# Join us for KubeCon + CloudNativeCon Virtual



**Event dates:** August 17-20, 2020
**Schedule:** Now available!
**Cost:** $75

**Register now!**

# Securing and Accelerating Kubernetes CNI

## Integrating Project Antrea and NVIDIA Mellanox ConnectX SmartNICS

**Antonin Bas**
Maintainer of Project Antrea and Staff Engineer at VMware

**Moshe Levi**
Sr. Staff Engineer at NVIDIA

July 14, 2020

# Antonin Bas

**Maintainer of Project Antrea and Staff Engineer at VMware**
abas@vmware.com

# Moshe Levi

**Sr. Staff Engineer at NVIDIA**
moshele@nvidia.com

# Cody McCain

**Sr. Product Manager Container Networking at VMware**
mmccain@vmware.com

# Itay Ozery

**Director, Product Management for Networking at NVIDIA**
itayo@nvidia.com

# Agenda

Securing and Accelerating the Kubernetes CNI Data Plane

Kubernetes Cluster Networking

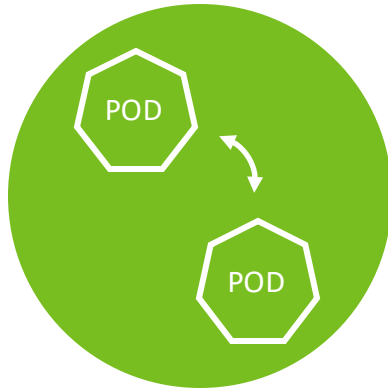Project Antrea Deep Dive

Hardware Acceleration
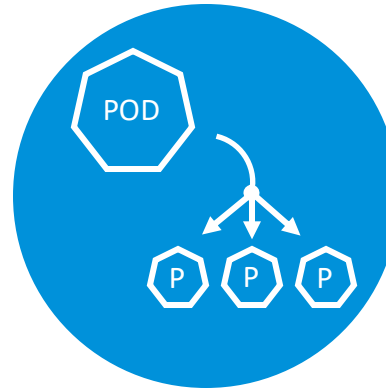
Demo

Roadmap

Get Involved

# Kubernetes Cluster Networking
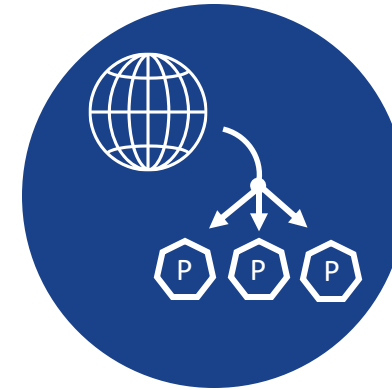
# Kubernetes Cluster Networking

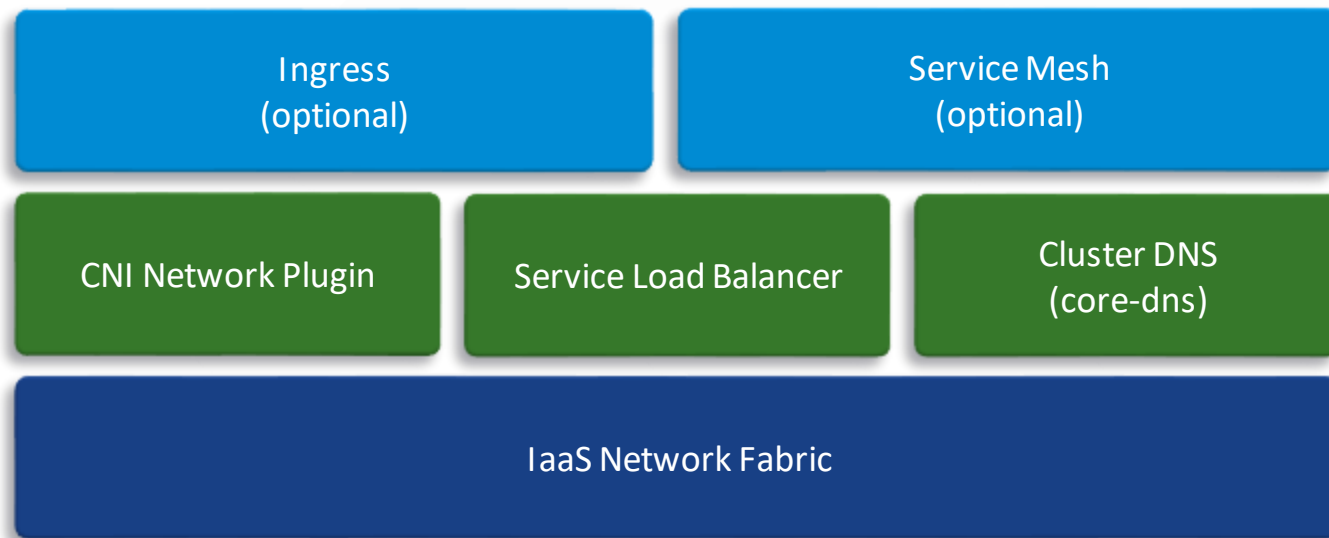Three connectivity scenarios must be enabled.



Pod
-to-
Pod

Pod
-to-
Service

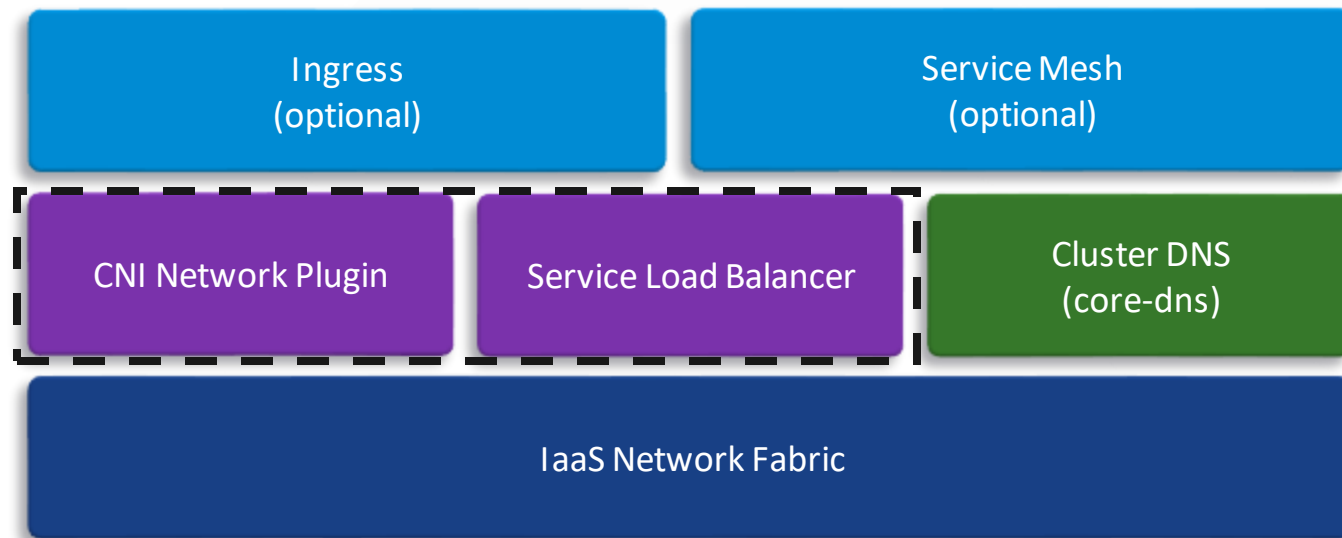External
-to-
Service

# Kubernetes Networking
## in Layers

| Ingress (optional) | | Service Mesh (optional) |
|---|---|---|
| CNI Network Plugin | Service Load Balancer | Cluster DNS (core-dns) |
| IaaS Network Fabric | | |

# Kubernetes Networking
# in Layers

| Ingress (optional) | Service Mesh (optional) |
|---|---|

| CNI Network Plugin | Service Load Balancer | Cluster DNS (core-dns) |
|---|---|---|

| IaaS Network Fabric |
|---|

## What is a Kubernetes CNI Network Plugin

responsible for?

