxx` and PCP `y`.

NXM: `z` is `((y << 1) | 1)`.

OF1.2: Presumably other masks such that `(vlan_vid_mask & 0x1fff) == 0x1fff` would also work.

[6] Matches packets with no 802.1Q header or with an 802.1Q header with a VID of 0. Only possible with NXM.

[7] Matches packets with no 802.1Q header or with an 802.1Q header with a PCP of 0. Only possible with NXM.

[8] Matches packets with no 802.1Q header or with an 802.1Q header with both VID and PCP of 0. Only possible with NXM.

[9] Matches only packets that have an 802.1Q header with an odd-numbered VID (and any PCP). Only possible with NXM and OF1.2. (This is just an example; one can match on any desired VID bit pattern.)

[10] Matches only packets that have an 802.1Q header with an odd-numbered PCP (and any VID). Only possible with NXM. (This is just an example; one can match on any desired VID bit pattern.)

[11] Matches any packet with an 802.1Q header, regardless of VID or PCP.

Additional notes:

OF1.2: The top three bits of `OXM_OF_VLAN_VID` are fixed to zero, so bits 13, 14, and 15 in the masks listed in the table may be set to arbitrary values, as long as the corresponding value bits are also zero. The suggested `ffff` mask for [2], [3], and [5] allows a shorter OXM representation (the mask is omitted) than the minimal `1fff` mask.

Flow Cookies

OpenFlow 1.0 and later versions have the concept of a “flow cookie”, which is a 64-bit integer value attached to each flow. The treatment of the flow cookie has varied greatly across OpenFlow versions, however.

In OpenFlow 1.0:

- `OFPPC_ADD` set the cookie in the flow that it added.
- `OFPPC_MODIFY` and `OFPPC_MODIFY_STRICT` updated the cookie for the flow or flows that it modified.
- `OFPST_FLOW` messages included the flow cookie.
- `OFPT_FLOW_REMOVED` messages reported the cookie of the flow that was removed.

OpenFlow 1.1 made the following changes:

- Flow mod operations `OFFFC_MODIFY`, `OFFFC_MODIFY_STRICT`, `OFFFC_DELETE`, and `OFFFC_DELETE_STRICT`, plus flow stats requests and aggregate stats requests, gained the ability to match on flow cookies with an arbitrary mask.
- `OFFFC_MODIFY` and `OFFFC_MODIFY_STRICT` were changed to add a new flow, in the case of no match, only if the flow table modification operation did not match on the cookie field. (In OpenFlow 1.0, modify operations always added a new flow when there was no match.)
- `OFFFC_MODIFY` and `OFFFC_MODIFY_STRICT` no longer updated flow cookies.

OpenFlow 1.2 made the following changes:

- `OFFFC_MODIFY` and `OFFFC_MODIFY_STRICT` were changed to never add a new flow, regardless of whether the flow cookie was used for matching.

Open vSwitch support for OpenFlow 1.0 implements the OpenFlow 1.0 behavior with the following extensions:

- An NXM extension field `NXM_NX_COOKIE(_W)` allows the NXM versions of `OFFFC_MODIFY`, `OFFFC_MODIFY_STRICT`, `OFFFC_DELETE`, and `OFFFC_DELETE_STRICT` flow_mod calls, plus flow stats requests and aggregate stats requests, to match on flow cookies with arbitrary masks. This is much like the equivalent OpenFlow 1.1 feature.
- Like OpenFlow 1.1, `OFFFC_MODIFY` and `OFFFC_MODIFY_STRICT` add a new flow if there is no match and the mask is zero (or not given).
- The `cookie` field in `OFPT_FLOW_MOD` and `NXT_FLOW_MOD` messages is used as the cookie value for `OFFFC_ADD` commands, as described in OpenFlow 1.0. For `OFFFC_MODIFY` and `OFFFC_MODIFY_STRICT` commands, the `cookie` field is used as a new cookie for flows that match unless it is `UINT64_MAX`, in which case the flow's cookie is not updated.
- `NXT_PACKET_IN` (the Nicira extended version of `OFPT_PACKET_IN`) reports the cookie of the rule that generated the packet, or all-1-bits if no rule generated the packet. (Older versions of OVS used all-0-bits instead of all-1-bits.)

The following table shows the handling of different protocols when receiving `OFFFC_MODIFY` and `OFFFC_MODIFY_STRICT` messages. A mask of 0 indicates either an explicit mask of zero or an implicit one by not specifying the `NXM_NX_COOKIE(_W)` field.

OpenFlow 1.0	no	yes	(add on miss)	(add on miss)
OpenFlow 1.1	yes	no	no	yes
OpenFlow 1.2	yes	no	no	no
NXM	yes	yes*	no	yes

* Updates the flow's cookie unless the `cookie` field is `UINT64_MAX`.

Multiple Table Support

OpenFlow 1.0 has only rudimentary support for multiple flow tables. Notably, OpenFlow 1.0 does not allow the controller to specify the flow table to which a flow is to be added. Open vSwitch adds an extension for this purpose, which is enabled on a per-OpenFlow connection basis using the `NXT_FLOW_MOD_TABLE_ID` message. When the extension is enabled, the upper 8 bits of the `command` member in an `OFPT_FLOW_MOD` or `NXT_FLOW_MOD` message designates the table to which a flow is to be added.

The Open vSwitch software switch implementation offers 255 flow tables. On packet ingress, only the first flow table (table 0) is searched, and the contents of the remaining tables are not considered in any way. Tables other than table 0 only come into play when an `NXAST_RESUBMIT_TABLE` action specifies another table to search.

Tables 128 and above are reserved for use by the switch itself. Controllers should use only tables 0 through 127.

OFPTC_* Table Configuration

This section covers the history of the OFPTC_* table configuration bits across OpenFlow versions.

OpenFlow 1.0 flow tables had fixed configurations.

OpenFlow 1.1 enabled controllers to configure behavior upon flow table miss and added the OFPTC_MISS_* constants for that purpose. OFPTC_* did not control anything else but it was nevertheless conceptualized as a set of bit-fields instead of an enum. OF1.1 added the OFPT_TABLE_MOD message to set OFPTC_MISS_* for a flow table and added the config field to the OFPST_TABLE reply to report the current setting.

OpenFlow 1.2 did not change anything in this regard.

OpenFlow 1.3 switched to another means to changing flow table miss behavior and deprecated OFPTC_MISS_* without adding any more OFPTC_* constants. This meant that OFPT_TABLE_MOD now had no purpose at all, but OF1.3 kept it around “for backward compatibility with older and newer versions of the specification.” At the same time, OF1.3 introduced a new message OFPMP_TABLE_FEATURES that included a field config documented as reporting the OFPTC_* values set with OFPT_TABLE_MOD; of course this served no real purpose because no OFPTC_* values are defined. OF1.3 did remove the OFPTC_* field from OFPMP_TABLE (previously named OFPST_TABLE).

OpenFlow 1.4 defined two new OFPTC_* constants, OFPTC_EVICTION and OFPTC_VACANCY_EVENTS, using bits that did not overlap with OFPTC_MISS_* even though those bits had not been defined since OF1.2. OFPT_TABLE_MOD still controlled these settings. The field for OFPTC_* values in OFPMP_TABLE_FEATURES was renamed from config to capabilities and documented as reporting the flags that are supported in a OFPT_TABLE_MOD message. The OFPMP_TABLE_DESC message newly added in OF1.4 reported the OFPTC_* setting.

OpenFlow 1.5 did not change anything in this regard.

Table 7: Revisions

Open-Flow	OFPTC_* flags	TABLE_MOD	Statistics	TABLE_FEATURES	TABLE_DESC
OF1.0	none	no (*) (+)	no (*)	nothing (*) (+)	no (*) (+)
OF1.1/1.2	MISS_*	yes	yes	nothing (+)	no (+)
OF1.3	none	yes (*)	no (*)	config (*)	no (*) (+)
OF1.4/1.5	EVICTION/VACANCY_EVENTS	yes	no	capabilities	yes

where:

OpenFlow: The OpenFlow version(s).

OFPTC_* flags: The OFPTC_* flags defined in those versions.

TABLE_MOD: Whether OFPT_TABLE_MOD can modify OFPTC_* flags.

Statistics: Whether OFPST_TABLE/OFPMP_TABLE reports the OFPTC_* flags.

TABLE_FEATURES: What OFPMP_TABLE_FEATURES reports (if it exists): either the current configuration or the switch’s capabilities.

TABLE_DESC: Whether OFPMP_TABLE_DESC reports the current configuration.

(*): Nothing to report/change anyway.

(+): No such message.

IPv6

Open vSwitch supports stateless handling of IPv6 packets. Flows can be written to support matching TCP, UDP, and ICMPv6 headers within an IPv6 packet. Deeper matching of some Neighbor Discovery messages is also supported.

IPv6 was not designed to interact well with middle-boxes. This, combined with Open vSwitch's stateless nature, have affected the processing of IPv6 traffic, which is detailed below.

Extension Headers

The base IPv6 header is incredibly simple with the intention of only containing information relevant for routing packets between two endpoints. IPv6 relies heavily on the use of extension headers to provide any other functionality. Unfortunately, the extension headers were designed in such a way that it is impossible to move to the next header (including the layer-4 payload) unless the current header is understood.

Open vSwitch will process the following extension headers and continue to the next header:

- Fragment (see the next section)
- AH (Authentication Header)
- Hop-by-Hop Options
- Routing
- Destination Options

When a header is encountered that is not in that list, it is considered “terminal”. A terminal header's IPv6 protocol value is stored in `nw_proto` for matching purposes. If a terminal header is TCP, UDP, or ICMPv6, the packet will be further processed in an attempt to extract layer-4 information.

Fragments

IPv6 requires that every link in the internet have an MTU of 1280 octets or greater (RFC 2460). As such, a terminal header (as described above in “Extension Headers”) in the first fragment should generally be reachable. In this case, the terminal header's IPv6 protocol type is stored in the `nw_proto` field for matching purposes. If a terminal header cannot be found in the first fragment (one with a fragment offset of zero), the `nw_proto` field is set to 0. Subsequent fragments (those with a non-zero fragment offset) have the `nw_proto` field set to the IPv6 protocol type for fragments (44).

Jumbograms

An IPv6 jumbogram (RFC 2675) is a packet containing a payload longer than 65,535 octets. A jumbogram is only relevant in subnets with a link MTU greater than 65,575 octets, and are not required to be supported on nodes that do not connect to link with such large MTUs. Currently, Open vSwitch doesn't process jumbograms.

In-Band Control

Motivation

An OpenFlow switch must establish and maintain a TCP network connection to its controller. There are two basic ways to categorize the network that this connection traverses: either it is completely separate from the one that the switch is otherwise controlling, or its path may overlap the network that the switch controls. We call the former case “out-of-band control”, the latter case “in-band control”.

Out-of-band control has the following benefits:

- **Simplicity:** Out-of-band control slightly simplifies the switch implementation.
- **Reliability:** Excessive switch traffic volume cannot interfere with control traffic.
- **Integrity:** Machines not on the control network cannot impersonate a switch or a controller.
- **Confidentiality:** Machines not on the control network cannot snoop on control traffic.

In-band control, on the other hand, has the following advantages:

- **No dedicated port:** There is no need to dedicate a physical switch port to control, which is important on switches that have few ports (e.g. wireless routers, low-end embedded platforms).
- **No dedicated network:** There is no need to build and maintain a separate control network. This is important in many environments because it reduces proliferation of switches and wiring.

Open vSwitch supports both out-of-band and in-band control. This section describes the principles behind in-band control. See the description of the Controller table in `ovs-vsitchd.conf.db(5)` to configure OVS for in-band control.

Principles

The fundamental principle of in-band control is that an OpenFlow switch must recognize and switch control traffic without involving the OpenFlow controller. All the details of implementing in-band control are special cases of this principle.

The rationale for this principle is simple. If the switch does not handle in-band control traffic itself, then it will be caught in a contradiction: it must contact the controller, but it cannot, because only the controller can set up the flows that are needed to contact the controller.

The following points describe important special cases of this principle.

- In-band control must be implemented regardless of whether the switch is connected.

It is tempting to implement the in-band control rules only when the switch is not connected to the controller, using the reasoning that the controller should have complete control once it has established a connection with the switch.

This does not work in practice. Consider the case where the switch is connected to the controller. Occasionally it can happen that the controller forgets or otherwise needs to obtain the MAC address of the switch. To do so, the controller sends a broadcast ARP request. A switch that implements the in-band control rules only when it is disconnected will then send an `OFPT_PACKET_IN` message up to the controller. The controller will be unable to respond, because it does not know the MAC address of the switch. This is a deadlock situation that can only be resolved by the switch noticing that its connection to the controller has hung and reconnecting.

- In-band control must override flows set up by the controller.

It is reasonable to assume that flows set up by the OpenFlow controller should take precedence over in-band control, on the basis that the controller should be in charge of the switch.

Again, this does not work in practice. Reasonable controller implementations may set up a “last resort” fallback rule that wildcards every field and, e.g., sends it up to the controller or discards it. If a controller does that, then it will isolate itself from the switch.

- The switch must recognize all control traffic.

The fundamental principle of in-band control states, in part, that a switch must recognize control traffic without involving the OpenFlow controller. More specifically, the switch must recognize *all* control traffic. “False negatives”, that is, packets that constitute control traffic but that the switch does not recognize as control traffic, lead to control traffic storms.

Consider an OpenFlow switch that only recognizes control packets sent to or from that switch. Now suppose that two switches of this type, named A and B, are connected to ports on an Ethernet hub (not a switch) and that an OpenFlow controller is connected to a third hub port. In this setup, control traffic sent by switch A will be seen by switch B, which will send it to the controller as part of an `OFPT_PACKET_IN` message. Switch A will then see the `OFPT_PACKET_IN` message's packet, re-encapsulate it in another `OFPT_PACKET_IN`, and send it to the controller. Switch B will then see that `OFPT_PACKET_IN`, and so on in an infinite loop.

Incidentally, the consequences of “false positives”, where packets that are not control traffic are nevertheless recognized as control traffic, are much less severe. The controller will not be able to control their behavior, but the network will remain in working order. False positives do constitute a security problem.

- The switch should use echo-requests to detect disconnection.

TCP will notice that a connection has hung, but this can take a considerable amount of time. For example, with default settings the Linux kernel TCP implementation will retransmit for between 13 and 30 minutes, depending on the connection's retransmission timeout, according to kernel documentation. This is far too long for a switch to be disconnected, so an OpenFlow switch should implement its own connection timeout. OpenFlow `OFPT_ECHO_REQUEST` messages are the best way to do this, since they test the OpenFlow connection itself.

Implementation

This section describes how Open vSwitch implements in-band control. Correctly implementing in-band control has proven difficult due to its many subtleties, and has thus gone through many iterations. Please read through and understand the reasoning behind the chosen rules before making modifications.

Open vSwitch implements in-band control as “hidden” flows, that is, flows that are not visible through OpenFlow, and at a higher priority than wildcarded flows can be set up through OpenFlow. This is done so that the OpenFlow controller cannot interfere with them and possibly break connectivity with its switches. It is possible to see all flows, including in-band ones, with the `ovs-appctl “bridge/dump-flows”` command.

The Open vSwitch implementation of in-band control can hide traffic to arbitrary “remotes”, where each remote is one TCP port on one IP address. Currently the remotes are automatically configured as the in-band OpenFlow controllers plus the OVSDB managers, if any. (The latter is a requirement because OVSDB managers are responsible for configuring OpenFlow controllers, so if the manager cannot be reached then OpenFlow cannot be reconfigured.)

The following rules (with the `OFPP_NORMAL` action) are set up on any bridge that has any remotes:

1. DHCP requests sent from the local port.
2. ARP replies to the local port's MAC address.
3. ARP requests from the local port's MAC address.

In-band also sets up the following rules for each unique next-hop MAC address for the remotes' IPs (the “next hop” is either the remote itself, if it is on a local subnet, or the gateway to reach the remote):

4. ARP replies to the next hop's MAC address.
5. ARP requests from the next hop's MAC address.

In-band also sets up the following rules for each unique remote IP address:

6. ARP replies containing the remote's IP address as a target.
7. ARP requests containing the remote's IP address as a source.

In-band also sets up the following rules for each unique remote (IP,port) pair:

8. TCP traffic to the remote's IP and port.
9. TCP traffic from the remote's IP and port.

The goal of these rules is to be as narrow as possible to allow a switch to join a network and be able to communicate with the remotes. As mentioned earlier, these rules have higher priority than the controller's rules, so if they are too broad, they may prevent the controller from implementing its policy. As such, in-band actively monitors some aspects of flow and packet processing so that the rules can be made more precise.

In-band control monitors attempts to add flows into the datapath that could interfere with its duties. The datapath only allows exact match entries, so in-band control is able to be very precise about the flows it prevents. Flows that miss in the datapath are sent to userspace to be processed, so preventing these flows from being cached in the "fast path" does not affect correctness. The only type of flow that is currently prevented is one that would prevent DHCP replies from being seen by the local port. For example, a rule that forwarded all DHCP traffic to the controller would not be allowed, but one that forwarded to all ports (including the local port) would.

As mentioned earlier, packets that miss in the datapath are sent to the userspace for processing. The userspace has its own flow table, the "classifier", so in-band checks whether any special processing is needed before the classifier is consulted. If a packet is a DHCP response to a request from the local port, the packet is forwarded to the local port, regardless of the flow table. Note that this requires L7 processing of DHCP replies to determine whether the 'chaddr' field matches the MAC address of the local port.

It is interesting to note that for an L3-based in-band control mechanism, the majority of rules are devoted to ARP traffic. At first glance, some of these rules appear redundant. However, each serves an important role. First, in order to determine the MAC address of the remote side (controller or gateway) for other ARP rules, we must allow ARP traffic for our local port with rules (b) and (c). If we are between a switch and its connection to the remote, we have to allow the other switch's ARP traffic to through. This is done with rules (d) and (e), since we do not know the addresses of the other switches a priori, but do know the remote's or gateway's. Finally, if the remote is running in a local guest VM that is not reached through the local port, the switch that is connected to the VM must allow ARP traffic based on the remote's IP address, since it will not know the MAC address of the local port that is sending the traffic or the MAC address of the remote in the guest VM.

With a few notable exceptions below, in-band should work in most network setups. The following are considered "supported" in the current implementation:

- Locally Connected. The switch and remote are on the same subnet. This uses rules (a), (b), (c), (h), and (i).
- Reached through Gateway. The switch and remote are on different subnets and must go through a gateway. This uses rules (a), (b), (c), (h), and (i).
- Between Switch and Remote. This switch is between another switch and the remote, and we want to allow the other switch's traffic through. This uses rules (d), (e), (h), and (i). It uses (b) and (c) indirectly in order to know the MAC address for rules (d) and (e). Note that DHCP for the other switch will not work unless an OpenFlow controller explicitly lets this switch pass the traffic.
- Between Switch and Gateway. This switch is between another switch and the gateway, and we want to allow the other switch's traffic through. This uses the same rules and logic as the "Between Switch and Remote" configuration described earlier.
- Remote on Local VM. The remote is a guest VM on the system running in-band control. This uses rules (a), (b), (c), (h), and (i).
- Remote on Local VM with Different Networks. The remote is a guest VM on the system running in-band control, but the local port is not used to connect to the remote. For example, an IP address is configured on eth0 of the switch. The remote's VM is connected through eth1 of the switch, but an IP address has not been configured for that port on the switch. As such, the switch will use eth0 to connect to the remote, and eth1's rules about the local port will not work. In the example, the switch attached to eth0 would use rules (a), (b), (c), (h), and (i) on eth0. The switch attached to eth1 would use rules (f), (g), (h), and (i).

The following are explicitly *not* supported by in-band control:

- Specify Remote by Name. Currently, the remote must be identified by IP address. A naive approach would be to permit all DNS traffic. Unfortunately, this would prevent the controller from defining any policy over DNS. Since switches that are located behind us need to connect to the remote, in-band cannot simply add a rule that

allows DNS traffic from the local port. The “correct” way to support this is to parse DNS requests to allow all traffic related to a request for the remote’s name through. Due to the potential security problems and amount of processing, we decided to hold off for the time-being.

- Differing Remotes for Switches. All switches must know the L3 addresses for all the remotes that other switches may use, since rules need to be set up to allow traffic related to those remotes through. See rules (f), (g), (h), and (i).
- Differing Routes for Switches. In order for the switch to allow other switches to connect to a remote through a gateway, it allows the gateway’s traffic through with rules (d) and (e). If the routes to the remote differ for the two switches, we will not know the MAC address of the alternate gateway.

Action Reproduction

It seems likely that many controllers, at least at startup, use the OpenFlow “flow statistics” request to obtain existing flows, then compare the flows’ actions against the actions that they expect to find. Before version 1.8.0, Open vSwitch always returned exact, byte-for-byte copies of the actions that had been added to the flow table. The current version of Open vSwitch does not always do this in some exceptional cases. This section lists the exceptions that controller authors must keep in mind if they compare actual actions against desired actions in a bitwise fashion:

- Open vSwitch zeros padding bytes in action structures, regardless of their values when the flows were added.
- Open vSwitch “normalizes” the instructions in OpenFlow 1.1 (and later) in the following way:
 - OVS sorts the instructions into the following order: Apply-Actions, Clear-Actions, Write-Actions, Write-Metadata, Goto-Table.
 - OVS drops Apply-Actions instructions that have empty action lists.
 - OVS drops Write-Actions instructions that have empty action sets.

Please report other discrepancies, if you notice any, so that we can fix or document them.

Suggestions

Suggestions to improve Open vSwitch are welcome at discuss@openvswitch.org.

4.1.2 Open vSwitch Datapath Development Guide

The Open vSwitch kernel module allows flexible userspace control over flow-level packet processing on selected network devices. It can be used to implement a plain Ethernet switch, network device bonding, VLAN processing, network access control, flow-based network control, and so on.

The kernel module implements multiple “datapaths” (analogous to bridges), each of which can have multiple “vports” (analogous to ports within a bridge). Each datapath also has associated with it a “flow table” that userspace populates with “flows” that map from keys based on packet headers and metadata to sets of actions. The most common action forwards the packet to another vport; other actions are also implemented.

When a packet arrives on a vport, the kernel module processes it by extracting its flow key and looking it up in the flow table. If there is a matching flow, it executes the associated actions. If there is no match, it queues the packet to userspace for processing (as part of its processing, userspace will likely set up a flow to handle further packets of the same type entirely in-kernel).

Flow Key Compatibility

Network protocols evolve over time. New protocols become important and existing protocols lose their prominence. For the Open vSwitch kernel module to remain relevant, it must be possible for newer versions to parse additional protocols as part of the flow key. It might even be desirable, someday, to drop support for parsing protocols that have become obsolete. Therefore, the Netlink interface to Open vSwitch is designed to allow carefully written userspace applications to work with any version of the flow key, past or future.

To support this forward and backward compatibility, whenever the kernel module passes a packet to userspace, it also passes along the flow key that it parsed from the packet. Userspace then extracts its own notion of a flow key from the packet and compares it against the kernel-provided version:

- If userspace's notion of the flow key for the packet matches the kernel's, then nothing special is necessary.
- If the kernel's flow key includes more fields than the userspace version of the flow key, for example if the kernel decoded IPv6 headers but userspace stopped at the Ethernet type (because it does not understand IPv6), then again nothing special is necessary. Userspace can still set up a flow in the usual way, as long as it uses the kernel-provided flow key to do it.
- If the userspace flow key includes more fields than the kernel's, for example if userspace decoded an IPv6 header but the kernel stopped at the Ethernet type, then userspace can forward the packet manually, without setting up a flow in the kernel. This case is bad for performance because every packet that the kernel considers part of the flow must go to userspace, but the forwarding behavior is correct. (If userspace can determine that the values of the extra fields would not affect forwarding behavior, then it could set up a flow anyway.)

How flow keys evolve over time is important to making this work, so the following sections go into detail.

Flow Key Format

A flow key is passed over a Netlink socket as a sequence of Netlink attributes. Some attributes represent packet metadata, defined as any information about a packet that cannot be extracted from the packet itself, e.g. the vport on which the packet was received. Most attributes, however, are extracted from headers within the packet, e.g. source and destination addresses from Ethernet, IP, or TCP headers.

The `<linux/openvswitch.h>` header file defines the exact format of the flow key attributes. For informal explanatory purposes here, we write them as comma-separated strings, with parentheses indicating arguments and nesting. For example, the following could represent a flow key corresponding to a TCP packet that arrived on vport 1:

```
in_port(1), eth(src=e0:91:f5:21:d0:b2, dst=00:02:e3:0f:80:a4),
eth_type(0x0800), ipv4(src=172.16.0.20, dst=172.18.0.52, proto=6, tos=0,
frag=no), tcp(src=49163, dst=80)
```

Often we ellipsize arguments not important to the discussion, e.g.:

```
in_port(1), eth(...), eth_type(0x0800), ipv4(...), tcp(...)
```

Wildcarded Flow Key Format

A wildcarded flow is described with two sequences of Netlink attributes passed over the Netlink socket. A flow key, exactly as described above, and an optional corresponding flow mask.

A wildcarded flow can represent a group of exact match flows. Each 1 bit in the mask specifies an exact match with the corresponding bit in the flow key. A 0 bit specifies a don't care bit, which will match either a 1 or 0 bit of an incoming packet. Using a wildcarded flow can improve the flow set up rate by reducing the number of new flows that need to be processed by the user space program.

Support for the mask Netlink attribute is optional for both the kernel and user space program. The kernel can ignore the mask attribute, installing an exact match flow, or reduce the number of don't care bits in the kernel to less than what was specified by the user space program. In this case, variations in bits that the kernel does not implement will simply result in additional flow setups. The kernel module will also work with user space programs that neither support nor supply flow mask attributes.

Since the kernel may ignore or modify wildcard bits, it can be difficult for the userspace program to know exactly what matches are installed. There are two possible approaches: reactively install flows as they miss the kernel flow table (and therefore not attempt to determine wildcard changes at all) or use the kernel's response messages to determine the installed wildcards.

When interacting with userspace, the kernel should maintain the match portion of the key exactly as originally installed. This will provide a handle to identify the flow for all future operations. However, when reporting the mask of an installed flow, the mask should include any restrictions imposed by the kernel.

The behavior when using overlapping wildcarded flows is undefined. It is the responsibility of the user space program to ensure that any incoming packet can match at most one flow, wildcarded or not. The current implementation performs best-effort detection of overlapping wildcarded flows and may reject some but not all of them. However, this behavior may change in future versions.

Unique Flow Identifiers

An alternative to using the original match portion of a key as the handle for flow identification is a unique flow identifier, or "UFID". UFIDs are optional for both the kernel and user space program.

User space programs that support UFID are expected to provide it during flow setup in addition to the flow, then refer to the flow using the UFID for all future operations. The kernel is not required to index flows by the original flow key if a UFID is specified.

Basic Rule for Evolving Flow Keys

Some care is needed to really maintain forward and backward compatibility for applications that follow the rules listed under "Flow key compatibility" above.

The basic rule is obvious:

New network protocol support must only supplement existing flow key attributes. It must not change the meaning of already defined flow key attributes.

This rule does have less-obvious consequences so it is worth working through a few examples. Suppose, for example, that the kernel module did not already implement VLAN parsing. Instead, it just interpreted the 802.1Q TPID (0x8100) as the Ethertype then stopped parsing the packet. The flow key for any packet with an 802.1Q header would look essentially like this, ignoring metadata:

```
eth(...), eth_type(0x8100)
```

Naively, to add VLAN support, it makes sense to add a new "vlan" flow key attribute to contain the VLAN tag, then continue to decode the encapsulated headers beyond the VLAN tag using the existing field definitions. With this change, a TCP packet in VLAN 10 would have a flow key much like this:

```
eth(...), vlan(vid=10, pcp=0), eth_type(0x0800), ip(proto=6, ...), tcp(...)
```

But this change would negatively affect a userspace application that has not been updated to understand the new "vlan" flow key attribute. The application could, following the flow compatibility rules above, ignore the "vlan" attribute that it does not understand and therefore assume that the flow contained IP packets. This is a bad assumption (the flow only contains IP packets if one parses and skips over the 802.1Q header) and it could cause the application's behavior to change across kernel versions even though it follows the compatibility rules.

The solution is to use a set of nested attributes. This is, for example, why 802.1Q support uses nested attributes. A TCP packet in VLAN 10 is actually expressed as:

```
eth(...), eth_type(0x8100), vlan(vid=10, pcp=0), encap(eth_type(0x0800),
ip(proto=6, ...), tcp(...)))
```

Notice how the `eth_type`, `ip`, and `tcp` flow key attributes are nested inside the `encap` attribute. Thus, an application that does not understand the `vlan` key will not see either of those attributes and therefore will not misinterpret them. (Also, the outer `eth_type` is still `0x8100`, not changed to `0x0800`)

Handling Malformed Packets

Don't drop packets in the kernel for malformed protocol headers, bad checksums, etc. This would prevent userspace from implementing a simple Ethernet switch that forwards every packet.

Instead, in such a case, include an attribute with "empty" content. It doesn't matter if the empty content could be valid protocol values, as long as those values are rarely seen in practice, because userspace can always forward all packets with those values to userspace and handle them individually.

For example, consider a packet that contains an IP header that indicates protocol 6 for TCP, but which is truncated just after the IP header, so that the TCP header is missing. The flow key for this packet would include a `tcp` attribute with all-zero `src` and `dst`, like this:

```
eth(...), eth_type(0x0800), ip(proto=6, ...), tcp(src=0, dst=0)
```

As another example, consider a packet with an Ethernet type of `0x8100`, indicating that a VLAN TCI should follow, but which is truncated just after the Ethernet type. The flow key for this packet would include an all-zero-bits `vlan` and an empty `encap` attribute, like this:

```
eth(...), eth_type(0x8100), vlan(0), encap()
```

Unlike a TCP packet with source and destination ports 0, an all-zero-bits VLAN TCI is not that rare, so the CFI bit (aka `VLAN_TAG_PRESENT` inside the kernel) is ordinarily set in a `vlan` attribute expressly to allow this situation to be distinguished. Thus, the flow key in this second example unambiguously indicates a missing or malformed VLAN TCI.

Other Rules

The other rules for flow keys are much less subtle:

- Duplicate attributes are not allowed at a given nesting level.
- Ordering of attributes is not significant.
- When the kernel sends a given flow key to userspace, it always composes it the same way. This allows userspace to hash and compare entire flow keys that it may not be able to fully interpret.

Coding Rules

Implement the headers and codes for compatibility with older kernel in `linux/compat/` directory. All public functions should be exported using `EXPORT_SYMBOL` macro. Public function replacing the same-named kernel function should be prefixed with `rpl_`. Otherwise, the function should be prefixed with `ovs_`. For special case when it is not possible to follow this rule (e.g., the `pskb_expand_head()` function), the function name must be added to `linux/compat/build-aux/export-check-whitelist`, otherwise, the compilation check `check-export-symbol` will fail.

4.1.3 Integration Guide for Centralized Control

This document describes how to integrate Open vSwitch onto a new platform to expose the state of the switch and attached devices for centralized control. (If you are looking to port the switching components of Open vSwitch to a new platform, refer to *Porting Open vSwitch to New Software or Hardware*) The focus of this guide is on hypervisors, but many of the interfaces are useful for hardware switches, as well. The XenServer integration is the most mature implementation, so most of the examples are drawn from it.

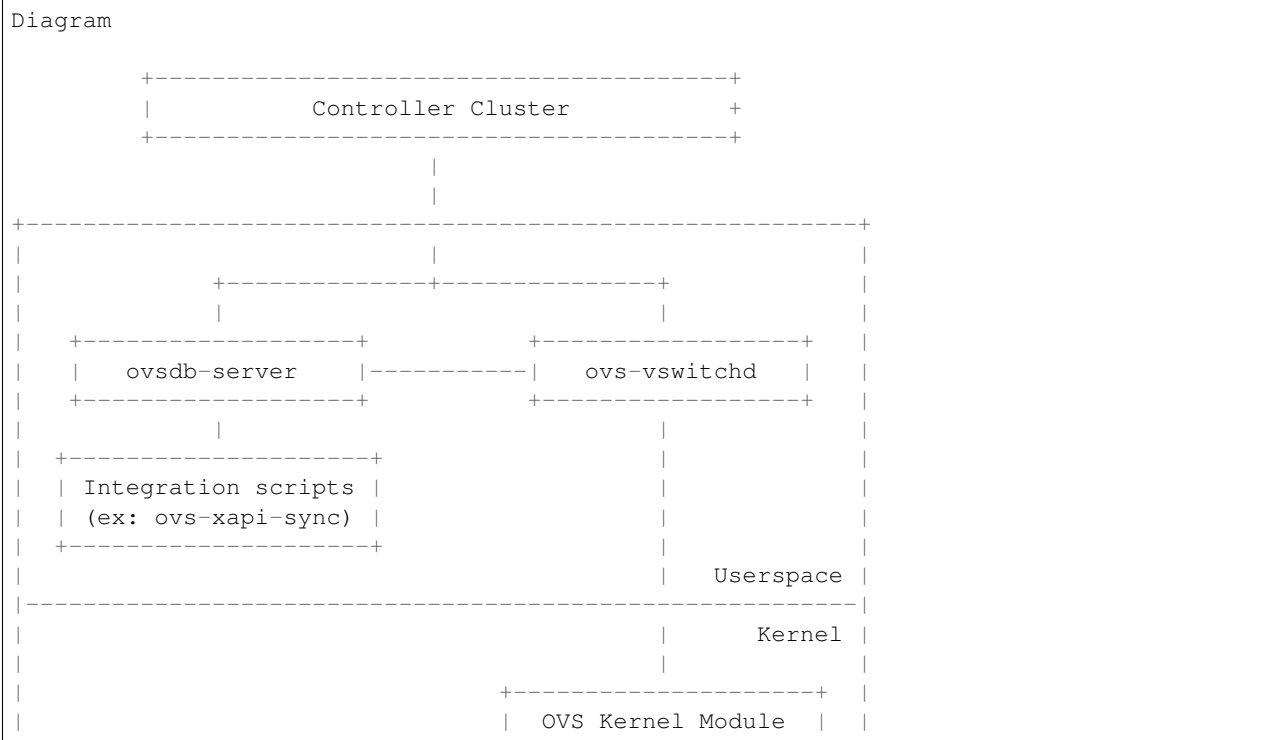
The externally visible interface to this integration is platform-agnostic. We encourage anyone who integrates Open vSwitch to use the same interface, because keeping a uniform interface means that controllers require less customization for individual platforms (and perhaps no customization at all).

Integration centers around the Open vSwitch database and mostly involves the `external_ids` columns in several of the tables. These columns are not interpreted by Open vSwitch itself. Instead, they provide information to a controller that permits it to associate a database record with a more meaningful entity. In contrast, the `other_config` column is used to configure behavior of the switch. The main job of the integrator, then, is to ensure that these values are correctly populated and maintained.

An integrator sets the columns in the database by talking to the `ovsdb-server` daemon. A few of the columns can be set during startup by calling the `ovs-ctl` tool from inside the startup scripts. The `xenserver/etc_init.d/openvswitch` script provides examples of its use, and the `ovs-ctl(8)` manpage contains complete documentation. At runtime, `ovs-vsctl` can be used to set columns in the database. The script `xenserver/etc_xensource_scripts_vif` contains examples of its use, and `ovs-vsctl(8)` manpage contains complete documentation.

Python and C bindings to the database are provided if deeper integration with a program are needed. The XenServer `ovs-xapi-sync` daemon (`xenserver/usr_share_openvswitch_scripts_ovs-xapi-sync`) provides an example of using the Python bindings. More information on the python bindings is available at `python/ovs/db/idl.py`. Information on the C bindings is available at `lib/ovsdb-idl.h`.

The following diagram shows how integration scripts fit into the Open vSwitch architecture:



(continues on next page)

(continued from previous page)

	+-----+	
+-----+		+-----+

A description of the most relevant fields for integration follows. By setting these values, controllers are able to understand the network and manage it more dynamically and precisely. For more details about the database and each individual column, please refer to the `ovs-vswitchd.conf.db(5)` manpage.

Open_vSwitch table

The `Open_vSwitch` table describes the switch as a whole. The `system_type` and `system_version` columns identify the platform to the controller. The `external_ids:system-id` key uniquely identifies the physical host. In XenServer, the `system-id` will likely be the same as the UUID returned by `xe host-list`. This key allows controllers to distinguish between multiple hypervisors.

Most of this configuration can be done with the `ovs-ctl` command at startup. For example:

```
$ ovs-ctl --system-type="XenServer" --system-version="6.0.0-50762p" \
  --system-id="${UUID}" "${other_options}" start
```

Alternatively, the `ovs-vsctl` command may be used to set a particular value at runtime. For example:

```
$ ovs-vsctl set open_vswitch . external_ids:system-id="${UUID}"
```

The `other_config:enable-statistics` key may be set to `true` to have OVS populate the database with statistics (e.g., number of CPUs, memory, system load) for the controller's use.

Bridge table

The `Bridge` table describes individual bridges within an Open vSwitch instance. The `external_ids:bridge-id` key uniquely identifies a particular bridge. In XenServer, this will likely be the same as the UUID returned by `xe network-list` for that particular bridge.

For example, to set the identifier for bridge “br0”, the following command can be used:

```
$ ovs-vsctl set Bridge br0 external_ids:bridge-id="${UUID}"
```

The MAC address of the bridge may be manually configured by setting it with the `other_config:hwaddr` key. For example:

```
$ ovs-vsctl set Bridge br0 other_config:hwaddr="12:34:56:78:90:ab"
```

Interface table

The `Interface` table describes an interface under the control of Open vSwitch. The `external_ids` column contains keys that are used to provide additional information about the interface:

`attached-mac`

This field contains the MAC address of the device attached to the interface. On a hypervisor, this is the MAC address of the interface as seen inside a VM. It does not necessarily correlate to the host-side MAC address. For example, on XenServer, the MAC address on a VIF in the hypervisor is always `FE:FF:FF:FF:FF:FF`, but inside the VM a normal MAC address is seen.

`iface-id`

This field uniquely identifies the interface. In hypervisors, this allows the controller to follow VM network interfaces as VMs migrate. A well-chosen identifier should also allow an administrator or a controller to associate the interface with the corresponding object in the VM management system. For example, the Open vSwitch integration with XenServer by default uses the XenServer assigned UUID for a VIF record as the `iface-id`.

`iface-status`

In a hypervisor, there are situations where there are multiple interface choices for a single virtual ethernet interface inside a VM. Valid values are “active” and “inactive”. A complete description is available in the `ovs-vswitchd.conf.db(5)` manpage.

`vm-id`

This field uniquely identifies the VM to which this interface belongs. A single VM may have multiple interfaces attached to it.

As in the previous tables, the `ovs-vsctl` command may be used to configure the values. For example, to set the `iface-id` on `eth0`, the following command can be used:

```
$ ovs-vsctl set Interface eth0 external-ids:iface-id='${UUID}'
```

HA for OVN DB servers using pacemaker

The `ovsdb` servers can work in either active or backup mode. In backup mode, db server will be connected to an active server and replicate the active servers contents. At all times, the data can be transacted only from the active server. When the active server dies for some reason, entire OVN operations will be stalled.

Pacemaker is a cluster resource manager which can manage a defined set of resource across a set of clustered nodes. Pacemaker manages the resource with the help of the resource agents. One among the resource agent is **OCF**

OCF is nothing but a shell script which accepts a set of actions and returns an appropriate status code.

With the help of the OCF resource agent `ovn/utilities/ovndb-servers.ocf`, one can defined a resource for the pacemaker such that pacemaker will always maintain one running active server at any time.

After creating a pacemaker cluster, use the following commands to create one active and multiple backup servers for OVN databases:

```
$ pcs resource create ovndb_servers ocf:ovn:ovndb-servers \
  master_ip=x.x.x.x \
  ovn_ctl=<path of the ovn-ctl script> \
  op monitor interval="10s" \
  op monitor role=Master interval="15s"
$ pcs resource master ovndb_servers-master ovndb_servers \
  meta notify="true"
```

The `master_ip` and `ovn_ctl` are the parameters that will be used by the OCF script. `ovn_ctl` is optional, if not given, it assumes a default value of `/usr/share/openvswitch/scripts/ovn-ctl`. `master_ip` is the IP address on which the active database server is expected to be listening, the slave node uses it to connect to the master node. You can add the optional parameters ‘`nb_master_port`’, ‘`nb_master_protocol`’, ‘`sb_master_port`’, ‘`sb_master_protocol`’ to set the protocol and port.

Whenever the active server dies, pacemaker is responsible to promote one of the backup servers to be active. Both `ovn-controller` and `ovn-northd` needs the ip-address at which the active server is listening. With pacemaker changing the node at which the active server is run, it is not efficient to instruct all the `ovn-controllers` and the `ovn-northd` to listen to the latest active server’s ip-address.

This problem can be solved by two ways:

1. By using a native ocf resource agent `ocf:heartbeat:IPAddr2`. The `IPAddr2` resource agent is just a resource with an ip-address. When we colocate this resource with the active server, pacemaker will enable the active server to be connected with a single ip-address all the time. This is the ip-address that needs to be given as the parameter while creating the `ovndb_servers` resource.

Use the following command to create the `IPAddr2` resource and colocate it with the active server:

```
$ pcs resource create VirtualIP ocf:heartbeat:IPAddr2 ip=x.x.x.x \
    op monitor interval=30s
$ pcs constraint order promote ovndb_servers-master then VirtualIP
$ pcs constraint colocation add VirtualIP with master ovndb_servers-master \
    score=INFINITY
```

2. Using load balancer vip ip as a master_ip. In order to use this feature, one needs to use `listen_on_master_ip_only` to no. Current code for load balancer have been tested to work with tcp protocol and needs to be tested/enhanced for ssl. Using load balancer, standby nodes will not listen on nb and sb db ports so that load balancer will always communicate to the active node and all the traffic will be sent to active node only. Standby will continue to sync using LB VIP IP in this case.

Use the following command to create pcs resource using LB VIP IP:

```
$ pcs resource create ovndb_servers ocf:ovn:ovndb-servers \
    master_ip="<load_balance_vip_ip>" \
    listen_on_master_ip_only="no" \
    ovn_ctl=<path of the ovn-ctl script> \
    op monitor interval="10s" \
    op monitor role=Master interval="15s"
$ pcs resource master ovndb_servers-master ovndb_servers \
    meta notify="true"
```

4.1.4 Porting Open vSwitch to New Software or Hardware

Open vSwitch (OVS) is intended to be easily ported to new software and hardware platforms. This document describes the types of changes that are most likely to be necessary in porting OVS to Unix-like platforms. (Porting OVS to other kinds of platforms is likely to be more difficult.)

Vocabulary

For historical reasons, different words are used for essentially the same concept in different areas of the Open vSwitch source tree. Here is a concordance, indexed by the area of the source tree:

datapath/	vport	---
vswitchd/	iface	port
ofproto/	port	bundle
ofproto/bond.c	slave	bond
lib/lacp.c	slave	lacp
lib/netdev.c	netdev	---
database	Interface	Port

Open vSwitch Architectural Overview

The following diagram shows the very high-level architecture of Open vSwitch from a porter's perspective.

```
+-----+
|   ovs-vswitchd   |<-->ovsdb-server
+-----+
|   ofproto        |<-->OpenFlow controllers
+-----+
| netdev | | ofproto |
+-----+ |provider|
| netdev | +-----+
|provider|
+-----+
```

Some of the components are generic. Modulo bugs or inadequacies, these components should not need to be modified as part of a port:

ovs-vswitchd The main Open vSwitch userspace program, in `vswitchd/`. It reads the desired Open vSwitch configuration from the `ovsdb-server` program over an IPC channel and passes this configuration down to the “ofproto” library. It also passes certain status and statistical information from `ofproto` back into the database.

ofproto The Open vSwitch library, in `ofproto/`, that implements an OpenFlow switch. It talks to OpenFlow controllers over the network and to switch hardware or software through an “ofproto provider”, explained further below.

netdev The Open vSwitch library, in `lib/netdev.c`, that abstracts interacting with network devices, that is, Ethernet interfaces. The `netdev` library is a thin layer over “netdev provider” code, explained further below.

The other components may need attention during a port. You will almost certainly have to implement a “netdev provider”. Depending on the type of port you are doing and the desired performance, you may also have to implement an “ofproto provider” or a lower-level component called a “dpif” provider.

The following sections talk about these components in more detail.

Writing a netdev Provider

A “netdev provider” implements an operating system and hardware specific interface to “network devices”, e.g. `eth0` on Linux. Open vSwitch must be able to open each port on a switch as a `netdev`, so you will need to implement a “netdev provider” that works with your switch hardware and software.

`struct netdev_class`, in `lib/netdev-provider.h`, defines the interfaces required to implement a `netdev`. That structure contains many function pointers, each of which has a comment that is meant to describe its behavior in detail. If the requirements are unclear, report this as a bug.

The `netdev` interface can be divided into a few rough categories:

- Functions required to properly implement OpenFlow features. For example, OpenFlow requires the ability to report the Ethernet hardware address of a port. These functions must be implemented for minimally correct operation.
- Functions required to implement optional Open vSwitch features. For example, the Open vSwitch support for in-band control requires `netdev` support for inspecting the TCP/IP stack’s ARP table. These functions must be implemented if the corresponding OVS features are to work, but may be omitted initially.
- Functions needed in some implementations but not in others. For example, most kinds of ports (see below) do not need functionality to receive packets from a network device.

The existing `netdev` implementations may serve as useful examples during a port:

- `lib/netdev-linux.c` implements `netdev` functionality for Linux network devices, using Linux kernel calls. It may be a good place to start for full-featured `netdev` implementations.
- `lib/netdev-vport.c` provides support for “virtual ports” implemented by the Open vSwitch datapath module for the Linux kernel. This may serve as a model for minimal `netdev` implementations.

- `lib/netdev-dummy.c` is a fake netdev implementation useful only for testing.

Porting Strategies

After a netdev provider has been implemented for a system's network devices, you may choose among three basic porting strategies.

The lowest-effort strategy is to use the “userspace switch” implementation built into Open vSwitch. This ought to work, without writing any more code, as long as the netdev provider that you implemented supports receiving packets. It yields poor performance, however, because every packet passes through the `ovs-vswitchd` process. Refer to *Open vSwitch without Kernel Support* for instructions on how to configure a userspace switch.

If the userspace switch is not the right choice for your port, then you will have to write more code. You may implement either an “ofproto provider” or a “dpif provider”. Which you should choose depends on a few different factors:

- Only an ofproto provider can take full advantage of hardware with built-in support for wildcards (e.g. an ACL table or a TCAM).
- A dpif provider can take advantage of the Open vSwitch built-in implementations of bonding, LACP, 802.1ag, 802.1Q VLANs, and other features. An ofproto provider has to provide its own implementations, if the hardware can support them at all.
- A dpif provider is usually easier to implement, but most appropriate for software switching. It “explodes” wildcard rules into exact-match entries (with an optional wildcard mask). This allows fast hash lookups in software, but makes inefficient use of TCAMs in hardware that support wildcarding.

The following sections describe how to implement each kind of port.

ofproto Providers

An “ofproto provider” is what ofproto uses to directly monitor and control an OpenFlow-capable switch. `struct ofproto_class`, in `ofproto/ofproto-provider.h`, defines the interfaces to implement an ofproto provider for new hardware or software. That structure contains many function pointers, each of which has a comment that is meant to describe its behavior in detail. If the requirements are unclear, report this as a bug.

The ofproto provider interface is preliminary. Let us know if it seems unsuitable for your purpose. We will try to improve it.

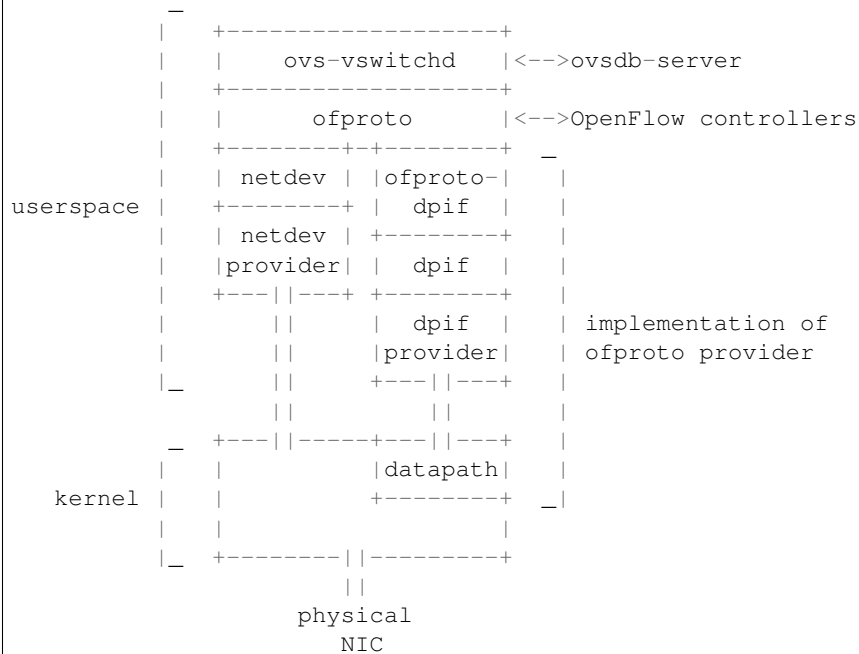
Writing a dpif Provider

Open vSwitch has a built-in ofproto provider named “ofproto-dpif”, which is built on top of a library for manipulating datapaths, called “dpif”. A “datapath” is a simple flow table, one that is only required to support exact-match flows, that is, flows without wildcards. When a packet arrives on a network device, the datapath looks for it in this table. If there is a match, then it performs the associated actions. If there is no match, the datapath passes the packet up to ofproto-dpif, which maintains the full OpenFlow flow table. If the packet matches in this flow table, then ofproto-dpif executes its actions and inserts a new entry into the dpif flow table. (Otherwise, ofproto-dpif passes the packet up to ofproto to send the packet to the OpenFlow controller, if one is configured.)

When calculating the dpif flow, ofproto-dpif generates an exact-match flow that describes the missed packet. It makes an effort to figure out what fields can be wildcarded based on the switch's configuration and OpenFlow flow table. The dpif is free to ignore the suggested wildcards and only support the exact-match entry. However, if the dpif supports wildcarding, then it can use the masks to match multiple flows with fewer entries and potentially significantly reduce the number of flow misses handled by ofproto-dpif.

The “dpif” library in turn delegates much of its functionality to a “dpif provider”. The following diagram shows how dpif providers fit into the Open vSwitch architecture:

Architecture



struct `dpif_class`, in `lib/dpif-provider.h`, defines the interfaces required to implement a `dpif` provider for new hardware or software. That structure contains many function pointers, each of which has a comment that is meant to describe its behavior in detail. If the requirements are unclear, report this as a bug.

There are two existing `dpif` implementations that may serve as useful examples during a port:

- `lib/dpif-netlink.c` is a Linux-specific `dpif` implementation that talks to an Open vSwitch-specific kernel module (whose sources are in the “datapath” directory). The kernel module performs all of the switching work, passing packets that do not match any flow table entry up to userspace. This `dpif` implementation is essentially a wrapper around calls into the kernel module.
- `lib/dpif-netdev.c` is a generic `dpif` implementation that performs all switching internally. This is how the Open vSwitch userspace switch is implemented.

Miscellaneous Notes

Open vSwitch source code uses `uint16_t`, `uint32_t`, and `uint64_t` as fixed-width types in host byte order, and `ovs_be16`, `ovs_be32`, and `ovs_be64` as fixed-width types in network byte order. Each of the latter is equivalent to the one of the former, but the difference in name makes the intended use obvious.

The default “fail-mode” for Open vSwitch bridges is “standalone”, meaning that, when the OpenFlow controllers cannot be contacted, Open vSwitch acts as a regular MAC-learning switch. This works well in virtualization environments where there is normally just one uplink (either a single physical interface or a bond). In a more general environment, it can create loops. So, if you are porting to a general-purpose switch platform, you should consider changing the default “fail-mode” to “secure”, which does not behave this way. See documentation for the “fail-mode” column in the Bridge table in `ovs-vswitchd.conf.db(5)` for more information.

`lib/entropy.c` assumes that it can obtain high-quality random number seeds at startup by reading from `/dev/urandom`. You will need to modify it if this is not true on your platform.

`vswitchd/system-stats.c` only knows how to obtain some statistics on Linux. Optionally you may implement them for your platform as well.

Why OVS Does Not Support Hybrid Providers

The *porting strategies* section above describes the “ofproto provider” and “dpif provider” porting strategies. Only an ofproto provider can take advantage of hardware TCAM support, and only a dpif provider can take advantage of the OVS built-in implementations of various features. It is therefore tempting to suggest a hybrid approach that shares the advantages of both strategies.

However, Open vSwitch does not support a hybrid approach. Doing so may be possible, with a significant amount of extra development work, but it does not yet seem worthwhile, for the reasons explained below.

First, user surprise is likely when a switch supports a feature only with a high performance penalty. For example, one user questioned why adding a particular OpenFlow action to a flow caused a 1,058x slowdown on a hardware OpenFlow implementation¹. The action required the flow to be implemented in software.

Given that implementing a flow in software on the slow management CPU of a hardware switch causes a major slowdown, software-implemented flows would only make sense for very low-volume traffic. But many of the features built into the OVS software switch implementation would need to apply to every flow to be useful. There is no value, for example, in applying bonding or 802.1Q VLAN support only to low-volume traffic.

Besides supporting features of OpenFlow actions, a hybrid approach could also support forms of matching not supported by particular switching hardware, by sending all packets that might match a rule to software. But again this can cause an unacceptable slowdown by forcing bulk traffic through software in the hardware switch’s slow management CPU. Consider, for example, a hardware switch that can match on the IPv6 Ethernet type but not on fields in IPv6 headers. An OpenFlow table that matched on the IPv6 Ethernet type would perform well, but adding a rule that matched only UDPv6 would force every IPv6 packet to software, slowing down not just UDPv6 but all IPv6 processing.

Questions

Direct porting questions to dev@openvswitch.org. We will try to use questions to improve this porting guide.

4.1.5 OpenFlow Support in Open vSwitch

Open vSwitch support for OpenFlow 1.1 and beyond is a work in progress. This file describes the work still to be done.

The Plan

OpenFlow version support is not a build-time option. A single build of Open vSwitch must be able to handle all supported versions of OpenFlow. Ideally, even at runtime it should be able to support all protocol versions at the same time on different OpenFlow bridges (and perhaps even on the same bridge).

At the same time, it would be a shame to litter the core of the OVS code with lots of ugly code concerned with the details of various OpenFlow protocol versions.

The primary approach to compatibility is to abstract most of the details of the differences from the core code, by adding a protocol layer that translates between OF1.x and a slightly higher-level abstract representation. The core of this approach is the many `struct ofputil_*` structures in `include/openvswitch/ofp-*.h`.

As a consequence of this approach, OVS cannot use OpenFlow protocol definitions that closely resemble those in the OpenFlow specification, because `openflow.h` in different versions of the OpenFlow specification defines the same identifier with different values. Instead, `openflow-common.h` contains definitions that are common to all the specifications and separate protocol version-specific headers contain protocol-specific definitions renamed so as not to conflict, e.g. `OFPAT10_ENQUEUE` and `OFPAT11_ENQUEUE` for the OpenFlow 1.0 and 1.1 values for

¹ Aaron Rosen, “Modify packet fields extremely slow”, openflow-discuss mailing list, June 26, 2011, archived at <https://mailman.stanford.edu/pipermail/openflow-discuss/2011-June/002386.html>.

OFPAT_ENQUEUE. Generally, in cases of conflict, the protocol layer will define a more abstract OFPUTIL_* or struct ofputil_*.

Here are the current approaches in a few tricky areas:

- Port numbering.

OpenFlow 1.0 has 16-bit port numbers and later OpenFlow versions have 32-bit port numbers. For now, OVS support for later protocol versions requires all port numbers to fall into the 16-bit range, translating the reserved OFPP_* port numbers.

- Actions.

OpenFlow 1.0 and later versions have very different ideas of actions. OVS reconciles by translating all the versions' actions (and instructions) to and from a common internal representation.

OpenFlow 1.1

OpenFlow 1.1 support is complete.

OpenFlow 1.2

OpenFlow 1.2 support is complete.

OpenFlow 1.3

OpenFlow 1.3 support requires OpenFlow 1.2 as a prerequisite, plus the following additional work. (This is based on the change log at the end of the OF1.3 spec, reusing most of the section titles directly. I didn't compare the specs carefully yet.)

- Add support for multipart requests.

Currently we always report OFPBRC_MULTIPART_BUFFER_OVERFLOW.

(optional for OF1.3+)

- IPv6 extension header handling support.

Fully implementing this requires kernel support. This likely will take some careful and probably time-consuming design work. The actual coding, once that is all done, is probably 2 or 3 days work.

(optional for OF1.3+)

- Auxiliary connections.

An implementation in generic code might be a week's worth of work. The value of an implementation in generic code is questionable, though, since much of the benefit of auxiliary connections is supposed to be to take advantage of hardware support. (We could make the kernel module somehow send packets across the auxiliary connections directly, for some kind of "hardware" support, if we judged it useful enough.)

(optional for OF1.3+)

- Provider Backbone Bridge tagging.

I don't plan to implement this (but we'd accept an implementation).

(optional for OF1.3+)

- On-demand flow counters.

I think this might be a real optimization in some cases for the software switch.

(optional for OF1.3+)

OpenFlow 1.4 & ONF Extensions for 1.3.X Pack1

The following features are both defined as a set of ONF Extensions for 1.3 and integrated in 1.4.

When defined as an ONF Extension for 1.3, the feature is using the Experimenter mechanism with the ONF Experimenter ID.

When defined integrated in 1.4, the feature use the standard OpenFlow structures (for example defined in openflow-1.4.h).

The two definitions for each feature are independent and can exist in parallel in OVS.

- Flow entry notifications

This seems to be modelled after OVS's NXST_FLOW_MONITOR.

(EXT-187) (optional for OF1.4+)

- Flow entry eviction

OVS has flow eviction functionality. `table_mod OFPTC_EVICTON`, `flow_mod 'importance'`, and `table_desc ofp_table_mod_prop_eviction` need to be implemented.

(EXT-192-e)

(optional for OF1.4+)

- Vacancy events

(EXT-192-v)

(optional for OF1.4+)

- Table synchronisation

Probably not so useful to the software switch.

(EXT-232)

(optional for OF1.4+)

- Group and Meter change notifications

(EXT-235)

(optional for OF1.4+)

- PBB UCA header field

See comment on Provider Backbone Bridge in section about OpenFlow 1.3.

(EXT-256)

(optional for OF1.4+)

OpenFlow 1.4 only

Those features are those only available in OpenFlow 1.4, other OpenFlow 1.4 features are listed in the previous section.

- More extensible wire protocol

Many on-wire structures got TLVs.

All required features are now supported. Remaining optional: table desc, table-status

(EXT-262)

(required for OF1.4+)

- Optical port properties

(EXT-154)

(optional for OF1.4+)

OpenFlow 1.5 & ONF Extensions for 1.3.X Pack2

The following features are both defined as a set of ONF Extensions for 1.3 and integrated in 1.5. Note that this list is not definitive as those are not yet published.

When defined as an ONF Extension for 1.3, the feature is using the Experimenter mechanism with the ONF Experimenter ID. When defined integrated in 1.5, the feature use the standard OpenFlow structures (for example defined in openflow-1.5.h).

The two definitions for each feature are independent and can exist in parallel in OVS.

- Time scheduled bundles

(EXT-340)

(optional for OF1.5+)

OpenFlow 1.5 only

Those features are those only available in OpenFlow 1.5, other OpenFlow 1.5 features are listed in the previous section. Note that this list is not definitive as OpenFlow 1.5 is not yet published.

- Egress Tables

(EXT-306)

(optional for OF1.5+)

- Extensible Flow Entry Statistics

(EXT-334)

(required for OF1.5+)

- Flow Entry Statistics Trigger

(EXT-335)

(optional for OF1.5+)

- Controller connection status

Prototype for OVS was done during specification.

(EXT-454)

(optional for OF1.5+)

- Meter action

(EXT-379)

(required for OF1.5+ if metering is supported)

- Port properties for pipeline fields

Prototype for OVS was done during specification.

(EXT-388)

(optional for OF1.5+)

- Port property for recirculation

Prototype for OVS was done during specification.

(EXT-399)

(optional for OF1.5+)

General

- ovs-ofctl(8) often lists as Nicira extensions features that later OpenFlow versions support in standard ways.

How to contribute

If you plan to contribute code for a feature, please let everyone know on ovs-dev before you start work. This will help avoid duplicating work.

Consider the following:

- Testing.

Please test your code.

- Unit tests.

Consider writing some. The tests directory has many examples that you can use as a starting point.

- ovs-ofctl.

If you add a feature that is useful for some ovs-ofctl command then you should add support for it there.

- Documentation.

If you add a user-visible feature, then you should document it in the appropriate manpage and mention it in NEWS as well.

Refer to *Contributing to Open vSwitch* for more information.

4.1.6 Bonding

Bonding allows two or more interfaces (the “slaves”) to share network traffic. From a high-level point of view, bonded interfaces act like a single port, but they have the bandwidth of multiple network devices, e.g. two 1 GB physical interfaces act like a single 2 GB interface. Bonds also increase robustness: the bonded port does not go down as long as at least one of its slaves is up.

In `vswitchd`, a bond always has at least two slaves (and may have more). If a configuration error, etc. would cause a bond to have only one slave, the port becomes an ordinary port, not a bonded port, and none of the special features of bonded ports described in this section apply.

There are many forms of bonding of which `ovs-vswitchd` implements only a few. The most complex bond `ovs-vswitchd` implements is called “source load balancing” or SLB bonding. SLB bonding divides traffic among the slaves based on the Ethernet source address. This is useful only if the traffic over the bond has multiple Ethernet source addresses, for example if network traffic from multiple VMs are multiplexed over the bond.

Note: Most of the `ovs-vswitchd` implementation is in `vswitchd/bridge.c`, so code references below should be assumed to refer to that file except as otherwise specified.

Enabling and Disabling Slaves

When a bond is created, a slave is initially enabled or disabled based on whether carrier is detected on the NIC (see `iface_create()`). After that, a slave is disabled if its carrier goes down for a period of time longer than the `down-delay`, and it is enabled if carrier comes up for longer than the `updelay` (see `bond_link_status_update()`). There is one exception where the `updelay` is skipped: if no slaves at all are currently enabled, then the first slave on which carrier comes up is enabled immediately.

The `updelay` should be set to a time longer than the STP forwarding delay of the physical switch to which the bond port is connected (if STP is enabled on that switch). Otherwise, the slave will be enabled, and load may be shifted to it, before the physical switch starts forwarding packets on that port, which can cause some data to be “blackholed” for a time. The exception for a single enabled slave does not cause any problem in this regard because when no slaves are enabled all output packets are blackholed anyway.

When a slave becomes disabled, the `vswitch` immediately chooses a new output port for traffic that was destined for that slave (see `bond_enable_slave()`). It also sends a “gratuitous learning packet”, specifically a RARP, on the bond port (on the newly chosen slave) for each MAC address that the `vswitch` has learned on a port other than the bond (see `bundle_send_learning_packets()`), to teach the physical switch that the new slave should be used in place of the one that is now disabled. (This behavior probably makes sense only for a `vswitch` that has only one port (the bond) connected to a physical switch; `vswitchd` should probably provide a way to disable or configure it in other scenarios.)

Bond Packet Input

Bonding accepts unicast packets on any bond slave. This can occasionally cause packet duplication for the first few packets sent to a given MAC, if the physical switch attached to the bond is flooding packets to that MAC because it has not yet learned the correct slave for that MAC.

Bonding only accepts multicast (and broadcast) packets on a single bond slave (the “active slave”) at any given time. Multicast packets received on other slaves are dropped. Otherwise, every multicast packet would be duplicated, once for every bond slave, because the physical switch attached to the bond will flood those packets.

Bonding also drops received packets when the `vswitch` has learned that the packet’s MAC is on a port other than the bond port itself. This is because it is likely that the `vswitch` itself sent the packet out the bond port on a different slave and is now receiving the packet back. This occurs when the packet is multicast or the physical switch has not yet learned the MAC and is flooding it. However, the `vswitch` makes an exception to this rule for broadcast ARP replies, which indicate that the MAC has moved to another switch, probably due to VM migration. (ARP replies are normally unicast, so this exception does not match normal ARP replies. It will match the learning packets sent on bond fail-over.)

The active slave is simply the first slave to be enabled after the bond is created (see `bond_choose_active_slave()`). If the active slave is disabled, then a new active slave is chosen among the

slaves that remain active. Currently due to the way that configuration works, this tends to be the remaining slave whose interface name is first alphabetically, but this is by no means guaranteed.

Bond Packet Output

When a packet is sent out a bond port, the bond slave actually used is selected based on the packet's source MAC and VLAN tag (see `bond_choose_output_slave()`). In particular, the source MAC and VLAN tag are hashed into one of 256 values, and that value is looked up in a hash table (the “bond hash”) kept in the `bond_hash` member of struct port. The hash table entry identifies a bond slave. If no bond slave has yet been chosen for that hash table entry, vswitchd chooses one arbitrarily.

Every 10 seconds, vswitchd rebalances the bond slaves (see `bond_rebalance()`). To rebalance, vswitchd examines the statistics for the number of bytes transmitted by each slave over approximately the past minute, with data sent more recently weighted more heavily than data sent less recently. It considers each of the slaves in order from most-loaded to least-loaded. If highly loaded slave H is significantly more heavily loaded than the least-loaded slave L, and slave H carries at least two hashes, then vswitchd shifts one of H's hashes to L. However, vswitchd will only shift a hash from H to L if it will decrease the ratio of the load between H and L by at least 0.1.

Currently, “significantly more loaded” means that H must carry at least 1 Mbps more traffic, and that traffic must be at least 3% greater than L's.

Bond Balance Modes

Each bond balancing mode has different considerations, described below.

LACP Bonding

LACP bonding requires the remote switch to implement LACP, but it is otherwise very simple in that, after LACP negotiation is complete, there is no need for special handling of received packets.

Several of the physical switches that support LACP block all traffic for ports that are configured to use LACP, until LACP is negotiated with the host. When configuring a LACP bond on a OVS host (eg: XenServer), this means that there will be an interruption of the network connectivity between the time the ports on the physical switch and the bond on the OVS host are configured. The interruption may be relatively long, if different people are responsible for managing the switches and the OVS host.

Such network connectivity failure can be avoided if LACP can be configured on the OVS host before configuring the physical switch, and having the OVS host fall back to a bond mode (active-backup) till the physical switch LACP configuration is complete. An option “lacp-fallback-ab” exists to provide such behavior on Open vSwitch.

Active Backup Bonding

Active Backup bonds send all traffic out one “active” slave until that slave becomes unavailable. Since they are significantly less complicated than SLB bonds, they are preferred when LACP is not an option. Additionally, they are the only bond mode which supports attaching each slave to a different upstream switch.

SLB Bonding

SLB bonding allows a limited form of load balancing without the remote switch's knowledge or cooperation. The basics of SLB are simple. SLB assigns each source MAC+VLAN pair to a link and transmits all packets from that MAC+VLAN through that link. Learning in the remote switch causes it to send packets to that MAC+VLAN through the same link.

SLB bonding has the following complications:

0. When the remote switch has not learned the MAC for the destination of a unicast packet and hence floods the packet to all of the links on the SLB bond, Open vSwitch will forward duplicate packets, one per link, to each other switch port.

Open vSwitch does not solve this problem.

1. When the remote switch receives a multicast or broadcast packet from a port not on the SLB bond, it will forward it to all of the links in the SLB bond. This would cause packet duplication if not handled specially.

Open vSwitch avoids packet duplication by accepting multicast and broadcast packets on only the active slave, and dropping multicast and broadcast packets on all other slaves.

2. When Open vSwitch forwards a multicast or broadcast packet to a link in the SLB bond other than the active slave, the remote switch will forward it to all of the other links in the SLB bond, including the active slave. Without special handling, this would mean that Open vSwitch would forward a second copy of the packet to each switch port (other than the bond), including the port that originated the packet.

Open vSwitch deals with this case by dropping packets received on any SLB bonded link that have a source MAC+VLAN that has been learned on any other port. (This means that SLB as implemented in Open vSwitch relies critically on MAC learning. Notably, SLB is incompatible with the “flood_vlans” feature.)

3. Suppose that a MAC+VLAN moves to an SLB bond from another port (e.g. when a VM is migrated from this hypervisor to a different one). Without additional special handling, Open vSwitch will not notice until the MAC learning entry expires, up to 60 seconds later as a consequence of rule #2.

Open vSwitch avoids a 60-second delay by listening for gratuitous ARPs, which VMs commonly emit upon migration. As an exception to rule #2, a gratuitous ARP received on an SLB bond is not dropped and updates the MAC learning table in the usual way. (If a move does not trigger a gratuitous ARP, or if the gratuitous ARP is lost in the network, then a 60-second delay still occurs.)

4. Suppose that a MAC+VLAN moves from an SLB bond to another port (e.g. when a VM is migrated from a different hypervisor to this one), that the MAC+VLAN emits a gratuitous ARP, and that Open vSwitch forwards that gratuitous ARP to a link in the SLB bond other than the active slave. The remote switch will forward the gratuitous ARP to all of the other links in the SLB bond, including the active slave. Without additional special handling, this would mean that Open vSwitch would learn that the MAC+VLAN was located on the SLB bond, as a consequence of rule #3.

Open vSwitch avoids this problem by “locking” the MAC learning table entry for a MAC+VLAN from which a gratuitous ARP was received from a non-SLB bond port. For 5 seconds, a locked MAC learning table entry will not be updated based on a gratuitous ARP received on a SLB bond.

4.1.7 Open vSwitch Networking Namespaces on Linux

The Open vSwitch has networking namespaces basic support on Linux. That allows ovs-vswitchd daemon to continue tracking status and statistics after moving a port to another networking namespace.

How It Works

The daemon ovs-vswitchd runs on what is called parent network namespace. It listens to netlink event messages from all networking namespaces (netns) with an identifier on the parent. Each netlink message contains the network namespace identifier (netnsid) as ancillary data which is used to match the event to the corresponding port.

The ovs-vswitchd uses an extended openvswitch kernel API¹ to get the current netnsid (stored in struct netdev_linux) and statistics from a specific port. The netnsid remains cached in userspace until a changing event is received, for example, when the port is moved to another network namespace.

¹ Request cmd: OVS_VPORT_CMD_GET, attribute: OVS_VPORT_ATTR_NETNSID

Using another extended kernel API², the daemon gets port’s information such as flags, MTU, MAC address and ifindex from a port already in another namespace.

The upstream kernel 4.15 includes the necessary changes for the basic support. In case of the running kernel doesn’t provide the APIs, the daemon falls back to the previous behavior.

Limitations

Currently it is only possible to retrieve the information listed in the above section. Most of other operations, for example querying MII or setting MTU, lacks the proper API in the kernel, so they remain unsupported.

In most use cases that needs to move ports to another networking namespaces should use veth pairs instead because it offers a cleaner and more robust solution with no noticeable performance penalty.

4.1.8 OVSDb Replication Implementation

Given two Open vSwitch databases with the same schema, OVSDb replication keeps these databases in the same state, i.e. each of the databases have the same contents at any given time even if they are not running in the same host. This document elaborates on the implementation details to provide this functionality.

Terminology

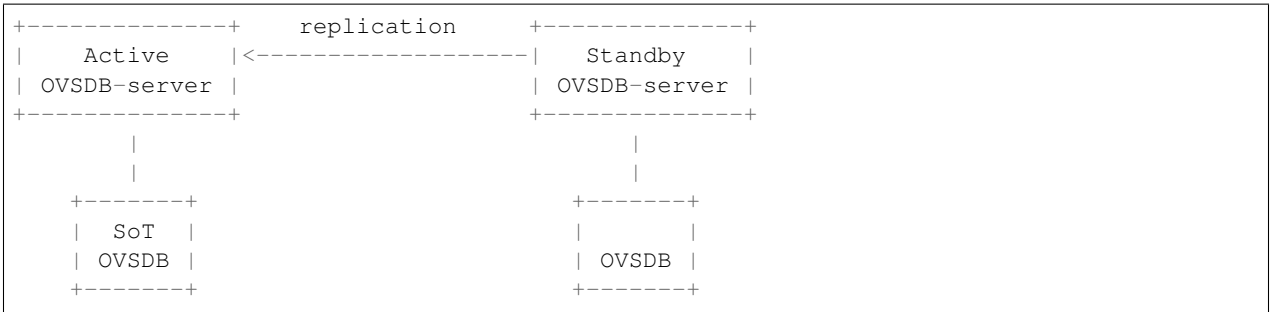
Source of truth database database whose content will be replicated to another database.

Active server ovssdb-server providing RPC interface to the source of truth database.

Standby server ovssdb-server providing RPC interface to the database that is not the source of truth.

Design

The overall design of replication consists of one ovssdb-server (active server) communicating the state of its databases to another ovssdb-server (standby server) so that the latter keep its own databases in that same state. To achieve this, the standby server acts as a client of the active server, in the sense that it sends a monitor request to keep up to date with the changes in the active server databases. When a notification from the active server arrives, the standby server executes the necessary set of operations so its databases reach the same state as the active server databases. Below is the design represented as a diagram.:



² Request cmd: RTM_GETLINK passing IFLA_IF_NETNSID attribute.

Setting Up The Replication

To initiate the replication process, the standby server must be executed indicating the location of the active server via the command line option `--sync-from=server`, where `server` can take any form described in the `ovsdb-client` manpage and it must specify an active connection type (`tcp`, `unix`, `ssl`). This option will cause the standby server to attempt to send a monitor request to the active server in every main loop iteration, until the active server responds.

