



Agilio OvS 2.2 User's Guide

Netronome Intelligent Server Adapters

- Proprietary and Confidential -

Agilio OvS 2.2 User's Guide: Netronome Intelligent Server Adapters

Copyright © 2011-2016 Netronome

COPYRIGHT

No part of this publication or documentation accompanying this Product may be reproduced in any form or by any means or used to make any derivative work by any means including but not limited to by translation, transformation or adaptation without permission from Netronome Systems, Inc., as stipulated by the United States Copyright Act of 1976. Contents are subject to change without prior notice.

WARRANTY

Netronome warrants that any media on which this documentation is provided will be free from defects in materials and workmanship under normal use for a period of ninety (90) days from the date of shipment. If a defect in any such media should occur during this 90-day period, the media may be returned to Netronome for a replacement.

NETRONOME DOES NOT WARRANT THAT THE DOCUMENTATION SHALL BE ERROR-FREE. THIS LIMITED WARRANTY SHALL NOT APPLY IF THE DOCUMENTATION OR MEDIA HAS BEEN (I) ALTERED OR MODIFIED; (II) SUBJECTED TO NEGLIGENCE, COMPUTER OR ELECTRICAL MALFUNCTION; OR (III) USED, ADJUSTED, OR INSTALLED OTHER THAN IN ACCORDANCE WITH INSTRUCTIONS FURNISHED BY NETRONOME OR IN AN ENVIRONMENT OTHER THAN THAT INTENDED OR RECOMMENDED BY NETRONOME.

EXCEPT FOR WARRANTIES SPECIFICALLY STATED IN THIS SECTION, NETRONOME HEREBY DISCLAIMS ALL EXPRESS OR IMPLIED WARRANTIES OF ANY KIND, INCLUDING, WITHOUT LIMITATION, ANY IMPLIED WARRANTIES OF MERCHANTABILITY, NON-INFRINGEMENT AND FITNESS FOR A PARTICULAR PURPOSE.

Some jurisdictions do not allow the exclusion of implied warranties, so the above exclusion may not apply to some users of this documentation. This limited warranty gives users of this documentation specific legal rights, and users of this documentation may also have other rights which vary from jurisdiction to jurisdiction.

LIABILITY

Regardless of the form of any claim or action, Netronome's total liability to any user of this documentation for all occurrences combined, for claims, costs, damages or liability based on any cause whatsoever and arising from or in connection with this documentation shall not exceed the purchase price (without interest) paid by such user.

IN NO EVENT SHALL NETRONOME OR ANYONE ELSE WHO HAS BEEN INVOLVED IN THE CREATION, PRODUCTION, OR DELIVERY OF THE DOCUMENTATION BE LIABLE FOR ANY LOSS OF DATA, LOSS OF PROFITS OR LOSS OF USE OF THE DOCUMENTATION OR FOR ANY SPECIAL, INCIDENTAL, CONSEQUENTIAL, EXEMPLARY, PUNITIVE, MULTIPLE OR OTHER DAMAGES, ARISING FROM OR IN CONNECTION WITH THE DOCUMENTATION EVEN IF NETRONOME HAS BEEN MADE AWARE OF THE POSSIBILITY OF SUCH DAMAGES. IN NO EVENT SHALL NETRONOME OR ANYONE ELSE WHO HAS BEEN INVOLVED IN THE CREATION, PRODUCTION, OR DELIVERY OF THE DOCUMENTATION BE LIABLE TO ANYONE FOR ANY CLAIMS, COSTS, DAMAGES OR LIABILITIES CAUSED BY IMPROPER USE OF THE DOCUMENTATION OR USE WHERE ANY PARTY HAS SUBSTITUTED PROCEDURES NOT SPECIFIED BY NETRONOME.

Revision History

Date	Revision	Description
24 February 2016	000	First DRAFT for Agilio OvS 2.2.
6 April 2016	001	Added virtiorelay configuration variables, VF rate limiting, VM QoS and PXE boot feature descriptions.
19 July 2016	002	Updated documentation for added features in the final release. Also added 'Use Case' section with examples.
29 July 2016	003	Addressed review comments.

Table of Contents

1. Introduction	8
1.1. Scope	8
1.2. Terminology	8
1.3. Intended Audience	9
2. Overview	10
2.1. Features at a Glance	11
2.2. Supported Platforms	11
2.3. Match fields	12
2.4. Actions	13
2.5. Ports	14
2.5.1. Physical Ports	14
2.5.2. Virtual VLAN Ports	15
2.5.3. Virtual Function (VF) Ports	16
2.5.4. OpenFlow Ports	17
2.6. Trivial NIC Mode	17
2.7. Link Aggregation	17
2.8. Load Balance with Group Select	17
2.9. Counters	18
2.10. Checksum Offload	18
2.11. VM QoS	19
3. Functionality and Interfaces	20
3.1. Configuration	20
3.2. ovs-ctl Command Usage	21
3.2.1. Main Startup/Shutdown Command for Agilio OvS Software and Firmware	21
3.2.2. Display Agilio OvS Software Version	23
3.2.3. Display Agilio OvS System Health Status	24
3.2.4. Gather System Troubleshooting Information	26
3.2.5. Fast Path (Wire) Configuration	28
3.2.6. Tap/Mirror Configuration	29
3.3. Configuring Open vSwitch	30
3.3.1. Representative Netdev	30
3.3.2. Example Bridge Setup	35
3.3.3. Layer 2 - Learning Bridge	35
3.3.4. Layer 3 - Routing	36
3.3.5. Tags and Labels	36
3.3.6. Running Packet Analyzer Tools (tcpdump) on Representative Netdevs	37
3.4. Tunneling	37
3.5. Link Aggregation	39
3.6. Load Balance with Group Select	40
3.7. Virtual Function Configuration	41
3.7.1. Creating VF Network Devices	41
3.7.2. Running a DPDK Application on NFP VFs	42
3.7.3. Virtualization Configuration	43
3.7.4. Checksum Offload Configuration	46
3.7.5. Quality of Service Configuration	47
3.8. VirtIO Relay	52

3.8.1. Requirements	52
3.8.2. libvirt and apparmor	52
3.8.3. hugepages	53
3.8.4. VirtIO relay configuration	54
3.8.5. Adding VF ports to VirtIO relay	55
3.8.6. VirtIO relay CPU affinity	56
3.8.7. Running virtual machines	56
3.8.8. Running DPDK on VirtIO inside virtual machines	57
3.8.9. Using VirtIO 1.0	57
3.8.10. VM live migration with libvirt	58
3.8.11. Inspecting traffic with tcpdump	60
3.8.12. Debug counters	60
3.9. Health Monitoring	61
3.10. Trivial NIC Mode	62
3.11. PXE Booting	64
3.12. Connection Tracking (Conntrack)	64
3.12.1. Module Requirements	65
3.12.2. Helper Functionality	66
3.12.3. Conntrack in Use (Recirculation)	66
3.12.4. Conntrack Match Fields	66
3.12.5. Conntrack Action Fields	68
3.12.6. OvS Rule Configuration for Conntrack	69
3.12.7. Viewing Conntrack	70
4. Use Case Examples	71
4.1. Tunnel over Aggregated Links	71
4.1.1. Configuration Example	71
4.2. Load balancing tunnels	72
4.2.1. Configuration Example	73
4.3. Fast dataplane processing combined with standard network services	73
4.3.1. Configuration Example	73
5. Technical Support	76
5.1. Related Documents	76

List of Figures

- 2.1. Agilio OvS Software Architecture and Interfaces 10
- 3.1. VM Rate limiter example setup 48
- 3.2. VM QoS overview 50
- 4.1. Tunnel over Aggregated Links 71

List of Tables

3.1. Agilio OvS Configuration Variables 20

3.2. Port counters available via ethtool 31

3.3. Port counters available via ethtool 34

1. Introduction

1.1 Scope

This manual provides an overview, how to configure, and how to use the Agilio OvS software.

The manual is organized as follows:

- Chapter 2 provides an overview of the Agilio OvS software architecture.
- Chapter 3 describes the configuration, usage, features and restrictions of the Agilio OvS software.
- Chapter 4 provides some configuration use cases for Agilio OvS.
- Chapter 5 provides information on technical support.

1.2 Terminology

Acronym	Description
API	Application Programming Interface
ARP	Address Resolution Protocol
BDF	Bus/Device/Function
BSP	Board Support Package
CDP	Customer Development Platform
CPU	Central Processing Unit
DPDK	Intel© Dataplane Development Kit
ECC	Error-Correcting Code
GSG	Getting Started Guide
GUI	Graphical User Interface
LA	Link Aggregation
LACP	Link Aggregation Control Protocol
LED	Light Emitting Diode
LTS	Long Term Support
MAC	Media Access Control
MPLS	MultiProtocol Label Switching
Mpps	Millions of packets per second
MTU	Maximum Transmission Unit
NFD	Netronome Flow Driver
NFP	Netronome Network Flow Processor

Acronym	Description
NIC	Network Interface Card
NVGRE	Network Virtualization using Generic Routing Encapsulation
OvS	Open vSwitch
OS	Operating System
PCIe	Peripheral Component Interconnect Express
PMD	Poll Mode Driver
QEMU	Generic and open source machine emulator and virtualizer (http://qemu.org)
RX	Receive
SDN	Software Defined Networking
SR-IOV	Single Root I/O Virtualization
TTL	Time To Live
TX	Transmit
VF	SR-IOV Virtual Function
VirtIO	Specification for family of devices used in virtual machine environments (http://docs.oasis-open.org/virtio/virtio/v1.0/virtio-v1.0.html)
VLAN	Virtual Local Area Network
VM	Virtual Machine
VNI	VXLAN Network Identifier
VNIC	Virtual Network Interface Card
VSID	Virtual Subnet Identifier
VXLAN	Virtual eXtensible Local Area Network
VXLAN-GPE	Virtual eXtensible Local Area Network Generic Protocol Extension

1.3 Intended Audience

This User's Guide is written for developers, integrators, and system administrators who are tasked with research, installation, configuration, integration, and testing of Netronome's Agilio OvS and Intelligent Server Adapters.

2. Overview

The Agilio OvS software provides host drivers for network connectivity to virtual machines (VM) and bare metal host applications, as well as firmware to accelerate and offload Open vSwitch (OvS) to the Netronome Intelligent Server Adapters. The Agilio OvS software runs on 64-bit x86 platforms with 64-bit Linux operating systems. In this revision, Ubuntu Server 14.04 LTS and CentOS/RHEL 7.1 distributions have been formally tested and certified by Netronome. The Agilio OvS software and firmware supports OvS version 2.5.1. The OvS kernel data path is offloaded to the Intelligent Server Adapter, improving the performance of the vSwitch as well as relieving x86 CPU cycles from vSwitch processing. Depending on workload and CPU capacity, a user can expect on the order of 5x performance gain over a non-accelerated vSwitch with CPU savings of up to 12 cores. By using OvS for the control plane, Agilio OvS can inherit traditional configuration methods present in OvS today. These include the use of `ovs-ofctl` utilities, as well as remote control from a controller using OVSDb for configuration and OpenFlow for flow table population. These configuration methods are also compatible with higher level software and orchestration such as OpenStack. Network connectivity to the host is supplied via Linux network devices (`netdev`) for terminating applications, as well as the DPDK Poll Mode Driver (PMD) for userspace applications. Both methods of host networking I/O are available to VMs as well as the host OS (a.k.a. Bare Metal Server). There are 60 virtual functions (VFs) that are available to the host or guest operating environments. The diagram in Figure 2.1 depicts a high-level view of the complete system, inclusive of Linux userspace and kernel software, VMs, drivers, and adapter firmware.

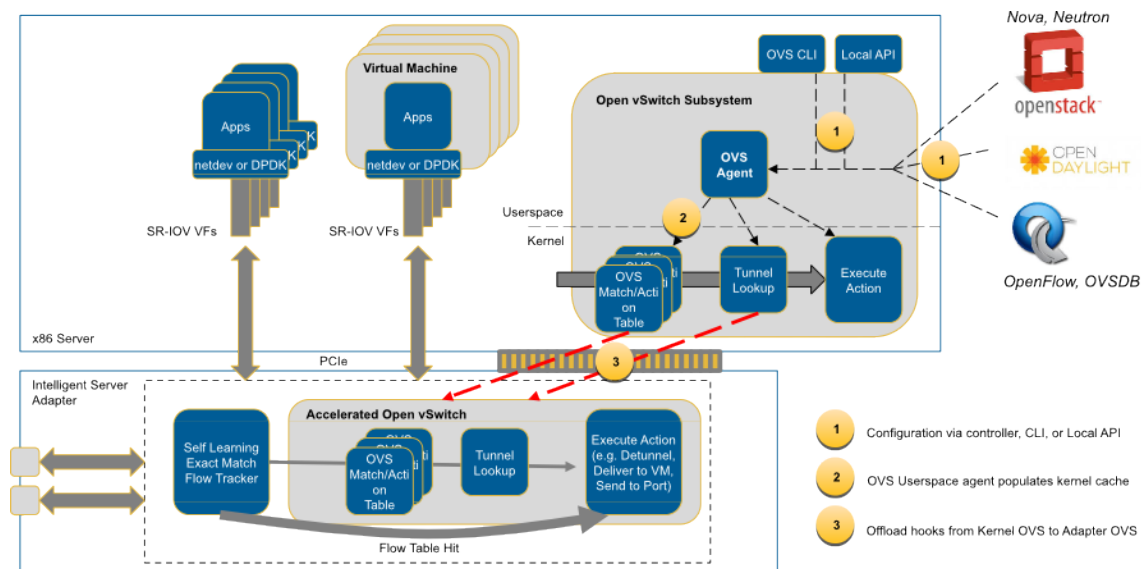


Figure 2.1. Agilio OvS Software Architecture and Interfaces

The Agilio OvS subsystem was designed to provide acceleration and offload for OvS. Therefore, configuration and integration of Agilio OvS with OvS shall be identical to an off-the-shelf, standard installation of OvS. Familiarity with OvS is a prerequisite for using Agilio OvS. In addition to this User's Guide, the OvS documentation can also be used for general OvS configuration.

2.1 Features at a Glance

The following list provides the high level features of the Agilio OvS software:

- Many Open vSwitch datapath traffic matching and action processing functions are accelerated by offload to the Netronome Intelligent Server Adapters. Where acceleration is not supported (e.g. certain control plane traffic), traffic is directed to the Open vSwitch host kernel and/or user mode code for processing in order to retain compatibility with the standard Open vSwitch behavior. This is referred to as fallback processing, and guarantees compatibility with pre-existing server networking configurations.
- Up to 60 SR-IOV VFs are available to users for each PCIe interface on the Intelligent Server Adapters. The VFs can optionally be selectively relayed into virtio connections for VMs that cannot take advantage of native SR-IOV performance.
- Accelerated Link Aggregation for a tcp-balance mode bond, with LACP enabled.
- Accelerated Load Balance functionality via OvS Group Select.
- Accelerated tunnel termination and origination, supporting VXLAN, VXLAN-GPE and NVGRE.
- Configurable 'wire' modes for fastpath forwarding of traffic between specified vPorts on the vSwitch.
- Traffic mirroring at ingress and egress.
- RX and TX checksum offloading, configurable per VF.
- VM QoS: Rate limiting and network bandwidth guarantees

2.2 Supported Platforms

The following Netronome Intelligent Server Adapters are supported with this release of Agilio OvS:

- 2x10G Lithium - Half height, half length.
- 1x40G Beryllium - Half height, half length.
- 2x40G Beryllium - Half height, half length.
- 1x40G Hydrogen - Half height, half length.
- 2x40G Starfighter - Full height, half length.
- 1x100G Starfighter - Full height, half length.
- PCIe expansion through Star Shuttle for each of the Starfighter adapters.
- 4x10G Breakout cables supported on each of the adapters supporting 40G interfaces.
- 10x10G Breakout cables supported on each of the adapters supporting 100G interfaces.

2.3 Match fields

The Agilio OvS software supports many of the match fields supported by Open vSwitch. Refer to the `ovs-ofctl` manual page for details on how to use these match fields. Traffic containing fields that cannot be matched by the accelerated datapath running on the Intelligent Server Adapter is directed to the Open vSwitch host software for fallback processing.

The following match fields are offloaded to the Intelligent Server Adapter:

- Tunnel ID
- Tunnel IPv4 Source
- Tunnel IPv4 Destination
- Tunnel Flags
- Tunnel IPv4 TOS
- Tunnel IPv4 TTL
- Input port
- Ethernet source address
- Ethernet destination address
- Ethernet TCI
- Ethernet Type
- MPLS top label stack entry
- IPv4 source address
- IPv4 destination address
- IPv6 source address
- IPv6 destination address
- IPv6 flow label
- IP protocol
- IP TOS
- IP TTL
- IP Fragmentation
- Transport layer SRC
- Transport layer DST
- Transport layer flags

2.4 Actions

The following basic actions are supported by Agilio OvS and can be offloaded from the Open vSwitch kernel datapath (refer to the OvS documentation for more detailed information on the usage of these actions):

- Output to port

This action is used to output to a specified OpenFlow port. It should be noted that this action need not be the last action in the action list. A copy of the packet will be made if this action is not the last action in the list.

- Fallback to userspace

The action used when a packet is sent to userspace. Examples of such a case could include the situation where the datapath is unable to handle a specific packet/rule combination.

- Set header

This action sets an entire packet header. The following set actions are supported:

- Set tunnel header

Sets tunnel fields `tun_id`, `ipv4_src`, `ipv4_dst`, `tun_flags`, `ipv4_tos`, and `ipv4_ttl`.

- Set ethernet header

Sets ethernet fields `eth_src`, and `eth_dst`.

- Set IPv4 header

Sets IPv4 fields `ipv4_src`, `ipv4_dst`, `ip_tos`, and `ip_ttl`.

- Set IPv6 header

Sets IPv6 fields `ipv6_src`, `ipv6_dst`, `ip_tos`, and `ip_ttl`.

- Set TCP header

Sets transport layer fields `tp_src`, and `tp_dst`.

- Set UDP header

Sets transport layer fields `tp_src`, and `tp_dst`.

- Set MPLS header

Sets MPLS top label stack entry.

- Push VLAN header

Add a VLAN header to a packet.

- Pop VLAN header

Strip outer VLAN header from a packet.

- Push MPLS header

Add an MPLS header to a packet.

- Pop MPLS header

Strip outer MPLS header from a packet.

2.5 Ports

The Agilio OvS software uses the concept of representative netdevs: Linux kernel netdevs are created to represent the various ports available within the Agilio OvS system, so that these can be attached to OvS in a familiar way. Each netdev corresponds to either a physical or virtual port available in the system.

The representative netdevs process only traffic that has not been offloaded to the hardware, for example initial packets of new flows, or flows that have match fields or actions that cannot be offloaded.

The following list illustrates examples of representative netdevs:

- `sdn_p0`: representative netdev corresponding to the first physical port.
- `sdn_v0.1`: representative netdev corresponding to VF.1 on PCIe 0.
- `sdn_v1.59`: representative netdev corresponding to VF.59 on PCIe 1.
- `sdn_p0.1000`: representative netdev corresponding to VLAN ID 1000 on physical port `sdn_p0`.

Refer to Section 3.3.1 for additional details with respect to representative netdevs.

2.5.1 Physical Ports

The Agilio OvS software supports adapters with multiple 10G, 40G, and 100G ports. Each of these physical ports will be represented by a kernel netdev of the format `sdn_pX`, where X is an integer from 0 to 47.

2.5.1.1 Configuration

On the Intelligent Server Adapter, each of the physical network interfaces can be configured to support fiber breakout cables and thereby support more lower rate physical interfaces. The physical interfaces can be configured in such a way by executing the following command:

```
/opt/netronome/bin/nfp-media phy0=CONFIG0 phy1=CONFIG1
```

where `CONFIG1` and `CONFIG2` are 4x10G, 10x10G, 40G, 2x40G or 100G respectively. Note that `phy1=CONFIG1` should be omitted on systems with Intelligent Server Adapters that only have a single physical interface. After the physical port configuration has been modified, the host system needs to be rebooted for the changes to take effect.