### Pod Connectivity

Plumbing eth0 (network interface) into Pod network (encapsulated or non-encapsulated)
Pod egress to world – SNAT

### IP Address Management (IPAM)

### Service Load Balancing

Make traffic available to upstream kube-proxy, or
Implement native service load balancing – VIP DNAT

### NetworkPolicy Enforcement (optional)

Enforcing Kubernetes Network Policy
Source Spoof Prevention
Connection Tracking (Stateful Firewall)

### hostPort Support

### Traffic Shaping Support

(experimental)

# Project Antrea
# Deep Dive

# Where can I run Antrea?
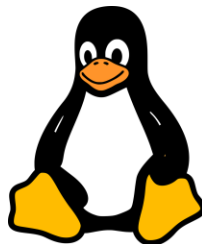
Our goal is to run anywhere Kubernetes runs.

Private Cloud

Public Cloud

Edge

Linux

Windows

# What is Open vSwitch (OVS)?

And why use it for K8s networking?

A high-performance programmable virtual switch
- Connects to VMs (tap) and containers (veth)

**Linux foundation project**, very active

**Portable**: Works out of the box on all Linux distributions and supports Windows
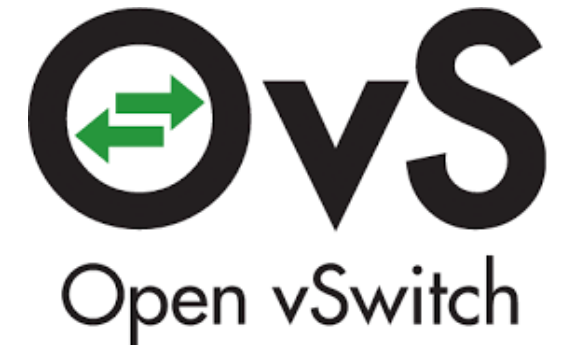
**Programmability**: Supports many protocols, build your own forwarding pipeline

**High-performance**
- DPDK, AF_XDP
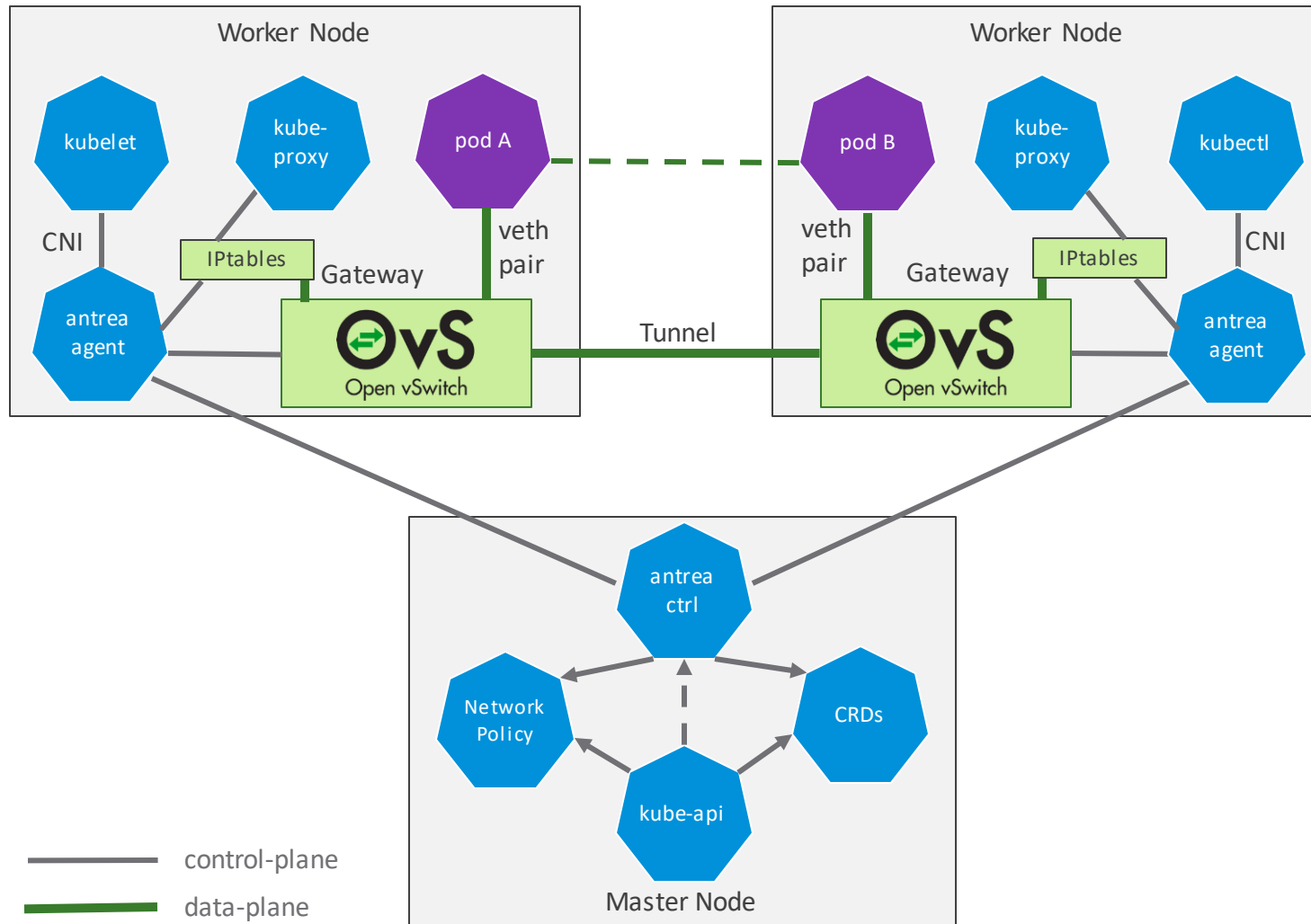- Hardware offload available across multiple vendors

**Rich feature set**:
- Advanced CLI tools
- Statistics, QoS
- Packet tracing

# Project Antrea Architecture

Open vSwitch provides a flexible and performant data plane.



**Supports K8S cluster networking**

**Antrea Agent**

- Manages Pod network interfaces and OVS bridge.
- Creates overlay tunnels across Nodes.
- Implements NetworkPolicies with OVS.