When sending a monitor request the standby server is doing the following:

1. Erase the content of the databases for which it is providing a RPC interface.
2. Open the jsonrpc channel to communicate with the active server.
3. Fetch all the databases located in the active server.
4. For each database with the same schema in both the active and standby servers: construct and send a monitor request message specifying the tables that will be monitored (i.e all the tables on the database except the ones blacklisted [*]).
5. Set the standby database to the current state of the active database.

Once the monitor request message is sent, the standby server will continuously receive notifications of changes occurring to the tables specified in the request. The process of handling this notifications is detailed in the next section.

[*] A set of tables that will be excluded from replication can be configure as a blacklist of tables via the command line option `--sync-exclude-tables=db:table[,db:table]...`, where `db` corresponds to the database where the table resides.

Replication Process

The replication process consists on handling the update notifications received in the standby server caused by the monitor request that was previously sent to the active server. In every loop iteration, the standby server attempts to receive a message from the active server which can be an error, an echo message (used to keep the connection alive) or an update notification. In case the message is a fatal error, the standby server will disconnect from the active without dropping the replicated data. If it is an echo message, the standby server will reply with an echo message as well. If the message is an update notification, the following process occurs:

1. Create a new transaction.
2. Get the `<table-updates>` object from the `params` member of the notification.
3. For each `<table-update>` in the `<table-updates>` object do:
 - (a) For each `<row-update>` in `<table-update>` check what kind of operation should be executed according to the following criteria about the presence of the object members:
 - If `old` member is not present, execute an insert operation using `<row>` from the `new` member.
 - If `old` member is present and `new` member is not present, execute a delete operation using `<row>` from the `old` member
 - If both `old` and `new` members are present, execute an update operation using `<row>` from the `new` member.
4. Commit the transaction.

If an error occurs during the replication process, all replication is restarted by resending a new monitor request as described in the section “Setting up the replication”.

Runtime Management Commands

Runtime management commands can be sent to a running standby server via `ovs-appctl` in order to configure the replication functionality. The available commands are the following.

ovsdb-server/set-remote-ovsdb-server {server} sets the name of the active server

ovsdb-server/get-remote-ovsdb-server gets the name of the active server

ovsdb-server/connect-remote-ovsdb-server causes the server to attempt to send a monitor request every main loop iteration

ovsdb-server/disconnect-remote-ovsdb-server closes the jsonrpc channel between the active server and frees the memory used for the replication configuration.

ovsdb-server/set-sync-exclude-tables {db:table, ...} sets the tables list that will be excluded from being replicated

ovsdb-server/get-sync-excluded-tables gets the tables list that is currently excluded from replication

4.1.9 The DPDK Datapath

DPDK Bridges

The DPDK datapath requires specially configured bridge(s) in order to utilize DPDK-backed *physical* and *virtual* ports.

Quick Example

This example demonstrates how to add a bridge using the DPDK datapath:

```
$ ovs-vsctl add-br br0 -- set bridge br0 datapath_type=netdev
```

This assumes Open vSwitch has been built with DPDK support. Refer to *Open vSwitch with DPDK* for more information.

Extended & Custom Statistics

The DPDK Extended Statistics API allows PMDs to expose a unique set of statistics. The Extended Statistics are implemented and supported only for DPDK physical and vHost ports. Custom statistics are a dynamic set of counters which can vary depending on the driver. Those statistics are implemented for DPDK physical ports and contain all “dropped”, “error” and “management” counters from XSTATS. A list of all XSTATS counters can be found [here](#).

Note: vHost ports only support RX packet size-based counters. TX packet size counters are not available.

To enable statistics, you have to enable OpenFlow 1.4 support for OVS. To configure a bridge, `br0`, to support OpenFlow version 1.4, run:

```
$ ovs-vsctl set bridge br0 datapath_type=netdev \
  protocols=OpenFlow10,OpenFlow11,OpenFlow12,OpenFlow13,OpenFlow14
```

Once configured, check the OVSDDB protocols column in the bridge table to ensure OpenFlow 1.4 support is enabled:

```
$ ovsdb-client dump Bridge protocols
```

You can also query the port statistics by explicitly specifying the `-O OpenFlow14` option:

```
$ ovs-ofctl -O OpenFlow14 dump-ports br0
```

EMC Insertion Probability

By default 1 in every 100 flows is inserted into the Exact Match Cache (EMC). It is possible to change this insertion probability by setting the `emc-insert-inv-prob=N` option:

```
$ ovs-vsctl --no-wait set Open_vSwitch . other_config:emc-insert-inv-prob=N
```

where:

N A positive integer representing the inverse probability of insertion, i.e. on average 1 in every N packets with a unique flow will generate an EMC insertion.

If N is set to 1, an insertion will be performed for every flow. If set to 0, no insertions will be performed and the EMC will effectively be disabled.

With default N set to 100, higher megaflow hits will occur initially as observed with `pmd stats`:

```
$ ovs-appctl dpif-netdev/pmd-stats-show
```

For certain traffic profiles with many parallel flows, it's recommended to set N to '0' to achieve higher forwarding performance.

It is also possible to enable/disable EMC on per-port basis using:

```
$ ovs-vsctl set interface <iface> other_config:emc-enable={true,false}
```

Note: This could be useful for cases where different number of flows expected on different ports. For example, if one of the VMs encapsulates traffic using additional headers, it will receive large number of flows but only few flows will come out of this VM. In this scenario it's much faster to use EMC instead of classifier for traffic from the VM, but it's better to disable EMC for the traffic which flows to the VM.

For more information on the EMC refer to [Open vSwitch with DPDK](#).

SMC cache (experimental)

SMC cache or signature match cache is a new cache level after EMC cache. The difference between SMC and EMC is SMC only stores a signature of a flow thus it is much more memory efficient. With same memory space, EMC can store 8k flows while SMC can store 1M flows. When traffic flow count is much larger than EMC size, it is generally beneficial to turn off EMC and turn on SMC. It is currently turned off by default and an experimental feature.

To turn on SMC:

```
$ ovs-vsctl --no-wait set Open_vSwitch . other_config:smc-enable=true
```

DPDK Physical Ports

The netdev datapath allows attaching of DPDK-backed physical interfaces in order to provide high-performance ingress/egress from the host.

Important: To use any DPDK-backed interface, you must ensure your bridge is configured correctly. For more information, refer to [DPDK Bridges](#).

Changed in version 2.7.0: Before Open vSwitch 2.7.0, it was necessary to prefix port names with a `dpdk` prefix. Starting with 2.7.0, this is no longer necessary.

Quick Example

This example demonstrates how to bind two `dpdk` ports, bound to physical interfaces identified by hardware IDs `0000:01:00.0` and `0000:01:00.1`, to an existing bridge called `br0`:

```
$ ovs-vsctl add-port br0 dpdk-p0 \
  -- set Interface dpdk-p0 type=dpdk options:dpdk-devargs=0000:01:00.0
$ ovs-vsctl add-port br0 dpdk-p1 \
  -- set Interface dpdk-p1 type=dpdk options:dpdk-devargs=0000:01:00.1
```

For the above example to work, the two physical interfaces must be bound to the DPDK poll-mode drivers in userspace rather than the traditional kernel drivers. See the *binding NIC drivers* <[dpdk-binding-nics](#)> section for details.

Binding NIC Drivers

DPDK operates entirely in userspace and, as a result, requires use of its own poll-mode drivers in user space for physical interfaces and a passthrough-style driver for the devices in kernel space.

There are two different tools for binding drivers: **driverctl** which is a generic tool for persistently configuring alternative device drivers, and **dpdk-devbind** which is a DPDK-specific tool and whose changes do not persist across reboots. In addition, there are two options available for this kernel space driver - VFIO (Virtual Function I/O) and UIO (Userspace I/O) - along with a number of drivers for each option. We will demonstrate examples of both tools and will use the `vfio-pci` driver, which is the more secure, robust driver of those available. More information can be found in the [DPDK documentation](#).

To list devices using **driverctl**, run:

```
$ driverctl -v list-devices | grep -i net
0000:07:00.0 igb (I350 Gigabit Network Connection (Ethernet Server Adapter I350-T2))
0000:07:00.1 igb (I350 Gigabit Network Connection (Ethernet Server Adapter I350-T2))
```

You can then bind one or more of these devices using the same tool:

```
$ driverctl set-override 0000:07:00.0 vfio-pci
```

Alternatively, to list devices using **dpdk-devbind**, run:

```
$ dpdk-devbind --status
Network devices using DPDK-compatible driver
=====
<none>
```

(continues on next page)

(continued from previous page)

```
Network devices using kernel driver
=====
0000:07:00.0 'I350 Gigabit Network Connection 1521' if=enp7s0f0 drv=igb unused=igb_uio
0000:07:00.1 'I350 Gigabit Network Connection 1521' if=enp7s0f1 drv=igb unused=igb_uio

Other Network devices
=====
...
```

Once again, you can then bind one or more of these devices using the same tool:

```
$ dpdk-devbind --bind=vfio-pci 0000:07:00.0
```

Changed in version 2.6.0: Open vSwitch 2.6.0 added support for DPDK 16.07, which in turn renamed the former `dpdk_nic_bind` tool to `dpdk-devbind`.

For more information, refer to the [DPDK documentation](#).

Multiqueue

Poll Mode Driver (PMD) threads are the threads that do the heavy lifting for the DPDK datapath. Correct configuration of PMD threads and the Rx queues they utilize is a requirement in order to deliver the high-performance possible with DPDK acceleration. It is possible to configure multiple Rx queues for dpdk ports, thus ensuring this is not a bottleneck for performance. For information on configuring PMD threads, refer to [PMD Threads](#).

Flow Control

Flow control can be enabled only on DPDK physical ports. To enable flow control support at Tx side while adding a port, run:

```
$ ovs-vsctl add-port br0 dpdk-p0 -- set Interface dpdk-p0 type=dpdk \
    options:dpdk-devargs=0000:01:00.0 options:tx-flow-ctrl=true
```

Similarly, to enable Rx flow control, run:

```
$ ovs-vsctl add-port br0 dpdk-p0 -- set Interface dpdk-p0 type=dpdk \
    options:dpdk-devargs=0000:01:00.0 options:rx-flow-ctrl=true
```

To enable flow control auto-negotiation, run:

```
$ ovs-vsctl add-port br0 dpdk-p0 -- set Interface dpdk-p0 type=dpdk \
    options:dpdk-devargs=0000:01:00.0 options:flow-ctrl-autoneg=true
```

To turn on the Tx flow control at run time for an existing port, run:

```
$ ovs-vsctl set Interface dpdk-p0 options:tx-flow-ctrl=true
```

The flow control parameters can be turned off by setting `false` to the respective parameter. To disable the flow control at Tx side, run:

```
$ ovs-vsctl set Interface dpdk-p0 options:tx-flow-ctrl=false
```

Rx Checksum Offload

By default, DPDK physical ports are enabled with Rx checksum offload.

Rx checksum offload can offer performance improvement only for tunneling traffic in OVS-DPDK because the checksum validation of tunnel packets is offloaded to the NIC. Also enabling Rx checksum may slightly reduce the performance of non-tunnel traffic, specifically for smaller size packet.

Hotplugging

OVS supports port hotplugging, allowing the use of physical ports that were not bound to DPDK when ovs-vswitchd was started.

Warning: This feature is not compatible with all NICs. Refer to vendor documentation for more information.

Important: Ports must be bound to DPDK. Refer to *Binding NIC Drivers* for more information.

To *hotplug* a port, simply add it like any other port:

```
$ ovs-vsctl add-port br0 dpdkx -- set Interface dpdkx type=dpdk \
    options:dpdk-devargs=0000:01:00.0
```

Ports can be detached using the `del-port` command:

```
$ ovs-vsctl del-port dpdkx
```

This should both delete the port and detach the device. If successful, you should see an INFO log. For example:

```
INFO|Device '0000:04:00.1' has been detached
```

If the log is not seen then the port can be detached like so:

```
$ ovs-appctl netdev-dpdk/detach 0000:01:00.0
```

Warning: Detaching should not be done if a device is known to be non-detachable, as this may cause the device to behave improperly when added back with `add-port`. The Chelsio Terminator adapters which use the `cxgbe` driver seem to be an example of this behavior; check the driver documentation if this is suspected.

For more information please refer to the [DPDK Port Hotplug Framework](#).

Representors

DPDK representors enable configuring a phy port to a guest (VM) machine.

OVS resides in the hypervisor which has one or more physical interfaces also known as the physical functions (PFs). If a PF supports SR-IOV it can be used to enable communication with the VMs via Virtual Functions (VFs). The VFs are virtual PCIe devices created from the physical Ethernet controller.

DPDK models a physical interface as a rte device on top of which an eth device is created. DPDK (version 18.xx) introduced the representors eth devices. A representor device represents the VF eth device (VM side) on the hypervisor side and operates on top of a PF. Representors are multi devices created on top of one PF.

For more information, refer to the [DPDK documentation](#).

Prior to port representors there was a one-to-one relationship between the PF and the eth device. With port representors the relationship becomes one PF to many eth devices. In case of two representors ports, when one of the ports is closed - the PCI bus cannot be detached until the second representor port is closed as well.

When configuring a PF-based port, OVS traditionally assigns the device PCI address in devargs. For an existing bridge called `br0` and PCI address `0000:08:00.0` an `add-port` command is written as:

```
$ ovs-vsctl add-port br0 dpdk-pf -- set Interface dpdk-pf type=dpdk \
    options:dpdk-devargs=0000:08:00.0
```

When configuring a VF-based port, DPDK uses an extended devargs syntax which has the following format:

```
BDBF,representor=[<representor id>]
```

This syntax shows that a representor is an enumerated eth device (with a representor ID) which uses the PF PCI address. The following commands add representors 3 and 5 using PCI device address `0000:08:00.0`:

```
$ ovs-vsctl add-port br0 dpdk-rep3 -- set Interface dpdk-rep3 type=dpdk \
    options:dpdk-devargs=0000:08:00.0,representor=[3]

$ ovs-vsctl add-port br0 dpdk-rep5 -- set Interface dpdk-rep5 type=dpdk \
    options:dpdk-devargs=0000:08:00.0,representor=[5]
```

Important: Representors ports are configured prior to OVS invocation and independently of it, or by other means as well. Please consult a NIC vendor instructions on how to establish representors. To verify their correct configuration, execute:

```
$ ovs-vsctl show
```

and make sure no errors are indicated.

Port representors are an example of multi devices. There are NICs which support multi devices by other methods than representors for which a generic devargs syntax is used. The generic syntax is based on the device mac address:

```
class=eth,mac=<MAC address>
```

For example, the following command adds a port to a bridge called `br0` using an eth device whose mac address is `00:11:22:33:44:55`:

```
$ ovs-vsctl add-port br0 dpdk-mac -- set Interface dpdk-mac type=dpdk \
    options:dpdk-devargs="class=eth,mac=00:11:22:33:44:55"
```

Jumbo Frames

DPDK physical ports can be configured to use Jumbo Frames. For more information, refer to [Jumbo Frames](#).

Link State Change (LSC) detection configuration

There are two methods to get the information when Link State Change (LSC) happens on a network interface: by polling or interrupt.

Configuring the lsc detection mode has no direct effect on OVS itself, instead it configures the NIC how it should handle link state changes. Processing the link state update request triggered by OVS takes less time using interrupt mode, since the NIC updates its link state in the background, while in polling mode the link state has to be fetched from the firmware every time to fulfil this request.

Note that not all PMD drivers support LSC interrupts.

The default configuration is polling mode. To set interrupt mode, option `dpdk-lsc-interrupt` has to be set to `true`.

Command to set interrupt mode for a specific interface:: `$ ovs-vsctl set interface <iface_name> options:dpdk-lsc-interrupt=true`

Command to set polling mode for a specific interface:: `$ ovs-vsctl set interface <iface_name> options:dpdk-lsc-interrupt=false`

DPDK vHost User Ports

The DPDK datapath provides DPDK-backed vHost user ports as a primary way to interact with guests. For more information on vHost User, refer to the [QEMU documentation](#) on same.

Important: To use any DPDK-backed interface, you must ensure your bridge is configured correctly. For more information, refer to [DPDK Bridges](#).

Quick Example

This example demonstrates how to add two `dpdkvhostuserclient` ports to an existing bridge called `br0`:

```
$ ovs-vsctl add-port br0 dpdkvhostclient0 \
  -- set Interface dpdkvhostclient0 type=dpdkvhostuserclient \
  options:vhost-server-path=/tmp/dpdkvhostclient0
$ ovs-vsctl add-port br0 dpdkvhostclient1 \
  -- set Interface dpdkvhostclient1 type=dpdkvhostuserclient \
  options:vhost-server-path=/tmp/dpdkvhostclient1
```

For the above examples to work, an appropriate server socket must be created at the paths specified (`/tmp/dpdkvhostclient0` and `/tmp/dpdkvhostclient1`). These sockets can be created with QEMU; see the [vhost-user client](#) section for details.

vhost-user vs. vhost-user-client

Open vSwitch provides two types of vHost User ports:

- `vhost-user` (`dpdkvhostuser`)
- `vhost-user-client` (`dpdkvhostuserclient`)

vHost User uses a client-server model. The server creates/manages/destroys the vHost User sockets, and the client connects to the server. Depending on which port type you use, `dpdkvhostuser` or `dpdkvhostuserclient`, a different configuration of the client-server model is used.

For vhost-user ports, Open vSwitch acts as the server and QEMU the client. This means if OVS dies, all VMs **must** be restarted. On the other hand, for vhost-user-client ports, OVS acts as the client and QEMU the server. This means OVS can die and be restarted without issue, and it is also possible to restart an instance itself. For this reason, vhost-user-client ports are the preferred type for all known use cases; the only limitation is that vhost-user client mode ports require QEMU version 2.7. Ports of type vhost-user are currently deprecated and will be removed in a future release.

vhost-user

Important: Use of vhost-user ports requires QEMU ≥ 2.2 ; vhost-user ports are *deprecated*.

To use vhost-user ports, you must first add said ports to the switch. DPDK vhost-user ports can have arbitrary names with the exception of forward and backward slashes, which are prohibited. For vhost-user, the port type is `dpdkvhostuser`:

```
$ ovs-vsctl add-port br0 vhost-user-1 -- set Interface vhost-user-1 \
    type=dpdkvhostuser
```

This action creates a socket located at `/usr/local/var/run/openvswitch/vhost-user-1`, which you must provide to your VM on the QEMU command line.

Note: If you wish for the vhost-user sockets to be created in a sub-directory of `/usr/local/var/run/openvswitch`, you may specify this directory in the `ovsdb` like so:

```
$ ovs-vsctl --no-wait \
    set Open_vSwitch . other_config:vhost-sock-dir=subdir
```

Once the vhost-user ports have been added to the switch, they must be added to the guest. There are two ways to do this: using QEMU directly, or using `libvirt`.

Note: IOMMU is not supported with vhost-user ports.

Adding vhost-user ports to the guest (QEMU)

To begin, you must attach the vhost-user device sockets to the guest. To do this, you must pass the following parameters to QEMU:

```
-chardev socket,id=char1,path=/usr/local/var/run/openvswitch/vhost-user-1
-netdev type=vhost-user,id=mynet1,chardev=char1,vhostforce
-device virtio-net-pci,mac=00:00:00:00:00:01,netdev=mynet1
```

where `vhost-user-1` is the name of the vhost-user port added to the switch.

Repeat the above parameters for multiple devices, changing the chardev path and id as necessary. Note that a separate and different chardev path needs to be specified for each vhost-user device. For example you have a second vhost-user port named `vhost-user-2`, you append your QEMU command line with an additional set of parameters:


```
-chardev socket,id=char2,path=/usr/local/var/run/openvswitch/vhost-user-2
-netdev type=vhost-user,id=mynet2,chardev=char2,vhostforce
-device virtio-net-pci,mac=00:00:00:00:00:02,netdev=mynet2
```

In addition, QEMU must allocate the VM's memory on hugetlbfs. vhost-user ports access a virtio-net device's virtual rings and packet buffers mapping the VM's physical memory on hugetlbfs. To enable vhost-user ports to map the VM's memory into their process address space, pass the following parameters to QEMU:

```
-object memory-backend-file,id=mem,size=4096M,mem-path=/dev/hugepages,share=on
-numa node,memdev=mem -mem-prealloc
```

Finally, you may wish to enable multiqueue support. This is optional but, should you wish to enable it, run:

```
-chardev socket,id=char2,path=/usr/local/var/run/openvswitch/vhost-user-2
-netdev type=vhost-user,id=mynet2,chardev=char2,vhostforce,queues=$q
-device virtio-net-pci,mac=00:00:00:00:00:02,netdev=mynet2,mq=on,vectors=$v
```

where:

\$q The number of queues

\$v The number of vectors, which is $\$q * 2 + 2$

The vhost-user interface will be automatically reconfigured with required number of Rx and Tx queues after connection of virtio device. Manual configuration of `n_rxq` is not supported because OVS will work properly only if `n_rxq` will match number of queues configured in QEMU.

A least two PMDs should be configured for the vswitch when using multiqueue. Using a single PMD will cause traffic to be enqueued to the same vhost queue rather than being distributed among different vhost queues for a vhost-user interface.

If traffic destined for a VM configured with multiqueue arrives to the vswitch via a physical DPDK port, then the number of Rx queues should also be set to at least two for that physical DPDK port. This is required to increase the probability that a different PMD will handle the multiqueue transmission to the guest using a different vhost queue.

If one wishes to use multiple queues for an interface in the guest, the driver in the guest operating system must be configured to do so. It is recommended that the number of queues configured be equal to `$q`.

For example, this can be done for the Linux kernel virtio-net driver with:

```
$ ethtool -L <DEV> combined <$q>
```

where:

-L Changes the numbers of channels of the specified network device

combined Changes the number of multi-purpose channels.

Adding vhost-user ports to the guest (libvirt)

To begin, you must change the user and group that qemu runs under, and restart libvirtd.

- In `/etc/libvirt/qemu.conf` add/edit the following lines:

```
user = "root"
group = "root"
```

- Finally, restart the libvirtd process, For example, on Fedora:

```
$ systemctl restart libvirtd.service
```

Once complete, instantiate the VM. A sample XML configuration file is provided at the [end of this file](#). Save this file, then create a VM using this file:

```
$ virsh create demovm.xml
```

Once created, you can connect to the guest console:

```
$ virsh console demovm
```

The demovm xml configuration is aimed at achieving out of box performance on VM. These enhancements include:

- The vcpus are pinned to the cores of the CPU socket 0 using `vcpupin`.
- Configure NUMA cell and memory shared using `memAccess='shared'`.
- Disable `mrg_rxbuf='off'`

Refer to the [libvirt documentation](#) for more information.

vhost-user-client

Important: Use of vhost-user ports requires QEMU >= 2.7

To use vhost-user-client ports, you must first add said ports to the switch. Like DPDK vhost-user ports, DPDK vhost-user-client ports can have mostly arbitrary names. However, the name given to the port does not govern the name of the socket device. Instead, this must be configured by the user by way of a `vhost-server-path` option. For vhost-user-client, the port type is `dpdkvhostuserclient`:

```
$ VHOST_USER_SOCKET_PATH=/path/to/socket
$ ovs-vsctl add-port br0 vhost-client-1 \
    -- set Interface vhost-client-1 type=dpdkvhostuserclient \
        options:vhost-server-path=$VHOST_USER_SOCKET_PATH
```

Once the vhost-user-client ports have been added to the switch, they must be added to the guest. Like vhost-user ports, there are two ways to do this: using QEMU directly, or using libvirt. Only the QEMU case is covered here.

Adding vhost-user-client ports to the guest (QEMU)

Attach the vhost-user device sockets to the guest. To do this, you must pass the following parameters to QEMU:

```
-chardev socket,id=char1,path=$VHOST_USER_SOCKET_PATH,server
-netdev type=vhost-user,id=mynet1,chardev=char1,vhostforce
-device virtio-net-pci,mac=00:00:00:00:00:01,netdev=mynet1
```

where `vhost-user-1` is the name of the vhost-user port added to the switch.

If the corresponding `dpdkvhostuserclient` port has not yet been configured in OVS with `vhost-server-path=/path/to/socket`, QEMU will print a log similar to the following:

```
QEMU waiting for connection on: disconnected:unix:/path/to/socket,server
```

QEMU will wait until the port is created successfully in OVS to boot the VM. One benefit of using this mode is the ability for vHost ports to ‘reconnect’ in event of the switch crashing or being brought down. Once it is brought back up, the vHost ports will reconnect automatically and normal service will resume.

vhost-user-client IOMMU Support

vhost IOMMU is a feature which restricts the vhost memory that a virtio device can access, and as such is useful in deployments in which security is a concern.

IOMMU support may be enabled via a global config value, ``vhost-iommu-support``. Setting this to true enables vhost IOMMU support for all vhost ports when/where available:

```
$ ovs-vsctl set Open_vSwitch . other_config:vhost-iommu-support=true
```

The default value is false.

Important: Changing this value requires restarting the daemon.

Important: Enabling the IOMMU feature also enables the vhost user reply-ack protocol; this is known to work on QEMU v2.10.0, but is buggy on older versions (2.7.0 - 2.9.0, inclusive). Consequently, the IOMMU feature is disabled by default (and should remain so if using the aforementioned versions of QEMU). Starting with QEMU v2.9.1, vhost-iommu-support can safely be enabled, even without having an IOMMU device, with no performance penalty.

DPDK in the Guest

The DPDK `testpmd` application can be run in guest VMs for high speed packet forwarding between vhostuser ports. DPDK and `testpmd` application has to be compiled on the guest VM. Below are the steps for setting up the `testpmd` application in the VM.

Note: Support for DPDK in the guest requires QEMU ≥ 2.2

To begin, instantiate a guest as described in *vhost-user* or *vhost-user-client*. Once started, connect to the VM, download the DPDK sources to VM and build DPDK:

```
$ cd /root/dpdk/
$ wget http://fast.dpdk.org/rel/dpdk-18.11.tar.xz
$ tar xf dpdk-18.11.tar.xz
$ export DPDK_DIR=/root/dpdk/dpdk-18.11
$ export DPDK_TARGET=x86_64-native-linuxapp-gcc
$ export DPDK_BUILD=$DPDK_DIR/$DPDK_TARGET
$ cd $DPDK_DIR
$ make install T=$DPDK_TARGET DESTDIR=install
```

Build the test-pmd application:

```
$ cd app/test-pmd
$ export RTE_SDK=$DPDK_DIR
$ export RTE_TARGET=$DPDK_TARGET
$ make
```

Setup huge pages and DPDK devices using UIO:

```
$ sysctl vm.nr_hugepages=1024
$ mkdir -p /dev/hugepages
$ mount -t hugetlbfs hugetlbfs /dev/hugepages # only if not already mounted
$ modprobe uio
$ insmod $DPDK_BUILD/kmod/igb_uio.ko
$ $DPDK_DIR/usertools/dpdk-devbind.py --status
$ $DPDK_DIR/usertools/dpdk-devbind.py -b igb_uio 00:03.0 00:04.0
```

Note: vhost ports pci ids can be retrieved using:

```
lspci | grep Ethernet
```

Finally, start the application:

```
# TODO
```

Sample XML

```
<domain type='kvm'>
  <name>demovm</name>
  <uuid>4a9b3f53-fa2a-47f3-a757-dd87720d9d1d</uuid>
  <memory unit='KiB'>4194304</memory>
  <currentMemory unit='KiB'>4194304</currentMemory>
  <memoryBacking>
    <hugepages>
      <page size='2' unit='M' nodeset='0' />
    </hugepages>
  </memoryBacking>
  <vcpu placement='static'>2</vcpu>
  <cputune>
    <shares>4096</shares>
    <vcpupin vcpu='0' cpuset='4' />
    <vcpupin vcpu='1' cpuset='5' />
    <emulatorpin cpuset='4,5' />
  </cputune>
  <os>
    <type arch='x86_64' machine='pc'>hvm</type>
    <boot dev='hd' />
  </os>
  <features>
    <acpi />
    <apic />
  </features>
  <cpu mode='host-model'>
    <model fallback='allow' />
    <topology sockets='2' cores='1' threads='1' />
    <numa>
      <cell id='0' cpus='0-1' memory='4194304' unit='KiB' memAccess='shared' />
    </numa>
  </cpu>
  <on_poweroff>destroy</on_poweroff>
  <on_reboot>restart</on_reboot>
```

(continues on next page)

(continued from previous page)

```

<on_crash>destroy</on_crash>
<devices>
  <emulator>/usr/bin/qemu-system-x86_64</emulator>
  <disk type='file' device='disk'>
    <driver name='qemu' type='qcow2' cache='none'>
    <source file='/root/CentOS7_x86_64.qcow2'>
    <target dev='vda' bus='virtio'>
  </disk>
  <interface type='vhostuser'>
    <mac address='00:00:00:00:00:01'>
    <source type='unix' path='/usr/local/var/run/openvswitch/dpdkvhostuser0' mode=
↪ 'client'>
    <model type='virtio'>
    <driver queues='2'>
    <host mrg_rxbuf='on'>
    </driver>
  </interface>
  <interface type='vhostuser'>
    <mac address='00:00:00:00:00:02'>
    <source type='unix' path='/usr/local/var/run/openvswitch/dpdkvhostuser1' mode=
↪ 'client'>
    <model type='virtio'>
    <driver queues='2'>
    <host mrg_rxbuf='on'>
    </driver>
  </interface>
  <serial type='pty'>
    <target port='0'>
  </serial>
  <console type='pty'>
    <target type='serial' port='0'>
  </console>
</devices>
</domain>

```

Jumbo Frames

DPDK vHost User ports can be configured to use Jumbo Frames. For more information, refer to [Jumbo Frames](#).

vhost-user Dequeue Zero Copy (experimental)

Normally when dequeuing a packet from a vHost User device, a memcpy operation must be used to copy that packet from guest address space to host address space. This memcpy can be removed by enabling dequeue zero-copy like so:

```

$ ovs-vsctl add-port br0 dpdkvhostuserclient0 -- set Interface \
  dpdkvhostuserclient0 type=dpdkvhostuserclient \
  options:vhost-server-path=/tmp/dpdkvhostclient0 \
  options:dq-zero-copy=true

```

With this feature enabled, a reference (pointer) to the packet is passed to the host, instead of a copy of the packet. Removing this memcpy can give a performance improvement for some use cases, for example switching large packets between different VMs. However additional packet loss may be observed.

Note that the feature is disabled by default and must be explicitly enabled by setting the `dq-zero-copy` option

to `true` while specifying the `vhost-server-path` option as above. If you wish to split out the command into multiple commands as below, ensure `dq-zero-copy` is set before `vhost-server-path`:

```
$ ovs-vsctl set Interface dpdkvhostuserclient0 options:dq-zero-copy=true
$ ovs-vsctl set Interface dpdkvhostuserclient0 \
    options:vhost-server-path=/tmp/dpdkvhostclient0
```

The feature is only available to `dpdkvhostuserclient` port types.

A limitation exists whereby if packets from a vHost port with `dq-zero-copy=true` are destined for a `dpdk` type port, the number of tx descriptors (`n_txq_desc`) for that port must be reduced to a smaller number, 128 being the recommended value. This can be achieved by issuing the following command:

```
$ ovs-vsctl set Interface dpdkport options:n_txq_desc=128
```

Note: The sum of the tx descriptors of all `dpdk` ports the VM will send to should not exceed 128. For example, in case of a bond over two physical ports in `balance-tcp` mode, one must divide 128 by the number of links in the bond.

Refer to *DPDK Physical Port Queue Sizes* for more information.

The reason for this limitation is due to how the zero copy functionality is implemented. The vHost device's 'tx used vring', a virtio structure used for tracking used ie. sent descriptors, will only be updated when the NIC frees the corresponding mbuf. If we don't free the mbufs frequently enough, that vring will be starved and packets will no longer be processed. One way to ensure we don't encounter this scenario, is to configure `n_txq_desc` to a small enough number such that the 'mbuf free threshold' for the NIC will be hit more often and thus free mbufs more frequently. The value of 128 is suggested, but values of 64 and 256 have been tested and verified to work too, with differing performance characteristics. A value of 512 can be used too, if the virtio queue size in the guest is increased to 1024 (available to configure in QEMU versions v2.10 and greater). This value can be set like so:

```
$ qemu-system-x86_64 ... -chardev socket,id=char1,path=<sockpath>,server
-netdev type=vhost-user,id=mynet1,chardev=char1,vhostforce
-device virtio-net-pci,mac=00:00:00:00:00:01,netdev=mynet1,
tx_queue_size=1024
```

Because of this limitation, this feature is considered 'experimental'.

Further information can be found in the [DPDK documentation](#)

DPDK Ring Ports

Warning: DPDK ring interfaces cannot be used for guest communication and are retained mainly for backwards compatibility purposes. In nearly all cases, *vhost-user ports* are a better choice and should be used instead.

The DPDK datapath provides DPDK-backed ring ports that are implemented using DPDK's `librte_ring` library. For more information on this library, refer to the [DPDK documentation](#).

Important: To use any DPDK-backed interface, you must ensure your bridge is configured correctly. For more information, refer to *DPDK Bridges*.

Quick Example

This example demonstrates how to add a `dpdkr` port to an existing bridge called `br0`:

```
$ ovs-vsctl add-port br0 dpdkr0 -- set Interface dpdkr0 type=dpdkr
```

dpdkr

To use ring ports, you must first add said ports to the switch. Unlike *vhost-user ports*, ring port names must take a specific format, `dpdkrNN`, where `NN` is the port ID. For example:

```
$ ovs-vsctl add-port br0 dpdkr0 -- set Interface dpdkr0 type=dpdkr
```

Once the port has been added to the switch, they can be used by host processes. A sample loopback application - `test-dpdkr` - is included with Open vSwitch. To use this, run the following:

```
$ ./tests/test-dpdkr -c 1 -n 4 --proc-type=secondary -- -n 0
```

Further functionality would require developing your own application. Refer to the [DPDK documentation](#) for more information on how to do this.

Adding dpdkr ports to the guest

It is **not** recommended to use ring ports from guests. Historically, this was possible using a patched version of QEMU and the IVSHMEM feature provided with DPDK. However, this functionality was removed because:

- The IVSHMEM library was removed from DPDK in DPDK 16.11
- Support for IVSHMEM was never upstreamed to QEMU and has been publicly rejected by the QEMU community
- *vhost-user interfaces* are the de facto DPDK-based path to guests

DPDK Virtual Devices

DPDK provides drivers for both physical and virtual devices. Physical DPDK devices are added to OVS by specifying a valid PCI address in `dpdk-devargs`. Virtual DPDK devices which do not have PCI addresses can be added using a different format for `dpdk-devargs`.

Important: To use any DPDK-backed interface, you must ensure your bridge is configured correctly. For more information, refer to [DPDK Bridges](#).

Note: Not all DPDK virtual PMD drivers have been tested and verified to work.

New in version 2.7.0.

Quick Example

To add a virtual dpdk devices, the `dpdk-devargs` argument should be of the format `eth_<driver_name><x>`, where `x` is a unique identifier of your choice for the given port. For example to add a dpdk port that uses the `null` DPDK PMD driver, run:

```
$ ovs-vsctl add-port br0 null0 -- set Interface null0 type=dtpdk \
    options:dtpdk-devargs=eth_null0
```

Similarly, to add a dtpdk port that uses the af_packet DPDK PMD driver, run:

```
$ ovs-vsctl add-port br0 myeth0 -- set Interface myeth0 type=dtpdk \
    options:dtpdk-devargs=eth_af_packet0,iface=eth0
```

More information on the different types of virtual DPDK PMDs can be found in the [DPDK documentation](#).

PMD Threads

Poll Mode Driver (PMD) threads are the threads that do the heavy lifting for the DPDK datapath and perform tasks such as continuous polling of input ports for packets, classifying packets once received, and executing actions on the packets once they are classified.

PMD threads utilize Receive (Rx) and Transmit (Tx) queues, commonly known as *rxqs* and *txqs*. While Tx queue configuration happens automatically, Rx queues can be configured by the user. This can happen in one of two ways:

- For physical interfaces, configuration is done using the **ovs-appctl** utility.
- For virtual interfaces, configuration is done using the **ovs-appctl** utility, but this configuration must be reflected in the guest configuration (e.g. QEMU command line arguments).

The **ovs-appctl** utility also provides a number of commands for querying PMD threads and their respective queues. This, and all of the above, is discussed here.

PMD Thread Statistics

To show current stats:

```
$ ovs-appctl dpif-netdev/pmd-stats-show
```

To clear previous stats:

```
$ ovs-appctl dpif-netdev/pmd-stats-clear
```

Port/Rx Queue Assignment to PMD Threads

Correct configuration of PMD threads and the Rx queues they utilize is a requirement in order to achieve maximum performance. This is particularly true for enabling things like multiqueue for *physical* and *vhost-user* interfaces.

To show port/Rx queue assignment:

```
$ ovs-appctl dpif-netdev/pmd-rxq-show
```

Rx queues may be manually pinned to cores. This will change the default Rx queue assignment to PMD threads:

```
$ ovs-vsctl set Interface <iface> \
    other_config:pmd-rxq-affinity=<rxq-affinity-list>
```

where:

- **<rxq-affinity-list>** is a CSV list of **<queue-id>:<core-id>** values

For example:

```
$ ovs-vsctl set interface dpdk-p0 options:n_rxq=4 \
    other_config:pmd-rxq-affinity="0:3,1:7,3:8"
```

This will ensure there are 4 Rx queues and that these queues are configured like so:

- Queue #0 pinned to core 3
- Queue #1 pinned to core 7
- Queue #2 not pinned
- Queue #3 pinned to core 8

PMD threads on cores where Rx queues are *pinned* will become *isolated*. This means that this thread will only poll the *pinned* Rx queues.

Warning: If there are no *non-isolated* PMD threads, *non-pinned* RX queues will not be polled. Also, if the provided <core-id> is not available (e.g. the <core-id> is not in pmd-cpu-mask), the RX queue will not be polled by any PMD thread.

If pmd-rxq-affinity is not set for Rx queues, they will be assigned to PMDs (cores) automatically.

The algorithm used to automatically assign Rxqs to PMDs can be set by:

```
$ ovs-vsctl set Open_vSwitch . other_config:pmd-rxq-assign=<assignment>
```

By default, *cycles* assignment is used where the Rxqs will be ordered by their measured processing cycles, and then be evenly assigned in descending order to PMDs based on an up/down walk of the PMDs. For example, where there are five Rx queues and three cores - 3, 7, and 8 - available and the measured usage of core cycles per Rx queue over the last interval is seen to be:

- Queue #0: 30%
- Queue #1: 80%
- Queue #3: 60%
- Queue #4: 70%
- Queue #5: 10%

The Rx queues will be assigned to the cores in the following order:

```
Core 3: Q1 (80%) |
Core 7: Q4 (70%) | Q5 (10%)
Core 8: Q3 (60%) | Q0 (30%)
```

Alternatively, *roundrobin* assignment can be used, where the Rxqs are assigned to PMDs in a round-robin fashion. This algorithm was used by default prior to OVS 2.9. For example, given the following ports and queues:

- Port #0 Queue #0 (P0Q0)
- Port #0 Queue #1 (P0Q1)
- Port #1 Queue #0 (P1Q0)
- Port #1 Queue #1 (P1Q1)
- Port #1 Queue #2 (P1Q2)

The Rx queues may be assigned to the cores in the following order:

```
Core 3: P0Q0 | P1Q1  
Core 7: P0Q1 | P1Q2  
Core 8: P1Q0 |
```

To see the current measured usage history of PMD core cycles for each Rx queue:

```
$ ovs-appctl dpif-netdev/pmd-rxq-show
```

Note: A history of one minute is recorded and shown for each Rx queue to allow for traffic pattern spikes. Any changes in the Rx queue's PMD core cycles usage, due to traffic pattern or reconfig changes, will take one minute to be fully reflected in the stats.

Rx queue to PMD assignment takes place whenever there are configuration changes or can be triggered by using:

```
$ ovs-appctl dpif-netdev/pmd-rxq-rebalance
```

Changed in version 2.6.0: The `pmd-rxq-show` command was added in OVS 2.6.0.

Changed in version 2.9.0: Utilization-based allocation of Rx queues to PMDs and the `pmd-rxq-rebalance` command were added in OVS 2.9.0. Prior to this, allocation was round-robin and processing cycles were not taken into consideration.

In addition, the output of `pmd-rxq-show` was modified to include Rx queue utilization of the PMD as a percentage. Prior to this, tracking of stats was not available.

Automatic assignment of Port/Rx Queue to PMD Threads (experimental)

Cycle or utilization based allocation of Rx queues to PMDs gives efficient load distribution but it is not adaptive to change in traffic pattern occurring over the time. This causes uneven load among the PMDs which results in overall lower throughput.

To address this automatic load balancing of PMDs can be set by:

```
$ ovs-vsctl set open_vswitch . other_config:pmd-auto-lb="true"
```

If `pmd-auto-lb` is set to true AND cycle based assignment is enabled then auto load balancing of PMDs is enabled provided there are 2 or more non-isolated PMDs and at least one of these PMDs is polling more than one RX queue. So, following conditions need to be met to have Auto Load balancing enabled:

1. cycle based assignment of RX queues to PMD is enabled.
2. `pmd-auto-lb` is set to true.
3. There are two or more non-isolated PMDs present.
4. And at least one of the non-isolated PMD has more than one queue polling.

If any of above is not met PMD Auto Load Balancing is disabled.

Once auto load balancing is set, each non-isolated PMD measures the processing load for each of its associated queues every 10 seconds. If the aggregated PMD load reaches 95% for 6 consecutive intervals then PMD considers itself to be overloaded.

If any PMD is overloaded, a dry-run of the PMD assignment algorithm is performed by OVS main thread. The dry-run does NOT change the existing queue to PMD assignments.

If the resultant mapping of dry-run indicates an improved distribution of the load then the actual reassignment will be performed.

Note: PMD Auto Load Balancing doesn't currently work if queues are assigned cross NUMA as actual processing load could get worse after assignment as compared to what dry run predicts.

The minimum time between 2 consecutive PMD auto load balancing iterations can also be configured by:

```
$ ovs-vsctl set open_vswitch .\
    other_config:pmd-auto-lb-rebal-interval="<interval>"
```

where **<interval>** is a value in minutes. The default interval is 1 minute and setting it to 0 will also result in default value i.e. 1 min.

A user can use this option to avoid frequent trigger of Auto Load Balancing of PMDs. For e.g. set this (in min) such that it occurs once in few hours or a day or a week.

Note: In some scenarios it may not be desired to have Auto Load Balancing triggered. For example, if traffic profile for specific RX queue is changing dramatically very frequently which in turn thrashes CPU cache due to changes required in dpctl flows and EMC for newly added flows. In such scenarios user should configure rebalance interval accordingly to avoid frequent rebalancing happening.

Quality of Service (QoS)

It is possible to apply both ingress and egress limiting when using the DPDK datapath. These are referred to as *QoS* and *Rate Limiting*, respectively.

New in version 2.7.0.

QoS (Egress Policing)

Assuming you have a *vhost-user port* transmitting traffic consisting of packets of size 64 bytes, the following command would limit the egress transmission rate of the port to ~1,000,000 packets per second:

```
$ ovs-vsctl set port vhost-user0 qos=@newqos -- \
    --id=@newqos create qos type=egress-policer other_config:cir=46000000 \
    other_config:cbs=2048`
```

To examine the QoS configuration of the port, run:

```
$ ovs-appctl -t ovs-vswitchd qos/show vhost-user0
```

To clear the QoS configuration from the port and ovsdb, run:

```
$ ovs-vsctl destroy QoS vhost-user0 -- clear Port vhost-user0 qos
```

Refer to `vswitch.xml` for more details on egress policer.

Rate Limiting (Ingress Policing)

Assuming you have a *vhost-user port* receiving traffic consisting of packets of size 64 bytes, the following command would limit the reception rate of the port to ~1,000,000 packets per second:

```
$ ovs-vsctl set interface vhost-user0 ingress_policing_rate=368000 \
    ingress_policing_burst=1000`
```

To examine the ingress policer configuration of the port:

```
$ ovs-vsctl list interface vhost-user0
```

To clear the ingress policer configuration from the port:

```
$ ovs-vsctl set interface vhost-user0 ingress_policing_rate=0
```

Refer to `vswitch.xml` for more details on ingress policer.

Flow Control

Flow control is available for *DPDK physical ports*. For more information, refer to *Flow Control*.

pdump

New in version 2.6.0.

pdump allows you to listen on DPDK ports and view the traffic that is passing on them. To use this utility, one must have libpcap installed on the system. Furthermore, DPDK must be built with `CONFIG_RTE_LIBRTE_PDUMP=y` and `CONFIG_RTE_LIBRTE_PMD_PCAP=y`.

Warning: A performance decrease is expected when using a monitoring application like the DPDK pdump app.

To use pdump, simply launch OVS as usual, then navigate to the `app/pdump` directory in DPDK, make the application and run like so:

```
$ sudo ./build/app/dpdk-pdump -- \
    --pdump port=0,queue=0,rx-dev=/tmp/pkts.pcap \
    --server-socket-path=/usr/local/var/run/openvswitch
```

The above command captures traffic received on queue 0 of port 0 and stores it in `/tmp/pkts.pcap`. Other combinations of port numbers, queues numbers and pcap locations are of course also available to use. For example, to capture all packets that traverse port 0 in a single pcap file:

```
$ sudo ./build/app/dpdk-pdump -- \
    --pdump 'port=0,queue=*,rx-dev=/tmp/pkts.pcap,tx-dev=/tmp/pkts.pcap' \
    --server-socket-path=/usr/local/var/run/openvswitch
```

`server-socket-path` must be set to the value of `ovs_rundir()` which typically resolves to `/usr/local/var/run/openvswitch`.

Many tools are available to view the contents of the pcap file. One example is `tcpdump`. Issue the following command to view the contents of `pkts.pcap`:

```
$ tcpdump -r pkts.pcap
```

More information on the pdump app and its usage can be found in the [DPDK documentation](#).

Jumbo Frames

New in version 2.6.0.

By default, DPDK ports are configured with standard Ethernet MTU (1500B). To enable Jumbo Frames support for a DPDK port, change the Interface's `mtu_request` attribute to a sufficiently large value. For example, to add a *DPDK physical port* with an MTU of 9000, run:

```
$ ovs-vsctl add-port br0 dpdk-p0 -- set Interface dpdk-p0 type=dpdk \
    options:dpdk-devargs=0000:01:00.0 mtu_request=9000
```

Similarly, to change the MTU of an existing port to 6200, run:

```
$ ovs-vsctl set Interface dpdk-p0 mtu_request=6200
```

Some additional configuration is needed to take advantage of jumbo frames with *vHost User ports*:

- *Mergeable buffers* must be enabled for vHost User ports, as demonstrated in the QEMU command line snippet below:

```
-netdev type=vhost-user,id=mynet1,chardev=char0,vhostforce \
-device virtio-net-pci,mac=00:00:00:00:00:01,netdev=mynet1,mrg_rxbuf=on
```

- Where virtio devices are bound to the Linux kernel driver in a guest environment (i.e. interfaces are not bound to an in-guest DPDK driver), the MTU of those logical network interfaces must also be increased to a sufficiently large value. This avoids segmentation of Jumbo Frames received in the guest. Note that 'MTU' refers to the length of the IP packet only, and not that of the entire frame.

To calculate the exact MTU of a standard IPv4 frame, subtract the L2 header and CRC lengths (i.e. 18B) from the max supported frame size. So, to set the MTU for a 9018B Jumbo Frame:

```
$ ip link set eth1 mtu 9000
```

When Jumbo Frames are enabled, the size of a DPDK port's mbuf segments are increased, such that a full Jumbo Frame of a specific size may be accommodated within a single mbuf segment.

Jumbo frame support has been validated against 9728B frames, which is the largest frame size supported by Fortville NIC using the DPDK i40e driver, but larger frames and other DPDK NIC drivers may be supported. These cases are common for use cases involving East-West traffic only.

DPDK Device Memory Models

DPDK device memory can be allocated in one of two ways in OVS DPDK, **shared memory** or **per port memory**. The specifics of both are detailed below.

Shared Memory

By default OVS DPDK uses a shared memory model. This means that multiple ports can share the same mempool. For example when a port is added it will have a given MTU and socket ID associated with it. If a mempool has been created previously for an existing port that has the same MTU and socket ID, that mempool is used for both ports. If there is no existing mempool supporting these parameters then a new mempool is created.

Per Port Memory

In the per port memory model, mempools are created per device and are not shared. The benefit of this is a more transparent memory model where mempools will not be exhausted by other DPDK devices. However this comes at a potential increase in cost for memory dimensioning for a given deployment. Users should be aware of the memory requirements for their deployment before using this model and allocate the required hugepage memory.

Per port mempool support may be enabled via a global config value, ``per-port-memory``. Setting this to true enables the per port memory model for all DPDK devices in OVS:

```
$ ovs-vsctl set Open_vSwitch . other_config:per-port-memory=true
```

Important: This value should be set before setting `dpdk-init=true`. If set after `dpdk-init=true` then the daemon must be restarted to use `per-port-memory`.

Calculating Memory Requirements

The amount of memory required for a given mempool can be calculated by the **number mbufs in the mempool * mbuf size**.

Users should be aware of the following:

- The **number of mbufs** per mempool will differ between memory models.
- The **size of each mbuf** will be affected by the requested **MTU** size.

Important: An mbuf size in bytes is always larger than the requested MTU size due to alignment and rounding needed in OVS DPDK.

Below are a number of examples of memory requirement calculations for both shared and per port memory models.

Shared Memory Calculations

In the shared memory model the number of mbufs requested is directly affected by the requested MTU size as described in the table below.

MTU Size	Num MBUFS
1500 or greater	262144
Less than 1500	16384

Important: If a deployment does not have enough memory to provide 262144 mbufs then the requested amount is halved up until 16384.

Example 1

```
MTU = 1500 Bytes
Number of mbufs = 262144
Mbuf size = 3008 Bytes
Memory required = 262144 * 3008 = 788 MB
```

Example 2

```
MTU = 1800 Bytes
Number of mbufs = 262144
Mbuf size = 3008 Bytes
Memory required = 262144 * 3008 = 788 MB
```

Note: Assuming the same socket is in use for example 1 and 2 the same mempool would be shared.

Example 3

```
MTU = 6000 Bytes
Number of mbufs = 262144
Mbuf size = 7104 Bytes
Memory required = 262144 * 7104 = 1862 MB
```

Example 4

```
MTU = 9000 Bytes
Number of mbufs = 262144
Mbuf size = 10176 Bytes
Memory required = 262144 * 10176 = 2667 MB
```

Per Port Memory Calculations

The number of mbufs requested in the per port model is more complicated and accounts for multiple dynamic factors in the datapath and device configuration.

A rough estimation of the number of mbufs required for a port is:

```
packets required to fill the device rxqs +
packets that could be stuck on other ports txqs +
packets on the pmd threads +
additional corner case memory.
```

The algorithm in OVS used to calculate this is as follows:

```
requested number of rxqs * requested rxq size +
requested number of txqs * requested txq size +
min(RTE_MAX_LCORE, requested number of rxqs) * netdev_max_burst +
MIN_NB_MBUF.
```

where:

- **requested number of rxqs:** Number of requested receive queues for a device.
- **requested rxq size:** The number of descriptors requested for a rx queue.
- **requested number of txqs:** Number of requested transmit queues for a device. Calculated as the number of PMDs configured +1.
- **requested txq size:** the number of descriptors requested for a tx queue.
- **min(RTE_MAX_LCORE, requested number of rxqs):** Compare the maximum number of lcores supported by DPDK to the number of requested receive queues for the device and use the variable of lesser value.
- **NETDEV_MAX_BURST:** Maximum number of of packets in a burst, defined as 32.
- **MIN_NB_MBUF:** Additional memory for corner case, defined as 16384.

For all examples below assume the following values:

- requested_rxq_size = 2048
- requested_txq_size = 2048
- RTE_MAX_LCORE = 128
- netdev_max_burst = 32
- MIN_NB_MBUF = 16384

Example 1: (1 rxq, 1 PMD, 1500 MTU)

```
MTU = 1500
Number of mbufs = (1 * 2048) + (2 * 2048) + (1 * 32) + (16384) = 22560
Mbuf size = 3008 Bytes
Memory required = 22560 * 3008 = 67 MB
```

Example 2: (1 rxq, 2 PMD, 6000 MTU)

```
MTU = 6000
Number of mbufs = (1 * 2048) + (3 * 2048) + (1 * 32) + (16384) = 24608
Mbuf size = 7104 Bytes
Memory required = 24608 * 7104 = 175 MB
```

Example 3: (2 rxq, 2 PMD, 9000 MTU)

```
MTU = 9000
Number of mbufs = (2 * 2048) + (3 * 2048) + (1 * 32) + (16384) = 26656
Mbuf size = 10176 Bytes
Memory required = 26656 * 10176 = 271 MB
```

4.1.10 OVS-on-Hyper-V Design

This document provides details of the effort to develop Open vSwitch on Microsoft Hyper-V. This document should give enough information to understand the overall design.

Note: The userspace portion of the OVS has been ported to Hyper-V in a separate effort, and committed to the openvswitch repo. This document will mostly emphasize on the kernel driver, though we touch upon some of the aspects of userspace as well.

Background Info

Microsoft’s hypervisor solution - Hyper-V¹ implements a virtual switch that is extensible and provides opportunities for other vendors to implement functional extensions². The extensions need to be implemented as NDIS drivers that bind within the extensible switch driver stack provided. The extensions can broadly provide the functionality of monitoring, modifying and forwarding packets to destination ports on the Hyper-V extensible switch. Correspondingly, the extensions can be categorized into the following types and provide the functionality noted:

- Capturing extensions: monitoring packets
- Filtering extensions: monitoring, modifying packets
- Forwarding extensions: monitoring, modifying, forwarding packets

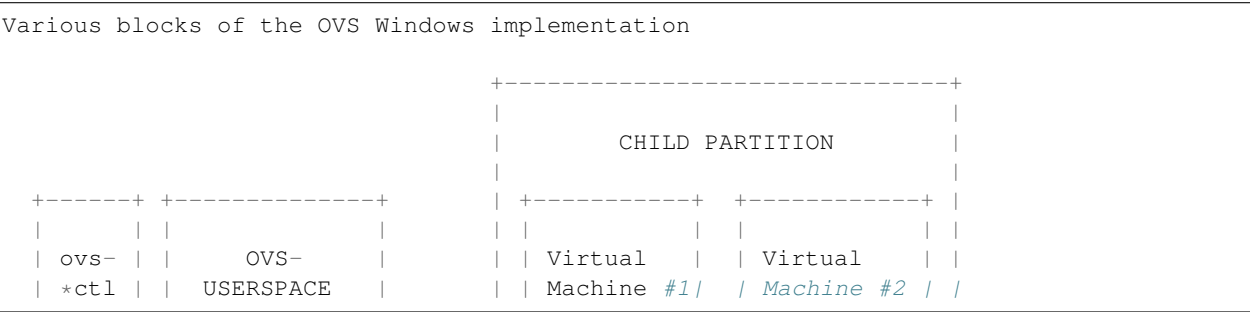
As can be expected, the kernel portion (datapath) of OVS on Hyper-V solution will be implemented as a forwarding extension.

In Hyper-V, the virtual machine is called the Child Partition. Each VIF or physical NIC on the Hyper-V extensible switch is attached via a port. Each port is both on the ingress path or the egress path of the switch. The ingress path is used for packets being sent out of a port, and egress is used for packet being received on a port. By design, NDIS provides a layered interface. In this layered interface, higher level layers call into lower level layers, in the ingress path. In the egress path, it is the other way round. In addition, there is a object identifier (OID) interface for control operations Eg. addition of a port. The workflow for the calls is similar in nature to the packets, where higher level layers call into the lower level layers. A good representational diagram of this architecture is in⁴.

Windows Filtering Platform (WFP)⁵ is a platform implemented on Hyper-V that provides APIs and services for filtering packets. WFP has been utilized to filter on some of the packets that OVS is not equipped to handle directly. More details in later sections.

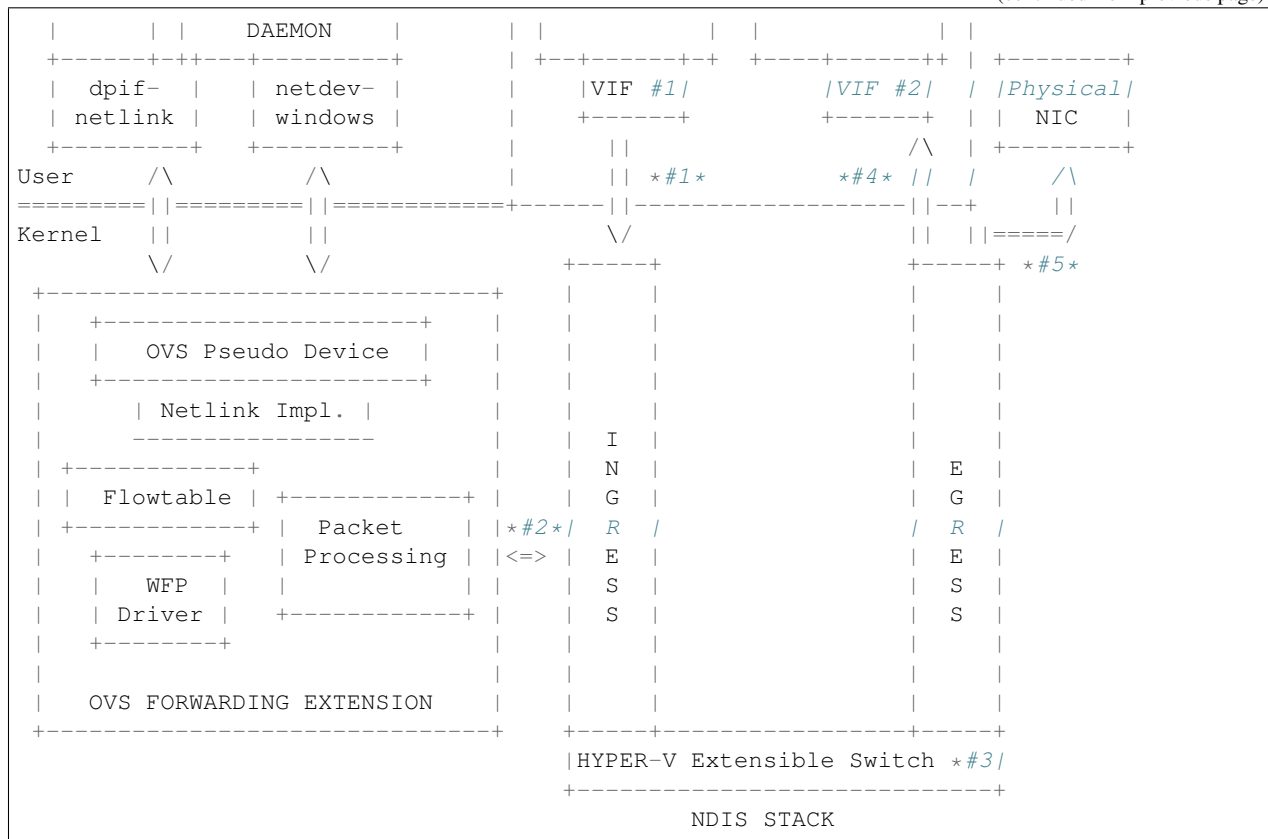
IP Helper⁶ is a set of API available on Hyper-V to retrieve information related to the network configuration information on the host machine. IP Helper has been used to retrieve some of the configuration information that OVS needs.

Design



¹ Hyper-V Extensible Switch <https://msdn.microsoft.com/windows/hardware/drivers/network/hyper-v-extensible-switch>
² Hyper-V Extensible Switch Extensions <https://msdn.microsoft.com/windows/hardware/drivers/network/hyper-v-extensible-switch-extensions>
⁴ Hyper-V Extensible Switch Components <https://msdn.microsoft.com/windows/hardware/drivers/network/hyper-v-extensible-switch-components>
⁵ Windows Filtering Platform [https://msdn.microsoft.com/en-us/library/windows/desktop/aa366510\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/desktop/aa366510(v=vs.85).aspx)
⁶ IP Helper <https://msdn.microsoft.com/windows/hardware/drivers/network/ip-helper>

(continued from previous page)



This diagram shows the various blocks involved in the OVS Windows implementation, along with some of the components available in the NDIS stack, and also the virtual machines. The workflow of a packet being transmitted from a VIF out and into another VIF and to a physical NIC is also shown. Later on in this section, we will discuss the flow of a packet at a high level.

The figure gives a general idea of where the OVS userspace and the kernel components fit in, and how they interface with each other.

The kernel portion (datapath) of OVS on Hyper-V solution has been implemented as a forwarding extension roughly implementing the following sub-modules/functionality. Details of each of these sub-components in the kernel are contained in later sections:

- Interfacing with the NDIS stack
- Netlink message parser
- Netlink sockets
- Switch/Datapath management
- Interfacing with userspace portion of the OVS solution to implement the necessary functionality that userspace needs
- Port management
- Flowtable/Actions/packet forwarding
- Tunneling
- Event notifications

The datapath for the OVS on Linux is a kernel module, and cannot be directly ported since there are significant differences in architecture even though the end functionality provided would be similar. Some examples of the differences are:

- Interfacing with the NDIS stack to hook into the NDIS callbacks for functionality such as receiving and sending packets, packet completions, OIDs used for events such as a new port appearing on the virtual switch.
- Interface between the userspace and the kernel module.
- Event notifications are significantly different.
- The communication interface between DPIF and the kernel module need not be implemented in the way OVS on Linux does. That said, it would be advantageous to have a similar interface to the kernel module for reasons of readability and maintainability.
- Any licensing issues of using Linux kernel code directly.

Due to these differences, it was a straightforward decision to develop the datapath for OVS on Hyper-V from scratch rather than porting the one on Linux. A re-development focused on the following goals:

- Adhere to the existing requirements of userspace portion of OVS (such as `ovs-vswitchd`), to minimize changes in the userspace workflow.
- Fit well into the typical workflow of a Hyper-V extensible switch forwarding extension.

The userspace portion of the OVS solution is mostly POSIX code, and not very Linux specific. Majority of the userspace code does not interface directly with the kernel datapath and was ported independently of the kernel datapath effort.

As explained in the OVS porting design document⁷, DPIF is the portion of userspace that interfaces with the kernel portion of the OVS. The interface that each DPIF provider has to implement is defined in `dpif-provider.h`³. Though each platform is allowed to have its own implementation of the DPIF provider, it was found, via community feedback, that it is desired to share code whenever possible. Thus, the DPIF provider for OVS on Hyper-V shares code with the DPIF provider on Linux. This interface is implemented in `dpif-netlink.c`.

We'll elaborate more on kernel-userspace interface in a dedicated section below. Here it suffices to say that the DPIF provider implementation for Windows is netlink-based and shares code with the Linux one.

Kernel Module (Datapath)

Interfacing with the NDIS Stack

For each virtual switch on Hyper-V, the OVS extensible switch extension can be enabled/disabled. We support enabling the OVS extension on only one switch. This is consistent with using a single datapath in the kernel on Linux. All the physical adapters are connected as external adapters to the extensible switch.

When the OVS switch extension registers itself as a filter driver, it also registers callbacks for the switch/port management and datapath functions. In other words, when a switch is created on the Hyper-V root partition (host), the extension gets an activate callback upon which it can initialize the data structures necessary for OVS to function. Similarly, there are callbacks for when a port gets added to the Hyper-V switch, and an External Network adapter or a VM Network adapter is connected/disconnected to the port. There are also callbacks for when a VIF (NIC of a child partition) send out a packet, or a packet is received on an external NIC.

As shown in the figures, an extensible switch extension gets to see a packet sent by the VM (VIF) twice - once on the ingress path and once on the egress path. Forwarding decisions are to be made on the ingress path. Correspondingly, we will be hooking onto the following interfaces:

⁷ How to Port Open vSwitch to New Software or Hardware *Porting Open vSwitch to New Software or Hardware*

³ DPIF Provider http://openvswitch.sourcearchive.com/documentation/1.1.0-1/dpif-provider_8h_source.html

- Ingress send indication: intercept packets for performing flow based forwarding. This includes straight forwarding to output ports. Any packet modifications needed to be performed are done here either inline or by creating a new packet. A forwarding action is performed as the flow actions dictate.
- Ingress completion indication: cleanup and free packets that we generated on the ingress send path, pass-through for packets that we did not generate.
- Egress receive indication: pass-through.
- Egress completion indication: pass-through.

Interfacing with OVS Userspace

We have implemented a pseudo device interface for letting OVS userspace talk to the OVS kernel module. This is equivalent to the typical character device interface on POSIX platforms where we can register custom functions for read, write and ioctl functionality. The pseudo device supports a whole bunch of ioctls that netdev and DPF on OVS userspace make use of.

Netlink Message Parser

The communication between OVS userspace and OVS kernel datapath is in the form of Netlink messages^{1, 8}. More details about this are provided below. In the kernel, a full fledged netlink message parser has been implemented along the lines of the netlink message parser in OVS userspace. In fact, a lot of the code is ported code.

On the lines of `struct ofpbuf` in OVS userspace, a managed buffer has been implemented in the kernel datapath to make it easier to parse and construct netlink messages.

Netlink Sockets

On Linux, OVS userspace utilizes netlink sockets to pass back and forth netlink messages. Since much of userspace code including DPF provider in `dpif-netlink.c` (formerly `dpif-linux.c`) has been reused, pseudo-netlink sockets have been implemented in OVS userspace. As it is known, Windows lacks native netlink socket support, and also the socket family is not extensible either. Hence it is not possible to provide a native implementation of netlink socket. We emulate netlink sockets in `lib/netlink-socket.c` and support all of the `nl_*` APIs to higher levels. The implementation opens a handle to the pseudo device for each netlink socket. Some more details on this topic are provided in the userspace section on netlink sockets.

Typical netlink semantics of read message, write message, dump, and transaction have been implemented so that higher level layers are not affected by the netlink implementation not being native.

Switch/Datapath Management

As explained above, we hook onto the management callback functions in the NDIS interface for when to initialize the OVS data structures, flow tables etc. Some of this code is also driven by OVS userspace code which sends down ioctls for operations like creating a tunnel port etc.

⁸ Netlink <https://en.wikipedia.org/wiki/Netlink>

Port Management

As explained above, we hook onto the management callback functions in the NDIS interface to know when a port is added/connected to the Hyper-V switch. We use these callbacks to initialize the port related data structures in OVS. Also, some of the ports are tunnel ports that don't exist on the Hyper-V switch and get added from OVS userspace.

In order to identify a Hyper-V port, we use the value of 'FriendlyName' field in each Hyper-V port. We call this the "OVS-port-name". The idea is that OVS userspace sets 'OVS-port-name' in each Hyper-V port to the same value as the 'name' field of the 'Interface' table in OVSDb. When OVS userspace calls into the kernel datapath to add a port, we match the name of the port with the 'OVS-port-name' of a Hyper-V port.

We maintain separate hash tables, and separate counters for ports that have been added from the Hyper-V switch, and for ports that have been added from OVS userspace.

Flowtable/Actions/Packet Forwarding

The flowtable and flow actions based packet forwarding is the core of the OVS datapath functionality. For each packet on the ingress path, we consult the flowtable and execute the corresponding actions. The actions can be limited to simple forwarding to a particular destination port(s), or more commonly involves modifying the packet to insert a tunnel context or a VLAN ID, and thereafter forwarding to the external port to send the packet to a destination host.

Tunneling

We make use of the Internal Port on a Hyper-V switch for implementing tunneling. The Internal Port is a virtual adapter that is exposed on the Hyper-V host, and connected to the Hyper-V switch. Basically, it is an interface between the host and the virtual switch. The Internal Port acts as the Tunnel end point for the host (aka VTEP), and holds the VTEP IP address.

Tunneling ports are not actual ports on the Hyper-V switch. These are virtual ports that OVS maintains and while executing actions, if the output is a tunnel port, we short circuit by performing the encapsulation action based on the tunnel context. The encapsulated packet gets forwarded to the external port, and appears to the outside world as though it was set from the VTEP.

Similarly, when a tunneled packet enters the OVS from the external port bound to the internal port (VTEP), and if yes, we short circuit the path, and directly forward the inner packet to the destination port (mostly a VIF, but dictated by the flow). We leverage the Windows Filtering Platform (WFP) framework to be able to receive tunneled packets that cannot be decapsulated by OVS right away. Currently, fragmented IP packets fall into that category, and we leverage the code in the host IP stack to reassemble the packet, and performing decapsulation on the reassembled packet.

We'll also be using the IP helper library to provide us IP address and other information corresponding to the Internal port.

Event Notifications

The pseudo device interface described above is also used for providing event notifications back to OVS userspace. A shared memory/overlapped IO model is used.

Userspace Components

The userspace portion of the OVS solution is mostly POSIX code, and not very Linux specific. Majority of the userspace code does not interface directly with the kernel datapath and was ported independently of the kernel datapath effort.

In this section, we cover the userspace components that interface with the kernel datapath.

As explained earlier, OVS on Hyper-V shares the DPIF provider implementation with Linux. The DPIF provider on Linux uses netlink sockets and netlink messages. Netlink sockets and messages are extensively used on Linux to exchange information between userspace and kernel. In order to satisfy these dependencies, netlink socket (pseudo and non-native) and netlink messages are implemented on Hyper-V.

The following are the major advantages of sharing DPIF provider code:

1. Maintenance is simpler:

Any change made to the interface defined in `dpif-provider.h` need not be propagated to multiple implementations. Also, developers familiar with the Linux implementation of the DPIF provider can easily ramp on the Hyper-V implementation as well.

2. Netlink messages provides inherent advantages:

Netlink messages are known for their extensibility. Each message is versioned, so the provided data structures offer a mechanism to perform version checking and forward/backward compatibility with the kernel module.

Netlink Sockets

As explained in other sections, an emulation of netlink sockets has been implemented in `lib/netlink-socket.c` for Windows. The implementation creates a handle to the OVS pseudo device, and emulates netlink socket semantics of receive message, send message, dump, and transact. Most of the `nl_*` functions are supported.

The fact that the implementation is non-native manifests in various ways. One example is that PID for the netlink socket is not automatically assigned in userspace when a handle is created to the OVS pseudo device. There's an extra command (defined in `OvsDpInterfaceExt.h`) that is used to grab the PID generated in the kernel.

DPIF Provider

As has been mentioned in earlier sections, the netlink socket and netlink message based DPIF provider on Linux has been ported to Windows.

Most of the code is common. Some divergence is in the code to receive packets. The Linux implementation uses `epoll()`⁹ which is not natively supported on Windows.

netdev-windows

We have a Windows implementation of the interface defined in `lib/netdev-provider.h`. The implementation provides functionality to get extended information about an interface. It is limited in functionality compared to the Linux implementation of the netdev provider and cannot be used to add any interfaces in the kernel such as a tap interface or to send/receive packets. The netdev-windows implementation uses the datapath interface extensions defined in `datapath-windows/include/OvsDpInterfaceExt.h`.

Powershell Extensions to Set `OVS-port-name`

As explained in the section on “Port management”, each Hyper-V port has a ‘FriendlyName’ field, which we call as the “OVS-port-name” field. We have implemented powershell command extensions to be able to set the “OVS-port-name” of a Hyper-V port.

⁹ `epoll` <https://en.wikipedia.org/wiki/Epoll>

Kernel-Userspace Interface

openvswitch.h and OvsDpInterfaceExt.h

Since the DPF provider is shared with Linux, the kernel datapath provides the same interface as the Linux datapath. The interface is defined in `datapath/linux/compat/include/linux/openvswitch.h`. Derivatives of this interface file are created during OVS userspace compilation. The derivative for the kernel datapath on Hyper-V is provided in `datapath-windows/include/OvsDpInterface.h`.

That said, there are Windows specific extensions that are defined in the interface file `datapath-windows/include/OvsDpInterfaceExt.h`.

Flow of a Packet

Figure 2 shows the numbered steps in which a packets gets sent out of a VIF and is forwarded to another VIF or a physical NIC. As mentioned earlier, each VIF is attached to the switch via a port, and each port is both on the ingress and egress path of the switch, and depending on whether a packet is being transmitted or received, one of the paths gets used. In the figure, each step *n* is annotated as #*n*

The steps are as follows:

1. When a packet is sent out of a VIF or an physical NIC or an internal port, the packet is part of the ingress path.
2. The OVS kernel driver gets to intercept this packet.
 - (a) OVS looks up the flows in the flowtable for this packet, and executes the corresponding action.
 - (b) If there is not action, the packet is sent up to OVS userspace to examine the packet and figure out the actions.
 - (c) Userspace executes the packet by specifying the actions, and might also insert a flow for such a packet in the future.
 - (d) The destination ports are added to the packet and sent down to the Hyper- V switch.
3. The Hyper-V forwards the packet to the destination ports specified in the packet, and sends it out on the egress path.
4. The packet gets forwarded to the destination VIF.
5. It might also get forwarded to a physical NIC as well, if the physical NIC has been added as a destination port by OVS.

Build/Deployment

The userspace components added as part of OVS Windows implementation have been integrated with autoconf, and can be built using the steps mentioned in the BUILD.Windows file. Additional targets need to be specified to make.

The OVS kernel code is part of a Visual Studio 2013 solution, and is compiled from the IDE. There are plans in the future to move this to a compilation mode such that we can compile it without an IDE as well.

Once compiled, we have an install script that can be used to load the kernel driver.

References

4.1.11 Language Bindings

Bindings exist for Open vSwitch in a variety of languages.

Official Bindings

Python

The Python bindings are part of the [Open vSwitch package](#). You can install the bindings using `pip`:

```
$ pip install ovs
```

Third-Party Bindings

Lua

- [LJIT2ovs](#): LuaJIT binding for Open vSwitch

Go

- [go-odp](#): A Go library to control the Open vSwitch in-kernel datapath

4.1.12 Testing

It is possible to test Open vSwitch using both tooling provided with Open vSwitch and using a variety of third party tooling.