The following physical network interfaces are supported on the 4xxx series Intelligent Server Adapters:

- 2x10GE

- 1x40GE
- 2x40GE
- 4x10GE

The following physical network interfaces are supported on the 6xxx series 2x40GE Intelligent Server Adapters:

- 1x40GE and 1x40GE
- 4x10GE and 1x40GE
- 4x10GE and 4x10GE



Note

When a mixed configuration is used (4x10GbE on one port, 1x40GbE on the other port), the highest numbered port must be the 40GbE port. phy0=1x40GE and phy1=4x10GE is not supported.

The following physical network interfaces are supported on the 6xxx series 1x100GE Intelligent Server Adapters:

- 1x100GE
- 2x40GE
- 10x10GE

2.5.1.2 LED Management

Please refer to the Hardware User Guide for pictures and LEDs-to-port mapping as this depends on the specific Intelligent Server Adapter used.

The following LEDs are used in the Intelligent Server Adapters:

- Activity LED (Yellow): Blinks on activity and is updated every 100ms.
- Link LED (Green): Active when MAC detects link from its peer.



Note

When a port is configured to be used with a breakout cable, e.g. a 4x10G port configuration on a 40G interface, the LEDs will only be updated for the first port in the set of ports on the breakout cable. This is applicable to all Intelligent Server Adapters that support breakout cables.

2.5.2 Virtual VLAN Ports

Offloadable 802.1Q VLAN ports can be created on any of the representative physical netdevs. Creating a logical VLAN port will create a new kernel netdev that can be attached to an OvS bridge like any other Agilio OvS

representative netdev. For example, it can be used in configurations with VLAN ports as tunnel endpoints, or as part of a load balancing configuration etc.

To create a VLAN port, one can enter:

```
ip link add link name sdn_p0 name sdn_p0.1000 type vlan id 1000
```

There will now be a kernel netdev with the name `sdn_p0.1000` available on the host. Any traffic sent to this VLAN netdev, via OvS or directly via the kernel netdev, will have a VLAN tag with a VID of 1000 inserted when the packet is sent out of physical port 0. Similarly, when traffic is received with a VLAN tag and a VID of 1000, the VLAN tag will be stripped before the packet enters the OvS pipeline.

Statistics, link state and MTU are available in the same way as any other 802.1Q VLAN interface. It must also be noted that when such an interface is created, the initial netdev state is 'down'. VLAN sub-interfaces have their own MTU setting independent of the parent netdev, however this MTU must be less than or equal to the parent netdev MTU.

This Agilio OvS feature aligns with the generic 802.1Q VLAN features available in the Linux kernel. All utilities supported by the kernel to create such netdevs are supported with Agilio OvS VLAN ports, e.g. the `ip` utility or the Debian based `interfaces` file.

A maximum of 4096 such interfaces can be offloaded. One may configure any range of VLAN tag VIDs on any physical interface up to a global maximum of 4096 offloaded VLAN ports. An appropriate error message will be displayed once the limit has been reached.

Note the following restrictions:

- Native Linux 802.1AD VLAN tagging is not supported.
- Creating a VLAN interface on a VF port is not supported¹.
- Nested VLAN tagging, i.e. creating a VLAN interface on an existing VLAN interface will not be offloaded to the NFP.

2.5.3 Virtual Function (VF) Ports

60 VFs per physical function are available to the user. The number of physical functions (PCIe units) is dependent on the platform. These VFs can be bound to DPDK applications (i.e. Poll Mode Drivers running in userspace), or to kernel netdevs. The DPDK applications and netdevs can be deployed in VMs or can be instantiated on the host itself. The VFs are represented by host kernel representative netdevs named `sdn_vX.Y`, where X is the NFP PCIe unit number (0-3) and Y is the VF number on that PCIe unit. It is important to understand that these kernel representative netdevs are NOT the VF traffic endpoints themselves. They merely provide a representative netdev that can be connected to the Open vSwitch software running on the host to facilitate non-accelerated processing of traffic. They can furthermore be used to obtain traffic statistics of each of the fallback and fastpathed traffic.



Note

Only 60 VFs per PCIe unit are available to the user. VFs 60 to 63 are reserved.

¹VLAN tags can be configured with OvS for VF ports.

2.5.4 OpenFlow Ports

Any representative netdev, either a physical port netdev, a VLAN netdev or a VF netdev, may be added to an Open vSwitch bridge². The OpenFlow port configuration requested by the user when configuring Open vSwitch will be honored with respect to all Open vSwitch commands. The OpenFlow port number may thus be used as an user configurable port number.

2.6 Trivial NIC Mode

The Trivial NIC mode refers to the mode of the Agilio OvS software when Open vSwitch itself is not running, but the NFP firmware is loaded, configured and able to pass traffic to the host. In this mode, all the representative netdevs are active in the kernel, but only the physical port representative netdevs are operational. Although these representative netdevs can be used for basic networking functions, no offloading or additional features are available in this mode.

The Trivial NIC feature is described in more detail in Section 3.10.

2.7 Link Aggregation

Link Aggregation (also known as bonding, port trunking or link bundling) allows using multiple network connections in parallel to:

- Provide redundancy (active / backup links)
- Provide increased throughput (active / active links)

In the case of increased throughput, for each packet sent to the bond, one output port is selected from the set of ports included in the bond.

The accelerated implementation of Link Aggregation supports 32 independent bonds, each containing up to 128 ports, using the tcp-balance bond mode. Each port in a bond is selected with equal likelihood³. A bond can contain a combination of physical ports and VFs. LACP (Link Aggregation Control Protocol) is supported, in both active / active and active / passive modes.

2.8 Load Balance with Group Select

The Agilio OvS software supports accelerated load balancing using group select in OvS. Group Select selects and then executes one bucket in a group. The selection algorithm results in equal load sharing between buckets.

²Any 3rd party NIC netdev or tap netdev may also be added to an OvS bridge. OpenFlow flows containing match fields or actions related to such ports will simply not be offloaded to the NFP, but will still function correctly.

³Traffic that does not contain layer 3 or 4 information will not be spread over multiple bonds equally. Examples of such traffic include MPLS and ARP.

Agilio OvS 2.2 supports the ability to combine the group select-based load balancer with other groups such that the offloaded group select is the last group in the chain. Load balancing is achieved by applying a Jenkins hash to the source and destination address information layer 3 and layer 4. Traffic that does not contain layer 3 or 4 will not be spread over multiple buckets. Examples of such traffic include MPLS and ARP.

2.9 Counters

Agilio OvS software supports all counters made available by Open vSwitch. Refer to the appropriate Open vSwitch documentation for instructions w.r.t. how to extract counters from the system. The software ensures that the Open vSwitch counters maintained on the host accurately reflect all traffic, whether traffic was handled by the accelerated datapath on the NFP or whether it was handed to the Open vSwitch software on the host for fallback processing.

2.10 Checksum Offload

Agilio OvS software supports L3 and L4 checksum offload on all VF ports. Both the Linux kernel network driver and the Poll Mode (DPDK) drivers are supported. Furthermore RX and TX checksum offloading is supported for traffic flowing from the network to a VF, from a VF to the network and from a VF to a VF. Layer 3 protocols IPv4 and IPv6 are supported. Layer 4 protocols TCP and UDP are supported.

The IP header checksum is checked as per RFC 791, where the packet meets the following criteria:

- If present, the DSA-tag header must be four or 8 octets long.
- If present, the VLAN header must be four octets long.
- Encapsulation must be RFC 894 Ethernet Type Encoding.
- IPv4 packet - (ethertype == 0x800) AND (ver == 4) AND (hdr_len >= 5).
- IP header is of a valid length.
- The packet can have at most 2 levels of VLAN tags.
- The packet can have at most 8 levels of MPLS tags.
- If both MPLS and VLAN are present, the packet can have at most 8 MPLS and 2 VLAN tags.
- The layer 3 protocols that are supported are IPv4 and IPv6.

The TCP/UDP header checksum is calculated as per RFC 793, where the packet meets the following criteria:

- IPv4 or IPv6.
- No IP fragmentation.
- No ESP header.
- No Authorization header.
- No Mobility header.
- TCP or UDP packet.

2.11 VM QoS

Agilio OvS software supports setting the maximum bit rate a VM can send or receive. This is enforced on the VF(s) that are bound to the VM. It also supports the prioritisation of specific VF traffic over physical network ports.

3. Functionality and Interfaces

3.1 Configuration

The Agilio OvS software can be configured by using the `/etc/netronome.conf` file. This file takes the form of `variable=value` entries, one per line. Lines starting with the `"#"` character are treated as comments and are ignored. The table below lists the variables that are recognized in this file. For each variable, if applicable, the default value and the permissible values or the acceptable value range are specified.

Table 3.1. Agilio OvS Configuration Variables

Name / Description	Unit	Valid values	Default
SDN_VERBOSE: Determines whether to perform verbose logging. This only affects the logging performed when the software is started and stopped.	Text	"y" or "n"	"n"
SDN_VIRTIORELAY_2M_HUGEPAGES: Selects the number of 2Mb hugepages which will be configured automatically when starting Agilio OvS. The startup procedure will only increase the number of hugepages if the currently configured number of hugepages is less than the specified number. Setting the number of hugepages to 0, prevents the startup script from making any modification to the hugepage configuration on the system.	Integer	0 - available memory on host system	768
SDN_VIRTIORELAY_ENABLE: Enables or disables the virtiorelayd process from starting up along with Agilio OvS.	Text	"y" or "n"	"n"
SDN_BIND_VF_DRIVER: Select the default driver VFs are bound to when they are created. Only applicable when the VFs are recreated, i.e. when trivial NIC mode is bypassed.	Text	"nfp_uio" / "nfp_netvf" / "vfio-pci" / "pci-stub" / "none"	"nfp_uio"
SDN_FORCE_NFP_RESET: Option to reset the NFP when starting or stopping Agilio OvS. Enabling this option bypasses the Trivial NIC mode. Refer to Section 3.10 for more details about the Trivial NIC feature.	Text	"y" or "n"	"n"
SDN_VIRTIORELAY_HUGEDIR:	Text	Refer to virtiorelayd usage.	"/mnt/aovs-huge-2M"

Name / Description	Unit	Valid values	Default
Provides the user the ability to modify the mount path to the hugepages for virtiorelayd.			
SDN_VIRTIORELAY_VHOSTDIR: Provides the user the ability to modify the vhost-user unix socket directory path for virtiorelayd.	Text	Refer to virtiorelayd usage.	"/tmp"
SDN_VIRTIORELAY_PARAM: Provides the user the ability to modify the commandline configuration options to the virtiorelayd process. This variable is set to the default value during installation unless there is already an existing value defined.	Text	Refer to virtiorelayd usage.	"--cpus=1,2 --vhost-username=libvirt-qemu --vhost-groupname=kvm --huge-dir=/mnt/aovs-huge-2M --vhost-path=/tmp"
SDN_CONNTRACK: Option to enable or disable the experimental Agilio OvS Conntrack mode. Although this option is configurable by the user, it should only be enabled when the experimental Agilio OvS Conntrack version was installed.	Text	"y" or "n"	Set according to installation option.

None of the variables are mandatory. If they are not specified the system will use the respective default value.

The following example illustrates a valid configuration file.

```
SDN_VIRTIORELAY_2M_HUGEPAGES=128
SDN_BIND_VF_DRIVER=nfp_uio
```

3.2 ovs-ctl Command Usage

The `ovs-ctl` command is used to control the Agilio OvS software. The Agilio OvS software is directly hooked into Open vSwitch, therefore the `ovs-ctl` command used to manage the Open vSwitch software will also manage the Agilio OvS software. Refer to the manual page of the `ovs-ctl` for all the details.

The following sections provide more detail on the specific Netronome extensions to `ovs-ctl`, along with examples of the usage of the particular command and expected output. The examples below are not indicative of the complete scope of all inputs and outputs, but merely illustrate the normal usage of the command.

3.2.1 Main Startup/Shutdown Command for Agilio OvS Software and Firmware

- **Command:**

```
/opt/netronome/bin/ovs-ctl [start|stop|restart]
```

- **Syntax Description:**

Parameter	Description
start	Start the Agilio OvS host software including Open vSwitch software. If the Trivial NIC (see Section 3.10) was not active, then all Agilio OvS kernel modules as well as the Network Flow Processor firmware will also be loaded.
stop	Stop the Agilio OvS host software. Kernel modules no longer required will be unloaded. This will also stop Open vSwitch software. If the Trivial NIC mode is bypassed (see Section 3.10), then the Netronome Flow Processor firmware will also be unloaded at this point.
restart	Similar to <code>ovs-ctl stop</code> followed by <code>ovs-ctl start</code> , except that the OpenFlow flows are preserved if possible. If the Trivial NIC mode is bypassed (see Section 3.10), the Netronome Flow Processor firmware will be reconfigured and reloaded along with all the Agilio OvS kernel modules.

- **Command Mode:**

Privileged Execution

- **Usage Guidelines:**

It is recommended that no other commands be executed while waiting for the Agilio OvS software to start or stop.

After the system has been started, the following processes should be running for normal system operation:

- `ovs-vswitchd`
- `ovsdb-server`
- `virtiorelayd` (if enabled in `/etc/netronome.conf`)

The required processes and their current status can be viewed with the `ovs-ctl status` command. This command is described in more detail below.

All options provided by the standard OvS `ovs-ctl` are available, for example to prevent the following log message

```
* system ID not configured, please use --system-id
```

one can specify `ovs-ctl start --system-id=random`. Refer to the man page for the standard OvS options and parameters applicable to the `ovs-ctl`.

Additionally, the `'--nfp-reset'` option can be used to bypass the Trivial NIC mode. Refer to Section 3.10 for more information about this feature.

The `ovs-ctl stop` command will not execute when the Trivial NIC mode is bypassed and there are any DPDK applications or VMs connected to any VFs. A warning with the respective process which is still running will be reported to the user. The user may attempt to force this operation by specifying the '`--force`' switch as follow:

```
/opt/netronome/bin/ovs-ctl --nfp-reset --force stop
```

- **Examples:**

Examples of typical commands and the expected output are supplied below. There may be some minor variation of the output depending on the platform the software is used on and options selected by the user.

Command:

```
/opt/netronome/bin/ovs-ctl start
```

Output:

```
* Starting Agilio OvS 2.2
* Starting ovssdb-server
* system ID not configured, please use --system-id
* Configuring Open vSwitch system IDs
* Starting ovs-vswitchd
* Enabling remote OVSSDB managers
* Agilio OvS 2.2 started successfully
* Running firmware image /lib/firmware/nfp_firmware2.2_4001.nffw
```

Command:

```
/opt/netronome/bin/ovs-ctl stop
```

Output:

```
* Stopping Agilio OvS 2.2
* Exiting ovs-vswitchd (70074)
* Exiting ovssdb-server (70003)
* Agilio OvS 2.2 stopped successfully
```

3.2.2 Display Agilio OvS Software Version

- **Command:**

```
/opt/netronome/bin/ovs-ctl version
```

- **Syntax Description:**

Parameter	Description
version	Display the current version of the Agilio OvS software installed on the system.

- **Command Mode:**

User Execution

- **Usage Guidelines:**

This version text should accompany any queries logged with Netronome support.

- **Example:**

Examples of typical commands and the expected output are supplied below.

Command:

```
/opt/netronome/bin/ovs-ctl version
```

Output:

```
ovsdb-server (Open vSwitch) 2.5.1
Compiled Jul 18 2016 21:12:42
ovs-vswitchd (Open vSwitch) 2.5.1
Compiled Jul 18 2016 21:12:52
Netronome Agilio OvS version 2.2-r4518:xxxxxxxxxxxxx
NFP Firmware Loaded: /lib/firmware/nfp_firmware2.2_40011.nffw
```

3.2.3 Display Agilio OvS System Health Status

- **Command:**

```
/opt/netronome/bin/ovs-ctl status
```

- **Syntax Description:**

Parameter	Description
status	Display the current health status of the Agilio OvS software and hardware.

- **Command Mode:**

Privileged Execution

- **Usage Guidelines:**

The health monitor utility will display the information for the NFP installed in the system and the current software status.

This command can be invoked routinely or whenever there is a suspicion that the system may have failed. It is recommended that this command should not be invoked while starting or stopping the Agilio OvS software.

Invoking this command will log any failures detected to syslog. Refer to Section 3.9 for the exact messages that will be logged to syslog for each type of failure. These syslog messages also describe the recommended action to take to resolve the failure.

If the syslog messages contain "restart", e.g. "restart.process", "restart.drivers", etc., restart the system by entering the command `/opt/netronome/bin/ovs-ctl --nfp-reset restart`. If the problem occurs again, power cycle the system. If the problem persists, enter the command `/opt/netronome/bin/ovs-ctl status troubleshoot` and contact Netronome support. Similarly, if the syslog messages contain "replace", e.g. "replace.nfe", the hardware may need to be replaced by contacting Netronome support.

The return code of this command indicates the number of failures detected. Warnings could also be logged, but although these warnings must be brought to the attention of the user, the system is not in a failed state and the return code of the command is zero.

- **Examples:**

Examples of typical commands and the expected output are supplied below.

Command: `/opt/netronome/bin/ovs-ctl status` on a system currently running Agilio OvS software.

Output:

```

    Detected Configurator version: 20160211132710

SDN Health Report:
-----
Userspace Process: ovsdb-server          ... [PASS]
Userspace Process: ovs-vswitchd         ... [PASS]
Userspace Process: virtiorelayd         ... [PASS]
Kernel Module: nfp                      ... [PASS]
Kernel Module: nfp_cmsg                 ... [PASS]
Kernel Module: nfp_fallback             ... [PASS]
Kernel Module: nfp_offloads             ... [PASS]
Kernel Module: openvswitch              ... [PASS]
NFP: Firmware Loaded                   ... [PASS]
NFP: Flow Processing Cores Responsive   ... [PASS]
NFP: Control Message Channel Responsive ... [PASS]
NFP: Fallback Traffic Divert Disabled   ... [PASS]
NFP: Ingress NBI Backed Up             ... [PASS]
NFP: Card Detected                     ... [PASS]
NFP: ECC Errors                        ... [PASS]
NFP: Parity Errors                     ... [PASS]
NFP: Configurator Version              ... [PASS]
Host: ERR47 Workaround                  ... [PASS]

Overall System State      [PASS]

virtiorelayd is running with pid 7957
ovsdb-server is running with pid 7816
ovs-vswitchd is running with pid 7826
    
```

Command: `/opt/netronome/bin/ovs-ctl status` on a system with Agilio OvS software currently stopped with Trivial NIC mode activated.

Output:

```

    Detected Configurator version: 20160211132710

SDN Health Report:
    
```

```

-----
Userspace Process: ovsdb-server          ... [FAIL]
Userspace Process: ovs-vswitchd         ... [FAIL]
Userspace Process: virtiorelayd         ... [SKIPPED]
Kernel Module: nfp                      ... [PASS]
Kernel Module: nfp_cmsg                 ... [PASS]
Kernel Module: nfp_fallback             ... [PASS]
Kernel Module: nfp_offloads             ... [FAIL]
Kernel Module: openvswitch              ... [PASS]
NFP: Firmware Loaded                    ... [PASS]
NFP: Flow Processing Cores Responsive    ... [PASS]
NFP: Control Message Channel Responsive ... [PASS]
NFP: Fallback Traffic Divert Disabled    ... [WARN]
NFP: Ingress NBI Backed Up              ... [PASS]
NFP: Card Detected                      ... [PASS]
NFP: ECC Errors                         ... [PASS]
NFP: Parity Errors                      ... [PASS]
NFP: Configurator Version               ... [PASS]
Host: ERR47 Workaround                  ... [PASS]

Overall System State      [FAIL]
    
```

Command: /opt/netronome/bin/ovs-ctl status on a system with Agilio OvS software currently stopped with Trivial NIC mode deactivated.

Output:

```

    Detected Configurator version: 20160211132710

SDN Health Report:
-----
Userspace Process: ovsdb-server          ... [FAIL]
Userspace Process: ovs-vswitchd         ... [FAIL]
Userspace Process: virtiorelayd         ... [SKIPPED]
Kernel Module: nfp                      ... [PASS]
Kernel Module: nfp_cmsg                 ... [FAIL]
Kernel Module: nfp_fallback             ... [FAIL]
Kernel Module: nfp_offloads             ... [FAIL]
Kernel Module: openvswitch              ... [PASS]
NFP: Firmware Loaded                    ... [FAIL]
NFP: Flow Processing Cores Responsive    ... [WARN]
NFP: Control Message Channel Responsive ... [FAIL]
NFP: Fallback Traffic Divert Disabled    ... [SKIPPED]
    Control channel unreliable
NFP: Ingress NBI Backed Up              ... [PASS]
NFP: Card Detected                      ... [PASS]
NFP: ECC Errors                         ... [PASS]
NFP: Parity Errors                      ... [PASS]
NFP: Configurator Version               ... [PASS]
Host: ERR47 Workaround                  ... [PASS]

Overall System State      [FAIL]
    
```

3.2.4 Gather System Troubleshooting Information

- **Command:**

```
/opt/netronome/bin/ovs-ctl status troubleshoot
```

- **Syntax Description:**

Parameter	Description
troubleshoot	Compiles and archives various troubleshooting information about the system which should be supplied to Netronome support when system diagnostics are required.