**Antrea Controller**

- Computes K8s NetworkPolicies and publishes the results to Antrea Agents.
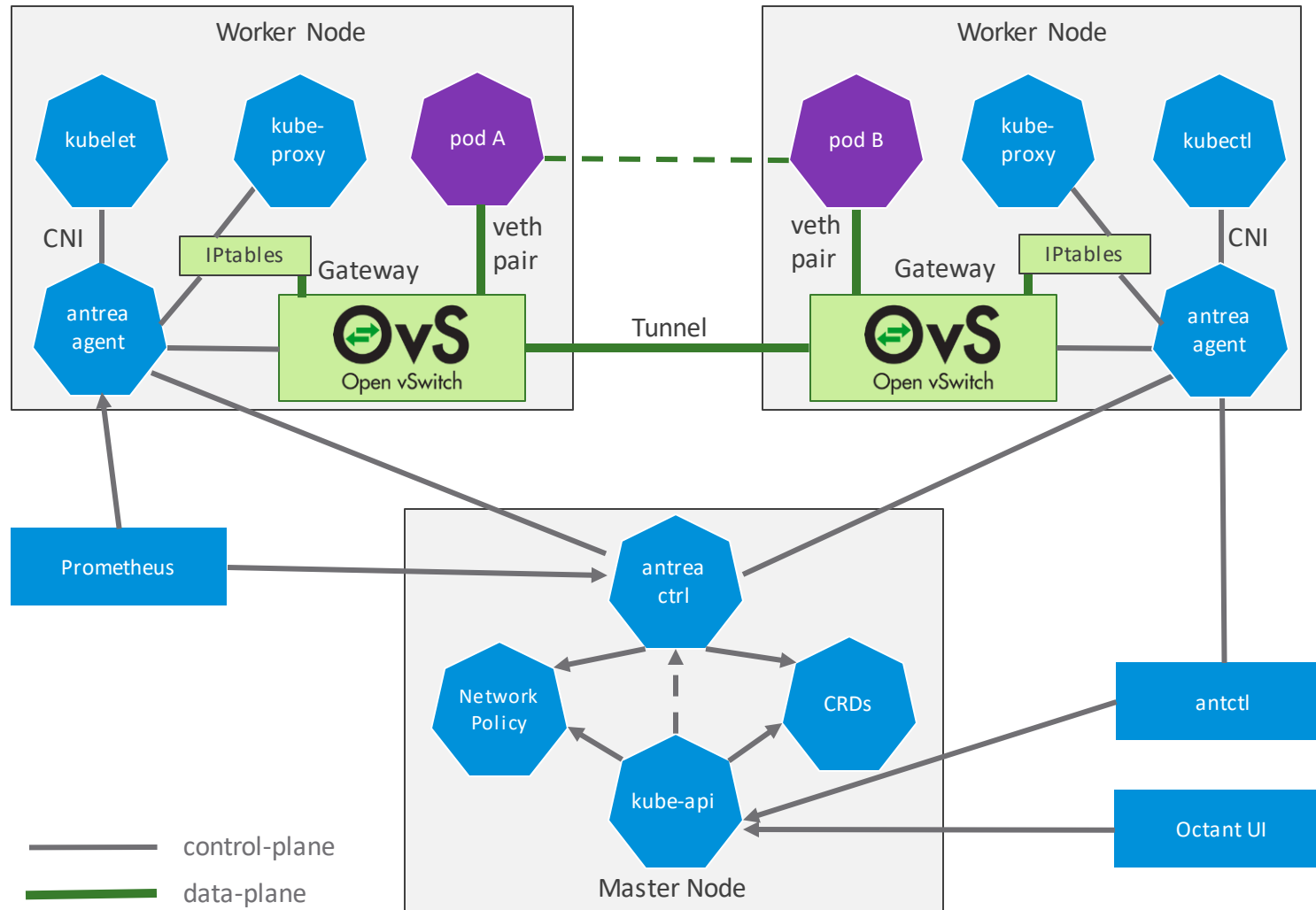
**Open vSwitch as dataplane**

- Antrea Agent programs Open vSwitch with OpenFlow flows.
- Geneve, VXLAN, GRE, or STT tunnel between nodes
- Also supports policy-only and no-encap modes

**Built with K8S technologies**

- Leverages K8S and K8S solutions for API, UI, deployment, control plane, and CLI.
- Antrea Controller and Agent are based on K8S controller and apiserver libs.

# Project Antrea Architecture
## Component Review



**Octant UI Plugin**

- Shows Antrea runtime information (CRDs).
- Diagnostic traceflow visualization.

**antctl – CLI for debugging**

- Connects to Agent Agent or Controller.
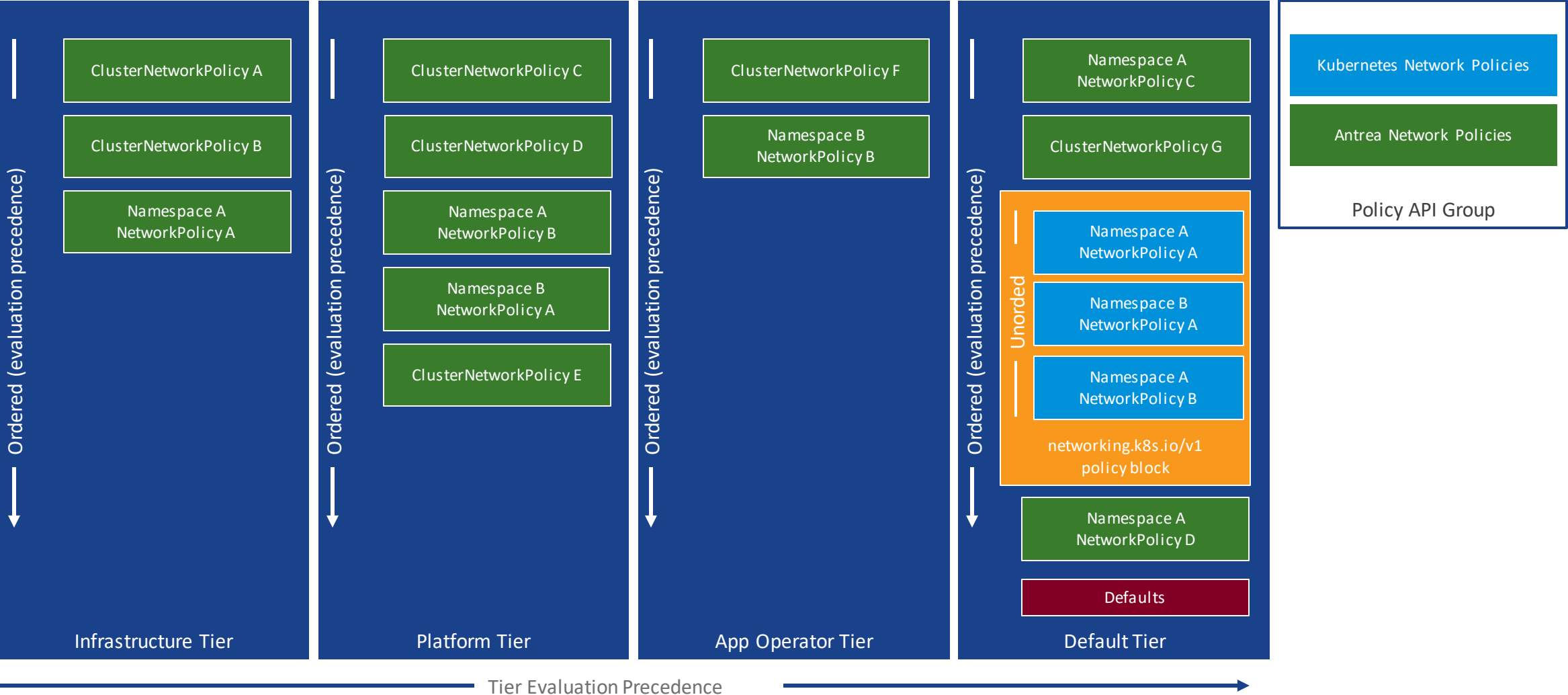- Packet tracing / Support bundle / etc.