Built-in Tooling

Open vSwitch provides a number of different test suites and other tooling for validating basic functionality of OVS. Before running any of the tests described here, you must bootstrap, configure and build Open vSwitch as described in *[Open vSwitch on Linux, FreeBSD and NetBSD](#)*. You do not need to install Open vSwitch or to build or load the kernel module to run these test suites. You do not need supervisor privilege to run these test suites.

Unit Tests

Open vSwitch includes a suite of self-tests. Before you submit patches upstream, we advise that you run the tests and ensure that they pass. If you add new features to Open vSwitch, then adding tests for those features will ensure your features don't break as developers modify other areas of Open vSwitch.

To run all the unit tests in Open vSwitch, one at a time, run:

```
$ make check
```

This takes under 5 minutes on a modern desktop system.

To run all the unit tests in Open vSwitch in parallel, run:

```
$ make check TESTSUITEFLAGS=-j8
```

You can run up to eight threads. This takes under a minute on a modern 4-core desktop system.

To see a list of all the available tests, run:


```
$ make check TESTSUITEFLAGS=--list
```

To run only a subset of tests, e.g. test 123 and tests 477 through 484, run:

```
$ make check TESTSUITEFLAGS='123 477-484'
```

Tests do not have inter-dependencies, so you may run any subset.

To run tests matching a keyword, e.g. `ovsdb`, run:

```
$ make check TESTSUITEFLAGS='-k ovsdb'
```

To see a complete list of test options, run:

```
$ make check TESTSUITEFLAGS=--help
```

The results of a testing run are reported in `tests/testsuite.log`. Report test failures as bugs and include the `testsuite.log` in your report.

Note: Sometimes a few tests may fail on some runs but not others. This is usually a bug in the testsuite, not a bug in Open vSwitch itself. If you find that a test fails intermittently, please report it, since the developers may not have noticed. You can make the testsuite automatically rerun tests that fail, by adding `RECHECK=yes` to the make command line, e.g.:

```
$ make check TESTSUITEFLAGS=-j8 RECHECK=yes
```

Coverage

If the build was configured with `--enable-coverage` and the `lcov` utility is installed, you can run the testsuite and generate a code coverage report by using the `check-lcov` target:

```
$ make check-lcov
```

All the same options are available via `TESTSUITEFLAGS`. For example:

```
$ make check-lcov TESTSUITEFLAGS='-j8 -k ovn'
```

Valgrind

If you have `valgrind` installed, you can run the testsuite under `valgrind` by using the `check-valgrind` target:

```
$ make check-valgrind
```

When you do this, the “valgrind” results for test `<N>` are reported in files named `tests/testsuite.dir/<N>/valgrind.*`.

To test the testsuite of kernel datapath under `valgrind`, you can use the `check-kernel-valgrind` target and find the “valgrind” results under directory `tests/system-kmod-testsuite.dir/`.

All the same options are available via `TESTSUITEFLAGS`.

Hint: You may find that the valgrind results are easier to interpret if you put `-q` in `~/ .valgrindrc`, since that reduces the amount of output.

OFTest

OFTest is an OpenFlow protocol testing suite. Open vSwitch includes a Makefile target to run OFTest with Open vSwitch in “dummy mode”. In this mode of testing, no packets travel across physical or virtual networks. Instead, Unix domain sockets stand in as simulated networks. This simulation is imperfect, but it is much easier to set up, does not require extra physical or virtual hardware, and does not require supervisor privileges.

To run OFTest with Open vSwitch, you must obtain a copy of OFTest and install its prerequisites. You need a copy of OFTest that includes commit 406614846c5 (make ovs-dummy platform work again). This commit was merged into the OFTest repository on Feb 1, 2013, so any copy of OFTest more recent than that should work. Testing OVS in dummy mode does not require root privilege, so you may ignore that requirement.

Optionally, add the top-level OFTest directory (containing the `oft` program) to your `$PATH`. This slightly simplifies running OFTest later.

To run OFTest in dummy mode, run the following command from your Open vSwitch build directory:

```
$ make check-oftest OFT=<oft-binary>
```

where `<oft-binary>` is the absolute path to the `oft` program in OFTest. If you added “oft” to your `$PATH`, you may omit the `OFT` variable assignment

By default, `check-oftest` passes `oft` just enough options to enable dummy mode. You can use `OFTFLAGS` to pass additional options. For example, to run just the `basic.Echo` test instead of all tests (the default) and enable verbose logging, run:

```
$ make check-oftest OFT=<oft-binary> OFTFLAGS='--verbose -T basic.Echo'
```

If you use OFTest that does not include commit 4d1f3eb2c792 (oft: change default port to 6653), merged into the OFTest repository in October 2013, then you need to add an option to use the IETF-assigned controller port:

```
$ make check-oftest OFT=<oft-binary> OFTFLAGS='--port=6653'
```

Interpret OFTest results cautiously. Open vSwitch can fail a given test in OFTest for many reasons, including bugs in Open vSwitch, bugs in OFTest, bugs in the “dummy mode” integration, and differing interpretations of the OpenFlow standard and other standards.

Note: Open vSwitch has not been validated against OFTest. Report test failures that you believe to represent bugs in Open vSwitch. Include the precise versions of Open vSwitch and OFTest in your bug report, plus any other information needed to reproduce the problem.

Ryu

Ryu is an OpenFlow controller written in Python that includes an extensive OpenFlow testsuite. Open vSwitch includes a Makefile target to run Ryu in “dummy mode”. See [OFTest](#) above for an explanation of dummy mode.

To run Ryu tests with Open vSwitch, first read and follow the instructions under **Testing** above. Second, obtain a copy of Ryu, install its prerequisites, and build it. You do not need to install Ryu (some of the tests do not get installed, so it does not help).

To run Ryu tests, run the following command from your Open vSwitch build directory:

```
$ make check-ryu RYUDIR=<ryu-source-dir>
```

where `<ryu-source-dir>` is the absolute path to the root of the Ryu source distribution. The default `<ryu-source-dir>` is `$srcdir/../../ryu` where `$srcdir` is your Open vSwitch source directory. If this is correct, omit `RYUDIR`

Note: Open vSwitch has not been validated against Ryu. Report test failures that you believe to represent bugs in Open vSwitch. Include the precise versions of Open vSwitch and Ryu in your bug report, plus any other information needed to reproduce the problem.

Datapath testing

Open vSwitch includes a suite of tests specifically for datapath functionality, which can be run against the userspace or kernel datapaths. If you are developing datapath features, it is recommended that you use these tests and build upon them to verify your implementation.

The datapath tests make some assumptions about the environment. They must be run under root privileges on a Linux system with support for network namespaces. For ease of use, the OVS source tree includes a vagrant box to invoke these tests. Running the tests inside Vagrant provides kernel isolation, protecting your development host from kernel panics or configuration conflicts in the testsuite. If you wish to run the tests without using the vagrant box, there are further instructions below.

Vagrant

Important: Requires Vagrant (version 1.7.0 or later) and a compatible hypervisor

Note: You must bootstrap and configure the sources (see `doc:intro/install/general`) before you run the steps described here.

A Vagrantfile is provided allowing to compile and provision the source tree as found locally in a virtual machine using the following command:

```
$ vagrant up
```

This will bring up a Fedora 23 VM by default. If you wish to use a different box or a vagrant backend not supported by the default box, the `Vagrantfile` can be modified to use a different box as base.

The VM can be reprovisioned at any time:

```
$ vagrant provision
```

OVS out-of-tree compilation environment can be set up with:

```
$ ./boot.sh
$ vagrant provision --provision-with configure_ovs,build_ovs
```

This will set up an out-of-tree build environment inside the VM in `/root/build`. The source code can be found in `/vagrant`.

To recompile and reinstall OVS in the VM using RPM:

```
$ ./boot.sh
$ vagrant provision --provision-with configure_ovs,install_rpm
```

Two provisioners are included to run system tests with the OVS kernel module or with a userspace datapath. These tests are different from the self-tests mentioned above. To run them:

```
$ ./boot.sh
$ vagrant provision --provision-with \
    configure_ovs,test_ovs_kmod,test_ovs_system_userspace
```

The results of the testsuite reside in the VM root user's home directory:

```
$ vagrant ssh
$ sudo -s
$ cd /root/build
$ ls tests/system*
```

Native

The datapath testsuite as invoked by Vagrant above may also be run manually on a Linux system with root privileges. Make sure, no other Open vSwitch instance is running on the test suite. These tests may take several minutes to complete, and cannot be run in parallel.

Userspace datapath

To invoke the datapath testsuite with the userspace datapath, run:

```
$ make check-system-userspace
```

The results of the testsuite are in `tests/system-userspace-testsuite.dir`.

DPDK datapath

To test *Open vSwitch with DPDK* (i.e., the build was configured with `--with-dpdk`, the DPDK is installed), run the testsuite and generate a report by using the `check-dpdk` target:

```
# make check-dpdk
```

or if you are not a root, but a sudo user:

```
$ sudo -E make check-dpdk
```

To see a list of all the available tests, run:

```
# make check-dpdk TESTSUITEFLAGS=--list
```

These tests support a [DPDK supported NIC](#). The tests operate on a wider set of environments, for instance, when a virtual port is used. They do require proper DPDK variables (`DPDK_DIR` and `DPDK_BUILD`). Moreover you need to have root privileges to load the required modules and to bind the NIC to the DPDK-compatible driver.

All tests are skipped if no hugepages are configured. User must look into the DPDK manual to figure out how to [Configure hugepages](#). The phy test will skip if no compatible physical device is available.

Kernel datapath

Make targets are also provided for testing the Linux kernel module. Note that these tests operate by inserting modules into the running Linux kernel, so if the tests are able to trigger a bug in the OVS kernel module or in the upstream kernel then the kernel may panic.

To run the testsuite against the kernel module which is currently installed on your system, run:

```
$ make check-kernel
```

To install the kernel module from the current build directory and run the testsuite against that kernel module:

```
$ make check-kmod
```

The results of the testsuite are in `tests/system-kmod-testsuite.dir`.

Static Code Analysis

Static Analysis is a method of debugging Software by examining code rather than actually executing it. This can be done through ‘scan-build’ commandline utility which internally uses clang (or) gcc to compile the code and also invokes a static analyzer to do the code analysis. At the end of the build, the reports are aggregated in to a common folder and can later be analyzed using ‘scan-view’.

Open vSwitch includes a Makefile target to trigger static code analysis:

```
$ ./boot.sh
$ ./configure CC=clang # clang
# or
$ ./configure CC=gcc CFLAGS="-std=gnu99" # gcc
$ make clang-analyze
```

You should invoke scan-view to view analysis results. The last line of output from `clang-analyze` will list the command (containing results directory) that you should invoke to view the results on a browser.

Continuous Integration with Travis CI

A `.travis.yml` file is provided to automatically build Open vSwitch with various build configurations and run the testsuite using Travis CI. Builds will be performed with gcc, sparse and clang with the `-Werror` compiler flag included, therefore the build will fail if a new warning has been introduced.

The CI build is triggered via git push (regardless of the specific branch) or pull request against any Open vSwitch GitHub repository that is linked to travis-ci.

Instructions to setup travis-ci for your GitHub repository:

1. Go to <https://travis-ci.org/> and sign in using your GitHub ID.
2. Go to the “Repositories” tab and enable the ovs repository. You may disable builds for pushes or pull requests.
3. In order to avoid forks sending build failures to the upstream mailing list, the notification email recipient is encrypted. If you want to receive email notification for build failures, replace the the encrypted string:
 - (a) Install the travis-ci CLI (Requires ruby >=2.0): `gem install travis`
 - (b) In your Open vSwitch repository: `travis encrypt mylist@mydomain.org`
 - (c) Add/replace the notifications section in `.travis.yml` and fill in the secure string as returned by travis encrypt:

```
notifications:
  email:
    recipients:
      - secure: "....."
```

Note: You may remove/omit the notifications section to fall back to default notification behaviour which is to send an email directly to the author and committer of the failing commit. Note that the email is only sent if the author/committer have commit rights for the particular GitHub repository.

4. Pushing a commit to the repository which breaks the build or the testsuite will now trigger a email sent to mylist@mydomain.org

vsperf

The vsperf project aims to develop a vSwitch test framework that can be used to validate the suitability of different vSwitch implementations in a telco deployment environment. More information can be found on the [OPNFV wiki](#).

Proof of Concepts

Proof of Concepts are documentation materialized into Ansible recipes executed in VirtualBox or Libvirt environments orchestrated by Vagrant. Proof of Concepts allow developers to create small virtualized setups that demonstrate how certain Open vSwitch features are intended to work avoiding user introduced errors by overlooking instructions. Proof of Concepts are also helpful when integrating with thirdparty software, because standard unit tests with make check are limited.

Vagrant by default uses VirtualBox provider. However, if Libvirt is your choice of virtualization technology, then you can use it by installing Libvirt plugin:

```
$ vagrant plugin install vagrant-libvirt
```

And then appending `--provider=libvirt` flag to vagrant commands.

The host where Vagrant runs does not need to have any special software installed besides vagrant, virtualbox (or libvirt and libvirt-dev) and ansible.

The following Proof of Concepts are supported:

Builders

This particular Proof of Concept demonstrates integration with Debian and RPM packaging tools:

```
$ cd ./poc/builders
$ vagrant up
```

Once that command finished you can get packages from `/var/www/html` directory. Since those hosts are also configured as repositories then you can add them to `/etc/apt/sources.list.d` or `/etc/yum.repos.d` configuration files on another host to retrieve packages with yum or apt-get.

When you have made changes to OVS source code and want to rebuild packages run:

```
$ git commit -a
$ vagrant rsync && vagrant provision
```

Whenever packages are rebuilt the Open vSwitch release number increases by one and you can simply upgrade Open vSwitch by running `yum` or `apt-get update` commands.

Once you are done with experimenting you can tear down setup with:

```
$ vagrant destroy
```

Sometimes deployment of Proof of Concept may fail, if, for example, VMs don't have network reachability to the Internet.

4.1.13 Tracing packets inside Open vSwitch

Open vSwitch (OVS) is a programmable software switch that can execute actions at per packet level. This document explains how to use the tracing tool to know what is happening with packets as they go through the data plane processing.

The `ovs-vsitchd(8)` manpage describes basic usage of the `ofproto/trace` command used for tracing in Open vSwitch. For a tool with a goal similar to `ofproto/trace` for tracing packets through OVN logical switches, see [ovn-trace\(8\)](#).

Packet Tracing

In order to understand the tool, let's use the following flows as an example:

```
table=3,ip,tcp,tcp_dst=80,action=output:2
table=2,ip,tcp,tcp_dst=22,action=output:1
table=0,in_port=3,ip,nw_src=192.0.2.0/24,action=resubmit(,2)
table=0,in_port=3,ip,nw_src=198.51.100.0/24,action=resubmit(,3)
```

Note: If you can't use a "real" OVS setup you can use `ovs-sandbox`, as described in [Open vSwitch Advanced Features](#), which also provides additional tracing examples.

The first line adds a rule in table 3 matching on TCP/IP packet with destination port 80 (HTTP). If a packet matches, the action is to output the packet on OpenFlow port 2.

The second line is similar but matches on destination port 22. If a packet matches, the action is to output the packet on OpenFlow port 1.

The next two lines matches on source IP addresses. If there is a match, the packet is submitted to table indicated as parameter to the `resubmit()` action.

Now let's see if a packet from IP address 192.0.2.1 and destination port 22 would really go to OpenFlow port 1:

```
$ ovs-appctl ofproto/trace br0 in_port=3,tcp,nw_src=192.0.2.2,tcp_dst=22
Flow: tcp,in_port=3,vlan_tci=0x0000,dl_src=00:00:00:00:00:00,dl_dst=00:00:00:00:00:00,
↪nw_src=192.0.2.2,nw_dst=0.0.0.0,nw_tos=0,nw_ecn=0,nw_ttl=0,tp_src=0,tp_dst=22,tcp_
↪flags=0

bridge("br0")
-----
0. ip,in_port=3,nw_src=192.0.2.0/24, priority 32768
   resubmit(,2)
2. tcp,tp_dst=22, priority 32768
   output:1

Final flow: unchanged
```

(continues on next page)

(continued from previous page)

```
Megaflow: recirc_id=0,tcp,in_port=3,nw_src=192.0.2.0/24,nw_frag=no,tp_dst=22
Datapath actions: 1
```

The first line is the trace command. The br0 is the bridge where the packet is going through. The next arguments describe the packet itself. For instance, the nw_src matches with the IP source address. All the packet fields are well documented in the [ovs-fields\(7\)](#) man-page.

The second line shows the flow extracted from the packet described in the command line. Unspecified packet fields are zeroed.

The second group of lines shows the packet's trip through bridge br0. We see, in table 0, the OpenFlow flow that the fields matched, along with its priority, followed by its actions, one per line. In this case, we see that this packet matches the flow that resubmit those packets to table 2. The “resubmit” causes a second lookup in OpenFlow table 2, described by the block of text that starts with “2.” In the second lookup we see that this packet matches the rule that outputs those packets to OpenFlow port #1.

In summary, it is possible to follow the flow entries and actions until the final decision is made. At the end, the trace tool shows the Megaflow which matches on all relevant fields followed by the data path actions.

Let's see what happens with the same packet but with another TCP destination port:

```
$ ovs-appctl ofproto/trace br0 in_port=3,tcp,nw_src=192.0.2.2,tcp_dst=80
Flow: tcp,in_port=3,vlan_tci=0x0000,dl_src=00:00:00:00:00:00,dl_dst=00:00:00:00:00:00,
↪nw_src=192.0.2.2,nw_dst=0.0.0.0,nw_tos=0,nw_ecn=0,nw_ttl=0,tp_src=0,tp_dst=80,tcp_
↪flags=0

bridge("br0")
-----
 0. ip,in_port=3,nw_src=192.0.2.0/24, priority 32768
    resubmit(,2)
 2. No match.
    drop

Final flow: unchanged
Megaflow: recirc_id=0,tcp,in_port=3,nw_src=192.0.2.0/24,nw_frag=no,tp_dst=0x40/0xffc0
Datapath actions: drop
```

In the second group of lines, in table 0, you can see that the packet matches with the rule because of the source IP address, so it is resubmitted to the table 2 as before. However, it doesn't match any rule there. When the packet doesn't match any rule in the flow tables, it is called a table miss. The virtual switch table miss behavior can be configured and it depends on the OpenFlow version being used. In this example the default action was to drop the packet.

Credits

This document is heavily based on content from Flavio Bruno Leitner at Red Hat:

- <https://developers.redhat.com/blog/2016/10/12/tracing-packets-inside-open-vswitch/>

4.1.14 C IDL Compound Indexes

Introduction

This document describes the design and usage of the C IDL Compound Indexes feature, which allows OVSDb client applications to efficiently search table contents using arbitrary sets of column values in a generic way.

This feature is implemented entirely in the client IDL, requiring no changes to the OVSDb Server, OVSDb Protocol (OVSDb RFC (RFC 7047)) or additional interaction with the OVSDb server.

Please note that in this document, the term “index” refers to the common database term defined as “a data structure that facilitates data retrieval”. Unless stated otherwise, the definition for index from the OVSDb RFC (RFC 7047) is not used.

Typical Use Cases

Fast lookups

Depending on the topology, the route table of a network device could manage thousands of routes. Commands such as “show ip route <specific route>” would need to do a sequential lookup of the routing table to find the specific route. With an index created, the lookup time could be faster.

This same scenario could be applied to other features such as Access List rules and even interfaces lists.

Lexicographic order

There are a number of cases in which retrieving data in a particular lexicographic order is needed. For example, SNMP. When an administrator or even a NMS would like to retrieve data from a specific device, it’s possible that they will request data from full tables instead of just specific values. Also, they would like to have this information displayed in lexicographic order. This operation could be done by the SNMP daemon or by the CLI, but it would be better if the database could provide the data ready for consumption. Also, duplicate efforts by different processes will be avoided. Another use case for requesting data in lexicographic order is for user interfaces (web or CLI) where it would be better and quicker if the DB sends the data sorted instead of letting each process to sort the data by itself.

Implementation Design

This feature maintains a collection of indexes per table. The application can create any number of indexes per table.

An index can be defined over any number of columns, and supports the following options:

- Add a column with type string, boolean, uuid, integer or real (using default comparators).
- Select ordering direction of a column (ascending or descending, must be selected when creating the index).
- Use a custom ordering comparator (eg: treat a string column like a IP, or sort by the value of the “config” key in a map column).

Indexes can be searched for matches based on the key. They can also be iterated across a range of keys or in full.

For lookups, the user needs to provide a key to be used for locating the specific rows that meet his criteria. This key could be an IP address, a MAC address, an ACL rule, etc. If several rows match the query then the user can easily iterate over all of the matches.

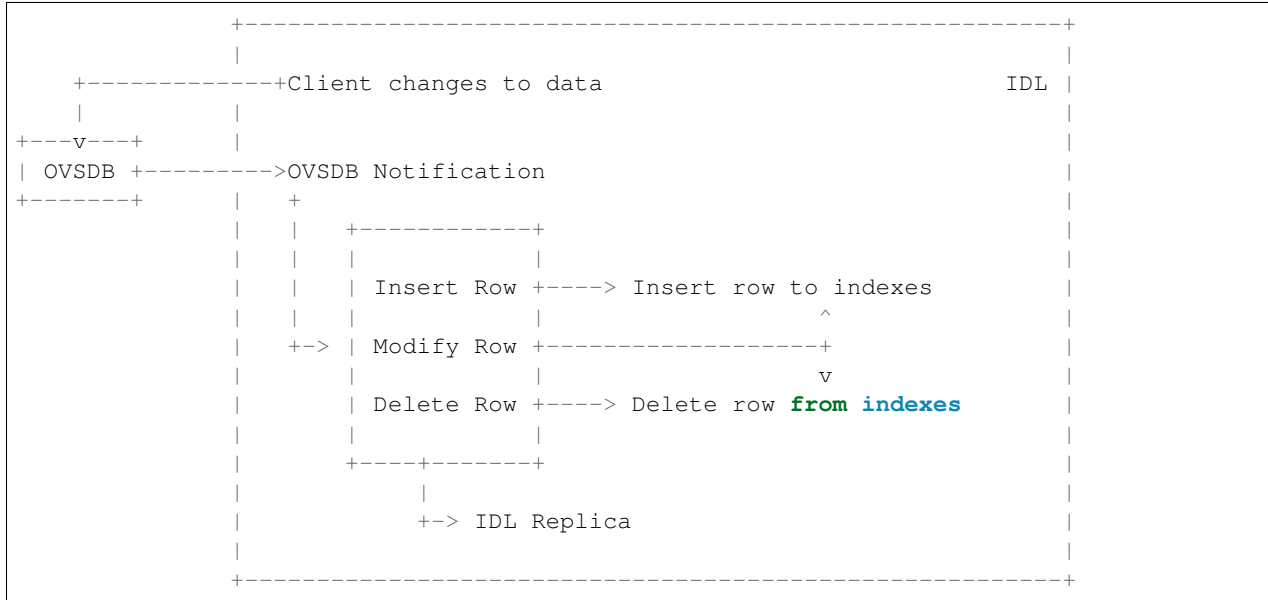
For accessing data in lexicographic order, the user can use the ranged iterators, which use “from” and “to” values to define a range.

The indexes maintain a pointer to the row in the local replica, avoiding the need to make additional copies of the data and thereby minimizing any additional memory and CPU overhead for their maintenance. It is intended that creating and maintaining indexes should be very cheap.

Another potential issue is the time needed to create the data structure and the time needed to add/remove elements. The indexes are always synchronized with the replica. For this reason is VERY IMPORTANT that the comparison functions (built-in and user provided) are FAST.

Skiplists are used as the primary data structure for the implementation of indexes. Indexes therefore have an expected $O(\log(n))$ cost when inserting, deleting or modifying a row, $O(\log(n))$ when retrieving a row by key, and $O(1)$ when retrieving the first or next row.

Indexes are maintained incrementally in the replica as notifications of database changes are received from the OVSDb server, as shown in the following diagram.



C IDL API

Index Creation

Each index must be created with the function `ovsdb_idl_index_create()` or one of the simpler convenience functions `ovsdb_idl_index_create1()` or `ovsdb_idl_index_create2()`. All indexes must be created before the first call to `ovsdb_idl_run()`.

Index Creation Example

```

/* Define a custom comparator for the column "stringField" in table
 * "Test". (Note that custom comparison functions are not often
 * necessary.)
 */
int stringField_comparator(const void *a, const void *b)
{
    struct ovsrec_test *AAA, *BBB;
    AAA = (struct ovsrec_test *)a;
    BBB = (struct ovsrec_test *)b;
    return strcmp(AAA->stringField, BBB->stringField);
}

void init_idl(struct ovsdb_idl **, char *remote)
{
    /* Add the columns to the IDL */
    *idl = ovsdb_idl_create(remote, &ovsrec_idl_class, false, true);
  
```

(continues on next page)

(continued from previous page)

```

ovsdb_idl_add_table(*idl, &ovsrec_table_test);
ovsdb_idl_add_column(*idl, &ovsrec_test_col_stringField);
ovsdb_idl_add_column(*idl, &ovsrec_test_col_numericField);
ovsdb_idl_add_column(*idl, &ovsrec_test_col_enumField);
ovsdb_idl_add_column(*idl, &ovsrec_test_col_boolField);

struct ovsdb_idl_index_column columns[] = {
    { .column = &ovsrec_test_col_stringField,
      .comparer = stringField_comparator },
    { .column = &ovsrec_test_col_numericField,
      .order = OVSDb_INDEX_DESC },
};
struct ovsdb_idl_index *index = ovsdb_idl_create_index(
    *idl, columns, ARRAY_SIZE(columns));
...
}

```

Index Usage

Iterators

The recommended way to do queries is using a “ranged foreach”, an “equal foreach” or a “full foreach” over an index. The mechanism works as follows:

1. Create index row objects with index columns set to desired search key values (one is needed for equality iterators, two for range iterators, a search key is not needed for the full index iterator).
2. Pass the index, an iteration variable, and the index row object to the iterator.
3. Use the values within iterator loop.

The library implements three different iterators: a range iterator, an equality iterator and a full index iterator. The range iterator receives two values and iterates over all rows with values that are within that range (inclusive of the two values defining the range). The equality iterator iterates over all rows that exactly match the value passed. The full index iterator iterates over all rows in the index, in an order determined by the comparison function and configured direction (ascending or descending).

Note that indexes are *sorted by the “concatenation” of the values in all indexed columns*, so the ranged iterator returns all the values between “from.col1 from.col2 ... from.coln” and “to.col1 to.col2 ... to.coln”, *NOT the rows with a value in column 1 between from.col1 and to.col1, and so on.*

The iterators are macros specific to each table. An example of the use of these iterators follows:

```

/*
 * Equality iterator; iterates over all the records equal to "value".
 */
struct ovsrec_test *target = ovsrec_test_index_init_row(index);
ovsrec_test_index_set_stringField(target, "hello world");
struct ovsrec_test *record;
OVSREC_TEST_FOR_EACH_EQUAL (record, target, index) {
    /* Can return zero, one or more records */
    assert(strcmp(record->stringField, "hello world") == 0);
    printf("Found one record with %s", record->stringField);
}
ovsrec_test_index_destroy_row(target);

```

(continues on next page)

(continued from previous page)

```

/*
 * Range iterator; iterates over all records between two values
 * (inclusive).
 */
struct ovsrec_test *from = ovsrec_test_index_init_row(index);
struct ovsrec_test *to = ovsrec_test_index_init_row(index);

ovsrec_test_index_set_stringField(from, "aaa");
ovsrec_test_index_set_stringField(to, "mmm");
OVSREC_TEST_FOR_EACH_RANGE (record, from, to, index) {
    /* Can return zero, one or more records */
    assert(strcmp("aaa", record->stringField) <= 0);
    assert(strcmp(record->stringField, "mmm") <= 0);
    printf("Found one record with %s", record->stringField);
}

ovsrec_test_index_destroy_row(from);
ovsrec_test_index_destroy_row(to);

/*
 * Index iterator; iterates over all nodes in the index, in order
 * determined by comparison function and configured order (ascending
 * or descending).
 */
OVSREC_TEST_FOR_EACH_BYINDEX (record, index) {
    /* Can return zero, one or more records */
    printf("Found one record with %s", record->stringField);
}

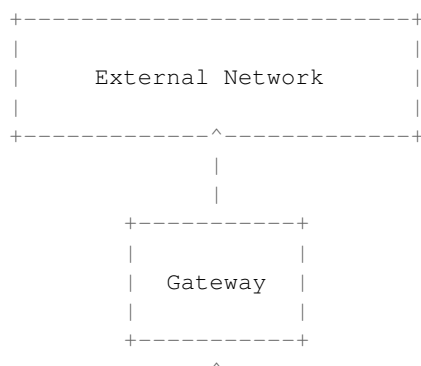
```

General Index Access

While the currently defined iterators are suitable for many use cases, it is also possible to create custom iterators using the more general API on which the existing iterators have been built. See `ovsdb-idl.h` for the details.

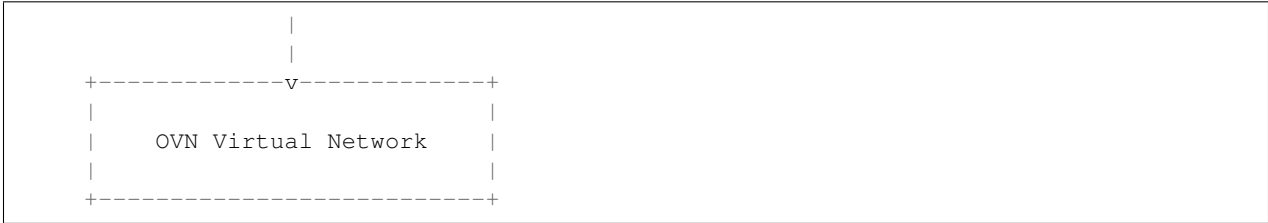
4.2 OVN

4.2.1 OVN Gateway High Availability Plan



(continues on next page)

(continued from previous page)



The OVN gateway is responsible for shuffling traffic between the tunneled overlay network (governed by `ovn-northd`), and the legacy physical network. In a naive implementation, the gateway is a single x86 server, or hardware VTEP. For most deployments, a single system has enough forwarding capacity to service the entire virtualized network, however, it introduces a single point of failure. If this system dies, the entire OVN deployment becomes unavailable. To mitigate this risk, an HA solution is critical – by spreading responsibility across multiple systems, no single server failure can take down the network.

An HA solution is both critical to the manageability of the system, and extremely difficult to get right. The purpose of this document, is to propose a plan for OVN Gateway High Availability which takes into account our past experience building similar systems. It should be considered a fluid changing proposal, not a set-in-stone decree.

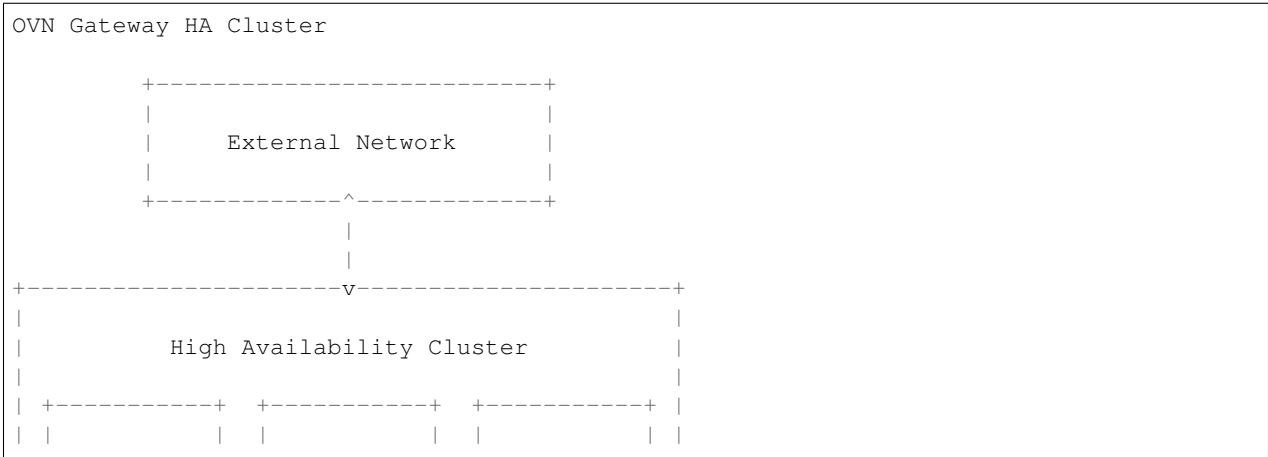
Note: This document describes a range of options OVN could take to provide high availability for gateways. The current implementation provides L3 gateway high availability by the “Router Specific Active/Backup” approach described in this document.

Basic Architecture

In an OVN deployment, the set of hypervisors and network elements operating under the guidance of `ovn-northd` are in what’s called “logical space”. These servers use VXLAN, STT, or Geneve to communicate, oblivious to the details of the underlying physical network. When these systems need to communicate with legacy networks, traffic must be routed through a Gateway which translates from OVN controlled tunnel traffic, to raw physical network traffic.

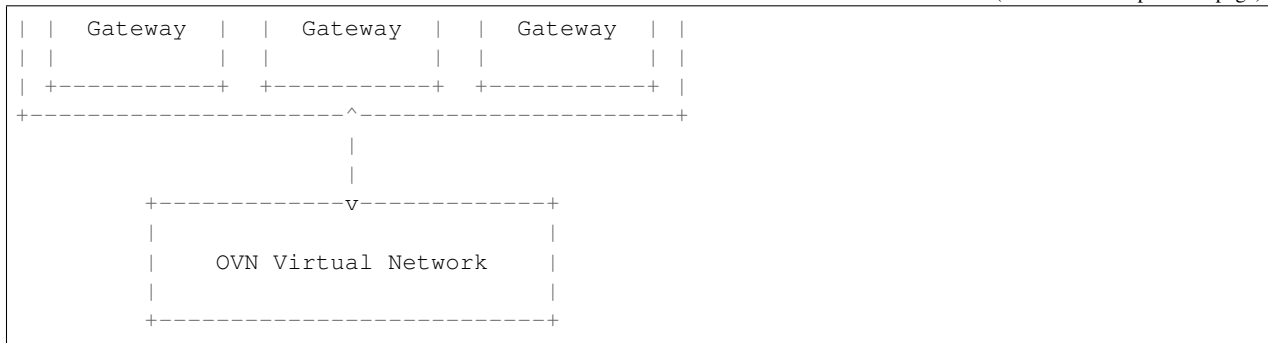
Since the gateway is typically the only system with a connection to the physical network all traffic between logical space and the WAN must travel through it. This makes it a critical single point of failure – if the gateway dies, communication with the WAN ceases for all systems in logical space.

To mitigate this risk, multiple gateways should be run in a “High Availability Cluster” or “HA Cluster”. The HA cluster will be responsible for performing the duties of a gateways, while being able to recover gracefully from individual member failures.



(continues on next page)

(continued from previous page)



L2 vs L3 High Availability

In order to achieve this goal, there are two broad approaches one can take. The HA cluster can appear to the network like a giant Layer 2 Ethernet Switch, or like a giant IP Router. These approaches are called L2HA, and L3HA respectively. L2HA allows ethernet broadcast domains to extend into logical space, a significant advantage, but this comes at a cost. The need to avoid transient L2 loops during failover significantly complicates their design. On the other hand, L3HA works for most use cases, is simpler, and fails more gracefully. For these reasons, it is suggested that OVN supports an L3HA model, leaving L2HA for future work (or third party VTEP providers). Both models are discussed further below.

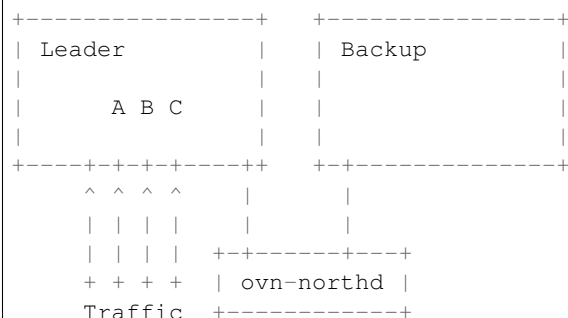
L3HA

In this section, we'll work through a basic simple L3HA implementation, on top of which we'll gradually build more sophisticated features explaining their motivations and implementations as we go.

Naive active-backup

Let's assume that there are a collection of logical routers which a tenant has asked for, our task is to schedule these logical routers on one of N gateways, and gracefully redistribute the routers on gateways which have failed. The absolute simplest way to achieve this is what we'll call "naive-active-backup".

Naive Active Backup HA Implementation



In a naive active-backup, one of the Gateways is chosen (arbitrarily) as a leader. All logical routers (A, B, C in the figure), are scheduled on this leader gateway and all traffic flows through it. ovn-northd monitors this gateway via OpenFlow echo requests (or some equivalent), and if the gateway dies, it recreates the routers on one of the backups.

This approach basically works in most cases and should likely be the starting point for OVN – it’s strictly better than no HA solution and is a good foundation for more sophisticated solutions. That said, it’s not without it’s limitations. Specifically, this approach doesn’t coordinate with the physical network to minimize disruption during failures, and it tightly couples failover to ovn-northd (we’ll discuss why this is bad in a bit), and wastes resources by leaving backup gateways completely unutilized.

Router Failover

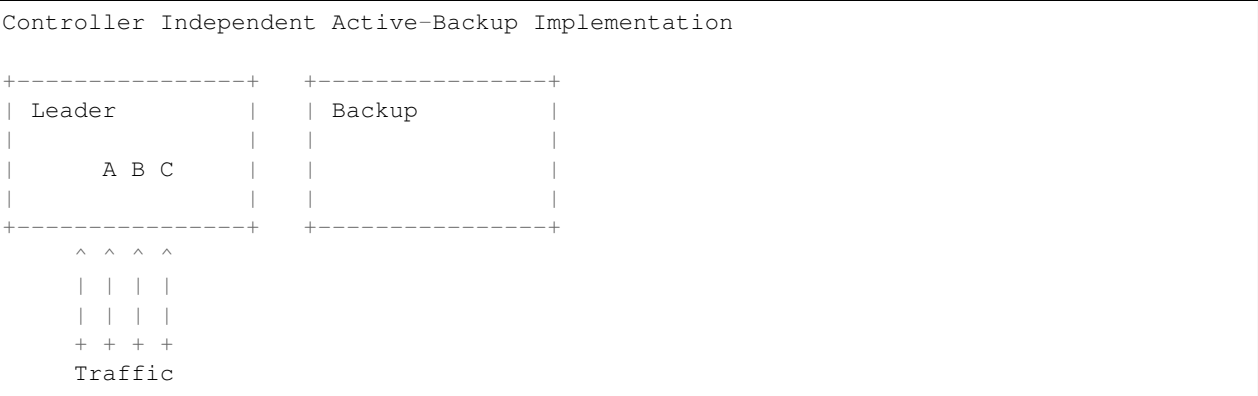
When ovn-northd notices the leader has died and decides to migrate routers to a backup gateway, the physical network has to be notified to direct traffic to the new gateway. Otherwise, traffic could be blackholed for longer than necessary making failovers worse than they need to be.

For now, let’s assume that OVN requires all gateways to be on the same IP subnet on the physical network. If this isn’t the case, gateways would need to participate in routing protocols to orchestrate failovers, something which is difficult and out of scope of this document.

Since all gateways are on the same IP subnet, we simply need to worry about updating the MAC learning tables of the Ethernet switches on that subnet. Presumably, they all have entries for each logical router pointing to the old leader. If these entries aren’t updated, all traffic will be sent to the (now defunct) old leader, instead of the new one.

In order to mitigate this issue, it’s recommended that the new gateway sends a Reverse ARP (RARP) onto the physical network for each logical router it now controls. A Reverse ARP is a benign protocol used by many hypervisors when virtual machines migrate to update L2 forwarding tables. In this case, the ethernet source address of the RARP is that of the logical router it corresponds to, and its destination is the broadcast address. This causes the RARP to travel to every L2 switch in the broadcast domain, updating forwarding tables accordingly. This strategy is recommended in all failover mechanisms discussed in this document – when a router newly boots on a new leader, it should RARP its MAC address.

Controller Independent Active-backup



The fundamental problem with naive active-backup, is it tightly couples the failover solution to ovn-northd. This can significantly increase downtime in the event of a failover as the (often already busy) ovn-northd controller has to recompute state for the new leader. Worse, if ovn-northd goes down, we can’t perform gateway failover at all. This violates the principle that control plane outages should have no impact on dataplane functionality.

In a controller independent active-backup configuration, ovn-northd is responsible for initial configuration while the HA cluster is responsible for monitoring the leader, and failing over to a backup if necessary. ovn-northd sets HA policy, but doesn’t actively participate when failovers occur.

Of course, in this model, ovn-northd is not without some responsibility. Its role is to pre-plan what should happen in the event of a failure, leaving it to the individual switches to execute this plan. It does this by assigning each gateway

a unique leadership priority. Once assigned, it communicates this priority to each node it controls. Nodes use the leadership priority to determine which gateway in the cluster is the active leader by using a simple metric: the leader is the gateway that is healthy, with the highest priority. If that gateway goes down, leadership falls to the next highest priority, and conversely, if a new gateway comes up with a higher priority, it takes over leadership.

Thus, in this model, leadership of the HA cluster is determined simply by the status of its members. Therefore if we can communicate the status of each gateway to each transport node, they can individually figure out which is the leader, and direct traffic accordingly.

Tunnel Monitoring

Since in this model leadership is determined exclusively by the health status of member gateways, a key problem is how do we communicate this information to the relevant transport nodes. Luckily, we can do this fairly cheaply using tunnel monitoring protocols like BFD.

The basic idea is pretty straightforward. Each transport node maintains a tunnel to every gateway in the HA cluster (not just the leader). These tunnels are monitored using the BFD protocol to see which are alive. Given this information, hypervisors can trivially compute the highest priority live gateway, and thus the leader.

In practice, this leadership computation can be performed trivially using the bundle or group action. Rather than using OpenFlow to simply output to the leader, all gateways could be listed in an active-backup bundle action ordered by their priority. The bundle action will automatically take into account the tunnel monitoring status to output the packet to the highest priority live gateway.

Inter-Gateway Monitoring

One somewhat subtle aspect of this model, is that failovers are not globally atomic. When a failover occurs, it will take some time for all hypervisors to notice and adjust accordingly. Similarly, if a new high priority Gateway comes up, it may take some time for all hypervisors to switch over to the new leader. In order to avoid confusing the physical network, under these circumstances it's important for the backup gateways to drop traffic they've received erroneously. In order to do this, each Gateway must know whether or not it is, in fact active. This can be achieved by creating a mesh of tunnels between gateways. Each gateway monitors the other gateways its cluster to determine which are alive, and therefore whether or not that gateway happens to be the leader. If leading, the gateway forwards traffic normally, otherwise it drops all traffic.

We should note that this method works well under the assumption that there are no inter-gateway connectivity failures, in such case this method would fail to elect a single master. The simplest example is two gateways which stop seeing each other but can still reach the hypervisors. Protocols like VRRP or CARP have the same issue. A mitigation for this type of failure mode could be achieved by having all network elements (hypervisors and gateways) periodically share their link status to other endpoints.

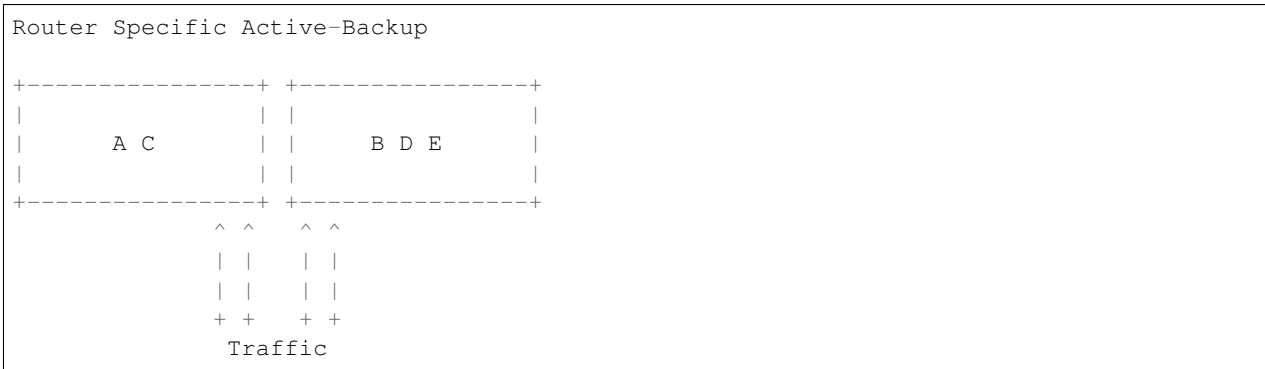
Gateway Leadership Resignation

Sometimes a gateway may be healthy, but still may not be suitable to lead the HA cluster. This could happen for several reasons including:

- The physical network is unreachable
- BFD (or ping) has detected the next hop router is unreachable
- The Gateway recently booted and isn't fully configured

In this case, the Gateway should resign leadership by holding its tunnels down using the `other_config:cpath_down` flag. This indicates to participating hypervisors and Gateways that this gateway should be treated as if it's down, even though its tunnels are still healthy.

Router Specific Active-Backup



Controller independent active-backup is a great advance over naive active-backup, but it still has one glaring problem – it under-utilizes the backup gateways. In ideal scenario, all traffic would split evenly among the live set of gateways. Getting all the way there is somewhat tricky, but as a step in the direction, one could use the “Router Specific Active-Backup” algorithm. This algorithm looks a lot like active-backup on a per logical router basis, with one twist. It chooses a different active Gateway for each logical router. Thus, in situations where there are several logical routers, all with somewhat balanced load, this algorithm performs better.

Implementation of this strategy is quite straightforward if built on top of basic controller independent active-backup. On a per logical router basis, the algorithm is the same, leadership is determined by the liveness of the gateways. The key difference here is that the gateways must have a different leadership priority for each logical router. These leadership priorities can be computed by `ovn-northd` just as they had been in the controller independent active-backup model.

Once we have these per logical router priorities, they simply need be communicated to the members of the gateway cluster and the hypervisors. The hypervisors in particular, need simply have an active-backup bundle action (or group action) per logical router listing the gateways in priority order for *that router*, rather than having a single bundle action shared for all the routers.

Additionally, the gateways need to be updated to take into account individual router priorities. Specifically, each gateway should drop traffic of backup routers it’s running, and forward traffic of active gateways, instead of simply dropping or forwarding everything. This should likely be done by having `ovn-controller` recompute OpenFlow for the gateway, though other options exist.

The final complication is that `ovn-northd`’s logic must be updated to choose these per logical router leadership priorities in a more sophisticated manner. It doesn’t matter much exactly what algorithm it chooses to do this, beyond that it should provide good balancing in the common case. I.E. each logical routers priorities should be different enough that routers balance to different gateways even when failures occur.

Preemption

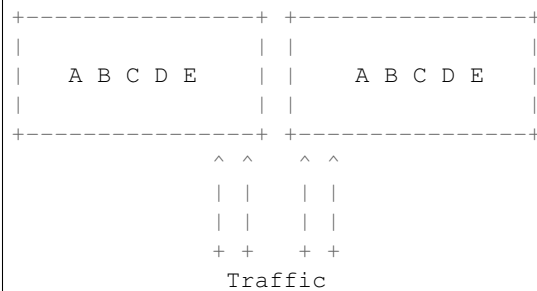
In an active-backup setup, one issue that users will run into is that of gateway leader preemption. If a new Gateway is added to a cluster, or for some reason an existing gateway is rebooted, we could end up in a situation where the newly activated gateway has higher priority than any other in the HA cluster. In this case, as soon as that gateway appears, it will preempt leadership from the currently active leader causing an unnecessary failover. Since failover can be quite expensive, this preemption may be undesirable.

The controller can optionally avoid preemption by cleverly tweaking the leadership priorities. For each router, new gateways should be assigned priorities that put them second in line or later when they eventually come up. Furthermore, if a gateway goes down for a significant period of time, its old leadership priorities should be revoked and new ones should be assigned as if it’s a brand new gateway. Note that this should only happen if a gateway has been down for a while (several minutes), otherwise a flapping gateway could have wide ranging, unpredictable, consequences.

Note that preemption avoidance should be optional depending on the deployment. One necessarily sacrifices optimal load balancing to satisfy these requirements as new gateways will get no traffic on boot. Thus, this feature represents a trade-off which must be made on a per installation basis.

Fully Active-Active HA

Fully Active-Active HA



The final step in L3HA is to have true active-active HA. In this scenario each router has an instance on each Gateway, and a mechanism similar to ECMP is used to distribute traffic evenly among all instances. This mechanism would require Gateways to participate in routing protocols with the physical network to attract traffic and alert of failures. It is out of scope of this document, but may eventually be necessary.

L2HA

L2HA is very difficult to get right. Unlike L3HA, where the consequences of problems are minor, in L2HA if two gateways are both transiently active, an L2 loop triggers and a broadcast storm results. In practice to get around this, gateways end up implementing an overly conservative “when in doubt drop all traffic” policy, or they implement something like MLAG.

MLAG has multiple gateways work together to pretend to be a single L2 switch with a large LACP bond. In principle, it’s the right solution to the problem as it solves the broadcast storm problem, and has been deployed successfully in other contexts. That said, it’s difficult to get right and not recommended.

4.2.2 Role Based Access Control

Where SSL provides authentication when connecting to an OVS database, role based access control (RBAC) provides authorization to operations performed by clients connecting to an OVS database. RBAC allows for administrators to restrict the database operations a client may perform and thus enhance the security already provided by SSL.

In theory, any OVS database could define RBAC roles and permissions, but at present only the OVN southbound database has the appropriate tables defined to facilitate RBAC.

Mechanics

RBAC is intended to supplement SSL. In order to enable RBAC, the connection to the database must use SSL. Some permissions in RBAC are granted based on the certificate common name (CN) of the connecting client.

RBAC is controlled with two database tables, RBAC_Role and RBAC_Permission. The RBAC_Permission table contains records that describe a set of permissions for a given table in the database.

The RBAC_Permission table contains the following columns:

table The table in the database for which permissions are being described.

insert_delete Describes whether insertion and deletion of records is allowed.

update A list of columns that are allowed to be updated.

authorization A list of column names. One of the listed columns must match the SSL certificate CN in order for the attempted operation on the table to succeed. If a key-value pair is provided, then the key is the column name, and the value is the name of a key in that column. An empty string gives permission to all clients to perform operations.

The RBAC_Role table contains the following columns:

name The name of the role being defined

permissions A list of key-value pairs. The key is the name of a table in the database, and the value is a UUID of a record in the RBAC_Permission table that describes the permissions the role has for that table.

Note: All tables not explicitly referenced in an RBAC_Role record are read-only

In order to enable RBAC, specify the role name as an argument to the set-connection command for the database. As an example, to enable the “ovn-controller” role on the OVN southbound database, use the following command:

```
$ ovn-sbctl set-connection role=ovn-controller ssl:192.168.0.1:6642
```

Pre-defined Roles

This section describes roles that have been defined internally by OVS/OVN.

ovn-controller

The ovn-controller role is specified in the OVN southbound database and is intended for use by hypervisors running the ovn-controller daemon. ovn-controller connects to the OVN southbound database mostly to read information, but there are a few cases where ovn-controller also needs to write. The ovn-controller role was designed to allow for ovn-controllers to write to the southbound database only in places where it makes sense to do so. This way, if an intruder were to take over a hypervisor running ovn-controller, it is more difficult to compromise the entire overlay network.

It is strongly recommended to set the ovn-controller role for the OVN southbound database to enhance security.

4.2.3 What’s New with OVS and OVN 2.8

This document is about what was added in Open vSwitch 2.8, which was released at the end of August 2017, concentrating on the new features in OVN. It also covers some of what is coming up in Open vSwitch and OVN 2.9, which is due to be released in February 2018. OVN has many features, and this document does not cover every new or enhanced feature (but contributions are welcome).

This document assumes a basic familiarity with Open vSwitch, OVN, and their associated tools. For more information, please refer to the Open vSwitch and OVN documentation, such as the `ovn-architecture(7)` manpage.

Debugging and Troubleshooting

Before version 2.8, Open vSwitch command-line tools were far more painful to use than they needed to be. This section covers the improvements made to the CLI in the 2.8 release.

User-Hostile UUIDs

The OVN CLI, through `ovn-nbctl`, `ovn-nbctl`, and `ovn-trace`, used full-length UUIDs almost everywhere. It didn't even provide any assistance with completion, etc., which in practice meant always cutting and pasting UUIDs from one command or window to another. This problem wasn't limited to the places where one would expect to have to see or use a UUID, either. In many places where one would expect to be able to use a network, router, or port name, a UUID was required instead. In many places where one would want to see a name, the UUID was displayed instead. More than anything else, these shortcomings made the CLI user-hostile.

There was an underlying problem that the southbound database didn't actually contain all the information needed to provide a decent user interface. In some cases, for example, the human-friendly names that one would want to use for entities simply weren't part of the database. These names weren't necessary for correctness, only for usability.

OVN 2.8 eased many of these problems. Most parts of the CLI now allow the user to abbreviate UUIDs, as long as the abbreviations are unique within the database. Some parts of the CLI where full-length UUIDs make output hard to read now abbreviate them themselves. Perhaps more importantly, in many places the OVN CLI now displays and accepts human-friendly names for networks, routers, ports, and other entities. In the places where the names were not previously available, OVN (through `ovn-northd`) now copies the names into the southbound database.

The CLIs for layers below OVN, at the OpenFlow and datapath layers with `ovs-ofctl` and `ovs-dpctl`, respectively, had some similar problems in which numbers were used for entities that had human-friendly names. Open vSwitch 2.8 also solves some of those problems. Other than that, the most notable enhancement in this area was the `--no-stats` option to `ovs-ofctl dump-flows`, which made that command's output more readable for the cases where per-flow statistics were not interesting to the reader.

Connections Between Levels

OVN and Open vSwitch work almost like a stack of compilers: the OVN Neutron plugin translates Neutron configuration into OVN northbound configuration, which `ovn-northd` translates into logical flows, which `ovn-controller` translates into OpenFlow flows, which `ovs-vswitchd` translates into datapath flows. For debugging and troubleshooting it is often necessary to understand exactly how these translations work. The relationship from a logical flow to its OpenFlow flows, or in the other direction, from an OpenFlow flow back to the logical flow that produced it, was often of particular interest, but OVN didn't provide good tools for the job.

OVN 2.8 added some new features that ease these jobs. `ovn-sbctl lflow-list` has a new option `--ovs` that lists the OpenFlow flows on a particular chassis that were generated from the logical flows that it lists. `ovn-trace` also added a similar `--ovs` option that applies to the logical flows it traces.

In the other direction, OVN 2.8 added a new utility `ovn-detrace` that, given an Open vSwitch trace of OpenFlow flows, annotates it with the logical flows that yielded those OpenFlow flows.

Distributed Firewall

OVN supports a distributed firewall with stateful connection tracking to ensure that only packets for established connections, or those that the configuration explicitly allows, can ingress a given VM or container. Neutron uses this feature by default. Most packets in an OpenStack environment pass through it twice, once after egress from the packet's source VM and once before ingress into its destination VM. Before OVN 2.8, the `ovn-trace` program, which shows the path of a packet through an OVN logical network, did not support the logical firewall, which in practice made it almost useless for Neutron.

In OVN 2.8, `ovn-trace` adds support for the logical firewall. By default it assumes that packets are part of an established connection, which is usually what the user wants as part of the trace. It also accepts command-line options to override that assumption, which allows the user to discover the treatment of packets that the firewall should drop.

At the next level deeper, prior to Open vSwitch 2.8, the OpenFlow tracing command `ofproto/trace` also supported neither the connection tracking feature underlying the OVN distributed firewall nor the “recirculation” feature that accompanied it. This meant that, even if the user tried to look deeper into the distributed firewall mechanism, he or she would encounter a further roadblock. Open vSwitch 2.8 added support for both of these features as well.

Summary Display

`ovn-nbctl show` and `ovn-sbctl show`, for showing an overview of the OVN configuration, didn’t show a lot of important information. OVN adds some more useful information here.

DNS, and IPAM

OVN 2.8 adds a built-in DNS server designed for assigning names to VMs and containers within an OVN logical network. DNS names are assigned using records in the OVN northbound database and, like other OVN features, translated into logical flows at the OVN southbound layer. DNS requests directed to the OVN DNS server never leave the hypervisor from which the request is sent; instead, OVN processes and replies to the request from its `ovn-controller` local agent. The OVN DNS server is not a general-purpose DNS server and cannot be used for that purpose.

OVN includes simple built-in support for IP address management (IPAM), in which OVN assigns IP addresses to VMs or containers from a pool or pools of IP addresses delegated to it by the administrator. Before OVN 2.8, OVN IPAM only supported IPv4 addresses; OVN 2.8 adds support for IPv6. OVN 2.8 also enhances the address pool support to allow specific addresses to be excluded. Neutron assigns IP addresses itself and does not use OVN IPAM.

High Availability

As a distributed system, in OVN a lot can go wrong. As OVN advances, it adds redundancy in places where currently a single failure could disrupt the functioning of the system as a whole. OVN 2.8 adds two new kinds of high availability.

ovn-northd HA

The `ovn-northd` program sits between the OVN northbound and southbound databases and translates from a logical network configuration into logical flows. If `ovn-northd` itself or the host on which it runs fails, then updates to the OVN northbound configuration will not propagate to the hypervisors and the OVN configuration freezes in place until `ovn-northd` restarts.

OVN 2.8 adds support for active-backup HA to `ovn-northd`. When more than one `ovn-northd` instance runs, it uses an OVSDB locking feature to automatically choose a single active instance. When that instance dies or becomes nonresponsive, the OVSDB server automatically choose one of the remaining instance(s) to take over.

L3 Gateway HA

In OVN 2.8, multiple chassis may now be specified for L3 gateways. When more than one chassis is specified, OVN manages high availability for that gateway. Each hypervisor uses the BFD protocol to keep track of the gateway nodes that are currently up. At any given time, a hypervisor uses the highest-priority gateway node that is currently up.

OVSDB

The OVN architecture relies heavily on OVSDB, the Open vSwitch database, for hosting the northbound and southbound databases. OVSDB was originally selected for this purpose because it was already used in Open vSwitch

for configuring OVS itself and, thus, it was well integrated with OVS and well supported in C and Python, the two languages that are used in Open vSwitch.

OVSDb was well designed for its original purpose of configuring Open vSwitch. It supports ACID transactions, has a small, efficient server, a flexible schema system, and good support for troubleshooting and debugging. However, it lacked several features that are important for OVN but not for Open vSwitch. As OVN advances, these missing features have become more and more of a problem. One option would be to switch to a different database that already has many of these features, but despite a careful search, no ideal existing database was identified, so the project chose instead to improve OVSDb where necessary to bring it up to speed. The following sections talk more about recent and future improvements.

High Availability

When `ovsdb-server` was only used for OVS configuration, high availability was not important. `ovsdb-server` was capable of restarting itself automatically if it crashed, and if the whole system went down then Open vSwitch itself was dead too, so the database server's failure was not important.

In contrast, the northbound and southbound databases are centralized components of a distributed system, so it is important that they not be a single point of failure for the system as a whole. In released versions of OVN, `ovsdb-server` supports only “active-backup replication” across a pair of servers. This means that if one server goes down, the other can pick it back up approximately where the other one left off. The servers do not have built-in support for deciding at any given time which is the active and which the backup, so the administrator must configure an external agent to do this management.

Active-backup replication is not entirely satisfactory, for multiple reasons. Replication is only approximate. Configuring the external agent requires extra work. There is no benefit from the backup server except when the active server fails. At most two servers can be used.

A new form of high availability for OVSDb is under development for the OVN 2.9 release, based on the Raft algorithm for distributed consensus. Whereas replication uses two servers, clustering using Raft requires three or more (typically an odd number) and continues functioning as long as more than half of the servers are up. The clustering implementation is built into `ovsdb-server` and does not require an external agent. Clustering preserves the ACID properties of the database, so that a transaction that commits is guaranteed to persist. Finally, reads (which are the bulk of the OVN workload) scale with the size of the cluster, so that adding more servers should improve performance as the number of hypervisors in an OVN deployment increases. As of this writing, OVSDb support for clustering is undergoing development and early deployment testing.

RBAC security

Until Open vSwitch 2.8, `ovsdb-server` had little support for access control within a database. If an OVSDb client could modify the database at all, it could make arbitrary changes. This was sufficient for most use case to that point.

Hypervisors in an OVN deployment need access to the OVN southbound database. Most of their access is reads, to find out about the OVN configuration. Hypervisors do need some write access to the southbound database, primarily to let the other hypervisors know what VMs and containers they are running and how to reach them. Thus, OVN gives all of the hypervisors in the OVN deployment write access to the OVN southbound database. This is fine when all is well, but if any of the hypervisors were compromised then they could disrupt the entire OVN deployment by corrupting the database.

The OVN developers considered a few ways to solve this problem. One way would be to introduce a new central service (perhaps in `ovn-northd`) that provided only the kinds of writes that the hypervisors legitimately need, and then grant hypervisors direct access to the southbound database only for reads. But ultimately the developers decided to introduce a new form of more access control for OVSDb, called the OVSDb RBAC (role-based access control) feature. OVSDb RBAC allows for granular enough control over access that hypervisors can be granted only the

ability to add, modify, and delete the records that relate to themselves, preventing them from corrupting the database as a whole.

Further Directions

For more information about new features in OVN and Open vSwitch, please refer to the NEWS file distributed with the source tree. If you have questions about Open vSwitch or OVN features, please feel free to write to the Open vSwitch discussion mailing list at ovs-discuss@openvswitch.org.

ovn-architecture(7)	(pdf)	(html)	(plain text)
-------------------------------------	-----------------------	------------------------	------------------------------

Answers to common “How do I?”-style questions. For more information on the topics covered herein, refer to *Deep Dive*.

5.1 OVS

5.1.1 Open vSwitch with KVM

This document describes how to use Open vSwitch with the Kernel-based Virtual Machine (KVM).

Note: This document assumes that you have Open vSwitch set up on a Linux system.

Setup

KVM uses `tunctl` to handle various bridging modes, which you can install with the Debian/Ubuntu package `uml-utilities`:

```
$ apt-get install uml-utilities
```

Next, you will need to modify or create custom versions of the `qemu-ifup` and `qemu-ifdown` scripts. In this guide, we’ll create custom versions that make use of example Open vSwitch bridges that we’ll describe in this guide.

Create the following two files and store them in known locations. For example:

```
$ cat << 'EOF' > /etc/ovs-ifup
#!/bin/sh

switch='br0'
ip link set $1 up
```

(continues on next page)

(continued from previous page)

```
ovs-vsctl add-port ${switch} $1
EOF
```

```
$ cat << 'EOF' > /etc/ovs-ifdown
#!/bin/sh

switch='br0'
ip addr flush dev $1
ip link set $1 down
ovs-vsctl del-port ${switch} $1
EOF
```

The basic usage of Open vSwitch is described at the end of *Open vSwitch on Linux, FreeBSD and NetBSD*. If you haven't already, create a bridge named `br0` with the following command:

```
$ ovs-vsctl add-br br0
```

Then, add a port to the bridge for the NIC that you want your guests to communicate over (e.g. `eth0`):

```
$ ovs-vsctl add-port br0 eth0
```

Refer to `ovs-vsctl(8)` for more details.

Next, we'll start a guest that will use our `ifup` and `ifdown` scripts:

```
$ kvm -m 512 -net nic,macaddr=00:11:22:EE:EE:EE -net \
    tap,script=/etc/ovs-ifup,downscript=/etc/ovs-ifdown -drive \
    file=/path/to/disk-image,boot=on
```

This will start the guest and associate a tap device with it. The `ovs-ifup` script will add a port on the `br0` bridge so that the guest will be able to communicate over that bridge.

To get some more information and for debugging you can use Open vSwitch utilities such as `ovs-dpctl` and `ovs-ofctl`. For example:

```
$ ovs-dpctl show
$ ovs-ofctl show br0
```

You should see tap devices for each KVM guest added as ports to the bridge (e.g. `tap0`)

Refer to `ovs-dpctl(8)` and `ovs-ofctl(8)` for more details.

Bug Reporting

Please report problems to bugs@openvswitch.org.

5.1.2 Encrypt Open vSwitch Tunnels with IPsec

This document gives detailed description on the OVS IPsec tunnel and its configuration modes. If you want to follow a step-by-step guide to run and test IPsec tunnel, please refer to *OVS IPsec Tutorial*.

Overview

Why do encryption?

OVS tunnel packets are transported from one machine to another. Along the path, the packets are processed by physical routers and physical switches. There are risks that these physical devices might read or write the contents of the tunnel packets. IPsec encrypts IP payload and prevents the malicious party sniffing or manipulating the tunnel traffic.

OVS IPsec

OVS IPsec aims to provide a simple interface for user to add encryption on OVS tunnels. It supports GRE, GENEVE, VXLAN, and STT tunnel. The IPsec configuration is done by setting options of the tunnel interface and other_config of Open_vSwitch. You can choose different authentication methods and plaintext tunnel policies based on your requirements.

OVS does not currently provide any support for IPsec encryption for traffic not encapsulated in a tunnel.

Configuration

Authentication Methods

Hosts of the IPsec tunnel need to authenticate each other to build a secure channel. There are three authentication methods:

1. You can use a pre-shared key (PSK) to do authentication. In both hosts, set the same PSK value. This PSK is like your password. You should never reveal it to untrusted parties. This method is easier to use but less secure than the certificate-based methods:

```
$ ovs-vsctl add-port br0 ipsec_gre0 -- \
    set interface ipsec_gre0 type=gre \
        options:remote_ip=2.2.2.2 \
        options:psk=swordfish
```

2. You can use a self-signed certificate to do authentication. In each host, generate a certificate and the paired private key. Copy the certificate of the remote host to the local host and configure the OVS as following:

```
$ ovs-vsctl set Open_vSwitch . \
    other_config:certificate=/path/to/local_cert.pem \
    other_config:private_key=/path/to/priv_key.pem
$ ovs-vsctl add-port br0 ipsec_gre0 -- \
    set interface ipsec_gre0 type=gre \
        options:remote_ip=2.2.2.2 \
        options:remote_cert=/path/to/remote_cert.pem
```

local_cert.pem is the certificate of the local host. *priv_key.pem* is the private key of the local host. *priv_key.pem* needs to be stored in a secure location. *remote_cert.pem* is the certificate of the remote host.

Note: OVS IPsec requires x.509 version 3 certificate with the subjectAltName DNS field setting the same string as the common name (CN) field. You can follow the tutorial in [OVS IPsec Tutorial](#) and use ovs-pki(8) to generate compatible certificate and key.

(Before OVS version 2.10.90, ovs-pki(8) did not generate x.509 v3 certificates, so if your existing PKI was generated by an older version, it is not suitable for this purpose.)

3. You can also use CA-signed certificate to do authentication. First, you need to create a CA certificate and sign each host certificate with the CA key (please see [OVS IPsec Tutorial](#)). Copy the CA certificate to each host and configure the OVS as following:

```
$ ovs-vsctl set Open_vSwitch . \
    other_config:certificate=/path/to/local_cert.pem \
    other_config:private_key=/path/to/priv_key.pem \
    other_config:ca_cert=/path/to/ca_cert.pem
$ ovs-vsctl add-port br0 ipsec_gre0 -- \
    set interface ipsec_gre0 type=gre \
        options:remote_ip=2.2.2.2 \
        options:remote_name=remote_cn
```

ca_cert.pem is the CA certificate. You need to set *remote_cn* as the common name (CN) of the remote host's certificate so that only the certificate with the expected CN can be trusted in this connection. It is preferable to use this method than 2) if there are many remote hosts since you don't have to copy every remote certificate to the local host.

Note: When using certificate-based authentication, you should not set *psk* in the interface options. When using psk-based authentication, you should not set *certificate*, *private_key*, *ca_cert*, *remote_cert*, and *remote_name*.

Plaintext Policies

When an IPsec tunnel is configured in this database, multiple independent components take responsibility for implementing it. *ovs-vswitchd* and its datapath handle packet forwarding to the tunnel and a separate daemon pushes the tunnel's IPsec policy configuration to the kernel or other entity that implements it. There is a race: if the former configuration completes before the latter, then packets sent by the local host over the tunnel can be transmitted in plaintext. Using this setting, OVS users can avoid this undesirable situation.

1. The default setting allows unencrypted packets to be sent before IPsec completes negotiation:

```
$ ovs-vsctl add-port br0 ipsec_gre0 -- \
    set interface ipsec_gre0 type=gre \
        options:remote_ip=2.2.2.2 \
        options:psk=swordfish
```

This setting should be used only and only if tunnel configuration is static and/or if there is firewall that can drop the plain packets that occasionally leak the tunnel unencrypted on OVSDB (re)configuration events.

2. Setting *ipsec_skb_mark* drops unencrypted packets by using *skb_mark* of tunnel packets:

```
$ ovs-vsctl set Open_vSwitch . other_config:ipsec_skb_mark=0/1
$ ovs-vsctl add-port br0 ipsec_gre0 -- \
    set interface ipsec_gre0 type=gre \
        options:remote_ip=2.2.2.2 \
        options:psk=swordfish
```

OVS IPsec drops unencrypted packets which carry the same *skb_mark* as *ipsec_skb_mark*. By setting the *ipsec_skb_mark* as 0/1, OVS IPsec prevents all unencrypted tunnel packets leaving the host since the default *skb_mark* value for tunnel packets are 0. This affects all OVS tunnels including those without IPsec being set up. You can install OpenFlow rules to whitelist those non-IPsec tunnels by setting the *skb_mark* of the tunnel traffic as non-zero value.

3. Setting *ipsec_skb_mark* as 1/1 only drops tunnel packets with *skb_mark* value being 1:

```
$ ovs-vsctl set Open_vSwitch . other_config:ipsec_skb_mark=1/1
$ ovs-vsctl add-port br0 ipsec_gre0 -- \
    set interface ipsec_gre0 type=gre \
        options:remote_ip=2.2.2.2 \
        options:psk=swordfish
```

Opposite to 2), this setting passes through unencrypted tunnel packets by default. To drop unencrypted IPsec tunnel traffic, you need to explicitly set `skb_mark` to a non-zero value for those tunnel traffic by installing OpenFlow rules.

Bug Reporting

If you think you may have found a bug with security implications, like

1. IPsec protected tunnel accepted packets that came unencrypted; OR
2. IPsec protected tunnel allowed packets to leave unencrypted

then please report such bugs according to [Open vSwitch's Security Process](#).

If the bug does not have security implications, then report it according to instructions in [Reporting Bugs in Open vSwitch](#).

5.1.3 Open vSwitch with SELinux

Security-Enhanced Linux (SELinux) is a Linux kernel security module that limits “the malicious things” that certain processes, including OVS, can do to the system in case they get compromised. In our case SELinux basically serves as the “second line of defense” that limits the things that OVS processes are allowed to do. The “first line of defense” is proper input validation that eliminates code paths that could be used by attacker to do any sort of “escape attacks”, such as file name escape, shell escape, command line argument escape, buffer escape. Since developers don’t always implement proper input validation, then SELinux Access Control’s goal is to confine damage of such attacks, if they turned out to be possible.

Besides Type Enforcement there are other SELinux features, but they are out of scope for this document.

Currently there are two SELinux policies for Open vSwitch:

- the one that ships with your Linux distribution (i.e. `selinux-policy-targeted` package)
- the one that ships with OVS (i.e. `openvswitch-selinux-policy` package)

Limitations

If Open vSwitch is directly started from command line, then it will run under `unconfined_t` SELinux domain that basically lets daemon to do whatever it likes. This is very important for developers to understand, because they might introduced code in OVS that invokes new system calls that SELinux policy did not anticipate. This means that their feature may have worked out just fine for them. However, if someone else would try to run the same code when Open vSwitch is started through `systemctl`, then Open vSwitch would get Permission Denied errors.

Currently the only distributions that enforce SELinux on OVS by default are RHEL, CentOS and Fedora. While Ubuntu and Debian also have some SELinux support, they run Open vSwitch under the unrestricted `unconfined` domain. Also, it seems that Ubuntu is leaning towards Apparmor that works slightly differently than SELinux.

SELinux and Open vSwitch are moving targets. What this means is that, if you solely rely on your Linux distribution’s SELinux policy, then this policy might not have correctly anticipated that a newer Open vSwitch version needs extra white list rules. However, if you solely rely on SELinux policy that ships with Open vSwitch, then Open vSwitch developers might not have correctly anticipated the feature set that your SELinux implementation supports.

Installation

Refer to *Fedora, RHEL 7.x Packaging for Open vSwitch* for instructions on how to build all Open vSwitch rpm packages.