- **Command Mode:**

Privileged Execution

- **Usage Guidelines:**

The troubleshoot utility will gather information of the system components and archive this in a single file stored in the `/opt/netronome/sysinfo/archives/` directory.

The troubleshoot utility should be invoked as soon as possible after a failure has been observed to obtain the most usable data for diagnostics later on.

While this command is executing, each of the component groups is displayed while information is extracted from it. The only important information displayed to the user is the name and location of the archive file where the collected data will be stored. The messages displayed on the screen merely indicate progress. This command may take a long time to complete.

If an error occurs while trying to extract information from a system component, the "done" status printed at end of each line will be replaced with a "failed" message. If this does occur, some of the debugging information contained within the archive may not be complete. However it may also be a consequence of the actual system failure and the remainder of the information will still be useful to Netronome Support.



Note

After the execution of the troubleshoot utility the host system will need to be rebooted due to the destructive nature of gathering the internal debugging data.

- **Example:**

Examples of typical commands and the expected output are supplied below.

Command:

```
/opt/netronome/bin/ovs-ctl status troubleshoot
```

Output:

```
Netronome Systems - System Troubleshoot Information

Saving last logs... done
Saving kernel ring buffer (dmesg)... done
Saving device info (lspci)... done
```

```

Saving NFE HWINFO and SDN config files...done
Saving versions... done
Saving status of SDN components... done
Saving information about host system... done
Saving `top' util output... done
Saving current process information... done
Saving ARP/Routing tables... done
Saving ME Information (this may take a while)... done
Saving current OVS configuration information... done
Saving flow entries... done
Saving flow cache diagnostics (this may take up to 10 minutes)... done
Saving NFD info...done
Saving control channel info...done
Saving GRO dump...done
Saving Port Information... done
Saving NFP symbols (this may take a while)... done
Saving NFP Debug Information (this may take up to 5 minutes)... done
Saving symbol list dump...done
Saving worker states...done
Saving Error Correct Code info... done
Dumping QINFO ... done
Dumping GPRs and Local Memory ... done

Archiving system information... done

Please review the contents of /opt/netronome/sysinfo and e-mail
/opt/netronome/sysinfo/archives/sysinfo-2016-04-06.1053.tgz
to support@netronome.com along with a description of the problem you have
encountered. If possible, also provide steps to reproduce the problem.

Please reboot the system. Due to the troubleshoot information gathered,
the system will need to be rebooted to operate normally again.

```

3.2.5 Fast Path (Wire) Configuration

- **Commands:**

```
/opt/netronome/bin/ovs-ctl configure wire from IN_PORT to OUT_PORT
```

```
/opt/netronome/bin/ovs-ctl configure wire clear IN_PORT
```

- **Syntax Description:**

Parameter	Description
IN_PORT, OUT_PORT	The PORT parameters refer to representative netdevs as described in Section 3.3.1. However, representative VLAN netdevs are not supported with this feature.

- **Command Mode:**

Privileged Execution

- **Usage Guidelines:**

There are no specific usage guidelines for this command.

- **Example:**

Examples of typical commands are supplied below.

To bidirectionally connect physical ports sdn_p0 and sdn_p1:

```
ovs-ctl configure wire from sdn_p0 to sdn_p1
ovs-ctl configure wire from sdn_p1 to sdn_p0
```

To bidirectionally connect physical port sdn_p1 to the DPDK application / netdev attached to VF 1 on NFP PCIe unit 0:

```
ovs-ctl configure wire from sdn_p1 to sdn_v0.1
ovs-ctl configure wire from sdn_v0.1 to sdn_p1
```

To remove a wire (i.e. resume normal operations):

```
ovs-ctl configure wire clear sdn_p1
ovs-ctl configure wire clear sdn_v0.1
```

3.2.6 Tap/Mirror Configuration

- **Commands:**

```
/opt/netronome/bin/ovs-ctl configure mirror rx from IN_PORT to OUT_PORT
```

```
/opt/netronome/bin/ovs-ctl configure mirror tx from IN_PORT to OUT_PORT
```

```
/opt/netronome/bin/ovs-ctl configure mirror rx clear IN_PORT
```

```
/opt/netronome/bin/ovs-ctl configure mirror tx clear IN_PORT
```

- **Syntax Description:**

Parameter	Description
IN_PORT, OUT_PORT	The PORT parameters refer to representative netdevs as described in Section 3.3.1. However, representative VLAN netdevs are not supported with this feature.

- **Command Mode:**

Privileged Execution

- **Usage Guidelines:**

Configure a tap or mirror on a physical port or a virtual function or another physical port. Packet will be copied on reception (mirror rx) or transmission (mirror tx) and the copy will be sent to the mirror destination. Expect a performance impact under high packet rates.

Note that TX mirrors will show only traffic for flows that are offloaded to the NIC. Run tcpdump on the appropriate representative netdev to inspect the fallback traffic

- **Example:**

Examples of typical commands are supplied below.

To mirror traffic received on sdn_p0 to sdn_p1:

```
ovs-ctl configure mirror rx from sdn_p0 to sdn_p1
```

To mirror received on physical port sdn_p0 to the DPDK application / netdev attached to VF 1 on NFP PCIe unit 0:

```
ovs-ctl configure mirror rx from sdn_p0 to sdn_v0.1
```

To mirror traffic transmitted on sdn_p0 to the DPDK application / netdev attached to VF 1 on NFP PCIe unit 0:

```
ovs-ctl configure mirror tx from sdn_p0 to sdn_v0.1
```

To remove a mirror (i.e. resume normal operations):

```
ovs-ctl configure mirror rx clear sdn_p0
ovs-ctl configure mirror tx clear sdn_p0
ovs-ctl configure mirror rx clear sdn_v0.1
```

3.3 Configuring Open vSwitch

3.3.1 Representative Netdev

The Agilio OvS software uses the concept of representative netdevs. Kernel netdevs are automatically created by the `nfp_fallback.ko` driver, and represent the physical and virtual function interfaces. When these netdevs are attached to an Open vSwitch bridge, traffic on them will automatically trigger the offloading of flows to the NFP. Examples of representative netdevs are:

- `sdn_p0`: representative netdev corresponding to the first physical port.
- `sdn_v0.1`: representative netdev corresponding to VF.1 on NFP PCIe unit 0.
- `sdn_v1.59`: representative netdev corresponding to VF.59 on NFP PCIe unit 1.
- `sdn_p0.1000`: representative netdev corresponding to VLAN ID 1000 on physical port `sdn_p0`.

Offloading flows to the NFP eventually results in fewer fallback packets being emitted from the representative netdevs. Running `ethtool -S` on these representative netdevs will give two distinct sets of packet and byte counters: one for the underlying device (physical port or VF), and one set for the actual fallback traffic received on the host representative netdev. Refer to the table below for a complete reference of all the port statistics supported by the software. The fallback traffic counters will ordinarily be far less than the underlying device counters when flows have been successfully offloaded to the NFP. The counters displayed by `ifconfig` (also seen in OvS port stats) represent the underlying device, not the fallback traffic.

Note that representative VLAN netdevs does not have separate counters for fallback and port counters, instead the counters displayed for these netdevs is always the NFP logical port counters.

Virtual Function devices can have two netdevs associated with them. The representative netdev representing the VF, and the netdev running on the actual VF (Refer to Section 3.7.1). The important distinction is that the representative netdev should only emit packets from flows that have not been offloaded to the NFP. The netdev running directly on the VF will receive and transmit all packets: fallback packets and offloaded packets will both be routed correctly. The actual VF can of course also be assigned to a VM, and/or be used by a DPDK poll mode driver instead or a kernel netdev driver, and the VF is where user applications will receive and transmit traffic. The VF netdev is optimized for bulk traffic performance, whereas the representative netdevs are considered to be a low bandwidth traffic path.

Consequently, to allow for proper offloading, the representative device is added to an OvS bridge and the netdev (or DPDK PMD, or VM etc.) associated with the VF is where user application traffic is directed.

The netdev running directly on the VF is named by the Linux distribution. The Agilio OvS software does not override this and the standard rules implemented by udev, for example, are applied.

The following table provides detailed descriptions of all the port statistics supported by physical ports on the NFP. These statistics are available via the ethtool utility. VF netdevs will provide a subset of these statistics.

Table 3.2. Port counters available via ethtool

Counter	Description
fallback rx bytes	Bytes received by representative netdev.
fallback rx packets	Packets received by representative netdev.
fallback tx bytes	Bytes transmitted by representative netdev.
fallback tx packets	Packets transmitted by representative netdev.
fallback tx errors	Transmit errors detected by representative netdev.
fallback tx copies	Packets copied on transmission by representative netdev.
device rx bytes	See description of <code>RxPifInOctets</code> below.
device rx packets	See description of <code>RxFramesReceivedOK</code> below.
device tx bytes	See description of <code>TxPifOutOctets</code> below.
device tx packets	See description of <code>TxFramesTransmittedOK</code> below.
device rx frame errors	See description of <code>RxFrameCheckSequenceErrors</code> below.
device rx errors	See description of <code>RxPifInErrors</code> below.
device rx multicast	See description of <code>RxPifInMultiCastPkts</code> below.
device rx dropped	See description of <code>RxPStatsDropEvents</code> below.
device tx errors	See description of <code>TxPifOutErrors</code> below.
rx rate packets/s	Packet rate received determined from the update rate of the 'device rx packets' counter.
rx rate bytes/s	Byte rate received determined from the update rate of the 'device rx bytes' counter.

Counter	Description
tx_rate_packets/s	Packet rate transmitted determined from the update rate of the 'device tx packets' counter.
tx_rate_bytes/s	Byte rate transmitted determined from the update rate of the 'device tx bytes' counter.
RxAlignmentErrors	Frames received with an alignment error.
RxCBFCPauseFramesReceived0	CBFC (Class Based Flow Control) pause frames received for class 0.
RxCBFCPauseFramesReceived1	CBFC (Class Based Flow Control) pause frames received for class 1.
RxCBFCPauseFramesReceived2	CBFC (Class Based Flow Control) pause frames received for class 2.
RxCBFCPauseFramesReceived3	CBFC (Class Based Flow Control) pause frames received for class 3.
RxCBFCPauseFramesReceived4	CBFC (Class Based Flow Control) pause frames received for class 4.
RxCBFCPauseFramesReceived5	CBFC (Class Based Flow Control) pause frames received for class 5.
RxCBFCPauseFramesReceived6	CBFC (Class Based Flow Control) pause frames received for class 6.
RxCBFCPauseFramesReceived7	CBFC (Class Based Flow Control) pause frames received for class 7.
RxFrameCheckSequenceErrors	CRC-32 Error is detected but the frame is otherwise of correct length.
RxFrameTooLongErrors	Frame received exceeded the maximum length.
RxFramesReceivedOK	Frame received without error (including pause frames).
RxInRangeLengthErrors	A count of frames with a length / type field value between 46 (VLAN: 42) and less than 0x0600, that does not match the number of payload data octets received. Should count also if length / type field is less than 46 (VLAN: 42) and the frame is longer than 64 bytes.
RxPIfInBroadCastPkts	Incremented with each valid frame received on the Receive FIFO interface and all bits of the destination address were set to '1'.
RxPIfInErrors	Number of frames received with error: FIFO Overflow Error, CRC Error, Frame Too Long Error, Alignment Error, The dedicated Error Code (0xfe, not a code error) was received.
RxPIfInMultiCastPkts	Incremented with each valid frame received on the Receive FIFO interface and bit 0 of the destination address was '1' but not the broadcast address (all bits set to '1'). Pause frames are not counted.
RxPIfInUniCastPkts	Incremented with each valid frame received on the Receive FIFO interface and bit 0 of the destination address was '0'.
RxPStatsDropEvents	Counts the number of dropped packets due to internal errors of the MAC Client. Occurs when a Receive FIFO overflow condition persists.
RxPStatsFragments	Total number of packets that were less than 64 octets long with a wrong CRC. Note: Fragments are not delivered to the FIFO interface.
RxPStatsJabbers	Total number of packets longer than the valid maximum length programmed in register FRM_LENGTH (excluding framing bits, but including FCS octets), and with a bad Frame Check Sequence.

Counter	Description
RxPStatsOversizePkts	Total number of packets longer than the valid maximum length programmed in register FRM_LENGTH (excluding framing bits, but including FCS octets), and with a good Frame Check Sequence.
RxPStatsPkts	Total number of packets received. Good and bad packets.
RxPStatsPkts1024to1518octets	Frames (good and bad) with 1024 to 1518 octets.
RxPStatsPkts128to255octets	Frames (good and bad) with 128 to 255 octets.
RxPStatsPkts1519toMaxoctets	Proprietary RMON extension counter that counts the number of frames with 1519 bytes to the maximum length programmed in register FRM_LENGTH.
RxPStatsPkts256to511octets	Frames (good and bad) with 256 to 511 octets.
RxPStatsPkts512to1023octets	Frames (good and bad) with 512 to 1023 octets.
RxPStatsPkts64octets	Incremented when a packet of 64 octets length is received (good and bad frames are counted).
RxPStatsPkts65to127octets	Frames (good and bad) with 65 to 127 octets.
RxPStatsUndersizePkts	Total number of packets that were less than 64 octets long with a good CRC. Note: Undersized packets are not delivered to the FIFO interface.
RxPauseMacCtlFramesReceived	Valid pause frame received.
RxVlanReceivedOK	VLAN frame received without error.
TxCBFCPauseFramesTransmitted0	CBFC (Class Based Flow Control) pause frames transmitted for class 0.
TxCBFCPauseFramesTransmitted1	CBFC (Class Based Flow Control) pause frames transmitted for class 1.
TxCBFCPauseFramesTransmitted2	CBFC (Class Based Flow Control) pause frames transmitted for class 2.
TxCBFCPauseFramesTransmitted3	CBFC (Class Based Flow Control) pause frames transmitted for class 3.
TxCBFCPauseFramesTransmitted4	CBFC (Class Based Flow Control) pause frames transmitted for class 4.
TxCBFCPauseFramesTransmitted5	CBFC (Class Based Flow Control) pause frames transmitted for class 5.
TxCBFCPauseFramesTransmitted6	CBFC (Class Based Flow Control) pause frames transmitted for class 6.
TxCBFCPauseFramesTransmitted7	CBFC (Class Based Flow Control) pause frames transmitted for class 7.
TxFramesTransmittedOK	Frame transmitted without error (including pause frames).
TxPIfOutBroadCastPkts	Incremented with each frame written to the FIFO interface and all bits of the destination address set to '1'.

Counter	Description
<code>TxPIfOutErrors</code>	Number of frames transmitted with error: FIFO Overflow Error, FIFO Underflow Error, or User application defined error (<code>ff_tx_err</code> asserted together with <code>ff_tx_eop</code>).
<code>TxPIfOutMultiCastPkts</code>	Incremented with each frame written to the FIFO interface and bit 0 of the destination address set to '1' but not the broadcast address (all bits set to '1').
<code>TxPIfOutUniCastPkts</code>	Incremented with each frame written to the FIFO interface and bit 0 of the destination address set to '0'.
<code>TxPStatsPkts1024to1518octets</code>	Frames (good and bad) with 1024 to 1518 octets.
<code>TxPStatsPkts128to255octets</code>	Frames (good and bad) with 128 to 255 octets.
<code>TxPStatsPkts1518toMAXoctets</code>	Frames (good and bad) with 1518 to MAX octets.
<code>TxPStatsPkts256to511octets</code>	Frames (good and bad) with 256 to 511 octets.
<code>TxPStatsPkts512to1023octets</code>	Frames (good and bad) with 512 to 1023 octets.
<code>TxPStatsPkts64octets</code>	Incremented when a packet of 64 octets length is received (good and bad frames are counted).
<code>TxPStatsPkts65to127octets</code>	Frames (good and bad) with 65 to 127 octets.
<code>TxPauseMacCtlFramesTransmitted</code>	Valid pause frame transmitted.
<code>TxVlanTransmittedOK</code>	VLAN frame transmitted without error.
<code>RxPIfInOctets</code>	All octets received except preamble (i.e. Header, Payload, Padding and FCS) for all frames and pause frames received.
<code>TxPIfOutOctets</code>	All octets transmitted except preamble (i.e. Header, Payload, Padding and FCS) for all valid frames and valid pause frames transmitted.

Additionally, the following counters will be available for representative VF netdevs only:

Table 3.3. Port counters available via ethtool

Counter	Description
<code>rx rate limited packets</code>	Packets dropped by ingress rate limiter from the vSwitch point of view.
<code>rx rate limited bytes</code>	Bytes dropped by ingress rate limiter from the vSwitch point of view.
<code>tx rate limited packets</code>	Packets dropped by egress rate limiter from the vSwitch point of view.
<code>tx rate limited bytes</code>	Bytes dropped by egress rate limiter from the vSwitch point of view.



Note

The RX and TX statistics displayed by the representative VF netdev is inverted from the actual VF statistics. This convention is consistent with that of the representative physical port netdev.

3.3.2 Example Bridge Setup

Assuming br0 for these examples, create the bridge:

```
ovs-vsctl add-br br0
```

Set the bridge OpenFlow protocol version by entering on the host, logged in as root:

```
ovs-vsctl set Bridge br0 protocols=OpenFlow13
```

and add ports to the bridge specifying the OpenFlow port numbers by using:

```
ovs-vsctl add-port br0 $PORT -- set Interface $PORT ofport_request=$PORT_NUM
```

\$PORT, in this case, corresponds to any representative netdev associated with the NFP. These are generally in the form sdn_pX for physical ports, or sdn_vY.Z for Virtual Functions. \$PORT_NUM is the OpenFlow number that is requested to be associated with that port, and is a user configurable integer with a minimum value of 1. This OpenFlow port number is used when specifying ingress port match field or output actions with OpenFlow flows.

The OpenFlow/Port mapping can be inspected by using:

```
ovs-ofctl -OOpenFlow13 show br0
```

Setting the OpenFlow protocol version of the bridge implies that OpenFlow operations such as those performed with the ovs-ofctl utility must be performed with compatible OpenFlow versions, for example adding flow entries:

```
ovs-ofctl -OOpenFlow13 add-flow br0 in_port=$INPORT,actions=output:$PORT_NUM
```

Set the bridge failure mode (behaviour when connection to the controller is lost). By default (usually at startup) Open vSwitch enters standalone failure mode. This may be changed by using:

```
ovs-vsctl set-fail-mode br0 secure
```

Set the datapath maximum idle timeout by using:

```
ovs-vsctl set Open_vSwitch . other_config:max-idle=$TIMEOUT
```

Setting TIMEOUT=-1 (specified in milliseconds) is interpreted as infinity. This can lead to large repeated hash tables in OvS and on the NFP and is therefore not recommended outside of troubleshooting scenarios. The default timeout is 1500ms (1.5s), it can be handy to increase this timeout (or set to infinity) when trying to inspect datapath flows (ovs-dpctl dump-flows) to ensure they are not timed out before you have a chance to see them.

3.3.3 Layer 2 - Learning Bridge

Configure Open vSwitch (explicitly) to act as a learning bridge by using the following command:

```
ovs-ofctl -OOpenFlow13 add-flow br0 actions=NORMAL
```

In Open vSwitch this is simply a flow rule (not strictly a slice of the switch resources) and thus may be combined with other flow rules in the OpenFlow pipeline. The flow rule with action NORMAL is installed by default when the failure mode is standalone and no OpenFlow controller is connected.