**Prometheus metrics available**

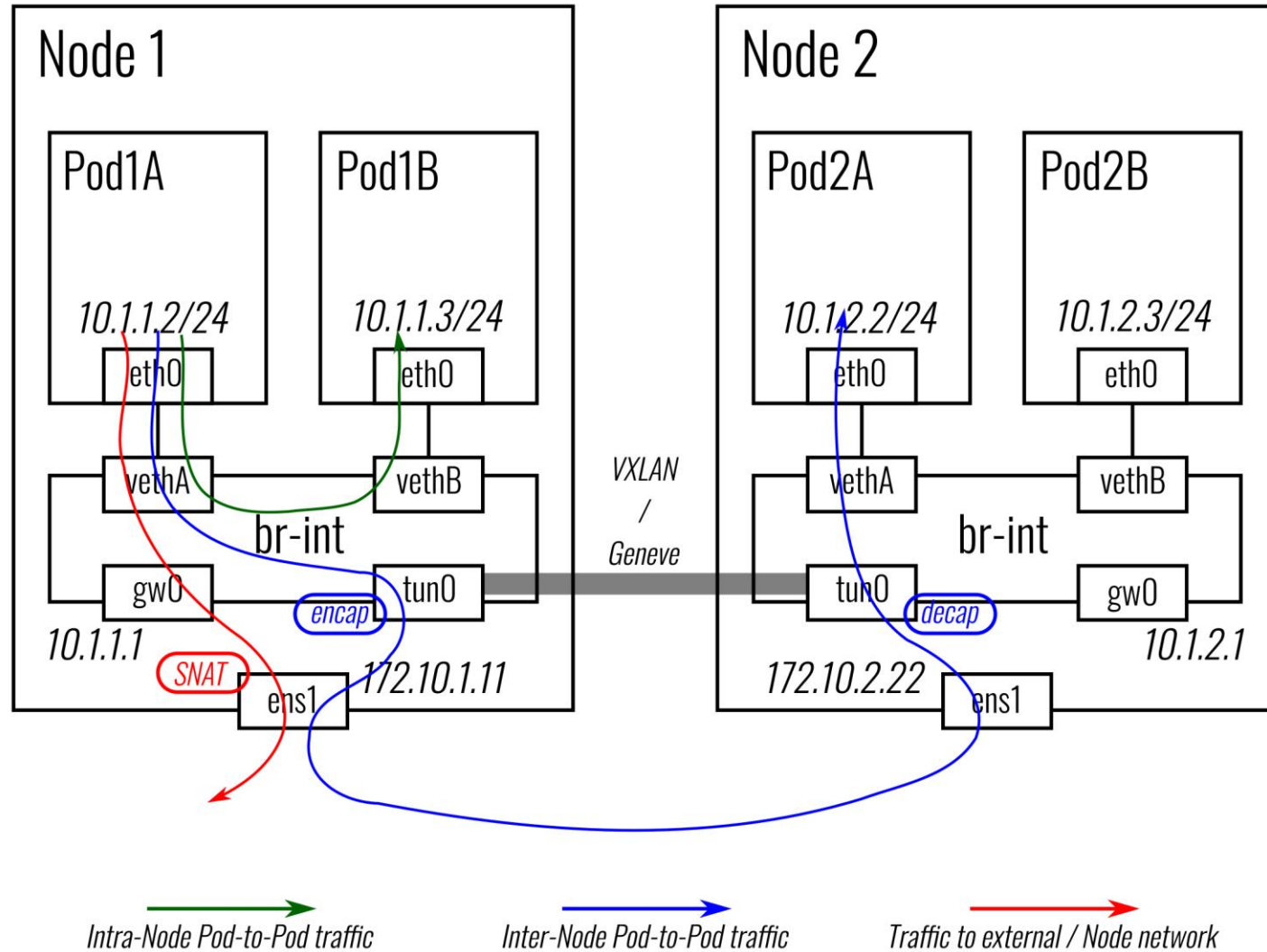**All bits (including OVS daemons) in a Docker image.**

**All components are deployed using K8S manifests.**

# Policy Model

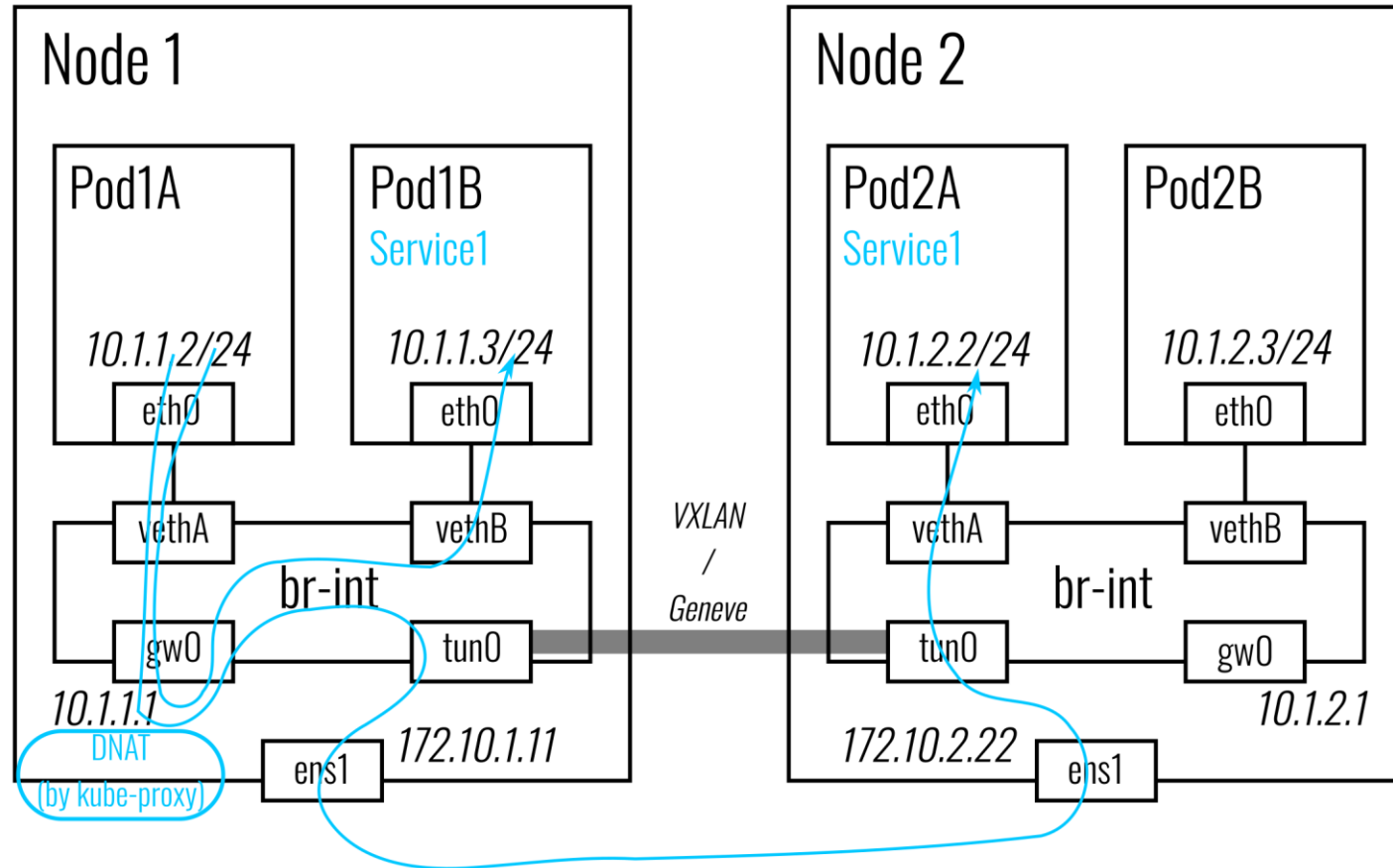Antrea will allow native and Kubernetes policies to co-exist.

**Infrastructure Tier** — Ordered (evaluation precedence)
- ClusterNetworkPolicy A
- ClusterNetworkPolicy B
- Namespace A NetworkPolicy A

**Platform Tier** — Ordered (evaluation precedence)
- ClusterNetworkPolicy C
- ClusterNetworkPolicy D
- Namespace A NetworkPolicy B
- Namespace B NetworkPolicy A
- ClusterNetworkPolicy E

**App Operator Tier** — Ordered (evaluation precedence)
- ClusterNetworkPolicy F
- Namespace B NetworkPolicy B

**Default Tier** — Ordered (evaluation precedence)
- Namespace A NetworkPolicy C
- ClusterNetworkPolicy G
- **Unordered — networking.k8s.io/v1 policy block**
  - Namespace A NetworkPolicy A
  - Namespace B NetworkPolicy A
  - Namespace A NetworkPolicy B
- Namespace A NetworkPolicy D
- Defaults