Once the package is built, install it on your Linux distribution:

```
$ dnf install openvswitch-selinux-policy-2.4.1-1.el7.centos.noarch.rpm
```

Restart Open vSwitch:

```
$ systemctl restart openvswitch
```

Troubleshooting

When SELinux was implemented some of the standard system utilities acquired `-Z` flag (e.g. `ps -Z`, `ls -Z`). For example, to find out under which SELinux security domain process runs, use:

```
$ ps -AZ | grep ovs-vswitchd
system_u:system_r:openvswitch_t:s0 854 ?      ovs-vswitchd
```

To find out the SELinux label of file or directory, use:

```
$ ls -Z /etc/openvswitch/conf.db
system_u:object_r:openvswitch_rw_t:s0 /etc/openvswitch/conf.db
```

If, for example, SELinux policy for Open vSwitch is too strict, then you might see in Open vSwitch log files “Permission Denied” errors:

```
$ cat /var/log/openvswitch/ovs-vswitchd.log
vlog|INFO|opened log file /var/log/openvswitch/ovs-vswitchd.log
ovs_numa|INFO|Discovered 2 CPU cores on NUMA node 0
ovs_numa|INFO|Discovered 1 NUMA nodes and 2 CPU cores
reconnect|INFO|unix:/var/run/openvswitch/db.sock: connecting...
reconnect|INFO|unix:/var/run/openvswitch/db.sock: connected
netlink_socket|ERR|fcntl: Permission denied
dpif_netlink|ERR|Generic Netlink family 'ovs_datapath' does not exist.
      The Open vSwitch kernel module is probably not loaded.
dpif|WARN|failed to enumerate system datapaths: Permission denied
dpif|WARN|failed to create datapath ovs-system: Permission denied
```

However, not all “Permission denied” errors are caused by SELinux. So, before blaming too strict SELinux policy, make sure that indeed SELinux was the one that denied OVS access to certain resources, for example, run:

```
$ grep “openvswitch_t” /var/log/audit/audit.log | tail type=AVC msg=audit(1453235431.640:114671):
avc:      denied      {      getopt      }      for      pid=4583      comm=”ovs-vswitchd”      scon-
text=system_u:system_r:openvswitch_t:s0      tcontext=system_u:system_r:openvswitch_t:s0
tclass=netlink_generic_socket permissive=0
```

If SELinux denied OVS access to certain resources, then make sure that you have installed our SELinux policy package that “loosens” up distribution’s SELinux policy:

```
$ rpm -qa | grep openvswitch-selinux
openvswitch-selinux-policy-2.4.1-1.el7.centos.noarch
```

Then verify that this module was indeed loaded:

```
# semodule -l | grep openvswitch
openvswitch-custom    1.0
openvswitch           1.1.1
```

If you still see Permission denied errors, then take a look into `selinux/openvswitch.te.in` file in the OVS source tree and try to add white list rules. This is really simple, just run SELinux `audit2allow` tool:

```
$ grep "openvswitch_t" /var/log/audit/audit.log | audit2allow -M ovslocal
```

Contributing SELinux policy patches

Here are few things to consider before proposing SELinux policy patches to Open vSwitch developer mailing list:

1. The SELinux policy that resides in Open vSwitch source tree amends SELinux policy that ships with your distributions.

Implications of this are that it is assumed that the distribution's Open vSwitch SELinux module must be already loaded to satisfy dependencies.

2. The SELinux policy that resides in Open vSwitch source tree must work on all currently relevant Linux distributions.

Implications of this are that you should use only those SELinux policy features that are supported by the lowest SELinux version out there. Typically this means that you should test your SELinux policy changes on the oldest RHEL or CentOS version that this OVS version supports. Refer to *Fedora*, *RHEL 7.x Packaging for Open vSwitch* to find out this.

3. The SELinux policy is enforced only when state transition to `openvswitch_t` domain happens.

Implications of this are that perhaps instead of loosening SELinux policy you can do certain things at the time rpm package is installed.

Reporting Bugs

Report problems to bugs@openvswitch.org.

5.1.4 Open vSwitch with Libvirt

This document describes how to use Open vSwitch with Libvirt 0.9.11 or later. This document assumes that you followed *Open vSwitch on Linux, FreeBSD and NetBSD* or installed Open vSwitch from distribution packaging such as a .deb or .rpm. The Open vSwitch support is included by default in Libvirt 0.9.11. Consult www.libvirt.org for instructions on how to build the latest Libvirt, if your Linux distribution by default comes with an older Libvirt release.

Limitations

Currently there is no Open vSwitch support for networks that are managed by libvirt (e.g. NAT). As of now, only bridged networks are supported (those where the user has to manually create the bridge).

Setup

First, create the Open vSwitch bridge by using the `ovs-vsctl` utility (this must be done with administrative privileges):

```
$ ovs-vsctl add-br ovsbr
```

Once that is done, create a VM, if necessary, and edit its Domain XML file:

```
$ virsh edit <vm>
```

Lookup in the Domain XML file the `<interface>` section. There should be one such XML section for each interface the VM has:

```
<interface type='network'>
  <mac address='52:54:00:71:b1:b6' />
  <source network='default' />
  <address type='pci' domain='0x0000' bus='0x00' slot='0x03' function='0x0' />
</interface>
```

And change it to something like this:

```
<interface type='bridge'>
  <mac address='52:54:00:71:b1:b6' />
  <source bridge='ovsbr' />
  <virtualport type='openvswitch' />
  <address type='pci' domain='0x0000' bus='0x00' slot='0x03' function='0x0' />
</interface>
```

The interface type must be set to `bridge`. The `<source>` XML element specifies to which bridge this interface will be attached to. The `<virtualport>` element indicates that the bridge in `<source>` element is an Open vSwitch bridge.

Then (re)start the VM and verify if the guest's vnet interface is attached to the ovsbr bridge:

```
$ ovs-vsctl show
```

Troubleshooting

If the VM does not want to start, then try to run the `libvirtd` process either from the terminal, so that all errors are printed in console, or inspect Libvirt/Open vSwitch log files for possible root cause.

Bug Reporting

Report problems to bugs@openvswitch.org.

5.1.5 Open vSwitch with SSL

If you plan to configure Open vSwitch to connect across the network to an OpenFlow controller, then we recommend that you build Open vSwitch with OpenSSL. SSL support ensures integrity and confidentiality of the OpenFlow connections, increasing network security.

This document describes how to configure an Open vSwitch to connect to an OpenFlow controller over SSL. Refer to *Open vSwitch on Linux, FreeBSD and NetBSD* for instructions on building Open vSwitch with SSL support.

Open vSwitch uses TLS version 1.0 or later (TLSv1), as specified by RFC 2246, which is very similar to SSL version 3.0. TLSv1 was released in January 1999, so all current software and hardware should implement it.

This document assumes basic familiarity with public-key cryptography and public-key infrastructure.

SSL Concepts for OpenFlow

This section is an introduction to the public-key infrastructure architectures that Open vSwitch supports for SSL authentication.

To connect over SSL, every Open vSwitch must have a unique private/public key pair and a certificate that signs that public key. Typically, the Open vSwitch generates its own public/private key pair. There are two common ways to obtain a certificate for a switch:

- Self-signed certificates: The Open vSwitch signs its certificate with its own private key. In this case, each switch must be individually approved by the OpenFlow controller(s), since there is no central authority.

This is the only switch PKI model currently supported by NOX (<http://noxrepo.org>).

- Switch certificate authority: A certificate authority (the “switch CA”) signs each Open vSwitch’s public key. The OpenFlow controllers then check that any connecting switches’ certificates are signed by that certificate authority.

This is the only switch PKI model supported by the simple OpenFlow controller included with Open vSwitch.

Each Open vSwitch must also have a copy of the CA certificate for the certificate authority that signs OpenFlow controllers’ keys (the “controller CA” certificate). Typically, the same controller CA certificate is installed on all of the switches within a given administrative unit. There are two common ways for a switch to obtain the controller CA certificate:

- Manually copy the certificate to the switch through some secure means, e.g. using a USB flash drive, or over the network with “scp”, or even FTP or HTTP followed by manual verification.
- Open vSwitch “bootstrap” mode, in which Open vSwitch accepts and saves the controller CA certificate that it obtains from the OpenFlow controller on its first connection. Thereafter the switch will only connect to controllers signed by the same CA certificate.

Establishing a Public Key Infrastructure

Open vSwitch can make use of your existing public key infrastructure. If you already have a PKI, you may skip forward to the next section. Otherwise, if you do not have a PKI, the `ovs-pki` script included with Open vSwitch can help. To create an initial PKI structure, invoke it as:

```
$ ovs-pki init
```

This will create and populate a new PKI directory. The default location for the PKI directory depends on how the Open vSwitch tree was configured (to see the configured default, look for the `--dir` option description in the output of `ovs-pki --help`).

The `pki` directory contains two important subdirectories. The *controllerca* subdirectory contains controller CA files, including the following:

cacert.pem Root certificate for the controller certificate authority. Each Open vSwitch must have a copy of this file to allow it to authenticate valid controllers.

private/cakey.pem Private signing key for the controller certificate authority. This file must be kept secret. There is no need for switches or controllers to have a copy of it.

The *switchca* subdirectory contains switch CA files, analogous to those in the *controllerca* subdirectory:

cacert.pem Root certificate for the switch certificate authority. The OpenFlow controller must have this file to enable it to authenticate valid switches.

private/cakey.pem Private signing key for the switch certificate authority. This file must be kept secret. There is no need for switches or controllers to have a copy of it.

After you create the initial structure, you can create keys and certificates for switches and controllers with `ovs-pki`. Refer to the `ovs-pki(8)` manpage for complete details. A few examples of its use follow:

Controller Key Generation

To create a controller private key and certificate in files named `ctl-privkey.pem` and `ctl-cert.pem`, run the following on the machine that contains the PKI structure:

```
$ ovs-pki req+sign ctl controller
```

`ctl-privkey.pem` and `ctl-cert.pem` would need to be copied to the controller for its use at runtime. If, for testing purposes, you were to use `ovs-testcontroller`, the simple OpenFlow controller included with Open vSwitch, then the `-private-key` and `-certificate` options, respectively, would point to these files.

It is very important to make sure that no stray copies of `ctl-privkey.pem` are created, because they could be used to impersonate the controller.

Switch Key Generation with Self-Signed Certificates

If you are using self-signed certificates (see “SSL Concepts for OpenFlow”), this is one way to create an acceptable certificate for your controller to approve.

1. Run the following command on the Open vSwitch itself:

```
$ ovs-pki self-sign sc
```

Note: This command does not require a copy of any of the PKI files generated by `ovs-pki init`, and you should not copy them to the switch because some of them have contents that must remain secret for security.)

The `ovs-pki self-sign` command has the following output:

sc-privkey.pem the switch private key file. For security, the contents of this file must remain secret. There is ordinarily no need to copy this file off the Open vSwitch.

sc-cert.pem the switch certificate, signed by the switch’s own private key. Its contents are not a secret.

2. Optionally, copy `controllerca/cacert.pem` from the machine that has the OpenFlow PKI structure and verify that it is correct. (Otherwise, you will have to use CA certificate bootstrapping when you configure Open vSwitch in the next step.)
3. Configure Open vSwitch to use the keys and certificates (see “Configuring SSL Support”, below).

Switch Key Generation with a Switch PKI (Easy Method)

If you are using a switch PKI (see “SSL Concepts for OpenFlow”, above), this method of switch key generation is a little easier than the alternate method described below, but it is also a little less secure because it requires copying a sensitive private key from file from the machine hosting the PKI to the switch.

1. Run the following on the machine that contains the PKI structure:

```
$ ovs-pki req+sign sc switch
```

This command has the following output:

sc-privkey.pem the switch private key file. For security, the contents of this file must remain secret.

sc-cert.pem the switch certificate. Its contents are not a secret.

2. Copy `sc-privkey.pem` and `sc-cert.pem`, plus `controllerca/cacert.pem`, to the Open vSwitch.
3. Delete the copies of `sc-privkey.pem` and `sc-cert.pem` on the PKI machine and any other copies that may have been made in transit. It is very important to make sure that there are no stray copies of `sc-privkey.pem`, because they could be used to impersonate the switch.

Warning: Don't delete `controllerca/cacert.pem`! It is not security-sensitive and you will need it to configure additional switches.

4. Configure Open vSwitch to use the keys and certificates (see “Configuring SSL Support”, below).

Switch Key Generation with a Switch PKI (More Secure)

If you are using a switch PKI (see “SSL Concepts for OpenFlow”, above), then, compared to the previous method, the method described here takes a little more work, but it does not involve copying the private key from one machine to another, so it may also be a little more secure.

1. Run the following command on the Open vSwitch itself:

```
$ ovs-pki req sc
```

Note: This command does not require a copy of any of the PKI files generated by “`ovs-pki init`”, and you should not copy them to the switch because some of them have contents that must remain secret for security.

The “`ovs-pki req`” command has the following output:

sc-privkey.pem the switch private key file. For security, the contents of this file must remain secret. There is ordinarily no need to copy this file off the Open vSwitch.

sc-req.pem the switch “certificate request”, which is essentially the switch’s public key. Its contents are not a secret.

a fingerprint this is output on stdout.

2. Write the fingerprint down on a slip of paper and copy `sc-req.pem` to the machine that contains the PKI structure.
3. On the machine that contains the PKI structure, run:

```
$ ovs-pki sign sc switch
```

This command will output a fingerprint to stdout and request that you verify it. Check that it is the same as the fingerprint that you wrote down on the slip of paper before you answer “yes”.

`ovs-pki sign` creates a file named `sc-cert.pem`, which is the switch certificate. Its contents are not a secret.

4. Copy the generated `sc-cert.pem`, plus `controllerca/cacert.pem` from the PKI structure, to the Open vSwitch, and verify that they were copied correctly.

You may delete `sc-cert.pem` from the machine that hosts the PKI structure now, although it is not important that you do so.

Warning: Don't delete *controllerca/cacert.pem*! It is not security-sensitive and you will need it to configure additional switches.

5. Configure Open vSwitch to use the keys and certificates (see “Configuring SSL Support”, below).

Configuring SSL Support

SSL configuration requires three additional configuration files. The first two of these are unique to each Open vSwitch. If you used the instructions above to build your PKI, then these files will be named *sc-privkey.pem* and *sc-cert.pem*, respectively:

- A private key file, which contains the private half of an RSA or DSA key.

This file can be generated on the Open vSwitch itself, for the greatest security, or it can be generated elsewhere and copied to the Open vSwitch.

The contents of the private key file are secret and must not be exposed.

- A certificate file, which certifies that the private key is that of a trustworthy Open vSwitch.

This file has to be generated on a machine that has the private key for the switch certification authority, which should not be an Open vSwitch; ideally, it should be a machine that is not networked at all.

The certificate file itself is not a secret.

The third configuration file is typically the same across all the switches in a given administrative unit. If you used the instructions above to build your PKI, then this file will be named *cacert.pem*:

- The root certificate for the controller certificate authority. The Open vSwitch verifies it that is authorized to connect to an OpenFlow controller by verifying a signature against this CA certificate.

Once you have these files, configure `ovs-vswitchd` to use them using the `ovs-vsctl set-ssl` command, e.g.:

```
$ ovs-vsctl set-ssl /etc/openvswitch/sc-privkey.pem \  
    /etc/openvswitch/sc-cert.pem /etc/openvswitch/cacert.pem
```

Substitute the correct file names, of course, if they differ from the ones used above. You should use absolute file names (ones that begin with `/`), because `ovs-vswitchd`'s current directory is unrelated to the one from which you run `ovs-vsctl`.

If you are using self-signed certificates (see “SSL Concepts for OpenFlow”) and you did not copy *controllerca/cacert.pem* from the PKI machine to the Open vSwitch, then add the `--bootstrap` option, e.g.:

```
$ ovs-vsctl -- --bootstrap set-ssl /etc/openvswitch/sc-privkey.pem \  
    /etc/openvswitch/sc-cert.pem /etc/openvswitch/cacert.pem
```

After you have added all of these configuration keys, you may specify `ssl:` connection methods elsewhere in the configuration database. `tcp:` connection methods are still allowed even after SSL has been configured, so for security you should use only `ssl:` connections.

Reporting Bugs

Report problems to bugs@openvswitch.org.

5.1.6 Using LISP tunneling

LISP is a layer 3 tunneling mechanism, meaning that encapsulated packets do not carry Ethernet headers, and ARP requests shouldn't be sent over the tunnel. Because of this, there are some additional steps required for setting up LISP tunnels in Open vSwitch, until support for L3 tunnels will improve.

This guide assumes tunneling between two VMs connected to OVS bridges on different hypervisors reachable over IPv4. Of course, more than one VM may be connected to any of the hypervisors, and a hypervisor may communicate with several different hypervisors over the same lisp tunneling interface. A LISP “map-cache” can be implemented using flows, see example at the bottom of this file.

There are several scenarios:

1. the VMs have IP addresses in the same subnet and the hypervisors are also in a single subnet (although one different from the VM's);
2. the VMs have IP addresses in the same subnet but the hypervisors are separated by a router;
3. the VMs are in different subnets.

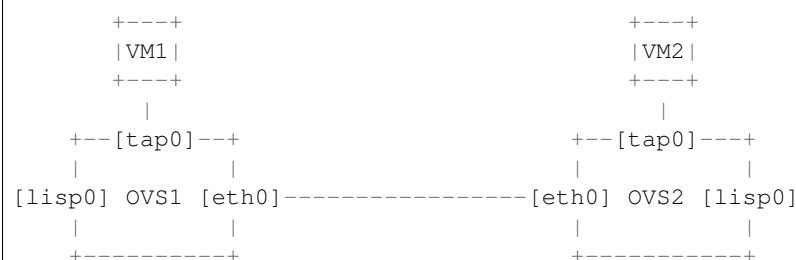
In cases 1) and 3) ARP resolution can work as normal: ARP traffic is configured not to go through the LISP tunnel. For case 1) ARP is able to reach the other VM, if both OVS instances default to MAC address learning. Case 3) requires the hypervisor be configured as the default router for the VMs.

In case 2) the VMs expect ARP replies from each other, but this is not possible over a layer 3 tunnel. One solution is to have static MAC address entries preconfigured on the VMs (e.g., `arp -f /etc/ethers` on startup on Unix based VMs), or have the hypervisor do proxy ARP. In this scenario, the `eth0` interfaces need not be added to the `br0` bridge in the examples below.

On the receiving side, the packet arrives without the original MAC header. The LISP tunneling code attaches a header with hard-coded source and destination MAC address `02:00:00:00:00:00`. This address has all bits set to 0, except the locally administered bit, in order to avoid potential collisions with existing allocations. In order for packets to reach their intended destination, the destination MAC address needs to be rewritten. This can be done using the flow table.

See below for an example setup, and the associated flow rules to enable LISP tunneling.

Diagram



On each hypervisor, interfaces `tap0`, `eth0`, and `lisp0` are added to a single bridge instance, and become numbered 1, 2, and 3 respectively:

```

$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 tap0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 lisp0 \
  -- set Interface lisp0 type=lisp options:remote_ip=flow options:key=flow
  
```

The last command sets up flow based tunneling on the `lisp0` interface. From the LISP point of view, this is like having the Tunnel Router map cache implemented as flow rules.

Flows on br0 should be configured as follows:

```
priority=3,dl_dst=02:00:00:00:00:00,action=mod_dl_dst:<VMx_MAC>,output:1
priority=2,in_port=1,dl_type=0x0806,action=NORMAL
priority=1,in_port=1,dl_type=0x0800,vlan_tci=0,nw_src=<EID_prefix>,action=set_field:
-><OVSx_IP>->tun_dst,output:3
priority=0,action=NORMAL
```

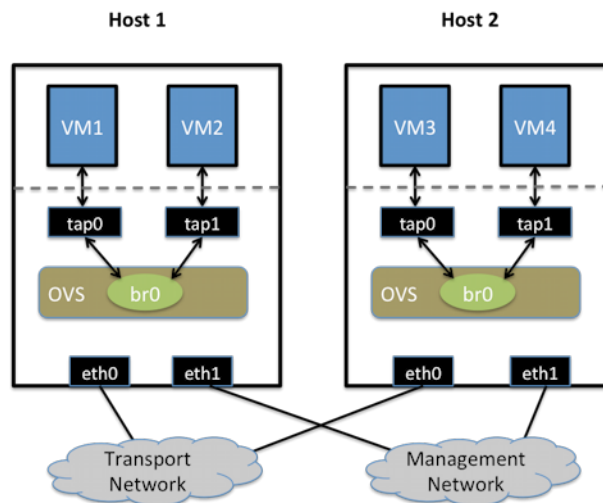
The third rule is like a map cache entry: the `<EID_prefix>` specified by the `nw_src` match field is mapped to the RLOC `<OVSx_IP>`, which is set as the tunnel destination for this particular flow.

Optionally, if you want to use Instance ID in a flow, you can add `set_tunnel:<IID>` to the action list.

5.1.7 Connecting VMs Using Tunnels

This document describes how to use Open vSwitch to allow VMs on two different hosts to communicate over port-based GRE tunnels.

Note: This guide covers the steps required to configure GRE tunneling. The same approach can be used for any of the other tunneling protocols supported by Open vSwitch.



Setup

This guide assumes the environment is configured as described below.

Two Physical Networks

- Transport Network

Ethernet network for tunnel traffic between hosts running OVS. Depending on the tunneling protocol being used (this cookbook uses GRE), some configuration of the physical switches may be required (for example, it may be necessary to adjust the MTU). Configuration of the physical switching hardware is outside the scope of this cookbook entry.

- Management Network

Strictly speaking this network is not required, but it is a simple way to give the physical host an IP address for remote access since an IP address cannot be assigned directly to a physical interface that is part of an OVS bridge.

Two Physical Hosts

The environment assumes the use of two hosts, named *host1* and *host2*. Both hosts are hypervisors running Open vSwitch. Each host has two NICs, *eth0* and *eth1*, which are configured as follows:

- *eth0* is connected to the Transport Network. *eth0* has an IP address that is used to communicate with Host2 over the Transport Network.
- *eth1* is connected to the Management Network. *eth1* has an IP address that is used to reach the physical host for management.

Four Virtual Machines

Each host will run two virtual machines (VMs). *vm1* and *vm2* are running on *host1*, while *vm3* and *vm4* are running on *host2*.

Each VM has a single interface that appears as a Linux device (e.g., *tap0*) on the physical host.

Note: For Xen/XenServer, VM interfaces appears as Linux devices with names like *vif1.0*. Other Linux systems may present these interfaces as *vnet0*, *vnet1*, etc.

Configuration Steps

Before you begin, you'll want to ensure that you know the IP addresses assigned to *eth0* on both *host1* and *host2*, as they will be needed during the configuration.

Perform the following configuration on *host1*.

1. Create an OVS bridge:

```
$ ovs-vsctl add-br br0
```

Note: You will *not* add *eth0* to the OVS bridge.

2. Boot *vm1* and *vm2* on *host1*. If the VMs are not automatically attached to OVS, add them to the OVS bridge you just created (the commands below assume *tap0* is for *vm1* and *tap1* is for *vm2*):

```
$ ovs-vsctl add-port br0 tap0
$ ovs-vsctl add-port br0 tap1
```

3. Add a port for the GRE tunnel:

```
$ ovs-vsctl add-port br0 gre0 \
  -- set interface gre0 type=gre options:remote_ip=<IP of eth0 on host2>
```

Create a mirrored configuration on *host2* using the same basic steps:

1. Create an OVS bridge, but do not add any physical interfaces to the bridge:

```
$ ovs-vsctl add-br br0
```

2. Launch *vm3* and *vm4* on *host2*, adding them to the OVS bridge if needed (again, *tap0* is assumed to be for *vm3* and *tap1* is assumed to be for *vm4*):

```
$ ovs-vsctl add-port br0 tap0
$ ovs-vsctl add-port br0 tap1
```

3. Create the GRE tunnel on *host2*, this time using the IP address for *eth0* on *host1* when specifying the *remote_ip* option:

```
$ ovs-vsctl add-port br0 gre0 -- set interface gre0 type=gre options:remote_ip=<IP of eth0 on host1>
```

Testing

Pings between any of the VMs should work, regardless of whether the VMs are running on the same host or different hosts.

Using `ip route show` (or equivalent command), the routing table of the operating system running inside the VM should show no knowledge of the IP subnets used by the hosts, only the IP subnet(s) configured within the VM's operating system. To help illustrate this point, it may be preferable to use very different IP subnet assignments within the guest VMs than what is used on the hosts.

Troubleshooting

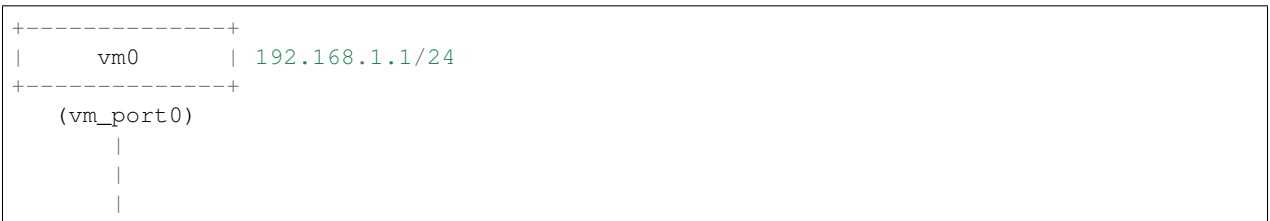
If connectivity between VMs on different hosts isn't working, check the following items:

- Make sure that *host1* and *host2* have full network connectivity over *eth0* (the NIC attached to the Transport Network). This may necessitate the use of additional IP routes or IP routing rules.
- Make sure that *gre0* on *host1* points to *eth0* on *host2*, and that *gre0* on *host2* points to *eth0* on *host1*.
- Ensure that all the VMs are assigned IP addresses on the same subnet; there is no IP routing functionality in this configuration.

5.1.8 Connecting VMs Using Tunnels (Userspace)

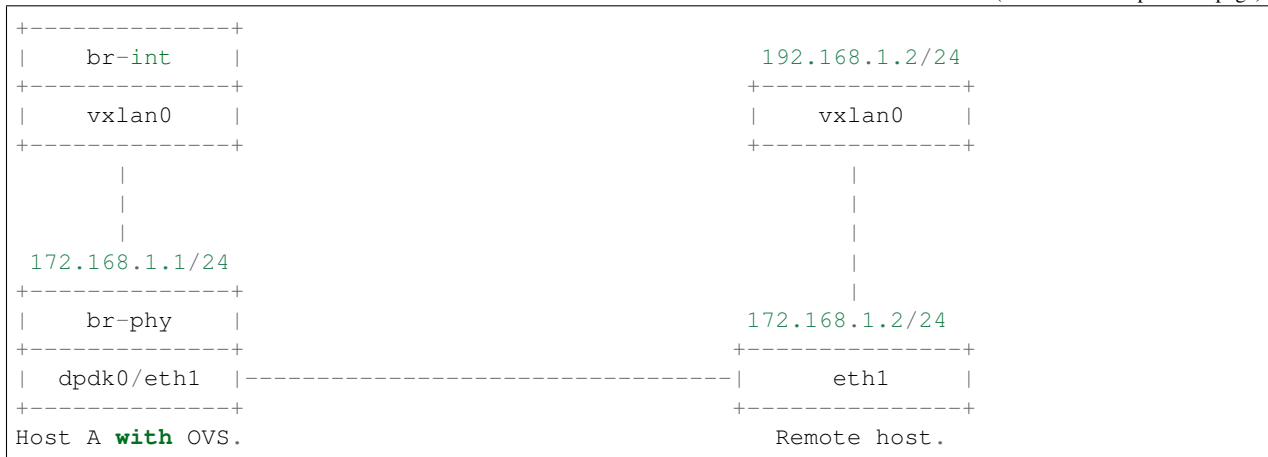
This document describes how to use Open vSwitch to allow VMs on two different hosts to communicate over VXLAN tunnels. Unlike *Connecting VMs Using Tunnels*, this configuration works entirely in userspace.

Note: This guide covers the steps required to configure VXLAN tunneling. The same approach can be used for any of the other tunneling protocols supported by Open vSwitch.



(continues on next page)

(continued from previous page)



Setup

This guide assumes the environment is configured as described below.

Two Physical Hosts

The environment assumes the use of two hosts, named *host1* and *host2*. We only detail the configuration of *host1* but a similar configuration can be used for *host2*. Both hosts should be configured with Open vSwitch (with or without the DPDK datapath), QEMU/KVM and suitable VM images. Open vSwitch should be running before proceeding.

Configuration Steps

Perform the following configuration on *host1*:

1. Create a `br-int` bridge:

```
$ ovs-vsctl --may-exist add-br br-int \
-- set Bridge br-int datapath_type=netdev \
-- br-set-external-id br-int bridge-id br-int \
-- set bridge br-int fail-mode=standalone
```

2. Add a port to this bridge. If using tap ports, first boot a VM and then add the port to the bridge:

```
$ ovs-vsctl add-port br-int tap0
```

If using DPDK vhost-user ports, add the port and then boot the VM accordingly, using `vm_port0` as the interface name:

```
$ ovs-vsctl add-port br-int vm_port0 \
    -- set Interface vm_port0 type=dpdkvhostuser
```

3. Configure the IP address of the VM interface *in the VM itself*:

```
$ ip addr add 192.168.1.1/24 dev eth0
$ ip link set eth0 up
```

4. On *host1*, add a port for the VXLAN tunnel:

```
$ ovs-vsctl add-port br-int vxlan0 \  
-- set interface vxlan0 type=vxlan options:remote_ip=172.168.1.2
```

Note: 172.168.1.2 is the remote tunnel end point address. On the remote host this will be 172.168.1.1

5. Create a br-phy bridge:

```
$ ovs-vsctl --may-exist add-br br-phy \  
-- set Bridge br-phy datapath_type=netdev \  
-- br-set-external-id br-phy bridge-id br-phy \  
-- set bridge br-phy fail-mode=standalone \  
   other_config:hwaddr=<mac address of eth1 interface>
```

Note: This additional bridge is required when running Open vSwitch in userspace rather than kernel-based Open vSwitch. The purpose of this bridge is to allow use of the kernel network stack for routing and ARP resolution. The datapath needs to look-up the routing table and ARP table to prepare the tunnel header and transmit data to the output port.

Note: eth1 is used rather than eth0. This is to ensure network connectivity is retained.

6. Attach eth1/dpdk0 to the br-phy bridge.

If the physical port eth1 is operating as a kernel network interface, run:

```
$ ovs-vsctl --timeout 10 add-port br-phy eth1  
$ ip addr add 172.168.1.1/24 dev br-phy  
$ ip link set br-phy up  
$ ip addr flush dev eth1 2>/dev/null  
$ ip link set eth1 up  
$ iptables -F
```

If instead the interface is a DPDK interface and bound to the igb_uio or vfio driver, run:

```
$ ovs-vsctl --timeout 10 add-port br-phy dpdk0 \  
-- set Interface dpdk0 type=dpdk options:dpdk-devargs=0000:06:00.0  
$ ip addr add 172.168.1.1/24 dev br-phy  
$ ip link set br-phy up  
$ iptables -F
```

The commands are different as DPDK interfaces are not managed by the kernel, thus, the port details are not visible to any ip commands.

Important: Attempting to use the kernel network commands for a DPDK interface will result in a loss of connectivity through eth1. Refer to *Basic Configuration* for more details.

Once complete, check the cached routes using ovs-appctl command:

```
$ ovs-appctl ovs/route/show
```

If the tunnel route is missing, adding it now:

```
$ ovs-appctl ovs/route/add 172.168.1.1/24 br-eth1
```

Repeat these steps if necessary for *host2*, but using 192.168.1.1 and 172.168.1.2 for the VM and tunnel interface IP addresses, respectively.

Testing

With this setup, ping to VXLAN target device (192.168.1.2) should work. Traffic will be VXLAN encapsulated and sent over the `eth1/dpdk0` interface.

Tunneling-related Commands

Tunnel routing table

To add route:

```
$ ovs-appctl ovs/route/add <IP address>/<prefix length> <output-bridge-name> <gw>
```

To see all routes configured:

```
$ ovs-appctl ovs/route/show
```

To delete route:

```
$ ovs-appctl ovs/route/del <IP address>/<prefix length>
```

To look up and display the route for a destination:

```
$ ovs-appctl ovs/route/lookup <IP address>
```

ARP

To see arp cache content:

```
$ ovs-appctl tnl/arp/show
```

To flush arp cache:

```
$ ovs-appctl tnl/arp/flush
```

To set a specific arp entry:

```
$ ovs-appctl tnl/arp/set <bridge> <IP address> <MAC address>
```

Ports

To check tunnel ports listening in `ovs-vswitchd`:

```
$ ovs-appctl tnl/ports/show
```

To set range for VxLan UDP source port:

```
$ ovs-appctl tnl/egress_port_range <num1> <num2>
```

To show current range:

```
$ ovs-appctl tnl/egress_port_range
```

Datapath

To check datapath ports:

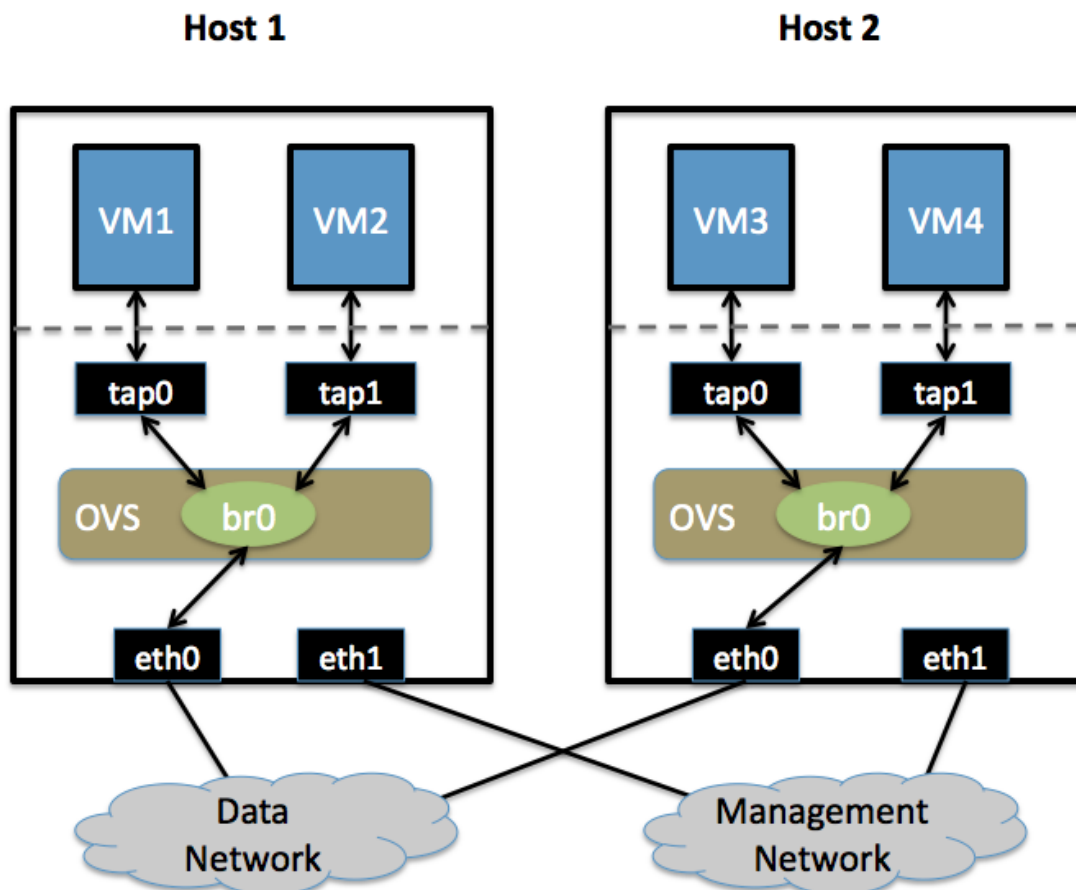
```
$ ovs-appctl dpif/show
```

To check datapath flows:

```
$ ovs-appctl dpif/dump-flows
```

5.1.9 Isolating VM Traffic Using VLANs

This document describes how to use Open vSwitch to isolate VM traffic using VLANs.



Setup

This guide assumes the environment is configured as described below.

Two Physical Networks

- Data Network

Ethernet network for VM data traffic, which will carry VLAN-tagged traffic between VMs. Your physical switch(es) must be capable of forwarding VLAN-tagged traffic and the physical switch ports should operate as VLAN trunks. (Usually this is the default behavior. Configuring your physical switching hardware is beyond the scope of this document.)

- Management Network

This network is not strictly required, but it is a simple way to give the physical host an IP address for remote access, since an IP address cannot be assigned directly to `eth0` (more on that in a moment).

Two Physical Hosts

The environment assumes the use of two hosts: *host1* and *host2*. Both hosts are running Open vSwitch. Each host has two NICs, `eth0` and `eth1`, which are configured as follows:

- `eth0` is connected to the Data Network. No IP address is assigned to `eth0`.
- `eth1` is connected to the Management Network (if necessary). `eth1` has an IP address that is used to reach the physical host for management.

Four Virtual Machines

Each host will run two virtual machines (VMs). *vm1* and *vm2* are running on *host1*, while *vm3* and *vm4* are running on *host2*.

Each VM has a single interface that appears as a Linux device (e.g., `tap0`) on the physical host.

Note: For Xen/XenServer, VM interfaces appears as Linux devices with names like `vif1.0`. Other Linux systems may present these interfaces as `vnet0`, `vnet1`, etc.

Configuration Steps

Perform the following configuration on *host1*:

1. Create an OVS bridge:

```
$ ovs-vsctl add-br br0
```

2. Add `eth0` to the bridge:

```
$ ovs-vsctl add-port br0 eth0
```

Note: By default, all OVS ports are VLAN trunks, so eth0 will pass all VLANs

Note: When you add eth0 to the OVS bridge, any IP addresses that might have been assigned to eth0 stop working. IP address assigned to eth0 should be migrated to a different interface before adding eth0 to the OVS bridge. This is the reason for the separate management connection via eth1.

3. Add *vm1* as an “access port” on VLAN 100. This means that traffic coming into OVS from VM1 will be untagged and considered part of VLAN 100:

```
$ ovs-vsctl add-port br0 tap0 tag=100
```

Add VM2 on VLAN 200:

```
$ ovs-vsctl add-port br0 tap1 tag=200
```

Repeat these steps on *host2*:

1. Setup a bridge with eth0 as a VLAN trunk:

```
$ ovs-vsctl add-br br0  
$ ovs-vsctl add-port br0 eth0
```

2. Add VM3 to VLAN 100:

```
$ ovs-vsctl add-port br0 tap0 tag=100
```

3. Add VM4 to VLAN 200:

```
$ ovs-vsctl add-port br0 tap1 tag=200
```

Validation

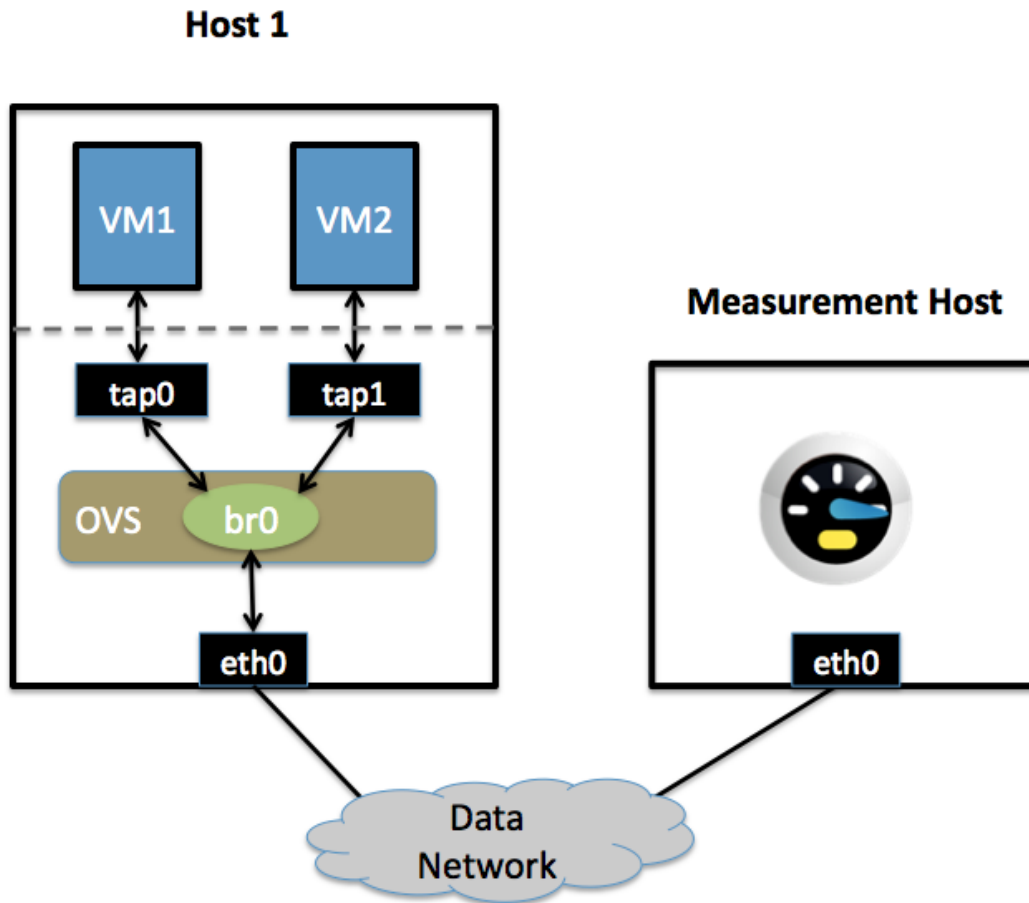
Pings from *vm1* to *vm3* should succeed, as these two VMs are on the same VLAN.

Pings from *vm2* to *vm4* should also succeed, since these VMs are also on the same VLAN as each other.

Pings from *vm1/vm3* to *vm2/vm4* should not succeed, as these VMs are on different VLANs. If you have a router configured to forward between the VLANs, then pings will work, but packets arriving at *vm3* should have the source MAC address of the router, not of *vm1*.

5.1.10 Quality of Service (QoS) Rate Limiting

This document explains how to use Open vSwitch to rate-limit traffic by a VM to either 1 Mbps or 10 Mbps.



Setup

This guide assumes the environment is configured as described below.

One Physical Network

- Data Network

Ethernet network for VM data traffic. This network is used to send traffic to and from an external host used for measuring the rate at which a VM is sending. For experimentation, this physical network is optional; you can instead connect all VMs to a bridge that is not connected to a physical interface and use a VM as the measurement host.

There may be other networks (for example, a network for management traffic), but this guide is only concerned with the Data Network.

Two Physical Hosts

The first host, named *host1*, is a hypervisor that runs Open vSwitch and has one NIC. This single NIC, *eth0*, is connected to the Data Network. Because it is participating in an OVS bridge, no IP address can be assigned on *eth0*.

The second host, named Measurement Host, can be any host capable of measuring throughput from a VM. For this guide, we use `netperf`, a free tool for testing the rate at which one host can send to another. The Measurement Host has only a single NIC, `eth0`, which is connected to the Data Network. `eth0` has an IP address that can reach any VM on `host1`.

Two VMs

Both VMs (`vm1` and `vm2`) run on `host1`.

Each VM has a single interface that appears as a Linux device (e.g., `tap0`) on the physical host.

Note: For Xen/XenServer, VM interfaces appears as Linux devices with names like `vif1.0`. Other Linux systems may present these interfaces as `vnet0`, `vnet1`, etc.

Configuration Steps

For both VMs, we modify the Interface table to configure an ingress policing rule. There are two values to set:

`ingress_policing_rate` the maximum rate (in Kbps) that this VM should be allowed to send

`ingress_policing_burst` a parameter to the policing algorithm to indicate the maximum amount of data (in Kb) that this interface can send beyond the policing rate.

To rate limit VM1 to 1 Mbps, use these commands:

```
$ ovs-vsctl set interface tap0 ingress_policing_rate=1000
$ ovs-vsctl set interface tap0 ingress_policing_burst=100
```

Similarly, to limit `vm2` to 10 Mbps, enter these commands on `host1`:

```
$ ovs-vsctl set interface tap1 ingress_policing_rate=10000
$ ovs-vsctl set interface tap1 ingress_policing_burst=1000
```

To see the current limits applied to VM1, run this command:

```
$ ovs-vsctl list interface tap0
```

Testing

To test the configuration, make sure `netperf` is installed and running on both VMs and on the Measurement Host. `netperf` consists of a client (`netperf`) and a server (`netserver`). In this example, we run `netserver` on the Measurement Host (installing Netperf usually starts `netserver` as a daemon, meaning this is running by default).

For this example, we assume that the Measurement Host has an IP of 10.0.0.100 and is reachable from both VMs.

From `vm1`, run this command:

```
$ netperf -H 10.0.0.100
```

This will cause VM1 to send TCP traffic as quickly as it can to the Measurement Host. After 10 seconds, this will output a series of values. We are interested in the “Throughput” value, which is measured in Mbps (10⁶ bits/sec). For VM1 this value should be near 1. Running the same command on VM2 should give a result near 10.

Troubleshooting

Open vSwitch uses the Linux [traffic-control](#) capability for rate-limiting. If you are not seeing the configured rate-limit have any effect, make sure that your kernel is built with “ingress qdisc” enabled, and that the user-space utilities (e.g., `/sbin/tc`) are installed.

Additional Information

Open vSwitch’s rate-limiting uses policing, which does not queue packets. It drops any packets beyond the specified rate. Specifying a larger burst size lets the algorithm be more forgiving, which is important for protocols like TCP that react severely to dropped packets. Setting a burst size of less than the MTU (e.g., 10 kb) should be avoided.

For TCP traffic, setting a burst size to be a sizeable fraction (e.g., > 10%) of the overall policy rate helps a flow come closer to achieving the full rate. If a burst size is set to be a large fraction of the overall rate, the client will actually experience an average rate slightly higher than the specific policing rate.

For UDP traffic, set the burst size to be slightly greater than the MTU and make sure that your performance tool does not send packets that are larger than your MTU (otherwise these packets will be fragmented, causing poor performance). For example, you can force netperf to send UDP traffic as 1000 byte packets by running:

```
$ netperf -H 10.0.0.100 -t UDP_STREAM -- -m 1000
```

5.1.11 How to Use the VTEP Emulator

This document explains how to use ovs-vtep, a VXLAN Tunnel Endpoint (VTEP) emulator that uses Open vSwitch for forwarding. VTEPs are the entities that handle VXLAN frame encapsulation and decapsulation in a network.

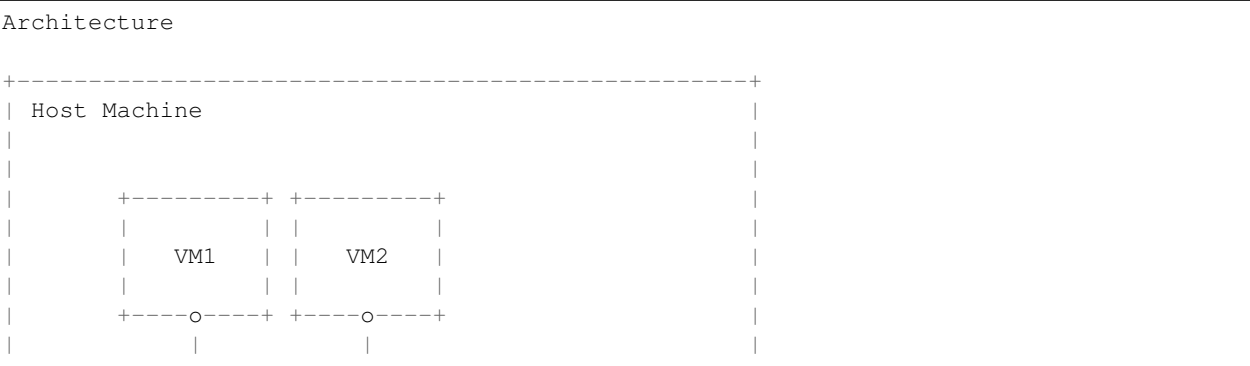
Requirements

The VTEP emulator is a Python script that invokes calls to tools like vtep-ctl and ovs-vsctl. It is only useful when Open vSwitch daemons like ovsdb-server and ovs-vswitchd are running and installed. To do this, either:

- Follow the instructions in [Open vSwitch on Linux, FreeBSD and NetBSD](#) (don’t start any daemons yet).
- Follow the instructions in [Debian Packaging for Open vSwitch](#) and then install the `openvswitch-vtep` package (if operating on a debian based machine). This will automatically start the daemons.

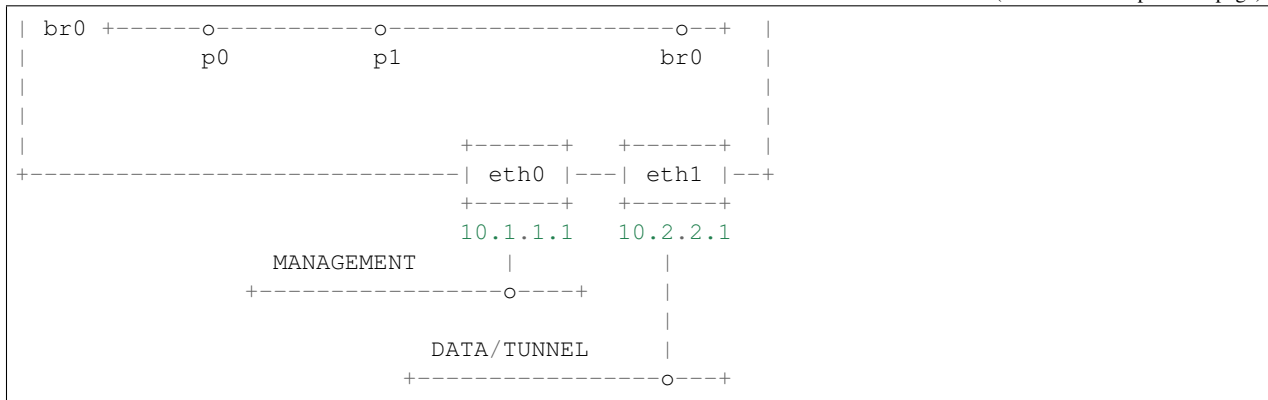
Design

At the end of this process, you should have the following setup:



(continues on next page)

(continued from previous page)



Some important points.

- We will use Open vSwitch to create our “physical” switch labeled `br0`
 - Our “physical” switch `br0` will have one internal port also named `br0` and two “physical” ports, namely `p0` and `p1`.
 - The host machine may have two external interfaces. We will use `eth0` for management traffic and `eth1` for tunnel traffic (One can use a single interface to achieve both). Please take note of their IP addresses in the diagram. You do not have to use exactly the same IP addresses. Just know that the above will be used in the steps below.
 - You can optionally connect physical machines instead of virtual machines to switch `br0`. In that case:
 - Make sure you have two extra physical interfaces in your host machine, `eth2` and `eth3`.
 - In the rest of this doc, replace `p0` with `eth2` and `p1` with `eth3`.
5. In addition to implementing `p0` and `p1` as physical interfaces, you can also optionally implement them as standalone TAP devices, or VM interfaces for simulation.
 6. Creating and attaching the VMs is outside the scope of this document and is included in the diagram for reference purposes only.

Startup

These instructions describe how to run with a single `ovsdb-server` instance that handles both the OVS and VTEP schema. You can skip steps 1-3 if you installed using the debian packages as mentioned in step 2 of the “Requirements” section.

1. Create the initial OVS and VTEP schemas:

```
$ ovsdb-tool create /etc/openvswitch/ovs.db vswitchd/vswitch.ovsschema
$ ovsdb-tool create /etc/openvswitch/vtep.db vtep/vtep.ovsschema
```

```

2. Start `ovsdb-server` and have it handle both databases:

```
$ ovsdb-server --pidfile --detach --log-file \
 --remote punix:/var/run/openvswitch/db.sock \
 --remote=db:hardware_vtep,Global,managers \
 /etc/openvswitch/ovs.db /etc/openvswitch/vtep.db
```

3. Start `ovs-vswitchd` as normal:

```
$ ovs-vswitchd --log-file --detach --pidfile \
 unix:/var/run/openvswitch/db.sock
```

4. Create a “physical” switch and its ports in OVS:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 p0
$ ovs-vsctl add-port br0 p1
```

5. Configure the physical switch in the VTEP database:

```
$ vtep-ctl add-ps br0
$ vtep-ctl set Physical_Switch br0 tunnel_ips=10.2.2.1
```

6. Start the VTEP emulator. If you installed the components following *Open vSwitch on Linux, FreeBSD and NetBSD*, run the following from the `vtep` directory:

```
$./ovs-vtep --log-file=/var/log/openvswitch/ovs-vtep.log \
 --pidfile=/var/run/openvswitch/ovs-vtep.pid \
 --detach br0
```

If the installation was done by installing the `openvswitch-vtep` package, you can find `ovs-vtep` at `/usr/share/openvswitch/scripts`.

7. Configure the VTEP database’s manager to point at an NVC:

```
$ vtep-ctl set-manager tcp:<CONTROLLER IP>:6640
```

Where `<CONTROLLER IP>` is your controller’s IP address that is accessible via the Host Machine’s `eth0` interface.

## Simulating an NVC

A VTEP implementation expects to be driven by a Network Virtualization Controller (NVC), such as NSX. If one does not exist, it’s possible to use `vtep-ctl` to simulate one:

1. Create a logical switch:

```
$ vtep-ctl add-ls ls0
```

2. Bind the logical switch to a port:

```
$ vtep-ctl bind-ls br0 p0 0 ls0
$ vtep-ctl set Logical_Switch ls0 tunnel_key=33
```

3. Direct unknown destinations out a tunnel.

For handling L2 broadcast, multicast and unknown unicast traffic, packets can be sent to all members of a logical switch referenced by a physical switch. The “unknown-dst” address below is used to represent these packets. There are different modes to replicate the packets. The default mode of replication is to send the traffic to a service node, which can be a hypervisor, server or appliance, and let the service node handle replication to other transport nodes (hypervisors or other VTEP physical switches). This mode is called *service node* replication. An alternate mode of replication, called *source node* replication, involves the source node sending to all other transport nodes. Hypervisors are always responsible for doing their own replication for locally attached VMs in both modes. Service node mode is the default. Service node replication mode is considered a basic requirement because it only requires sending the packet to a single transport node. The following configuration is for service node replication mode as only a single transport node destination is specified for the unknown-dst address:

```
$ vtep-ctl add-mcast-remote ls0 unknown-dst 10.2.2.2
```

4. Optionally, change the replication mode from a default of `service_node` to `source_node`, which can be done at the logical switch level:

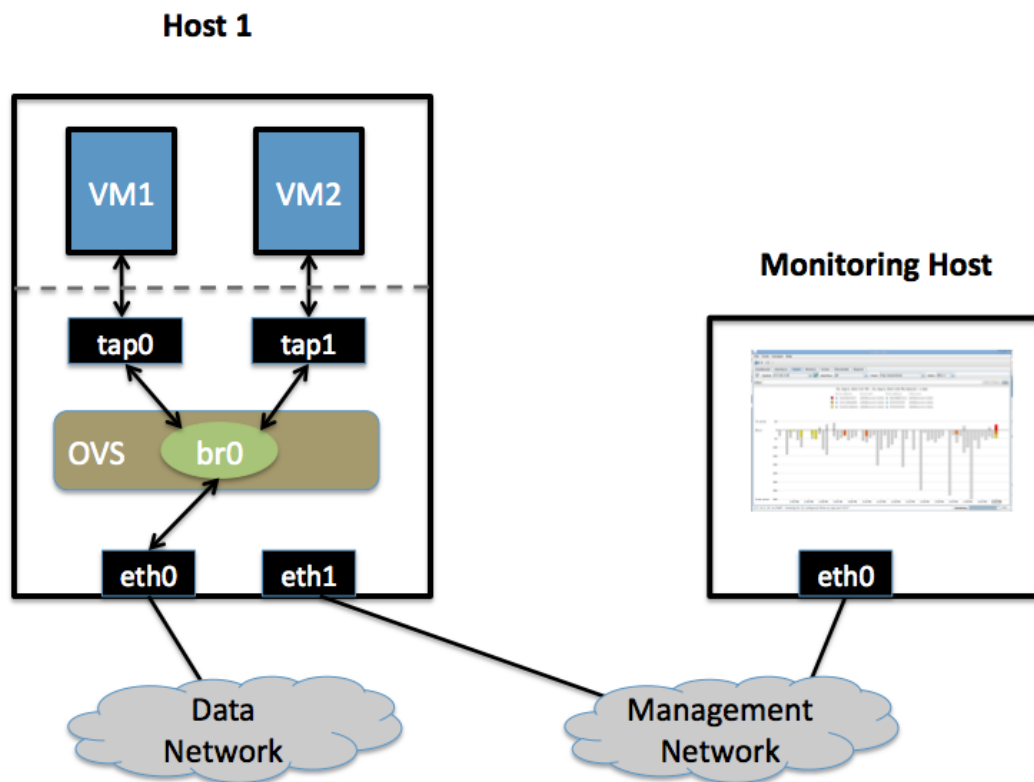
```
$ vtep-ctl set-replication-mode ls0 source_node
```

5. Direct unicast destinations out a different tunnel:

```
$ vtep-ctl add-ucast-remote ls0 00:11:22:33:44:55 10.2.2.3
```

## 5.1.12 Monitoring VM Traffic Using sFlow

This document describes how to use Open vSwitch to monitor traffic sent between two VMs on the same host using an sFlow collector. VLANs.



### Setup

This guide assumes the environment is configured as described below.

## Two Physical Networks

- Data Network

Ethernet network for VM data traffic. For experimentation, this physical network is optional. You can instead connect all VMs to a bridge that is not connected to a physical interface.

- Management Network

This network must exist, as it is used to send sFlow data from the agent to the remote collector.

## Two Physical Hosts

The environment assumes the use of two hosts: *host1* and *hostMon*. *host* is a hypervisor that run Open vSwitch and has two NICs:

- *eth0* is connected to the Data Network. No IP address can be assigned on *eth0* because it is part of an OVS bridge.
- *eth1* is connected to the Management Network. *eth1* has an IP address for management traffic, including sFlow.

*hostMon* can be any computer that can run the sFlow collector. For this cookbook entry, we use [sFlowTrend](#), a free sFlow collector that is a simple cross-platform Java download. Other sFlow collectors should work equally well. *hostMon* has a single NIC, *eth0*, that is connected to the Management Network. *eth0* has an IP address that can reach *eth1* on *host1*.

## Two Virtual Machines

This guide uses two virtual machines - *vm1* and *vm2*- running on *host1*.

---

**Note:** For Xen/XenServer, VM interfaces appears as Linux devices with names like *vif1.0*. Other Linux systems may present these interfaces as *vnet0*, *vnet1*, etc.

---

## Configuration Steps

On *host1*, define the following configuration values in your shell environment:

```
COLLECTOR_IP=10.0.0.1
COLLECTOR_PORT=6343
AGENT_IP=eth1
HEADER_BYTES=128
SAMPLING_N=64
POLLING_SECS=10
```

Port 6343 (*COLLECTOR\_PORT*) is the default port number for sFlowTrend. If you are using an sFlow collector other than sFlowTrend, set this value to the appropriate port for your particular collector. Set your own IP address for the collector in the place of 10.0.0.1 (*COLLECTOR\_IP*). Setting the *AGENT\_IP* value to *eth1* indicates that the sFlow agent should send traffic from *eth1*'s IP address. The other values indicate settings regarding the frequency and type of packet sampling that sFlow should perform.

Still on *host1*, run the following command to create an sFlow configuration and attach it to bridge *br0*:

```
$ ovs-vsctl -- --id=@sflow create sflow agent=${AGENT_IP} \
 target="\${COLLECTOR_IP}:${COLLECTOR_PORT}" header=${HEADER_BYTES} \
 sampling=${SAMPLING_N} polling=${POLLING_SECS} \
 -- set bridge br0 sflow=@sflow
```

Make note of the UUID that is returned by this command; this value is necessary to remove the sFlow configuration.

On *hostMon*, go to the [sFlowTrend](#) and click “Install” in the upper right-hand corner. If you have Java installed, this will download and start the sFlowTrend application. Once sFlowTrend is running, the light in the lower right-hand corner of the sFlowTrend application should blink green to indicate that the collector is receiving traffic.

The sFlow configuration is now complete, and sFlowTrend on *hostMon* should be receiving sFlow data from OVS on *host1*.

To configure sFlow on additional bridges, just replace `br0` in the above command with a different bridge name.

To remove sFlow configuration from a bridge (in this case, `br0`), run this command, where “sFlow UUID” is the UUID returned by the command used to set the sFlow configuration initially:

```
$ ovs-vsctl remove bridge br0 sflow <sFlow UUID>
```

To see all current sets of sFlow configuration parameters, run:

```
$ ovs-vsctl list sflow
```

## Troubleshooting

If sFlow data isn’t being collected and displayed by sFlowTrend, check the following items:

- Make sure the VMs are sending/receiving network traffic over bridge `br0`, preferably to multiple other hosts and using a variety of protocols.
- To confirm that the agent is sending traffic, check that running the following command shows that the agent on the physical server is sending traffic to the collector IP address (change the port below to match the port your collector is using):

```
$ tcpdump -ni eth1 udp port 6343
```

If no traffic is being sent, there is a problem with the configuration of OVS. If traffic is being sent but nothing is visible in the sFlowTrend user interface, this may indicate a configuration problem with the collector.

Check to make sure the host running the collector (*hostMon*) does not have a firewall that would prevent UDP port 6343 from reaching the collector.

## Credit

This document is heavily based on content from Neil McKee at InMon:

- <https://mail.openvswitch.org/pipermail/ovs-dev/2010-July/165245.html>
- <https://blog.sflow.com/2010/01/open-vswitch.html>

---

**Note:** The configuration syntax is out of date, but the high-level descriptions are correct.

---

### 5.1.13 Using Open vSwitch with DPDK

This document describes how to use Open vSwitch with DPDK datapath.

---

**Important:** Using the DPDK datapath requires building OVS with DPDK support. The mapping of OVS version to DPDK can vary between releases. For version mapping information refer to [releases FAQ](#). For build instructions refer to [Open vSwitch with DPDK](#).

---

#### Ports and Bridges

ovs-vsctl can be used to set up bridges and other Open vSwitch features. Bridges should be created with a `datapath_type=netdev`:

```
$ ovs-vsctl add-br br0 -- set bridge br0 datapath_type=netdev
```

ovs-vsctl can also be used to add DPDK devices. ovs-vswitchd should print the number of dpdk devices found in the log file:

```
$ ovs-vsctl add-port br0 dpdk-p0 -- set Interface dpdk-p0 type=dpdk \
 options:dpdk-devargs=0000:01:00.0
$ ovs-vsctl add-port br0 dpdk-p1 -- set Interface dpdk-p1 type=dpdk \
 options:dpdk-devargs=0000:01:00.1
```

Some NICs (i.e. Mellanox ConnectX-3) have only one PCI address associated with multiple ports. Using a PCI device like above won't work. Instead, below usage is suggested:

```
$ ovs-vsctl add-port br0 dpdk-p0 -- set Interface dpdk-p0 type=dpdk \
 options:dpdk-devargs="class=eth,mac=00:11:22:33:44:55"
$ ovs-vsctl add-port br0 dpdk-p1 -- set Interface dpdk-p1 type=dpdk \
 options:dpdk-devargs="class=eth,mac=00:11:22:33:44:56"
```

---

**Important:** Hotplugging physical interfaces is not supported using the above syntax. This is expected to change with the release of DPDK v18.05. For information on hotplugging physical interfaces, you should instead refer to [Hotplugging](#).

---

After the DPDK ports get added to switch, a polling thread continuously polls DPDK devices and consumes 100% of the core, as can be checked from `top` and `ps` commands:

```
$ top -H
$ ps -eLo pid,psr,comm | grep pmd
```

Creating bonds of DPDK interfaces is slightly different to creating bonds of system interfaces. For DPDK, the interface type and devargs must be explicitly set. For example:

```
$ ovs-vsctl add-bond br0 dpdkbond p0 p1 \
 -- set Interface p0 type=dpdk options:dpdk-devargs=0000:01:00.0 \
 -- set Interface p1 type=dpdk options:dpdk-devargs=0000:01:00.1
```

To stop ovs-vswitchd & delete bridge, run:

```
$ ovs-appctl -t ovs-vswitchd exit
$ ovs-appctl -t ovsdb-server exit
$ ovs-vsctl del-br br0
```

## OVS with DPDK Inside VMs

Additional configuration is required if you want to run `ovs-vswitchd` with DPDK backend inside a QEMU virtual machine. `ovs-vswitchd` creates separate DPDK TX queues for each CPU core available. This operation fails inside QEMU virtual machine because, by default, VirtIO NIC provided to the guest is configured to support only single TX queue and single RX queue. To change this behavior, you need to turn on `mq` (multiqueue) property of all `virtio-net-pci` devices emulated by QEMU and used by DPDK. You may do it manually (by changing QEMU command line) or, if you use Libvirt, by adding the following string to `<interface>` sections of all network devices used by DPDK:

```
<driver name='vhost' queues='N' />
```

where:

**N** determines how many queues can be used by the guest.

This requires QEMU `>= 2.2`.

## PHY-PHY

Add a userspace bridge and two dpdk (PHY) ports:

```
Add userspace bridge
$ ovs-vsctl add-br br0 -- set bridge br0 datapath_type=netdev

Add two dpdk ports
$ ovs-vsctl add-port br0 phy0 -- set Interface phy0 type=dpdk \
 options:dpdk-devargs=0000:01:00.0 ofport_request=1

$ ovs-vsctl add-port br0 phy1 -- set Interface phy1 type=dpdk \
 options:dpdk-devargs=0000:01:00.1 ofport_request=2
```

Add test flows to forward packets between DPDK port 0 and port 1:

```
Clear current flows
$ ovs-ofctl del-flows br0

Add flows between port 1 (phy0) to port 2 (phy1)
$ ovs-ofctl add-flow br0 in_port=1,action=output:2
$ ovs-ofctl add-flow br0 in_port=2,action=output:1
```

Transmit traffic into either port. You should see it returned via the other.

## PHY-VM-PHY (vHost Loopback)

Add a userspace bridge, two dpdk (PHY) ports, and two `dpdkvhostuser` ports:

```
Add userspace bridge
$ ovs-vsctl add-br br0 -- set bridge br0 datapath_type=netdev

Add two dpdk ports
$ ovs-vsctl add-port br0 phy0 -- set Interface phy0 type=dpdk \
 options:dpdk-devargs=0000:01:00.0 ofport_request=1

$ ovs-vsctl add-port br0 phy1 -- set Interface phy1 type=dpdk \
 options:dpdk-devargs=0000:01:00.1 ofport_request=2
```

(continues on next page)



(continued from previous page)

```
Add two dpdkvhostuser ports
$ ovs-vsctl add-port br0 dpdkvhostuser0 \
 -- set Interface dpdkvhostuser0 type=dpdkvhostuser ofport_request=3
$ ovs-vsctl add-port br0 dpdkvhostuser1 \
 -- set Interface dpdkvhostuser1 type=dpdkvhostuser ofport_request=4
```

Add test flows to forward packets between DPDK devices and VM ports:

```
Clear current flows
$ ovs-ofctl del-flows br0

Add flows
$ ovs-ofctl add-flow br0 in_port=1,action=output:3
$ ovs-ofctl add-flow br0 in_port=3,action=output:1
$ ovs-ofctl add-flow br0 in_port=4,action=output:2
$ ovs-ofctl add-flow br0 in_port=2,action=output:4

Dump flows
$ ovs-ofctl dump-flows br0
```

Create a VM using the following configuration:

Configuration	Values	Comments
QEMU version	2.2.0	n/a
QEMU thread affinity	core 5	taskset 0x20
Memory	4GB	n/a
Cores	2	n/a
Qcow2 image	CentOS7	n/a
mrg_rxbuf	off	n/a

You can do this directly with QEMU via the `qemu-system-x86_64` application:

```
$ export VM_NAME=vhost-vm
$ export GUEST_MEM=3072M
$ export QCOW2_IMAGE=/root/CentOS7_x86_64.qcow2
$ export VHOST_SOCK_DIR=/usr/local/var/run/openvswitch

$ taskset 0x20 qemu-system-x86_64 -name $VM_NAME -cpu host -enable-kvm \
-m $GUEST_MEM -drive file=$QCOW2_IMAGE --nographic -snapshot \
-numa node,memdev=mem -mem-prealloc -smp sockets=1,cores=2 \
-object memory-backend-file,id=mem,size=$GUEST_MEM,mem-path=/dev/hugepages,share=on,
↪ \
-chardev socket,id=char0,path=$VHOST_SOCK_DIR/dpdkvhostuser0 \
-netdev type=vhost-user,id=mynet1,chardev=char0,vhostforce \
-device virtio-net-pci,mac=00:00:00:00:00:01,netdev=mynet1,mrg_rxbuf=off \
-chardev socket,id=char1,path=$VHOST_SOCK_DIR/dpdkvhostuser1 \
-netdev type=vhost-user,id=mynet2,chardev=char1,vhostforce \
-device virtio-net-pci,mac=00:00:00:00:00:02,netdev=mynet2,mrg_rxbuf=off
```

For a explanation of this command, along with alternative approaches such as booting the VM via libvirt, refer to [DPDK vHost User Ports](#).

Once the guest is configured and booted, configure DPDK packet forwarding within the guest. To accomplish this, build the `testpmd` application as described in [DPDK in the Guest](#). Once compiled, run the application:

```
$ cd $DPDK_DIR/app/test-pmd;
$./testpmd -c 0x3 -n 4 --socket-mem 1024 -- \
 --burst=64 -i --txqflags=0xf00 --disable-hw-vlan
$ set fwd mac retry
$ start
```

When you finish testing, bind the vNICs back to kernel:

```
$ $DPDK_DIR/usertools/dpdk-devbind.py --bind=virtio-pci 0000:00:03.0
$ $DPDK_DIR/usertools/dpdk-devbind.py --bind=virtio-pci 0000:00:04.0
```

---

**Note:** Valid PCI IDs must be passed in above example. The PCI IDs can be retrieved like so:

```
$ $DPDK_DIR/usertools/dpdk-devbind.py --status
```

---

More information on the dpdkvhostuser ports can be found in *DPDK vHost User Ports*.

### PHY-VM-PHY (vHost Loopback) (Kernel Forwarding)

*PHY-VM-PHY (vHost Loopback)* details steps for PHY-VM-PHY loopback testcase and packet forwarding using DPDK testpmd application in the Guest VM. For users wishing to do packet forwarding using kernel stack below, you need to run the below commands on the guest:

```
$ ip addr add 1.1.1.2/24 dev eth1
$ ip addr add 1.1.2.2/24 dev eth2
$ ip link set eth1 up
$ ip link set eth2 up
$ systemctl stop firewalld.service
$ systemctl stop iptables.service
$ sysctl -w net.ipv4.ip_forward=1
$ sysctl -w net.ipv4.conf.all.rp_filter=0
$ sysctl -w net.ipv4.conf.eth1.rp_filter=0
$ sysctl -w net.ipv4.conf.eth2.rp_filter=0
$ route add -net 1.1.2.0/24 eth2
$ route add -net 1.1.1.0/24 eth1
$ arp -s 1.1.2.99 DE:AD:BE:EF:CA:FE
$ arp -s 1.1.1.99 DE:AD:BE:EF:CA:EE
```

### PHY-VM-PHY (vHost Multiqueue)

vHost Multiqueue functionality can also be validated using the PHY-VM-PHY configuration. To begin, follow the steps described in *PHY-PHY* to create and initialize the database, start ovs-vswitchd and add dpdk-type devices to bridge br0. Once complete, follow the below steps:

#### 1. Configure PMD and RXQs.

For example, set the number of dpdk port rx queues to at least 2. The number of rx queues at vhost-user interface gets automatically configured after virtio device connection and doesn't need manual configuration:

```
$ ovs-vsctl set Open_vSwitch . other_config:pmd-cpu-mask=0xc
$ ovs-vsctl set Interface phy0 options:n_rxq=2
$ ovs-vsctl set Interface phy1 options:n_rxq=2
```

## 2. Instantiate Guest VM using QEMU cmdline

We must configure with appropriate software versions to ensure this feature is supported.

Table 1: VM Configuration

Setting	Value
QEMU version	2.5.0
QEMU thread affinity	2 cores (taskset 0x30)
Memory	4 GB
Cores	2
Distro	Fedora 22
Multiqueue	Enabled

To do this, instantiate the guest as follows:

```
$ export VM_NAME=vhost-vm
$ export GUEST_MEM=4096M
$ export QCOW2_IMAGE=/root/Fedora22_x86_64.qcow2
$ export VHOST_SOCK_DIR=/usr/local/var/run/openvswitch
$ taskset 0x30 qemu-system-x86_64 -cpu host -smp 2,cores=2 -m 4096M \
 -drive file=$QCOW2_IMAGE --enable-kvm -name $VM_NAME \
 -nographic -numa node,memdev=mem -mem-prealloc \
 -object memory-backend-file,id=mem,size=$GUEST_MEM,mem-path=/dev/hugepages,
 →share=on \
 -chardev socket,id=char1,path=$VHOST_SOCK_DIR/dpdkvhostuser0 \
 -netdev type=vhost-user,id=mynet1,chardev=char1,vhostforce,queues=2 \
 -device virtio-net-pci,mac=00:00:00:00:00:01,netdev=mynet1,mq=on,vectors=6 \
 -chardev socket,id=char2,path=$VHOST_SOCK_DIR/dpdkvhostuser1 \
 -netdev type=vhost-user,id=mynet2,chardev=char2,vhostforce,queues=2 \
 -device virtio-net-pci,mac=00:00:00:00:00:02,netdev=mynet2,mq=on,vectors=6
```

**Note:** Queue value above should match the queues configured in OVS, The vector value should be set to “number of queues x 2 + 2”

## 3. Configure the guest interface

Assuming there are 2 interfaces in the guest named eth0, eth1 check the channel configuration and set the number of combined channels to 2 for virtio devices:

```
$ ethtool -l eth0
$ ethtool -L eth0 combined 2
$ ethtool -L eth1 combined 2
```

More information can be found in vHost walkthrough section.