Note

Setting the OvS 'no-forward' option on the LOCAL port will prevent traffic being forwarded to the bridge netdev.

```
ovs-ofctl -OOpenFlow13 mod-port br0 br0 no-forward
```

3.3.4 Layer 3 - Routing

Matching on L3 addresses and then modifying L2 addresses as well as TTL/Hop count results in datapath actions that set both the L2 and L3 headers, for example:

```
ovs-ofctl -OOpenFlow13 add-flow br0 \
ip,nw_dst=192.168.1.0/24,actions=set_field:00:1a:2b:3c:4d:5e->dl_src,\
set_field:00:11:aa:22:bb:33->dl_dst,dec_ttl,output:2
```

This rule matches the 192.168.1.0/24 subnet, sets the source and destination ethernet addresses and decrements the IP TTL before outputting to port 2.

3.3.5 Tags and Labels

Adding and removing of tags/labels may be achieved as follows using flow rules.

3.3.5.1 VLAN Tags

To push a VLAN tag and set the VID and PCP:

```
ovs-ofctl -OOpenFlow13 add-flow br0 actions=push_vlan:0x8100,set_field:4196->\
vlan_vid,set_field:3->vlan_pcp,output:2
```



Note

Only an ethertype of 0x8100 is supported and therefore only a single VLAN tag may be pushed onto a packet.

In order to pop a VLAN tag from a packet use the following:

```
ovs-ofctl -OOpenFlow13 add-flow br0 dl_vlan=1,actions=strip_vlan,output:2
```

3.3.5.2 MPLS Labels

Pushing MPLS labels can be done using the following flow rule, which simultaneously sets the label, TC and TTL:

```
ovs-ofctl -OOpenFlow13 add-flow br0 actions=push_mpls:$ETHERTYPE,set_field:100->\>\>
mpls_label,set_field:3->\>mpls_tc,set_mpls_ttl:32,output:2
```

\$ETHERTYPE may be one of either 0x8847 or 0x8848 for unicast or multicast MPLS respectively.

Removing MPLS labels is much the same as popping VLAN tags:

```
ovs-ofctl -OOpenFlow13 add-flow br0 dl_type=$ETHERTYPE_MATCH,actions=pop_mpls:\
$ETHERTYPE_OUT,output:2
```

however \$ETHERTYPE_OUT in this case specifies the ethernet type of the resulting packet and \$ETHERTYPE_MATCH the ethernet type of the ingress packet to be matched.

3.3.6 Running Packet Analyzer Tools (tcpdump) on Representative Netdevs

Due to the architecture of the software, one cannot always run packet analyzer tools like tcpdump on a representative netdev. As soon as all the flows are offloaded to the NFP, the traffic may no longer traverse the host at all. An example of this is simple match port/set port rules. All traffic will be forwarded to the opposite port, and only stats updates will reach the host. In this case, something like `tcpdump -i sdn_p0` will not show any traffic.

To effectively run analysis tools on a representative netdev in this case, one needs to add a tap/mirror on the physical port to send traffic to an unused VF:

```
ovs-ctl configure mirror rx from sdn_p0 to sdn_v0.1
```

Now run tcpdump on the netdev attached to VF 1 on NFP PCIe unit 0.

For more information on configuring mirrors, please refer to Section 3.2.6

3.4 Tunneling

The Agilio OvS software supports accelerated termination and origination of VXLAN, VXLAN-GPE and NVGRE tunnels. The Agilio OvS tunnel configuration is the same as configuring tunnels with unaccelerated Open vSwitch.

Open vSwitch provides a multitude of methods to configure tunnel endpoints, however to make full use of the acceleration of tunnels with Agilio OvS software, the following configuration should be used.

The port on which tunneled traffic ingresses and egresses the system should not be added to the Open vSwitch bridge. Instead, this port (or ports) must be configured with the tunnel endpoint IP address using Linux system tools like `ifconfig` or `ip`. Any other system ports should be added to Open vSwitch and a tunnel configuration created in the Open vSwitch bridge.

The following commands can be used to configure such a tunnel. This example assumes that entunneled traffic will egress the system on netdev `sdn_p0` to a remote peer with IP address `10.1.1.2`. The rules are configured in such a manner that all traffic ingressing from `sdn_v0.0` will be entunneled and sent out of `sdn_p0`. Conversely,

all entunneled traffic ingressing on `sdn_p0`, will be detunneled and sent out of `sdn_v0.0` again. The `sdn_p0` netdev is not added to the Open vSwitch bridge since it is configured with the remote IP address of the tunnel endpoint.

```
ovs-vsctl add-br br0
ovs-vsctl add-port br0 sdn_v0.0 -- set interface sdn_v0.0 ofport_request=1
ovs-vsctl add-port br0 tun0 -- set Interface tun0 type=vxlan options:remote_ip=10.1.1.2 \
options:local_ip=10.1.1.1 options:key=flow ofport_request=2
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1,actions=set_tunnel:0xff,output:2
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=2,actions=1
ifconfig sdn_p0 10.1.1.1/24 up
```

Configuring a NVGRE tunnel can be achieved using the same process described above, where `type=vxlan` is replaced with `type=gre`:

```
ovs-vsctl add-port br0 tun0 -- set Interface tun0 type=gre options:remote_ip=10.1.1.2 \
options:local_ip=10.1.1.1 options:key=flow ofport_request=2
```

Configuring a VXLAN-GPE tunnel can be achieved using the same process described above, where `type=vxlan` and an additional `ext=gpe` parameter is added:

```
ovs-vsctl add-port br0 tun0 -- set interface tun0 type=vxlan options:remote_ip=10.1.1.2 \
options:local_ip=10.1.1.1 options:dst_port=4790 \
options:exts=gpe options:key=flow ofport_request=2
```

In the VXLAN-GPE case the flow rule dictates the next protocol header and `gpe-flags`, this is achieved by following flow rule:

```
ovs-ofctl -O OpenFlow13 add-flow br0 ip,in_port=1,actions=set_field:0x0d->tun_gpe_flags,\
set_field:0x03->tun_gpe_np,set_field:0xff->tun_id,output:2
```

For more information about the Open vSwitch commands used in this example configuration, refer to the Open vSwitch documentation and manual pages for the respective utilities. Also refer to the Open vSwitch documentation to configure other supported tunnels with the Agilio OvS software.

The Agilio OvS software does support header modification actions along with tunnels to the same extent as supported by Open vSwitch. The Agilio OvS software also supports tunnel conversion over all permutations of VXLAN and NVGRE. For example, VXLAN entunneled traffic can be detunneled and entunneled again using NVGRE.

The NFP can support terminating tunnels on up to 32 local IP addresses. The first 32 IP addresses for netdevs that have active link are programmed into the NFP for this purpose. Note that IP addresses in the 127.0.0.0/8 subnet are ignored for this purpose. If there are more than 32 IP addresses configured on the host system, the kernel may return the IP addresses in an undeterministic order resulting in the required IP addresses to be absent from the NFP configuration.

Agilio OvS provides the feature to prioritize the IP address list programmed on the NFP. To set the prioritized IP address list, execute the following command (note that the list is space separated):

```
ovs-vsctl set Open_vSwitch . other_config:priority-offloaded-ip-list=\
"192.168.0.1 192.168.0.2 192.168.0.3 192.168.0.4 192.168.0.5 192.168.0.6"
```

This will result in the first six IP addresses on the NFP to be populated with the specified values. The IP addresses does not necessarily have to be active on the system at this point. Up to 32 IPs can be specified in this list. If more are specified, only the first 32 IPs will be used. If less are specified, all of them will be used and supplimented with any IP addresses configured on the host system that are not already in the priority list.

To view the current configuration, enter the following command:

```
ovs-vsctl get Open_vSwitch . other_config:priority-offloaded-ip-list
```

To remove the priority list, enter the following command:

```
ovs-vsctl remove Open_vSwitch . other_config priority-offloaded-ip-list
```

3.5 Link Aggregation

The Agilio OvS software supports accelerated Link Aggregation, using the tcp-balance bond mode, to up to 32 independent bonds, where each bond consists of up to 128 ports, where each port can be a physical port or a VF.

A bond is typically created between two machines that are connected via multiple physical links. To configure each machine, proceed as follows:

To associate a specific port with a bond, the OpenFlow port bound to the representative netdev is used. Refer to Section 3.3.1 for more information on the representative netdevs.

To create a bond (named `bondbr0` in this example) over four physical ports (`sdn_p0`, `sdn_p1`, `sdn_p2` and `sdn_p3`) to provide redundancy (active / backup links) where one link is used, and another is selected when the current link fails (default behaviour):

```
ovs-vsctl add-bond bondbr0 bond0 sdn_p0 sdn_p1 sdn_p2 sdn_p3
```

To use all available links to achieve increased throughput (active / active links), the bond mode needs to be changed to TCP load balancing, where the link used for a particular packet is selected based on L3 addresses (IPv4 and IPv6) and L4 ports:

```
ovs-vsctl set port bond0 bond_mode=balance-tcp
```



Note

Packets are approximately evenly spread over all ports on the bridge; the capacity of the links is not considered and is assumed equal for all ports.

To enable LACP on a bond, so that the availability of links is actively monitored (needs to be done on at least one of the two machines involved in the bond):

```
ovs-vsctl set port bond0 lacp=active
```

The interval at which LACP packets are transmitted can be configured as being fast (1 s) or slow (30 s) by using one of the following:

```
ovs-vsctl set port bond0 other_config:lacp-time=slow
ovs-vsctl set port bond0 other_config:lacp-time=fast
```

To show the bond configuration:

```
ovs-appctl bond/show bond0
```

When LACP has been enabled, OvS rebalances flows amongst the available slaves. The rebalance interval can be configured by:

```
ovs-vsctl set port bond0 other_config:bond-rebalance-interval=100
```



Note

OvS sends packets to the LOCAL port for each bond. This results in the flow rule including output to the LOCAL port. This cannot be accelerated by Agilio OvS. To prevent this from occurring, and to achieve acceleration, packets must **not** be sent to the LOCAL port, which can be achieved by:

```
ovs-ofctl -OOpenflow13 mod-port bondbr0 bondbr0 no-forward
```



Note

Link Aggregation in OvS is based on the NORMAL rule; links will not be aggregated when the bond bridge does not contain a NORMAL rule.

Should match/actions be required, an additional bridge (named `br0` in this example) is required on which the match/actions are performed, allowing the bond bridge to only have the NORMAL rule. This additional bridge can be connected to the bond bridge using a patch port:

```
ovs-vsctl add-port br0 bond_patch1 -- set Interface bond_patch1 \
    type=patch options:peer=bond_patch2
ovs-vsctl add-port bondbr0 bond_patch2 -- set Interface bond_patch2 \
    type=patch options:peer=bond_patch1
```

3.6 Load Balance with Group Select

The Agilio OvS software supports accelerated load balancing using group select in OvS. Group Select selects and then executes one bucket in a group. The selection algorithm results in equal load sharing between buckets. An accelerated group can be instantiated by using, for example, the following command:

```
ovs-ofctl -O OpenFlow13 add-group br0 group_id=0,type=select,bucket=actions=10, \
bucket=actions=11,bucket=actions=12
```

Buckets may contain any action set allowed by OvS.



Note

Agilio OvS 2.2 supports the ability to combine the group select-based load balancer with other groups such that the offloaded group select is the last group in the chain.

The group could then be applied to traffic by configuring a flow that outputs to that group:

```
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1,actions=group:0
```


**Note**

Load balancing is achieved by applying a Jenkins hash to the source and destination address information layer 3 and layer 4. Traffic that does not contain layer 3 or 4 will not be spread over multiple buckets. Examples of such traffic include MPLS and ARP. This behaviour differs from previous versions of Agilio OvS, which would have balanced on Layer 2.

3.7 Virtual Function Configuration

The following sections describe different Virtual Function (VF) configurations available in the Agilio OvS software. For each of these configurations, the Bus/Device/Function (BDF) address will be required for each individual VF.

Each NFP PCIe unit has a different PCIe physical function on the PCIe bus and each such physical function can have a corresponding set of "slave" VFs. One can find a list of the BDF prefixes for all the VFs on the NFP PCIe units by performing:

```
lspci -d 19ee:6003
```

This should produce output containing multiple lines of the form:

```
XX:0Y.Z Ethernet controller: Netronome Systems, Inc. Device 6003
```

The VFs on any PCIe unit that are available on the system will be at BDF addresses of the form XX:0Y:Z with Y ranging from 8 to f and Z ranging from 0 to 7. For example, if the output shows c4:08.3 Ethernet controller... then VF 3 on NFP PCIe unit 0 will be at c4:08.3. The PCIe BDF address of a particular VF can most easily be found by running `ethtool -i` on the relevant representative netdev, for example:

```
ethtool -i sdn_v0.2 | grep bus-info
bus-info: NFP 0 VF0.2 0000:c4:08.2
```

One can use the `/opt/netronome/libexec/dpdk_nic_bind.py` command to bind specific drivers to VFs. The Agilio OvS software provides two different drivers which may be bound to VFs, the Netronome DPDK driver `nfp_uio` and the Netronome netdev driver `nfp_netvf`.

The user may use VFs 0 to 59 on each of the NFP's PCIe units. VFs 60 to 63 on each NFP PCIe unit are reserved for internal use. Three of these VFs (VFs 60, 61, and 62 on NFP PCIe unit 0) are bound to the `nfp_netvf` driver when the software is started. The resulting netdevs, named `sdn_ctl`, `sdn_pkt`, and `sdn_pri` respectively, are used for communication between the kernel and the NFP. **These netdevs should not be used or reconfigured in any way by the user. Ensure that third party applications (e.g. GNOME NetworkManager) do not try to manage these netdevs.**

3.7.1 Creating VF Network Devices

The Agilio OvS software VF may be used like any other netdev made available by the Linux kernel. In order to create such a netdev, one must bind the `nfp_netvf` driver to a particular VF via its BDF identifier. Refer to Section 3.7 on how to determine the BDF identifier for a particular VF.

Creating a VF netdev can be accomplished by executing the following command with the BDF identifier for the particular VF:

```
/opt/netronome/libexec/dpdk_nic_bind.py -b nfp_netvf c4:08.2
```

Standard Linux netdev naming rules, for example udev rules, will apply when the netdev is created. The name of the netdev can be determined by observing the output in the kernel ring buffer for the specific BDF identifier. For example the following output may be seen after binding the netdev driver to a VF:

```
[ 5947.941387] nfp_netvf 0000:c4:08.2 eth3: Netronome NFP-6xxx VF Netdev: ...
```

This indicates the netdev name is `eth3` for the VF with BDF identifier of `c4:08.2`. Alternatively the netdev name can be observed in the output of the `/opt/netronome/libexec/dpdk_nic_status.py --status` command, e.g.

```
0000:c4:08.2 'Device 6003' if=eth3 drv=nfp_netvf unused=
```

3.7.2 Running a DPDK Application on NFP VFs

Running a DPDK application while binding it to Agilio OvS VFs is usually accomplished by passing command line arguments to the DPDK application. Some of these arguments go to the poll-mode driver (PMD) library used to interface with the VFs. To produce said arguments, one must first know the BDF identifiers that correspond to each VF. Refer to Section 3.7 on how to determine the BDF identifier for a particular VF.

The VFs are required to be bound to the Netronome DPDK driver `nfp_uio` before starting the DPDK application. This can be accomplished by executing the following command with the BDF identifier of the particular VF:

```
/opt/netronome/libexec/dpdk_nic_bind.py -b nfp_uio c4:08.2
```

Here is an example running the DPDK `l2fwd` application on three VFs to bump traffic inline.

```
# Build the application
cd /opt/netronome/srcpkg/dpdk-ns/examples/l2fwd
make RTE_SDK=/opt/netronome/srcpkg/dpdk-ns RTE_TARGET=x86_64-native-linuxapp-gcc

# Configure hugepages (see DPDK manuals)
mkdir -p /mnt/hugetlbfs
echo 1200 > /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages
mount -thugetlbfs none /mnt/hugetlbfs

# Start the l2fwd application on VFs c4:08.1, c4:08.2 and c4:08.3
build/l2fwd -c f -n 4 -d /opt/netronome/lib/librte_pmd_nfp_net.so \
  --log-level=7 -w c4:08.1 -w c4:08.2 -w c4:08.3 -- -p 7
```

It should be noted that the number of hugepages will be allocated per NUMA-Node (two nodes will allocate twice the requested amount). Furthermore, at least 2 GiB of RAM should be left free for the kernel, the OvS datapath and Agilio OvS. Failing to adhere to this requirement could lead to the Linux Out-of-Memory (OOM) killer being invoked.

Most of the DPDK command line parameters are documented in DPDK itself. However, a few are worth noting in the example above:

1. `-d /opt/netronome/lib/librte_pmd_nfp_net.so` tells the application to use the NFP PMD. This is required.
2. `--log-level=7` significantly reduces the debug output in the system, crucial to achieving better performance.
3. Each `-w XX:YY.Z` adds another VF that the application will consider using, as a DPDK port numbered incrementally from 0. However the `-p N` specifies a bitmask of those ports that will actually be used by the application. In this example, `-p 7` indicates 0b111, indicating that all three ports (numbered 0, 1 and 2) are used.

See the DPDK documentation for more details.



Note

The l2fwd example application forwards traffic between pairs of adjacent VFs when transmitting. If the application opens an odd number of VFs, then the last VF will send traffic out the VF it came in on. Note also that the l2fwd application modifies the MAC addresses of packets that traverse it.

3.7.3 Virtualization Configuration

3.7.3.1 Host Setup for Virtualization

Intel Virtualization Technology (VT) and Intel Virtualization Technology for directed I/O (VT-d) must be enabled in the host's BIOS. Afterward, the kernel command line parameters must be set to enable IOMMU. This can be done as follow:

Edit grub config `/etc/default/grub`. Edit either `GRUB_CMDLINE_LINUX` or `GRUB_CMDLINE_LINUX_DEFAULT` and include the following parameters:

- `intel_iommu=on`
- `iommu=pt`

For example:

```
GRUB_CMDLINE_LINUX_DEFAULT="intel_iommu=on iommu=pt"
```

Update grub and reboot by entering:

- `update-grub2`
- `reboot`

After the system has rebooted, verify the kernel command line with the following command:

```
cat /proc/cmdline.
```

For example, the output would look similar to this:

```
BOOT_IMAGE=/boot/vmlinuz-3.19.0-65-generic \
```

```
root=UUID=baf2673a-84d1-4dc0-9f08-0b4e4d757de2 ro \
console=ttyS0,115200n81 intel_iommu=on iommu=pt
```

Using virtualization with the Agilio OvS software requires `kvm` and `virt-manager` installed on the host system. These packages can be installed with the following commands on Ubuntu systems (similar packages should be installed on other Linux distributions using the relevant package management utility):

```
apt-get install kvm virt-manager
```

In order to allow `virt-manager` to memory map the Netronome libraries, the following changes must be made to the `apparmor` profile:

The following line

```
/opt/netronome/lib/** rm,
```

must be added to the `/etc/apparmor.d/usr.lib.libvirt.virt-aa-helper` file before

```
{media,mnt,opt,src}/** r,
```

Save the changes to the file and reload `apparmor` for the changes to take effect with the following commands:

```
service apparmor reload
```

```
service libvirt-bin restart
```

Start `virt-manager` and use the X11 GUI¹ to create a VM with the following recommended options:

- CPU: at least 4 CPU's (recommended for DPDK)
- CPU Configuration: Copy from host (this is required to compile DPDK)
- Memory (at least 2GB recommended)

Refer to the online documentation for `virt-manager` for more details to install and configure VMs.

3.7.3.2 Guest Setup for Virtualization

Depending on the guest OS distribution, there are three available options for Agilio OvS guest tools installation. In order to setup virtualization on the guest system, all the Agilio OvS required packages listed in the Getting Started Guide must be installed in the guest. Installation packages are available on the Netronome Support site.

- The following command can be used for a Debian package based installation:

```
dpkg --install nfp-bsp-release-2015.11-dkms_2016.7.13.1116-1_all.deb \
netronome-dpdk_16.04-0000_amd64.deb nfp-uio-dkms_0.0.1_amd64.deb
```

- The following command can be used for a RHEL/CentOS package based installation:

```
yum install agilio-ovs-common-2.2-1.el7.centos.x86_64.rpm \
openvswitch-dkms-2.5.1-1.el7.centos.x86_64.rpm \
agilio-ovs-dpdk1604-2.2-1.el7.centos.x86_64.rpm
```

- Alternatively, to use the Agilio OvS source installation, copy the installation tarball to the guest and extract the tarball. Change to the extracted directory and enter `make vm_install`.