**Policy API Group**
- Kubernetes Network Policies
- Antrea Network Policies

Tier Evaluation Precedence

# Traffic Walk (in "encap" mode)



Intra-Node Pod-to-Pod traffic  ·  Inter-Node Pod-to-Pod traffic  ·  Traffic to external / Node network
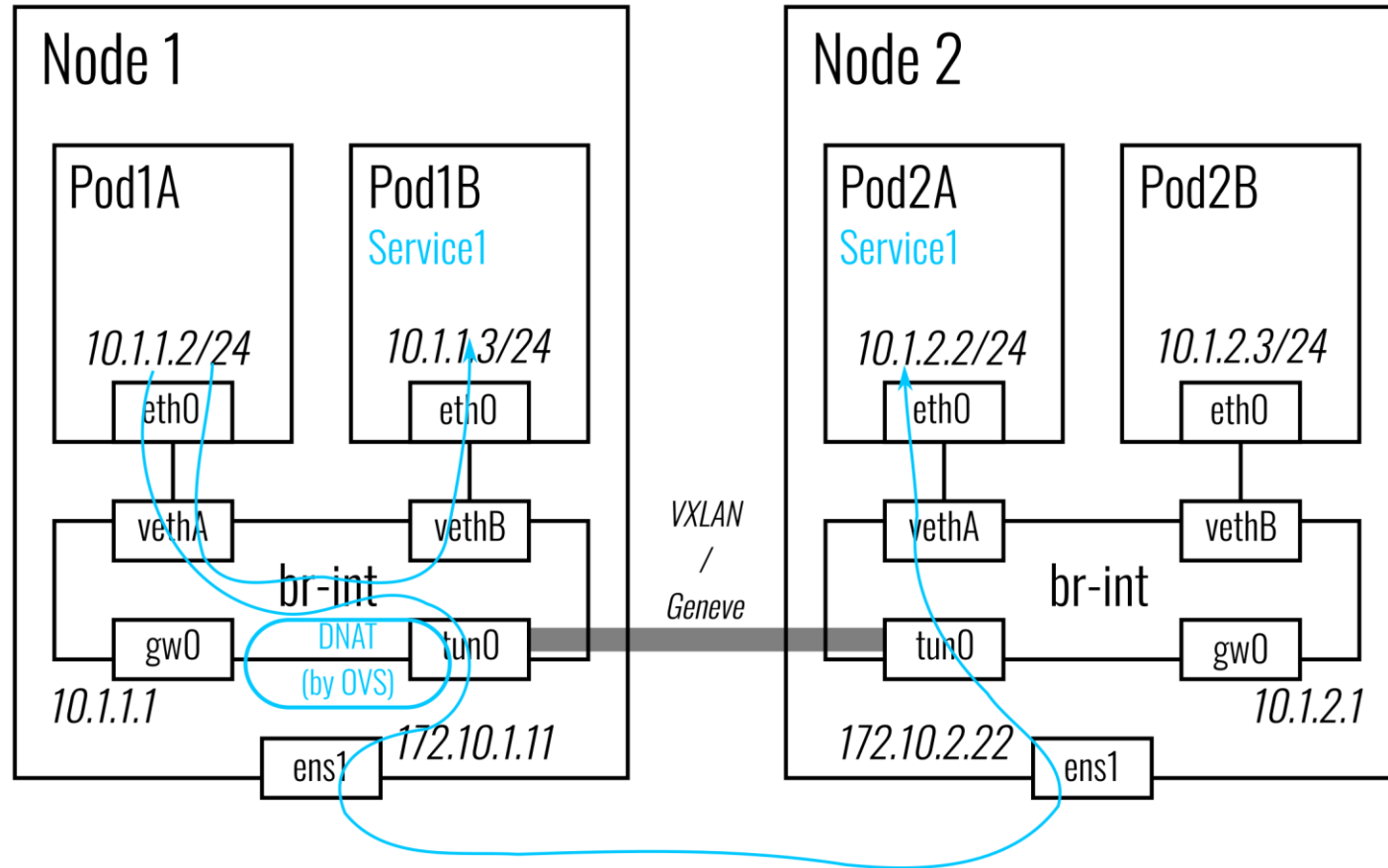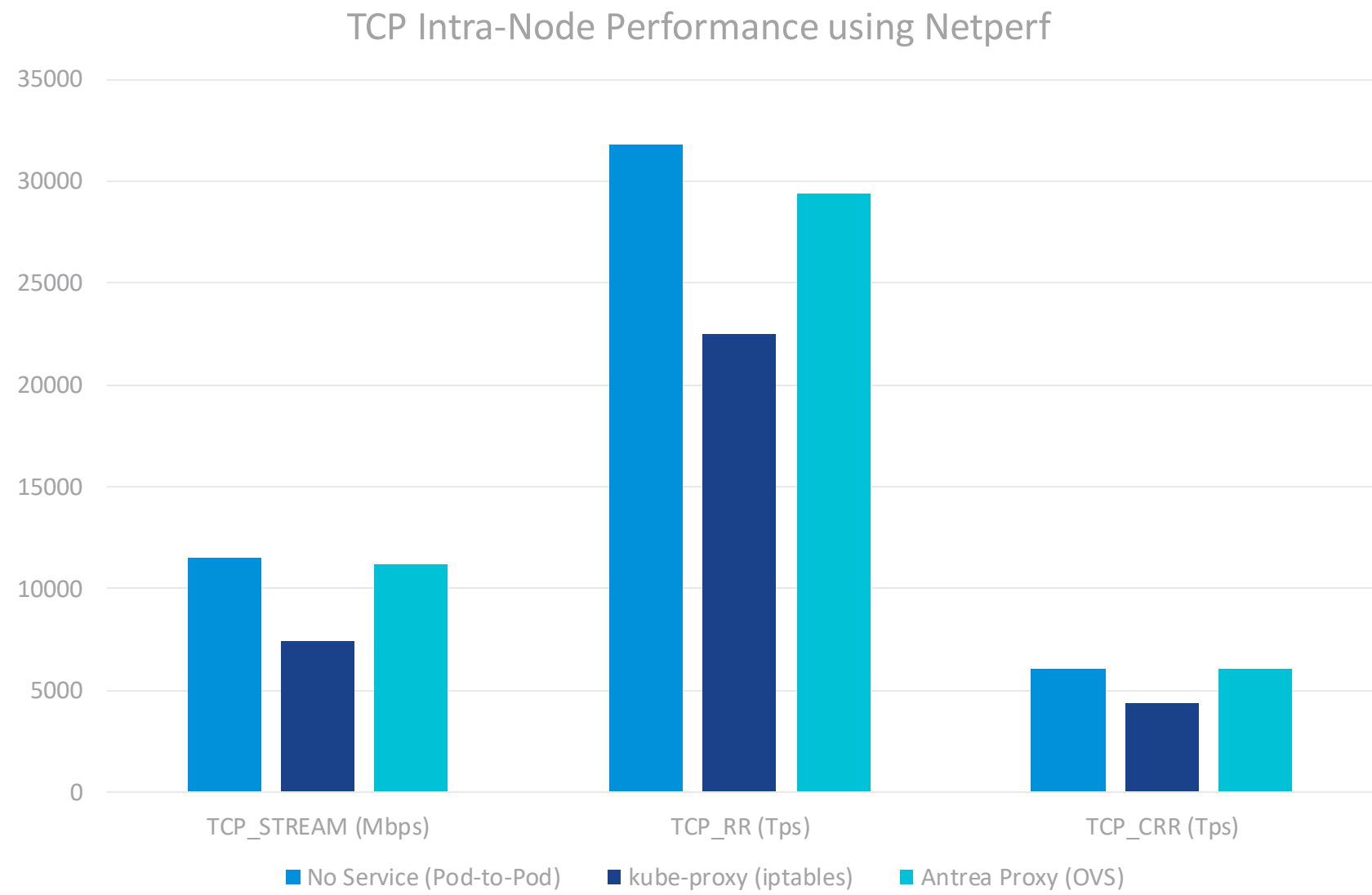
# Traffic Walk: ClusterIP Services