## 4. Configure kernel packet forwarding

Configure IP and enable interfaces:

```
$ ip addr add 5.5.5.1/24 dev eth0
$ ip addr add 90.90.90.1/24 dev eth1
$ ip link set eth0 up
$ ip link set eth1 up
```

Configure IP forwarding and add route entries:

```
$ sysctl -w net.ipv4.ip_forward=1
$ sysctl -w net.ipv4.conf.all.rp_filter=0
$ sysctl -w net.ipv4.conf.eth0.rp_filter=0
$ sysctl -w net.ipv4.conf.eth1.rp_filter=0
$ ip route add 2.1.1.0/24 dev eth1
$ route add default gw 2.1.1.2 eth1
$ route add default gw 90.90.90.90 eth1
$ arp -s 90.90.90.90 DE:AD:BE:EF:CA:FE
$ arp -s 2.1.1.2 DE:AD:BE:EF:CA:FA
```

Check traffic on multiple queues:

```
$ cat /proc/interrupts | grep virtio
```

### Flow Hardware Offload (Experimental)

The flow hardware offload is disabled by default and can be enabled by:

```
$ ovs-vsctl set Open_vSwitch . other_config:hw-offload=true
```

So far only partial flow offload is implemented. Moreover, it only works with PMD drivers have the `rte_flow` action “MARK + RSS” support.

The validated NICs are:

- Mellanox (ConnectX-4, ConnectX-4 Lx, ConnectX-5)
- Napatech (NT200B01)

Supported protocols for hardware offload are: - L2: Ethernet, VLAN - L3: IPv4, IPv6 - L4: TCP, UDP, SCTP, ICMP

### Further Reading

More detailed information can be found in the [DPDK topics section](#) of the documentation. These guides are listed below.

## 5.2 OVN

### 5.2.1 Open Virtual Networking With Docker

This document describes how to use Open Virtual Networking with Docker 1.9.0 or later.

---

**Important:** Requires Docker version 1.9.0 or later. Only Docker 1.9.0+ comes with support for multi-host networking. Consult [www.docker.com](http://www.docker.com) for instructions on how to install Docker.

---

---

**Note:** You must build and install Open vSwitch before proceeding with the below guide. Refer to [Installing Open vSwitch](#) for more information.

---

## Setup

For multi-host networking with OVN and Docker, Docker has to be started with a distributed key-value store. For example, if you decide to use consul as your distributed key-value store and your host IP address is `$HOST_IP`, start your Docker daemon with:

```
$ docker daemon --cluster-store=consul://127.0.0.1:8500 \
 --cluster-advertise=$HOST_IP:0
```

OVN provides network virtualization to containers. OVN's integration with Docker currently works in two modes - the "underlay" mode or the "overlay" mode.

In the "underlay" mode, OVN requires a OpenStack setup to provide container networking. In this mode, one can create logical networks and can have containers running inside VMs, standalone VMs (without having any containers running inside them) and physical machines connected to the same logical network. This is a multi-tenant, multi-host solution.

In the "overlay" mode, OVN can create a logical network amongst containers running on multiple hosts. This is a single-tenant (extendable to multi-tenants depending on the security characteristics of the workloads), multi-host solution. In this mode, you do not need a pre-created OpenStack setup.

For both the modes to work, a user has to install and start Open vSwitch in each VM/host that they plan to run their containers on.

## The "overlay" mode

**Note:** OVN in "overlay" mode needs a minimum Open vSwitch version of 2.5.

### 1. Start the central components.

OVN architecture has a central component which stores your networking intent in a database. On one of your machines, with an IP Address of `$CENTRAL_IP`, where you have installed and started Open vSwitch, you will need to start some central components.

Start `ovn-northd` daemon. This daemon translates networking intent from Docker stored in the `OVN_Northbound` database to logical flows in `OVN_Southbound` database. For example:

```
$ /usr/share/openvswitch/scripts/ovn-ctl start_northd
```

With Open vSwitch version of 2.7 or greater, you need to run the following additional commands (Please read the manpages of `ovn-nb` for more control on the types of connection allowed.)

```
$ ovn-nbctl set-connection tcp:6641
$ ovn-sbctl set-connection tcp:6642
```

### 2. One time setup

On each host, where you plan to spawn your containers, you will need to run the below command once. You may need to run it again if your OVS database gets cleared. It is harmless to run it again in any case:

```
$ ovs-vsctl set Open_vSwitch . \
 external_ids:ovn-remote="tcp:$CENTRAL_IP:6642" \
 external_ids:ovn-nb="tcp:$CENTRAL_IP:6641" \
 external_ids:ovn-encap-ip=$LOCAL_IP \
 external_ids:ovn-encap-type="$ENCAP_TYPE"
```

where:

**\$LOCAL\_IP** is the IP address via which other hosts can reach this host. This acts as your local tunnel endpoint.

**\$ENCAP\_TYPE** is the type of tunnel that you would like to use for overlay networking. The options are `geneve` or `stt`. Your kernel must have support for your chosen `$ENCAP_TYPE`. Both `geneve` and `stt` are part of the Open vSwitch kernel module that is compiled from this repo. If you use the Open vSwitch kernel module from upstream Linux, you will need a minimum kernel version of 3.18 for `geneve`. There is no `stt` support in upstream Linux. You can verify whether you have the support in your kernel as follows:

```
$ lsmod | grep $ENCAP_TYPE
```

In addition, each Open vSwitch instance in an OVN deployment needs a unique, persistent identifier, called the `system-id`. If you install OVS from distribution packaging for Open vSwitch (e.g. `.deb` or `.rpm` packages), or if you use the `ovs-ctl` utility included with Open vSwitch, it automatically configures a `system-id`. If you start Open vSwitch manually, you should set one up yourself. For example:

```
$ id_file=/etc/openvswitch/system-id.conf
$ test -e $id_file || uuidgen > $id_file
$ ovs-vsctl set Open_vSwitch . external_ids:system-id=$(cat $id_file)
```

### 3. Start the `ovn-controller`.

You need to run the below command on every boot:

```
$ /usr/share/openvswitch/scripts/ovn-ctl start_controller
```

### 4. Start the Open vSwitch network driver.

By default Docker uses Linux bridge for networking. But it has support for external drivers. To use Open vSwitch instead of the Linux bridge, you will need to start the Open vSwitch driver.

The Open vSwitch driver uses the Python's flask module to listen to Docker's networking api calls. So, if your host does not have Python's flask module, install it:

```
$ sudo pip install Flask
```

Start the Open vSwitch driver on every host where you plan to create your containers. Refer to the note on `$OVS_PYTHON_LIBS_PATH` that is used below at the end of this document:

```
$ PYTHONPATH=$OVS_PYTHON_LIBS_PATH ovn-docker-overlay-driver --detach
```

---

**Note:** The `$OVS_PYTHON_LIBS_PATH` variable should point to the directory where Open vSwitch Python modules are installed. If you installed Open vSwitch Python modules via the Debian package of `python-openvswitch` or via `pip` by running `pip install ovs`, you do not need to specify the `PATH`. If you installed it by following the instructions in [Open vSwitch on Linux, FreeBSD and NetBSD](#), then you should specify the `PATH`. In this case, the `PATH` depends on the options passed to `./configure`. It is usually either `/usr/share/openvswitch/python` or `/usr/local/share/openvswitch/python`

---

Docker has inbuilt primitives that closely match OVN's logical switches and logical port concepts. Consult Docker's documentation for all the possible commands. Here are some examples.

## Create a logical switch

To create a logical switch with name 'foo', on subnet '192.168.1.0/24', run:

```
$ NID=`docker network create -d openvswitch --subnet=192.168.1.0/24 foo`
```

### List all logical switches

```
$ docker network ls
```

You can also look at this logical switch in OVN's northbound database by running the following command:

```
$ ovn-nbctl --db=tcp:$CENTRAL_IP:6640 ls-list
```

### Delete a logical switch

```
$ docker network rm bar
```

### Create a logical port

Docker creates your logical port and attaches it to the logical network in a single step. For example, to attach a logical port to network `foo` inside container `busybox`, run:

```
$ docker run -itd --net=foo --name=busybox busybox
```

### List all logical ports

Docker does not currently have a CLI command to list all logical ports but you can look at them in the OVN database by running:

```
$ ovn-nbctl --db=tcp:$CENTRAL_IP:6640 lsp-list $NID
```

### Create and attach a logical port to a running container

```
$ docker network create -d openvswitch --subnet=192.168.2.0/24 bar
$ docker network connect bar busybox
```

### Detach and delete a logical port from a running container

You can delete your logical port and detach it from a running container by running:

```
$ docker network disconnect bar busybox
```

### The “underlay” mode

---

**Note:** This mode requires that you have a OpenStack setup pre-installed with OVN providing the underlay networking.

---

### 1. One time setup

A OpenStack tenant creates a VM with a single network interface (or multiple) that belongs to management logical networks. The tenant needs to fetch the port-id associated with the interface via which he plans to send the container traffic inside the spawned VM. This can be obtained by running the below command to fetch the 'id' associated with the VM:

```
$ nova list
```

and then by running:

```
$ neutron port-list --device_id=$id
```

Inside the VM, download the OpenStack RC file that contains the tenant information (henceforth referred to as `openrc.sh`). Edit the file and add the previously obtained port-id information to the file by appending the following line:

```
$ export OS_VIF_ID=$port_id
```

After this edit, the file will look something like:

```
#!/bin/bash
export OS_AUTH_URL=http://10.33.75.122:5000/v2.0
export OS_TENANT_ID=fab106b215d943c3bad519492278443d
export OS_TENANT_NAME="demo"
export OS_USERNAME="demo"
export OS_VIF_ID=e798c371-85f4-4f2d-ad65-d09dd1d3c1c9
```

### 2. Create the Open vSwitch bridge

If your VM has one ethernet interface (e.g.: 'eth0'), you will need to add that device as a port to an Open vSwitch bridge 'breth0' and move its IP address and route related information to that bridge. (If it has multiple network interfaces, you will need to create and attach an Open vSwitch bridge for the interface via which you plan to send your container traffic.)

If you use DHCP to obtain an IP address, then you should kill the DHCP client that was listening on the physical Ethernet interface (e.g. eth0) and start one listening on the Open vSwitch bridge (e.g. breth0).

Depending on your VM, you can make the above step persistent across reboots. For example, if your VM is Debian/Ubuntu-based, read *openvswitch-switch.README.Debian* found in *debian* folder. If your VM is RHEL-based, refer to *RHEL 5.6, 6.x Packaging for Open vSwitch*.

### 3. Start the Open vSwitch network driver

The Open vSwitch driver uses the Python's flask module to listen to Docker's networking api calls. The driver also uses OpenStack's `python-neutronclient` libraries. If your host does not have Python's flask module or `python-neutronclient` you must install them. For example:

```
$ pip install python-neutronclient
$ pip install Flask
```

Once installed, source the `openrc` file:

```
$. ./openrc.sh
```

Start the network driver and provide your OpenStack tenant password when prompted:

```
$ PYTHONPATH=$OVS_PYTHON_LIBS_PATH ovn-docker-underlay-driver \
--bridge breth0 --detach
```



From here-on you can use the same Docker commands as described in [docker-overlay](#).

Refer to the ovs-architecture man pages (`man ovn-architecture`) to understand OVN's architecture in detail.

## 5.2.2 Integration of Containers with OVN and OpenStack

Isolation between containers is weaker than isolation between VMs, so some environments deploy containers for different tenants in separate VMs as an additional security measure. This document describes creation of containers inside VMs and how they can be made part of the logical networks securely. The created logical network can include VMs, containers and physical machines as endpoints. To better understand the proposed integration of containers with OVN and OpenStack, this document describes the end to end workflow with an example.

- A OpenStack tenant creates a VM (say VM-A) with a single network interface that belongs to a management logical network. The VM is meant to host containers. OpenStack Nova chooses the hypervisor on which VM-A is created.
- A Neutron port may have been created in advance and passed in to Nova with the request to create a new VM. If not, Nova will issue a request to Neutron to create a new port. The ID of the logical port from Neutron will also be used as the vif-id for the virtual network interface (VIF) of VM-A.
- When VM-A is created on a hypervisor, its VIF gets added to the Open vSwitch integration bridge. This creates a row in the Interface table of the Open\_vSwitch database. As explained in the [integration guide](#), the vif-id associated with the VM network interface gets added in the `external_ids:iface-id` column of the newly created row in the Interface table.
- Since VM-A belongs to a logical network, it gets an IP address. This IP address is used to spawn containers (either manually or through container orchestration systems) inside that VM and to monitor the health of the created containers.
- The vif-id associated with the VM's network interface can be obtained by making a call to Neutron using tenant credentials.
- This flow assumes a component called a "container network plugin". If you take Docker as an example for containers, you could envision the plugin to be either a wrapper around Docker or a feature of Docker itself that understands how to perform part of this workflow to get a container connected to a logical network managed by Neutron. The rest of the flow refers to this logical component that does not yet exist as the "container network plugin".
- All the calls to Neutron will need tenant credentials. These calls can either be made from inside the tenant VM as part of a container network plugin or from outside the tenant VM (if the tenant is not comfortable using temporary Keystone tokens from inside the tenant VMs). For simplicity, this document explains the work flow using the former method.
- The container hosting VM will need Open vSwitch installed in it. The only work for Open vSwitch inside the VM is to tag network traffic coming from containers.
- When a container needs to be created inside the VM with a container network interface that is expected to be attached to a particular logical switch, the network plugin in that VM chooses any unused VLAN (This VLAN tag only needs to be unique inside that VM. This limits the number of container interfaces to 4096 inside a single VM). This VLAN tag is stripped out in the hypervisor by OVN and is only useful as a context (or metadata) for OVN.
- The container network plugin then makes a call to Neutron to create a logical port. In addition to all the inputs that a call to create a port in Neutron that are currently needed, it sends the vif-id and the VLAN tag as inputs.
- Neutron in turn will verify that the vif-id belongs to the tenant in question and then uses the OVN specific plugin to create a new row in the Logical\_Switch\_Port table of the OVN Northbound Database. Neutron responds back with an IP address and MAC address for that network interface. So Neutron becomes the IPAM system and provides unique IP and MAC addresses across VMs and containers in the same logical network.

- The Neutron API call above to create a logical port for the container could add a relatively significant amount of time to container creation. However, an optimization is possible here. Logical ports could be created in advance and reused by the container system doing container orchestration. Additional Neutron API calls would only be needed if the port needs to be attached to a different logical network.
- When a container is eventually deleted, the network plugin in that VM may make a call to Neutron to delete that port. Neutron in turn will delete the entry in the `Logical_Switch_Port` table of the OVN Northbound Database.

As an example, consider Docker containers. Since Docker currently does not have a network plugin feature, this example uses a hypothetical wrapper around Docker to make calls to Neutron.

- Create a Logical switch:

```
$ ovn-docker --cred=cca86bd13a564ac2a63ddf14bf45d37f create network LS1
```

The above command will make a call to Neutron with the credentials to create a logical switch. The above is optional if the logical switch has already been created from outside the VM.

- List networks available to the tenant:

```
$ ovn-docker --cred=cca86bd13a564ac2a63ddf14bf45d37f list networks
```

- Create a container and attach a interface to the previously created switch as a logical port:

```
$ ovn-docker --cred=cca86bd13a564ac2a63ddf14bf45d37f --vif-id=$VIF_ID \
 --network=LS1 run -d --net=none ubuntu:14.04 /bin/sh -c \
 "while true; do echo hello world; sleep 1; done"
```

The above command will make a call to Neutron with all the inputs it currently needs to create a logical port. In addition, it passes the `$VIF_ID` and a unused VLAN. Neutron will add that information in OVN and return back a MAC address and IP address for that interface. `ovn-docker` will then create a veth pair, insert one end inside the container as `eth0` and the other end as a port of a local OVS bridge as an access port of the chosen VLAN.

## 5.2.3 Open Virtual Network With firewalld

`firewalld` is a service that allows for easy administration of firewalls. OVN ships with a set of service files that can be used with `firewalld` to allow for remote connections to the northbound and southbound databases.

This guide will describe how you can use these files with your existing `firewalld` setup. Setup and administration of `firewalld` is outside the scope of this document.

### Installation

If you have installed OVN from an RPM, then the service files for `firewalld` will automatically be installed in `/usr/lib/firewalld/services`. Installation from RPM includes installation from the `yum` or `dnf` package managers.

If you have installed OVN from source, then from the top level source directory, issue the following commands to copy the `firewalld` service files:

```
$ cp rhel/usr_lib_firewalld_services_ovn-central-firewall-service.xml \
/etc/firewalld/services/
$ cp rhel/usr_lib_firewalld_services_ovn-host-firewall-service.xml \
/etc/firewalld/services/
```

## Activation

Assuming you are already running `firewalld`, you can issue the following commands to enable the OVN services.

On the central server (the one running `ovn-northd`), issue the following:

```
$ firewall-cmd --zone=public --add-service=ovn-central-firewall-service
```

This will open TCP ports 6641 and 6642, allowing for remote connections to the northbound and southbound databases.

On the OVN hosts (the ones running `ovn-controller`), issue the following:

```
$ firewall-cmd --zone=public --add-service=ovn-host-firewall-service
```

This will open UDP port 6081, allowing for geneve traffic to flow between the controllers.

## Variations

When installing the XML service files, you have the choice of copying them to `/etc/firewalld/services` or `/usr/lib/firewalld/services`. The former is recommended since the latter can be overwritten if `firewalld` is upgraded.

The above commands assumed your underlay network interfaces are in the “public” `firewalld` zone. If your underlay network interfaces are in a separate zone, then adjust the above commands accordingly.

The `--permanent` option may be passed to the above `firewall-cmd` invocations in order for the services to be permanently added to the `firewalld` configuration. This way it is not necessary to re-issue the commands each time the `firewalld` service restarts.

The `ovn-host-firewall-service` only opens port 6081. This is because the default protocol for OVN tunnels is geneve. If you are using a different encapsulation protocol, you will need to modify the XML service file to open the appropriate port(s). For VXLAN, open port 4789. For STT, open port 7471.

## Recommendations

The `firewalld` service files included with the OVS repo are meant as a convenience for `firewalld` users. All that the service files do is to open the common ports used by OVN. No additional security is provided. To ensure a more secure environment, it is a good idea to do the following

- Use tools such as `iptables` or `nftables` to restrict access to known hosts.
- Use SSL for all remote connections to OVN databases.
- Use role-based access control for connections to the OVN southbound database.



## 6.1 Man Pages

The following man pages are written in rST and converted to roff at compile time:

### 6.1.1 ovs-sim

#### Synopsis

`ovs-sim` [*option*...]... [*script*]...

#### Description

`ovs-sim` provides a convenient environment for running one or more Open vSwitch instances and related software in a sandboxed simulation environment.

To use `ovs-sim`, first build Open vSwitch, then invoke it directly from the build directory, e.g.:

```
git clone https://github.com/openvswitch/ovs.git
cd ovs
./configure
make
utilities/ovs-sim
```

When invoked in the most ordinary way as shown above, `ovs-sim` does the following:

1. Creates a directory `sandbox` as a subdirectory of the current directory (first destroying such a directory if it already exists) and makes it the current directory.
2. Installs all of the Open vSwitch manpages into a `man` subdirectory of `sandbox` and adjusts the `MANPATH` environment variable so that `man` and other manpage viewers can find them.

3. Creates a simulated Open vSwitch named `main` and sets it up as the default target for OVS commands, as if the following `ovs-sim` commands had been run:

```
sim_add main
as main
```

See [Commands](#), below, for an explanation.

4. Runs any scripts specified on the command line (see [Options](#), below). The scripts can use arbitrary Bash syntax, plus the additional commands described under [Commands](#), below.
5. If no scripts were specified, or if `-i` or `--interactive` was specified, invokes an interactive Bash subshell. The user can use arbitrary Bash commands, plus the additional commands described under [Commands](#), below.

`ovs-sim` and the sandbox environment that it creates does not require superuser or other special privileges. Generally, it should not be run with such privileges.

## Options

**script** Runs *script*, which should be a Bash script, within a subshell after initializing. If multiple script arguments are given, then they are run in the order given. If any script exits with a nonzero exit code, then `ovs-sim` exits immediately with the same exit code.

**-i or --interactive** By default, if any script is specified, `ovs-sim` exits as soon as the scripts finish executing. With this option, or if no scripts are specified, `ovs-sim` instead starts an interactive Bash session.

## Commands

Scripts and interactive usage may use the following commands implemented by `ovs-sim`. They are implemented as Bash shell functions exported to subshells.

### Basic Commands

These are the basic commands for working with sandboxed Open vSwitch instances.

**sim\_add sandbox** Starts a new simulated Open vSwitch instance named *sandbox*. Files related to the instance, such as logs, databases, sockets, and pidfiles, are created in a subdirectory also named *sandbox*. Afterward, the `as` command (see below) can be used to run Open vSwitch utilities in the context of the new sandbox.

The new sandbox starts out without any bridges. Use `ovs-vsctl` in the context of the new sandbox to create a bridge, e.g.:

```
sim_add hv0 # Create sandbox hv0.
as hv0 # Set hv0 as default sandbox.
ovs-vsctl add-br br0 # Add bridge br0 inside hv0.
```

The Open vSwitch instances that `sim_add` creates enable dummy devices. This means that bridges and interfaces can be created with type `dummy` to indicate that they should be totally simulated, without any reference to system entities. In fact, `ovs-sim` also configures Open vSwitch so that the default system type of bridges and interfaces are replaced by dummy devices. Other types of devices, however, retain their usual functions, which means that, e.g., vxlan tunnels still act as tunnels (refer to the documentation).

**as sandbox** Sets *sandbox* as the default simulation target for Open vSwitch commands (e.g. `ovs-vsctl`, `ovs-ofctl`, `ovs-appctl`).

This command updates the beginning of the shell prompt to indicate the new default target.

**as sandbox command arg...** Runs the given command with *sandbox* as the simulation target, e.g. `as hv0 ovs-vsctl add-br br0` runs `ovs-vsctl add-br br0` within sandbox `hv0`. The default target is unchanged.

## Interconnection Network Commands

When multiple sandboxed Open vSwitch instances exist, one will inevitably want to connect them together. These commands allow for that. Conceptually, an interconnection network is a switch that `ovs-sim` makes it easy to plug into other switches in other sandboxed Open vSwitch instances. Interconnection networks are implemented as bridges in the `main` switch that `ovs-sim` creates by default, so to use interconnection networks please avoid working with `main` directly.

**net\_add network** Creates a new interconnection network named *network*.

**net\_attach network bridge** Adds a new port to *bridge* in the default sandbox (as set with `as`) and plugs it into interconnection network *network*, which must already have been created by a previous invocation of `net_add`. The default sandbox must not be `main`.

## OVN Commands

These commands interact with OVN, the Open Virtual Network.

**ovn\_start [options]** Creates and initializes the central OVN databases (both `ovn-sb(5)` and `ovn-nb(5)`) and starts an instance of `ovsdb-server` for each one. Also starts an instance of `ovn-northd`.

The following options are available:

**--nbdb-model model** Uses the given database model for the northbound database. The *model* may be `standalone` (the default), `backup`, or `clustered`.

**--nbdb-servers n** For a clustered northbound database, the number of servers in the cluster. The default is 3.

**--sbdb-model model** Uses the given database model for the southbound database. The *model* may be `standalone` (the default), `backup`, or `clustered`.

**--sbdb-servers n** For a clustered southbound database, the number of servers in the cluster. The default is 3.

**ovn\_attach network bridge ip [masklen]** First, this command attaches *bridge* to interconnection network *network*, just like `net_attach network bridge`. Second, it configures (simulated) IP address *ip* (with network mask length *masklen*, which defaults to 24) on *bridge*. Finally, it configures the Open vSwitch database to work with OVN and starts `ovn-controller`.

## Examples

The following creates a pair of Open vSwitch instances `hv0` and `hv1`, adds a port named `vif0` or `vif1`, respectively, to each one, and then connects the two through an interconnection network `n1`:

```
net_add n1
for i in 0 1; do
 sim_add hv$i
 as hv$i ovs-vsctl add-br br0 -- add-port br0 vif$i
 as hv$i net_attach n1 br0
done
```

Here's an extended version that also starts OVN:

```
ovn_start
ovn-nbctl ls-add lsw0
net_add n1
for i in 0 1; do
 sim_add hv$i
 as hv$i
 ovs-vsctl add-br br-phys
 ovn_attach n1 br-phys 192.168.0.`expr $i + 1`
 ovs-vsctl add-port br-int vif$i -- set Interface vif$i external-ids:iface-id=lp$i
 ovn-nbctl lsp-add lsw0 lp$i
 ovn-nbctl lsp-set-addresses lp$i f0:00:00:00:00:0$i
done
```

Here's a primitive OVN “scale test” (adjust the scale by changing *n* in the first line:

```
n=200; export n
ovn_start --sbdb-model=clustered
net_add n1
ovn-nbctl ls-add br0
for i in `seq $n`; do
 (sim_add hv$i
 as hv$i
 ovs-vsctl add-br br-phys
 y=$(expr $i / 256)
 x=$(expr $i % 256)
 ovn_attach n1 br-phys 192.168.$y.$x
 ovs-vsctl add-port br-int vif$i -- set Interface vif$i external-ids:iface-id=lp
 ↪$i) &
 case $i in
 *50|*00) echo $i; wait ;;
 esac
done
wait
for i in `seq $n`; do
 yy=$(printf %02x $(expr $i / 256))
 xx=$(printf %02x $(expr $i % 256))
 ovn-nbctl lsp-add br0 lp$i
 ovn-nbctl lsp-set-addresses lp$i f0:00:00:00:$yy:$xx
done
```

When the scale test has finished initializing, you can watch the logical ports come up with a command like this:

```
watch 'for i in `seq $n`; do if test `ovn-nbctl lsp-get-up lp$i` != up; then echo $i; ↪
 ↪fi; done'
```

## 6.1.2 ovs-test

### Synopsis

**ovs-test** -s *port*

**ovs-test** -c *server1 server2* [-b *targetbandwidth*] [-i *testinterval*] [-d] [-l *vlantag*] [-t *tunnelmodes*]



## Description

The **ovs-test** program may be used to check for problems sending 802.1Q or GRE traffic that Open vSwitch may uncover. These problems, for example, can occur when Open vSwitch is used to send 802.1Q traffic through physical interfaces running certain drivers of certain Linux kernel versions. To run a test, configure IP addresses on *server1* and *server2* for interfaces you intended to test. These interfaces could also be already configured OVS bridges that have a physical interface attached to them. Then, on one of the nodes, run **ovs-test** in server mode and on the other node run it in client mode. The client will connect to **ovs-test** server and schedule tests between both of them. The **ovs-test** client will perform UDP and TCP tests.

UDP tests can report packet loss and achieved bandwidth for various datagram sizes. By default target bandwidth for UDP tests is 1Mbit/s.

TCP tests report only achieved bandwidth, because kernel TCP stack takes care of flow control and packet loss. TCP tests are essential to detect potential TSO related issues.

To determine whether Open vSwitch is encountering any problems, the user must compare packet loss and achieved bandwidth in a setup where traffic is being directly sent and in one where it is not. If in the 802.1Q or L3 tunneled tests both **ovs-test** processes are unable to communicate or the achieved bandwidth is much lower compared to direct setup, then, most likely, Open vSwitch has encountered a pre-existing kernel or driver bug.

Some examples of the types of problems that may be encountered are:

- When NICs use VLAN stripping on receive they must pass a pointer to a *vlan\_group* when reporting the stripped tag to the networking core. If no *vlan\_group* is in use then some drivers just drop the extracted tag. Drivers are supposed to only enable stripping if a *vlan\_group* is registered but not all of them do that.
- On receive, some drivers handle priority tagged packets specially and don't pass the tag onto the network stack at all, so Open vSwitch never has a chance to see it.
- Some drivers size their receive buffers based on whether a *vlan\_group* is enabled, meaning that a maximum size packet with a VLAN tag will not fit if no *vlan\_group* is configured.
- On transmit, some drivers expect that VLAN acceleration will be used if it is available, which can only be done if a *vlan\_group* is configured. In these cases, the driver may fail to parse the packet and correctly setup checksum offloading or TSO.

**Client Mode** An **ovs-test** client will connect to two **ovs-test** servers and will ask them to exchange test traffic. It is also possible to spawn an **ovs-test** server automatically from the client.

**Server Mode** To conduct tests, two **ovs-test** servers must be running on two different hosts where the client can connect. The actual test traffic is exchanged only between both **ovs-test** servers. It is recommended that both servers have their IP addresses in the same subnet, otherwise one would have to make sure that routing is set up correctly.

## Options

**-s** <port>, **--server** <port>

Run in server mode and wait for the client to establish XML RPC Control Connection on this TCP port. It is recommended to have *ethtool(8)* installed on the server so that it could retrieve information about the NIC driver.

**-c** <server1> <server2>, **--client** <server1> <server2>

Run in client mode and schedule tests between *server1* and *server2*, where each server must be given in the following format:

```
OuterIP[:OuterPort], InnerIP[/Mask][:InnerPort].
```

The *OuterIP* must be already assigned to the physical interface which is going to be tested. This is the IP address where client will try to establish XML RPC connection. If *OuterIP* is 127.0.0.1 then client will automatically

spawn a local instance of **ovs-test** server. OuterPort is TCP port where server is listening for incoming XML/RPC control connections to schedule tests (by default it is 15531). The **ovs-test** will automatically assign *InnerIP[/Mask]* to the interfaces that will be created on the fly for testing purposes. It is important that *InnerIP[/Mask]* does not interfere with already existing IP addresses on both **ovs-test** servers and client. InnerPort is port which will be used by server to listen for test traffic that will be encapsulated (by default it is 15532).

- b** <targetbandwidth>, **--bandwidth** <targetbandwidth>  
Target bandwidth for UDP tests. The targetbandwidth must be given in bits per second. It is possible to use postfix *M* or *K* to alter the target bandwidth magnitude.
- i** <testinterval>, **--interval** <testinterval>  
How long each test should run. By default 5 seconds.
- h**, **--help**  
Prints a brief help message to the console.
- V**, **--version**  
Prints version information to the console.

The following test modes are supported by **ovs-test**. It is possible to combine multiple of them in a single **ovs-test** invocation.

- d**, **--direct**  
Perform direct tests between both OuterIP addresses. These tests could be used as a reference to compare 802.1Q or L3 tunneling test results.
- l** <vlantag>, **--vlan-tag** <vlantag>  
Perform 802.1Q tests between both servers. These tests will create a temporary OVS bridge, if necessary, and attach a VLAN tagged port to it for testing purposes.
- t** <tunnelmodes>, **--tunnel-modes** <tunnelmodes>  
Perform L3 tunneling tests. The given argument is a comma separated string that specifies all the L3 tunnel modes that should be tested (e.g. gre). The L3 tunnels are terminated on interface that has the OuterIP address assigned.

## Examples

On host 1.2.3.4 start **ovs-test** in server mode:

```
ovs-test -s 15531
```

On host 1.2.3.5 start **ovs-test** in client mode and do direct, VLAN and GRE tests between both nodes:

```
ovs-test -c 127.0.0.1,1.1.1.1/30 1.2.3.4,1.1.1.2/30 -d -l 123 -t gre
```

## See Also

*ovs-vswitchd(8)*, *ovs-ofctl(8)*, *ovs-vsctl(8)*, **ovs-vlan-test**, *ethtool(8)*, *uname(1)*

### 6.1.3 ovs-vlan-test

#### Synopsis

**ovs-vlan-test** [-s | -server] control\_ip vlan\_ip

## Description

The **ovs-vlan-test** utility has some limitations, for example, it does not use TCP in its tests. Also it does not take into account MTU to detect potential edge cases. To overcome those limitations a new tool was developed - **ovs-test**. **ovs-test** is currently supported only on Debian so, if possible, try to use that on instead of **ovs-vlan-test**.

The **ovs-vlan-test** program may be used to check for problems sending 802.1Q traffic which may occur when running Open vSwitch. These problems can occur when Open vSwitch is used to send 802.1Q traffic through physical interfaces running certain drivers of certain Linux kernel versions. To run a test, configure Open vSwitch to tag traffic originating from *vlan\_ip* and forward it out the target interface. Then run the **ovs-vlan-test** in client mode connecting to an **ovs-vlan-test** server. **ovs-vlan-test** will display “OK” if it did not detect problems.

Some examples of the types of problems that may be encountered are:

- When NICs use VLAN stripping on receive they must pass a pointer to a *vlan\_group* when reporting the stripped tag to the networking core. If no *vlan\_group* is in use then some drivers just drop the extracted tag. Drivers are supposed to only enable stripping if a *vlan\_group* is registered but not all of them do that.
- On receive, some drivers handle priority tagged packets specially and don’t pass the tag onto the network stack at all, so Open vSwitch never has a chance to see it.
- Some drivers size their receive buffers based on whether a *vlan\_group* is enabled, meaning that a maximum size packet with a VLAN tag will not fit if no *vlan\_group* is configured.
- On transmit, some drivers expect that VLAN acceleration will be used if it is available, which can only be done if a *vlan\_group* is configured. In these cases, the driver may fail to parse the packet and correctly setup checksum offloading or TSO.

**Client Mode** An **ovs-vlan-test** client may be run on a host to check for VLAN connectivity problems. The client must be able to establish HTTP connections with an **ovs-vlan-test** server located at the specified *control\_ip* address. UDP traffic sourced at *vlan\_ip* should be tagged and directed out the interface whose connectivity is being tested.

**Server Mode** To conduct tests, an **ovs-vlan-test** server must be running on a host known not to have VLAN connectivity problems. The server must have a *control\_ip* on a non-VLAN network which clients can establish connectivity with. It must also have a *vlan\_ip* address on a VLAN network which clients will use to test their VLAN connectivity. Multiple clients may test against a single **ovs-vlan-test** server concurrently.

## Options

- s, --server**  
Run in server mode.
- h, --help**  
Prints a brief help message to the console.
- V, --version**  
Prints version information to the console.

## Examples

Display the Linux kernel version and driver of *eth1*:

```
uname -r
ethtool -i eth1
```

Set up a bridge which forwards traffic originating from *1.2.3.4* out *eth1* with VLAN tag 10:

```
ovs-vsctl -- add-br vlan-br \
-- add-port vlan-br eth1 \
-- add-port vlan-br vlan-br-tag tag=10 \
-- set Interface vlan-br-tag type=internal
ip addr add 1.2.3.4/8 dev vlan-br-tag
ip link set vlan-br-tag up
```

Run an **ovs-vlan-test** server listening for client control traffic on *172.16.0.142* port *8080* and VLAN traffic on the default port of *1.2.3.3*:

```
ovs-vlan-test -s 172.16.0.142:8080 1.2.3.3
```

Run an **ovs-vlan-test** client with a control server located at *172.16.0.142* port *8080* and a local VLAN IP of *1.2.3.4*:

```
ovs-vlan-test 172.16.0.142:8080 1.2.3.4
```

### See Also

*ovs-vswitchd(8)*, *ovs-ofctl(8)*, *ovs-vsctl(8)*, **ovs-test**, *ethtool(8)*, *uname(1)*

## 6.1.4 ovssdb-server

### Description

ovssdb-server implements the Open vSwitch Database (OVSSDB) protocol specified in RFC 7047. This document provides clarifications for how ovssdb-server implements the protocol and describes the extensions that it provides beyond RFC 7047. Numbers in section headings refer to corresponding sections in RFC 7047.

### 3.1 JSON Usage

RFC 4627 says that names within a JSON object should be unique. The Open vSwitch JSON parser discards all but the last value for a name that is specified more than once.

The definition of <error> allows for implementation extensions. Currently ovssdb-server uses the following additional error strings (which might change in later releases):

**syntax error or unknown column** The request could not be parsed as an OVSSDB request. An additional `syntax` member, whose value is a string that contains JSON, may narrow down the particular syntax that could not be parsed.

**internal error** The request triggered a bug in ovssdb-server.

**ovssdb error** A map or set contains a duplicate key.

**permission error** The request was denied by the role-based access control extension, introduced in version 2.8.

### 3.2 Schema Format

RFC 7047 requires the `version` field in <database-schema>. Current versions of ovssdb-server allow it to be omitted (future versions are likely to require it).

RFC 7047 allows columns that contain weak references to be immutable. This raises the issue of the behavior of the weak reference when the rows that it references are deleted. Since version 2.6, `ovsdb-server` forces columns that contain weak references to be mutable.

Since version 2.8, the table name `RBAC_Role` is used internally by the role-based access control extension to `ovsdb-server` and should not be used for purposes other than defining mappings of role names to table access permissions. This table has one row per role name and the following columns:

**name** The role name.

**permissions** A map of table name to a reference to a row in a separate permission table.

The separate RBAC permission table has one row per access control configuration and the following columns:

**name** The name of the table to which the row applies.

**authorization** The set of column names and column:key pairs to be compared with the client ID in order to determine the authorization status of the requested operation.

**insert\_delete** A boolean value, true if authorized insertions and deletions are allowed, false if no insertions or deletions are allowed.

**update** The set of columns and column:key pairs for which authorized update and mutate operations should be permitted.

## 4 Wire Protocol

The original OVSDb specifications included the following reasons, omitted from RFC 7047, to operate JSON-RPC directly over a stream instead of over HTTP:

- JSON-RPC is a peer-to-peer protocol, but HTTP is a client-server protocol, which is a poor match. Thus, JSON-RPC over HTTP requires the client to periodically poll the server to receive server requests.
- HTTP is more complicated than stream connections and doesn't provide any corresponding advantage.
- The JSON-RPC specification for HTTP transport is incomplete.

### 4.1.3 Transact

Since version 2.8, role-based access controls can be applied to operations within a transaction that would modify the contents of the database (these operations include row insert, row delete, column update, and column mutate). Role-based access controls are applied when the database schema contains a table with the name `RBAC_Role` and the connection on which the transaction request was received has an associated role name (from the `role` column in the remote connection table). When role-based access controls are enabled, transactions that are otherwise well-formed may be rejected depending on the client's role, ID, and the contents of the `RBAC_Role` table and associated permissions table.

### 4.1.5 Monitor

For backward compatibility, `ovsdb-server` currently permits a single `<monitor-request>` to be used instead of an array; it is treated as a single-element array. Future versions of `ovsdb-server` might remove this compatibility feature.

Because the `<json-value>` parameter is used to match subsequent update notifications (see below) to the request, it must be unique among all active monitors. `ovsdb-server` rejects attempt to create two monitors with the same identifier.

### 4.1.7 Monitor Cancellation

When a database monitored by a session is removed, and database change awareness is enabled for the session (see Section 4.1.16), the database server spontaneously cancels all monitors (including conditional monitors described in Section 4.1.12) for the removed database. For each canceled monitor, it issues a notification in the following form:

```
"method": "monitor_canceled"
"params": [<json-value>]
"id": null
```

### 4.1.12 Monitor\_cond

A new monitor method added in Open vSwitch version 2.6. The `monitor_cond` request enables a client to replicate subsets of tables within an OVSDB database by requesting notifications of changes to rows matching one of the conditions specified in `where` by receiving the specified contents of these rows when table updates occur. `monitor_cond` also allows a more efficient update notifications by receiving `<table-updates2>` notifications (described below).

The `monitor` method described in Section 4.1.5 also applies to `monitor_cond`, with the following exceptions:

- RPC request method becomes `monitor_cond`.
- Reply result follows `<table-updates2>`, described in Section 4.1.14.
- Subsequent changes are sent to the client using the `update2` monitor notification, described in Section 4.1.14
- Update notifications are being sent only for rows matching [`<condition>*`].

The request object has the following members:

```
"method": "monitor_cond"
"params": [<db-name>, <json-value>, <monitor-cond-requests>]
"id": <nonnull-json-value>
```

The `<json-value>` parameter is used to match subsequent update notifications (see below) to this request. The `<monitor-cond-requests>` object maps the name of the table to an array of `<monitor-cond-request>`.

Each `<monitor-cond-request>` is an object with the following members:

```
"columns": [<column>*] optional
"where": [<condition>*] optional
"select": <monitor-select> optional
```

The `columns`, if present, define the columns within the table to be monitored that match conditions. If not present, all columns are monitored.

The `where`, if present, is a JSON array of `<condition>` and boolean values. If not present or condition is an empty array, implicit True will be considered and updates on all rows will be sent.

`<monitor-select>` is an object with the following members:

```
"initial": <boolean> optional
"insert": <boolean> optional
"delete": <boolean> optional
"modify": <boolean> optional
```

The contents of this object specify how the columns or table are to be monitored as explained in more detail below.

The response object has the following members:

```
"result": <table-updates2>
"error": null
"id": same "id" as request
```

The `<table-updates2>` object is described in detail in Section 4.1.14. It contains the contents of the tables for which initial rows are selected. If no tables initial contents are requested, then `result` is an empty object.

Subsequently, when changes to a specified table that match one of the conditions in `<monitor-cond-request>` are committed, the changes are automatically sent to the client using the `update2` monitor notification (see Section 4.1.14). This monitoring persists until the JSON-RPC session terminates or until the client sends a `monitor_cancel` JSON-RPC request.

Each `<monitor-cond-request>` specifies one or more conditions and the manner in which the rows that match the conditions are to be monitored. The circumstances in which an `update` notification is sent for a row within the table are determined by `<monitor-select>`:

- If `initial` is omitted or true, every row in the original table that matches one of the conditions is sent as part of the response to the `monitor_cond` request.
- If `insert` is omitted or true, update notifications are sent for rows newly inserted into the table that match conditions or for rows modified in the table so that their old version does not match the condition and new version does.
- If `delete` is omitted or true, update notifications are sent for rows deleted from the table that match conditions or for rows modified in the table so that their old version does match the conditions and new version does not.
- If `modify` is omitted or true, update notifications are sent whenever a row in the table that matches conditions in both old and new version is modified.

Both `monitor` and `monitor_cond` sessions can exist concurrently. However, `monitor` and `monitor_cond` shares the same `<json-value>` parameter space; it must be unique among all `monitor` and `monitor_cond` sessions.

#### 4.1.13 Monitor\_cond\_change

The `monitor_cond_change` request enables a client to change an existing `monitor_cond` replication of the database by specifying a new condition and columns for each replicated table. Currently changing the columns set is not supported.

The request object has the following members:

```
"method": "monitor_cond_change"
"params": [<json-value>, <json-value>, <monitor-cond-update-requests>]
"id": <nonnull-json-value>
```

The `<json-value>` parameter should have a value of an existing conditional monitoring session from this client. The second `<json-value>` in `params` array is the requested value for this session. This value is valid only after `monitor_cond_change` is committed. A user can use these values to distinguish between update messages before conditions update and after. The `<monitor-cond-update-requests>` object maps the name of the table to an array of `<monitor-cond-update-request>`. Monitored tables not included in `<monitor-cond-update-requests>` retain their current conditions.

Each `<monitor-cond-update-request>` is an object with the following members:

```
"columns": [<column>*] optional
"where": [<condition>*] optional
```

The `columns` specify a new array of columns to be monitored, although this feature is not yet supported.

The `where` specify a new array of conditions to be applied to this monitoring session.

The response object has the following members:

```
"result": null
"error": null
"id": same "id" as request
```

Subsequent `<table-updates2>` notifications are described in detail in Section 4.1.14 in the RFC. If insert contents are requested by original `monitor_cond` request, `<table-updates2>` will contain rows that match the new condition and do not match the old condition. If deleted contents are requested by origin monitor request, `<table-updates2>` will contain any matched rows by old condition and not matched by the new condition.

Changes according to the new conditions are automatically sent to the client using the `update2` monitor notification. An update, if any, as a result of a condition change, will be sent to the client before the reply to the `monitor_cond_change` request.

#### 4.1.14 Update2 notification

The `update2` notification is sent by the server to the client to report changes in tables that are being monitored following a `monitor_cond` request as described above. The notification has the following members:

```
"method": "update2"
"params": [<json-value>, <table-updates2>]
"id": null
```

The `<json-value>` in `params` is the same as the value passed as the `<json-value>` in `params` for the corresponding `monitor` request. `<table-updates2>` is an object that maps from a table name to a `<table-update2>`. A `<table-update2>` is an object that maps from row's UUID to a `<row-update2>` object. A `<row-update2>` is an object with one of the following members:

```
"initial": <row> present for initial updates
"insert": <row> present for insert updates
"delete": <row> present for delete updates
"modify": <row> present for modify updates
```

The format of `<row>` is described in Section 5.1.

`<row>` is always a null object for a delete update. In `initial` and `insert` updates, `<row>` omits columns whose values equal the default value of the column type.

For a `modify` update, `<row>` contains only the columns that are modified. `<row>` stores the difference between the old and new value for those columns, as described below.

For columns with single value, the difference is the value of the new column.

The difference between two sets are all elements that only belong to one of the sets.

The difference between two maps are all key-value pairs whose keys appears in only one of the maps, plus the key-value pairs whose keys appear in both maps but with different values. For the latter elements, `<row>` includes the value from the new column.

Initial views of rows are not presented in `update2` notifications, but in the response object to the `monitor_cond` request. The formatting of the `<table-updates2>` object, however, is the same in either case.



#### 4.1.15 Get Server ID

A new RPC method added in Open vSwitch version 2.7. The request contains the following members:

```
"method": "get_server_id"
"params": null
"id": <nonnull-json-value>
```

The response object contains the following members:

```
"result": "<server_id>"
"error": null
"id": same "id" as request
```

<server\_id> is JSON string that contains a UUID that uniquely identifies the running OVSDb server process. A fresh UUID is generated when the process restarts.

#### 4.1.16 Database Change Awareness

RFC 7047 does not provide a way for a client to find out about some kinds of configuration changes, such as about databases added or removed while a client is connected to the server, or databases changing between read/write and read-only due to a transition between active and backup roles. Traditionally, `ovsdb-server` disconnects all of its clients when this happens, because this prompts a well-written client to reassess what is available from the server when it reconnects.

OVS 2.9 provides a way for clients to keep track of these kinds of changes, by monitoring the `Database` table in the `_Server` database introduced in this release (see `ovsdb-server(5)` for details). By itself, this does not suppress `ovsdb-server` disconnection behavior, because a client might monitor this database without understanding its special semantics. Instead, `ovsdb-server` provides a special request:

```
"method": "set_db_change_aware"
"params": [<boolean>]
"id": <nonnull-json-value>
```

If the boolean in the request is true, it suppresses the connection-closing behavior for the current connection, and false restores the default behavior. The reply is always the same:

```
"result": {}
"error": null
"id": same "id" as request
```

#### 4.1.17 Schema Conversion

Open vSwitch 2.9 adds a new JSON-RPC request to convert an online database from one schema to another. The request contains the following members:

```
"method": "convert"
"params": [<db-name>, <database-schema>]
"id": <nonnull-json-value>
```

Upon receipt, the server converts database <db-name> to schema <database-schema>. The schema's name must be <db-name>. The conversion is atomic, consistent, isolated, and durable. The data in the database must be valid when interpreted under <database-schema>, with only one exception: data for tables and columns that do not exist in the

new schema are ignored. Columns that exist in <database-schema> but not in the database are set to their default values. All of the new schema's constraints apply in full.

If the conversion is successful, the server notifies clients that use the `set_db_change_aware` RPC introduced in Open vSwitch 2.9 and cancels their outstanding transactions and monitors. The server disconnects other clients, enabling them to notice the change when they reconnect. The server sends the following reply:

```
"result": {}
"error": null
"id": same "id" as request
```

If the conversion fails, then the server sends an error reply in the following form:

```
"result": null
"error": [<error>]
"id": same "id" as request
```

## 5.1 Notation

For <condition>, RFC 7047 only allows the use of `!=`, `==`, `includes`, and `excludes` operators with set types. Open vSwitch 2.4 and later extend <condition> to allow the use of `<`, `<=`, `>=`, and `>` operators with a column with type “set of 0 or 1 integer” and an integer argument, and with “set of 0 or 1 real” and a real argument. These conditions evaluate to false when the column is empty, and otherwise as described in RFC 7047 for integer and real types.

<condition> is specified in Section 5.1 in the RFC with the following change: A condition can be either a 3-element JSON array as described in the RFC or a boolean value. In case of an empty array an implicit true boolean value will be considered.

### 5.2.6 Wait, 5.2.7 Commit, 5.2.9 Comment

RFC 7047 says that the `wait`, `commit`, and `comment` operations have no corresponding result object. This is not true. Instead, when such an operation is successful, it yields a result object with no members.

## 6.1.5 ovssdb

### Description

OVSSDB, the Open vSwitch Database, is a database system whose network protocol is specified by RFC 7047. The RFC does not specify an on-disk storage format. The OVSSDB implementation in Open vSwitch implements two storage formats: one for standalone (and active-backup) databases, and the other for clustered databases. This manpage documents both of these formats.

Most users do not need to be concerned with this specification. Instead, to manipulate OVSSDB files, refer to *ovssdb-tool(1)*. For an introduction to OVSSDB as a whole, read *ovssdb(7)*.

OVSSDB files explicitly record changes that are implied by the database schema. For example, the OVSSDB “garbage collection” feature means that when a client removes the last reference to a garbage-collected row, the database server automatically removes that row. The database file explicitly records the deletion of the garbage-collected row, so that the reader does not need to infer it.

OVSSDB files do not include the values of ephemeral columns.

Standalone and clustered database files share the common structure described here. They are text files encoded in UTF-8 with LF (U+000A) line ends, organized as append-only series of records. Each record consists of 2 lines of text.

The first line in each record has the format `OVSDB <magic> <length> <hash>`, where `<magic>` is `JSON` for standalone databases or `CLUSTER` for clustered databases, `<length>` is a positive decimal integer, and `<hash>` is a SHA-1 checksum expressed as 40 hexadecimal digits. Words in the first line must be separated by exactly one space.

The second line must be exactly *length* bytes long (including the LF) and its SHA-1 checksum (including the LF) must match *hash* exactly. The line's contents must be a valid JSON object as specified by RFC 4627. Strings in the JSON object must be valid UTF-8. To ensure that the second line is exactly one line of text, the OVSDB implementation expresses any LF characters within a JSON string as `\n`. For the same reason, and to save space, the OVSDB implementation does not “pretty print” the JSON object with spaces and LFs. (The OVSDB implementation tolerates LFs when reading an OVSDB database file, as long as *length* and *hash* are correct.)

## JSON Notation

We use notation from RFC 7047 here to describe the JSON data in records. In addition to the notation defined there, we add the following:

**<raw-uuid>** A 36-character JSON string that contains a UUID in the format described by RFC 4122, e.g. `"550e8400-e29b-41d4-a716-446655440000"`

## Standalone Format

The first record in a standalone database contains the JSON schema for the database, as specified in RFC 7047. Only this record is mandatory (a standalone file that contains only a schema represents an empty database).

The second and subsequent records in a standalone database are transaction records. Each record may have the following optional special members, which do not have any semantics but are often useful to administrators looking through a database log with `ovsdb-tool show-log`:

**"\_date": <integer>** The time at which the transaction was committed, as an integer number of milliseconds since the Unix epoch. Early versions of OVSDB counted seconds instead of milliseconds; these can be detected by noticing that their values are less than  $2^{32}$ .

OVSDB always writes a `_date` member.

**"\_comment": <string>** A JSON string that specifies the comment provided in a transaction `comment` operation. If a transaction has multiple `comment` operations, OVSDB concatenates them into a single `_comment` member, separated by a new-line.

OVSDB only writes a `_comment` member if it would be a nonempty string.

Each of these records also has one or more additional members, each of which maps from the name of a database table to a `<table-txn>`:

**<table-txn>** A JSON object that describes the effects of a transaction on a database table. Its names are `<raw-uuid>`s for rows in the table and its values are `<row-txn>`s.

**<row-txn>** Either `null`, which indicates that the transaction deleted this row, or a JSON object that describes how the transaction inserted or modified the row, whose names are the names of columns and whose values are `<value>`s that give the column's new value.

For new rows, the OVSDB implementation omits columns whose values have the default values for their types defined in RFC 7047 section 5.2.1; for modified rows, the OVSDB implementation omits columns whose values are unchanged.

## Clustered Format

The clustered format has the following additional notation:

**<uint64>** A JSON integer that represents a 64-bit unsigned integer. The OVS JSON implementation only supports integers in the range  $-2^{63}$  through  $2^{63}-1$ , so 64-bit unsigned integer values from  $2^{63}$  through  $2^{64}-1$  are expressed as negative numbers.

**<address>** A JSON string that represents a network address to support clustering, in the `<protocol>:<ip>:<port>` syntax described in `ovsdb-tool(1)`.

**<servers>** A JSON object whose names are `<raw-uuid>`s that identify servers and whose values are `<address>`es that specify those servers' addresses.

**<cluster-txn>** A JSON array with two elements:

1. The first element is either a `<database-schema>` or `null`. A `<database-schema>` element is always present in the first record of a clustered database to indicate the database's initial schema. If it is not `null` in a later record, it indicates a change of schema for the database.
2. The second element is either a transaction record in the format described under `Standalone Format` above, or `null`.

When a schema is present, the transaction record is relative to an empty database. That is, a schema change effectively resets the database to empty and the transaction record represents the full database contents. This allows readers to be ignorant of the full semantics of schema change.

The first record in a clustered database contains the following members, all of which are required:

**"server\_id": <raw-uuid>** The server's own UUID, which must be unique within the cluster.

**"local\_address": <address>** The address on which the server listens for connections from other servers in the cluster.

**name": <id>** The database schema name. It is only important when a server is in the process of joining a cluster: a server will only join a cluster if the name matches. (If the database schema name were unique, then we would not also need a cluster ID.)

**"cluster\_id": <raw-uuid>** The cluster's UUID. The all-zeros UUID is not a valid cluster ID.

**"prev\_term": <uint64> and "prev\_index": <uint64>** The Raft term and index just before the beginning of the log.

**"prev\_servers": <servers>** The set of one or more servers in the cluster at index "prev\_index" and term "prev\_term". It might not include this server, if it was not the initial server in the cluster.

**"prev\_data": <json-value> and "prev\_eid": <raw-uuid>** A snapshot of the data in the database at index "prev\_index" and term "prev\_term", and the entry ID for that data. The snapshot must contain a schema.

The second and subsequent records, if present, in a clustered database represent changes to the database, to the cluster state, or both. There are several types of these records. The most important types of records directly represent persistent state described in the Raft specification:

**Entry** A Raft log entry.

**Term** The start of a new term.

**Vote** The server's vote for a leader in the current term.

The following additional types of records aid debugging and troubleshooting, but they do not affect correctness.

**Leader** Identifies a newly elected leader for the current term.

**Commit Index** An update to the server's `commit_index`.

**Note** A human-readable description of some event.

The table below identifies the members that each type of record contains. "yes" indicates that a member is required, "?" that it is optional, blank that it is forbidden, and [1] that `data` and `eid` must be either both present or both absent.

member	Entry	Term	Vote	Leader	Commit Index	Note
comment	?	?	?	?	?	?
term	yes	yes	yes	yes		
index	yes					
servers	?					
data	[1]					
eid	[1]					
vote			yes			
leader				yes		
commit_index					yes	
note						yes

The members are:

**"comment":** **<string>** A human-readable string giving an administrator more information about the reason a record was emitted.

**"term":** **<uint64>** The term in which the activity occurred.

**"index":** **<uint64>** The index of a log entry.

**"servers":** **<servers>** Server configuration in a log entry.

**"data":** **<json-value>** The data in a log entry.

**"eid":** **<raw-uuid>** Entry ID in a log entry.

**"vote":** **<raw-uuid>** The server ID for which this server voted.

**"leader":** **<raw-uuid>** The server ID of the server. Emitted by both leaders and followers when a leader is elected.

**"commit\_index":** **<uint64>** Updated `commit_index` value.

**"note":** **<string>** One of a few special strings indicating important events. The currently defined strings are:

**"transfer leadership"** This server transferred leadership to a different server (with details included in `comment`).

**"left"** This server finished leaving the cluster. (This lets subsequent readers know that the server is not part of the cluster and should not attempt to connect to it.)

## Joining a Cluster

In addition to general format for a clustered database, there is also a special case for a database file created by `ovsdb-tool join-cluster`. Such a file contains exactly one record, which conveys the information passed to the `join-cluster` command. It has the following members:

**"server\_id":** **<raw-uuid>** and **"local\_address":** **<address>** and **"name":** **<id>** These have the same semantics described above in the general description of the format.

**"cluster\_id":** **<raw-uuid>** This is provided only if the user gave the `--cid` option to `join-cluster`. It has the same semantics described above.

**"remote\_addresses";** [**<address>\***] One or more remote servers to contact for joining the cluster.

When the server successfully joins the cluster, the database file is replaced by one described in *Clustered Format*.

## 6.1.6 ovssdb

### Description

OVSSDB, the Open vSwitch Database, is a network-accessible database system. Schemas in OVSSDB specify the tables in a database and their columns' types and can include data, uniqueness, and referential integrity constraints. OVSSDB offers atomic, consistent, isolated, durable transactions. RFC 7047 specifies the JSON-RPC based protocol that OVSSDB clients and servers use to communicate.

The OVSSDB protocol is well suited for state synchronization because it allows each client to monitor the contents of a whole database or a subset of it. Whenever a monitored portion of the database changes, the server tells the client what rows were added or modified (including the new contents) or deleted. Thus, OVSSDB clients can easily keep track of the newest contents of any part of the database.

While OVSSDB is general-purpose and not particularly specialized for use with Open vSwitch, Open vSwitch does use it for multiple purposes. The leading use of OVSSDB is for configuring and monitoring `ovs-vswitchd`(8), the Open vSwitch switch daemon, using the schema documented in `ovs-vswitchd.conf.db`(5). The Open Virtual Network (OVN) sub-project of OVS uses two OVSSDB schemas, documented in `ovn-nb`(5) and `ovn-sb`(5). Finally, Open vSwitch includes the "VTEP" schema, documented in `vtep`(5) that many third-party hardware switches support for configuring VXLAN, although OVS itself does not directly use this schema.

The OVSSDB protocol specification allows independent, interoperable implementations of OVSSDB to be developed. Open vSwitch includes an OVSSDB server implementation named `ovssdb-server`(1), which supports several protocol extensions documented in its manpage, and a basic command-line OVSSDB client named `ovssdb-client`(1), as well as OVSSDB client libraries for C and for Python. Open vSwitch documentation often speaks of these OVSSDB implementations in Open vSwitch as simply "OVSSDB," even though that is distinct from the OVSSDB protocol; we make the distinction explicit only when it might otherwise be unclear from the context.

In addition to these generic OVSSDB server and client tools, Open vSwitch includes tools for working with databases that have specific schemas: `ovs-vsctl` works with the `ovs-vswitchd` configuration database, `vtepctl` works with the VTEP database, `ovn-nbctl` works with the OVN Northbound database, and so on.

RFC 7047 specifies the OVSSDB protocol but it does not specify an on-disk storage format. Open vSwitch includes `ovssdb-tool`(1) for working with its own on-disk database formats. The most notable feature of this format is that `ovssdb-tool`(1) makes it easy for users to print the transactions that have changed a database since the last time it was compacted. This feature is often useful for troubleshooting.

### Schemas

Schemas in OVSSDB have a JSON format that is specified in RFC 7047. They are often stored in files with an extension `.ovsschema`. An on-disk database in OVSSDB includes a schema and data, embedding both into a single file. The Open vSwitch utility `ovssdb-tool` has commands that work with schema files and with the schemas embedded in database files.

An Open vSwitch schema has three important identifiers. The first is its name, which is also the name used in JSON-RPC calls to identify a database based on that schema. For example, the schema used to configure Open vSwitch has the name `Open_vSwitch`. Schema names begin with a letter or an underscore, followed by any number of letters, underscores, or digits. The `ovssdb-tool` commands `schema-name` and `db-name` extract the schema name from a schema or database file, respectively.

An OVSSDB schema also has a version of the form `x.y.z` e.g. `1.2.3`. Schemas managed within the Open vSwitch project manage version numbering in the following way (but OVSSDB does not mandate this approach). Whenever we change the database schema in a non-backward compatible way (e.g. when we delete a column or a table), we increment `<x>` and set `<y>` and `<z>` to 0. When we change the database schema in a backward compatible way (e.g. when we add a new column), we increment `<y>` and set `<z>` to 0. When we change the database schema cosmetically (e.g. we reindent its syntax), we increment `<z>`. The `ovssdb-tool` commands `schema-version` and `db-version` extract the schema version from a schema or database file, respectively.

Very old OVSDb schemas do not have a version, but RFC 7047 mandates it.

An OVSDb schema optionally has a “checksum.” RFC 7047 does not specify the use of the checksum and recommends that clients ignore it. Open vSwitch uses the checksum to remind developers to update the version: at build time, if the schema’s embedded checksum, ignoring the checksum field itself, does not match the schema’s content, then it fails the build with a recommendation to update the version and the checksum. Thus, a developer who changes the schema, but does not update the version, receives an automatic reminder. In practice this has been an effective way to ensure compliance with the version number policy. The `ovsdb-tool` commands `schema-cksum` and `db-cksum` extract the schema checksum from a schema or database file, respectively.

## Service Models

OVSDb supports three service models for databases: **standalone**, **active-backup**, and **clustered**. The service models provide different compromises among consistency, availability, and partition tolerance. They also differ in the number of servers required and in terms of performance. The standalone and active-backup database service models share one on-disk format, and clustered databases use a different format, but the OVSDb programs work with both formats. `ovsdb(5)` documents these file formats.

RFC 7047, which specifies the OVSDb protocol, does not mandate or specify any particular service model.

The following sections describe the individual service models.

### Standalone Database Service Model

A **standalone** database runs a single server. If the server stops running, the database becomes inaccessible, and if the server’s storage is lost or corrupted, the database’s content is lost. This service model is appropriate when the database controls a process or activity to which it is linked via “fate-sharing.” For example, an OVSDb instance that controls an Open vSwitch virtual switch daemon, `ovs-vswitchd`, is a standalone database because a server failure would take out both the database and the virtual switch.

To set up a standalone database, use `ovsdb-tool create` to create a database file, then run `ovsdb-server` to start the database service.

To configure a client, such as `ovs-vswitchd` or `ovs-vsctl`, to use a standalone database, configure the server to listen on a “connection method” that the client can reach, then point the client to that connection method. See [Connection Methods](#) below for information about connection methods.

### Active-Backup Database Service Model

An **active-backup** database runs two servers (on different hosts). At any given time, one of the servers is designated with the **active** role and the other the **backup** role. An active server behaves just like a standalone server. A backup server makes an OVSDb connection to the active server and uses it to continuously replicate its content as it changes in real time. OVSDb clients can connect to either server but only the active server allows data modification or lock transactions.

Setup for an active-backup database starts from a working standalone database service, which is initially the active server. On another node, to set up a backup server, create a database file with the same schema as the active server. The initial contents of the database file do not matter, as long as the schema is correct, so `ovsdb-tool create` will work, as will copying the database file from the active server. Then use `ovsdb-server --sync-from=<active>` to start the backup server, where `<active>` is an OVSDb connection method (see [Connection Methods](#) below) that connects to the active server. At that point, the backup server will fetch a copy of the active database and keep it up-to-date until it is killed.

When the active server in an active-backup server pair fails, an administrator can switch the backup server to an active role with the `ovs-appctl` command `ovsdb-server/disconnect-active-ovsdb-server`. Clients then



have read/write access to the now-active server. Of course, administrators are slow to respond compared to software, so in practice external management software detects the active server's failure and changes the backup server's role. For example, the "Integration Guide for Centralized Control" in the Open vSwitch documentation describes how to use Pacemaker for this purpose in OVN.

Suppose an active server fails and its backup is promoted to active. If the failed server is revived, it must be started as a backup server. Otherwise, if both servers are active, then they may start out of sync, if the database changed while the server was down, and they will continue to diverge over time. This also happens if the software managing the database servers cannot reach the active server and therefore switches the backup to active, but other hosts can reach both servers. These "split-brain" problems are unsolvable in general for server pairs.

Compared to a standalone server, the active-backup service model somewhat increases availability, at a risk of split-brain. It adds generally insignificant performance overhead. On the other hand, the clustered service model, discussed below, requires at least 3 servers and has greater performance overhead, but it avoids the need for external management software and eliminates the possibility of split-brain.

Open vSwitch 2.6 introduced support for the active-backup service model.

### Clustered Database Service Model

A **clustered** database runs across 3 or 5 or more database servers (the **cluster**) on different hosts. Servers in a cluster automatically synchronize writes within the cluster. A 3-server cluster can remain available in the face of at most 1 server failure; a 5-server cluster tolerates up to 2 failures. Clusters larger than 5 servers will also work, with every 2 added servers allowing the cluster to tolerate 1 more failure, but write performance decreases. The number of servers should be odd: a 4- or 6-server cluster cannot tolerate more failures than a 3- or 5-server cluster, respectively.

To set up a clustered database, first initialize it on a single node by running `ovsdb-tool create-cluster`, then start `ovsdb-server`. Depending on its arguments, the `create-cluster` command can create an empty database or copy a standalone database's contents into the new database.

To configure a client, such as `ovn-controller` or `ovn-sbctl`, to use a clustered database, first configure all of the servers to listen on a connection method that the client can reach, then point the client to all of the servers' connection methods, comma-separated. See [Connection Methods](#), below, for more detail.

Open vSwitch 2.9 introduced support for the clustered service model.

### How to Maintain a Clustered Database

To add a server to a cluster, run `ovsdb-tool join-cluster` on the new server and start `ovsdb-server`. To remove a running server from a cluster, use `ovs-appctl` to invoke the `cluster/leave` command. When a server fails and cannot be recovered, e.g. because its hard disk crashed, or to otherwise remove a server that is down from a cluster, use `ovs-appctl` to invoke `cluster/kick` to make the remaining servers kick it out of the cluster.

The above methods for adding and removing servers only work for healthy clusters, that is, for clusters with no more failures than their maximum tolerance. For example, in a 3-server cluster, the failure of 2 servers prevents servers joining or leaving the cluster (as well as database access). To prevent data loss or inconsistency, the preferred solution to this problem is to bring up enough of the failed servers to make the cluster healthy again, then if necessary remove any remaining failed servers and add new ones. If this cannot be done, though, use `ovs-appctl` to invoke `cluster/leave --force` on a running server. This command forces the server to which it is directed to leave its cluster and form a new single-node cluster that contains only itself. The data in the new cluster may be inconsistent with the former cluster: transactions not yet replicated to the server will be lost, and transactions not yet applied to the cluster may be committed. Afterward, any servers in its former cluster will regard the server to have failed.

Once a server leaves a cluster, it may never rejoin it. Instead, create a new server and join it to the cluster.



The servers in a cluster synchronize data over a cluster management protocol that is specific to Open vSwitch; it is not the same as the OVSDb protocol specified in RFC 7047. For this purpose, a server in a cluster is tied to a particular IP address and TCP port, which is specified in the `ovsdb-tool` command that creates or joins the cluster. The TCP port used for clustering must be different from that used for OVSDb clients. To change the port or address of a server in a cluster, first remove it from the cluster, then add it back with the new address.

To upgrade the `ovsdb-server` processes in a cluster from one version of Open vSwitch to another, upgrading them one at a time will keep the cluster healthy during the upgrade process. (This is different from upgrading a database schema, which is covered later under *Upgrading or Downgrading a Database*.)

Clustered OVSDb does not support the OVSDb “ephemeral columns” feature. `ovsdb-tool` and `ovsdb-client` change ephemeral columns into persistent ones when they work with schemas for clustered databases. Future versions of OVSDb might add support for this feature.

## Understanding Cluster Consistency

To ensure consistency, clustered OVSDb uses the Raft algorithm described in Diego Ongaro’s Ph.D. thesis, “Consensus: Bridging Theory and Practice”. In an operational Raft cluster, at any given time a single server is the “leader” and the other nodes are “followers”. Only the leader processes transactions, but a transaction is only committed when a majority of the servers confirm to the leader that they have written it to persistent storage.

In most database systems, read and write access to the database happens through transactions. In such a system, Raft allows a cluster to present a strongly consistent transactional interface. OVSDb uses conventional transactions for writes, but clients often effectively do reads a different way, by asking the server to “monitor” a database or a subset of one on the client’s behalf. Whenever monitored data changes, the server automatically tells the client what changed, which allows the client to maintain an accurate snapshot of the database in its memory. Of course, at any given time, the snapshot may be somewhat dated since some of it could have changed without the change notification yet being received and processed by the client.

Given this unconventional usage model, OVSDb also adopts an unconventional clustering model. Each server in a cluster acts independently for the purpose of monitors and read-only transactions, without verifying that data is up-to-date with the leader. Servers forward transactions that write to the database to the leader for execution, ensuring consistency. This has the following consequences:

- Transactions that involve writes, against any server in the cluster, are linearizable if clients take care to use correct prerequisites, which is the same condition required for linearizability in a standalone OVSDb. (Actually, “at-least-once” consistency, because OVSDb does not have a session mechanism to drop duplicate transactions if a connection drops after the server commits it but before the client receives the result.)
- Read-only transactions can yield results based on a stale version of the database, if they are executed against a follower. Transactions on the leader always yield fresh results. (With monitors, as explained above, a client can always see stale data even without clustering, so clustering does not change the consistency model for monitors.)
- Monitor-based (or read-heavy) workloads scale well across a cluster, because clustering OVSDb adds no additional work or communication for reads and monitors.
- A write-heavy client should connect to the leader, to avoid the overhead of followers forwarding transactions to the leader.
- When a client conducts a mix of read and write transactions across more than one server in a cluster, it can see inconsistent results because a read transaction might read stale data whose updates have not yet propagated from the leader. By default, `ovn-sbctl` and similar utilities connect to the cluster leader to avoid this issue.

The same might occur for transactions against a single follower except that the OVSDb server ensures that the results of a write forwarded to the leader by a given server are visible at that server before it replies to the requesting client.

- If a client uses a database on one server in a cluster, then another server in the cluster (perhaps because the first server failed), the client could observe stale data. Clustered OVSDB clients, however, can use a column in the `_Server` database to detect that data on a server is older than data that the client previously read. The OVSDB client library in Open vSwitch uses this feature to avoid servers with stale data.

## Database Replication

OVSDB can layer **replication** on top of any of its service models. Replication, in this context, means to make, and keep up-to-date, a read-only copy of the contents of a database (the `replica`). One use of replication is to keep an up-to-date backup of a database. A replica used solely for backup would not need to support clients of its own. A set of replicas that do serve clients could be used to scale out read access to the primary database.

A database replica is set up in the same way as a backup server in an active-backup pair, with the difference that the replica is never promoted to an active role.

A database can have multiple replicas.

Open vSwitch 2.6 introduced support for database replication.

## Connection Methods

An OVSDB **connection method** is a string that specifies how to make a JSON-RPC connection between an OVSDB client and server. Connection methods are part of the Open vSwitch implementation of OVSDB and not specified by RFC 7047. `ovsdb-server` uses connection methods to specify how it should listen for connections from clients and `ovsdb-client` uses them to specify how it should connect to a server. Connections in the opposite direction, where `ovsdb-server` connects to a client that is configured to listen for an incoming connection, are also possible.

Connection methods are classified as **active** or **passive**. An active connection method makes an outgoing connection to a remote host; a passive connection method listens for connections from remote hosts. The most common arrangement is to configure an OVSDB server with passive connection methods and clients with active ones, but the OVSDB implementation in Open vSwitch supports the opposite arrangement as well.

OVSDB supports the following active connection methods:

**ssl:<host>:<port>** The specified SSL or TLS <port> on the given <host>.

**tcp:<host>:<port>** The specified TCP <port> on the given <host>.

**unix:<file>** On Unix-like systems, connect to the Unix domain server socket named <file>.

On Windows, connect to a local named pipe that is represented by a file created in the path <file> to mimic the behavior of a Unix domain socket.

**<method1>,<method2>,...,<methodN>** For a clustered database service to be highly available, a client must be able to connect to any of the servers in the cluster. To do so, specify connection methods for each of the servers separated by commas (and optional spaces).

In theory, if machines go up and down and IP addresses change in the right way, a client could talk to the wrong instance of a database. To avoid this possibility, add `cid:<uuid>` to the list of methods, where <uuid> is the cluster ID of the desired database cluster, as printed by `ovsdb-tool db-cid`. This feature is optional.

OVSDB supports the following passive connection methods:

**pssl:<port>[:<ip>]** Listen on the given TCP <port> for SSL or TLS connections. By default, connections are not bound to a particular local IP address. Specifying <ip> limits connections to those from the given IP.

**ptcp:<port>[:<ip>]** Listen on the given TCP <port>. By default, connections are not bound to a particular local IP address. Specifying <ip> limits connections to those from the given IP.

**punix:<file>** On Unix-like systems, listens for connections on the Unix domain socket named <file>.

On Windows, listens on a local named pipe, creating a named pipe <file> to mimic the behavior of a Unix domain socket.

All IP-based connection methods accept IPv4 and IPv6 addresses. To specify an IPv6 address, wrap it in square brackets, e.g. `ssl::[::1]:6640`. Passive IP-based connection methods by default listen for IPv4 connections only; use `[::]` as the address to accept both IPv4 and IPv6 connections, e.g. `pssl:6640::[::]`. DNS names are also accepted if built with unbound library. On Linux, use `%<device>` to designate a scope for IPv6 link-level addresses, e.g. `ssl:[fe80::1234%eth0]:6653`.

The <port> may be omitted from connection methods that use a port number. The default <port> for TCP-based connection methods is 6640, e.g. `pssl:` is equivalent to `pssl:6640`. In Open vSwitch prior to version 2.4.0, the default port was 6632. To avoid incompatibility between older and newer versions, we encourage users to specify a port number.

The `ssl` and `pssl` connection methods requires additional configuration through `--private-key`, `--certificate`, and `--ca-cert` command line options. Open vSwitch can be built without SSL support, in which case these connection methods are not supported.

## Database Life Cycle

This section describes how to handle various events in the life cycle of a database using the Open vSwitch implementation of OVSDb.

### Creating a Database

Creating and starting up the service for a new database was covered separately for each database service model in the *Service Models* section, above.

### Backing Up and Restoring a Database

OVSDb is often used in contexts where the database contents are not particularly valuable. For example, in many systems, the database for configuring `ovs-vswitchd` is essentially rebuilt from scratch at boot time. It is not worthwhile to back up these databases.

When OVSDb is used for valuable data, a backup strategy is worth considering. One way is to use database replication, discussed above in *Database Replication* which keeps an online, up-to-date copy of a database, possibly on a remote system. This works with all OVSDb service models.