¹VMs can also be created with commandline tools by knowledgeable users.

Installation is complete once the following message is displayed on screen:

```
*** Agilio OvS Software 2.2 Virtual Machine Installation Complete!
Build Number: 'build number'
```

The next step is to shutdown the VM and configure PCI passthrough.



Note

The package based installations only provide DPDK 16.04 in `/opt/netronome/src/dpdk-ns.16.04/`.

The Agilio OvS 2.2 source installation provides various DPDK versions, including DPDK 1.8, DPDK 2.2 and DPDK 16.04. The last built DPDK version is symlinked to `/opt/netronome/srcpkg/dpdk-ns`. If one wants to build another version of DPDK, simply change to the `/opt/netronome/srcpkg` path and execute one of the following commands to build another version of DPDK:

- `make build_dpdk1.8`
- `make build_dpdk2.2`
- `make build_dpdk16.04`

The same process can be used on the host installation of Agilio OvS 2.2 if the software was installed with the source tarball. DPDK 2.2 is built by default for any source installation of Agilio OvS.

3.7.3.3 PCI Passthrough using VFIO Setup

To configure PCI passthrough, Agilio OvS software should be running. To start Agilio OvS software enter the following command: `ovs-ctl start`

Once Agilio OvS software is running, the PCI addresses of the VFs can be determined from the corresponding representative netdevs in the BDF format. Refer to Section 3.7 to determine the BDF identifier for a particular VF. The selected BDF addresses must be replaced in the commands below.

Enter the following commands as root on the host system to install the VFIO module and set it to accept Netronome VFs as valid passthrough devices:

- `modprobe vfio-pci`
- `echo "19ee 6003" > /sys/bus/pci/drivers/vfio-pci/new_id`

To bind the VFs to the vfio-pci passthrough driver, execute the following command for each VF:

- `/opt/netronome/libexec/dpdk_nic_bind.py -b vfio-pci <PCI_BDF>`

(where `<PCI_BDF>` is a bus:device.function address, for example 05:08.1)

On the VM system, using `virt-manager`, click on "Add Hardware", choose PCI Host Device and pick the correct VFs. After starting the VM, the VFs should be visible as PCI devices. Either the `nfp_uio` (the Netronome DPDK module) or the `nfp_netvf` (the Netronome netdev module) can be bound to the devices

individually. Refer to Section 3.7.1 and Section 3.7.2 covering netdev and DPDK configuration respectively for VFs. To determine the BDF addresses of all the VFs inside the VM the following command may be used:

```
lspci -d 19ee:6003
```

The following VM configuration has been observed to cause IOMMU mapping errors when rebooting the VM and should be avoided. If such a host/guest configuration is required, the recommended solution is to power off the VM and start it up again instead of rebooting the guest. This problem is not specific to Netronome hardware.

The IOMMU mapping errors have been observed on the following host systems:

- CPU Architecture: Ivy Bridge
- Linux Kernel: 3.13.0-46 and 3.13.0-49
- Linux Distributions: Ubuntu 14.04
- Kernel Parameters:
 - intel_iommu=on
 - iommu=pt
 - default_hugepagesz=1G
 - hugepagesz=1G
 - hugepages=8

The guest system must be a QEMU VM with the following parameters passed to QEMU:

- -enable-kvm
- -device pci-assign
- -mem-path /mnt/huge

The following errors are expected to be seen on the host system in the kernel ring buffer:

```
[ 1798.594046] dmar: DRHD: handling fault status reg 202
[ 1798.596735] dmar: DMAR:[DMA Read] Request device [05:08.1] fault addr 12800000
[ 1798.596735] DMAR:[fault reason 06] PTE Read access is not set
```

This may cause long latencies and possibly cause the host system to become unstable.

3.7.4 Checksum Offload Configuration

3.7.4.1 Linux Kernel Netdev (nfp_netvf driver)

The Agilio OvS checksum offload capabilities are enabled by default when a VF is bound to the `nfp_netvf` driver. Both RX and TX checksum offloading is supported.

`ethtool` can be used to view the status of RX and TX checksum offloading:

```
ethtool -k eth0 | grep checksum
rx-checksumming: on
tx-checksumming: on
    tx-checksum-ipv4: on
```

```
tx-checksum-ip-generic: off [fixed]
tx-checksum-ipv6: on
tx-checksum-fcoe-crc: off [fixed]
tx-checksum-sctp: off [fixed]
```

To disable RX checksum offloading:

```
ethtool -K eth0 rx-checksum off
Actual changes:
rx-checksumming: off
```

To disable TX checksum offloading:

```
ethtool -K eth0 tx-checksum-ipv4 off tx-checksum-ipv6 off
Actual changes:
tx-checksumming: off
    tx-checksum-ipv4: off
    tx-checksum-ipv6: off
```

In both cases, replacing "off" with "on" will enable the feature. The Agilio OvS software makes no distinction between IPv4 and IPv6 where TX checksum offloading is concerned. Thus if either is enabled, the feature will be enabled for both. If the user desires to disable the feature, both need to be disabled.

3.7.4.2 DPDK (nfp_uio driver)

Agilio OvS supports the same checksum offload features for DPDK as the Linux kernel network driver. Refer to the DPDK documentation on how to configure the offloads programmatically.

3.7.5 Quality of Service Configuration

3.7.5.1 VM rate limiter configuration

The rate limiters can be used to limit traffic to and/or from a VF(Virtual Function). Since a VF is used by a Virtual Machine (VM), this will effectively rate limit the traffic to and from the VM. Egress rate limiters limit traffic going from the virtual switch to a VF. Ingress rate limiters limit traffic coming from a VF to the virtual switch. The following commands can be used to configure VF rate limiters.

```
ovs-vsctl set interface <intf> ingress_policing_rate=<rate>
ovs-vsctl set interface <intf> ingress_policing_burst=<burst>

ovs-vsctl set interface <intf> egress_policing_rate=<rate>
ovs-vsctl set interface <intf> egress_policing_burst=<burst>
```

Where <intf> is sdn_v0.0 - sdn_v0.59, or sdn_v1.0 - sdn_v1.59 depending on the amount of PCIE busses in the system. <rate> is the maximum allowed rate, in kilobits per second (Kbps) <burst> is the maximum burst size, in kilobits.

The <rate> parameter can be between 10 (10Kbps) and 40000000 (40Gbps).

Note that a too large burst size can cause the actual rate to be higher than the desired rate. For very low rate and small burst size configurations, the burst size is adjusted to allow at least one jumbo packet to be processed by the rate limiter.

Rate limiters are refreshed every 5 to 10ms. Measuring the rate limiter output rate over small time intervals (e.g. 5 to 10ms) will show inaccuracies. However, these will average out over longer intervals (1s or more).

3.7.5.2 VM rate limiter example:

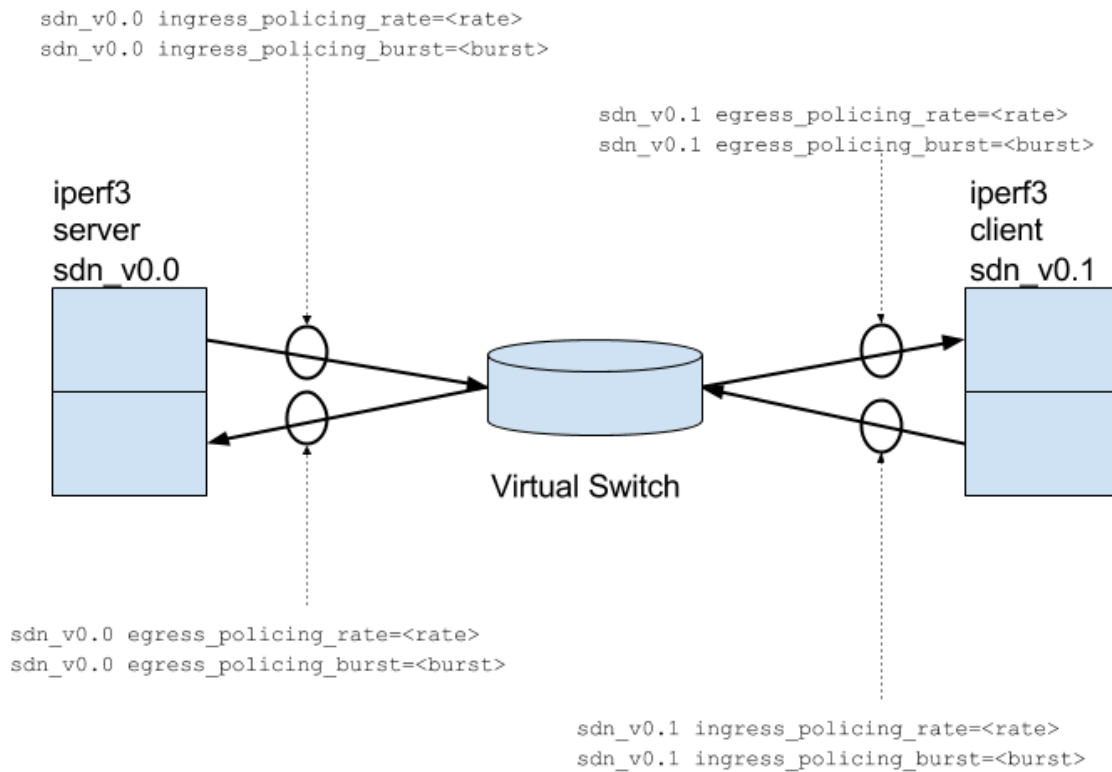


Figure 3.1. VM Rate limiter example setup

The image above shows the setup used to test the rate limiters using iperf3. The server runs on VF0 (sdn_v0.0) and the client runs on VF1 (sdn_v0.1). In this type of test the performance is measured by the amount of data the client sends to the server, i.e. the bottom arrow, going from right to left. To effectively rate limit traffic in this scenario, one only needs the sdn_v0.0 egress limiter or the sdn_v0.1 ingress rate limiter. If both are used, the resultant rate will be determined by the limiter with the lowest rate.

Server setup:

```
SDN_VF_0=sdn_v0.0
# This setup assumes sdn_v0.0 is already bound to nfp_netvf,
# and that the netdev name is eth0
if_0=eth0
# create bridge with 2 ports
ovs-vsctl --if-exists del-br $BR
ovs-vsctl add-br $BR
ovs-ofctl del-flows -Oopenflow13 br0
ovs-vsctl set bridge $BR protocols=OpenFlow13
#create ns1
ip netns add ns1
ovs-vsctl add-port br0 $SDN_VF_0 -- set Interface $SDN_VF_0 ofport_request=1
```



```
ip link set $if_0 netns ns1
ip netns exec ns1 ifconfig $if_0 12.0.0.1/24 mtu 9000 up
ip netns exec ns1 ifconfig
#add flows
ovs-ofctl add-flow -Oopenflow13 br0 in_port=1,actions=output:2
ovs-ofctl add-flow -Oopenflow13 br0 in_port=2,actions=output:1
ip netns exec ns1 iperf3 -s
```

Client setup:

```
SDN_VF_1=sdn_v0.1
# This setup assumes sdn_v0.1 is already bound to nfp_netvf,
# and that the netdev name is eth1
if_1=eth1
ip netns add ns2
ovs-vsctl add-port br0 $SDN_VF_1 -- set Interface $SDN_VF_1 ofport_request=2
ip link set $if_1 netns ns2
ip netns exec ns2 ifconfig $if_1 12.0.0.2/24 mtu 9000 up
echo "ip netns exec ns2 iperf3 -c 12.0.0.1 -i2 -O2 -t30"
```

Rate limiter setup:

```
ovs-vsctl set interface sdn_v0.1 ingress_policing_rate=2000 #2Mbit/sec rate
ovs-vsctl set interface sdn_v0.1 ingress_policing_burst=100000 #100KB burst
```

Result:

```
ip netns exec ns2 iperf3 -c 12.0.0.1 -i3 -O2 -t30
Connecting to host 12.0.0.1, port 5201
[ 4] local 12.0.0.2 port 60101 connected to 12.0.0.1 port 5201
[ ID] Interval           Transfer     Bandwidth       Retr   Cwnd
[ 4]  0.00-3.00        sec    769 KBytes  2.10 Mbits/sec   56    26.2 KBytes
[ 4]  3.00-6.00        sec    769 KBytes  2.10 Mbits/sec   56    26.2 KBytes
[ 4]  6.00-9.00        sec    734 KBytes  2.00 Mbits/sec   56    26.2 KBytes
[ 4]  9.00-12.00       sec    612 KBytes  1.67 Mbits/sec   53    26.2 KBytes
[ 4] 12.00-15.00       sec    734 KBytes  2.00 Mbits/sec   55    26.2 KBytes
[ 4] 15.00-18.00       sec    786 KBytes  2.15 Mbits/sec   56    26.2 KBytes
[ 4] 18.00-21.00       sec    769 KBytes  2.10 Mbits/sec   56    26.2 KBytes
[ 4] 21.00-24.00       sec    734 KBytes  2.00 Mbits/sec   57    26.2 KBytes
[ 4] 24.00-27.00       sec    769 KBytes  2.10 Mbits/sec   56    26.2 KBytes
[ 4] 27.00-30.00       sec    786 KBytes  2.15 Mbits/sec   55    26.2 KBytes
- - - - -
[ ID] Interval           Transfer     Bandwidth       Retr
[ 4]  0.00-30.00       sec   7.29 MBytes  2.04 Mbits/sec   556
[ 4]  0.00-30.00       sec   7.29 MBytes  2.04 Mbits/sec
                                sender
                                receiver

iperf Done.
```

Note the -O2 parameter when invoking the iperf3 client. This is to ignore the initial TCP slow start and the burst that follows due to the rate limiter allowing a 100KB burst.

3.7.5.3 VM QoS

The VM QoS feature allows the user to define up to 8 different traffic classes for use in egress scheduling. The class is determined by the VF or network port. Each physical egress network port can have relative bandwidths assigned to each of the 8 classes. The scheduling algorithm used is Deficit Weighted Round Robin (DWRR).

The DWRR scheduling algorithm works as follows: Bandwidth will be allocated to each class in a relative fashion. Each class will be programmed with a unique weight. Weights are assigned in a relative manner across the classes according to the desired allocation of intended bandwidth for each class. For example, in an oversubscribed application if it is desired that one class of service receives 100 times the bandwidth as compared to a different class, the weight for that class would be programmed to 100x the other class.

- DWRR Example, Sum of all weights = 402,000
- Class 0, large weight = 200,000 : Relative bandwidth = 49.75%
- Class 1, medium weight = 100,000 : Relative bandwidth = 24.88%
- Class 2, medium weight = 100,000 : Relative bandwidth = 24.88%
- Class 3, small port weight = 2,000 : Relative bandwidth = 0.50%

Note that if no traffic of a certain class is scheduled, the relative bandwidth of that class can be used by the other classes. E.g in the example above, if there is only traffic assigned to class 3, it can use 100% of the bandwidth. If at the same time, traffic assigned to classes 0, 1, and 2 appear, the traffic assigned class 3 will be dropped up until the point where it reaches its assigned relative bandwidth.

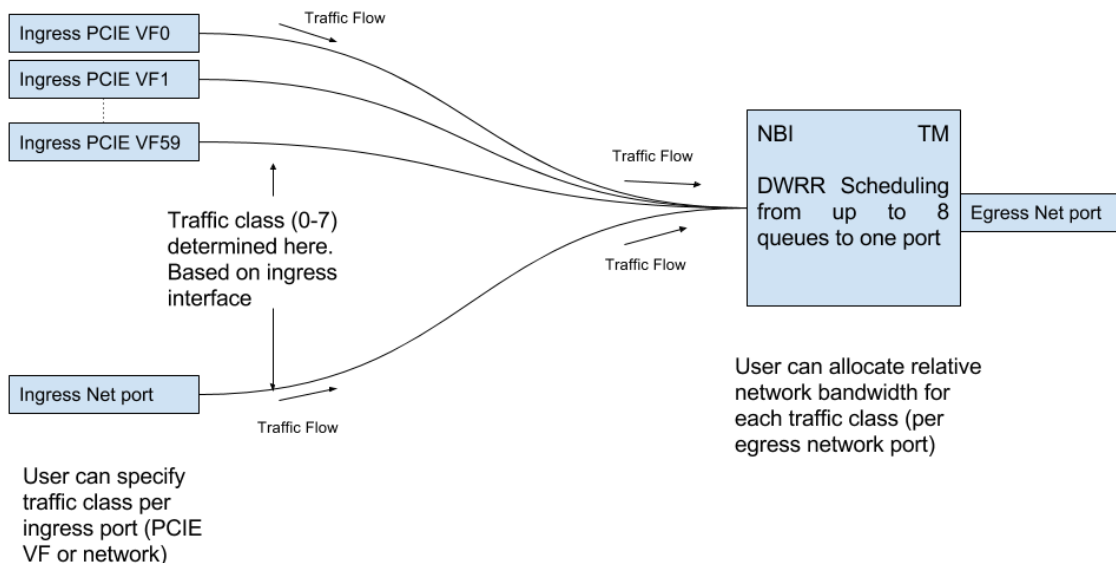


Figure 3.2. VM QoS overview

3.7.5.4 VM QoS example

Clear QoS (per port or all):

```
ovs-vsctl clear port sdn_p0 qos
ovs-vsctl clear port sdn_p1 qos
ovs-vsctl --all destroy qos
```

Add new QoS class to port (sdn_p0 is assigned class 2, sdn_p1 is assigned class 0). Also applicable to VF's (sdn_vx.y)

```
ovs-vsctl set port sdn_p0 qos=@newqos -- \
--id=@newqos create qos type=nfp-dwrr-qos other-config:ingress-class=2
```

```
ovs-vsctl set port sdn_p1 qos=@newqos -- \
--id=@newqos create qos type=nfp-dwrr-qos other-config:ingress-class=0
```

Modify Ingress Traffic Class per port:

```
ovs-vsctl list port sdn_p0 # Observe _UUID for QoS entry
ovs-vsctl set qos 0e5544fb-d90a-4128-8a0b-431ab8ff0233 other-config:ingress-class=1
```

View all QoS entries:

```
ovs-vsctl list qos
```

Change the DWRR weight of a class at an egress port:

```
ovs-vsctl set port sdn_p1 qos=@newqos -- \
--id=@newqos create qos type=nfp-dwrr-qos other-config:ingress-class=2 \
queues:0=@q0 queues:1=@q1 queues:2=@q2 queues:3=@q3 \
queues:4=@q4 queues:5=@q5 queues:6=@q6 queues:7=@q7 -- \
--id=@q0 create Queue other-config:dwrr-relative-bandwidth=1500 -- \
--id=@q1 create Queue other-config:dwrr-relative-bandwidth=10000000 -- \
--id=@q2 create Queue other-config:dwrr-relative-bandwidth=5000000 -- \
--id=@q3 create Queue other-config:dwrr-relative-bandwidth=2000000 -- \
--id=@q4 create Queue other-config:dwrr-relative-bandwidth=40000 -- \
--id=@q5 create Queue other-config:dwrr-relative-bandwidth=5000 -- \
--id=@q6 create Queue other-config:dwrr-relative-bandwidth=4000 -- \
--id=@q7 create Queue other-config:dwrr-relative-bandwidth=1500
```

In the above example, classes 0 and 7 will be allocated 0.0088% of the relative bandwidth. In the case where a 10Gbps port is oversubscribed, classes 0 and 7 will each be allocated 880Kbps. Note that it is not possible to set relative bandwidth weight to absolute 0.

The minimum value of "dwrr-relative-bandwidth" is 1. However, the user should not set the minimum lower than the largest expected frame size. As a guideline, values smaller than the interface MTU should not be used. The maximum value is 16777215. Attempting to set a value of 0 will have no effect.

The ovs-vsctl command above is just an example to configure all 8 queues at once. Queues can be added one by one or in any other fashion. See below.

To modify only a single queue after it has been configured:

```
ovs-vsctl set queue <UUID> other-config:dwrr-relative-bandwidth=1000
```

UUID can be determined with 'ovs-vsctl list queue'

It is important to note that if a user disconnects the QoS class from the port without clearing the queues first (as shown above), the code cannot reset the queues to its default state and the last QoS configuration will remain active on the particular port.