Delegating to kube-proxy

# Traffic Walk: ClusterIP Services in OVS

New in v0.8.0: ClusterIP without kube-proxy

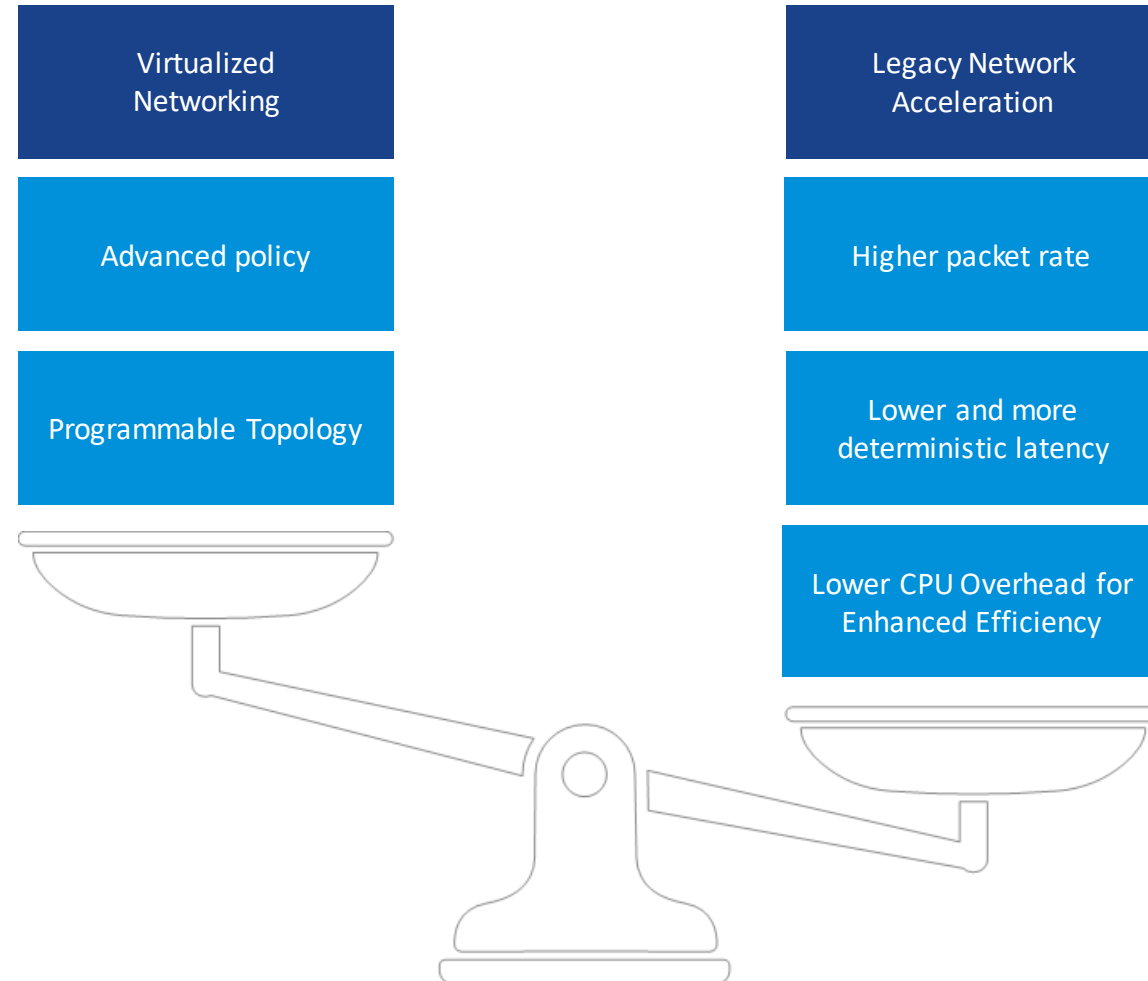# ClusterIP Services in OVS
## "Antrea Proxy"

**TCP Intra-Node Performance using Netperf**



Legend: No Service (Pod-to-Pod), kube-proxy (iptables), Antrea Proxy (OVS)

Categories: TCP_STREAM (Mbps), TCP_RR (Tps), TCP_CRR (Tps)

# Hardware Acceleration

# No Tradeoff between Virtualized and Accelerated Networking
## Decision used to be Either/Or

| Virtualized Networking | Legacy Network Acceleration |
|---|---|
| Advanced policy | Higher packet rate |
| Programmable Topology | Lower and more deterministic latency |
| | Lower CPU Overhead for Enhanced Efficiency |

# Introducing OVS Hardware Offload

## Now we can have Both/And

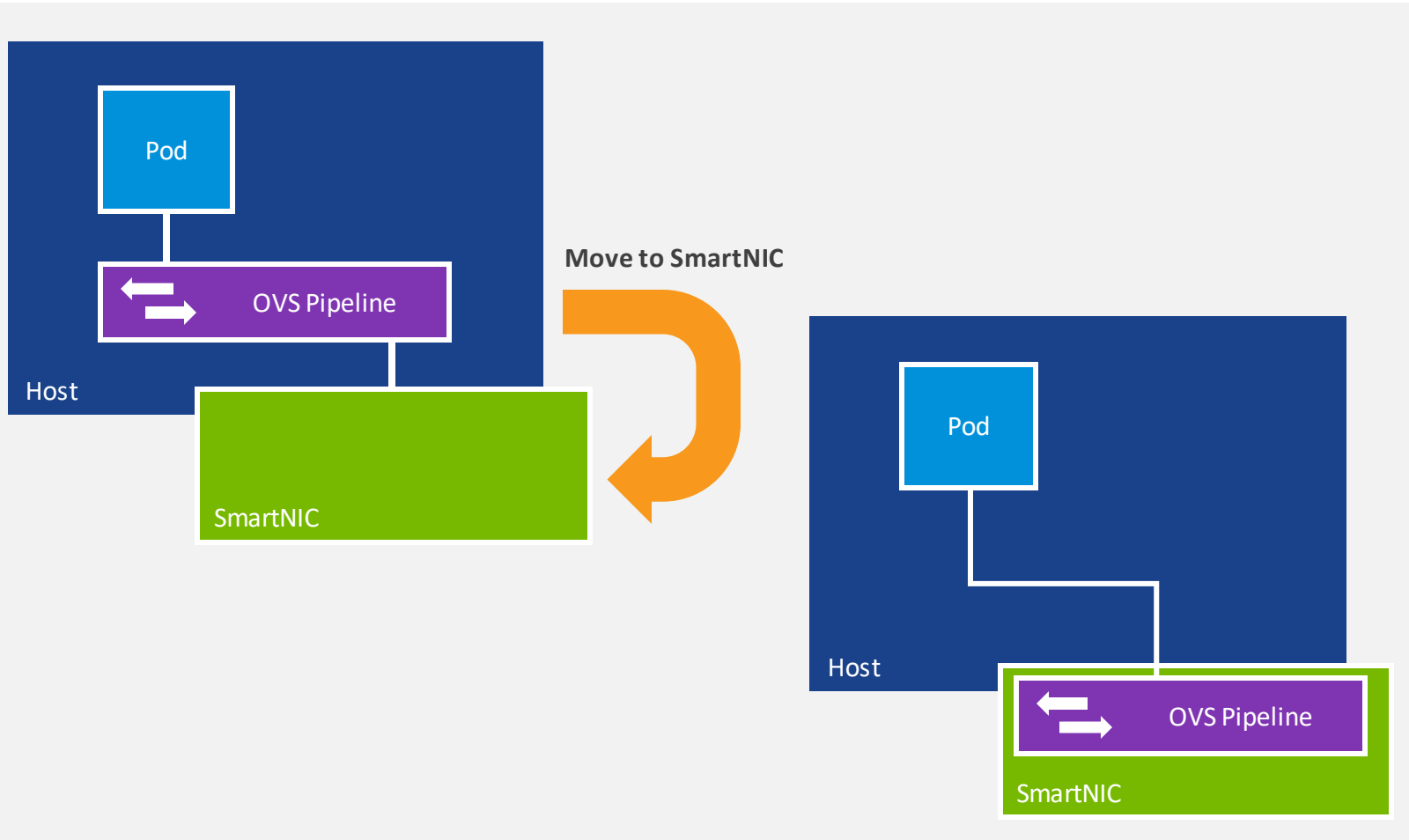| Virtual Switch Control Plane | + | Hardware Accelerated Data Plane | + | Standard Hardware Abstraction Interface | = | OVS Hardware Offload |
|---|---|---|---|---|---|---|