A more common backup strategy is to periodically take and store a snapshot. For the standalone and active-backup service models, making a copy of the database file, e.g. using `cp`, effectively makes a snapshot, and because OVSDb database files are append-only, it works even if the database is being modified when the snapshot takes place. This approach does not work for clustered databases.

Another way to make a backup, which works with all OVSDb service models, is to use `ovsdb-client backup`, which connects to a running database server and outputs an atomic snapshot of its schema and content, in the same format used for standalone and active-backup databases.

Multiple options are also available when the time comes to restore a database from a backup. For the standalone and active-backup service models, one option is to stop the database server or servers, overwrite the database file with the backup (e.g. with `cp`), and then restart the servers. Another way, which works with any service model, is to use `ovsdb-client restore`, which connects to a running database server and replaces the data in one of its databases by a provided snapshot. The advantage of `ovsdb-client restore` is that it causes zero downtime for

the database and its server. It has the downside that UUIDs of rows in the restored database will differ from those in the snapshot, because the OVSDb protocol does not allow clients to specify row UUIDs.

None of these approaches saves and restores data in columns that the schema designates as ephemeral. This is by design: the designer of a schema only marks a column as ephemeral if it is acceptable for its data to be lost when a database server restarts.

Clustering and backup serve different purposes. Clustering increases availability, but it does not protect against data loss if, for example, a malicious or malfunctioning OVSDb client deletes or tampers with data.

## Changing Database Service Model

Use `ovsdb-tool create-cluster` to create a clustered database from the contents of a standalone database. Use `ovsdb-tool backup` to create a standalone database from the contents of a clustered database.

## Upgrading or Downgrading a Database

The evolution of a piece of software can require changes to the schemas of the databases that it uses. For example, new features might require new tables or new columns in existing tables, or conceptual changes might require a database to be reorganized in other ways. In some cases, the easiest way to deal with a change in a database schema is to delete the existing database and start fresh with the new schema, especially if the data in the database is easy to reconstruct. But in many other cases, it is better to convert the database from one schema to another.

The OVSDb implementation in Open vSwitch has built-in support for some simple cases of converting a database from one schema to another. This support can handle changes that add or remove database columns or tables or that eliminate constraints (for example, changing a column that must have exactly one value into one that has one or more values). It can also handle changes that add constraints or make them stricter, but only if the existing data in the database satisfies the new constraints (for example, changing a column that has one or more values into a column with exactly one value, if every row in the column has exactly one value). The built-in conversion can cause data loss in obvious ways, for example if the new schema removes tables or columns, or indirectly, for example by deleting unreferenced rows in tables that the new schema marks for garbage collection.

Converting a database can lose data, so it is wise to make a backup beforehand.

To use OVSDb's built-in support for schema conversion with a standalone or active-backup database, first stop the database server or servers, then use `ovsdb-tool convert` to convert it to the new schema, and then restart the database server.

OVSDb also supports online database schema conversion for any of its database service models. To convert a database online, use `ovsdb-client convert`. The conversion is atomic, consistent, isolated, and durable. `ovsdb-server` disconnects any clients connected when the conversion takes place (except clients that use the `set_db_change_aware` Open vSwitch extension RPC). Upon reconnection, clients will discover that the schema has changed.

Schema versions and checksums (see [Schemas](#) above) can give hints about whether a database needs to be converted to a new schema. If there is any question, though, the `needs-conversion` command on `ovsdb-tool` and `ovsdb-client` can provide a definitive answer.

## Working with Database History

Both on-disk database formats that OVSDb supports are organized as a stream of transaction records. Each record describes a change to the database as a list of rows that were inserted or deleted or modified, along with the details. Therefore, in normal operation, a database file only grows, as each change causes another record to be appended at the end. Usually, a user has no need to understand this file structure. This section covers some exceptions.

## Compacting Databases

If OVSDb database files were truly append-only, then over time they would grow without bound. To avoid this problem, OVSDb can **compact** a database file, that is, replace it by a new version that contains only the current database contents, as if it had been inserted by a single transaction. From time to time, `ovsdb-server` automatically compacts a database that grows much larger than its minimum size.

Because `ovsdb-server` automatically compacts databases, it is usually not necessary to compact them manually, but OVSDb still offers a few ways to do it. First, `ovsdb-tool compact` can compact a standalone or active-backup database that is not currently being served by `ovsdb-server` (or otherwise locked for writing by another process). To compact any database that is currently being served by `ovsdb-server`, use `ovs-appctl` to send the `ovsdb-server/compact` command. Each server in an active-backup or clustered database maintains its database file independently, so to compact all of them, issue this command separately on each server.

## Viewing History

The `ovsdb-tool` utility's `show-log` command displays the transaction records in an OVSDb database file in a human-readable format. By default, it shows minimal detail, but adding the option `-m` once or twice increases the level of detail. In addition to the transaction data, it shows the time and date of each transaction and any “comment” added to the transaction by the client. The comments can be helpful for quickly understanding a transaction; for example, `ovs-vsctl` adds its command line to the transactions that it makes.

The `show-log` command works with both OVSDb file formats, but the details of the output format differ. For active-backup and clustered databases, the sequence of transactions in each server's log will differ, even at points when they reflect the same data.

## Truncating History

It may occasionally be useful to “roll back” a database file to an earlier point. Because of the organization of OVSDb records, this is easy to do. Start by noting the record number `<i>` of the first record to delete in `ovsdb-tool show-log` output. Each record is two lines of plain text, so trimming the log is as simple as running `head -n <j>`, where `<j> = 2 * <i>`.

## Corruption

When `ovsdb-server` opens an OVSDb database file, of any kind, it reads as many transaction records as it can from the file until it reaches the end of the file or it encounters a corrupted record. At that point it stops reading and regards the data that it has read to this point as the full contents of the database file, effectively rolling the database back to an earlier point.

Each transaction record contains an embedded SHA-1 checksum, which the server verifies as it reads a database file. It detects corruption when a checksum fails to verify. Even though SHA-1 is no longer considered secure for use in cryptography, it is acceptable for this purpose because it is not used to defend against malicious attackers.

The first record in a standalone or active-backup database file specifies the schema. `ovsdb-server` will refuse to work with a database where this record is corrupted, or with a clustered database file with corruption in the first few records. Delete and recreate such a database, or restore it from a backup.

When `ovsdb-server` adds records to a database file in which it detected corruption, it first truncates the file just after the last good record.

## See Also

RFC 7047, “The Open vSwitch Database Management Protocol.”

Open vSwitch implementations of generic OVSDb functionality: `ovsdb-server(1)`, `ovsdb-client(1)`, `ovsdb-tool(1)`.

Tools for working with databases that have specific OVSDb schemas: `ovs-vsctl(8)`, `vtep-ctl(8)`, `ovn-nbctl(8)`, `ovn-sbctl(8)`.

OVSDb schemas for Open vSwitch and related functionality: `ovs-vswitchd.conf.db(5)`, `vtep(5)`, `ovn-nb(5)`, `ovn-sb(5)`.

The remainder are still in roff format can be found below:

<code>ovs-actions(7)</code>	(pdf)	(html)	(plain text)
<code>ovn-architecture(7)</code>	(pdf)	(html)	(plain text)
<code>ovn-controller(8)</code>	(pdf)	(html)	(plain text)
<code>ovn-controller-vtep(8)</code>	(pdf)	(html)	(plain text)
<code>ovn-ctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovn-nb(5)</code>	(pdf)	(html)	(plain text)
<code>ovn-nbctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovn-northd(8)</code>	(pdf)	(html)	(plain text)
<code>ovn-sb(5)</code>	(pdf)	(html)	(plain text)
<code>ovn-sbctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovn-trace(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-appctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-bugtool(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-ctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovsdb-client(1)</code>	(pdf)	(html)	(plain text)
<code>ovsdb-server(1)</code>	(pdf)	(html)	(plain text)
<code>ovsdb-tool(1)</code>	(pdf)	(html)	(plain text)
<code>ovs-dpctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-dpctl-top(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-fields(7)</code>	(pdf)	(html)	(plain text)
<code>ovs-l3ping(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-ofctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-parse-backtrace(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-pcap(1)</code>	(pdf)	(html)	(plain text)
<code>ovs-pki(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-tcpdump(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-tcpundump(1)</code>	(pdf)	(html)	(plain text)
<code>ovs-test(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-testcontroller(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-vlan-bug-workaround(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-vlan-test(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-vsctl(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-vswitchd(8)</code>	(pdf)	(html)	(plain text)
<code>ovs-vswitchd.conf.db(5)</code>	(pdf)	(html)	(plain text)
<code>vtep(5)</code>	(pdf)	(html)	(plain text)
<code>vtep-ctl(8)</code>	(pdf)	(html)	(plain text)

## 7.1 Basic Configuration

Q: How do I configure a port as an access port?

A. Add `tag=VLAN` to your `ovs-vsctl add-port` command. For example, the following commands configure `br0` with `eth0` as a trunk port (the default) and `tap0` as an access port for VLAN 9:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0 tag=9
```

If you want to configure an already added port as an access port, use `ovs-vsctl set`, e.g.:

```
$ ovs-vsctl set port tap0 tag=9
```

Q: How do I configure a port as a SPAN port, that is, enable mirroring of all traffic to that port?

A. The following commands configure `br0` with `eth0` and `tap0` as trunk ports. All traffic coming in or going out on `eth0` or `tap0` is also mirrored to `tap1`; any traffic arriving on `tap1` is dropped:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0
$ ovs-vsctl add-port br0 tap1 \
 -- --id=@p get port tap1 \
 -- --id=@m create mirror name=m0 select-all=true output-port=@p \
 -- set bridge br0 mirrors=@m
```

To later disable mirroring, run:

```
$ ovs-vsctl clear bridge br0 mirrors
```

Q: Does Open vSwitch support configuring a port in promiscuous mode?

A: Yes. How you configure it depends on what you mean by “promiscuous mode”:



- Conventionally, “promiscuous mode” is a feature of a network interface card. Ordinarily, a NIC passes to the CPU only the packets actually destined to its host machine. It discards the rest to avoid wasting memory and CPU cycles. When promiscuous mode is enabled, however, it passes every packet to the CPU. On an old-style shared-media or hub-based network, this allows the host to spy on all packets on the network. But in the switched networks that are almost everywhere these days, promiscuous mode doesn’t have much effect, because few packets not destined to a host are delivered to the host’s NIC.

This form of promiscuous mode is configured in the guest OS of the VMs on your bridge, e.g. with “ip link set <device> promisc”.

- The VMware vSwitch uses a different definition of “promiscuous mode”. When you configure promiscuous mode on a VMware vNIC, the vSwitch sends a copy of every packet received by the vSwitch to that vNIC. That has a much bigger effect than just enabling promiscuous mode in a guest OS. Rather than getting a few stray packets for which the switch does not yet know the correct destination, the vNIC gets every packet. The effect is similar to replacing the vSwitch by a virtual hub.

This “promiscuous mode” is what switches normally call “port mirroring” or “SPAN”. For information on how to configure SPAN, see “How do I configure a port as a SPAN port, that is, enable mirroring of all traffic to that port?”

Q: How do I configure a DPDK port as an access port?

A: Firstly, you must have a DPDK-enabled version of Open vSwitch.

If your version is DPDK-enabled it may support the `dpdk_version` and `dpdk_initialized` keys in the configuration database. Earlier versions of Open vSwitch only supported the `other-config:dpdk-init` key in the configuration in the database. All versions will display lines with “EAL:...” during startup when `other_config:dpdk-init` is set to ‘true’.

Secondly, when adding a DPDK port, unlike a system port, the type for the interface and valid `dpdk-devargs` must be specified. For example:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 myportname -- set Interface myportname \
 type=dpdk options:dpdk-devargs=0000:06:00.0
```

Refer to *Open vSwitch with DPDK* for more information on enabling and using DPDK with Open vSwitch.

Q: How do I configure a VLAN as an RSPAN VLAN, that is, enable mirroring of all traffic to that VLAN?

A: The following commands configure `br0` with `eth0` as a trunk port and `tap0` as an access port for VLAN 10. All traffic coming in or going out on `tap0`, as well as traffic coming in or going out on `eth0` in VLAN 10, is also mirrored to VLAN 15 on `eth0`. The original tag for VLAN 10, in cases where one is present, is dropped as part of mirroring:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0 tag=10
$ ovs-vsctl \
 -- --id=@m create mirror name=m0 select-all=true select-vlan=10 \
 output-vlan=15 \
 -- set bridge br0 mirrors=@m
```

To later disable mirroring, run:

```
$ ovs-vsctl clear bridge br0 mirrors
```



Mirroring to a VLAN can disrupt a network that contains unmanaged switches. See `ovs-vswitchd.conf.db(5)` for details. Mirroring to a GRE tunnel has fewer caveats than mirroring to a VLAN and should generally be preferred.

**Q:** Can I mirror more than one input VLAN to an RSPAN VLAN?

**A:** Yes, but mirroring to a VLAN strips the original VLAN tag in favor of the specified output-vlan. This loss of information may make the mirrored traffic too hard to interpret.

To mirror multiple VLANs, use the commands above, but specify a comma-separated list of VLANs as the value for `select-vlan`. To mirror every VLAN, use the commands above, but omit `select-vlan` and its value entirely.

When a packet arrives on a VLAN that is used as a mirror output VLAN, the mirror is disregarded. Instead, in standalone mode, OVS floods the packet across all the ports for which the mirror output VLAN is configured. (If an OpenFlow controller is in use, then it can override this behavior through the flow table.) If OVS is used as an intermediate switch, rather than an edge switch, this ensures that the RSPAN traffic is distributed through the network.

Mirroring to a VLAN can disrupt a network that contains unmanaged switches. See `ovs-vswitchd.conf.db(5)` for details. Mirroring to a GRE tunnel has fewer caveats than mirroring to a VLAN and should generally be preferred.

**Q:** How do I configure mirroring of all traffic to a GRE tunnel?

**A:** The following commands configure `br0` with `eth0` and `tap0` as trunk ports. All traffic coming in or going out on `eth0` or `tap0` is also mirrored to `gre0`, a GRE tunnel to the remote host `192.168.1.10`; any traffic arriving on `gre0` is dropped:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0
$ ovs-vsctl add-port br0 gre0 \
 -- set interface gre0 type=gre options:remote_ip=192.168.1.10 \
 -- --id=@p get port gre0 \
 -- --id=@m create mirror name=m0 select-all=true output-port=@p \
 -- set bridge br0 mirrors=@m
```

To later disable mirroring and destroy the GRE tunnel:

```
$ ovs-vsctl clear bridge br0 mirrors
$ ovs-vsctl del-port br0 gre0
```

**Q:** Does Open vSwitch support ERSPAN?

**A:** Yes. ERSPAN version I and version II over IPv4 GRE and IPv6 GRE tunnel are supported. See `ovs-fields(7)` for matching and setting ERSPAN fields.

```
$ ovs-vsctl add-br br0
$ #For ERSPAN type 2 (version I)
$ ovs-vsctl add-port br0 at_erspan0 -- \
 set int at_erspan0 type=erspan options:key=1 \
 options:remote_ip=172.31.1.1 \
 options:erspan_ver=1 options:erspan_idx=1
$ #For ERSPAN type 3 (version II)
$ ovs-vsctl add-port br0 at_erspan0 -- \
 set int at_erspan0 type=erspan options:key=1 \
 options:remote_ip=172.31.1.1 \
 options:erspan_ver=2 options:erspan_dir=1 \
 options:erspan_hwid=4
```

Q: How do I connect two bridges?

A: First, why do you want to do this? Two connected bridges are not much different from a single bridge, so you might as well just have a single bridge with all your ports on it.

If you still want to connect two bridges, you can use a pair of patch ports. The following example creates bridges br0 and br1, adds eth0 and tap0 to br0, adds tap1 to br1, and then connects br0 and br1 with a pair of patch ports.

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0
$ ovs-vsctl add-br br1
$ ovs-vsctl add-port br1 tap1
$ ovs-vsctl \
 -- add-port br0 patch0 \
 -- set interface patch0 type=patch options:peer=patch1 \
 -- add-port br1 patch1 \
 -- set interface patch1 type=patch options:peer=patch0
```

Bridges connected with patch ports are much like a single bridge. For instance, if the example above also added eth1 to br1, and both eth0 and eth1 happened to be connected to the same next-hop switch, then you could loop your network just as you would if you added eth0 and eth1 to the same bridge (see the “Configuration Problems” section below for more information).

If you are using Open vSwitch 1.9 or an earlier version, then you need to be using the kernel module bundled with Open vSwitch rather than the one that is integrated into Linux 3.3 and later, because Open vSwitch 1.9 and earlier versions need kernel support for patch ports. This also means that in Open vSwitch 1.9 and earlier, patch ports will not work with the userspace datapath, only with the kernel module.

Q: How do I configure a bridge without an OpenFlow local port? (Local port in the sense of OFPP\_LOCAL)

A: Open vSwitch does not support such a configuration. Bridges always have their local ports.

## 7.2 Development

Q: How do I implement a new OpenFlow message?

A: Add your new message to enum ofpraw and enum ofptype in include/openvswitch/ofp-msgs.h, following the existing pattern. Then recompile and fix all of the new warnings, implementing new functionality for the new message as needed. (If you configure with --enable-Werror, as described in *Open vSwitch on Linux, FreeBSD and NetBSD*, then it is impossible to miss any warnings.)

To add an OpenFlow vendor extension message (aka experimenter message) for a vendor that doesn't yet have any extension messages, you will also need to edit build-aux/extract-ofp-msgs and at least ofphdrs\_decode() and ofpraw\_put\_\_() in lib/ofp-msgs.c. OpenFlow doesn't standardize vendor extensions very well, so it's hard to make the process simpler than that. (If you have a choice of how to design your vendor extension messages, it will be easier if you make them resemble the ONF and OVS extension messages.)

Q: How do I add support for a new field or header?

A: Add new members for your field to struct flow in include/openvswitch/flow.h, and add new enumerations for your new field to enum mf\_field\_id in include/openvswitch/meta-flow.h, following the existing pattern. If the field uses a new OXM class, add it to OXM\_CLASSES in build-aux/extract-ofp-fields. Also, add support to miniflow\_extract() in lib/flow.c for extracting your new field from a packet into struct miniflow, and to nx\_put\_raw() in lib/nx-match.c to output your new field in OXM matches. Then

recompile and fix all of the new warnings, implementing new functionality for the new field or header as needed. (If you configure with `--enable-Werror`, as described in *Open vSwitch on Linux, FreeBSD and NetBSD*, then it is impossible to miss any warnings.)

If you want kernel datapath support for your new field, you also need to modify the kernel module for the operating systems you are interested in. This isn't mandatory, since fields understood only by userspace work too (with a performance penalty), so it's reasonable to start development without it. If you implement kernel module support for Linux, then the Linux kernel "netdev" mailing list is the place to submit that support first; please read up on the Linux kernel development process separately. The Windows datapath kernel module support, on the other hand, is maintained within the OVS tree, so patches for that can go directly to ovs-dev.

Q: How do I add support for a new OpenFlow action?

A: Add your new action to `enum ofp_raw_action_type` in `lib/ofp-actions.c`, following the existing pattern. Then recompile and fix all of the new warnings, implementing new functionality for the new action as needed. (If you configure with `--enable-Werror`, as described in the *Open vSwitch on Linux, FreeBSD and NetBSD*, then it is impossible to miss any warnings.)

If you need to add an OpenFlow vendor extension action for a vendor that doesn't yet have any extension actions, then you will also need to add the vendor to `vendor_map` in `build-aux/extract-ofp-actions`. Also, you will need to add support for the vendor to `ofpact_decode_raw()` and `ofpact_put_raw()` in `lib/ofp-actions.c`. (If you have a choice of how to design your vendor extension actions, it will be easier if you make them resemble the ONF and OVS extension actions.)

Q: How do I add support for a new OpenFlow error message?

A: Add your new error to `enum ofperr` in `include/openvswitch/ofp-errors.h`. Read the large comment at the top of the file for details. If you need to add an OpenFlow vendor extension error for a vendor that doesn't yet have any, first add the vendor ID to the `<name>_VENDOR_ID` list in `include/openflow/openflow-common.h`.

Q: What's a Signed-off-by and how do I provide one?

A: Free and open source software projects usually require a contributor to provide some assurance that they're entitled to contribute the code that they provide. Some projects, for example, do this with a Contributor License Agreement (CLA) or a copyright assignment that is signed on paper or electronically.

For this purpose, Open vSwitch has adopted something called the Developer's Certificate of Origin (DCO), which is also used by the Linux kernel and originated there. Informally stated, agreeing to the DCO is the developer's way of attesting that a particular commit that they are contributing is one that they are allowed to contribute. You should visit <https://developercertificate.org/> to read the full statement of the DCO, which is less than 200 words long.

To certify compliance with the Developer's Certificate of Origin for a particular commit, just add the following line to the end of your commit message, properly substituting your name and email address:

Signed-off-by: Firstname Lastname <email@example.org>

Git has special support for adding a Signed-off-by line to a commit message: when you run "git commit", just add the `-s` option, as in "git commit -s". If you use the "git citool" GUI for commits, you can add a Signed-off-by line to the commit message by pressing Control+S. Other Git user interfaces may provide similar support.

## 7.3 Implementation Details

Q: I hear OVS has a couple of kinds of flows. Can you tell me about them?

A: Open vSwitch uses different kinds of flows for different purposes:

- OpenFlow flows are the most important kind of flow. OpenFlow controllers use these flows to define a switch's policy. OpenFlow flows support wildcards, priorities, and multiple tables.

When in-band control is in use, Open vSwitch sets up a few "hidden" flows, with priority higher than a controller or the user can configure, that are not visible via OpenFlow. (See the "Controller" section of the FAQ for more information about hidden flows.)

- The Open vSwitch software switch implementation uses a second kind of flow internally. These flows, called "datapath" or "kernel" flows, do not support priorities and comprise only a single table, which makes them suitable for caching. (Like OpenFlow flows, datapath flows do support wildcarding, in Open vSwitch 1.11 and later.) OpenFlow flows and datapath flows also support different actions and number ports differently.

Datapath flows are an implementation detail that is subject to change in future versions of Open vSwitch. Even with the current version of Open vSwitch, hardware switch implementations do not necessarily use this architecture.

Users and controllers directly control only the OpenFlow flow table. Open vSwitch manages the datapath flow table itself, so users should not normally be concerned with it.

Q: Why are there so many different ways to dump flows?

A: Open vSwitch has two kinds of flows (see the previous question), so it has commands with different purposes for dumping each kind of flow:

- `ovs-ofctl dump-flows <br>` dumps OpenFlow flows, excluding hidden flows. This is the most commonly useful form of flow dump. (Unlike the other commands, this should work with any OpenFlow switch, not just Open vSwitch.)
- `ovs-appctl bridge/dump-flows <br>` dumps OpenFlow flows, including hidden flows. This is occasionally useful for troubleshooting suspected issues with in-band control.
- `ovs-dpctl dump-flows [dp]` dumps the datapath flow table entries for a Linux kernel-based datapath. In Open vSwitch 1.10 and later, `ovs-vswitchd` merges multiple switches into a single datapath, so it will show all the flows on all your kernel-based switches. This command can occasionally be useful for debugging.
- `ovs-appctl dpif/dump-flows <br>`, new in Open vSwitch 1.10, dumps datapath flows for only the specified bridge, regardless of the type.

Q: How does multicast snooping works with VLANs?

A: Open vSwitch maintains snooping tables for each VLAN.

Q: Can OVS populate the kernel flow table in advance instead of in reaction to packets?

A: No. There are several reasons:

- Kernel flows are not as sophisticated as OpenFlow flows, which means that some OpenFlow policies could require a large number of kernel flows. The "conjunctive match" feature is an extreme example: the number of kernel flows it requires is the product of the number of flows in each dimension.
- With multiple OpenFlow flow tables and simple sets of actions, the number of kernel flows required can be as large as the product of the number of flows in each dimension. With more sophisticated actions, the number of kernel flows could be even larger.
- Open vSwitch is designed so that any version of OVS userspace interoperates with any version of the OVS kernel module. This forward and backward compatibility requires that userspace observe how the kernel module parses received packets. This is only possible in a straightforward way when userspace adds kernel flows in reaction to received packets.

For more relevant information on the architecture of Open vSwitch, please read “The Design and Implementation of Open vSwitch”, published in USENIX NSDI 2015.

Q: How many packets does OVS buffer?

A: Open vSwitch fast path packet processing uses a “run to completion” model in which every packet is completely handled in a single pass. Therefore, in the common case where a packet just passes through the fast path, Open vSwitch does not buffer packets itself. The operating system and the network drivers involved in receiving and later in transmitting the packet do often include buffering. Open vSwitch is only a middleman between these and does not have direct access or influence over their buffers.

Outside the common case, Open vSwitch does sometimes buffer packets. When the OVS fast path processes a packet that does not match any of the flows in its megaflow cache, it passes that packet to the Open vSwitch slow path. This procedure queues a copy of the packet to the Open vSwitch userspace which processes it and, if necessary, passes it back to the kernel module. Queuing the packet to userspace as part of this process involves buffering. (Going the opposite direction does not, because the kernel actually processes the request synchronously.) A few other exceptional cases also queue packets to userspace for processing; most of these are due to OpenFlow actions that the fast path cannot handle and that must therefore be handled by the slow path instead.

OpenFlow also has a concept of packet buffering. When an OpenFlow switch sends a packet to a controller, it may opt to retain a copy of the packet in an OpenFlow “packet buffer”. Later, if the controller wants to tell the switch to forward a copy of that packet, it can refer to the packet through its assigned buffer, instead of sending the whole packet back to the switch, thereby saving bandwidth in the OpenFlow control channel. Before Open vSwitch 2.7, OVS implemented such buffering; Open vSwitch 2.7 and later do not.

## 7.4 General

Q: What is Open vSwitch?

A: Open vSwitch is a production quality open source software switch designed to be used as a vswitch in virtualized server environments. A vswitch forwards traffic between different VMs on the same physical host and also forwards traffic between VMs and the physical network. Open vSwitch supports standard management interfaces (e.g. sFlow, NetFlow, IPFIX, RSPAN, CLI), and is open to programmatic extension and control using OpenFlow and the OVSDDB management protocol.

Open vSwitch is designed to be compatible with modern switching chipsets. This means that it can be ported to existing high-fanout switches allowing the same flexible control of the physical infrastructure as the virtual infrastructure. It also means that Open vSwitch will be able to take advantage of on-NIC switching chipsets as their functionality matures.

Q: What virtualization platforms can use Open vSwitch?

A: Open vSwitch can currently run on any Linux-based virtualization platform (kernel 3.10 and newer), including: KVM, VirtualBox, Xen, Xen Cloud Platform, XenServer. As of Linux 3.3 it is part of the mainline kernel. The bulk of the code is written in platform-independent C and is easily ported to other environments. We welcome inquiries about integrating Open vSwitch with other virtualization platforms.

Q: How can I try Open vSwitch?

A: The Open vSwitch source code can be built on a Linux system. You can build and experiment with Open vSwitch on any Linux machine. Packages for various Linux distributions are available on many platforms, including: Debian, Ubuntu, Fedora.

You may also download and run a virtualization platform that already has Open vSwitch integrated. For example, download a recent ISO for XenServer or Xen Cloud Platform. Be aware that the version integrated with a particular platform may not be the most recent Open vSwitch release.

Q: Does Open vSwitch only work on Linux?

A: No, Open vSwitch has been ported to a number of different operating systems and hardware platforms. Most of the development work occurs on Linux, but the code should be portable to any POSIX system. We've seen Open vSwitch ported to a number of different platforms, including FreeBSD, Windows, and even non-POSIX embedded systems.

By definition, the Open vSwitch Linux kernel module only works on Linux and will provide the highest performance. However, a userspace datapath is available that should be very portable.

Q: What's involved with porting Open vSwitch to a new platform or switching ASIC?

A: *Porting Open vSwitch to New Software or Hardware* describes how one would go about porting Open vSwitch to a new operating system or hardware platform.

Q: Why would I use Open vSwitch instead of the Linux bridge?

A: Open vSwitch is specially designed to make it easier to manage VM network configuration and monitor state spread across many physical hosts in dynamic virtualized environments. Refer to *Why Open vSwitch?* for a more detailed description of how Open vSwitch relates to the Linux Bridge.

Q: How is Open vSwitch related to distributed virtual switches like the VMware vNetwork distributed switch or the Cisco Nexus 1000V?

A: Distributed vswitch applications (e.g., VMware vNetwork distributed switch, Cisco Nexus 1000V) provide a centralized way to configure and monitor the network state of VMs that are spread across many physical hosts. Open vSwitch is not a distributed vswitch itself, rather it runs on each physical host and supports remote management in a way that makes it easier for developers of virtualization/cloud management platforms to offer distributed vswitch capabilities.

To aid in distribution, Open vSwitch provides two open protocols that are specially designed for remote management in virtualized network environments: OpenFlow, which exposes flow-based forwarding state, and the OVSDb management protocol, which exposes switch port state. In addition to the switch implementation itself, Open vSwitch includes tools (ovs-ofctl, ovs-vsctl) that developers can script and extend to provide distributed vswitch capabilities that are closely integrated with their virtualization management platform.

Q: Why doesn't Open vSwitch support distribution?

A: Open vSwitch is intended to be a useful component for building flexible network infrastructure. There are many different approaches to distribution which balance trade-offs between simplicity, scalability, hardware compatibility, convergence times, logical forwarding model, etc. The goal of Open vSwitch is to be able to support all as a primitive building block rather than choose a particular point in the distributed design space.

Q: How can I contribute to the Open vSwitch Community?

A: You can start by joining the mailing lists and helping to answer questions. You can also suggest improvements to documentation. If you have a feature or bug you would like to work on, send a mail to one of the *mailing lists*.

Q: Why can I no longer connect to my OpenFlow controller or OVSDb manager?

A: Starting in OVS 2.4, we switched the default ports to the IANA-specified port numbers for OpenFlow (6633->6653) and OVSDb (6632->6640). We recommend using these port numbers, but if you cannot, all the programs allow overriding the default port. See the appropriate man page.

## 7.5 Common Configuration Issues

Q: I created a bridge and added my Ethernet port to it, using commands like these:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
```

and as soon as I ran the “add-port” command I lost all connectivity through eth0. Help!

A: A physical Ethernet device that is part of an Open vSwitch bridge should not have an IP address. If one does, then that IP address will not be fully functional.

You can restore functionality by moving the IP address to an Open vSwitch “internal” device, such as the network device named after the bridge itself. For example, assuming that eth0’s IP address is 192.168.128.5, you could run the commands below to fix up the situation:

```
$ ip addr flush dev eth0
$ ip addr add 192.168.128.5/24 dev br0
$ ip link set br0 up
```

(If your only connection to the machine running OVS is through the IP address in question, then you would want to run all of these commands on a single command line, or put them into a script.) If there were any additional routes assigned to eth0, then you would also want to use commands to adjust these routes to go through br0.

If you use DHCP to obtain an IP address, then you should kill the DHCP client that was listening on the physical Ethernet interface (e.g. eth0) and start one listening on the internal interface (e.g. br0). You might still need to manually clear the IP address from the physical interface (e.g. with “ip addr flush dev eth0”).

There is no compelling reason why Open vSwitch must work this way. However, this is the way that the Linux kernel bridge module has always worked, so it’s a model that those accustomed to Linux bridging are already used to. Also, the model that most people expect is not implementable without kernel changes on all the versions of Linux that Open vSwitch supports.

By the way, this issue is not specific to physical Ethernet devices. It applies to all network devices except Open vSwitch “internal” devices.

Q: I created a bridge and added a couple of Ethernet ports to it, using commands like these:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 eth1
```

and now my network seems to have melted: connectivity is unreliable (even connectivity that doesn’t go through Open vSwitch), all the LEDs on my physical switches are blinking, wireshark shows duplicated packets, and CPU usage is very high.

A: More than likely, you’ve looped your network. Probably, eth0 and eth1 are connected to the same physical Ethernet switch. This yields a scenario where OVS receives a broadcast packet on eth0 and sends it out on eth1, then the physical switch connected to eth1 sends the packet back on eth0, and so on forever. More complicated scenarios, involving a loop through multiple switches, are possible too.

The solution depends on what you are trying to do:

- If you added eth0 and eth1 to get higher bandwidth or higher reliability between OVS and your physical Ethernet switch, use a bond. The following commands create br0 and then add eth0 and eth1 as a bond:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-bond br0 bond0 eth0 eth1
```

Bonds have tons of configuration options. Please read the documentation on the Port table in `ovs-vswitchd.conf.db(5)` for all the details.

Configuration for DPDK-enabled interfaces is slightly less straightforward. Refer to *Open vSwitch with DPDK* for more information.

- Perhaps you don't actually need `eth0` and `eth1` to be on the same bridge. For example, if you simply want to be able to connect each of them to virtual machines, then you can put each of them on a bridge of its own:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0

$ ovs-vsctl add-br br1
$ ovs-vsctl add-port br1 eth1
```

and then connect VMs to `br0` and `br1`. (A potential disadvantage is that traffic cannot directly pass between `br0` and `br1`. Instead, it will go out `eth0` and come back in `eth1`, or vice versa.)

- If you have a redundant or complex network topology and you want to prevent loops, turn on spanning tree protocol (STP). The following commands create `br0`, enable STP, and add `eth0` and `eth1` to the bridge. The order is important because you don't want have to have a loop in your network even transiently:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl set bridge br0 stp_enable=true
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 eth1
```

The Open vSwitch implementation of STP is not well tested. Report any bugs you observe, but if you'd rather avoid acting as a beta tester then another option might be your best shot.

**Q:** I can't seem to use Open vSwitch in a wireless network.

**A:** Wireless base stations generally only allow packets with the source MAC address of NIC that completed the initial handshake. Therefore, without MAC rewriting, only a single device can communicate over a single wireless link.

This isn't specific to Open vSwitch, it's enforced by the access point, so the same problems will show up with the Linux bridge or any other way to do bridging.

**Q:** I can't seem to add my PPP interface to an Open vSwitch bridge.

**A:** PPP most commonly carries IP packets, but Open vSwitch works only with Ethernet frames. The correct way to interface PPP to an Ethernet network is usually to use routing instead of switching.

**Q:** Is there any documentation on the database tables and fields?

**A:** Yes. `ovs-vswitchd.conf.db(5)` is a comprehensive reference.

**Q:** When I run `ovs-dpctl` I no longer see the bridges I created. Instead, I only see a datapath called "ovs-system". How can I see datapath information about a particular bridge?

**A:** In version 1.9.0, OVS switched to using a single datapath that is shared by all bridges of that type. The `ovs-appctl dpif/*` commands provide similar functionality that is scoped by the bridge.

**Q:** I created a GRE port using `ovs-vsctl` so why can't I send traffic or see the port in the datapath?

**A:** On Linux kernels before 3.11, the OVS GRE module and Linux GRE module cannot be loaded at the same time. It is likely that on your system the Linux GRE module is already loaded and blocking OVS (to confirm, check `dmesg` for errors regarding GRE registration). To fix this, unload all GRE modules that appear in `lsmod` as well as the OVS kernel module. You can then reload the OVS module following



the directions in *Open vSwitch on Linux, FreeBSD and NetBSD*, which will ensure that dependencies are satisfied.

Q: Open vSwitch does not seem to obey my packet filter rules.

A: It depends on mechanisms and configurations you want to use.

You cannot usefully use typical packet filters, like iptables, on physical Ethernet ports that you add to an Open vSwitch bridge. This is because Open vSwitch captures packets from the interface at a layer lower below where typical packet-filter implementations install their hooks. (This actually applies to any interface of type “system” that you might add to an Open vSwitch bridge.)

You can usefully use typical packet filters on Open vSwitch internal ports as they are mostly ordinary interfaces from the point of view of packet filters.

For example, suppose you create a bridge br0 and add Ethernet port eth0 to it. Then you can usefully add iptables rules to affect the internal interface br0, but not the physical interface eth0. (br0 is also where you would add an IP address, as discussed elsewhere in the FAQ.)

For simple filtering rules, it might be possible to achieve similar results by installing appropriate Open-Flow flows instead. The OVS conntrack feature (see the “ct” action in ovs-actions(7)) can implement a stateful firewall.

If the use of a particular packet filter setup is essential, Open vSwitch might not be the best choice for you. On Linux, you might want to consider using the Linux Bridge. (This is the only choice if you want to use ebtables rules.) On NetBSD, you might want to consider using the bridge(4) with BRIDGE\_IPF option.

Q: It seems that Open vSwitch does nothing when I removed a port and then immediately put it back. For example, consider that p1 is a port of type=internal:

```
$ ovs-vsctl del-port br0 p1 -- \
 add-port br0 p1 -- \
 set interface p1 type=internal
```

Any other type of port gets the same effect.

A: It's an expected behaviour.

If del-port and add-port happen in a single OVSDB transaction as your example, Open vSwitch always “skips” the intermediate steps. Even if they are done in multiple transactions, it's still allowed for Open vSwitch to skip the intermediate steps and just implement the overall effect. In both cases, your example would be turned into a no-op.

If you want to make Open vSwitch actually destroy and then re-create the port for some side effects like resetting kernel setting for the corresponding interface, you need to separate operations into multiple OVSDB transactions and ensure that at least the first one does not have `--no-wait`. In the following example, the first ovs-vsctl will block until Open vSwitch reloads the new configuration and removes the port:

```
$ ovs-vsctl del-port br0 p1
$ ovs-vsctl add-port br0 p1 -- \
 set interface p1 type=internal
```

Q: I want to add thousands of ports to an Open vSwitch bridge, but it takes too long (minutes or hours) to do it with ovs-vsctl. How can I do it faster?

A: If you add them one at a time with ovs-vsctl, it can take a long time to add thousands of ports to an Open vSwitch bridge. This is because every invocation of ovs-vsctl first reads the current configuration from OVSDB. As the number of ports grows, this starts to take an appreciable amount of time, and when it is repeated thousands of times the total time becomes significant.

The solution is to add the ports in one invocation of `ovs-vsctl` (or a small number of them). For example, using bash:

```
$ ovs-vsctl add-br br0
$ cmds=; for i in {1..5000}; do cmds+=" -- add-port br0 p$i"; done
$ ovs-vsctl $cmds
```

takes seconds, not minutes or hours, in the OVS sandbox environment.

Q: I created a bridge named `br0`. My bridge shows up in “`ovs-vsctl show`”, but “`ovs-ofctl show br0`” just prints “`br0` is not a bridge or a socket”.

A: Open vSwitch wasn’t able to create the bridge. Check the `ovs-vswitchd` log for details (Debian and Red Hat packaging for Open vSwitch put it in `/var/log/openvswitch/ovs-vswitchd.log`).

In general, the Open vSwitch database reflects the desired configuration state. `ovs-vswitchd` monitors the database and, when it changes, reconfigures the system to reflect the new desired state. This normally happens very quickly. Thus, a discrepancy between the database and the actual state indicates that `ovs-vswitchd` could not implement the configuration, and so one should check the log to find out why. (Another possible cause is that `ovs-vswitchd` is not running. This will make `ovs-vsctl` commands hang, if they change the configuration, unless one specifies `--no-wait`.)

Q: I have a bridge `br0`. I added a new port `vif1.0`, and it shows up in “`ovs-vsctl show`”, but “`ovs-vsctl list port`” says that it has OpenFlow port (“`ofport`”) -1, and “`ovs-ofctl show br0`” doesn’t show `vif1.0` at all.

A: Open vSwitch wasn’t able to create the port. Check the `ovs-vswitchd` log for details (Debian and Red Hat packaging for Open vSwitch put it in `/var/log/openvswitch/ovs-vswitchd.log`). Please see the previous question for more information.

You may want to upgrade to Open vSwitch 2.3 (or later), in which `ovs-vsctl` will immediately report when there is an issue creating a port.

Q: I created a tap device `tap0`, configured an IP address on it, and added it to a bridge, like this:

```
$ tuncctl -t tap0
$ ip addr add 192.168.0.123/24 dev tap0
$ ip link set tap0 up
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 tap0
```

I expected that I could then use this IP address to contact other hosts on the network, but it doesn’t work. Why not?

A: The short answer is that this is a misuse of a “tap” device. Use an “internal” device implemented by Open vSwitch, which works differently and is designed for this use. To solve this problem with an internal device, instead run:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 int0 -- set Interface int0 type=internal
$ ip addr add 192.168.0.123/24 dev int0
$ ip link set int0 up
```

Even more simply, you can take advantage of the internal port that every bridge has under the name of the bridge:

```
$ ovs-vsctl add-br br0
$ ip addr add 192.168.0.123/24 dev br0
$ ip link set br0 up
```

In more detail, a “tap” device is an interface between the Linux (or BSD) network stack and a user program that opens it as a socket. When the “tap” device transmits a packet, it appears in the socket opened by

the userspace program. Conversely, when the userspace program writes to the “tap” socket, the kernel TCP/IP stack processes the packet as if it had been received by the “tap” device.

Consider the configuration above. Given this configuration, if you “ping” an IP address in the 192.168.0.x subnet, the Linux kernel routing stack will transmit an ARP on the tap0 device. Open vSwitch userspace treats “tap” devices just like any other network device; that is, it doesn’t open them as “tap” sockets. That means that the ARP packet will simply get dropped.

You might wonder why the Open vSwitch kernel module doesn’t intercept the ARP packet and bridge it. After all, Open vSwitch intercepts packets on other devices. The answer is that Open vSwitch only intercepts *received* packets, but this is a packet being transmitted. The same thing happens for all other types of network devices, except for Open vSwitch “internal” ports. If you, for example, add a physical Ethernet port to an OVS bridge, configure an IP address on a physical Ethernet port, and then issue a “ping” to an address in that subnet, the same thing happens: an ARP gets transmitted on the physical Ethernet port and Open vSwitch never sees it. (You should not do that, as documented at the beginning of this section.)

It can make sense to add a “tap” device to an Open vSwitch bridge, if some userspace program (other than Open vSwitch) has opened the tap socket. This is the case, for example, if the “tap” device was created by KVM (or QEMU) to simulate a virtual NIC. In such a case, when OVS bridges a packet to the “tap” device, the kernel forwards that packet to KVM in userspace, which passes it along to the VM, and in the other direction, when the VM sends a packet, KVM writes it to the “tap” socket, which causes OVS to receive it and bridge it to the other OVS ports. Please note that in such a case no IP address is configured on the “tap” device (there is normally an IP address configured in the virtual NIC inside the VM, but this is not visible to the host Linux kernel or to Open vSwitch).

There is one special case in which Open vSwitch does directly read and write “tap” sockets. This is an implementation detail of the Open vSwitch userspace switch, which implements its “internal” ports as Linux (or BSD) “tap” sockets. In such a userspace switch, OVS receives packets sent on the “tap” device used to implement an “internal” port by reading the associated “tap” socket, and bridges them to the rest of the switch. In the other direction, OVS transmits packets bridged to the “internal” port by writing them to the “tap” socket, causing them to be processed by the kernel TCP/IP stack as if they had been received on the “tap” device. Users should not need to be concerned with this implementation detail.

Open vSwitch has a network device type called “tap”. This is intended only for implementing “internal” ports in the OVS userspace switch and should not be used otherwise. In particular, users should not configure KVM “tap” devices as type “tap” (use type “system”, the default, instead).

Q: I observe packet loss at the beginning of RFC2544 tests on a server running few hundred container apps bridged to OVS with traffic generated by HW traffic generator. How can I fix this?

A: This is expected behavior on virtual switches. RFC2544 tests were designed for hardware switches, which don’t have caches on the fastpath that need to be heated. Traffic generators in order to prime the switch use learning phase to heat the caches before sending the actual traffic in test phase. In case of OVS the cache is flushed quickly and to accommodate the traffic generator’s delay between learning and test phase, the max-idle timeout settings should be changed to 50000 ms.:

```
$ ovs-vsctl --no-wait set Open_vSwitch . other_config:max-idle=50000
```

Q: How can I configure the bridge internal interface MTU? Why does Open vSwitch keep changing internal ports MTU?

A: By default Open vSwitch overrides the internal interfaces (e.g. br0) MTU. If you have just an internal interface (e.g. br0) and a physical interface (e.g. eth0), then every change in MTU to eth0 will be reflected to br0. Any manual MTU configuration using *ip* on internal interfaces is going to be overridden by Open vSwitch to match the current bridge minimum.

Sometimes this behavior is not desirable, for example with tunnels. The MTU of an internal interface can be explicitly set using the following command:

```
$ ovs-vsctl set int br0 mtu_request=1450
```

After this, Open vSwitch will configure br0 MTU to 1450. Since this setting is in the database it will be persistent (compared to what happens with *ip*).

The MTU configuration can be removed to restore the default behavior with:

```
$ ovs-vsctl set int br0 mtu_request=[]
```

The `mtu_request` column can be used to configure MTU even for physical interfaces (e.g. `eth0`).

Q: I just upgraded and I see a performance drop. Why?

A: The OVS kernel datapath may have been updated to a newer version than the OVS userspace components. Sometimes new versions of OVS kernel module add functionality that is backwards compatible with older userspace components but may cause a drop in performance with them. Especially, if a kernel module from OVS 2.1 or newer is paired with OVS userspace 1.10 or older, there will be a performance drop for TCP traffic.

Updating the OVS userspace components to the latest released version should fix the performance degradation.

To get the best possible performance and functionality, it is recommended to pair the same versions of the kernel module and OVS userspace.

Q: I can't unload the `openvswitch` kernel module. Why?

A: The userspace might still hold the reference count. So `rmmod openvswitch` does not work, for example:

```
$ lsmod | grep openvswitch
openvswitch 155648 4
nf_conntrack 24576 1 openvswitch
```

Use the command below to drop the refcnt:

```
$ ovs-dpctl del-dp system@ovs-system
$ rmmod openvswitch
```

## 7.6 Using OpenFlow

Q: What versions of OpenFlow does Open vSwitch support?

A: The following table lists the versions of OpenFlow supported by each version of Open vSwitch:

Open vSwitch	OF1.0	OF1.1	OF1.2	OF1.3	OF1.4	OF1.5
1.9 and earlier	yes	—	—	—	—	—
1.10, 1.11	yes	—	(*)	(*)	—	—
2.0, 2.1	yes	(*)	(*)	(*)	—	—
2.2	yes	(*)	(*)	(*)	(%)	(*)
2.3, 2.4	yes	yes	yes	yes	(*)	(*)
2.5, 2.6, 2.7	yes	yes	yes	yes	(*)	(*)
2.8	yes	yes	yes	yes	yes	(*)

—Not supported. yes Supported and enabled by default (\*) Supported, but missing features, and must be enabled by user. (%) Experimental, unsafe implementation.

In any case, the user may override the default:

- To enable OpenFlow 1.0, 1.1, 1.2, and 1.3 on bridge br0:

```
$ ovs-vsctl set bridge br0 \
 protocols=OpenFlow10,OpenFlow11,OpenFlow12,OpenFlow13
```

- To enable OpenFlow 1.0, 1.1, 1.2, 1.3, 1.4, and 1.5 on bridge br0:

```
$ ovs-vsctl set bridge br0 \
 protocols=OpenFlow10,OpenFlow11,OpenFlow12,OpenFlow13,OpenFlow14,
 ↪OpenFlow15
```

- To enable only OpenFlow 1.0 on bridge br0:

```
$ ovs-vsctl set bridge br0 protocols=OpenFlow10
```

All current versions of `ovs-ofctl` enable only OpenFlow 1.0 by default. Use the `-O` option to enable support for later versions of OpenFlow in `ovs-ofctl`. For example:

```
$ ovs-ofctl -O OpenFlow13 dump-flows br0
```

(Open vSwitch 2.2 had an experimental implementation of OpenFlow 1.4 that could cause crashes. We don't recommend enabling it.)

*OpenFlow Support in Open vSwitch* tracks support for OpenFlow 1.1 and later features. When support for OpenFlow 1.5 is solidly implemented, Open vSwitch will enable it by default.

Q: Does Open vSwitch support MPLS?

A: Before version 1.11, Open vSwitch did not support MPLS. That is, these versions can match on MPLS Ethernet types, but they cannot match, push, or pop MPLS labels, nor can they look past MPLS labels into the encapsulated packet.

Open vSwitch versions 1.11, 2.0, and 2.1 have very minimal support for MPLS. With the userspace datapath only, these versions can match, push, or pop a single MPLS label, but they still cannot look past MPLS labels (even after popping them) into the encapsulated packet. Kernel datapath support is unchanged from earlier versions.

Open vSwitch version 2.3 can match, push, or pop a single MPLS label and look past the MPLS label into the encapsulated packet. Both userspace and kernel datapaths will be supported, but MPLS processing always happens in userspace either way, so kernel datapath performance will be disappointing.

Open vSwitch version 2.4 can match, push, or pop up to 3 MPLS labels and look past the MPLS label into the encapsulated packet. It will have kernel support for MPLS, yielding improved performance.

Q: I'm getting "error type 45250 code 0". What's that?

A: This is a Open vSwitch extension to OpenFlow error codes. Open vSwitch uses this extension when it must report an error to an OpenFlow controller but no standard OpenFlow error code is suitable.

Open vSwitch logs the errors that it sends to controllers, so the easiest thing to do is probably to look at the `ovs-vswitchd` log to find out what the error was.

If you want to dissect the extended error message yourself, the format is documented in `include/openflow/nicira-ext.h` in the Open vSwitch source distribution. The extended error codes are documented in `include/openvswitch/ofp-errors.h`.

Q: Some of the traffic that I'd expect my OpenFlow controller to see doesn't actually appear through the OpenFlow connection, even though I know that it's going through.

A: By default, Open vSwitch assumes that OpenFlow controllers are connected “in-band”, that is, that the controllers are actually part of the network that is being controlled. In in-band mode, Open vSwitch sets up special “hidden” flows to make sure that traffic can make it back and forth between OVS and the controllers. These hidden flows are higher priority than any flows that can be set up through OpenFlow, and they are not visible through normal OpenFlow flow table dumps.

Usually, the hidden flows are desirable and helpful, but occasionally they can cause unexpected behavior. You can view the full OpenFlow flow table, including hidden flows, on bridge br0 with the command:

```
$ ovs-appctl bridge/dump-flows br0
```

to help you debug. The hidden flows are those with priorities greater than 65535 (the maximum priority that can be set with OpenFlow).

The Documentation/topics/design doc describes the in-band model in detail.

If your controllers are not actually in-band (e.g. they are on localhost via 127.0.0.1, or on a separate network), then you should configure your controllers in “out-of-band” mode. If you have one controller on bridge br0, then you can configure out-of-band mode on it with:

```
$ ovs-vsctl set controller br0 connection-mode=out-of-band
```

Q: Some of the OpenFlow flows that my controller sets up don’t seem to apply to certain traffic, especially traffic between OVS and the controller itself.

A: See above.

Q: I configured all my controllers for out-of-band control mode but “ovs-appctl bridge/dump-flows” still shows some hidden flows.

A: You probably have a remote manager configured (e.g. with “ovs-vsctl set-manager”). By default, Open vSwitch assumes that managers need in-band rules set up on every bridge. You can disable these rules on bridge br0 with:

```
$ ovs-vsctl set bridge br0 other-config:disable-in-band=true
```

This actually disables in-band control entirely for the bridge, as if all the bridge’s controllers were configured for out-of-band control.

Q: My OpenFlow controller doesn’t see the VLANs that I expect.

A: See answer under “VLANs”, above.

Q: I ran `ovs-ofctl add-flow br0 nw_dst=192.168.0.1,actions=drop` but I got a funny message like this:

```
ofp_util|INFO|normalization changed ofp_match, details:
ofp_util|INFO| pre: nw_dst=192.168.0.1
ofp_util|INFO|post:
```

and when I ran `ovs-ofctl dump-flows br0` I saw that my `nw_dst` match had disappeared, so that the flow ends up matching every packet.

A: The term “normalization” in the log message means that a flow cannot match on an L3 field without saying what L3 protocol is in use. The “ovs-ofctl” command above didn’t specify an L3 protocol, so the L3 field match was dropped.

In this case, the L3 protocol could be IP or ARP. A correct command for each possibility is, respectively:

```
$ ovs-ofctl add-flow br0 ip,nw_dst=192.168.0.1,actions=drop
```

and:

```
$ ovs-ofctl add-flow br0 arp,nw_dst=192.168.0.1,actions=drop
```

Similarly, a flow cannot match on an L4 field without saying what L4 protocol is in use. For example, the flow match `tp_src=1234` is, by itself, meaningless and will be ignored. Instead, to match TCP source port 1234, write `tcp,tp_src=1234`, or to match UDP source port 1234, write `udp,tp_src=1234`.

Q: How can I figure out the OpenFlow port number for a given port?

A: The `OFPT_FEATURES_REQUEST` message requests an OpenFlow switch to respond with an `OFPT_FEATURES_REPLY` that, among other information, includes a mapping between OpenFlow port names and numbers. From a command prompt, `ovs-ofctl show br0` makes such a request and prints the response for switch `br0`.

The Interface table in the Open vSwitch database also maps OpenFlow port names to numbers. To print the OpenFlow port number associated with interface `eth0`, run:

```
$ ovs-vsctl get Interface eth0 ofport
```

You can print the entire mapping with:

```
$ ovs-vsctl -- --columns=name,ofport list Interface
```

but the output mixes together interfaces from all bridges in the database, so it may be confusing if more than one bridge exists.

In the Open vSwitch database, `ofport` value `-1` means that the interface could not be created due to an error. (The Open vSwitch log should indicate the reason.) `ofport` value `[]` (the empty set) means that the interface hasn't been created yet. The latter is normally an intermittent condition (unless `ovs-vswitchd` is not running).

Q: I added some flows with my controller or with `ovs-ofctl`, but when I run “`ovs-dpctl dump-flows`” I don't see them.

A: `ovs-dpctl` queries a kernel datapath, not an OpenFlow switch. It won't display the information that you want. You want to use `ovs-ofctl dump-flows` instead.

Q: It looks like each of the interfaces in my bonded port shows up as an individual OpenFlow port. Is that right?

A: Yes, Open vSwitch makes individual bond interfaces visible as OpenFlow ports, rather than the bond as a whole. The interfaces are treated together as a bond for only a few purposes:

- Sending a packet to the `OFPP_NORMAL` port. (When an OpenFlow controller is not configured, this happens implicitly to every packet.)
- Mirrors configured for output to a bonded port.

It would make a lot of sense for Open vSwitch to present a bond as a single OpenFlow port. If you want to contribute an implementation of such a feature, please bring it up on the Open vSwitch development mailing list at [dev@openvswitch.org](mailto:dev@openvswitch.org).

Q: I have a sophisticated network setup involving Open vSwitch, VMs or multiple hosts, and other components. The behavior isn't what I expect. Help!

A: To debug network behavior problems, trace the path of a packet, hop-by-hop, from its origin in one host to a remote host. If that's correct, then trace the path of the response packet back to the origin.

The open source tool called `plotnetcfg` can help to understand the relationship between the networking devices on a single host.

Usually a simple ICMP echo request and reply (ping) packet is good enough. Start by initiating an ongoing ping from the origin host to a remote host. If you are tracking down a connectivity problem,

the “ping” will not display any successful output, but packets are still being sent. (In this case the packets being sent are likely ARP rather than ICMP.)

Tools available for tracing include the following:

- `tcpdump` and `wireshark` for observing hops across network devices, such as Open vSwitch internal devices and physical wires.
- `ovs-appctl dpif/dump-flows` <br> in Open vSwitch 1.10 and later or `ovs-dpctl dump-flows` <br> in earlier versions. These tools allow one to observe the actions being taken on packets in ongoing flows.

See `ovs-vswitchd(8)` for `ovs-appctl dpif/dump-flows` documentation, `ovs-dpctl(8)` for `ovs-dpctl dump-flows` documentation, and “Why are there so many different ways to dump flows?” above for some background.

- `ovs-appctl ofproto/trace` to observe the logic behind how `ovs-vswitchd` treats packets. See `ovs-vswitchd(8)` for documentation. You can get more details about a given flow that `ovs-dpctl dump-flows` displays, by cutting and pasting a flow from the output into an `ovs-appctl ofproto/trace` command.
- SPAN, RSPAN, and ERSPAN features of physical switches, to observe what goes on at these physical hops.

Starting at the origin of a given packet, observe the packet at each hop in turn. For example, in one plausible scenario, you might:

1. `tcpdump` the `eth` interface through which an ARP egresses a VM, from inside the VM.
2. `tcpdump` the `vif` or `tap` interface through which the ARP ingresses the host machine.
3. Use `ovs-dpctl dump-flows` to spot the ARP flow and observe the host interface through which the ARP egresses the physical machine. You may need to use `ovs-dpctl show` to interpret the port numbers. If the output seems surprising, you can use `ovs-appctl ofproto/trace` to observe details of how `ovs-vswitchd` determined the actions in the `ovs-dpctl dump-flows` output.
4. `tcpdump` the `eth` interface through which the ARP egresses the physical machine.
5. `tcpdump` the `eth` interface through which the ARP ingresses the physical machine, at the remote host that receives the ARP.
6. Use `ovs-dpctl dump-flows` to spot the ARP flow on the remote host remote host that receives the ARP and observe the VM `vif` or `tap` interface to which the flow is directed. Again, `ovs-dpctl show` and `ovs-appctl ofproto/trace` might help.
7. `tcpdump` the `vif` or `tap` interface to which the ARP is directed.
8. `tcpdump` the `eth` interface through which the ARP ingresses a VM, from inside the VM.

It is likely that during one of these steps you will figure out the problem. If not, then follow the ARP reply back to the origin, in reverse.

**Q:** How do I make a flow drop packets?

**A:** To drop a packet is to receive it without forwarding it. OpenFlow explicitly specifies forwarding actions. Thus, a flow with an empty set of actions does not forward packets anywhere, causing them to be dropped. You can specify an empty set of actions with `actions=` on the `ovs-ofctl` command line. For example:

```
$ ovs-ofctl add-flow br0 priority=65535,actions=
```



would cause every packet entering switch br0 to be dropped.

You can write “drop” explicitly if you like. The effect is the same. Thus, the following command also causes every packet entering switch br0 to be dropped:

```
$ ovs-ofctl add-flow br0 priority=65535,actions=drop
```

drop is not an action, either in OpenFlow or Open vSwitch. Rather, it is only a way to say that there are no actions.

Q: I added a flow to send packets out the ingress port, like this:

```
$ ovs-ofctl add-flow br0 in_port=2,actions=2
```

but OVS drops the packets instead.

A: Yes, OpenFlow requires a switch to ignore attempts to send a packet out its ingress port. The rationale is that dropping these packets makes it harder to loop the network. Sometimes this behavior can even be convenient, e.g. it is often the desired behavior in a flow that forwards a packet to several ports (“floods” the packet).

Sometimes one really needs to send a packet out its ingress port (“hairpin”). In this case, output to `OFPP_IN_PORT`, which in `ovs-ofctl` syntax is expressed as just `in_port`, e.g.:

```
$ ovs-ofctl add-flow br0 in_port=2,actions=in_port
```

This also works in some circumstances where the flow doesn’t match on the input port. For example, if you know that your switch has five ports numbered 2 through 6, then the following will send every received packet out every port, even its ingress port:

```
$ ovs-ofctl add-flow br0 actions=2,3,4,5,6,in_port
```

or, equivalently:

```
$ ovs-ofctl add-flow br0 actions=all,in_port
```

Sometimes, in complicated flow tables with multiple levels of `resubmit` actions, a flow needs to output to a particular port that may or may not be the ingress port. It’s difficult to take advantage of `OFPP_IN_PORT` in this situation. To help, Open vSwitch provides, as an OpenFlow extension, the ability to modify the `in_port` field. Whatever value is currently in the `in_port` field is the port to which outputs will be dropped, as well as the destination for `OFPP_IN_PORT`. This means that the following will reliably output to port 2 or to ports 2 through 6, respectively:

```
$ ovs-ofctl add-flow br0 in_port=2,actions=load:0->NXM_OF_IN_PORT[],2
$ ovs-ofctl add-flow br0 actions=load:0->NXM_OF_IN_PORT[],2,3,4,5,6
```

If the input port is important, then one may save and restore it on the stack:

```
$ ovs-ofctl add-flow br0 actions=push:NXM_OF_IN_PORT[], load:0-
>NXM_OF_IN_PORT[],2,3,4,5,6,pop:NXM_OF_IN_PORT[]
```

Q: My bridge br0 has host 192.168.0.1 on port 1 and host 192.168.0.2 on port 2. I set up flows to forward only traffic destined to the other host and drop other traffic, like this:

```
priority=5,in_port=1,ip,nw_dst=192.168.0.2,actions=2
priority=5,in_port=2,ip,nw_dst=192.168.0.1,actions=1
priority=0,actions=drop
```

But it doesn’t work—I don’t get any connectivity when I do this. Why?

A: These flows drop the ARP packets that IP hosts use to establish IP connectivity over Ethernet. To solve the problem, add flows to allow ARP to pass between the hosts:

```
priority=5,in_port=1,arp,actions=2
priority=5,in_port=2,arp,actions=1
```

This issue can manifest other ways, too. The following flows that match on Ethernet addresses instead of IP addresses will also drop ARP packets, because ARP requests are broadcast instead of being directed to a specific host:

```
priority=5,in_port=1,dl_dst=54:00:00:00:00:02,actions=2
priority=5,in_port=2,dl_dst=54:00:00:00:00:01,actions=1
priority=0,actions=drop
```

The solution already described above will also work in this case. It may be better to add flows to allow all multicast and broadcast traffic:

```
priority=5,in_port=1,dl_dst=01:00:00:00:00:00/01:00:00:00:00:00,actions=2
priority=5,in_port=2,dl_dst=01:00:00:00:00:00/01:00:00:00:00:00,actions=1
```

Q: My bridge disconnects from my controller on add-port/del-port.

A: Reconfiguring your bridge can change your bridge's datapath-id because Open vSwitch generates datapath-id from the MAC address of one of its ports. In that case, Open vSwitch disconnects from controllers because there's no graceful way to notify controllers about the change of datapath-id.

To avoid the behaviour, you can configure datapath-id manually.:

```
$ ovs-vsctl set bridge br0 other-config:datapath-id=0123456789abcdef
```

Q: My controller complains that OVS is not buffering packets. What's going on?

A: "Packet buffering" is an optional OpenFlow feature, and controllers should detect how many "buffers" an OpenFlow switch implements. It was recently noticed that OVS implementation of the buffering feature was not compliant to OpenFlow specifications. Rather than fix it and risk controller incompatibility, the buffering feature is removed as of OVS 2.7. Controllers are already expected to work properly in cases where the switch can not buffer packets, but sends full packets in "packet-in" messages instead, so this change should not affect existing users. After the change OVS always sends the `buffer_id` as `0xffffffff` in "packet-in" messages and will send an error response if any other value of this field is included in a "packet-out" or a "flow mod" sent by a controller.

Packet buffers have limited usefulness in any case. Table-miss packet-in messages most commonly pass the first packet in a microflow to the OpenFlow controller, which then sets up an OpenFlow flow that handles remaining traffic in the microflow without further controller intervention. In such a case, the packet that initiates the microflow is in practice usually small (certainly for TCP), which means that the switch sends the entire packet to the controller and the buffer only saves a small number of bytes in the reverse direction.

Q: How does OVS divide flows among buckets in an OpenFlow "select" group?

A: In Open vSwitch 2.3 and earlier, Open vSwitch used the destination Ethernet address to choose a bucket in a select group.

Open vSwitch 2.4 and later by default hashes the source and destination Ethernet address, VLAN ID, Ethernet type, IPv4/v6 source and destination address and protocol, and for TCP and SCTP only, the source and destination ports. The hash is "symmetric", meaning that exchanging source and destination addresses does not change the bucket selection.

Select groups in Open vSwitch 2.4 and later can be configured to use a different hash function, using a Netronome extension to the OpenFlow 1.5+ `group_mod` message. For more information, see

Documentation/group-selection-method-property.txt in the Open vSwitch source tree. (OpenFlow 1.5 support in Open vSwitch is still experimental.)

Q: I added a flow to accept packets on VLAN 123 and output them on VLAN 456, like so:

```
$ ovs-ofctl add-flow br0 dl_vlan=123,actions=output:1,mod_vlan_vid:456
```

but the packets are actually being output in VLAN 123. Why?

A: OpenFlow actions are executed in the order specified. Thus, the actions above first output the packet, then change its VLAN. Since the output occurs before changing the VLAN, the change in VLAN will have no visible effect.

To solve this and similar problems, order actions so that changes to headers happen before output, e.g.:

```
$ ovs-ofctl add-flow br0 dl_vlan=123,actions=mod_vlan_vid:456,output:1
```

See also the following question.

Q: I added a flow to a redirect packets for TCP port 80 to port 443, like so:

```
$ ovs-ofctl add-flow br0 tcp,tcp_dst=123,actions=mod_tp_dst:443
```

but the packets are getting dropped instead. Why?

A: This set of actions does change the TCP destination port to 443, but then it does nothing more. It doesn't, for example, say to continue to another flow table or to output the packet. Therefore, the packet is dropped.

To solve the problem, add an action that does something with the modified packet. For example:

```
$ ovs-ofctl add-flow br0 tcp,tcp_dst=123,actions=mod_tp_dst:443,normal
```

See also the preceding question.

Q: When using the “ct” action with FTP connections, it doesn't seem to matter if I set the “alg=ftp” parameter in the action. Is this required?

A: It is advisable to use this option. Some platforms may automatically detect and apply ALGs in the “ct” action regardless of the parameters you provide, however this is not consistent across all implementations. The `ovs-ofctl(8)` man pages contain further details in the description of the ALG parameter.

## 7.7 Quality of Service (QoS)

Q: Does OVS support Quality of Service (QoS)?

A: Yes. For traffic that egresses from a switch, OVS supports traffic shaping; for traffic that ingresses into a switch, OVS support policing. Policing is a simple form of quality-of-service that simply drops packets received in excess of the configured rate. Due to its simplicity, policing is usually less accurate and less effective than egress traffic shaping, which queues packets.

Keep in mind that ingress and egress are from the perspective of the switch. That means that egress shaping limits the rate at which traffic is allowed to transmit from a physical interface, but not the rate at which traffic will be received on a virtual machine's VIF. For ingress policing, the behavior is the opposite.

Q: How do I configure egress traffic shaping?

A: Suppose that you want to set up bridge br0 connected to physical Ethernet port eth0 (a 1 Gbps device) and virtual machine interfaces vif1.0 and vif2.0, and that you want to limit traffic from vif1.0 to eth0 to 10 Mbps and from vif2.0 to eth0 to 20 Mbps. Then, you could configure the bridge this way:

```
$ ovs-vsctl -- \
 add-br br0 -- \
 add-port br0 eth0 -- \
 add-port br0 vif1.0 -- set interface vif1.0 ofport_request=5 -- \
 add-port br0 vif2.0 -- set interface vif2.0 ofport_request=6 -- \
 set port eth0 qos=@newqos -- \
 --id=@newqos create qos type=linux-htb \
 other-config:max-rate=1000000000 \
 queues:123=@vif10queue \
 queues:234=@vif20queue -- \
 --id=@vif10queue create queue other-config:max-rate=10000000 -- \
 --id=@vif20queue create queue other-config:max-rate=20000000
```

At this point, bridge br0 is configured with the ports and eth0 is configured with the queues that you need for QoS, but nothing is actually directing packets from vif1.0 or vif2.0 to the queues that we have set up for them. That means that all of the packets to eth0 are going to the “default queue”, which is not what we want.

We use OpenFlow to direct packets from vif1.0 and vif2.0 to the queues reserved for them:

```
$ ovs-ofctl add-flow br0 in_port=5,actions=set_queue:123,normal
$ ovs-ofctl add-flow br0 in_port=6,actions=set_queue:234,normal
```

Each of the above flows matches on the input port, sets up the appropriate queue (123 for vif1.0, 234 for vif2.0), and then executes the “normal” action, which performs the same switching that Open vSwitch would have done without any OpenFlow flows being present. (We know that vif1.0 and vif2.0 have OpenFlow port numbers 5 and 6, respectively, because we set their ofport\_request columns above. If we had not done that, then we would have needed to find out their port numbers before setting up these flows.)

Now traffic going from vif1.0 or vif2.0 to eth0 should be rate-limited.

By the way, if you delete the bridge created by the above commands, with:

```
$ ovs-vsctl del-br br0
```

then that will leave one unreferenced QoS record and two unreferenced Queue records in the Open vSwitch database. One way to clear them out, assuming you don’t have other QoS or Queue records that you want to keep, is:

```
$ ovs-vsctl -- --all destroy QoS -- --all destroy Queue
```

If you do want to keep some QoS or Queue records, or the Open vSwitch you are using is older than version 1.8 (which added the `--all` option), then you will have to destroy QoS and Queue records individually.

**Q:** How do I configure ingress policing?

**A:** A policing policy can be configured on an interface to drop packets that arrive at a higher rate than the configured value. For example, the following commands will rate-limit traffic that vif1.0 may generate to 10Mbps:

```
$ ovs-vsctl set interface vif1.0 ingress_policing_rate=10000 $ ovs-vsctl set interface vif1.0
 ingress_policing_burst=8000
```

Traffic policing can interact poorly with some network protocols and can have surprising results. The “Ingress Policing” section of `ovs-vswitchd.conf.db(5)` discusses the issues in greater detail.

**Q:** I configured Quality of Service (QoS) in my OpenFlow network by adding records to the QoS and Queue table, but the results aren’t what I expect.

A: Did you install OpenFlow flows that use your queues? This is the primary way to tell Open vSwitch which queues you want to use. If you don't do this, then the default queue will be used, which will probably not have the effect you want.

Refer to the previous question for an example.

Q: I'd like to take advantage of some QoS feature that Open vSwitch doesn't yet support. How do I do that?

A: Open vSwitch does not implement QoS itself. Instead, it can configure some, but not all, of the QoS features built into the Linux kernel. If you need some QoS feature that OVS cannot configure itself, then the first step is to figure out whether Linux QoS supports that feature. If it does, then you can submit a patch to support Open vSwitch configuration for that feature, or you can use "tc" directly to configure the feature in Linux. (If Linux QoS doesn't support the feature you want, then first you have to add that support to Linux.)

Q: I configured QoS, correctly, but my measurements show that it isn't working as well as I expect.

A: With the Linux kernel, the Open vSwitch implementation of QoS has two aspects:

- Open vSwitch configures a subset of Linux kernel QoS features, according to what is in OVSDB. It is possible that this code has bugs. If you believe that this is so, then you can configure the Linux traffic control (QoS) stack directly with the "tc" program. If you get better results that way, you can send a detailed bug report to [bugs@openvswitch.org](mailto:bugs@openvswitch.org).

It is certain that Open vSwitch cannot configure every Linux kernel QoS feature. If you need some feature that OVS cannot configure, then you can also use "tc" directly (or add that feature to OVS).

- The Open vSwitch implementation of OpenFlow allows flows to be directed to particular queues. This is pretty simple and unlikely to have serious bugs at this point.

However, most problems with QoS on Linux are not bugs in Open vSwitch at all. They tend to be either configuration errors (please see the earlier questions in this section) or issues with the traffic control (QoS) stack in Linux. The Open vSwitch developers are not experts on Linux traffic control. We suggest that, if you believe you are encountering a problem with Linux traffic control, that you consult the tc manpages (e.g. `tc(8)`, `tc-htb(8)`, `tc-hfsc(8)`), web resources (e.g. <http://lartc.org/>), or mailing lists (e.g. <http://vger.kernel.org/vger-lists.html#netdev>).

Q: Does Open vSwitch support OpenFlow meters?

A: Open vSwitch 2.0 added OpenFlow protocol support for OpenFlow meters. Open vSwitch 2.7 implemented meters in the userspace datapath. Open vSwitch 2.10 implemented meters in the Linux kernel datapath.

## 7.8 Releases

Q: What does it mean for an Open vSwitch release to be LTS (long-term support)?

A: All official releases have been through a comprehensive testing process and are suitable for production use. Planned releases occur twice a year. If a significant bug is identified in an LTS release, we will provide an updated release that includes the fix. Releases that are not LTS may not be fixed and may just be supplanted by the next major release. The current LTS release is 2.5.x.

For more information on the Open vSwitch release process, refer to *Open vSwitch Release Process*.

Q: What Linux kernel versions does each Open vSwitch release work with?

A: The following table lists the Linux kernel versions against which the given versions of the Open vSwitch kernel module will successfully build. The Linux kernel versions are upstream kernel versions, so Linux kernels modified from the upstream sources may not build in some cases even if they are based

on a supported version. This is most notably true of Red Hat Enterprise Linux (RHEL) kernels, which are extensively modified from upstream.

Open vSwitch	Linux kernel
1.4.x	2.6.18 to 3.2
1.5.x	2.6.18 to 3.2
1.6.x	2.6.18 to 3.2
1.7.x	2.6.18 to 3.3
1.8.x	2.6.18 to 3.4
1.9.x	2.6.18 to 3.8
1.10.x	2.6.18 to 3.8
1.11.x	2.6.18 to 3.8
2.0.x	2.6.32 to 3.10
2.1.x	2.6.32 to 3.11
2.3.x	2.6.32 to 3.14
2.4.x	2.6.32 to 4.0
2.5.x	2.6.32 to 4.3
2.6.x	3.10 to 4.7
2.7.x	3.10 to 4.9
2.8.x	3.10 to 4.12
2.9.x	3.10 to 4.13
2.10.x	3.10 to 4.17
2.11.x	3.10 to 4.18

Open vSwitch userspace should also work with the Linux kernel module built into Linux 3.3 and later.

Open vSwitch userspace is not sensitive to the Linux kernel version. It should build against almost any kernel, certainly against 2.6.32 and later.

Q: Are all features available with all datapaths?

A: Open vSwitch supports different datapaths on different platforms. Each datapath has a different feature set: the following tables try to summarize the status.