To view QoS configuration, it is recommended to use this utility:

```
ovs-appctl qos/show sdn_p0
```

If an invalid configuration is attempted, the above ovs-appctl command will show what the actual configuration is. OVSDb will unfortunately show the invalid configuration and users may think the configuration was successful.

Note that by default all ingress ports and VF's are set to class 0. On the egress network ports, class 0 has a DWRR weight of 200000 and classes 1-7 have weights of 1. If the user decides to give precedence to VM traffic over network-network traffic, he/she is responsible to expressly configure it as such.

3.8 VirtIO Relay

The optional VirtIO relay daemon (`virtiorelayd`) can relay packets between SR-IOV VFs on the host, and VirtIO devices in a QEMU guest VM. The benefit of this is that any guest OS supporting VirtIO can be used, at the expense of losing the native SR-IOV performance advantage. In a nutshell, VirtIO relay monitors OVSDb to discover which SR-IOV VF ports to relay, uses DPDK to interface with such VFs, and uses vhost-user unix sockets to exchange VirtIO information with QEMU virtual machines. VFs can currently only be relayed 1:1 with VirtIO instances going to VMs.²

To enable this daemon, add the following to `/etc/netronome.conf`:

```
SDN_VIRTIORELAY_ENABLE=y
```

This must be done prior to starting Agilio OvS with the `ovs-ctl start` command. The `ovs-ctl status` command will also show the status of the VirtIO relay. Further configuration details are supplied below.

3.8.1 Requirements

Some external requirements must be met in order to use VirtIO relay successfully:

- QEMU version 2.5 (or newer) must be used for the virtual machine hypervisor. The older QEMU 2.3 and 2.4 do work with `virtiorelayd`, though there are bugs³, less optimised performance and missing features⁴. Some specific QEMU settings are required, detailed in a subsection below.
- libvirt 1.2.6 or newer (if using libvirt to manage VMs, manually scripted QEMU commandline VMs don't require libvirt)
- 2M hugepages must be configured in Linux, a corresponding `hugetlbfs` mountpoint must exist, and at least 320 hugepages must be free for use by VirtIO relay.
- The SR-IOV VFs added to the relay must be bound to the `nfp_uio` driver on the host (all the VFs are bound to `nfp_uio` after Agilio startup by default, see also Section 3.7.2)

3.8.2 libvirt and apparmor

Libvirt's apparmor permissions might need to be modified to allow read/write access to the hugepages directory and library files for QEMU:

```
# in /etc/apparmor.d/abstractions/libvirt-qemu
```

²Future versions of VirtIO relay may make it possible to combine many VFs into fewer VirtIO pipes, or few VFs into many VirtIO instances.

³In particular, shared memory hugepages may be leaked when restarting VMs.

⁴For example, QEMU 2.5 is required to use VirtIO 1.0 successfully.

```
# for latest QEMU
/usr/lib/x86_64-linux-gnu/qemu/* rmix,
# for access to hugepages
owner "/mnt/huge/libvirt/qemu/*" rw,
owner "/mnt/huge-1G/libvirt/qemu/*" rw,
```

Be sure to substitute the `hugetlbfs` mountpoints that you use into the above. It may also be prudent to check for any deny lines in the `apparmor` configuration that may refer to paths used by `VirtIO` relay, such as `hugepage` mounts or `vhostuser` sockets (default `/tmp`).

3.8.3 hugepages

`VirtIO` relay requires 2M hugepages and `QEMU/KVM` can give better performance with 1G hugepages. To set up the system for use with `QEMU` and `VirtIO` relay, the following should be added to the Linux kernel command line parameters:

```
hugepagesz=2M hugepages=768 default_hugepagesz=1G hugepagesz=1G hugepages=8
```

The following could be done after each boot:

```
#Reserve at least 320 * 2M for virtio-relay:
echo 768 > /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages
#Reserve 8G for application hugepages (modify this as needed for VMs):
echo 8 > /sys/kernel/mm/hugepages/hugepages-1048576kB/nr_hugepages
```

Note that reserving hugepages after boot may fail if not enough contiguous free memory is available, and it is therefore recommended to reserve them at boot time with Linux kernel command line parameters. This is especially true for 1G hugepages.

`hugetlbfs` needs to be mounted on the filesystem to allow applications to create and allocate handles to the mapped memory. The following lines mount the two types of hugepages on `/mnt/huge` (2M) and `/mnt/huge-1G` (1G):

```
grep hugetlbfs /proc/mounts | grep -q "pagesize=2M" || \
( mkdir -p /mnt/huge && mount nodev -t hugetlbfs -o rw,pagesize=2M /mnt/huge/ )
grep hugetlbfs /proc/mounts | grep -q "pagesize=1G" || \
( mkdir -p /mnt/huge-1G && mount nodev -t hugetlbfs -o rw,pagesize=1G /mnt/huge-1G/ )
```

Finally, if `libvirt` is to be used, `libvirt` requires a special directory inside the hugepages mounts with the correct permissions in order to create the necessary per-VM handles:

```
mkdir /mnt/huge-1G/libvirt
mkdir /mnt/huge/libvirt
chown libvirt-qemu:kvm -R /mnt/huge-1G/libvirt
chown libvirt-qemu:kvm -R /mnt/huge/libvirt
```

On CentOS the `libvirt-qemu` username above needs to be substituted with the `qemu` username.



Note

After these mounts have been prepared, the `libvirt` daemon will probably need to be restarted.

It is also possible (and simpler) to have the VMs use 2M hugepages as well instead of 1G hugepages, then one does not need to manually reserve 1G hugepages at boot time. There is a performance benefit to using 1G hugepages for the VMs however, especially if they use a few gigabytes of RAM.

3.8.4 VirtIO relay configuration

To enable this daemon, add the following to `/etc/netronome.conf`:

```
SDN_VIRTIORELAY_ENABLE=y
```

The daemon can be further configured by command line parameters, which can be specified in the following variable:

```
SDN_VIRTIORELAY_PARAM="--cpus=1,2 --vhost-username=libvirt-qemu\  
--vhost-groupname=kvm --huge-dir=/mnt/huge --virtio-cpu=0:1,2"
```

By default, when this variable is not specified, VirtIO relay will use CPUs 1 and 2, the vhost username and group will be set according to the Linux distribution in use and the hugepage directory will be autodetected from the mount points. The available command line options are:

```
USAGE: virtiorelayd --cpus=<val> [--nodaemon] [--loglevel=<val>] [--pid-path=<val>] \  
[--huge-dir=<val>] [--ovsdb-sock=<val>] [--vhost-username[=<val>]] \  
[--vhost-groupname[=<val>]] [--vhost-path=<val>] [--vhost-socket=<val>]
```

VirtIO Relay daemon: forward packets between SR-IOV VFs (serviced by DPDK) and VirtIO network backend.

Options:

```
-C<val>    or --cpus=<val>  
    CPUs to use for worker threads (either comma separated integers or hex  
    bitmap starting with 0x)  
-n         or --nodaemon  
    Don't daemonize (default: run as daemon)  
-l<val>    or --loglevel=<val>  
    Set log threshold 0-7 (least to most verbose) (default: 6)  
-p<val>    or --pid-path=<val>  
    PID file (virtiorelayd.pid) will be written to this directory (default:  
    /var/run)  
-H<val>    or --huge-dir=<val>  
    The mount path to the hugepages (default: /mnt/huge)  
-O<val>    or --ovsdb-sock=<val>  
    OVSDb unix domain socket file (default:  
    /usr/local/var/run/openvswitch/db.sock)  
-I         or --ipc  
    Use IPC to add/remove VF&virtio instead of OVSDb (default: use OVSDb  
    monitoring)  
-u[<val>] or --vhost-username[=<val>]  
    vhost-user unix socket ownership username, omit value to inherit process  
    username (default: libvirt-qemu)  
-g[<val>] or --vhost-groupname[=<val>]  
    vhost-user unix socket ownership groupname, omit value to inherit process  
    groupname (default: kvm)  
-V<val>    or --vhost-path=<val>  
    vhost-user unix socket directory path (default: /tmp)  
-S<val>    or --vhost-socket=<val>
```

```

vhost-user unix socket file name, must contain exactly one %u to denote
VirtIO ID (default: virtiorelay%u.sock)
-c<val>    or --virtio-cpu=<val>
  CPU(s) to use for specified virtio, format is '<virtio>:<cpu>[,<cpu>]'.
  Can be specified multiple times to configure CPUs for each virtio ID.
-v         or --version
  Show version number and exit

```

3.8.5 Adding VF ports to VirtIO relay

SR-IOV VF ports are added to an Agilio OvS bridge using the representative netdevs and the `ovs-vsctl` command (see Section 3.3.2). Similarly, to add a VF to the VirtIO relay, the `ovs-vsctl` command can be used with a special `external_ids` value containing an indication to use the relay. The bridge name `br-virtio` in this example is arbitrary, any bridge name may be used:

```

ovs-vsctl add-port br-virtio sdn_v0.42 -- set interface sdn_v0.42 \
external_ids:virtio_relay=1

```

The `ovs-vsctl` command manipulates the OVSDb, which is monitored for changes by VirtIO relay. The preceding command can be interpreted as follows: the VF represented by `sdn_v0.42` will be configured to be relayed to the VirtIO ID 42. This ID is present in the vhost-user unix socket path, which is passed to the QEMU VM (see next section), by default this will be `/tmp/virtiorelay42.sock`. The ID can be any value from 0 to 59, and is determined from the VF used. Note that the ports in the OVSDb remain configured across OvS restarts, and when `virtiorelayd` starts it will find the initial list of ports with associated VirtIO relay indications and recreate the necessary associations.

Changing an interface with no VirtIO relay indication to one with a VirtIO relay indication, or changing one with a VirtIO relay indication to one without a VirtIO relay indication also works. e.g.

```

# add to OvS bridge without VirtIO relay (ignored by VirtIO relay)
ovs-vsctl add-port br-virtio sdn_v0.2
# add VirtIO relay (detected by VirtIO relay)
ovs-vsctl set interface sdn_v0.2 external_ids:virtio_relay=1
# remove VirtIO relay (detected by VirtIO relay and removed from relay,
# but remains on OvS bridge)
ovs-vsctl remove interface sdn_v0.2 external_ids virtio_relay

```

The `external_ids` of a particular interface can be viewed with `ovs-vsctl` as follows:

```

ovs-vsctl list interface sdn_v0.1 | grep external_ids

```

A list of all the interfaces with `external_ids` can be queried from OVSDb:

```

ovsdb-client --pretty -f list dump Interface name external_ids | \
grep -A2 -E "external_ids.*: {.+}"

```

This will result in output similar to:

```

external_ids      : {virtio_relay="1"}
name              : "sdn_v0.1"

external_ids      : {virtio_relay="1"}
name              : "sdn_v0.2"

```

3.8.6 VirtIO relay CPU affinity

The `--virtio-cpu` command line switch can be used to control VF relay CPU affinity. The option may be specified multiple times, once for each VF number. The format of the option is '`--virtio-cpu=<vf>:<cpu>[,<cpu>]`' where `<cpu>` must be a valid CPU enabled in the '`--cpus`' command line option, and VF must be 0 to 59. Specifying two CPUs for a particular VF allows the NFP-to-virtio and virtio-to-NFP relay directions to be serviced by separate CPUs, enabling higher performance to a particular virtio endpoint in a VM. If a given VF is not bound to a CPU (or CPUs), then that VF relay will be assigned to the least busy CPU in the list of CPUs provided in the configuration. It is recommended that all VFs which will be handled by the VirtIO relay be explicitly assigned to CPUs; this is especially important in multi-socket systems to ensure correct NUMA node memory allocation.

3.8.7 Running virtual machines

QEMU virtual machines can be run manually on the command line, or by using libvirt to manage them. To use QEMU manually with the vhost-user backed VirtIO which the VirtIO relay provides, the following example can be used:

```
-object memory-backend-file,id=mem,size=3584M,mem-path=/mnt/huge-1G,share=on,prealloc=on \
-numa node,memdev=mem -mem-prealloc \
-chardev socket,id=chr0,path=/tmp/virtiorelay1.sock \
-netdev type=vhost-user,id=guest3,chardev=chr0,vhostforce \
-device virtio-net-pci,netdev=guest3,csum=off,gso=off,guest_tso4=off,guest_tso6=off,\
guest_ecn=off,mac=00:03:02:03:04:01
```

It is important for the VM memory to be marked as shareable (`share=on`) and preallocated (`prealloc=on` and `mem-prealloc`), the `mem-path` must also be correctly specified to the hugepage mount point used on the system. The path of the socket must be set to the correct VirtIO relay vhost-user instance, and the MAC address may be configured as needed.

Virtual machines may also be managed using libvirt, and this requires some specific XML snippets in the libvirt VM domain specification file:

```
<memoryBacking>
  <hugepages>
    <page size='1048576' unit='KiB' nodeset='0' />
  </hugepages>
</memoryBacking>

<cpu mode='custom' match='exact'>
  <model fallback='allow'>SandyBridge</model>
  <feature policy='require' name='ssse3' />
  <numa>
    <cell id='0' cpus='0-1' memory='3670016' unit='KiB' memAccess='shared' />
  </numa>
</cpu>
```

If only 2M hugepages are in use on the system, the domain can be configured with the following page size:

```
<page size='2' unit='MiB' nodeset='0' />
```


Note, the emulated CPU requires SSSE3 instructions for DPDK support.

The following snippet illustrates how to add a vhost-user interface to the domain:

```
<devices>
  <interface type='vhostuser'>
    <source type='unix' path='/tmp/virtiorelayRELAYID.sock' mode='client' />
    <model type='virtio' />
    <alias name='net1' />
    <address type='pci' domain='0x0000' bus='0x00' slot='0x06' function='0x0' />
  </interface>
</devices>
```

Note: when starting the domain, make sure that the permissions are correctly set on the relay vhost-user socket, as well as adding the required permissions to the apparmor profile. The `--vhost-username` and `--vhost-groupname` command line parameters to the VirtIO relay can also be used to set the permissions on the vhost-user sockets, these default to `user:group of libvirt-qemu:kvm` on Ubuntu, or `qemu:kvm` for CentOS.

It is recommended to pin the virtual CPUs to dedicated host CPUs, for example:

```
<cputune>
  <vcpupin vcpu='0' cpuset='4' />
  <vcpupin vcpu='1' cpuset='5' />
</cputune>
```

This can also be set at runtime with the `virsh vcpupin` command.

3.8.8 Running DPDK on VirtIO inside virtual machines

Instead of using a VirtIO netdev, DPDK can be used with VirtIO inside the VM with a script similar to:

```
mkdir -p /mnt/huge
echo 256 > /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages
grep -q hugetlbfs /proc/mounts || mount -t hugetlbfs nodev /mnt/huge
modprobe uio
modprobe uio_pci_generic
/opt/netronome/libexec/dpdk_nic_bind.py -b uio_pci_generic 00:06.0
# Run a DPDK app with relevant parameters
./trafgen -c 0x3 -n 1 --proc-type auto --socket-mem 512 --file-prefix sink -w 0000:00:06.0 \
-- -b -p 0x1 -w 1a:46:0b:ca:bc:01 -x 1a:46:0b:ca:bc:02 -i 64 -y 64 -j 64 -t 32 -Q 900 -z 64
```

A handy way to test is to use VirtIO PMD inside the VM with `trafgen` sample app in benchmark mode, and set up an OvS flow to echo the VF traffic on the NFP using an OpenFlow in_port action, e.g.

```
ovs-vsctl add-br br-VF0
ovs-vsctl add-port br-VF0 sdn_v0.0 -- set interface sdn_v0.0 external_ids:virtio_relay=1
ovs-ofctl del-flows br-VF0
ovs-ofctl add-flow br-VF0 in_port=1,actions=in_port
```

3.8.9 Using VirtIO 1.0

To enable VirtIO 1.0 (as opposed to legacy VirtIO), the backend virtual PCI device provided by QEMU needs to be enabled. Using QEMU 2.5, you need to supply an extra cmdline parameter to prevent VirtIO 1.0 support from

being disabled (it is disabled by default, since there are apparently still known issues with performance, stability and live migration):

```
-global virtio-pci.disable_modern=off
```

This can be done in a libvirt domain by ensuring the domain spec starts with something like:

```
<domain type='kvm' xmlns:qemu='http://libvirt.org/schemas/domain/qemu/1.0'>
```

and just prior to the closing `</domain>` tag adding the following:

```
<qemu:commandline>
  <qemu:arg value='-global' />
  <qemu:arg value='virtio-pci.disable-modern=off' />
</qemu:commandline>
```

In addition to this, the vhost or vhost-user connected to the device in QEMU must support VirtIO 1.0. The vhost-user interface which VirtIO relay supplies does support this, but if the host is running a Linux kernel older than 4.0, you likely won't have vhost-net (kernel) support for any network interfaces in your QEMU VM which are not connected to VirtIO relay, for example if you have a bridged management network interface. Libvirt will by default use vhost-net for that, you can disable vhost-net by adding `<driver name='qemu' />` to the relevant bridge interface as follows:

```
<interface type='bridge'>
  ...
  <model type='virtio' />
  <driver name='qemu' />
  ...
</interface>
```

To use VirtIO 1.0 with DPDK inside a VM, you will need to use DPDK 16.04. To use a VirtIO 1.0 netdev in the VM, the VM must be running Linux kernel version 4.0 or newer.

3.8.10 VM live migration with libvirt

The VirtIO relay is compatible with QEMU VM live migration as abstracted by libvirt, and has been tested using QEMU 2.5 with libvirt 1.2.16. The VM configuration must conform to some requirements⁵ to allow live migration to take place. In short:

- VM disk image must be accessible over shared network storage accessible to the source and destination machines
- Same versions of QEMU must be available on both machines
- apparmor configuration must be correct on both machines
- VM disk cache must be disabled, e.g. `<driver name='qemu' type='qcow2' cache='none' />` (inside the disk element)
- The hugepages for both machines must be correctly configured
- Ensure both machines have Linux kernels new enough to support vhost-net live migration for any virtio network devices not using the vhostuser interface, or configure such interfaces to only use vanilla QEMU

⁵See for example documentation on RedHat's website, or <http://www.linux-kvm.org/page/Migration#Requirements> or various other articles that can be found in a Google search for "QEMU live migration requirements"

virtio backend support, e.g. `<model type='virtio' />` `<driver name='qemu' />` (inside the relevant interface elements)

The VM live migration can be initiated from the source machine by giving the VM name and target user&hostname as follows:

```
virsh migrate --live <vm_name> qemu+ssh://<user@host>/system
```

The `--verbose` argument can optionally be added for extra information. If all goes well, `virsh list` on the source machine should no longer show `<vm_name>` and instead it should appear in the output of `virsh list` on the destination machine. If anything goes wrong, the following log files often have additional details to help troubleshoot the problem:

```
/var/log/syslog
/var/log/libvirt/libvirt.log
/var/log/libvirt/qemu/<vm_name>.log
```

In the simplest scenario, the source and destination machines have the same VM configuration, particularly with respect to the vhostuser socket used on VirtIO relay. It may be handy to configure the vhostuser socket in the VM to point to a symlink file which links to one of the VirtIO relay sockets. This is one way to allow the source and destination machines to use different vhostuser sockets if necessary. For example, on the source machine one might be using a symlink called `/tmp/vm_abc.sock` linking to `/tmp/virtiorelay1.sock`, while on the destination machine `/tmp/vm_abc.sock` might link to `/tmp/virtiorelay13.sock`.

It is also possible to migrate between machines where one is using VirtIO relay, and the other is using a different virtio backend driver (could be a different vhostuser implementation, or could even be vhost-net or plain QEMU backend). The key to achieving this is the `--xml` parameter for the `virsh migrate` command (`virsh help migrate` reveals: "`--xml <string>` filename containing updated XML for the target").