✓ Best of both worlds: Enable hardware-accelerated networking data plane with programmable control plane

✓ Up to 10X network performance with practically zero CPU utilization
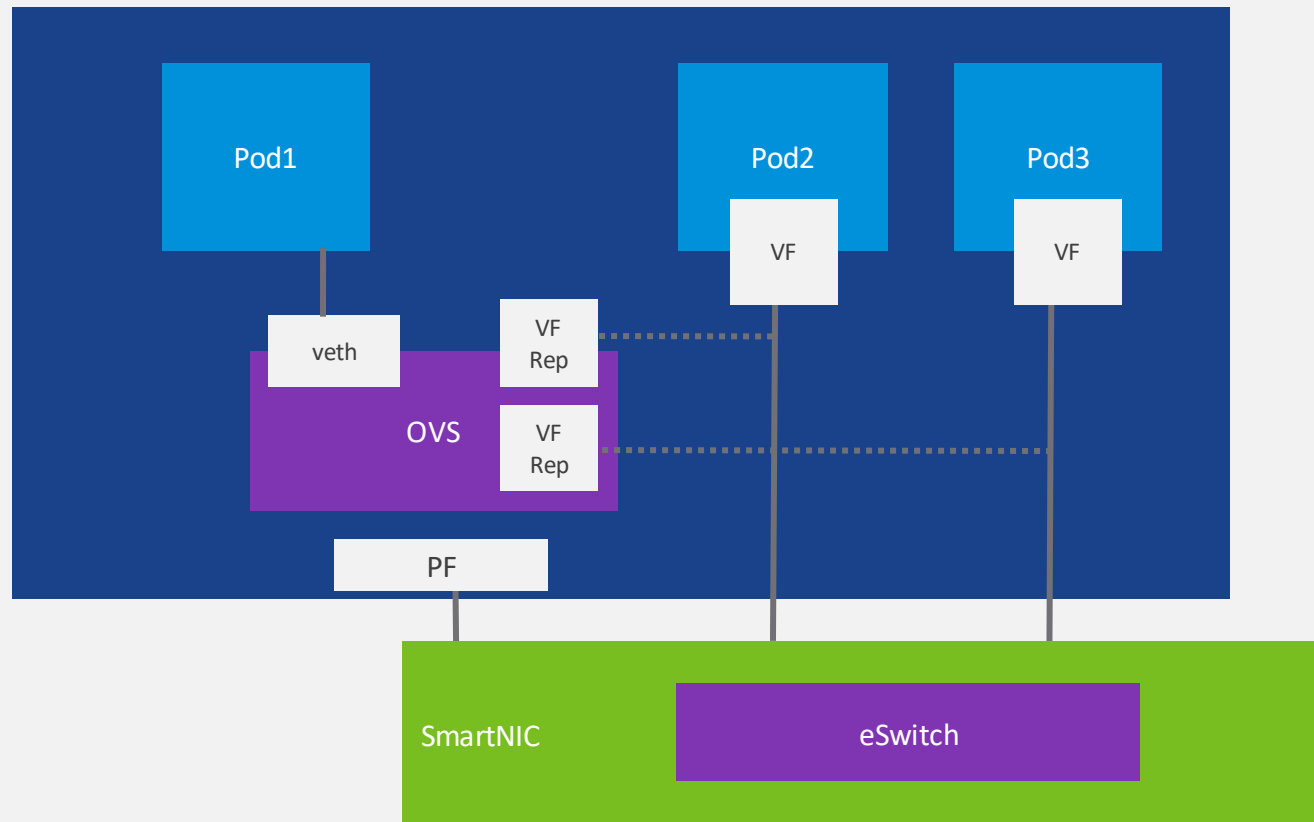
# OVS Hardware Offload

## Move OVS OpenFlow Processing to a SmartNIC

**Move to SmartNIC**

Pod

OVS Pipeline

Host

SmartNIC

Pod

Host

OVS Pipeline

SmartNIC

Typically, OVS flows are processed on a bare metal host, VM or hypervisor.

- The OVS kernel or user space component consumes CPU
- Less CPU resources available for apps
- Moving OVS processing to the SmartNIC frees up CPU

# SR-IOV Definitions



**SR-IOV** – Single Root
I/O Virtualization

**PF** – Physical Function. The
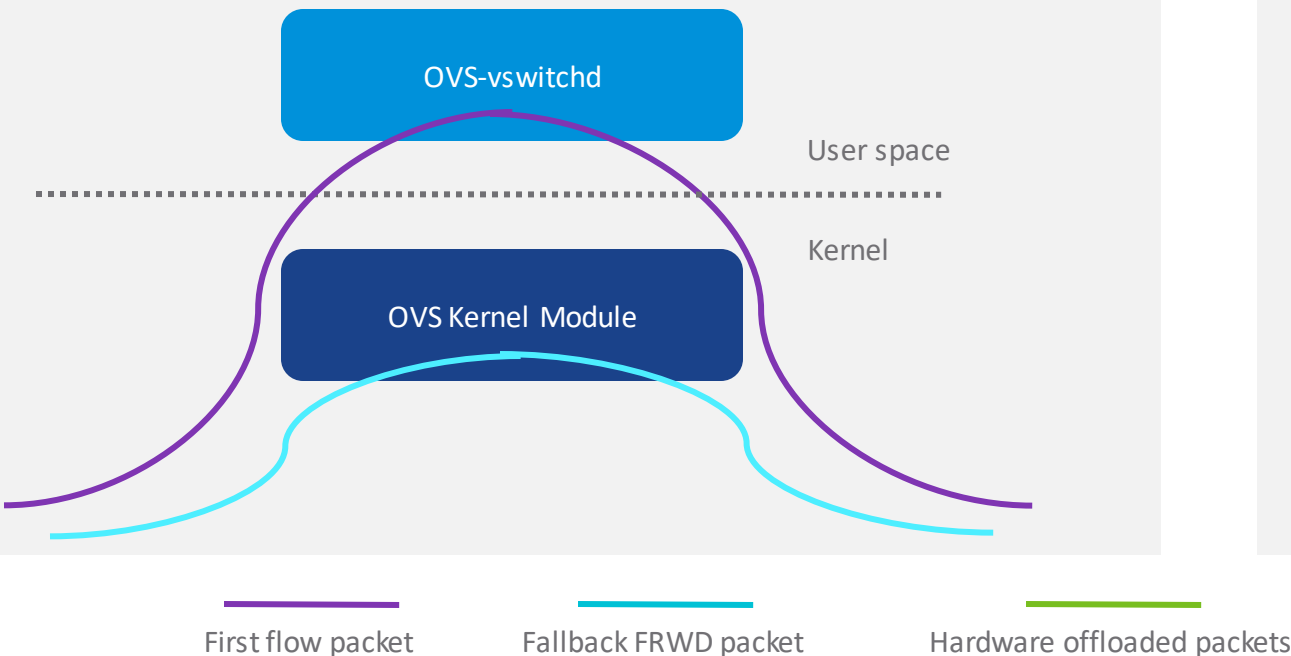physical Ethernet controller
that supports SR-IOV.