Supported datapaths:

**Linux upstream** The datapath implemented by the kernel module shipped with Linux upstream. Since features have been gradually introduced into the kernel, the table mentions the first Linux release whose OVS module supports the feature.

**Linux OVS tree** The datapath implemented by the Linux kernel module distributed with the OVS source tree.

**Userspace** Also known as DPDK, dpif-netdev or dummy datapath. It is the only datapath that works on NetBSD, FreeBSD and Mac OSX.

**Hyper-V** Also known as the Windows datapath.

The following table lists the datapath supported features from an Open vSwitch user's perspective.

Feature	Linux upstream	Linux OVS tree	Userspace	Hyper-V
NAT	4.6	YES	Yes	NO
Connection tracking	4.3	YES	PARTIAL	PARTIAL
Tunnel - LISP	NO	YES	NO	NO
Tunnel - STT	NO	YES	NO	YES
Tunnel - GRE	3.11	YES	YES	YES
Tunnel - VXLAN	3.12	YES	YES	YES
Tunnel - Geneve	3.18	YES	YES	YES
Tunnel - GRE-IPv6	4.18	YES	YES	NO
Tunnel - VXLAN-IPv6	4.3	YES	YES	NO
Tunnel - Geneve-IPv6	4.4	YES	YES	NO
Tunnel - ERSPAN	4.18	YES	YES	NO
Tunnel - ERSPAN-IPv6	4.18	YES	YES	NO
QoS - Policing	YES	YES	YES	NO
QoS - Shaping	YES	YES	NO	NO
sFlow	YES	YES	YES	NO
IPFIX	3.10	YES	YES	NO
Set action	YES	YES	YES	PARTIAL
NIC Bonding	YES	YES	YES	YES
Multiple VTEPs	YES	YES	YES	YES
Meters	4.15	YES	YES	NO
Conntrack zone limit	4.18	YES	NO	NO

Do note, however:

- Only a limited set of flow fields is modifiable via the set action by the Hyper-V datapath.
- Userspace datapath support, in some cases, is dependent on the associated interface types. For example, DPDK interfaces support ingress and egress policing, but not shaping.

The following table lists features that do not *directly* impact an Open vSwitch user, e.g. because their absence can be hidden by the ofproto layer (usually this comes with a performance penalty).

Feature	Linux upstream	Linux OVS tree	Userspace	Hyper-V
SCTP flows	3.12	YES	YES	YES
MPLS	3.19	YES	YES	YES
UFID	4.0	YES	YES	NO
Megaflows	3.12	YES	YES	NO
Masked set action	4.0	YES	YES	NO
Recirculation	3.19	YES	YES	YES
TCP flags matching	3.13	YES	YES	NO
Validate flow actions	YES	YES	N/A	NO
Multiple datapaths	YES	YES	YES	NO
Tunnel TSO - STT	N/A	YES	NO	YES

Q: What DPDK version does each Open vSwitch release work with?

A: The following table lists the DPDK version against which the given versions of Open vSwitch will successfully build.

Open vSwitch	DPDK
2.2.x	1.6
2.3.x	1.6
2.4.x	2.0
2.5.x	2.2
2.6.x	16.07.2
2.7.x	16.11.8
2.8.x	17.05.2
2.9.x	17.11.4
2.10.x	17.11.4

Q: Are all the DPDK releases that OVS versions work with maintained?

No. DPDK follows YY.MM.n (Year.Month.Number) versioning.

Typically, all DPDK releases get a stable YY.MM.1 update with bugfixes 3 months after the YY.MM.0 release. In some cases there may also be a YY.MM.2 release.

DPDK LTS releases start once a year at YY.11.0 and are maintained for two years, with YY.MM.n+1 releases around every 3 months.

The latest information about DPDK stable and LTS releases can be found at [DPDK stable](#).

Q: I get an error like this when I configure Open vSwitch:

```
configure: error: Linux kernel in <dir> is version <x>, but version newer than <y> is not
supported (please refer to the FAQ for advice)
```

What should I do?

A: You have the following options:

- Use the Linux kernel module supplied with the kernel that you are using. (See also the following FAQ.)
- If there is a newer released version of Open vSwitch, consider building that one, because it may support the kernel that you are building against. (To find out, consult the table in the previous FAQ.)
- The Open vSwitch “master” branch may support the kernel that you are using, so consider building the kernel module from “master”.

All versions of Open vSwitch userspace are compatible with all versions of the Open vSwitch kernel module, so you do not have to use the kernel module from one source along with the userspace programs from the same source.

Q: What features are not available in the Open vSwitch kernel datapath that ships as part of the upstream Linux kernel?

A: The kernel module in upstream Linux does not include support for LISP. Work is in progress to add support for LISP to the upstream Linux version of the Open vSwitch kernel module. For now, if you need this feature, use the kernel module from the Open vSwitch distribution instead of the upstream Linux kernel module.

Certain features require kernel support to function or to have reasonable performance. If the ovs-vsitchd log file indicates that a feature is not supported, consider upgrading to a newer upstream Linux release or using the kernel module paired with the userspace distribution.

Q: Why do tunnels not work when using a kernel module other than the one packaged with Open vSwitch?

A: Support for tunnels was added to the upstream Linux kernel module after the rest of Open vSwitch. As a result, some kernels may contain support for Open vSwitch but not tunnels. The minimum kernel version that supports each tunnel protocol is:



Protocol	Linux Kernel
GRE	3.11
VXLAN	3.12
Geneve	3.18
ERSPAN	4.18
LISP	not upstream
STT	not upstream

If you are using a version of the kernel that is older than the one listed above, it is still possible to use that tunnel protocol. However, you must compile and install the kernel module included with the Open vSwitch distribution rather than the one on your machine. If problems persist after doing this, check to make sure that the module that is loaded is the one you expect.

Q: Why are UDP tunnel checksums not computed for VXLAN or Geneve?

A: Generating outer UDP checksums requires kernel support that was not part of the initial implementation of these protocols. If using the upstream Linux Open vSwitch module, you must use kernel 4.0 or newer. The out-of-tree modules from Open vSwitch release 2.4 and later support UDP checksums.

Q: What features are not available when using the userspace datapath?

A: Tunnel virtual ports are not supported, as described in the previous answer. It is also not possible to use queue-related actions. On Linux kernels before 2.6.39, maximum-sized VLAN packets may not be transmitted.

Q: Should userspace or kernel be upgraded first to minimize downtime?

A: In general, the Open vSwitch userspace should be used with the kernel version included in the same release or with the version from upstream Linux. However, when upgrading between two releases of Open vSwitch it is best to migrate userspace first to reduce the possibility of incompatibilities.

Q: What happened to the bridge compatibility feature?

A: Bridge compatibility was a feature of Open vSwitch 1.9 and earlier. When it was enabled, Open vSwitch imitated the interface of the Linux kernel “bridge” module. This allowed users to drop Open vSwitch into environments designed to use the Linux kernel bridge module without adapting the environment to use Open vSwitch.

Open vSwitch 1.10 and later do not support bridge compatibility. The feature was dropped because version 1.10 adopted a new internal architecture that made bridge compatibility difficult to maintain. Now that many environments use OVS directly, it would be rarely useful in any case.

To use bridge compatibility, install OVS 1.9 or earlier, including the accompanying kernel modules (both the main and bridge compatibility modules), following the instructions that come with the release. Be sure to start the `ovs-brcompatd` daemon.

## 7.9 Terminology

Q: I thought Open vSwitch was a virtual Ethernet switch, but the documentation keeps talking about bridges. What’s a bridge?

A: In networking, the terms “bridge” and “switch” are synonyms. Open vSwitch implements an Ethernet switch, which means that it is also an Ethernet bridge.

Q: What’s a VLAN?

A: See [VLANs](#).

## 7.10 VLANs

Q: What's a VLAN?

A: At the simplest level, a VLAN (short for “virtual LAN”) is a way to partition a single switch into multiple switches. Suppose, for example, that you have two groups of machines, group A and group B. You want the machines in group A to be able to talk to each other, and you want the machine in group B to be able to talk to each other, but you don't want the machines in group A to be able to talk to the machines in group B. You can do this with two switches, by plugging the machines in group A into one switch and the machines in group B into the other switch.

If you only have one switch, then you can use VLANs to do the same thing, by configuring the ports for machines in group A as VLAN “access ports” for one VLAN and the ports for group B as “access ports” for a different VLAN. The switch will only forward packets between ports that are assigned to the same VLAN, so this effectively subdivides your single switch into two independent switches, one for each group of machines.

So far we haven't said anything about VLAN headers. With access ports, like we've described so far, no VLAN header is present in the Ethernet frame. This means that the machines (or switches) connected to access ports need not be aware that VLANs are involved, just like in the case where we use two different physical switches.

Now suppose that you have a whole bunch of switches in your network, instead of just one, and that some machines in group A are connected directly to both switches 1 and 2. To allow these machines to talk to each other, you could add an access port for group A's VLAN to switch 1 and another to switch 2, and then connect an Ethernet cable between those ports. That works fine, but it doesn't scale well as the number of switches and the number of VLANs increases, because you use up a lot of valuable switch ports just connecting together your VLANs.

This is where VLAN headers come in. Instead of using one cable and two ports per VLAN to connect a pair of switches, we configure a port on each switch as a VLAN “trunk port”. Packets sent and received on a trunk port carry a VLAN header that says what VLAN the packet belongs to, so that only two ports total are required to connect the switches, regardless of the number of VLANs in use. Normally, only switches (either physical or virtual) are connected to a trunk port, not individual hosts, because individual hosts don't expect to see a VLAN header in the traffic that they receive.

None of the above discussion says anything about particular VLAN numbers. This is because VLAN numbers are completely arbitrary. One must only ensure that a given VLAN is numbered consistently throughout a network and that different VLANs are given different numbers. (That said, VLAN 0 is usually synonymous with a packet that has no VLAN header, and VLAN 4095 is reserved.)

Q: VLANs don't work.

A: Many drivers in Linux kernels before version 3.3 had VLAN-related bugs. If you are having problems with VLANs that you suspect to be driver related, then you have several options:

- Upgrade to Linux 3.3 or later.
- Build and install a fixed version of the particular driver that is causing trouble, if one is available.
- Use a NIC whose driver does not have VLAN problems.
- Use “VLAN splinters”, a feature in Open vSwitch 1.4 upto 2.5 that works around bugs in kernel drivers. To enable VLAN splinters on interface eth0, use the command:

```
$ ovs-vsctl set interface eth0 other-config:enable-vlan-splinters=true
```

For VLAN splinters to be effective, Open vSwitch must know which VLANs are in use. See the “VLAN splinters” section in the Interface table in `ovs-vswitchd.conf.db(5)` for details on how Open vSwitch infers in-use VLANs.

VLAN splinters increase memory use and reduce performance, so use them only if needed.

- Apply the “vlan workaround” patch from the XenServer kernel patch queue, build Open vSwitch against this patched kernel, and then use `ovs-vlan-bug-workaround(8)` to enable the VLAN workaround for each interface whose driver is buggy.

(This is a nontrivial exercise, so this option is included only for completeness.)

It is not always easy to tell whether a Linux kernel driver has buggy VLAN support. The `ovs-vlan-test(8)` and `ovs-test(8)` utilities can help you test. See their manpages for details. Of the two utilities, `ovs-test(8)` is newer and more thorough, but `ovs-vlan-test(8)` may be easier to use.

Q: VLANs still don’t work. I’ve tested the driver so I know that it’s OK.

A: Do you have VLANs enabled on the physical switch that OVS is attached to? Make sure that the port is configured to trunk the VLAN or VLANs that you are using with OVS.

Q: Outgoing VLAN-tagged traffic goes through OVS to my physical switch and to its destination host, but OVS seems to drop incoming return traffic.

A: It’s possible that you have the VLAN configured on your physical switch as the “native” VLAN. In this mode, the switch treats incoming packets either tagged with the native VLAN or untagged as part of the native VLAN. It may also send outgoing packets in the native VLAN without a VLAN tag.

If this is the case, you have two choices:

- Change the physical switch port configuration to tag packets it forwards to OVS with the native VLAN instead of forwarding them untagged.
- Change the OVS configuration for the physical port to a native VLAN mode. For example, the following sets up a bridge with port `eth0` in “native-tagged” mode in VLAN 9:

```
$ ovs-vsctl add-br br0 $ ovs-vsctl add-port br0 eth0 tag=9
vlan_mode=native-tagged
```

In this situation, “native-untagged” mode will probably work equally well. Refer to the documentation for the Port table in `ovs-vswitchd.conf.db(5)` for more information.

Q: I added a pair of VMs on different VLANs, like this:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0 tag=9
$ ovs-vsctl add-port br0 tap1 tag=10
```

but the VMs can’t access each other, the external network, or the Internet.

A: It is to be expected that the VMs can’t access each other. VLANs are a means to partition a network. When you configured `tap0` and `tap1` as access ports for different VLANs, you indicated that they should be isolated from each other.

As for the external network and the Internet, it seems likely that the machines you are trying to access are not on VLAN 9 (or 10) and that the Internet is not available on VLAN 9 (or 10).

Q: I added a pair of VMs on the same VLAN, like this:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0 tag=9
$ ovs-vsctl add-port br0 tap1 tag=9
```

The VMs can access each other, but not the external network or the Internet.

A: It seems likely that the machines you are trying to access in the external network are not on VLAN 9 and that the Internet is not available on VLAN 9. Also, ensure VLAN 9 is set up as an allowed trunk VLAN on the upstream switch port to which eth0 is connected.

Q: Can I configure an IP address on a VLAN?

A: Yes. Use an “internal port” configured as an access port. For example, the following configures IP address 192.168.0.7 on VLAN 9. That is, OVS will forward packets from eth0 to 192.168.0.7 only if they have an 802.1Q header with VLAN 9. Conversely, traffic forwarded from 192.168.0.7 to eth0 will be tagged with an 802.1Q header with VLAN 9:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 vlan9 tag=9 \
 -- set interface vlan9 type=internal
$ ip addr add 192.168.0.7/24 dev vlan9
$ ip link set vlan9 up
```

See also the following question.

Q: I configured one IP address on VLAN 0 and another on VLAN 9, like this:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 eth0
$ ip addr add 192.168.0.5/24 dev br0
$ ip link set br0 up
$ ovs-vsctl add-port br0 vlan9 tag=9 -- set interface vlan9 type=internal
$ ip addr add 192.168.0.9/24 dev vlan9
$ ip link set vlan9 up
```

but other hosts that are only on VLAN 0 can reach the IP address configured on VLAN 9. What’s going on?

A: [RFC 1122 section 3.3.4.2 “Multihoming Requirements”](#) describes two approaches to IP address handling in Internet hosts:

- In the “Strong ES Model”, where an ES is a host (“End System”), an IP address is primarily associated with a particular interface. The host discards packets that arrive on interface A if they are destined for an IP address that is configured on interface B. The host never sends packets from interface A using a source address configured on interface B.
- In the “Weak ES Model”, an IP address is primarily associated with a host. The host accepts packets that arrive on any interface if they are destined for any of the host’s IP addresses, even if the address is configured on some interface other than the one on which it arrived. The host does not restrict itself to sending packets from an IP address associated with the originating interface.

Linux uses the weak ES model. That means that when packets destined to the VLAN 9 IP address arrive on eth0 and are bridged to br0, the kernel IP stack accepts them there for the VLAN 9 IP address, even though they were not received on vlan9, the network device for vlan9.

To simulate the strong ES model on Linux, one may add iptables rule to filter packets based on source and destination address and adjust ARP configuration with sysctls.

BSD uses the strong ES model.

Q: My OpenFlow controller doesn’t see the VLANs that I expect.

A: The configuration for VLANs in the Open vSwitch database (e.g. via ovs-vsctl) only affects traffic that goes through Open vSwitch’s implementation of the OpenFlow “normal switching” action. By default, when Open vSwitch isn’t connected to a controller and nothing has been manually configured in the flow table, all traffic goes through the “normal switching” action. But, if you set up OpenFlow flows on your

own, through a controller or using `ovs-ofctl` or through other means, then you have to implement VLAN handling yourself.

You can use “normal switching” as a component of your OpenFlow actions, e.g. by putting “normal” into the lists of actions on `ovs-ofctl` or by outputting to `OFPP_NORMAL` from an OpenFlow controller. In situations where this is not suitable, you can implement VLAN handling yourself, e.g.:

- If a packet comes in on an access port, and the flow table needs to send it out on a trunk port, then the flow can add the appropriate VLAN tag with the “`mod_vlan_vid`” action.
- If a packet comes in on a trunk port, and the flow table needs to send it out on an access port, then the flow can strip the VLAN tag with the “`strip_vlan`” action.

Q: I configured ports on a bridge as access ports with different VLAN tags, like this:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl set-controller br0 tcp:192.168.0.10:6653
$ ovs-vsctl add-port br0 eth0
$ ovs-vsctl add-port br0 tap0 tag=9
$ ovs-vsctl add-port br0 tap1 tag=10
```

but the VMs running behind `tap0` and `tap1` can still communicate, that is, they are not isolated from each other even though they are on different VLANs.

A: Do you have a controller configured on `br0` (as the commands above do)? If so, then this is a variant on the previous question, “My OpenFlow controller doesn’t see the VLANs that I expect,” and you can refer to the answer there for more information.

Q: How MAC learning works with VLANs?

A: Open vSwitch implements Independent VLAN Learning (IVL) for `OFPP_NORMAL` action, e.g. it logically has separate learning tables for each VLANs.

## 7.11 VXLANs

Q: What’s a VXLAN?

A: VXLAN stands for Virtual eXtensible Local Area Network, and is a means to solve the scaling challenges of VLAN networks in a multi-tenant environment. VXLAN is an overlay network which transports an L2 network over an existing L3 network. For more information on VXLAN, please see [RFC 7348](#).

Q: How much of the VXLAN protocol does Open vSwitch currently support?

A: Open vSwitch currently supports the framing format for packets on the wire. There is currently no support for the multicast aspects of VXLAN. To get around the lack of multicast support, it is possible to pre-provision MAC to IP address mappings either manually or from a controller.

Q: What destination UDP port does the VXLAN implementation in Open vSwitch use?

A: By default, Open vSwitch will use the assigned IANA port for VXLAN, which is 4789. However, it is possible to configure the destination UDP port manually on a per-VXLAN tunnel basis. An example of this configuration is provided below.:

```
$ ovs-vsctl add-br br0
$ ovs-vsctl add-port br0 vxlan1 -- set interface vxlan1 type=vxlan \
 options:remote_ip=192.168.1.2 options:key=flow options:dst_port=8472
```

## 7.12 OVN

Q: Why does OVN use STT and Geneve instead of VLANs or VXLAN (or GRE)?

A: OVN implements a fairly sophisticated packet processing pipeline in “logical datapaths” that can implement switching or routing functionality. A logical datapath has an ingress pipeline and an egress pipeline, and each of these pipelines can include logic based on packet fields as well as packet metadata such as the logical ingress and egress ports (the latter only in the egress pipeline).

The processing for a logical datapath can be split across hypervisors. In particular, when a logical ingress pipeline executes an “output” action, OVN passes the packet to the egress pipeline on the hypervisor (or, in the case of output to a logical multicast group, hypervisors) on which the logical egress port is located. If this hypervisor is not the same as the ingress hypervisor, then the packet has to be transmitted across a physical network.

This situation is where tunneling comes in. To send the packet to another hypervisor, OVN encapsulates it with a tunnel protocol and sends the encapsulated packet across the physical network. When the remote hypervisor receives the tunnel packet, it decapsulates it and passes it through the logical egress pipeline. To do so, it also needs the metadata, that is, the logical ingress and egress ports.

Thus, to implement OVN logical packet processing, at least the following metadata must pass across the physical network:

- Logical datapath ID, a 24-bit identifier. In Geneve, OVN uses the VNI to hold the logical datapath ID; in STT, OVN uses 24 bits of STT’s 64-bit context ID.
- Logical ingress port, a 15-bit identifier. In Geneve, OVN uses an option to hold the logical ingress port; in STT, 15 bits of the context ID.
- Logical egress port, a 16-bit identifier. In Geneve, OVN uses an option to hold the logical egress port; in STT, 16 bits of the context ID.

See `ovn-architecture(7)`, under “Tunnel Encapsulations”, for details.

Together, these metadata require  $24 + 15 + 16 = 55$  bits. GRE provides 32 bits, VXLAN provides 24, and VLAN only provides 12. Most notably, if logical egress pipelines do not match on the logical ingress port, thereby restricting the class of ACLs available to users, then this eliminates 15 bits, bringing the requirement down to 40 bits. At this point, one can choose to limit the size of the OVN logical network in various ways, e.g.:

- 16 bits of logical datapaths + 16 bits of logical egress ports. This combination fits within a 32-bit GRE tunnel key.
- 12 bits of logical datapaths + 12 bits of logical egress ports. This combination fits within a 24-bit VXLAN VNI.
- It’s difficult to identify an acceptable compromise for a VLAN-based deployment.

These compromises wouldn’t suit every site, since some deployments may need to allocate more bits to the datapath or egress port identifiers.

As a side note, OVN does support VXLAN for use with ASIC-based top of rack switches, using `ovn-controller-vtep(8)` and the OVSDB VTEP schema described in `vtep(5)`, but this limits the features available from OVN to the subset available from the VTEP schema.

---

## Open vSwitch Internals

---

Information for people who want to know more about the Open vSwitch project itself and how they might be involved.

### 8.1 Contributing to Open vSwitch

The below guides provide information on contributing to Open vSwitch itself.

#### 8.1.1 Submitting Patches

Send changes to Open vSwitch as patches to [dev@openvswitch.org](mailto:dev@openvswitch.org). One patch per email. More details are included below.

If you are using Git, then *git format-patch* takes care of most of the mechanics described below for you.

##### Before You Start

Before you send patches at all, make sure that each patch makes sense. In particular:

- A given patch should not break anything, even if later patches fix the problems that it causes. The source tree should still build and work after each patch is applied. (This enables *git bisect* to work best.)
- A patch should make one logical change. Don't make multiple, logically unconnected changes to disparate subsystems in a single patch.
- A patch that adds or removes user-visible features should also update the appropriate user documentation or manpages. Consider adding an item to NEWS for nontrivial changes. Check "Feature Deprecation Guidelines" section in this document if you intend to remove user-visible feature.

Testing is also important:

- Test a patch that modifies existing code with `make check` before submission. Refer to the "Unit Tests" in *Testing*, for more information. We also encourage running the kernel and userspace system tests.

- Consider testing a patch that adds or deletes files with `make distcheck` before submission.
- A patch that modifies Linux kernel code should be at least build-tested on various Linux kernel versions before submission. I suggest versions 3.10 and whatever the current latest release version is at the time.
- A patch that adds a new feature should add appropriate tests for the feature. A bug fix patch should preferably add a test that would fail if the bug recurs.

If you are using GitHub, then you may utilize the [travis-ci.org](https://travis-ci.org) CI build system by linking your GitHub repository to it. This will run some of the above tests automatically when you push changes to your repository. See the “Continuous Integration with Travis-CI” in [Testing](#) for details on how to set it up.

### Email Subject

The subject line of your email should be in the following format:

[PATCH <n>/<m>] <area>: <summary>

Where:

**[PATCH <n>/<m>]:** indicates that this is the nth of a series of m patches. It helps reviewers to read patches in the correct order. You may omit this prefix if you are sending only one patch.

**<area>:** indicates the area of the Open vSwitch to which the change applies (often the name of a source file or a directory). You may omit it if the change crosses multiple distinct pieces of code.

**<summary>:**

briefly describes the change. Use the imperative form, e.g. “Force SNAT for multiple gateway routers.” or “Fix daemon exit for bad datapaths or flows.” Try to keep the summary short, about 50 characters wide.

The subject, minus the [PATCH <n>/<m>] prefix, becomes the first line of the commit’s change log message.

### Description

The body of the email should start with a more thorough description of the change. This becomes the body of the commit message, following the subject. There is no need to duplicate the summary given in the subject.

Please limit lines in the description to 75 characters in width. That allows the description to format properly even when indented (e.g. by “git log” or in email quotations).

The description should include:

- The rationale for the change.
- Design description and rationale (but this might be better added as code comments).
- Testing that you performed (or testing that should be done but you could not for whatever reason).
- Tags (see below).

There is no need to describe what the patch actually changed, if the reader can see it for himself.

If the patch refers to a commit already in the Open vSwitch repository, please include both the commit number and the subject of the patch, e.g. ‘commit 632d136c (vswitch: Remove restriction on datapath names.)’.

If you, the person sending the patch, did not write the patch yourself, then the very first line of the body should take the form **From:** <author name> <author email>, followed by a blank line. This will automatically cause the named author to be credited with authorship in the repository.



## Tags

The description ends with a series of tags, written one to a line as the last paragraph of the email. Each tag indicates some property of the patch in an easily machine-parseable manner.

Please don't wrap a tag across multiple lines. If necessary, it's OK to have a tag extend beyond the customary maximum width of a commit message.

Examples of common tags follow.

Signed-off-by: Author Name <author.name@email.address...>

Informally, this indicates that Author Name is the author or submitter of a patch and has the authority to submit it under the terms of the license. The formal meaning is to agree to the Developer's Certificate of Origin (see below).

If the author and submitter are different, each must sign off. If the patch has more than one author, all must sign off.

Signed-off-by tags should be the last tags in the commit message. If the author (or authors) and submitter are different, the author tags should come first. More generally, occasionally a patch might pass through a chain of submitters, and in such a case the sign-offs should be arranged in chronological order.

```
Signed-off-by: Author Name <author.name@email.address...>
Signed-off-by: Submitter Name <submitter.name@email.address...>
```

Co-authored-by: Author Name <author.name@email.address...>

Git can only record a single person as the author of a given patch. In the rare event that a patch has multiple authors, one must be given the credit in Git and the others must be credited via Co-authored-by: tags. (All co-authors must also sign off.)

Acked-by: Reviewer Name <reviewer.name@email.address...>

Reviewers will often give an Acked-by: tag to code of which they approve. It is polite for the submitter to add the tag before posting the next version of the patch or applying the patch to the repository. Quality reviewing is hard work, so this gives a small amount of credit to the reviewer.

Not all reviewers give Acked-by: tags when they provide positive reviews. It's customary only to add tags from reviewers who actually provide them explicitly.

Tested-by: Tester Name <reviewer.name@email.address...>

When someone tests a patch, it is customary to add a Tested-by: tag indicating that. It's rare for a tester to actually provide the tag; usually the patch submitter makes the tag himself in response to an email indicating successful testing results.

Tested-at: <URL>

When a test report is publicly available, this provides a way to reference it. Typical <URL>s would be build logs from autobuilders or references to mailing list archives.

Some autobuilders only retain their logs for a limited amount of time. It is less useful to cite these because they may be dead links for a developer reading the commit message months or years later.

Reported-by: Reporter Name <reporter.name@email.address...>

When a patch fixes a bug reported by some person, please credit the reporter in the commit log in this fashion. Please also add the reporter's name and email address to the list of people who provided helpful bug reports in the AUTHORS file at the top of the source tree.

Fairly often, the reporter of a bug also tests the fix. Occasionally one sees a combined “Reported-and-tested-by:” tag used to indicate this. It is also acceptable, and more common, to include both tags separately.

(If a bug report is received privately, it might not always be appropriate to publicly credit the reporter. If in doubt, please ask the reporter.)

Requested-by: Requester Name <requester.name@email.address...>

When a patch implements a request or a suggestion made by some person, please credit that person in the commit log in this fashion. For a helpful suggestion, please also add the person’s name and email address to the list of people who provided suggestions in the AUTHORS file at the top of the source tree.

(If a suggestion or a request is received privately, it might not always be appropriate to publicly give credit. If in doubt, please ask.)

Suggested-by: Suggester Name <suggester.name@email.address...>

See Requested-by:.

CC: Person <name@email>

This is a way to tag a patch for the attention of a person when no more specific tag is appropriate. One use is to request a review from a particular person. It doesn’t make sense to include the same person in CC and another tag, so e.g. if someone who is CCed later provides an Acked-by, add the Acked-by and remove the CC at the same time.

Reported-at: <URL>

If a patch fixes or is otherwise related to a bug reported in a public bug tracker, please include a reference to the bug in the form of a URL to the specific bug, e.g.:

Reported-at: <https://bugs.debian.org/743635>

This is also an appropriate way to refer to bug report emails in public email archives, e.g.:

Reported-at: [https://mail.openvswitch.org/pipermail/ovs-dev/2014-June/284495.  
↪html](https://mail.openvswitch.org/pipermail/ovs-dev/2014-June/284495.html)

Submitted-at: <URL>

If a patch was submitted somewhere other than the Open vSwitch development mailing list, such as a GitHub pull request, this header can be used to reference the source.

Submitted-at: <https://github.com/openvswitch/ovs/pull/92>

VMware-BZ: #1234567

If a patch fixes or is otherwise related to a bug reported in a private bug tracker, you may include some tracking ID for the bug for your own reference. Please include some identifier to make the origin clear, e.g. “VMware-BZ” refers to VMware’s internal Bugzilla instance and “ONF-JIRA” refers to the Open Networking Foundation’s JIRA bug tracker.

ONF-JIRA: EXT-12345

See VMware-BZ:.

Bug #1234567.

These are obsolete forms of VMware-BZ: that can still be seen in old change log entries. (They are obsolete because they do not tell the reader what bug tracker is referred to.)

Issue: 1234567

See Bug:.

Fixes: 63bc9fb1c69f ("packets: Reorder CS\_\* flags to remove gap.")

If you would like to record which commit introduced a bug being fixed, you may do that with a “Fixes” header. This assists in determining which OVS releases have the bug, so the patch can be applied to all affected versions. The easiest way to generate the header in the proper format is with this git command. This command also CCs the author of the commit being fixed, which makes sense unless the author also made the fix or is already named in another tag:

```
$ git log -1 --pretty=format:"CC: %an <%ae>%nFixes: %h (%s)" \
--abbrev=12 COMMIT_REF
```

Vulnerability: CVE-2016-2074

Specifies that the patch fixes or is otherwise related to a security vulnerability with the given CVE identifier. Other identifiers in public vulnerability databases are also suitable.

If the vulnerability was reported publicly, then it is also appropriate to cite the URL to the report in a Reported-at tag. Use a Reported-by tag to acknowledge the reporters.

## Developer's Certificate of Origin

To help track the author of a patch as well as the submission chain, and be clear that the developer has authority to submit a patch for inclusion in Open vSwitch please sign off your work. The sign off certifies the following:

Developer's Certificate of Origin 1.1

By making a contribution to this project, I certify that:

- (a) The contribution was created **in** whole **or in** part by me **and** I have the right to submit it under the **open** source license indicated **in** the file; **or**
- (b) The contribution **is** based upon previous work that, to the best of my knowledge, **is** covered under an appropriate **open** source license **and** I have the right under that license to submit that work **with** modifications, whether created **in** whole **or in** part by me, under the same **open** source license (unless I am permitted to submit under a different license), **as** indicated **in** the file; **or**
- (c) The contribution was provided directly to me by some other person who certified (a), (b) **or** (c) **and** I have **not** modified it.
- (d) I understand **and** agree that this project **and** the contribution are public **and** that a record of the contribution (including **all** personal information I submit **with** it, including my sign-off) **is** maintained indefinitely **and** may be redistributed consistent **with** this project **or** the **open** source license(s) involved.

See also <http://developercertificate.org/>.

### Feature Deprecation Guidelines

Open vSwitch is intended to be user friendly. This means that under normal circumstances we don't abruptly remove features from OVS that some users might still be using. Otherwise, if we would, then we would possibly break our user setup when they upgrade and would receive bug reports.

Typical process to deprecate a feature in Open vSwitch is to:

1. Mention deprecation of a feature in the NEWS file. Also, mention expected release or absolute time when this feature would be removed from OVS altogether. Don't use relative time (e.g. "in 6 months") because that is not clearly interpretable.
2. If Open vSwitch is configured to use deprecated feature it should print a warning message to the log files clearly indicating that feature is deprecated and that use of it should be avoided.
3. If this feature is mentioned in man pages, then add "Deprecated" keyword to it.

Also, if there is alternative feature to the one that is about to be marked as deprecated, then mention it in (a), (b) and (c) as well.

Remember to follow-up and actually remove the feature from OVS codebase once deprecation grace period has expired and users had opportunity to use at least one OVS release that would have informed them about feature deprecation!

### Comments

If you want to include any comments in your email that should not be part of the commit's change log message, put them after the description, separated by a line that contains just ---. It may be helpful to include a diffstat here for changes that touch multiple files.

### Patch

The patch should be in the body of the email following the description, separated by a blank line.

Patches should be in `diff -up` format. We recommend that you use Git to produce your patches, in which case you should use the `-M -C` options to `git diff` (or other Git tools) if your patch renames or copies files. [Quilt](#) might be useful if you do not want to use Git.

Patches should be inline in the email message. Some email clients corrupt white space or wrap lines in patches. There are hints on how to configure many email clients to avoid this problem on [kernel.org](#). If you cannot convince your email client not to mangle patches, then sending the patch as an attachment is a second choice.

Follow the style used in the code that you are modifying. [Open vSwitch Coding Style](#) file describes the coding style used in most of Open vSwitch. Use Linux kernel coding style for Linux kernel code.

If your code is non-datapath code, you may use the `utilities/checkpatch.py` utility as a quick check for certain commonly occurring mistakes (improper leading/trailing whitespace, missing signoffs, some improper formatted patch files). For Linux datapath code, it is a good idea to use the Linux script `checkpatch.pl`.

### Example

```
From fa29a1c2c17682879e79a21bb0cdd5bbe67fa7c0 Mon Sep 17 00:00:00 2001
From: Jesse Gross <jesse@nicira.com>
Date: Thu, 8 Dec 2011 13:17:24 -0800
Subject: [PATCH] datapath: Alphabetize include/net/ipv6.h compat header.

Signed-off-by: Jesse Gross <jesse@nicira.com>
```

(continues on next page)

(continued from previous page)

```

datapath/linux/Modules.mk | 2 +-
1 files changed, 1 insertions(+), 1 deletions(-)

diff --git a/datapath/linux/Modules.mk b/datapath/linux/Modules.mk
index fdd952e..f6cb88e 100644
--- a/datapath/linux/Modules.mk
+++ b/datapath/linux/Modules.mk
@@ -56,11 +56,11 @@ openvswitch_headers += \
 linux/compat/include/net/dst.h \
 linux/compat/include/net/genetlink.h \
 linux/compat/include/net/ip.h \
+ linux/compat/include/net/ipv6.h \
 linux/compat/include/net/net_namespace.h \
 linux/compat/include/net/netlink.h \
 linux/compat/include/net/protocol.h \
 linux/compat/include/net/route.h \
- linux/compat/include/net/ipv6.h \
 linux/compat/genetlink.inc

both_modules += brcompat
--
1.7.7.3

```

## 8.1.2 Backporting patches

**Note:** This is an advanced topic for developers and maintainers. Readers should familiarize themselves with building and running Open vSwitch, with the git tool, and with the Open vSwitch patch submission process.

The backporting of patches from one git tree to another takes multiple forms within Open vSwitch, but is broadly applied in the following fashion:

- Contributors submit their proposed changes to the latest development branch
- Contributors and maintainers provide feedback on the patches
- When the change is satisfactory, maintainers apply the patch to the development branch.
- Maintainers backport changes from a development branch to release branches.

With regards to Open vSwitch user space code and code that does not comprise the Linux datapath and compat code, the development branch is *master* in the Open vSwitch repository. Patches are applied first to this branch, then to the most recent *branch-X.Y*, then earlier *branch-X.Z*, and so on. The most common kind of patch in this category is a bugfix which affects master and other branches.

For Linux datapath code, the primary development branch is in the *net-next* tree as described in the section below, and patch discussion occurs on the *netdev* mailing list. Patches are first applied to the upstream branch by the networking maintainer, then the contributor backports the patch to the Open vSwitch *master* development branch. Patches in this category may include features which have been applied upstream, or bugfixes to the Open vSwitch datapath code. For bugfixes, the patches subsequently follow the regular Open vSwitch process as described above to reach older branches.

### Changes to userspace components

Patches which are fixing bugs should be considered for backporting from *master* to release branches. Open vSwitch contributors submit their patches targeted to the *master* branch, using the `Fixes` tag described in [Submitting Patches](#). The maintainer first applies the patch to *master*, then backports the patch to each older affected tree, as far back as it goes or at least to all currently supported branches. This is usually each branch back to the most recent LTS release branch.

If the fix only affects a particular branch and not *master*, contributors should submit the change with the target branch listed in the subject line of the patch. Contributors should list all versions that the bug affects. The `git format-patch` argument `--subject-prefix` may be used when posting the patch, for example:

```
$ git format-patch HEAD --subject-prefix="PATCH branch-2.7"
```

If a maintainer is backporting a change to older branches and the backport is not a trivial cherry-pick, then the maintainer may opt to submit the backport for the older branch on the mailing list for further review. This should be done in the same manner as described above.

### Changes to Linux kernel components

The Linux kernel components in Open vSwitch go through initial review in the upstream Linux netdev community before they go into the Open vSwitch tree. As such, backports from upstream to the Open vSwitch tree may include bugfixes or new features. The [netdev-FAQ](#) describes the general process for merging patches to the upstream Linux tree.

To keep track of the changes which are made upstream against the changes which have been backported to the Open vSwitch tree, backports should be done in the order that they are applied to the upstream *net-next* tree. For example, if the git history in `linux/net/openvswitch/` in the *net-next* tree lists patches A, B and C that were applied (in that order), then the backports of these patches to `openvswitch/datapath/` should be done submitted in the order A, B, then C.

Patches that are proposed against the Open vSwitch tree, including backports, should follow the guidelines described in [Submitting Patches](#). Ideally, a series which backports new functionality would also include a series of patches for the userspace components which show how to use the new functionality, and include tests to validate the behaviour. However, in the interests of keeping the Open vSwitch tree in sync with upstream *net-next*, contributors may send Open vSwitch kernel module changes independently of userspace changes.

### How to backport kernel patches

First, the patch should be submitted upstream to *netdev*. When the patch has been applied to *net-next*, it is ready to be backported. Starting from the Linux tree, use `git format-patch` to format each patch that should be backported. For each of these patches, they may only include changes to `linux/net/openvswitch/`, or they may include changes to other directories. Depending on which files the patch touches, the backport may be easier or more difficult to undertake.

Start by formatting the relevant patches from the Linux tree. For example, to format the last 5 patches to `net/openvswitch`, going back from OVS commit `1234c0ffee5`, placing them into `/tmp/`:

```
$ git format-patch -5 1234c0ffee5 net/openvswitch/ -o /tmp/
```

Next, change into the Open vSwitch directory and apply the patch:

```
$ git am -p3 --reject --directory=datapath/ <patch>
```

If this is successful, proceed to the next patch:

```
$ git am --continue
```

If this is unsuccessful, the above command applies all changes that it can to the working tree, and leaves rejected hunks in corresponding \*.rej files. Proceed by using `git diff` to identify the changes, and edit the files so that the hunk matches what the file looks like when the corresponding commit is checked out in the linux tree. When all hunks are fixed, add the files to the index using `git add`.

If the patch only changes filepaths under `linux/net/openvswitch`, then most likely the patch is fully backported. At this point, review the patch's changes and compare with the latest upstream code for the modified functions. Occasionally, there may be bugs introduced in a particular patch which were fixed in a later patch upstream. To prevent breakage in the OVS tree, consider rolling later bugfixes into the current patch - particularly if they are small, clear bugfixes in the logic of this patch. Then proceed to the next patch using `git am --continue`. If you made any changes to the patch compared with the original version, describe the changes in the commit message.

If the changes affects other paths, then you may also need to backport function definitions from the upstream tree into the `datapath/linux/compat` directory. First, attempt to compile the datapath. If this is successful, then most likely there is no further work required. As per the previous paragraph, consider reviewing and backporting any minor fixes to this code if applicable, then proceed to the next patch using `git am --continue`.

If compilation fails, the compiler will show which functions are missing or broken. Typically this should match with some function definitions provided in the patch file. The following command will attempt to apply all such changes from the patch into the `openvswitch/datapath/linux/compat` directory; Like the previous `git am` command above, it may succeed or fail. If it succeeds, review the patch and proceed to the next patch using `git am --continue`.

```
$ git am -p3 --reject --directory='datapath/linux/compat/' <patch>
```

For each conflicting hunk, attempt to resolve the change so that the function reflects what the function looks like in the upstream Linux tree. After resolving these changes, compile the changes, add the modified files to the index using `git add`, review the patch, and proceed to the next patch using `git am --continue`.

## Submission

Once the patches are all assembled and working on the Open vSwitch tree, they need to be formatted again using `git format-patch`. The common format for commit messages for Linux backport patches is as follows:

```
datapath: Remove incorrect WARN_ONCE().

Upstream commit:
commit c6b2aafffc6934be72d96855c9a1d88970597fbc
Author: Jarno Rajahalme <jarno@ovn.org>
Date: Mon Aug 1 19:08:29 2016 -0700

openvswitch: Remove incorrect WARN_ONCE().

ovs_ct_find_existing() issues a warning if an existing conntrack entry
classified as IP_CT_NEW is found, with the premise that this should
not happen. However, a newly confirmed, non-expected conntrack entry
remains IP_CT_NEW as long as no reply direction traffic is seen. This
has resulted into somewhat confusing kernel log messages. This patch
removes this check and warning.

Fixes: 289f2253 ("openvswitch: Find existing conntrack entry after upcall.")
Suggested-by: Joe Stringer <joe@ovn.org>
Signed-off-by: Jarno Rajahalme <jarno@ovn.org>
```

(continues on next page)

(continued from previous page)

```
Acked-by: Joe Stringer <joe@ovn.org>
```

```
Signed-off-by: Jarno Rajahalme <jarno@ovn.org>
```

The upstream commit SHA should be the one that appears in Linus' tree so that reviewers can compare the backported patch with the one upstream. Note that the subject line for the backported patch replaces the original patch's `openvswitch` prefix with `datapath`. Patches which only affect the `datapath/linux/compat` directory should be prefixed with `compat`.

The contents of a backport should be equivalent to the changes made by the original patch; explain any variations from the original patch in the commit message - For instance if you rolled in a bugfix. Reviewers will verify that the changes made by the backport patch are the same as the changes made in the original commit which the backport is based upon. Patch submission should otherwise follow the regular steps described in [Submitting Patches](#). In particular, if performing kernel patch backports, pay attention to [Datapath testing](#).

### 8.1.3 Open vSwitch Coding Style

This file describes the coding style used in most C files in the Open vSwitch distribution. However, Linux kernel code `datapath` directory follows the Linux kernel's established coding conventions. For the Windows kernel `datapath` code, use the coding style described in [Open vSwitch Windows Datapath Coding Style](#).

The following GNU indent options approximate this style.

```
-npro -bad -bap -bbb -br -blf -brs -cdw -ce -fca -cli0 -npcs -i4 -l79 \
-lc79 -nbfd -nut -saf -sai -saw -sbi4 -sc -sob -st -ncdb -pi4 -cs -bs \
-dil -lp -il0 -hnl
```

#### Basics

- Limit lines to 79 characters.
- Use form feeds (control+L) to divide long source files into logical pieces. A form feed should appear as the only character on a line.
- Do not use tabs for indentation.
- Avoid trailing spaces on lines.

#### Naming

- Use names that explain the purpose of a function or object.
- Use underscores to separate words in an identifier: `multi_word_name`.
- Use lowercase for most names. Use uppercase for macros, macro parameters, and members of enumerations.
- Give arrays names that are plural.
- Pick a unique name prefix (ending with an underscore) for each module, and apply that prefix to all of that module's externally visible names. Names of macro parameters, struct and union members, and parameters in function prototypes are not considered externally visible for this purpose.
- Do not use names that begin with `_`. If you need a name for "internal use only", use `__` as a suffix instead of a prefix.
- Avoid negative names: `found` is a better name than `not_found`.



- In names, a `size` is a count of bytes, a `length` is a count of characters. A buffer has `size`, but a string has `length`. The length of a string does not include the null terminator, but the size of the buffer that contains the string does.

## Comments

Comments should be written as full sentences that start with a capital letter and end with a period. Put two spaces between sentences.

Write block comments as shown below. You may put the `/*` and `*/` on the same line as comment text if you prefer.

```
/*
 * We redirect stderr to /dev/null because we often want to remove all
 * traffic control configuration on a port so its in a known state. If
 * this done when there is no such configuration, tc complains, so we just
 * always ignore it.
 */
```

Each function and each variable declared outside a function, and each struct, union, and typedef declaration should be preceded by a comment. See [functions](#) below for function comment guidelines.

Each struct and union member should each have an inline comment that explains its meaning. structs and unions with many members should be additionally divided into logical groups of members by block comments, e.g.:

```
/* An event that will wake the following call to poll_block(). */
struct poll_waiter {
 /* Set when the waiter is created. */
 struct ovs_list node; /* Element in global waiters list. */
 int fd; /* File descriptor. */
 short int events; /* Events to wait for (POLLIN, POLLOUT). */
 poll_fd_func *function; /* Callback function, if any, or null. */
 void *aux; /* Argument to callback function. */
 struct backtrace *backtrace; /* Event that created waiter, or null. */

 /* Set only when poll_block() is called. */
 struct pollfd *pollfd; /* Pointer to element of the pollfds array
 (null if added from a callback). */
};
```

Use XXX or FIXME comments to mark code that needs work.

Don't use `//` comments.

Don't comment out or `#if 0` out code. Just remove it. The code that was there will still be in version control history.

## Functions

Put the return type, function name, and the braces that surround the function's code on separate lines, all starting in column 0.

Before each function definition, write a comment that describes the function's purpose, including each parameter, the return value, and side effects. References to argument names should be given in single-quotes, e.g. `'arg'`. The comment should not include the function name, nor need it follow any formal structure. The comment does not need to describe how a function does its work, unless this information is needed to use the function correctly (this is often better done with comments *inside* the function).

Simple static functions do not need a comment.

Within a file, non-static functions should come first, in the order that they are declared in the header file, followed by static functions. Static functions should be in one or more separate pages (separated by form feed characters) in logical groups. A commonly useful way to divide groups is by “level”, with high-level functions first, followed by groups of progressively lower-level functions. This makes it easy for the program’s reader to see the top-down structure by reading from top to bottom.

All function declarations and definitions should include a prototype. Empty parentheses, e.g. `int foo();`, do not include a prototype (they state that the function’s parameters are unknown); write `void` in parentheses instead, e.g. `int foo(void);`.

Prototypes for static functions should either all go at the top of the file, separated into groups by blank lines, or they should appear at the top of each page of functions. Don’t comment individual prototypes, but a comment on each group of prototypes is often appropriate.

In the absence of good reasons for another order, the following parameter order is preferred. One notable exception is that data parameters and their corresponding size parameters should be paired.

1. The primary object being manipulated, if any (equivalent to the `this` pointer in C++).
2. Input-only parameters.
3. Input/output parameters.
4. Output-only parameters.
5. Status parameter.

Example:

```
...
/* Stores the features supported by 'netdev' into each of '*current',
 * '*advertised', '*supported', and '*peer' that are non-null. Each value
 * is a bitmap of "enum ofp_port_features" bits, in host byte order.
 * Returns 0 if successful, otherwise a positive errno value. On failure,
 * all of the passed-in values are set to 0. */
int
netdev_get_features(struct netdev *netdev,
 uint32_t *current, uint32_t *advertised,
 uint32_t *supported, uint32_t *peer)
{
 ...
}
...
```

Functions that destroy an instance of a dynamically-allocated type should accept and ignore a null pointer argument. Code that calls such a function (including the C standard library function `free()`) should omit a null-pointer check. We find that this usually makes code easier to read.

Functions in `.c` files should not normally be marked `inline`, because it does not usually help code generation and it does suppress compiler warnings about unused functions. (Functions defined in `.h` usually should be marked `inline`.)

## Function Prototypes

Put the return type and function name on the same line in a function prototype:

```
static const struct option_class *get_option_class(int code);
```

Omit parameter names from function prototypes when the names do not give useful information, e.g.:

```
int netdev_get_mtu(const struct netdev *, int *mtup);
```

## Statements

Indent each level of code with 4 spaces. Use BSD-style brace placement:

```
if (a()) {
 b();
 d();
}
```

Put a space between `if`, `while`, `for`, etc. and the expressions that follow them.

Enclose single statements in braces:

```
if (a > b) {
 return a;
} else {
 return b;
}
```

Use comments and blank lines to divide long functions into logical groups of statements.

Avoid assignments inside `if` and `while` conditions.

Do not put gratuitous parentheses around the expression in a return statement, that is, write `return 0;` and not `return(0);`

Write only one statement per line.

Indent `switch` statements like this:

```
switch (conn->state) {
case S_RECV:
 error = run_connection_input(conn);
 break;

case S_PROCESS:
 error = 0;
 break;

case S_SEND:
 error = run_connection_output(conn);
 break;

default:
 OVS_NOT_REACHED();
}
```

`switch` statements with very short, uniform cases may use an abbreviated style:

```
switch (code) {
case 200: return "OK";
case 201: return "Created";
case 202: return "Accepted";
case 204: return "No Content";
default: return "Unknown";
}
```

Use `for (;;)`  to write an infinite loop.

In an `if/else` construct where one branch is the “normal” or “common” case and the other branch is the “uncommon” or “error” case, put the common case after the `if`, not the `else`. This is a form of documentation. It also places the most important code in sequential order without forcing the reader to visually skip past less important details. (Some compilers also assume that the `if` branch is the more common case, so this can be a real form of optimization as well.)

## Return Values

For functions that return a success or failure indication, prefer one of the following return value conventions:

- An `int` where 0 indicates success and a positive `errno` value indicates a reason for failure.
- A `bool` where `true` indicates success and `false` indicates failure.

## Macros

Don’t define an object-like macro if an `enum` can be used instead.

Don’t define a function-like macro if a `static inline` function can be used instead.

If a macro’s definition contains multiple statements, enclose them with `do { ... } while (0)` to allow them to work properly in all syntactic circumstances.

Do use macros to eliminate the need to update different parts of a single file in parallel, e.g. a list of `enums` and an array that gives the name of each `enum`. For example:

```
/* Logging importance levels. */
#define VLOG_LEVELS \
 VLOG_LEVEL(EMER, LOG_ALERT) \
 VLOG_LEVEL(ERR, LOG_ERR) \
 VLOG_LEVEL(WARN, LOG_WARNING) \
 VLOG_LEVEL(INFO, LOG_NOTICE) \
 VLOG_LEVEL(DBG, LOG_DEBUG)
enum vlog_level {
#define VLOG_LEVEL(NAME, SYSLOG_LEVEL) VLL_##NAME,
 VLOG_LEVELS
#undef VLOG_LEVEL
 VLL_N_LEVELS
};

/* Name for each logging level. */
static const char *level_names[VLL_N_LEVELS] = {
#define VLOG_LEVEL(NAME, SYSLOG_LEVEL) #NAME,
 VLOG_LEVELS
#undef VLOG_LEVEL
};
```

## Thread Safety Annotations

Use the macros in `lib/compiler.h` to annotate locking requirements. For example:

```
static struct ovs_mutex mutex = OVS_MUTEX_INITIALIZER;
static struct ovs_rwlock rwlock = OVS_RWLOCK_INITIALIZER;
```

(continues on next page)

(continued from previous page)

```
void function_require_plain_mutex(void) OVS_REQUIRES (mutex);
void function_require_rwlock(void) OVS_REQ_RDLOCK (rwlock);
```

Pass lock objects, not their addresses, to the annotation macros. (Thus we have `OVS_REQUIRES (mutex)` above, not `OVS_REQUIRES (&mutex)`.)

## Source Files

Each source file should state its license in a comment at the very top, followed by a comment explaining the purpose of the code that is in that file. The comment should explain how the code in the file relates to code in other files. The goal is to allow a programmer to quickly figure out where a given module fits into the larger system.

The first non-comment line in a `.c` source file should be:

```
#include <config.h>
```

`#include` directives should appear in the following order:

1. `#include <config.h>`
2. The module's own headers, if any. Including this before any other header (besides `<config.h>`) ensures that the module's header file is self-contained (see [header files](#) below).
3. Standard C library headers and other system headers, preferably in alphabetical order. (Occasionally one encounters a set of system headers that must be included in a particular order, in which case that order must take precedence.)
4. Open vSwitch headers, in alphabetical order. Use `" "`, not `<>`, to specify Open vSwitch header names.

## Header Files

Each header file should start with its license, as described under [source files](#) above, followed by a “header guard” to make the header file idempotent, like so:

```
#ifndef NETDEV_H
#define NETDEV_H 1

...

#endif /* netdev.h */
```

Header files should be self-contained; that is, they should `#include` whatever additional headers are required, without requiring the client to `#include` them for it.

Don't define the members of a struct or union in a header file, unless client code is actually intended to access them directly or if the definition is otherwise actually needed (e.g. inline functions defined in the header need them).

Similarly, don't `#include` a header file just for the declaration of a struct or union tag (e.g. just for `struct ;`). Just declare the tag yourself. This reduces the number of header file dependencies.

## Types

Use typedefs sparingly. Code is clearer if the actual type is visible at the point of declaration. Do not, in general, declare a typedef for a struct, union, or enum. Do not declare a typedef for a pointer type, because this can be very confusing to the reader.

A function type is a good use for a typedef because it can clarify code. The type should be a function type, not a pointer-to-function type. That way, the typedef name can be used to declare function prototypes. (It cannot be used for function definitions, because that is explicitly prohibited by C89 and C99.)

You may assume that `char` is exactly 8 bits and that `int` and `long` are at least 32 bits.

Don't assume that `long` is big enough to hold a pointer. If you need to cast a pointer to an integer, use `intptr_t` or `uintptr_t` from `.`

Use the `int_t` and `uint_t` types from `for exact-width integer types`. Use the `PRId`, `PRiU`, and `PRiX` macros from `for formatting them with printf()` and related functions.

For compatibility with antique `printf()` implementations:

- Instead of `"%zu"`, use `"%PRIuSIZE"`.
- Instead of `"%td"`, use `"%PRIdPTR"`.
- Instead of `"%ju"`, use `"%PRIuMAX"`.

Other variants exist for different radices. For example, use `"%PRIxSIZE"` instead of `"%zx"` or `"%x"` instead of `"%hhx"`.

Also, instead of `"%hhd"`, use `"%d"`. Be cautious substituting `"%u"`, `"%x"`, and `"%o"` for the corresponding versions with `"hh"`: cast the argument to unsigned char if necessary, because `printf("%hhu", -1)` prints 255 but `printf("%u", -1)` prints 4294967295.

Use bit-fields sparingly. Do not use bit-fields for layout of network protocol fields or in other circumstances where the exact format is important.

Declare bit-fields to be signed or unsigned integer types or `_Bool` (aka `bool`). Do *not* declare bit-fields of type `int`: C99 allows these to be either signed or unsigned according to the compiler's whim. (A 1-bit bit-field of type `int` may have a range of `-1...0!`)

Try to order structure members such that they pack well on a system with 2-byte `short`, 4-byte `int`, and 4- or 8-byte `long` and pointer types. Prefer clear organization over size optimization unless you are convinced there is a size or speed benefit.

Pointer declarators bind to the variable name, not the type name. Write `int *x`, not `int* x` and definitely not `int * x`.

## Expressions

Put one space on each side of infix binary and ternary operators:

```
* / %
+ -
<< >>
< <= > >=
== !=
&
^
|
&&
||
?:
= += -= *= /= %= &= ^= |= <<= >>=
```

Avoid comma operators.

Do not put any white space around postfix, prefix, or grouping operators:

```
() [] -> .
! ~ ++ -- + - * &
```

Exception 1: Put a space after (but not before) the “sizeof” keyword.

Exception 2: Put a space between the ( ) used in a cast and the expression whose type is cast: (void \*) 0.

Break long lines before the ternary operators ? and :, rather than after them, e.g.

```
return (out_port != VIGP_CONTROL_PATH
 ? alpheus_output_port(dp, skb, out_port)
 : alpheus_output_control(dp, skb, fwd_save_skb(skb),
 VIGR_ACTION));
```

Parenthesize the operands of && and || if operator precedence makes it necessary, or if the operands are themselves expressions that use && and ||, but not otherwise. Thus:

```
if (rule && (!best || rule->priority > best->priority)) {
 best = rule;
}
```

but:

```
if (!isdigit((unsigned char)s[0]) ||
 !isdigit((unsigned char)s[1]) ||
 !isdigit((unsigned char)s[2])) {
 printf("string %s does not start with 3-digit code\n", s);
}
```

Do parenthesize a subexpression that must be split across more than one line, e.g.:

```
*idxp = ((l1_idx << PORT_ARRAY_L1_SHIFT) |
 (l2_idx << PORT_ARRAY_L2_SHIFT) |
 (l3_idx << PORT_ARRAY_L3_SHIFT));
```

Breaking a long line after a binary operator gives its operands a more consistent look, since each operand has the same horizontal position. This makes the end-of-line position a good choice when the operands naturally resemble each other, as in the previous two examples. On the other hand, breaking before a binary operator better draws the eye to the operator, which can help clarify code by making it more obvious what’s happening, such as in the following example:

```
if (!ctx.freezing
 && xbridge->has_in_band
 && in_band_must_output_to_local_port(flow)
 && !actions_output_to_local_port(&ctx)) {
```

Thus, decide whether to break before or after a binary operator separately in each situation, based on which of these factors appear to be more important.

Try to avoid casts. Don’t cast the return value of malloc().

The sizeof operator is unique among C operators in that it accepts two very different kinds of operands: an expression or a type. In general, prefer to specify an expression, e.g. `int *x = xmalloc(sizeof *x);`. When the operand of sizeof is an expression, there is no need to parenthesize that operand, and please don’t.

Use the ARRAY\_SIZE macro from `lib/util.h` to calculate the number of elements in an array.

When using a relational operator like < or ==, put an expression or variable argument on the left and a constant argument on the right, e.g. `x == 0`, *not* `0 == x`.

## Blank Lines

Put one blank line between top-level definitions of functions and global variables.

## C DIALECT

Most C99 features are OK because they are widely implemented:

- Flexible array members (e.g. `struct { int foo[]; }`).
- `static inline` functions (but no other forms of `inline`, for which GCC and C99 have differing interpretations).
- `long long`
- `bool` and `<stdbool.h>`, but don't assume that `bool` or `_Bool` can only take on the values 0 or 1, because this behavior can't be simulated on C89 compilers.

Also, don't assume that a conversion to `bool` or `_Bool` follows C99 semantics, i.e. use `(bool)(some_value != 0)` rather than `(bool) some_value`. The latter might produce unexpected results on non-C99 environments. For example, if `bool` is implemented as a typedef of `char` and `some_value = 0x10000000`.

- Designated initializers (e.g. `struct foo foo = { .a = 1 };` and `int a[] = { [2] = 5 };`).
- Mixing of declarations and code within a block. Favor positioning that allows variables to be initialized at their point of declaration.
- Use of declarations in iteration statements (e.g. `for (int i = 0; i < 10; i++)`).
- Use of a trailing comma in an enum declaration (e.g. `enum { x = 1, };`).

As a matter of style, avoid `//` comments.

Avoid using GCC or Clang extensions unless you also add a fallback for other compilers. You can, however, use C99 features or GCC extensions also supported by Clang in code that compiles only on GNU/Linux (such as `lib/netdev-linux.c`), because GCC is the system compiler there.

## Python

When introducing new Python code, try to follow Python's [PEP 8](#) style. Consider running the `pep8` or `flake8` tool against your code to find issues.

## Libraries

When introducing a new library, follow *Open vSwitch Library ABI guide*

## 8.1.4 Open vSwitch Windows Datapath Coding Style

The *coding style* guide gives the flexibility for each platform to use its own coding style for the kernel datapath. This file describes the specific coding style used in most of the C files in the Windows kernel datapath of the Open vSwitch distribution.

Most of the coding conventions applicable for the Open vSwitch distribution are applicable to the Windows kernel datapath as well. There are some exceptions and new guidelines owing to the commonly followed practices in Windows kernel/driver code. They are noted as follows:



## Basics

- Limit lines to 79 characters.

Many times, this is not possible due to long names of functions and it is fine to go beyond the characters limit. One common example is when calling into NDIS functions.

## Types

Use data types defined by Windows for most of the code. This is a common practice in Windows driver code, and it makes integrating with the data structures and functions defined by Windows easier. Example: `DWORD` and `BOOLEAN`.

Use caution in portions of the code that interface with the OVS userspace. OVS userspace does not use Windows specific data types, and when copying data back and forth between kernel and userspace, care should be exercised.

## Naming

It is common practice to use camel casing for naming variables, functions and files in Windows. For types, especially structures, unions and enums, using all upper case letters with words separated by ‘\_’ is common. These practices can be used for OVS Windows datapath. However, use the following guidelines:

- Use lower case to begin the name of a variable.
- Do not use ‘\_’ to begin the name of the variable. ‘\_’ is to be used to begin the parameters of a pre-processor macro.
- Use upper case to begin the name of a function, enum, file name etc.
- Static functions whose scope is limited to the file they are defined in can be prefixed with ‘\_’. This is not mandatory though.
- For types, use all upper case for all letters with words separated by ‘\_’. If camel casing is preferred, use upper case for the first letter.
- It is a common practice to define a pointer type by prefixing the letter ‘P’ to a data type. The same practice can be followed here as well.

For example:

```
static __inline BOOLEAN
OvsDetectTunnelRxPkt (POVS_FORWARDING_CONTEXT ovsFwdCtx,
 POVS_FLOW_KEY flowKey)
{
 POVS_VPORT_ENTRY tunnelVport = NULL;

 if (!flowKey->ipKey.nwFrag &&
 flowKey->ipKey.nwProto == IPPROTO_UDP &&
 flowKey->ipKey.l4.tpDst == VXLAN_UDP_PORT_NBO) {
 tunnelVport = OvsGetTunnelVport (OVSWIN_VPORT_TYPE_VXLAN);
 ovsActionStats.rxVxlan++;
 } else {
 return FALSE;
 }

 if (tunnelVport) {
 ASSERT(ovsFwdCtx->tunnelRxNic == NULL);
 ovsFwdCtx->tunnelRxNic = tunnelVport;
 return TRUE;
 }
}
```

(continues on next page)

(continued from previous page)

```
}

 return FALSE;
}
```

For declaring variables of pointer type, use of the pointer data type prefixed with ‘P’ is preferred over using ‘\*’. This is not mandatory though, and is only prescribed since it is a common practice in Windows.

Example, #1 is preferred over #2 though #2 is also equally correct:

1. PNET\_BUFFER\_LIST curNbl;
2. NET\_BUFFER\_LIST \*curNbl;

## Comments

Comments should be written as full sentences that start with a capital letter and end with a period. Putting two spaces between sentences is not necessary.

// can be used for comments as long as the comment is a single line comment. For block comments, use /\* \*/ comments

## Functions

Put the return type, function name, and the braces that surround the function’s code on separate lines, all starting in column 0.

Before each function definition, write a comment that describes the function’s purpose, including each parameter, the return value, and side effects. References to argument names should be given in single-quotes, e.g. ‘arg’. The comment should not include the function name, nor need it follow any formal structure. The comment does not need to describe how a function does its work, unless this information is needed to use the function correctly (this is often better done with comments *inside* the function).

Mention any side effects that the function has that are not obvious based on the name of the function or based on the workflow it is called from.

In the interest of keeping comments describing functions similar in structure, use the following template.

```
/*
*-----
* Any description of the function, arguments, return types, assumptions and
* side effects.
*-----
*/
```

## Source Files

Each source file should state its license in a comment at the very top, followed by a comment explaining the purpose of the code that is in that file. The comment should explain how the code in the file relates to code in other files. The goal is to allow a programmer to quickly figure out where a given module fits into the larger system.

The first non-comment line in a .c source file should be:

```
#include <precomp.h>
```

`#include` directives should appear in the following order:

1. `#include <precomp.h>`
2. The module's own headers, if any. Including this before any other header (besides `<precomp.h>`) ensures that the module's header file is self-contained (see *Header Files*) below.
3. Standard C library headers and other system headers, preferably in alphabetical order. (Occasionally one encounters a set of system headers that must be included in a particular order, in which case that order must take precedence.)
4. Open vSwitch headers, in alphabetical order. Use `"`, not `<>`, to specify Open vSwitch header names.

### 8.1.5 Open vSwitch Documentation Style

This file describes the documentation style used in all documentation found in Open vSwitch. Documentation includes any documents found in `Documentation` along with any `README`, `MAINTAINERS`, or generally `rst` suffixed documents found in the project tree.

---

**Note:** This guide only applies to documentation for Open vSwitch v2.7. or greater. Previous versions of Open vSwitch used a combination of Markdown and raw plain text, and guidelines for these are not detailed here.

---

#### reStructuredText vs. Sphinx

reStructuredText (rST) is the syntax, while Sphinx is a documentation generator. Sphinx introduces a number of extensions to rST, like the `:ref:` role, which can and should be used in documentation, but these will not work correctly on GitHub. As such, these extensions should not be used in any documentation in the root level, such as the `README`.

#### rST Conventions

##### Basics

Many of the basic documentation guidelines match those of the *Open vSwitch Coding Style*.

- Use reStructuredText (rST) for all documentation.  
Sphinx extensions can be used, but only for documentation in the `Documentation` folder.
- Limit lines at 79 characters.

---

**Note:** An exception to this rule is text within code-block elements that cannot be wrapped and links within references.

---

- Use spaces for indentation.
- Match indentation levels.

A change in indentation level usually signifies a change in content nesting, by either closing the existing level or introducing a new level.

- Avoid trailing spaces on lines.
- Include a license (see this file) in all docs.

- Most importantly, always build and display documentation before submitting changes! Docs aren't unit testable, so visible inspection is necessary.

### File Names

- Use hyphens as space delimiters. For example: `my-readme-document.rst`

---

**Note:** An exception to this rule is any man pages, which take an trailing number corresponding to the number of arguments required. This number is preceded by an underscore.

---

- Use lowercase filenames.

---

**Note:** An exception to this rule is any documents found in the root-level of the project.

---

### Titles

- Use the following headers levels.

```
===== Heading 0 (reserved for the title in a document)
----- Heading 1
~~~~~ Heading 2
+++++++ Heading 3
'''''''' Heading 4
```

---

**Note:** Avoid using lower heading levels by rewriting and reorganizing the information.

---

- Under- and overlines should be of the same length as that of the heading text.
- Use “title case” for headers.

### Code

- Use `::` to prefix code.
- Don't use syntax highlighting such as `.. highlight:: <syntax>` or `code-block:: <syntax>` because it depends on external `pygments` library.
- Prefix commands with `$`.
- Where possible, include fully-working snippets of code. If there pre-requisites, explain what they are and how to achieve them.

### Admonitions

- Use admonitions to call attention to important information.:

```
.. note::

    This is a sample callout for some useful tip or trick.
```

Example admonitions include: warning, important, note, tip or seealso.

- Use notes sparingly. Avoid having more than one per subsection.

Tables

- Use either graphic tables, list tables or CSV tables.

Graphic tables

```
.. table:: OVS-Linux kernel compatibility

=====
Open vSwitch Linux kernel
=====
1.4.x      2.6.18 to 3.2
1.5.x      2.6.18 to 3.2
1.6.x      2.6.18 to 3.2
=====
```

```
.. table:: OVS-Linux kernel compatibility

+-----+-----+
| Open vSwitch | Linux kernel |
+=====+=====+
| 1.4.x      | 2.6.18 to 3.2 |
+-----+-----+
| 1.5.x      | 2.6.18 to 3.2 |
+-----+-----+
| 1.6.x      | 2.6.18 to 3.2 |
+-----+-----+
```

**Note:** The table role - `.. table:: <name>` - can be safely omitted.

List tables

```
.. list-table:: OVS-Linux kernel compatibility
:widths: 10 15
:header-rows: 1

* - Open vSwitch
  - Linux kernel
* - 1.4.x
  - 2.6.18 to 3.2
* - 1.5.x
  - 2.6.18 to 3.2
```

(continues on next page)

(continued from previous page)

```
* - 1.6.x
  - 2.6.18 to 3.2
```

### CSV tables

```
.. csv-table:: OVS-Linux kernel compatibility
   :header: Open vSwitch, Linux kernel
   :widths: 10 15

   1.4.x, 2.6.18 to 3.2
   1.5.x, 2.6.18 to 3.2
   1.6.x, 2.6.18 to 3.2
```

### Cross-referencing

- To link to an external file or document, include as a link.:

```
Here's a `link <http://openvswitch.org>`__ to the Open vSwitch website.

Here's a `link`_ in reference style.

.. _link: http://openvswitch.org
```

- You can also use citations.:

```
Refer to the Open vSwitch documentation [1]_.

References
-----

.. [1]: http://openvswitch.org
```

- To cross-reference another doc, use the `doc` role.:

```
Here is a link to the :doc:`/README.rst`
```

---

**Note:** This is a Sphinx extension. Do not use this in any top-level documents.

---

- To cross-reference an arbitrary location in a doc, use the `ref` role.:

```
.. _sample-crossref

Title
~~~~

Hello, world.

Another Title
~~~~~
```

(continues on next page)

(continued from previous page)

```
Here is a cross-reference to :ref:`sample-crossref`.
```

---

**Note:** This is a Sphinx extension. Do not use this in any top-level documents.

---

## Figures and Other Media

- All images should be in PNG format and compressed where possible. For PNG files, use OptiPNG and AdvanceCOMP's advpng:

```
$ optipng -o7 -zml-9 -i0 -strip all <path_to_png>
$ advpng -z4 <path_to_png>
```

- Any ASCII text “images” should be included in code-blocks to preserve formatting
- Include other reStructuredText verbatim in a current document

## Comments

- Comments are indicated by means of the `..` marker.:

```
.. TODO(stephenfin) This section needs some work. This TODO will not
   appear in the final generated document, however.
```

## Man Pages

In addition to the above, man pages have some specific requirements:

- You **must** define the following sections:
  - Synopsis
  - Description
  - Options

Note that *NAME* is not included - this is automatically generated by Sphinx and should not be manually defined. Also note that these do not need to be uppercase - Sphinx will do this automatically.

Additional sections are allowed. Refer to *man-pages(8)* for information on the sections generally allowed.

- You **must not** define a *NAME* section.

See above.

- The *OPTIONS* section must describe arguments and options using the `program` and `option` directives.

This ensures the output is formatted correctly and that you can cross-reference various programs and commands from the documentation. For example:

```
.. program:: ovs-do-something

.. option:: -f, --force
```

(continues on next page)

(continued from previous page)

```
Force the operation

.. option:: -b <bridge>, --bridge <bridge>

Name or ID of bridge
```

---

**Important:** Option argument names should be enclosed in angle brackets, as above.

---

- Any references to the application or any other Open vSwitch application must be marked up using the *program* role.

This allows for easy linking in the HTML output and correct formatting in the man page output. For example:

```
To do something, run :program:`ovs-do-something`.
```

- The man page must be included in the list of man page documents found in `conf.py`

Refer to existing man pages, such as *ovs-vlan-test* for a worked example.

## Writing Style

Follow these guidelines to ensure readability and consistency of the Open vSwitch documentation. These guidelines are based on the *IBM Style Guide*.

- Use standard US English
  - Use a spelling and grammar checking tool as necessary.
- Expand initialisms and acronyms on first usage.
  - Commonly used terms like CPU or RAM are allowed.

Do not use	Do use
OVS is a virtual switch. OVS has...	Open vSwitch (OVS) is a virtual switch. OVS has...
The VTEP emulator is...	The Virtual Tunnel Endpoint (VTEP) emulator is...

- Write in the active voice
  - The subject should do the verb's action, rather than be acted upon.

Do not use	Do use
A bridge is created by you	Create a bridge

- Write in the present tense

Do not use	Do use
Once the bridge is created, you can create a port	Once the bridge is created, create a port

- Write in second person

Do not use	Do use
To create a bridge, the user runs:	To create a bridge, run:



- Keep sentences short and concise
  - Eliminate needless politeness
- Avoid “please” and “thank you”

## Helpful Tools

There are a number of tools, online and offline, which can be used to preview documents are you edit them:

- [rst.ninjs.org](http://rst.ninjs.org)  
An online rST editor/previewer
- [ReText](#)  
A simple but powerful editor for Markdown and reStructuredText. ReText is written in Python.
- [restview](#)  
A viewer for ReStructuredText documents that renders them on the fly.

## Useful Links

- [Quick reStructuredText](#)
- [Sphinx Documentation](#)

## 8.1.6 Open vSwitch Library ABI Updates

This file describes the manner in which the Open vSwitch shared library manages different ABI and API revisions. This document aims to describe the background, goals, and concrete mechanisms used to export code-space functionality so that it may be shared between multiple applications.

### Definitions

Table 1: Definitions for terms appearing in this document

Term	Definition
ABI	Abbreviation of Application Binary Interface
API	Abbreviation of Application Programming Interface
Application Binary Interface	The low-level runtime interface exposed by an object file.
Application Programming Interface	The source-code interface descriptions intended for use in multiple translation units when compiling.
Code library	A collection of function implementations and definitions intended to be exported and called through a well-defined interface.
Shared Library	A code library which is imported at run time.

### Overview

C and C++ applications often use ‘external’ functionality, such as printing specialized data types or parsing messages, which has been exported for common use. There are many possible ways for applications to call such external functionality, for instance by including an appropriate inline definition which the compiler can emit as code in each

function it appears. One such way of exporting and importing such functionality is through the use of a library of code.

When a compiler builds object code from source files to produce object code, the results are binary data arranged with specific calling conventions, alignments, and order suitable for a run-time environment or linker. This result defines a specific ABI.

As library of code develops and its exported interfaces change over time, the resulting ABI may change as well. Therefore, care must be taken to ensure the changes made to libraries of code are effectively communicated to applications which use them. This includes informing the applications when incompatible changes are made.