Here is an example of the procedure to migrate from a vhostuser VM (connected to VirtIO relay) to a non-vhostuser VM:

On the destination machine, set up a libvirt network that you want to migrate the interface onto, e.g. named 'migrate', by passing the following XML file to `virsh net-define <xml_file>` and running it with `virsh net-start migrate`; `virsh net-autostart migrate`:

```
<network>
  <name>migrate</name>
  <bridge name='migratebr0' stp='off' delay='0' />
</network>
```

On the source machine (where the VM is defined to use vhostuser connected to VirtIO relay), dump the VM XML to a file by running `virsh dumpxml <vm_name> >domain.xml`. Edit the `domain.xml` file to change the vhostuser interfaces to be sourced by the migrate network, i.e. change these:

```
<interface type='vhostuser'>
  <mac address='00:0a:00:00:00:00' />
  <source type='unix' path='/tmp/virtiorelay0.sock' mode='client' />
  <model type='virtio' />
  <address type='pci' domain='0x0000' bus='0x00' slot='0x05' function='0x0' />
</interface>
```

to these:

```
<interface type='network'>
```

```
<mac address='00:0a:00:00:00:00' />
<source network='migrate'>
<model type='virtio' />
<address type='pci' domain='0x0000' bus='0x00' slot='0x05' function='0x0' />
</interface>
```

Finally, once you have this modified domain.xml file, the VM can be migrated as follows:

```
virsh migrate --live <vm_name> qemu+ssh://<user@host>/system --xml domain.xml
```

Migrating from a non-virtiorelay machine to a virtiorelay machine follows this same procedure in reverse; a new XML file is made where the migrate network interfaces are changed to vhostuser interfaces.

3.8.11 Inspecting traffic with tcpdump

To inspect traffic going to/from the relayed VFs, the port mirroring feature detailed in Section 3.2.6 can be used to mirror packets on a VF to another VF running a netdev. For example, if you have a virtiorelay using VF0 and wish to inspect the traffic going to and from the VM, you could do the following:

```
ovs-ctl configure mirror rx from sdn_v0.0 to sdn_v0.40
ovs-ctl configure mirror tx from sdn_v0.0 to sdn_v0.40
ethtool -i sdn_v0.40 (note PCI ID)
/opt/netronome/libexec/dpdk_nic_bind.py -b nfp_netvf <PCI ID>
dmesg|tail (note netdev name, e.g. eth4)
tcpdump -i <netdev>
```

3.8.12 Debug counters

Some debug counters (showing packet counts) from the running VirtIO relay can be retrieved using the virtio_relay_stats command, e.g.:

```
# /opt/netronome/bin/virtio_relay_stats
Netronome VirtIO relay
relay1.virtio_rx=204862176
relay1.dpdk_tx=172212896
relay1.dpdk_drop_full=32649344
relay1.dpdk_drop_unavail=0
relay1.dpdk_rx=172209976
relay1.virtio_tx=172152046
relay1.virtio_drop_full=57936
relay1.virtio_drop_unavail=0
relay2.virtio_rx=32550944
relay2.dpdk_tx=32550624
relay2.dpdk_drop_full=320
relay2.dpdk_drop_unavail=0
relay2.dpdk_rx=32550510
relay2.virtio_tx=32547760
relay2.virtio_drop_full=2718
relay2.virtio_drop_unavail=0
```

The drop_full counters increment when the queue in that direction becomes full and VirtIO relay needs to drop packets, e.g. virtio_drop_full means packets were received from the NFP VF, but could not be

delivered to the VM because the virtio queue going to the VM was full. The `drop_unavail` counters increment when that side of the relay is not connected, e.g. `dpdk_drop_unavail` will increment when packets were received from virtio but the VF to which they should be forwarded is not connected to the relay (this is possible because VFs can be dynamically added and removed via DPDK hotplug). In the ideal case when no drops are occurring in the relay, `virtio_rx` will match `dpdk_tx`, and `dpdk_rx` will match `virtio_tx`.

3.9 Health Monitoring

The health of the Agilio OvS software and the platform hardware can be checked by invoking the following command:

```
/opt/netronome/bin/ovs-ctl status
```

The output of this command will indicate the current status (`PASS/FAIL/WARN/SKIPPED`) of each hardware / software component as well as a consolidated result which reflects the health of the system as a whole. The consolidated result is only `PASS` if none of the respective system components is in a `FAIL` state. A syslog message is generated for each of the failures or warnings that are detected at the component level.

The health monitoring utility could potentially skip certain tests when the system is in a state where that test cannot be performed. For example, the Flow Processing Core test will be skipped when the firmware is not loaded. Similarly, many tests will be skipped while the Agilio OvS software is being started or stopped.

The utility can also log warnings for some tests. This will happen when a system component is not in a failed state, but there is a condition which the user should be made aware of.



Note

It is recommended to not run the health monitoring utility while the system is being started or stopped.

The following syslog failure messages could potentially be generated when invoking the Agilio OvS health checking utility:

- *SDN Failure=(Userspace Process Not Running: ovsdb-server) affected_component=(sw.host) action=(restart.software)*
- *SDN Failure=(Userspace Process Not Running: ovs-vswitchd) affected_component=(sw.host) action=(restart.software)*
- *SDN Failure=(Userspace Process Not Running: virtiorelayd) affected_component=(sw.host) action=(restart.software)*
- *SDN Failure=(Kernel Module Not Loaded: nfp) affected_component=(sw.driver) action=(restart.drivers)*
- *SDN Failure=(Kernel Module Not Loaded: nfp_cmsg) affected_component=(sw.driver) action=(restart.drivers)*
- *SDN Failure=(Kernel Module Not Loaded: nfp_fallback) affected_component=(sw.driver) action=(restart.drivers)*
- *SDN Failure=(Kernel Module Not Loaded: nfp_offloads) affected_component=(sw.driver) action=(restart.drivers)*

- *SDN Failure=(Kernel Module Not Loaded: openvswitch) affected_component=(sw.driver) action=(restart.drivers)*
- *SDN Failure=(Firmware Not Loaded) affected_component=(sw.nfp) action=(restart.nfp)*
- *SDN Failure=(FPC Unresponsive) affected_component=(sw.nfp) action=(restart.nfp)*
- *SDN Failure=(NFP Control Channel Unresponsive) affected_component=(sw.nfp) action=(restart.nfp)*
- *SDN Failure=(Fallback Traffic Divert Disabled) affected_component=(sw.nfp) action=(restart.nfp)*
- *SDN Failure=(Ingress NBI Backed Up) affected_component=(sw.nfp) action=(restart.nfp)*
- *SDN Failure=(NFP Not Detected) affected_component=(hw.nfp) action=(replace.nfp)*
- *SDN Failure=(ECC Error Detected) affected_component=(hw.nfp) action=(replace.nfp)*
- *SDN Failure=(Parity Error Detected) affected_component=(hw.nfp) action=(replace.nfp)*
- *SDN Failure=(`[warning] ERR47 Workaround not detected`) affected_component=(sw.host) action=(patch.os)*
- *SDN Failure=(`[warning] NFP Configurator outdated`) affected_component=(hw.nfp) action=(update.nfp)*

All syslog messages are generated with the `ns_sdn_health` tag, `local0` facility, and error or warning level.

The health utility is also provided in the form of a shared library which can be embedded into an user application. Refer to the Agilio OvS PRM for more details on the health utility library usage.

3.10 Trivial NIC Mode

The Trivial NIC mode refers to the mode of the Agilio OvS software when Open vSwitch itself is not running, but the NFP firmware is loaded, configured and able to pass traffic to the host. In this mode, all the representative netdevs are active in the kernel, but only the physical port representative netdevs are operational. Although these representative netdevs can be used for basic networking functions, no offloading or additional features are available in this mode.

For example, an IP address may be assigned to a physical port representative netdev and used for ssh connections to the host. The Trivial NIC mode is not intended to be a high performance mode of the Agilio OvS software, instead this mode can be used to maintaining basic network connections such as NFS mounts or ssh connections over the Agilio OvS physical ports while restarting Agilio OvS.

There are two ways in which the Trivial NIC mode is activated:

1. On host system boot:

During system bootup, the NFP firmware is loaded and configured in the Trivial NIC mode, making physical port representative netdevs available for use at the same time as any other standard NIC netdevs. This provides the ability to mount the root file system over NFS using one of the Intelligent Server Adapter physical ports. Refer to Section 3.11 on details for how this feature can be utilized. This also provides the means to use ssh for connections to the host over one of the Intelligent Server Adapter physical ports.

It is important to understand that the Agilio OvS kernel modules and firmware are added to the current initramfs image during installation. The Agilio OvS components will be stay in the initramfs image whenever it is refreshed, but if another initramfs image is used for booting the host, the Trivial NIC mode will not be activated on system boot.

To disable the Trivial NIC mode on bootup one can add the following to the kernel command line:

```
modprobe.blacklist=nfp
```

2. When Agilio OvS is stopped:

The Trivial NIC mode is activated whenever Open vSwitch is not running. For example, when executing the `ovs-ctl stop` command, the Trivial NIC mode will be activated and when executing the `ovs-ctl start` the Trivial NIC mode will be deactivated again.

Agilio OvS starts faster when transitioning from the Trivial NIC mode to fully operational as opposed to having the NFP firmware reloaded and configured when Agilio OvS is started.

If one wants to avoid using the Trivial NIC mode or wants to completely restart the Agilio OvS software, add the `--nfp-reset` switch to the `ovs-ctl` command to start, stop or restart the Agilio OvS software.

To see whether the Agilio OvS software is in the Trivial NIC mode, one can observe the output of the `ovs-ctl status` command. For example, the following output of the `ovs-ctl status` command indicates that the Trivial NIC is active:

```
Detected Configurator version: 20160211132710

SDN Health Report:
-----
Userspace Process: ovsdb-server          ... [FAIL]
Userspace Process: ovs-vswitchd          ... [FAIL]
Userspace Process: virtiorelayd          ... [SKIPPED]
Kernel Module: nfp                       ... [PASS]
Kernel Module: nfp_cmsg                  ... [PASS]
Kernel Module: nfp_fallback               ... [PASS]
Kernel Module: nfp_offloads               ... [FAIL]
Kernel Module: openvswitch                ... [FAIL]
NFP: Firmware Loaded                     ... [PASS]
NFP: Flow Processing Cores Responsive     ... [PASS]
NFP: Control Message Channel Responsive  ... [PASS]
NFP: Fallback Traffic Divert Disabled     ... [WARN]
NFP: Ingress NBI Backed Up                ... [PASS]
NFP: Card Detected                       ... [PASS]
NFP: ECC Errors                           ... [PASS]
NFP: Parity Errors                        ... [PASS]
NFP: Configurator Version                  ... [PASS]
Host: ERR47 Workaround                    ... [PASS]

Overall System State      [FAIL]
```

The `openvswitch` kernel module may or may not be loaded. The important check to determine whether the Trivial NIC mode is active, is:

```
NFP: Fallback Traffic Divert Disabled     ... [WARN]
```

This check will display 'WARN' when the Trivial NIC mode is active and 'PASS' when it is inactive.



Note

Open vSwitch cannot be used in Trivial NIC mode. VFs will be kept alive, but traffic received or transmitted on a VF will not be allowed to pass through the NFP. The Trivial

NIC mode is only intended to be a transient state of the Agilio OvS software and not intended to be used as a production mode of the software.

3.11 PXE Booting

The Intelligent Server Adapters along with the Agilio OvS software supports PXE booting a server via the physical Intelligent Server Adapter ports. However the PXE booting feature itself is independent of the Agilio OvS software as this occurs during the boot process of the server.

Only 64-bit UEFI boot mode is supported to PXE boot a server via the Intelligent Server Adapter ports. The server must be configured into UEFI boot mode, and the appropriate ports enabled in the BIOS. Refer to the server's BIOS documentation on how to configure this feature.

Agilio OvS does support mounting a NFS filesystem along with the PXE boot feature by means of the Trivial NIC mode. A system can be booted and the root filesystem can be mounted over any representative physical port netdev by using the appropriate PXE configuration. Note that the initramfs image provided by the PXE server must be one obtained from a system with Agilio OvS installed.

All the Intelligent Server Adapter ports are available to PXE boot from including support for breakout cables. By default only port 0 will be enabled for booting after the Agilio OvS installation. The ports available for booting can be observed by executing the following command:

```
# /opt/netronome/bin/nfp-hwinfo | grep eth.*boot
eth0.boot=1
eth1.boot=1
eth2.boot=0
eth3.boot=0
eth4.boot=0
```

To enable or disable any of the other ports for booting, execute the following command:

```
# /opt/netronome/bin/nfp-media --set-hwinfo ethN.boot=1
```

where N = ethernet instance as shown in the 'nfp-hwinfo' output above. The system must be rebooted for the changes to take effect.

If multiple bootable ports are available, the system will try to PXE boot on the first port with an active link starting from the lowest port number. If no bootable ports are configured, the PXE boot firmware will default to using port 0 for PXE booting.

3.12 Connection Tracking (Conntrack)

Conntrack support was introduced into the main branch of OvS in version 2.5. It allows the state tracking of individual packets with respect to their given flow. Giving OvS the ability to handle such functionality means that users can specify rules to replicate security group, access control list, and stateful firewall applications.

Within OvS, Conntrack functionality is offloaded at the kernel level to NetFilter kernel modules when a specific action is set within a flow rule. NetFilter Conntrack treats each packet as coming from a flow this is defined

by its five tuple value (source and destination IP, source and destination port, and transport layer protocol) in a bidirectional capacity. By storing information on each active flow, it can be determined how a packet fits into the flow. For example, in UDP, if a packet is the first viewed from a given flow then it can be considered to be starting a NEW connection, but if the packet is from a flow that has already witnessed packets in both directions then it can be considered to be from an ESTABLISHED flow. Once this Conntrack metadata has been associated with the packet, it can be used to determine how it should be handled.

Agilio OvS offers the ability to offload this Conntrack functionality to the NFP instead of NetFilter. In the NFP datapath, all packets that require Conntrack processing will be handled inline. If the packet travels via kernel space and is to be sent to Conntrack then, instead of forwarding to NetFilter modules, it will be forwarded to `nfp_conntrack.ko` which will relay the request to the NFP, receive a response, and continue kernel space processing.

Conntrack offloading within Agilio OvS requires a specialised firmware to be loaded. To utilise this, it must be selected at install time. This is done by running `make install_conntrack` or `make clean_install_conntrack` as opposed to the standard `make install`. Once the Conntrack version has been installed then it can be run in the same way as standard Agilio OvS, via the `ovs-ctl start/stop/restart` command. We refer to this separate install as Agilio OvS Conntrack.

Conntrack offloading introduces significant processing that is not required in standard (non-Conntrack install) Agilio OvS. Central to this is a new lookup table structure and memory to maintain active flow state, along with the code to replicate the NetFilter Conntrack functionality. Due to this, it should be noted that choosing to install Agilio OvS Conntrack will incur a slight performance decrease if handling non Conntrack rules.



Note

The Conntrack version of the software is in an experimental state which implies that this is not a production version yet. Customers can evaluate this software for deployment planning purposes, however bugs are expected at this stage. Please stay in close synchronization with Netronome support.

3.12.1 Module Requirements

Agilio OvS Conntrack will load a new kernel module called `nfp_conntrack` when it is run. This module will make use of some functions contained within NetFilter modules. These modules should be loaded automatically when the Netronome Conntrack module is loaded. If this does not happen and an ‘unknown symbol’ error appears in the dmesg logs then check if the following modules are loaded (`lsmod`) and if not, try loading them (`modprobe`):

- `nf_conntrack`
- `nf_conntrack_ipv4`
- `nf_conntrack_ipv6`

If helper functionality is to be used within the rules then the appropriate helper kernel modules should be loaded:

- `nf_conntrack_ftp`
- `nf_conntrack_sip`

3.12.2 Helper Functionality

Conntrack also includes the tracking of 'expected' flows. These are flows that have not yet been seen but are anticipated to arrive. An example of such a case is in a FTP Passive connection. Here an FTP control channel is set up and established over well known port 21. This control channel can then be used to negotiate further data channels over potentially random port numbers. Conntrack can examine the negotiation phase of the control channel by inspecting the application of layer of the packets to determine the ports that are agreed on for the data channel. It then 'expects' a flow to appear between these IP addresses and over the negotiated port. When a packet arrives that matches one of these expected flows its state is marked as 'related' meaning that the packets flow is related to, or set up by, an already established and accepted flow.

Agilio OvS Conntrack supports integration with NetFilter modules that perform the packet application layer analysis for both FTP and SIP. The handling of the output from these modules to give expected flows, and the lookup to determine related flows is offloaded to the NFP.

3.12.3 Conntrack in Use (Recirculation)

To start using Conntrack rules within Agilio OvS, the pipeline for rules processing should be considered. Conntrack state tracking, within OvS, is considered an action. This means that a packet must match a rule with a Conntrack action before its state information is known. Because of this, to determine how to handle a packet based on its Conntrack information, the packet needs to be matched a second time. This process is known as 'recirculation' and is key to the implementation of Conntrack. Recirculation means inserting the packet back into the rules matching block as if it is a new packet - the difference being that it will now have Conntrack metadata attached.

Recirculation without Conntrack is an action that is available in OvS 2.5, however, it is not offloaded by standard Agilio OvS. Agilio OvS Conntrack only offloads recirculation actions when used in conjunction with Conntrack.

3.12.4 Conntrack Match Fields

In addition to those match fields in Agilio OvS, several new fields are available in Conntrack mode.

3.12.4.1 Conntrack State

Using the ovs-ofctl app, the Conntrack state of a packet can be matched by calling `ct_state=` followed by one or more flags that are defined by three letter abbreviations. These are flags can either be checked (+), unchecked (-) or don't care (left out):

- `trk` - Tracked
- `new` - New
- `est` - Established

- `rel` - Related
- `rpl` - Reply
- `inv` - Invalid

The `trk` option indicates if a packet has been processed by Conntrack. Every packet that enters the system will, by default, match the `-trk` option, indicating that Conntrack has not been applied to it. If a packet is sent to Conntrack, via an action, and recirculated, then the `+trk` flag will be matched.

`new` and `est` are two core states for a flow. `new` indicates the first time a packet from that flow has been seen. Typically, a packet with this state is initiating a new connection. `est` indicates that the packet is from a well formed flow that is communicating bidirectionally between source and destination. A packet cannot be considered to have both the `new` and `est` flags marked at the same time. The `rpl` flag can, however, be either checked or unchecked along with the `est` flag. `rly` indicates the direction of the packet within the flow. If `rpl` is checked then it is in the opposite direction to the initiator (server to client) while it will remain unchecked if it is in the client to sever direction.

The `rel` flag determines if a packet is from a flow that is related to or spawned by a previously seen flow. The packets tuple will have been extracted from examination of the negotiation fields defined by the applied helper functions.

The `inv` flag is checked if the packet is determined to be invalid. In OvS (depending on NetFilter options selected), this may be due to a checksum error or witnessing a TCP packet that has a sequence number that falls outside the range of already viewed sequence numbers combined with window size. Currently Agilio OvS does not consider checksum calculation or TCP sequence number checking so does not utilise the `inv` flag.

3.12.4.2 Additional Conntrack Match Fields

As well as states, the following fields can be matched on:

- `ct_zone` (16-bit)
- `ct_mark` (32-bit)
- `ct_label` (128-bit)

The Conntrack zone (`ct_zone`) allows the separation of flows that are from different networks but on the same OvS bridge. Take, for example, the case where two private networks that both use the same private address space are connected to the same OvS switch. It is possible that two 5-tuple values on each of these networks could be the same. If Conntrack is applied at the OvS switch then there is scope for confusion between the two and the potential for incorrect state information to be returned. Adding a zone to connections from different networks effectively makes the Conntrack lookup a 6 tuple lookup. This way there can be no confusion caused, even if the 5 tuples are the same. Effectively, the zone splits the Conntrack lookup table into multiple tables (one per zone). When a packet has passed through Conntrack and recirculated, the zone it is applied to will then become part of the packets Conntrack metadata and can be used as a match field.

The Conntrack mark and label (`ct_mark/ct_label`) are both used to pass metadata between recirculations. For all intensive purposes, the only difference is that the mark is a 32-bit value while label is 128-bit. The values can be matched in a recirculated rule if they have been added within the Conntrack action. These match fields can have a mask applied to filter the bits we are concerned with. This is denoted with a `'` followed by the mask in a similar way to the subnet mask on an IP address.

Agilio OvS Conntrack supports the offloading of all of these additional Conntrack match fields to the NFP data path.

3.12.4.3 Recirculation ID

A further match field that is not strictly part of Conntrack but needs to be considered is the `recirc_id` or recirculation id. The `recirc_id` is a 32-bit metadata field that is set when a recirculation action is applied and is then matched on the recirculated packet. This metadata is in place to reduce the possibility of an infinite loop where a packet matches a rule, is recirculated, and matches the same rule again. From a user perspective, this field does not need to be considered as it will be added to the datapath rules automatically when they are converted from user-space rules. Agilio OvS Conntrack supports these automatically added `recirc_id` match fields on the NFP so will also support it as a match field if a user wishes to add it manually.