The Open vSwitch project exports much of its functionality through multiple such libraries of code. These libraries are intended for multiple applications to import and use. As the Open vSwitch project continues to evolve and change, its exported code will evolve as well. To ensure that applications linking to these libraries are aware of these changes, Open vSwitch employs libtool version stamps.

### ABI Policy

Open vSwitch will export the ABI version at the time of release, such that the library name will be the major.minor version, and the rest of the release version information will be conveyed with a libtool interface version.

The intent is for Open vSwitch to maintain an ABI stability for each minor revision only (so that Open vSwitch release 2.5 carries a guarantee for all 2.5.ZZ micro-releases). This means that any porting effort to stable branches must take not to disrupt the existing ABI.

In the event that a bug must be fixed in a backwards-incompatible way, developers must bump the libtool ‘current’ version to inform the linker of the ABI breakage. This will signal that libraries exposed by the subsequent release will not maintain ABI stability with the previous version.

### Coding

At build time, if building shared libraries by passing the *–enable-shared* arguments to *./configure*, version information is extracted from the `$PACKAGE_VERSION` automake variable and formatted into the appropriate arguments. These get exported for use in Makefiles as `$OVS_LTINFO`, and passed to each exported library along with other `LDFLAGS`.

Therefore, when adding a new library to the build system, these version flags should be included with the `$LDFLAGS` variable. Nothing else needs to be done.

Changing an exported function definition (from a file in, for instance *lib/\*.h*) is only permitted from minor release to minor release. Likewise changes to library data structures should only occur from minor release to minor release.

## 8.2 Mailing Lists

---

**Important:** Report security issues **only** to [security@openvswitch.org](mailto:security@openvswitch.org). For more information, refer to our [security policies](#).

---

### 8.2.1 ovs-announce

The [ovs-announce](#) mailing list is used to announce new versions of Open vSwitch and is extremely low-volume. ([subscribe](#)) ([archives](#))

### 8.2.2 ovs-discuss

The `ovs-discuss` mailing list is used to discuss plans and design decisions for Open vSwitch. It is also an appropriate place for user questions. ([subscribe](#)) ([archives](#))

### 8.2.3 ovs-dev

The `ovs-dev` mailing list is used to discuss development and review code before being committed. ([subscribe](#)) ([archives](#))

### 8.2.4 ovs-git

The `ovs-git` mailing list hooks into Open vSwitch's version control system to receive commits. ([subscribe](#)) ([archives](#))

### 8.2.5 ovs-build

The `ovs-build` mailing list hooks into Open vSwitch's continuous integration system to receive build reports. ([subscribe](#)) ([archives](#))

### 8.2.6 bugs

The `bugs` mailing list is an alias for the discuss mailing list.

### 8.2.7 security

The `security` mailing list is for submitting security vulnerabilities to the security team.

## 8.3 Patchwork

Open vSwitch uses `Patchwork` to track the status of patches sent to the *ovs-dev mailing list*. The Open vSwitch Patchwork instance can be found on [ozlabs.org](http://ozlabs.org).

Patchwork provides a number of useful features for developers working on Open vSwitch:

- Tracking the lifecycle of patches (accepted, rejected, under-review, ...)
- Assigning reviewers (delegates) to patches
- Downloading/applying patches, series, and bundles via the web UI or the REST API (see *git-pw*)
- A usable UI for viewing patch discussions

### 8.3.1 git-pw

The *git-pw* tool provides a way to download and apply patches, series, and bundles. You can install *git-pw* from PyPi like so:

```
$ pip install --user git-pw
```

To actually use *git-pw*, you must configure it with the Patchwork instance URL, Patchwork project, and your Patchwork user authentication token. The URL and project are provided below, but you must obtain your authentication token from your [Patchwork User Profile](#) page. If you do not already have a Patchwork user account, you should create one now.

Once your token is obtained, configure *git-pw* as below. Note that this must be run from within the Open vSwitch Git repository:

```
$ git config pw.server https://patchwork.ozlabs.org/
$ git config pw.project openvswitch
$ git config pw.token $PW_TOKEN # using the token obtained earlier
```

Once configured, run the following to get information about available commands:

```
$ git pw --help
```

### 8.3.2 pwclient

The *pwclient* is a legacy tool that provides some of the functionality of *git-pw* but uses the legacy XML-RPC API. It is considered deprecated in its current form and *git-pw* should be used instead.

## 8.4 Open vSwitch Release Process

This document describes the process ordinarily used for Open vSwitch development and release. Exceptions are sometimes necessary, so all of the statements here should be taken as subject to change through rough consensus of Open vSwitch contributors, obtained through public discussion on, e.g., ovs-dev or the #openvswitch IRC channel.

### 8.4.1 Release Strategy

Open vSwitch feature development takes place on the “master” branch. Ordinarily, new features are rebased against master and applied directly. For features that take significant development, sometimes it is more appropriate to merge a separate branch into master; please discuss this on ovs-dev in advance.

The process of making a release has the following stages. See [Release Scheduling](#) for the timing of each stage:

1. “Soft freeze” of the master branch.

During the freeze, we ask committers to refrain from applying patches that add new features unless those patches were already being publicly discussed and reviewed before the freeze began. Bug fixes are welcome at any time. Please propose and discuss exceptions on ovs-dev.

2. Fork a release branch from master, named for the expected release number, e.g. “branch-2.3” for the branch that will yield Open vSwitch 2.3.x.

Release branches are intended for testing and stabilization. At this stage and in later stages, they should receive only bug fixes, not new features. Bug fixes applied to release branches should be backports of corresponding bug fixes to the master branch, except for bugs present only on release branches (which are rare in practice).

At this stage, sometimes there can be exceptions to the rule that a release branch receives only bug fixes. Like bug fixes, new features on release branches should be backports of the corresponding commits on the master branch. Features to be added to release branches should be limited in scope and risk and discussed on ovs-dev before creating the branch.

3. When committers come to rough consensus that the release is ready, they release the .0 release on its branch, e.g. 2.3.0 for branch-2.3. To make the actual release, a committer pushes a signed tag named, e.g. v2.3.0, to the Open vSwitch repository, makes a release tarball available on [openvswitch.org](http://openvswitch.org), and posts a release announcement to [ovs-announce](mailto:ovs-announce).
4. As bug fixes accumulate, or after important bugs or vulnerabilities are fixed, committers may make additional releases from a branch: 2.3.1, 2.3.2, and so on. The process is the same for these additional release as for a .0 release.

At most two release branches are formally maintained at any given time: the latest release and the latest release designed as LTS. An LTS release is one that the OVS project has designated as being maintained for a longer period of time. Currently, an LTS release is maintained until the next LTS is chosen. There is not currently a strict guideline on how often a new LTS release is chosen, but so far it has been about every 2 years. That could change based on the current state of OVS development. For example, we do not want to designate a new release as LTS that includes disruptive internal changes, as that may make it harder to support for a longer period of time. Discussion about choosing the next LTS release occurs on the OVS development mailing list.

## 8.4.2 Release Numbering

The version number on master should normally end in .90. This indicates that the Open vSwitch version is “almost” the next version to branch.

Forking master into branch-x.y requires two commits to master. The first is titled “Prepare for x.y.0” and increments the version number to x.y. This is the initial commit on branch-x.y. The second is titled “Prepare for post-x.y.0 (x.y.90)” and increments the version number to x.y.90.

The version number on a release branch is x.y.z, where z is initially 0. Making a release requires two commits. The first is titled *Set release dates for x.y.z.* and updates NEWS and debian/changelog to specify the release date of the new release. This commit is the one made into a tarball and tagged. The second is titled *Prepare for x.y.(z+1).* and increments the version number and adds a blank item to NEWS with an unspecified date.

## 8.4.3 Release Scheduling

Open vSwitch makes releases at the following six-month cadence. All dates are approximate:

Time (months)	Dates	Stage
T	Mar 1, Sep 1	Begin x.y release cycle
T + 4	Jul 1, Jan 1	“Soft freeze” master for x.y release
T + 4.5	Jul 15, Jan 15	Fork branch-x.y from master
T + 5.5	Aug 15, Feb 15	Release version x.y.0

## 8.4.4 Contact

Use [dev@openvswitch.org](mailto:dev@openvswitch.org) to discuss the Open vSwitch development and release process.

## 8.5 Reporting Bugs in Open vSwitch

We are eager to hear from users about problems that they have encountered with Open vSwitch. This file documents how best to report bugs so as to ensure that they can be fixed as quickly as possible.

Please report bugs by sending email to [bugs@openvswitch.org](mailto:bugs@openvswitch.org).

For reporting security vulnerabilities, please read *Open vSwitch's Security Process*.

The most important parts of your bug report are the following:

- What you did that make the problem appear.
- What you expected to happen.
- What actually happened.

Please also include the following information:

- The Open vSwitch version number (as output by `ovs-vswitchd --version`).
- The Git commit number (as output by `git rev-parse HEAD`), if you built from a Git snapshot.
- Any local patches or changes you have applied (if any).

The following are also handy sometimes:

- The kernel version on which Open vSwitch is running (from `/proc/version`) and the distribution and version number of your OS (e.g. “Centos 5.0”).
- The contents of the vswitchd configuration database (usually `/etc/openvswitch/conf.db`).
- The output of `ovs-dpctl show`.
- If you have Open vSwitch configured to connect to an OpenFlow controller, the output of `ovs-ofctl show <bridge>` for each <bridge> configured in the vswitchd configuration database.
- A fix or workaround, if you have one.
- Any other information that you think might be relevant.

---

**Important:** [bugs@openvswitch.org](mailto:bugs@openvswitch.org) is a public mailing list, to which anyone can subscribe, so do not include confidential information in your bug report.

---

## 8.6 Open vSwitch's Security Process

This is a proposed security vulnerability reporting and handling process for Open vSwitch. It is based on the OpenStack vulnerability management process described at [https://wiki.openstack.org/wiki/Vulnerability\\_Management](https://wiki.openstack.org/wiki/Vulnerability_Management).

The OVS security team coordinates vulnerability management using the ovs-security mailing list. Membership in the security team and subscription to its mailing list consists of a small number of trustworthy people, as determined by rough consensus of the Open vSwitch committers on the ovs-committers mailing list. The Open vSwitch security team should include Open vSwitch committers, to ensure prompt and accurate vulnerability assessments and patch review.

We encourage everyone involved in the security process to GPG-sign their emails. We additionally encourage GPG-encrypting one-on-one conversations as part of the security process.

### 8.6.1 What is a vulnerability?

All vulnerabilities are bugs, but not every bug is a vulnerability. Vulnerabilities compromise one or more of:

- Confidentiality (personal or corporate confidential data).
- Integrity (trustworthiness and correctness).
- Availability (uptime and service).

Here are some examples of vulnerabilities to which one would expect to apply this process:

- A crafted packet that causes a kernel or userspace crash (Availability).
- A flow translation bug that misforwards traffic in a way likely to hop over security boundaries (Integrity).
- An OpenFlow protocol bug that allows a controller to read arbitrary files from the file system (Confidentiality).
- Misuse of the OpenSSL library that allows bypassing certificate checks (Integrity).
- A bug (memory corruption, overflow, ...) that allows one to modify the behaviour of OVS through external configuration interfaces such as OVSDDB (Integrity).
- Privileged information is exposed to unprivileged users (Confidentiality).

If in doubt, please do use the vulnerability management process. At worst, the response will be to report the bug through the usual channels.

## 8.6.2 Step 1: Reception

To report an Open vSwitch vulnerability, send an email to the ovs-security mailing list (see [contact](#) at the end of this document). A security team member should reply to the reporter acknowledging that the report has been received.

Consider reporting the information mentioned in *Reporting Bugs in Open vSwitch*, where relevant.

Reporters may ask for a GPG key while initiating contact with the security team to deliver more sensitive reports.

The Linux kernel has [its own vulnerability management process](#). Handling of vulnerabilities that affect both the Open vSwitch tree and the upstream Linux kernel should be reported through both processes. Send your report as a single email to both the kernel and OVS security teams to allow those teams to most easily coordinate among themselves.

## 8.6.3 Step 2: Assessment

The security team should discuss the vulnerability. The reporter should be included in the discussion (via “CC”) to an appropriate degree.

The assessment should determine which Open vSwitch versions are affected (e.g. every version, only the latest release, only unreleased versions), the privilege required to take advantage of the vulnerability (e.g. any network user, any local L2 network user, any local system user, connected OpenFlow controllers), the severity of the vulnerability, and how the vulnerability may be mitigated (e.g. by disabling a feature).

The treatment of the vulnerability could end here if the team determines that it is not a realistic vulnerability.

## 8.6.4 Step 3a: Document

The security team develops a security advisory document. The security team may, at its discretion, include the reporter (via “CC”) in developing the security advisory document, but in any case should accept feedback from the reporter before finalizing the document. When the document is final, the security team should obtain a CVE for the vulnerability from a CNA (<https://cve.mitre.org/cve/cna.html>).

The document credits the reporter and describes the vulnerability, including all of the relevant information from the assessment in step 2. Suitable sections for the document include:

```
* Title: The CVE identifier, a short description of the
  vulnerability. The title should mention Open vSwitch.

In email, the title becomes the subject. Pre-release advisories
are often passed around in encrypted email, which have plaintext
```

(continues on next page)

(continued from previous page)

subjects, so the title should not be too specific.

- \* **Description:** A few paragraphs describing the general characteristics of the vulnerability, including the versions of Open vSwitch that are vulnerable, the kind of attack that exposes the vulnerability, and potential consequences of the attack.

The description should re-state the CVE identifier, in case the subject is lost when an advisory is sent over email.

- \* **Mitigation:** How an Open vSwitch administrator can minimize the potential for exploitation of the vulnerability, before applying a fix. If no mitigation is possible or recommended, explain why, to reduce the chance that at-risk users believe they are not at risk.
- \* **Fix:** Describe how to fix the vulnerability, perhaps in terms of applying a source patch. The patch or patches themselves, if included in the email, should be at the very end of the advisory to reduce the risk that a reader would stop reading at this point.
- \* **Recommendation:** A concise description of the security team's recommendation to users.
- \* **Acknowledgments:** Thank the reporters.
- \* **Vulnerability Check:** A step-by-step procedure by which a user can determine whether an installed copy of Open vSwitch is vulnerable.

The procedure should clearly describe how to interpret the results, including expected results in vulnerable and not-vulnerable cases. Thus, procedures that produce clear and easily distinguished results are preferred.

The procedure should assume as little understanding of Open vSwitch as possible, to make it more likely that a competent administrator who does not specialize in Open vSwitch can perform it successfully.

The procedure should have minimal dependencies on tools that are not widely installed.

Given a choice, the procedure should be one that takes at least some work to turn into a useful exploit. For example, a procedure based on "ovs-appctl" commands, which require local administrator access, is preferred to one that sends test packets to a machine, which only requires network connectivity.

The section should say which operating systems it is designed for. If the procedure is likely to be specific to particular architectures (e.g. x86-64, i386), it should state on which ones it has been tested.

This section should state the risks of the procedure. For

(continues on next page)



(continued from previous page)

example, if it can crash Open vSwitch or disrupt packet forwarding, say so.

It is more useful to explain how to check an installed and running Open vSwitch than one built locally from source, but if it is easy to use the procedure from a sandbox environment, it can be helpful to explain how to do so.

- \* Patch: If a patch or patches are available, and it is practical to include them in the email, put them at the end. Format them as described in :doc:`contributing/submitting-patches`, that is, as output by "git format-patch".

The patch subjects should include the version for which they are suited, e.g. "[PATCH branch-2.3]" for a patch against Open vSwitch 2.3.x. If there are multiple patches for multiple versions of Open vSwitch, put them in separate sections with clear titles.

Multiple patches for a single version of Open vSwitch, that must be stacked on top of each other to fix a single vulnerability, are undesirable because users are less likely to apply all of them correctly and in the correct order.

Each patch should include a Vulnerability tag with the CVE identifier, a Reported-by tag or tags to credit the reporters, and a Signed-off-by tag to acknowledge the Developer's Certificate of Origin. It should also include other appropriate tags, such as Acked-by tags obtained during review.

CVE-2016-2074 is an example advisory document.

## 8.6.5 Step 3b: Fix

Steps 3a and 3b may proceed in parallel.

The security team develops and obtains (private) reviews for patches that fix the vulnerability. If necessary, the security team pulls in additional developers, who must agree to maintain confidentiality.

## 8.6.6 Step 4: Embargoed Disclosure

The security advisory and patches are sent to downstream stakeholders, with an embargo date and time set from the time sent. Downstream stakeholders are expected not to deploy or disclose patches until the embargo is passed.

A disclosure date is negotiated by the security team working with the bug submitter as well as vendors. However, the Open vSwitch security team holds the final say when setting a disclosure date. The timeframe for disclosure is from immediate (esp. if it's already publicly known) to a few weeks. As a basic default policy, we expect report date to disclosure date to be 10 to 15 business days.

Operating system vendors are obvious downstream stakeholders. It may not be necessary to be too choosy about who to include: any major Open vSwitch user who is interested and can be considered trustworthy enough could be included. To become a downstream stakeholder, email the ovs-security mailing list.

If the vulnerability is already public, skip this step.

### 8.6.7 Step 5: Public Disclosure

When the embargo expires, push the (reviewed) patches to appropriate branches, post the patches to the ovs-dev mailing list (noting that they have already been reviewed and applied), post the security advisory to appropriate mailing lists (ovs-announce, ovs-discuss), and post the security advisory on the Open vSwitch webpage.

When the patch is applied to LTS (long-term support) branches, a new version should be released.

The security advisory should be GPG-signed by a security team member with a key that is in a public web of trust.

#### Contact

Report security vulnerabilities to the ovs-security mailing list: [security@openvswitch.org](mailto:security@openvswitch.org)

Report problems with this document to the ovs-bugs mailing list: [bugs@openvswitch.org](mailto:bugs@openvswitch.org)

## 8.7 The Linux Foundation Open vSwitch Project Charter

Effective August 9, 2016

### 1. Mission of Open vSwitch Project (“OVS”).

The mission of OVS is to:

- (a) create an open source, production quality virtual networking platform, including a software switch, control plane, and related components, that supports standard management interfaces and opens the forwarding functions to programmatic extension and control; and
- (b) host the infrastructure for an OVS community, establishing a neutral home for community assets, infrastructure, meetings, events and collaborative discussions.

### 2. Technical Steering Committee (“TSC”)

- (a) A TSC shall be composed of the Committers for OVS. The list of Committers on the TSC are available at *Committers*.
- (b) TSC projects generally will involve Committers and Contributors:
  - i. Contributors: anyone in the technical community that contributes code, documentation or other technical artifacts to the OVS codebase.
  - ii. Committers: Contributors who have the ability to commit directly to a project’s main branch or repository on an OVS project.
- (c) Participation in as a Contributor and/or Committer is open to anyone under the terms of this Charter. The TSC may:
  - i. establish work flows and procedures for the submission, approval and closure or archiving of projects,
  - ii. establish criteria and processes for the promotion of Contributors to Committer status, available at *OVS Committer Grant/Revocation Policy*. and
  - iii. amend, adjust and refine the roles of Contributors and Committers listed in Section 2.b., create new roles and publicly document responsibilities and expectations for such roles, as it sees fit, available at *Expectations for Developers with Open vSwitch Repo Access*.
- (d) Responsibilities: The TSC is responsible for overseeing OVS activities and making decisions that impact the mission of OVS, including:
  - i. coordinating the technical direction of OVS;

- ii. approving project proposals (including, but not limited to, incubation, deprecation and changes to a project's charter or scope);
- iii. creating sub-committees or working groups to focus on cross-project technical issues and requirements;
- iv. communicating with external and industry organizations concerning OVS technical matters;
- v. appointing representatives to work with other open source or standards communities;
- vi. establishing community norms, workflows or policies including processes for contributing (available at *Contributing to Open vSwitch*), issuing releases, and security issue reporting policies;
- vii. discussing, seeking consensus, and where necessary, voting on technical matters relating to the code base that affect multiple projects; and
- viii. coordinate any marketing, events or communications with The Linux Foundation.

### 3. TSC Voting

- (a) While it is the goal of OVS to operate as a consensus based community, if any TSC decision requires a vote to move forward, the Committers shall vote on a one vote per Committer basis.
- (b) TSC votes should be conducted by email. In the case of a TSC meeting where a valid vote is taken, the details of the vote and any discussion should be subsequently documented for the community (e.g. to the appropriate email mailing list).
- (c) Quorum for TSC meetings shall require two-thirds of the TSC representatives. The TSC may continue to meet if quorum is not met, but shall be prevented from making any decisions requiring a vote at the meeting.
- (d) Except as provided in Section 8.d. and 9.a., decisions by electronic vote (e.g. email) shall require a majority of all voting TSC representatives. Decisions by electronic vote shall be made timely, and unless specified otherwise, within three (3) business days. Except as provided in Section 8.d. and 9.a., decisions by vote at a meeting shall require a majority vote, provided quorum is met.
- (e) In the event of a tied vote with respect to an action that cannot be resolved by the TSC, any TSC representative shall be entitled to refer the matter to the Linux Foundation for assistance in reaching a decision.

### 4. Antitrust Guidelines

- (a) All participants in OVS shall abide by The Linux Foundation Antitrust Policy available at <http://www.linuxfoundation.org/antitrust-policy>.
- (b) All members shall encourage open participation from any organization able to meet the participation requirements, regardless of competitive interests. Put another way, the community shall not seek to exclude any participant based on any criteria, requirements or reasons other than those that are reasonable and applied on a non-discriminatory basis to all participants.

### 5. Code of Conduct

- (a) The TSC may adopt a specific OVS Project code of conduct, with approval from the LF.

### 6. Budget and Funding

- (a) The TSC shall coordinate any budget or funding needs with The Linux Foundation. Companies participating may be solicited to sponsor OVS activities and infrastructure needs on a voluntary basis.
- (b) The Linux Foundation shall have custody of and final authority over the usage of any fees, funds and other cash receipts.
- (c) A General & Administrative (G&A) fee will be applied by the Linux Foundation to funds raised to cover Finance, Accounting, and operations. The G&A fee shall equal 9% of OVS's first \$1,000,000 of gross receipts and 6% of OVS's gross receipts over \$1,000,000.

- (d) Under no circumstances shall The Linux Foundation be expected or required to undertake any action on behalf of OVS that is inconsistent with the tax exempt purpose of The Linux Foundation.

#### 7. General Rules and Operations.

The OVS project shall be conducted so as to:

- (a) engage in the work of the project in a professional manner consistent with maintaining a cohesive community, while also maintaining the goodwill and esteem of The Linux Foundation in the open source software community;
- (b) respect the rights of all trademark owners, including any branding and usage guidelines;
- (c) engage The Linux Foundation for all OVS press and analyst relations activities;
- (d) upon request, provide information regarding Project participation, including information regarding attendance at Project-sponsored events, to The Linux Foundation; and
- (e) coordinate with The Linux Foundation in relation to any websites created directly for OVS.

#### 8. Intellectual Property Policy

- (a) Members agree that all new inbound code contributions to OVS shall be made under the Apache License, Version 2.0 (available at <http://www.apache.org/licenses/LICENSE-2.0>). All contributions shall be accompanied by a Developer Certificate of Origin sign-off (<http://developercertificate.org>) that is submitted through a TSC and LF-approved contribution process.
- (b) All outbound code will be made available under the Apache License, Version 2.0.
- (c) All documentation will be contributed to and made available by OVS under the Apache License, Version 2.0.
- (d) For any new project source code, if an alternative inbound or outbound license is required for compliance with the license for a leveraged open source project (e.g. GPLv2 for Linux kernel) or is otherwise required to achieve OVS's mission, the TSC may approve the use of an alternative license for specific inbound or outbound contributions on an exception basis. Any exceptions must be approved by a majority vote of the entire TSC and must be limited in scope to what is required for such purpose. Please email [tsc@openvswitch.org](mailto:tsc@openvswitch.org) to obtain exception approval.
- (e) Subject to available funds, OVS may engage The Linux Foundation to determine the availability of, and register, trademarks, service marks, which shall be owned by the LF.

#### 9. Amendments

- (a) This charter may be amended by a two-thirds vote of the entire TSC, subject to approval by The Linux Foundation.

## 8.8 Emeritus Status for OVS Committers

OVS committers are nominated and elected based on their impact on the Open vSwitch project. Over time, as committers' responsibilities change, some may become unable or uninterested in actively participating in project governance. Committer "emeritus" status provides a way for committers to take a leave of absence from OVS governance responsibilities. The following guidelines clarify the process around the emeritus status for committers:

- A committer may choose to transition from active to emeritus, or from emeritus to active, by sending an email to the committers mailing list.
- If a committer hasn't been heard from in 6 months, and does not respond to reasonable attempts to contact him or her, the other committers can vote as a majority to transition the committer from active to emeritus. (If the committer resurfaces, he or she can transition back to active by sending an email to the committers mailing list.)

- Emeritus committers may stay on the committers mailing list to continue to follow any discussions there.
- Emeritus committers do not nominate or vote in committer elections. From a governance perspective, they are equivalent to a non-committer.
- Emeritus committers cannot merge patches to the OVS repository.
- Emeritus committers will be listed in a separate section in the MAINTAINERS.rst file to continue to recognize their contributions to the project.

Emeritus status does not replace the procedures for forcibly removing a committer.

Note that just because a committer is not able to work on the project on a day-to-day basis, we feel they are still capable of providing input on the direction of the project. No committer should feel pressured to move themselves to this status. Again, it's just an option for those that do not currently have the time or interest.

## 8.9 Expectations for Developers with Open vSwitch Repo Access

### 8.9.1 Pre-requisites

Be familiar with the guidelines and standards defined in *Contributing to Open vSwitch*.

### 8.9.2 Review

Code (yours or others') must be reviewed publicly (by you or others) before you push it to the repository. With one exception (see below), every change needs at least one review.

If one or more people know an area of code particularly well, code that affects that area should ordinarily get a review from one of them.

The riskier, more subtle, or more complicated the change, the more careful the review required. When a change needs careful review, use good judgment regarding the quality of reviews. If a change adds 1000 lines of new code, and a review posted 5 minutes later says just "Looks good," then this is probably not a quality review.

(The size of a change is correlated with the amount of care needed in review, but it is not strictly tied to it. A search and replace across many files may not need much review, but one-line optimization changes can have widespread implications.)

Your own small changes to fix a recently broken build ("make") or tests ("make check"), that you believe to be visible to a large number of developers, may be checked in without review. If you are not sure, ask for review. If you do push a build fix without review, send the patch to ovs-dev afterward as usual, indicating in the email that you have already pushed it.

Regularly review submitted code in areas where you have expertise. Consider reviewing other code as well.

### 8.9.3 Git conventions

Do not push merge commits to the Git repository without prior discussion on ovs-dev.

If you apply a change (yours or another's) then it is your responsibility to handle any resulting problems, especially broken builds and other regressions. If it is someone else's change, then you can ask the original submitter to address it. Regardless, you need to ensure that the problem is fixed in a timely way. The definition of "timely" depends on the severity of the problem.

If a bug is present on master and other branches, fix it on master first, then backport the fix to other branches. Straight-forward backports do not require additional review (beyond that for the fix on master).

Feature development should be done only on master. Occasionally it makes sense to add a feature to the most recent release branch, before the first actual release of that branch. These should be handled in the same way as bug fixes, that is, first implemented on master and then backported.

Keep the authorship of a commit clear by maintaining a correct list of “Signed-off-by:”s. If a confusing situation comes up, as it occasionally does, bring it up on the mailing list. If you explain the use of “Signed-off-by:” to a new developer, explain not just how but why, since the intended meaning of “Signed-off-by:” is more important than the syntax. As part of your explanation, quote or provide a URL to the Developer’s Certificate of Origin in [Submitting Patches](#).

Use Reported-by: and Tested-by: tags in commit messages to indicate the source of a bug report.

Keep the `AUTHORS.rst` file up to date.

## 8.10 OVS Committer Grant/Revocation Policy

An OVS committer is a participant in the project with the ability to commit code directly to the master repository. Commit access grants a broad ability to affect the progress of the project as presented by its most important artifact, the code and related resources that produce working binaries of Open vSwitch. As such it represents a significant level of trust in an individual’s commitment to working with other committers and the community at large for the benefit of the project. It can not be granted lightly and, in the worst case, must be revocable if the trust placed in an individual was inappropriate.

This document suggests guidelines for granting and revoking commit access. It is intended to provide a framework for evaluation of such decisions without specifying deterministic rules that wouldn’t be sensitive to the nuance of specific situations. In the end the decision to grant or revoke committer privileges is a judgment call made by the existing set of committers.

### 8.10.1 Granting Commit Access

Granting commit access should be considered when a candidate has demonstrated the following in their interaction with the project:

- Contribution of significant new features through the patch submission process where:
  - Submissions are free of obvious critical defects
  - Submissions do not typically require many iterations of improvement to be accepted
- Consistent participation in code review of other’s patches, including existing committers, with comments consistent with the overall project standards
- Assistance to those in the community who are less knowledgeable through active participation in project forums such as the ovs-discuss mailing list.
- Plans for sustained contribution to the project compatible with the project’s direction as viewed by current committers.
- Commitment to meet the expectations described in the “Expectations of Developer’s with Open vSwitch Access”

The process to grant commit access to a candidate is simple:

- An existing committer nominates the candidate by sending an email to all existing committers with information substantiating the contributions of the candidate in the areas described above.
- All existing committers discuss the pros and cons of granting commit access to the candidate in the email thread.
- When the discussion has converged or a reasonable time has elapsed without discussion developing (e.g. a few business days) the nominator calls for a final decision on the candidate with a followup email to the thread.

- Each committer may vote yes, no, or abstain by replying to the email thread. A failure to reply is an implicit abstention.
- After votes from all existing committers have been collected or a reasonable time has elapsed for them to be provided (e.g. a couple of business days) the votes are evaluated. To be granted commit access the candidate must receive yes votes from a majority of the existing committers and zero no votes. Since a no vote is effectively a veto of the candidate it should be accompanied by a reason for the vote.
- The nominator summarizes the result of the vote in an email to all existing committers.
- If the vote to grant commit access passed, the candidate is contacted with an invitation to become a committer to the project which asks them to agree to the committer expectations documented on the project web site.
- If the candidate agrees access is granted by setting up commit access to the repos on github.

### 8.10.2 Revoking Commit Access

When a committer behaves in a manner that other committers view as detrimental to the future of the project, it raises a delicate situation with the potential for the creation of division within the greater community. These situations should be handled with care. The process in this case is:

- Discuss the behavior of concern with the individual privately and explain why you believe it is detrimental to the project. Stick to the facts and keep the email professional. Avoid personal attacks and the temptation to hypothesize about unknowable information such as the other's motivations. Make it clear that you would prefer not to discuss the behavior more widely but will have to raise it with other contributors if it does not change. Ideally the behavior is eliminated and no further action is required. If not,
- Start an email thread with all committers, including the source of the behavior, describing the behavior and the reason it is detrimental to the project. The message should have the same tone as the private discussion and should generally repeat the same points covered in that discussion. The person whose behavior is being questioned should not be surprised by anything presented in this discussion. Ideally the wider discussion provides more perspective to all participants and the issue is resolved. If not,
- Start an email thread with all committers except the source of the detrimental behavior requesting a vote on revocation of commit rights. Cite the discussion among all committers and describe all the reasons why it was not resolved satisfactorily. This email should be carefully written with the knowledge that the reasoning it contains may be published to the larger community to justify the decision.
- Each committer may vote yes, no, or abstain by replying to the email thread. A failure to reply is an implicit abstention.
- After all votes have been collected or a reasonable time has elapsed for them to be provided (e.g. a couple of business days) the votes are evaluated. For the request to revoke commit access for the candidate to pass it must receive yes votes from two thirds of the existing committers.
- anyone that votes no must provide their reasoning, and
- if the proposal passes then counter-arguments for the reasoning in no votes should also be documented along with the initial reasons the revocation was proposed. Ideally there should be no new counter-arguments supplied in a no vote as all concerns should have surfaced in the discussion before the vote.
- The original person to propose revocation summarizes the result of the vote in an email to all existing committers excepting the candidate for removal.
- If the vote to revoke commit access passes, access is removed and the candidate for revocation is informed of that fact and the reasons for it as documented in the email requesting the revocation vote.
- Ideally the revoked committer peacefully leaves the community and no further action is required. However, there is a distinct possibility that he/she will try to generate support for his/her point of view within the larger



community. In this case the reasoning for removing commit access as described in the request for a vote will be published to the community.

### 8.10.3 Changing the Policy

The process for changing the policy is:

- Propose the changes to the policy in an email to all current committers and request discussion.
- After an appropriate period of discussion (a few days) update the proposal based on feedback if required and resend it to all current committers with a request for a formal vote.
- After all votes have been collected or a reasonable time has elapsed for them to be provided (e.g. a couple of business days) the votes are evaluated. For the request to modify the policy to pass it must receive yes votes from two thirds of the existing committers.

#### Template Emails

### 8.10.4 Nomination to Grant Commit Access

I would like to nominate *[candidate]* for commit access. I believe *[he/she]* has met the conditions for commit access described in the committer grant policy on the project web site in the following ways:

*[list of requirements & evidence]*

Please reply to all in this message thread with your comments and questions. If that discussion concludes favorably I will request a formal vote on the nomination in a few days.

### 8.10.5 Vote to Grant Commit Access

I nominated *[candidate]* for commit access on *[date]*. Having allowed sufficient time for discussion it's now time to formally vote on the proposal.

Please reply to all in this thread with your vote of: YES, NO, or ABSTAIN. A failure to reply will be counted as an abstention. If you vote NO, by our policy you must include the reasons for that vote in your reply. The deadline for votes is *[date and time]*.

If a majority of committers vote YES and there are zero NO votes commit access will be granted.

### 8.10.6 Vote Results for Grant of Commit Access

The voting period for granting to commit access to *[candidate]* initiated at *[date and time]* is now closed with the following results:

YES: *[count of yes votes]* (*[% of voters]*)

NO: *[count of no votes]* (*[% of voters]*)

ABSTAIN: *[count of abstentions]* (*[% of voters]*)

Based on these results commit access *[is/is NOT]* granted.



### 8.10.7 Invitation to Accepted Committer

Due to your sustained contributions to the Open vSwitch (OVS) project we would like to provide you with commit access to the project repository. Developers with commit access must agree to fulfill specific responsibilities described in the source repository:

*/Documentation/internals/committer-responsibilities.rst*

Please let us know if you would like to accept commit access and if so that you agree to fulfill these responsibilities. Once we receive your response we'll set up access. We're looking forward continuing to work together to advance the Open vSwitch project.

### 8.10.8 Proposal to Revoke Commit Access for Detrimental Behavior

I regret that I feel compelled to propose revocation of commit access for *[candidate]*. I have privately discussed with *[him/her]* the following reasons I believe *[his/her]* actions are detrimental to the project and we have failed to come to a mutual understanding:

*[List of reasons and supporting evidence]*

Please reply to all in this thread with your thoughts on this proposal. I plan to formally propose a vote on the proposal on or after *[date and time]*.

It is important to get all discussion points both for and against the proposal on the table during the discussion period prior to the vote. Please make it a high priority to respond to this proposal with your thoughts.

### 8.10.9 Vote to Revoke Commit Access

I nominated *[candidate]* for revocation of commit access on *[date]*. Having allowed sufficient time for discussion it's now time to formally vote on the proposal.

Please reply to all in this thread with your vote of: YES, NO, or ABSTAIN. A failure to reply will be counted as an abstention. If you vote NO, by our policy you must include the reasons for that vote in your reply. The deadline for votes is *[date and time]*.

If 2/3rds of committers vote YES commit access will be revoked.

The following reasons for revocation have been given in the original proposal or during discussion:

*[list of reasons to remove access]*

The following reasons for retaining access were discussed:

*[list of reasons to retain access]*

The counter-argument for each reason for retaining access is:

*[list of counter-arguments for retaining access]*

### 8.10.10 Vote Results for Revocation of Commit Access

The voting period for revoking the commit access of *[candidate]* initiated at *[date and time]* is now closed with the following results:

- YES: *[count of yes votes]* (*[% of voters]*)
- NO: *[count of no votes]* (*[% of voters]*)

- ABSTAIN: *[count of abstentions] ([% of voters])*

Based on these results commit access *[is/is NOT]* revoked. The following reasons for retaining commit access were proposed in NO votes:

*[list of reasons]*

The counter-arguments for each of these reasons are:

*[list of counter-arguments]*

### 8.10.11 Notification of Commit Revocation for Detrimental Behavior

After private discussion with you and careful consideration of the situation, the other committers to the Open vSwitch (OVS) project have concluded that it is in the best interest of the project that your commit access to the project repositories be revoked and this has now occurred.

The reasons for this decision are:

*[list of reasons for removing access]*

While your goals and those of the project no longer appear to be aligned we greatly appreciate all the work you have done for the project and wish you continued success in your future work.

## 8.11 Authors

The following people authored or signed off on commits in the Open vSwitch source code or webpage version control repository.

Name	Email
Aaron Conole	<a href="mailto:aconole@redhat.com">aconole@redhat.com</a>
Aaron Rosen	<a href="mailto:arosen@clermson.edu">arosen@clermson.edu</a>
Alan Pevec	<a href="mailto:alan.pevec@redhat.com">alan.pevec@redhat.com</a>
Alexander Duyck	<a href="mailto:alexander.h.duyck@redhat.com">alexander.h.duyck@redhat.com</a>
Alexandru Copot	<a href="mailto:alex.mihai.c@gmail.com">alex.mihai.c@gmail.com</a>
Alexei Starovoitov	<a href="mailto:ast@plumgrid.com">ast@plumgrid.com</a>
Alexey I. Froloff	<a href="mailto:raorn@raorn.name">raorn@raorn.name</a>
Alex Wang	<a href="mailto:ee07b291@gmail.com">ee07b291@gmail.com</a>
Alfredo Finelli	<a href="mailto:alf@computationes.de">alf@computationes.de</a>
Alin Balutoiu	<a href="mailto:abalutoiu@cloudbasesolutions.com">abalutoiu@cloudbasesolutions.com</a>
Alin Serdean	<a href="mailto:aserdean@cloudbasesolutions.com">aserdean@cloudbasesolutions.com</a>
Ambika Arora	<a href="mailto:ambika.arora@tcs.com">ambika.arora@tcs.com</a>
Amit Bose	<a href="mailto:bose@noironetworks.com">bose@noironetworks.com</a>
Amitabha Biswas	<a href="mailto:azbiswas@gmail.com">azbiswas@gmail.com</a>
Anand Kumar	<a href="mailto:kumaranand@vmware.com">kumaranand@vmware.com</a>
Andrea Kao	<a href="mailto:eirinikos@gmail.com">eirinikos@gmail.com</a>
Andreas Karis	<a href="mailto:akaris@redhat.com">akaris@redhat.com</a>
Andrew Evans	
Andrew Beekhof	<a href="mailto:abeekhof@redhat.com">abeekhof@redhat.com</a>
Andrew Kampjes	<a href="mailto:a.kampjes@gmail.com">a.kampjes@gmail.com</a>
Andrew Lambeth	<a href="mailto:alambeth@vmware.com">alambeth@vmware.com</a>
Andre McCurdy	<a href="mailto:armccurdy@gmail.com">armccurdy@gmail.com</a>
Andy Hill	<a href="mailto:hillad@gmail.com">hillad@gmail.com</a>

Continued on next page

Table 2 – continued from previous page

Name	Email
Andy Southgate	andy.southgate@citrix.com
Andy Zhou	azhou@ovn.org
Ankur Sharma	ankursharma@vmware.com
Anoob Soman	anoob.soman@citrix.com
Ansis Atteka	aatteka@vmware.com
Antonio Fischetti	antonio.fischetti@intel.com
Anupam Chanda	
Ariel Tubaltsev	atubaltsev@vmware.com
Arnoldo Lutz	arnoldo.lutz.guevara@hpe.com
Arun Sharma	arun.sharma@calsoftinc.com
Aryan TaheriMonfared	aryan.taherimonfared@uis.no
Asaf Penso	asafp@mellanox.com
Ashish Varma	ashishvarma.ovs@gmail.com
Ashwin Swaminathan	ashwinds@arista.com
Babu Shanmugam	bschanmu@redhat.com
Bala Sankaran	bsankara@redhat.com
Ben Pfaff	blp@ovn.org
Ben Warren	ben@skyportsystems.com
Benli Ye	daniely@vmware.com
Bert Vermeulen	bert@biot.com
Bhanuprakash Bodireddy	bhanuprakash.bodireddy@intel.com
Billy O'Mahony	billy.o.mahony@intel.com
Binbin Xu	xu.binbin1@zte.com.cn
Brian Kruger	bkruger+ovsdev@gmail.com
Bruce Davie	bdavie@vmware.com
Bryan Phillippe	bp@toroki.com
Carlo Andreotti	c.andreotti@m3s.it
Casey Barker	crbarker@google.com
Chandra Sekhar Vejendla	csvejend@us.ibm.com
Christoph Jaeger	cj@linux.com
Chris Wright	chrisw@sous-sol.org
Chuck Short	zulcss@ubuntu.com
Ciara Loftus	ciara.loftus@intel.com
Clint Byrum	clint@fewbar.com
Cong Wang	amwang@redhat.com
Conner Herriges	conner.herriges@ibm.com
Damien Millescamp	damien.millescamp@6wind.com
Dan Carpenter	dan.carpenter@oracle.com
Dan McGregor	dan.mcgregor@usask.ca
Dan Wendlandt	
Dan Williams	dcbw@redhat.com
Daniel Alvarez	dalvarez@redhat.com
Daniel Borkmann	dborkman@redhat.com
Daniel Hiltgen	daniel@netkine.com
Daniel Roman	
Daniele Di Proietto	daniele.di.proietto@gmail.com
Daniele Venturino	venturino.daniele+ovs@gmail.com
Danny Kukawka	danny.kukawka@bisect.de
Darrell Ball	dlu998@gmail.com

Continued on next page

Table 2 – continued from previous page

Name	Email
Dave Tucker	dave@dtucker.co.uk
David Erickson	derickso@stanford.edu
David Hill	dhill@redhat.com
David Marchand	david.marchand@redhat.com
David S. Miller	davem@davemloft.net
David Yang	davidy@vmware.com
Dennis Sam	dsam@arista.com
Devendra Naga	devendra.aaru@gmail.com
Dmitry Krivenok	krivenok.dmitry@gmail.com
Dominic Curran	dominic.curran@citrix.com
Dongdong	dongdong1@huawei.com
Dongjun	dongj@dtream.com
Duan Jiong	djduanjiong@gmail.com
Duffie Cooley	
Dustin Lundquist	dustin@null-ptr.net
Ed Maste	emaste@freebsd.org
Ed Swierk	eswierk@skyportsystems.com
Edouard Bourguignon	madko@linuxed.net
Eelco Chaudron	echaudro@redhat.com
Eric Lapointe	elapointe@corsa.com
Esteban Rodriguez Betancourt	estebarb@hpe.com
Aymerich Edward	edward.aymerich@hpe.com
Edward Tomasz Napierała	trasz@freebsd.org
Eitan Eliahu	eliahue@vmware.com
Eohyung Lee	liquidnuker@gmail.com
Eric Dumazet	edumazet@google.com
Eric Garver	e@erig.me
Eric Sesterhenn	eric.sesterhenn@lsexperts.de
Ethan J. Jackson	ejj@eecs.berkeley.edu
Ethan Rahn	erahn@arista.com
Eziz Durdyev	ezizurdy@gmail.com
Flavio Fernandes	flavio@flaviof.com
Flavio Leitner	fbl@redhat.com
Francesco Fusco	ffusco@redhat.com
Frédéric Tobias Christ	fchrist@live.de
Frode Nordahl	frode.nordahl@gmail.com
FUJITA Tomonori	fujita.tomonori@lab.ntt.co.jp
Gabe Begeg-Dov	gabe@begegdov.com
Gaetano Catalli	gaetano.catalli@gmail.com
Gal Sagie	gal.sagie@gmail.com
Genevieve LEsperance	glesperance@pivotal.io
Geoffrey Wossum	gwossum@acm.org
Gianluca Merlo	gianluca.merlo@gmail.com
Giuseppe Lettieri	g.lettieri@iet.unipi.it
Glen Gibb	grg@stanford.edu
Guoshuai Li	ligs@dtream.com
Guolin Yang	gyang@vmware.com
Guru Chaitanya Perakam	gperakam@Brocade.com
Gurucharan Shetty	guru@ovn.org

Continued on next page

Table 2 – continued from previous page

Name	Email
Han Zhou	zhouhan@gmail.com
Henry Mai	
Hao Zheng	
Helmut Schaa	helmut.schaa@gmail.com
Hiteshi Kalra	hiteshi.kalra@tcs.com
Huanle Han	hanxueluo@gmail.com
Hui Kang	kangh@us.ibm.com
Hyong Youb Kim	hyonkim@cisco.com
Ian Campbell	Ian.Campbell@citrix.com
Ian Stokes	ian.stokes@intel.com
Ilya Maximets	i.maximets@samsung.com
Iman Tabrizian	tabrizian@outlook.com
Isaku Yamahata	yamahata@valinux.co.jp
Ivan Dyukov	i.dyukov@samsung.com
IWASE Yusuke	iwase.yusuke@gmail.com
Jakub Libosvar	libosvar@redhat.com
Jakub Sitnicki	jsitnicki@gmail.com
James P.	roampune@gmail.com
James Page	james.page@ubuntu.com
Jamie Lennox	jamielennox@gmail.com
Jan Scheurich	jan.scheurich@ericsson.com
Jan Vansteenkiste	jan@vstone.eu
Jarno Rajahalme	jarno@ovn.org
Jason Kölker	jason@koelker.net
Jason Wessel	jason.wessel@windriver.com
Jasper Capel	jasper@capel.tv
Jean Tourrilhes	jt@hpl.hp.com
Jeremy Stribling	
Jeroen van Bommel	jvb127@gmail.com
Jesse Gross	jesse@kernel.org
Jian Li	lijian@ooclab.com
Jianbo Liu	jianbol@mellanox.com
Jing Ai	jinga@google.com
Jiri Benc	jbenc@redhat.com
Joe Perches	joe@perches.com
Joe Stringer	joe@ovn.org
Jonathan Vestin	jonavest@kau.se
Jorge Arturo Sauma Vargas	jorge.sauma@hpe.com
Jun Nakajima	jun.nakajima@intel.com
Junhan Yan	juyan@redhat.com
JunoZhu	zhunatu@zhu@gmail.com
Justin Pettit	jpettit@ovn.org
Kaige Fu	fukaige@huawei.com
Keith Amidon	
Ken Ajiro	ajiro@mxw.nes.nec.co.jp
Ken Sanislo	ken@intherack.com
Kenneth Duda	kduda@arista.com
Kentaro Ebisawa	ebiken.g@gmail.com
Keshav Gupta	keshav.gupta@ericsson.com

Continued on next page

Table 2 – continued from previous page

Name	Email
Kevin Lo	kevlo@FreeBSD.org
Kevin Traynor	kevin.traynor@intel.com
Khem Raj	raj.khem@gmail.com
Kmindg G	kmindg@gmail.com
Kris Murphy	kriskend@linux.vnet.ibm.com
Krishna Kondaka	kkondaka@vmware.com
Kyle Mestery	mestery@mestery.com
Kyle Simpson	kyleandrew.simpson@gmail.com
Kyle Upton	kupton@baymicrosystems.com
Lance Richardson	lrichard@redhat.com
Lars Kellogg-Stedman	lars@redhat.com
Lei Huang	huang.f.lei@gmail.com
Leif Madsen	lmadsen@redhat.com
Leo Alterman	
Li RongQing	lirongqing@baidu.com
Lian-min Wang	liang-min.wang@intel.com
Lilijun	jerry.lilijun@huawei.com
Lili Huang	huanglili.huang@huawei.com
Linda Sun	lsun@vmware.com
Lior Neudorfer	lior@guardicore.com
Lorand Jakab	lojakab@cisco.com
Lorenzo Bianconi	lorenzo.bianconi@redhat.com
Luca Giraud	
Lucas Alvares Gomes	lucasagomes@gmail.com
Lucian Petrut	lpetrut@cloudbasesolutions.com
Luigi Rizzo	rizzo@iet.unipi.it
Luis E. P.	l31g@hotmail.com
Lukasz Rzasik	lukasz.rzasik@gmail.com
Madhu Challa	challa@noironetworks.com
Manohar K C	manukc@gmail.com
Marcin Mirecki	mmirecki@redhat.com
Mario Cabrera	mario.cabrera@hpe.com
Mark D. Gray	mark.d.gray@intel.com
Mark Hamilton	
Mark Kavanagh	mark.b.kavanagh81@gmail.com
Mark Maglana	mmaglana@gmail.com
Mark Michelson	mmichels@redhat.com
Markos Chandras	mchandras@suse.de
Martin Casado	casado@cs.stanford.edu
Martin Fong	mwfong@csl.sri.com
Martino Fornasa	mf@fornasa.it
Martin Xu	martinxu9.ovs@gmail.com
Maryam Tahhan	maryam.tahhan@intel.com
Matteo Croce	mcroce@redhat.com
Mauricio Vásquez	mauricio.vasquezbernal@studenti.polito.it
Maxime Coquelin	maxime.coquelin@redhat.com
Mehak Mahajan	
Michael Arnaldi	arnaldimichael@gmail.com
Michal Weglicki	michalx.weglicki@intel.com

Continued on next page

Table 2 – continued from previous page

Name	Email
Mickey Spiegel	mickeys.dev@gmail.com
Miguel Angel Ajo	majopela@redhat.com
Mijo Safradin	mijo@linux.vnet.ibm.com
Mika Vaisanen	mika.vaisanen@gmail.com
Minoru TAKAHASHI	takahashi.minoru7@gmail.com
Murphy McCauley	murphy.mccauley@gmail.com
Natasha Gude	
Neal Shrader	neal@digitalocean.com
Neil McKee	neil.mckee@inmon.com
Neil Zhu	zhuj@centecnetworks.com
Nimay Desai	nimaydesai1@gmail.com
Nithin Raju	nithin@vmware.com
Niti Rohilla	niti.rohilla@tcs.com
Nitin Katiyar	nitin.katiyar@ericsson.com
Numan Siddique	nusiddiq@redhat.com
Ofer Ben-Yacov	ofer.benyacov@gmail.com
Ophir Munk	ophirmu@mellanox.com
Or Gerlitz	ogerlitz@mellanox.com
Ori Shoshan	ori.shoshan@guardicore.com
Padmanabhan Krishnan	kprad1@yahoo.com
Panu Matilainen	pmatilai@redhat.com
Paraneetharan Chandrasekaran	paraneetharanc@gmail.com
Paul Boca	pboca@cloudbasesolutions.com
Paul Fazzone	pfazzone@vmware.com
Paul Ingram	
Paul-Emmanuel Raoul	skyper@skylabs.net
Pavithra Ramesh	paramesh@vmware.com
Peter Downs	padowns@gmail.com
Philippe Jung	phil.jung@free.fr
Pim van den Berg	pim@nethuis.nl
pritesh	pritesh.kothari@cisco.com
Pravin B Shelar	pshelar@ovn.org
Przemysław Szczerbik	przemyslawx.szczerbik@intel.com
Quentin Monnet	quentin.monnet@6wind.com
Qiuyu Xiao	qiuyu.xiao.qyx@gmail.com
Raju Subramanian	
Rami Rosen	ramirose@gmail.com
Ramu Ramamurthy	ramu.ramamurthy@us.ibm.com
Randall Sharo	andall.sharo@navy.mil
Ravi Kerur	Ravi.Kerur@telekom.com
Raymond Burkholder	ray@oneunified.net
Reid Price	
Remko Tronçon	git@el-tramo.be
Rich Lane	rlane@bigswitch.com
Richard Oliver	richard@richard-oliver.co.uk
Rishi Bamba	rishi.bamba@tcs.com
Rob Adams	readams@readams.net
Robert Åkerblom-Andersson	Robert.nr1@gmail.com
Robert Wojciechowicz	robertx.wojciechowicz@intel.com

Continued on next page

Table 2 – continued from previous page

Name	Email
Rob Hoes	rob.hoes@citrix.com
Rohith Basavaraja	rohith.basavaraja@gmail.com
Roi Dayan	roid@mellanox.com
Róbert Mulik	robert.mulik@ericsson.com
Romain Lenglet	romain.lenglet@berabera.info
Russell Bryant	russell@ovn.org
RYAN D. MOATS	rmoats@us.ibm.com
Ryan Wilson	
Sairam Venugopal	vsairam@vmware.com
Sajjad Lateef	
Saloni Jain	saloni.jain@tcs.com
Samuel Ghinet	sghinet@cloudbasesolutions.com
Sanjay Sane	
Saurabh Mohan	saurabh@cplanetworks.com
Saurabh Shah	
Saurabh Shrivastava	saurabh.shrivastava@nuagenetworks.net
Scott Cheloha	scottcheloha@gmail.com
Scott Lowe	scott.lowe@scottlowe.org
Scott Mann	sdmnix@gmail.com
Selvamuthukumar	smkumar@merunetworks.com
Sha Zhang	zhangsha.zhang@huawei.com
Shad Ansari	shad.ansari@hpe.com
Shan Wei	davidshan@tencent.com
Shashank Ram	rams@vmware.com
Shashwat Srivastava	shashwat.srivastava@tcs.com
Shih-Hao Li	shihli@vmware.com
Shu Shen	shu.shen@radisys.com
Simon Horman	horms@verge.net.au
Simon Horman	simon.horman@netronome.com
Sorin Vinturis	svinturis@cloudbasesolutions.com
Steffen Gebert	steffen.gebert@informatik.uni-wuerzburg.de
Sten Spans	sten@blinkerlights.nl
Stephane A. Sezer	sas@cd80.net
Stephen Finucane	stephen@that.guru
Steve Ruan	ruansx@cn.ibm.com
Stuart Cardall	developer@it-offshore.co.uk
Sugesh Chandran	sugesh.chandran@intel.com
SUGYO Kazushi	sugyo.org@gmail.com
Tadaaki Nagao	nagao@stratosphere.co.jp
Terry Wilson	twilson@redhat.com
Tetsuo NAKAGAWA	nakagawa@mxn.nes.nec.co.jp
Thadeu Lima de Souza Cascardo	cascardo@cascardo.eti.br
Thomas F. Herbert	thomasfherbert@gmail.com
Thomas Goirand	zigo@debian.org
Thomas Graf	tgraf@noironetworks.com
Thomas Lacroix	thomas.lacroix@citrix.com
Timo Puha	timox.puha@intel.com
Timothy Redaelli	tredaelli@redhat.com
Todd Deshane	deshantm@gmail.com

Continued on next page



Table 2 – continued from previous page

Name	Email
Tom Everman	teverman@google.com
Toms Atteka	cpp.code.lv@gmail.com
Torgny Lindberg	torgny.lindberg@ericsson.com
Tsvi Slonim	tsvi@toroki.com
Tuan Nguyen	tuan.nguyen@veriksystems.com
Tyler Coumbes	coumbes@gmail.com
Tony van der Peet	tony.vanderpeet@alliedtelesis.co.nz
Tonghao Zhang	xiangxia.m.yue@gmail.com
Valient Gough	vgough@pobox.com
Venkata Anil Kommaddi	vkommadi@redhat.com
Vishal Deep Ajmera	vishal.deep.ajmera@ericsson.com
Vivien Bernet-Rollande	vbr@soprive.net
wangqianyu	wang.qianyu@zte.com.cn
Wang Sheng-Hui	shhuiw@gmail.com
Wang Zhike	wangzhike@jd.com
Wei Li	liw@dt dream.com
Wei Yongjun	yjwei@cn.fujitsu.com
Wenyu Zhang	wenyuz@vmware.com
William Fulton	
William Tu	u9012063@gmail.com
Xiao Liang	shaw.leon@gmail.com
xu rong	xu.rong@zte.com.cn
YAMAMOTO Takashi	yamamoto@midokura.com
Yasuhito Takamiya	yasuhito@gmail.com
Yi-Hung Wei	yihung.wei@gmail.com
Yifeng Sun	pkusunyifeng@gmail.com
Yin Lin	linyi@vmware.com
Yu Zhiguo	yuzg@cn.fujitsu.com
Yuanhan Liu	yuanhan.liu@linux.intel.com
Yunjian Wang	wangyunjian@huawei.com
Yousong Zhou	yszhou4tech@gmail.com
Zak Whittington	zwhitt.vmware@gmail.com
ZhengLingyun	konghuarukhr@163.com
Zoltán Balogh	zoltan.balogh.eth@gmail.com
Zoltan Kiss	zoltan.kiss@citrix.com
Zongkai LI	zealokii@gmail.com
Zhi Yong Wu	zwu.kernel@gmail.com
Zang MingJie	zealot0630@gmail.com
Zhenyu Gao	sysugaozhenyu@gmail.com
ZhiPeng Lu	luzhipeng@uniudc.com
Zhou Yangchao	1028519445@qq.com
aginwala	amginwal@gmail.com
solomon	liwei.solomon@gmail.com
wenxu	wenxu@ucloud.cn
wisd0me	ak47izatoool@gmail.com
xushengping	shengping.xu@huawei.com
yinpeijun	yinpeijun@huawei.com
zangchuanqiang	zangchuanqiang@huawei.com
zhaojingjing	zhao.jingjing1@zte.com.cn

Continued on next page

Table 2 – continued from previous page

Name	Email
zhongbaisong	zhongbaisong@huawei.com
zhaozhanxu	zhaozhanxu@163.com

The following additional people are mentioned in commit logs as having provided helpful bug reports or suggestions.

Name	Email
Aaron M. Ucko	ucko@debian.org
Abhinav Singhal	Abhinav.Singhal@spirent.com
Adam Heath	doogie@brainfood.com
Ahmed Bilal	numan252@gmail.com
Alan Kayahan	hsykay@gmail.com
Alan Shieh	
Alban Browaeys	prahal@yahoo.com
Alex Yip	
Alexey I. Froloff	raorn@altlinux.org
Amar Padmanabhan	
Amey Bhide	
Amre Shakimov	ashakimov@vmware.com
André Ruß	andre.russ@hybris.com
Andreas Beckmann	debian@abeckmann.de
Andrei Andone	andrei.andone@softvision.ro
Andrey Korolyov	andrey@xdel.ru
Anil Jangam	anilj.mailing@gmail.com
Anshuman Manral	anshuman.manral@outlook.com
Anton Matsiuk	anton.matsiuk@gmail.com
Anup Khadka	khadka.py@gmail.com
Anuprem Chalvadi	achalvadi@vmware.com
Ariel Tubaltsev	atubaltsev@vmware.com
Arkajit Ghosh	arkajit.ghosh@tcs.com
Atzm Watanabe	atzm@stratosphere.co.jp
Aurélien Poulain	aurepoulain@viacesi.fr
Bastian Blank	waldi@debian.org
Ben Basler	
Bhargava Shastry	bshastry@sec.t-labs.tu-berlin.de
Bob Ball	bob.ball@citrix.com
Brad Hall	
Brad Cowie	brad@wand.net.nz
Brailey Josh	josh@faucet.nz
Brandon Heller	brandonh@stanford.edu
Brendan Kelley	
Brent Salisbury	brent.salisbury@gmail.com
Brian Field	Brian_Field@cable.comcast.com
Bryan Fulton	
Bryan Osoro	
Cedric Hobbs	
Chris Hydon	chydon@aristanetworks.com
Christian Stigen Larsen	cslarsen@gmail.com
Christopher Paggen	cpaggen@cisco.com
Chunhe Li	lichunhe@huawei.com

Continued on next page

Table 3 – continued from previous page

Name	Email
Daniel Badea	daniel.badea@windriver.com
Darragh O'Reilly	darragh.oreilly@hpe.com
Dave Walker	DaveWalker@ubuntu.com
David Evans	davidjoshuaevans@gmail.com
David Palma	palma@onesource.pt
David van Moolenbroek	dvmoelenbroek@aimvalley.nl
Derek Cormier	derek.cormier@lab.ntt.co.jp
Dhaval Badiani	dbadiani@vmware.com
DK Moon	
Ding Zhi	zhi.ding@6wind.com
Dong Jun	dongj@dtdream.com
Dustin Spinhirne	dspinhirne@vmware.com
Edwin Chiu	echiu@vmware.com
Eivind Bulie Haanaes	
Enas Ahmad	enas.ahmad@kaust.edu.sa
Eric Lopez	
Frido Roose	fr.roose@gmail.com
Gaetano Catalli	gaetano.catalli@gmail.com
Gavin Remaley	gavin_remaley@selinc.com
Georg Schmuecking	georg.schmuecking@ericsson.com
George Shuklin	amarao@desunote.ru
Gerald Rogers	gerald.rogers@intel.com
Ghanem Bahri	bahri.ghanem@gmail.com
Giuseppe de Candia	giuseppe.decandia@gmail.com
Gordon Good	ggood@vmware.com
Greg Dahlman	gdahlman@hotmail.com
Greg Rose	gvrose8192@gmail.com
Gregor Schaffrath	grsch@net.t-labs.tu-berlin.de
Gregory Smith	gasmith@nutanix.com
Guolin Yang	gyang@vmware.com
Gur Stavi	gstavi@mrv.com
Harish Kanakaraju	hkanakaraju@vmware.com
Hari Sasank Bhamidipalli	hbhamidi@cisco.com
Hassan Khan	hassan.khan@seecs.edu.pk
Hector Oron	hector.oron@gmail.com
Hemanth Kumar Mantri	mantri@nutanix.com
Henrik Amren	
Hiroshi Tanaka	
Hiroshi Miyata	miyahiro.dazu@gmail.com
Hsin-Yi Shen	shenh@vmware.com
Hui Xiang	xianghuir@gmail.com
Hyojoon Kim	joonk@gatech.edu
Igor Ganichev	
Igor Sever	igor@xorops.com
Jacob Cherkas	cherkasj@vmware.com
Jad Naous	jnaous@gmail.com
Jamal Hadi Salim	hadi@cyberus.ca
James Schmidt	jschmidt@vmware.com
Jan Medved	jmedved@juniper.net

Continued on next page

Table 3 – continued from previous page

Name	Email
Janis Hamme	janis.hamme@student.kit.edu
Jari Sundell	sundell.software@gmail.com
Javier Albornoz	javier.albornoz@hpe.com
Jed Daniels	openvswitch@jeddaniels.com
Jeff Merrick	jmerrick@vmware.com
Jeongkeun Lee	jkleee@hp.com
Jian Qiu	swordqiu@gmail.com
Joan Cirer	joan@ev0.net
John Darrington	john@darrington.wattle.id.au
John Galgay	john@galgay.net
John Hurley	john.hurley@netronome.com
John Reumann	nofutznetworks@gmail.com
Karthik Sundaravel	ksundara@redhat.com
Kashyap Thimmaraju	kashyap.thimmaraju@sec.t-labs.tu-berlin.de
Keith Holleman	hollemanietf@gmail.com
Kevin Lin	kevinlin@berkeley.edu
K	k940545@hotmail.com
Kevin Mancuso	kevin.mancuso@rackspace.com
Kiran Shanbhog	kiran@vmware.com
Kirill Kabardin	
Kirkland Spector	kspector@salesforce.com
Koichi Yagishita	yagishita.koichi@jrc.co.jp
Konstantin Khorenko	khorenko@openvz.org
Kris zhang	zhang.kris@gmail.com
Krishna Miriyala	miriyalak@vmware.com
Krishna Mohan Elluru	elluru.kri.mohan@hpe.com
László Süri	laszlo.suru@ericsson.com
Len Gao	leng@vmware.com
Logan Rosen	logatronico@gmail.com
Luca Falavigna	dktrkranz@debian.org
Luiz Henrique Ozaki	luiz.ozaki@gmail.com
Manpreet Singh	er.manpreet25@gmail.com
Marco d'Itri	md@Linux.IT
Martin Vizvary	vizvary@ics.muni.cz
Marvin Pascual	marvin@pascual.com.ph
Maxime Brun	m.brun@alphalink.fr
Madhu Venugopal	mavenugo@gmail.com
Michael A. Collins	mike.a.collins@ark-net.org
Michael Ben-Ami	mbenami@digitalocean.com
Michael Hu	humichael@vmware.com
Michael J. Smalley	michaeljsmalley@gmail.com
Michael Mao	
Michael Shigorin	mike@osdn.org.ua
Michael Stapelberg	stapelberg@debian.org
Mihir Gangar	gangarm@vmware.com
Mike Bursell	mike.bursell@citrix.com
Mike Kruze	
Mike Qing	mqing@vmware.com
Min Chen	ustcer.tonychan@gmail.com

Continued on next page

Table 3 – continued from previous page

Name	Email
Mikael Doverhag	
Mircea Ulinic	ping@mirceaulinic.net
Mrinmoy Das	mradas@ixiacom.com
Muhammad Shahbaz	mshahbaz@cs.princeton.edu
Murali R	muralirdev@gmail.com
Nagi Reddy Jonnala	njonnala@Brocade.com
Niels van Adrichem	N.L.M.vanAdrichem@tudelft.nl
Niklas Andersson	
Oscar Wilde	xdxiaobin@gmail.com
Pankaj Thakkar	pthakkar@vmware.com
Pasi Kärkkäinen	pasik@iki.fi
Patrik Andersson R	patrik.r.andersson@ericsson.com
Paul Greenberg	
Paulo Cravero	pcravero@as2594.net
Pawan Shukla	shuklap@vmware.com
Periyasamy Palanisamy	periyasamy.palanisamy@ericsson.com
Peter Amidon	peter@picnicpark.org
Peter Balland	
Peter Phaal	peter.phaal@inmon.com
Prabina Pattnaik	Prabina.Pattnaik@nechclst.in
Pratap Reddy	
Ralf Heiringhoff	ralf@frosty-geek.net
Ram Jothikumar	
Ramana Reddy	gtvrreddy@gmail.com
Ray Li	rayli1107@gmail.com
Richard Theis	rtheis@us.ibm.com
RishiRaj Maulick	rishi.raj2509@gmail.com
Rob Sherwood	rob.sherwood@bigswitch.com
Robert Strickler	anomalyst@gmail.com
Roger Leigh	rleigh@codelibre.net
Rogério Vinhal Nunes	
Roman Sokolkov	rsokolkov@gmail.com
Ronaldo A. Ferreira	ronaldof@CS.Princeton.EDU
Ronny L. Bull	bullrl@clarkson.edu
Sandeep Kumar	sandeep.kumar16@tcs.com
Sander Eikelenboom	linux@eikelenboom.it
Saul St. John	sstjohn@cs.wisc.edu
Scott Hendricks	
Sean Brady	sbrady@gtfsservices.com
Sebastian Andrzej Siewior	sebastian@breakpoint.cc
Sébastien RICCIO	sr@swisscenter.com
Shweta Seth	shwseth@cisco.com
Simon Jouet	simon.jouet@gmail.com
Spiro Kourtessis	spiro@vmware.com
Sridhar Samudrala	samudrala.sridhar@gmail.com
Srini Seetharaman	seethara@stanford.edu
Sabyasachi Sengupta	Sabyasachi.Sengupta@alcatel-lucent.com
Salvatore Cambria	salvatore.cambria@citrix.com
Soner Sevinc	sevincs@vmware.com

Continued on next page

Table 3 – continued from previous page

Name	Email
Stepan Andrushko	stepanx.andrushko@intel.com
Stephen Hemminger	shemminger@vyatta.com
Stuart Cardall	developer@it-offshore.co.uk
Suganya Ramachandran	suganyar@vmware.com
Sundar Nadathur	undar.nadathur@intel.com
Taekho Nam	thnam@smartx.kr
Takayuki HAMA	t-hama@cb.jp.nec.com
Teemu Koponen	
Thomas Morin	thomas.morin@orange.com
Timothy Chen	
Torbjorn Tornkvist	kruskakli@gmail.com
Tulio Ribeiro	tribeiro@lasige.di.fc.ul.pt
Tytus Kurek	Tytus.Kurek@pega.com
Valentin Bud	valentin@hackaserver.com
Vasiliy Tolstov	v.tolstov@selfip.ru
Vasu Dasari	vdasari@gmail.com
Vinllen Chen	cvinllen@gmail.com
Vishal Swarnkar	vishal.swarnkar@gmail.com
Vjekoslav Brajkovic	balkan@cs.washington.edu
Voravit T.	voravit@kth.se
Yeming Zhao	zhaoyeming@gmail.com
Yi Ba	yby.developer@yahoo.com
Ying Chen	yingchen@vmware.com
Yongqiang Liu	liuyq7809@gmail.com
ZHANG Zhiming	zhangzhiming@yunshan.net.cn
Zhangguanghui	zhang.guanghui@h3c.com
Ziyou Wang	ziyouw@vmware.com
ankur dwivedi	ankurengg2003@gmail.com
chen zhang	3zhangchen9211@gmail.com
james hopper	jameshopper@email.com
kk yap	yapkke@stanford.edu
likunyun	kunyunli@hotmail.com
meishengxin	meishengxin@huawei.com
neeraj mehta	mehtaneeraj07@gmail.com
rahim entezari	rahim.entezari@gmail.com
shaoke xi	xishaoke.xsk@gmail.com
shivani dommeti	shivani.dommeti@gmail.com
weizj	34965317@qq.com
	zhaojun12@outlook.com
(Crab)	fqs888@126.com
	fortitude.zhang@gmail.com
	hujingfei914@msn.com
	zhangwqh@126.com
	zhangqiang@meizu.com

Thanks to all Open vSwitch contributors. If you are not listed above but believe that you should be, please write to [dev@openvswitch.org](mailto:dev@openvswitch.org).

## 8.12 Committers

Open vSwitch committers are the people who have been granted access to push changes to the Open vSwitch git repository.

The responsibilities of an Open vSwitch committer are documented [here](#).

The process for adding or removing committers is documented [here](#).

This is the current list of active Open vSwitch committers:

Table 4: OVS Maintainers

Name	Email
Alex Wang	ee07b291@gmail.com
Alin Serdean	aserdean@cloudbasesolutions.com
Andy Zhou	azhou@ovn.org
Ansis Atteka	aatteka@nicira.com
Ben Pfaff	blp@ovn.org
Daniele Di Proietto	daniele.di.proietto@gmail.com
Gurucharan Shetty	guru@ovn.org
Ian Stokes	istokes@ovn.org
Jarno Rajahalme	jarno@ovn.org
Jesse Gross	jesse@kernel.org
Joe Stringer	joe@ovn.org
Justin Pettit	jpettit@ovn.org
Pravin B Shelar	pshelar@ovn.org
Russell Bryant	russell@ovn.org
Simon Horman	horms@ovn.org
Thomas Graf	tgraf@noironetworks.com
YAMAMOTO Takashi	yamamoto@midokura.com

The project also maintains a list of Emeritus Committers (or Maintainers). More information about Emeritus Committers can be found [here](#).

Table 5: OVS Emeritus Maintainers

Name	Email
Ethan J. Jackson	ejj@eecs.berkeley.edu

## 8.13 How Open vSwitch’s Documentation Works

This document provides a brief overview on how the documentation build system within Open vSwitch works. This is intended to maximize the “bus factor” and share best practices with other projects.

### 8.13.1 reStructuredText and Sphinx

Nearly all of Open vSwitch’s documentation is written in [reStructuredText](#), with man pages being the sole exception. Of this documentation, most of it is fed into [Sphinx](#), which provides not only the ability to convert rST to a variety of other output formats but also allows for things like cross-referencing and indexing. for more information on the two, refer to the [Open vSwitch Documentation Style](#).

### 8.13.2 ovs-sphinx-theme

The documentation uses its own theme, *ovs-sphinx-theme*, which can be found on [GitHub](#) and is published on [pypi](#). This is packaged separately from Open vSwitch itself to ensure all documentation gets the latest version of the theme (assuming there are no major version bumps in that package). If building locally and the package is installed, it will be used. If the package is not installed, Sphinx will fallback to the default theme.

The package is currently maintained by Stephen Finucane and Russell Bryant.

### 8.13.3 Read the Docs

The documentation is hosted on [readthedocs.org](#) and a CNAME redirect is in place to allow access from [docs.openvswitch.org](#). *Read the Docs* provides a couple of nifty features for us, such as automatic building of docs whenever there are changes and versioning of documentation.

The *Read the Docs* project is currently maintained by Stephen Finucane, Russell Bryant and Ben Pfaff.

### 8.13.4 openvswitch.org

The sources for [openvswitch.org](#) are maintained separately from [docs.openvswitch.org](#). For modifications to this site, refer to the [GitHub project](#).



## Symbols

-V, -version  
    ovs-test command line option, 282  
    ovs-vlan-test command line option, 283

-b <targetbandwidth>, -bandwidth <targetbandwidth>  
    ovs-test command line option, 282

-c <server1> <server2>, -client <server1> <server2>  
    ovs-test command line option, 281

-d, -direct  
    ovs-test command line option, 282

-h, -help  
    ovs-test command line option, 282  
    ovs-vlan-test command line option, 283

-i <testinterval>, -interval <testinterval>  
    ovs-test command line option, 282

-l <vlantag>, -vlan-tag <vlantag>  
    ovs-test command line option, 282

-s <port>, -server <port>  
    ovs-test command line option, 281

-s, -server  
    ovs-vlan-test command line option, 283

-t <tunnelmodes>, -tunnel-modes <tunnelmodes>  
    ovs-test command line option, 282

-V, -version, 283  
-h, -help, 283  
-s, -server, 283

## O

ovs-test command line option  
    -V, -version, 282  
    -b <targetbandwidth>, -bandwidth <targetbandwidth>, 282  
    -c <server1> <server2>, -client <server1> <server2>, 281  
    -d, -direct, 282  
    -h, -help, 282  
    -i <testinterval>, -interval <testinterval>, 282  
    -l <vlantag>, -vlan-tag <vlantag>, 282  
    -s <port>, -server <port>, 281  
    -t <tunnelmodes>, -tunnel-modes <tunnelmodes>, 282

ovs-vlan-test command line option