3.12.5 Conntrack Action Fields

Given the extra match fields required for OvS Conntrack, Agilio OvS Conntrack also supports the offloading of Conntrack related actions. In `ovs-oftl`, an action to send a packet to Conntrack and populate its state information is denoted as `ct()`. The brackets can contain a number of further Conntrack related actions. These are as follows:

- `exec(set_field:n->ct_mark)`
- `exec(set_field:n->ct_label)`
- `table=n`
- `zone=n`
- `alg=(ftp|sip)`

The two execute actions are similar to the set field actions in standard Agilio OvS. These are used to set the `ct_mark` or `ct_label` values within the packet metadata. As with their corresponding match fields, masks can be used to set specific bits only.

Recirculation is signalled by a `table=n` action where `n` indicates the table number that the packet should be recirculated to. This value will typically be set to 0 to recirculate as if it is a new packet entering the system. It should be remembered that the idea of tables does not exist in the datapath so a user-space rule containing `table=n` action will be translated to a recirculate action with a specified `recirc_id`.

The Conntrack zone that a packet is sent to should be specified in the action with a `zone=n`. The zone selected directly affects the state lookup for the packet as described previously. If the packet is to be recirculated then the zone will be added to its metadata and can be used a match field. If a zone is not entered as a Conntrack action parameter then a default zone of 0 will be used.

`alg=`, when applied as a Conntrack action parameter, requests that any matching packet be examined by a helper module to see if any port negotiation is taking place for further 'expected' flows. This action should only be applied to flows that are recognised control channels. The computational cost of sending packets to these helper functions is expensive so should be limited to those that are likely to contain port negotiation data (e.g. match on TCP port 21 for FTP or TCP/UDP port 5060 for SIP). Agilio OvS Conntrack supports `alg=ftp` or `alg=sip` actions (note that SIP is not currently supported in OvS 2.5 but is patched into Agilio). If an 'expected' flow is

found then this tuple will be written to the NFP automatically. Any new flow will be compared to these written expectations to see if it can also be classified as the 'related' state. The user does not need to add any additional command for this functionality to be applied.

A final action that can be used as a Conntrack parameter is `commit`. This is related to how NetFilter Conntrack kernel modules work. In OvS it tells NetFilter to insert the packet information into the main Conntrack lookup table. It, therefore, tends to be used when a packet is classified as the `new` state so that the entry will be stored and matched with any other packets from that flow that follow. In Agilio OvS Conntrack we ignore this action. The reason being that when a lookup is requested on the NFP Conntrack table, it is added by default if it does not already exist. Therefore, the `commit` action is always applied in Agilio. It should be noted that `commit` will still need to be used in some actions (for example, when sending to a `alg`). This is to satisfy some of syntax requirements in the `ovs-ofctl` app.

3.12.6 OvS Rule Configuration for Conntrack

The following are examples of programming Conntrack rules on Agilio OvS Conntrack using the `ovs-ofctl` app.

```
ovs-ofctl add-flow -O Openflow13 br0 "ct_state=-trk,tcp,action=ct(table=0)"
ovs-ofctl add-flow -O Openflow13 br0 "ct_state=+trk+new-est,tcp,action=2"
ovs-ofctl add-flow -O Openflow13 br0 "ct_state=+trk-new+est-rpl,tcp,action=2"
ovs-ofctl add-flow -O Openflow13 br0 "ct_state=+trk-new+est+rpl,tcp,action=2"
```

The first rule above matches any packet that comes into the system. It sends this packet to the Conntrack block and tells it to recirculate back to table 0. On recirculation the packet can then match any of the last 3 rules based on the state that is determined from Conntrack.

```
ovs-ofctl -O Openflow13 add-flow br0 "ct_state=-trk, action=ct(commit,table=0,zone=123,
exec(set_field:0xdead->ct_label, set_field:0xbeef->ct_mark))"

ovs-ofctl -O Openflow13 add-flow br0
"ct_mark=0xbeef,ct_zone=123,ct_label=0xdead,ct_state=+trk,action=0"
```

Here, the top rule sends to Conntrack setting the `ct_label`, `ct_mark`, and `ct_zone` before recirculating. The second rule should match any recirculating packet as it matches on the same mark, label and zone that have been set by the first rule.

```
ovs-ofctl -O Openflow13 add-flow br0 "tcp,ct_state=-trk,
action=ct(alg=ftp,table=0,commit)"

ovs-ofctl -O Openflow13 add-flow br0 "ct_state=+trk-rel,action=drop"

ovs-ofctl -O Openflow13 add-flow br0 "ct_state=+trk+rel,action=2"
```

In this final example, the first rule sends all matching packets to Conntrack and has them examined for FTP data channel negotiation, then recirculated. The final 2 rules check for any new data channels by looking for the inclusion or emission of a `rel` flag in the packet state.

3.12.7 Viewing Conntrack

Conntrack rules that have been added to the system can be viewed in the same way as standard OvS rule by using the command:

```
ovs-ofctl dump_flows BRIDGE_NAME
```

As Conntrack can not be applied in user-space, all rule matches should happen in the datapath. All datapath rules can then be viewed with:

```
ovs-dpctl dump-flows
```

For debugging purposes it may be useful to remove the timeouts on these rules with:

```
ovs-vsctl set Open_vSwitch . other_config:max-idle=-1
```

On the datapath it should be clear to see the Conntrack action and match fields that have filtered down to the NFP, along with the recirculation instructions and `recirc_id` match and action fields. The statistics on the datapath rules count the number of packets/bytes that match the rule. Because recirculation is active in Conntrack we may see one packet matching more than one rule. Therefore it is likely that the total rule matches will be greater than the number of packets that pass through the system.

4. Use Case Examples

4.1 Tunnel over Aggregated Links

Figure 4.1 demonstrates how to set up a tunnel over aggregated links. The setup consists of two bridges, the first establishes a tunnel (vxlan in this case), which will route packets to the second bridge via the IP stack. The second bridge instantiates a bond as described in section 3.5.

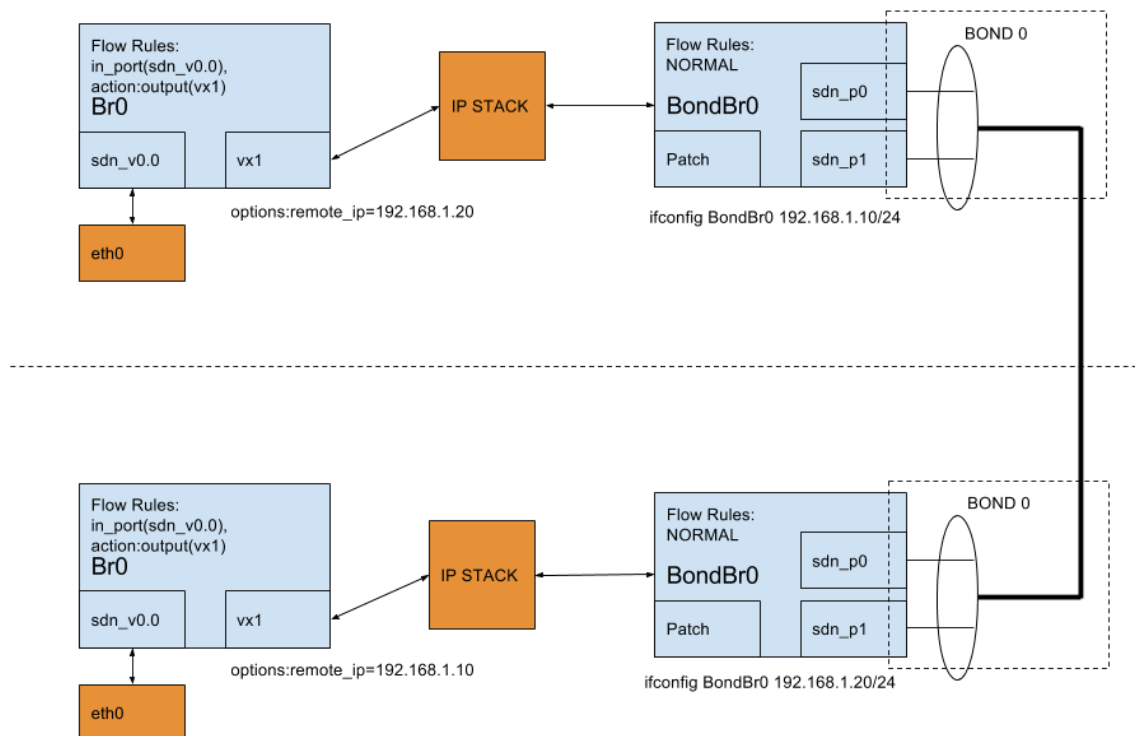


Figure 4.1. Tunnel over Aggregated Links

4.1.1 Configuration Example

The configuration assumes a symmetrical setup. The following commands can be used to configure the tunnel used on the first bridge:

```
Machine 1:
ovs-vsctl add-br br0
ovs-vsctl add-port br0 sdn_v0.0 -- set interface sdn_v0.0 ofport_request=1
ovs-vsctl add-port br0 vx1 -- set interface vx1 type=vxlan options:remote_ip=192.168.1.10 \
ofport_request=2
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1,actions=output:2
```

```
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=2,actions=1
Machine 2:
ovs-vsctl add-br br0
ovs-vsctl add-port br0 sdn_v0.0 -- set interface sdn_v0.0 ofport_request=1
ovs-vsctl add-port br0 vx1 -- set interface vx1 type=vxlan options:remote_ip=192.168.1.20 \
ofport_request=2
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1,actions=output:2
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=2,actions=1
```

The second bridge should be configured as follows:

```
Machine 1:
ovs-vsctl add-br bondbr0
ovs-vsctl add-bond bondbr0 bond0 sdn_p0 sdn_p1
ovs-vsctl set port bond0 bond_mode=balance-tcp
ovs-vsctl set port bond0 lacp=active
ovs-vsctl set port bond0 other_config:lacp-time=fast
Machine 2:
ovs-vsctl add-br bondbr0
ovs-vsctl add-bond bondbr0 bond0 sdn_p0 sdn_p1
ovs-vsctl set port bond0 bond_mode=balance-tcp
ovs-vsctl set port bond0 lacp=active
ovs-vsctl set port bond0 other_config:lacp-time=fast
```

This configuration relies on the fact that entunneled traffic will egress the system on aggregated links to a remote peer. This is achieved by configuring the IP address on the bond bridges:

```
Machine 1:
ifconfig bondbr0 192.168.1.20/24 up
Machine 2:
ifconfig bondbr0 192.168.1.10/24 up
```

The flow rule set on the second bridge should only contain a NORMAL rule. This should be the default flow rule on bridge creation.



Note

The OvS 'no-forward' option has not been applied to this setup as entunneled traffic is expected to ingress on the LOCAL port. The Openflow OvS show command can be used to observe port settings.

```
ovs-ofctl show bondbr0 -OOpenflow13
```

The 'forward' option can be enabled on a port using:

```
ovs-ofctl -OOpenflow13 mod-port bondbr0 bondbr0 forward
```

4.2 Load balancing tunnels

Configuring load balance over tunnels follows a similar approach to that described in section 3.6. The key differences here are setting up tunnel endpoints and configuring groups that could be used in combination with these tunnels.

4.2.1 Configuration Example

This configuration assumes the relevant remote peer IP addresses are setup correctly. The tunnel endpoint configuration follows:

```
ovs-vsctl add-port br0 vx1 -- set Interface vx1 type=vxlan options:key=flow \
options:remote_ip=192.168.1.20 ofport_request=100
ovs-vsctl add-port br0 gre1 -- set Interface gre1 type=gre options:key=flow \
options:remote_ip=192.168.2.20 ofport_request=200
```

Typically when specifying key=flow in the port setup, one would be required to specify the tunnel key using a flow rule. However, when load balancing tunnels the group configuration should contain the tunnel specification:

```
ovs-ofctl -O OpenFlow13 add-group br0 group_id=0,type=select,\
bucket=actions=set_field:204.167.145.118->nw_src,set_field:10->tun_id,output:100,\
bucket=actions=set_field:204.167.145.118->nw_src,set_field:20->tun_id,output:200,\
bucket=actions=set_field:204.167.145.118->nw_src,set_field:30->tun_id,output:100,\
bucket=actions=set_field:204.167.145.118->nw_src,set_field:40->tun_id,output:200
```

All that remain now is to configure a flow rule that would apply the group to traffic:

```
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1,actions=group:0
```

4.3 Fast dataplane processing combined with standard network services

It is a common misconception that a particular server with a single Intelligent Server Adapter must be executing a single network function only. However, in this example it is illustrated how with the use of OpenFlow rules and VF netdevs, traffic can be separated into high speed dataplane processing (via an arbitrary DPDK application), and at the same time running simple network services such as a TCP connection for ssh.

This example can be extended in many different ways, adding multiple DPDK applications, or routing network services traffic to a VM instead of the host.

For Intelligent Server Adapters with multiple physical ports, the example can be modified to run the network services on different physical ports. The advantage of such a configuration above this example, is that those network services will be available even if Agilio OvS is stopped.

4.3.1 Configuration Example

The following procedure details the steps to configure Agilio OvS for this example. The following preconditions are expected to be properly configured before attempting the commands below:

- Agilio OvS is running.

- There is no OvS configuration loaded in Agilio OvS at this time.
- Hugepages are configured for DPDK usage. Refer to Section 3.7 and specifically Section 3.7.2 for information about using DPDK applications with Agilio OvS.
- DPDK application compiled for the target system.

Allocate sdn_v0.0, sdn_v0.1, sdn_v0.2 and sdn_v0.3 for DPDK use by binding these VFs to the nfp_uio driver:

```
/opt/netronome/libexec/dpdk_nic_bind.py -b nfp_uio 02:08.0 02:08.1 02:08.2 02:08.3
```

Allocate sdn_v0.8 and sdn_v0.9 for network services use by binding these VFs to the nfp_netvf driver:

```
/opt/netronome/libexec/dpdk_nic_bind.py -b nfp_netvf 02:09.0 02:09.1
```

Configure the OvS Bridge:

```
ovs-vsctl add-br br0
ovs-vsctl set-fail-mode br0 secure
ovs-vsctl add-port br0 sdn_p0 -- set Interface sdn_p0 ofport_request=1
ovs-vsctl add-port br0 sdn_v0.0 -- set Interface sdn_v0.0 ofport_request=1000
ovs-vsctl add-port br0 sdn_v0.1 -- set Interface sdn_v0.1 ofport_request=1001
ovs-vsctl add-port br0 sdn_v0.2 -- set Interface sdn_v0.2 ofport_request=1002
ovs-vsctl add-port br0 sdn_v0.3 -- set Interface sdn_v0.3 ofport_request=1003
ovs-vsctl add-port br0 sdn_v0.8 -- set Interface sdn_v0.8 ofport_request=1008
ovs-vsctl add-port br0 sdn_v0.9 -- set Interface sdn_v0.9 ofport_request=1009
```

Configure the flow rules for load balancing UDP traffic to VFs 0-3. UDP traffic could be VoIP for example:

```
ovs-ofctl -O OpenFlow13 add-group br0 group_id=0,type=select,bucket=actions=1000,\
    bucket=actions=1001,bucket=actions=1002,bucket=actions=1003
ovs-ofctl -O OpenFlow13 add-flow br0 ip,udp,in_port=1,actions=group:0
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1000,actions=1
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1001,actions=1
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1002,actions=1
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1003,actions=1
```

Configure the flow rules for ssh connections to VF 8:

```
ovs-ofctl -O OpenFlow13 add-flow br0 priority=11,in_port=1,ip,tcp,\
    tcp_dst=22,nw_dst=123.10.0.0/24,actions=1008
ovs-ofctl -O OpenFlow13 add-flow br0 priority=11,in_port=1,arp,nw_dst=123.10.0.0/24,\
    actions=1008
ovs-ofctl -O OpenFlow13 add-flow br0 priority=11,in_port=1,icmp,nw_dst=123.10.0.0/24,\
    actions=1008
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1008,actions=1
```

Configure an IP address on the netdev associated with VF 8. Assuming that the netdev is named eth0 by the system, run the following command:

```
ifconfig eth0 123.10.0.200 up
```

Configure the flow rules for an arbitrary TCP connection on port 1000:

```
ovs-ofctl -O OpenFlow13 add-flow br0 priority=10,in_port=1,ip,tcp,\
    tcp_dst=1000,actions=1009
ovs-ofctl -O OpenFlow13 add-flow br0 priority=10,in_port=1,arp,actions=1009
ovs-ofctl -O OpenFlow13 add-flow br0 priority=10,in_port=1,icmp,actions=1009
ovs-ofctl -O OpenFlow13 add-flow br0 in_port=1009,actions=1
```

Now configure the IP address on the netdev associated with VF 9. Assuming that the netdev is named eth1 by the system, run the following command:

```
ifconfig eth1 192.168.0.1 up
```

Now an example of a DPDK application that can be run for VFs 0-3:

```
/opt/netronome/samples/dpdk/trafgen/build/trafgen -c 0x1e -n 4 \  
-w 02:08.0 -w 02:08.1 -w 02:08.2 -w 02:08.3 -- -p 0xf -T 5
```

Each of the three network connections can be verified with the nc Linux utility. Assuming traffic is sent on a peer netdev named eth8, configure the following IP address aliases on the netdev:

- ifconfig eth8:0 10.0.0.10
- ifconfig eth8:1 123.10.0.10
- ifconfig eth8:2 192.168.0.10

The IP address aliases are only required to configure multiple routes to the same destination and keep the example simple. Then execute the following commands on the peer host to test these network connections:

- UDP traffic: echo "RANDOM UDP DATA" | nc -u 10.0.0.1
- TCP port 22 traffic: echo "RANDOM SSH DATA" | nc 123.10.0.1 22
- TCP port 1000 traffic: echo "RANDOM TCP1000 DATA" | nc 192.168.0.1 1000

Traffic can be verified on the Agilio OvS host by executing the commands nc -l 22 and nc -l 1000 and observing the DPDK application statistics.

5. Technical Support

To obtain additional information, or to provide feedback, please email <support@netronome.com> or contact the nearest **Netronome** technical support representative.

5.1 Related Documents

Descriptive Name	Description
Netronome Network Flow Processor: Agilio OvS 2.2 Getting Started Guide	A guide to new users of Netronome's Agilio OvS software.
Netronome Network Flow Processor: Agilio OvS 2.2 Programmer's Reference Manual	A reference for programming and system design using the Agilio OvS software.
Intelligent Server Adapters: Hardware User Manual	Contains summary information on the Netronome Intelligent Server Adapter (ISA) PCIe card including card physical descriptions.
Netronome Network Flow Processor: Datasheet	Contains summary information on the Netronome Network Flow Processor NFP including a functional description, signal descriptions, electrical specifications, and mechanical specifications.
Netronome Network Flow Processor: Databook	Contains detailed reference information on the Netronome Network Flow Processor NFP.
Netronome Network Flow Processor: Development Tools User's Guide	Describes Programmer Studio and the development tools that can be accessed through Programmer Studio.
Netronome Network Flow Processor: Network Flow Assembler System User's Guide	Describes the syntax of the NFP's assembly language, supplies assembler usage information, and lists assembler warnings and errors.
Netronome Network Flow Processor: Microengine Programmer's Reference Manual	A reference for microcode programming of the Netronome Network Flow Processor NFP.
Netronome Network Flow Processor: Network Flow C Compiler User's Guide	Presents information, language structures and extensions to the language specific to the Netronome Network Flow C Compiler for Netronome NFP.
Netronome Network Flow Compiler LibC: Reference Manual	Specifies the subset and the extensions to the language that support the unique features of the Netronome Network Flow Processor NFP product line.
Open vSwitch Software Documentation	Agilio OvS software offers acceleration of Open vSwitch software. Refer to http://openvswitch.org/ for more details on Open vSwitch.
OpenFlow Specification	Open vSwitch (which is accelerated by Agilio OvS software) is an OpenFlow switch implementation. Refer to https://www.opennetworking.org/sdn-

Descriptive Name	Description
	resources/openflow for more details on this specification.
Data Plane Development Kit Documentation	DPDK related documentation is available at http://dpdk.org .