**VF** – Virtual Function. The
virtual PCIe device created
from a physical
Ethernet controller.

**VF Representor** – Port
representor of
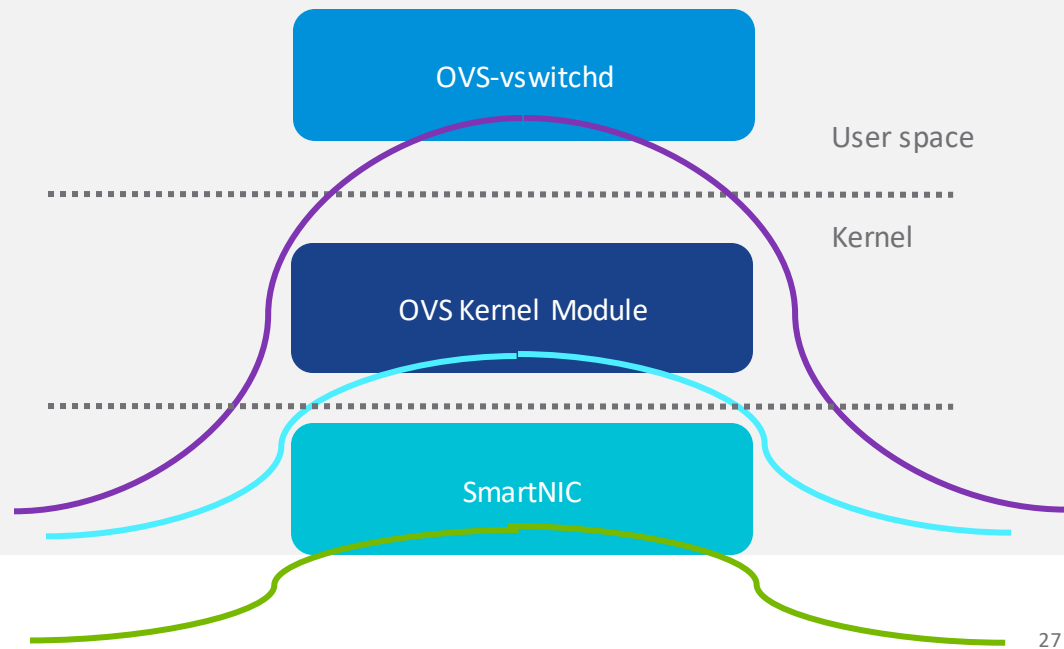the Virtual Function

# How OVS Hardware Offload Works

## Software only OVS Implementation

High latency, low bandwidth, CPU intensive

OVS-vswitchd

User space

Kernel

OVS Kernel Module

## Software-defined, Hardware-accelerated

Low latency, high bandwidth, CPU efficient

OVS-vswitchd

User space

Kernel

OVS Kernel Module

SmartNIC

First flow packet          Fallback FRWD packet          Hardware offloaded packets
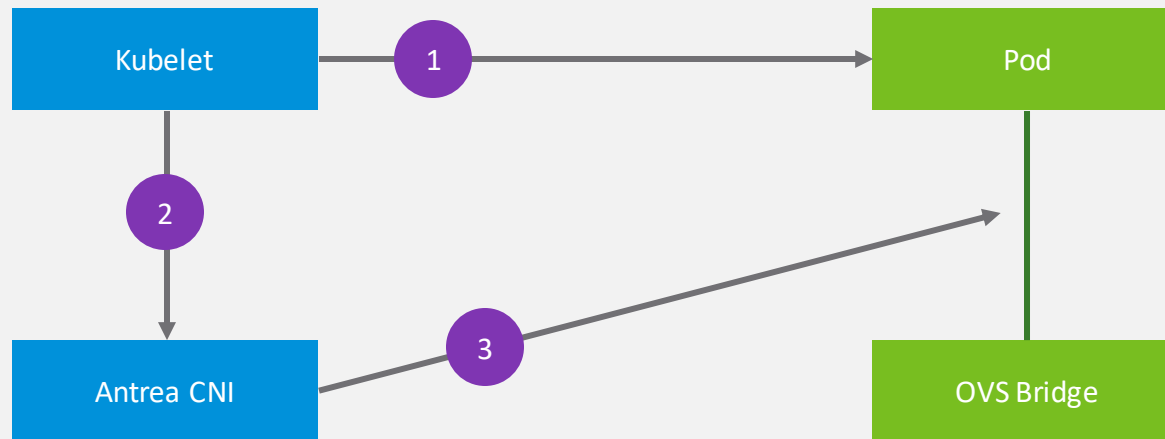
27

# OVS Hardware Offload

Requires additional CNI plugins and SR-IOV VF enablement on NIC

- Multus

- SR-IOV Network Device Plugin

- Antrea

# Antrea CNI Plumbing Without Offload

**Control Plane**

**Data Plane**

Kubelet —1→ Pod

Kubelet —2→ Antrea CNI

Antrea CNI —3→ OVS Bridge

Pod — OVS Bridge

1. Kubelet creates pod

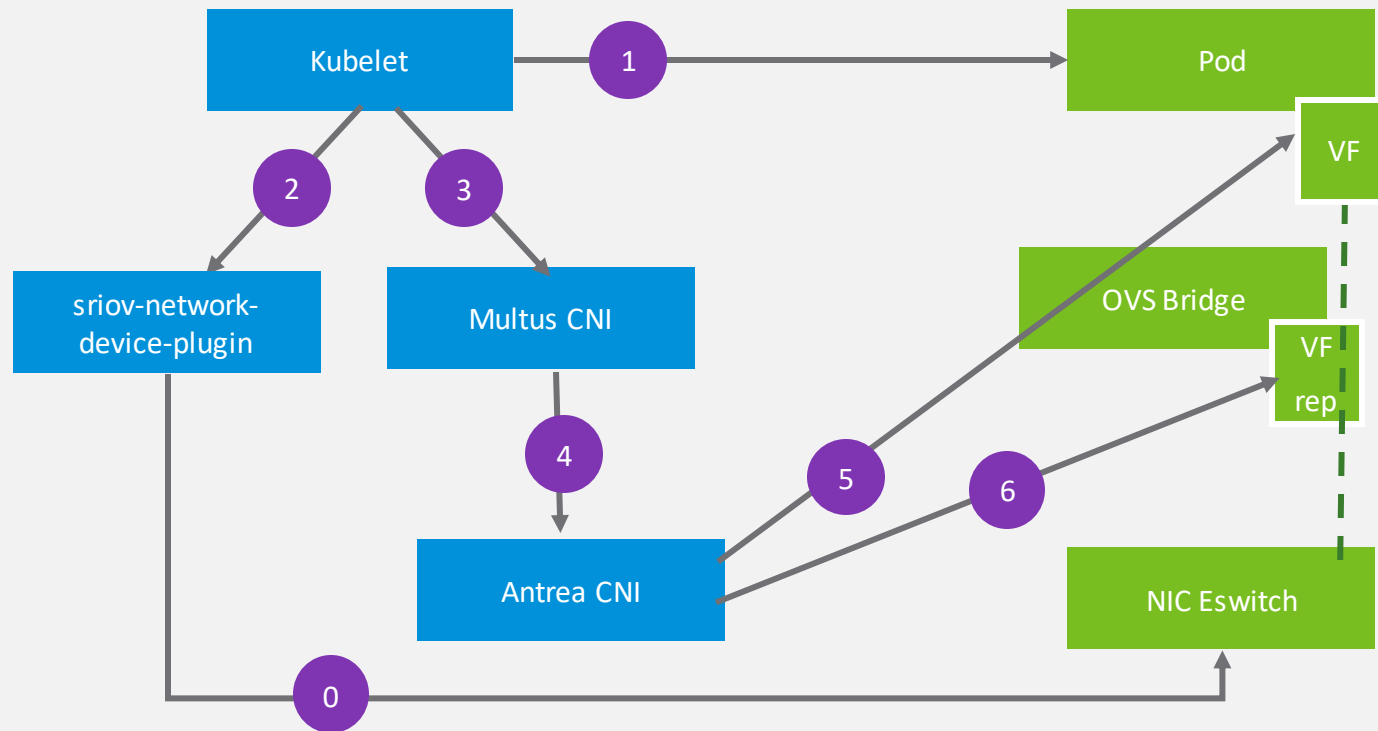2. Kubelet calls CNI to add pod to network

3. Antrea CNI provisions veth pair
   - eth0 in pod network namespace
   - connect other end to OVS bridge port

# Antrea CNI Plumbing With Offload

**Control Plane**

**Data Plane**

Kubelet

1 → Pod

VF

2

3

sriov-network-device-plugin

Multus CNI

OVS Bridge

VF rep

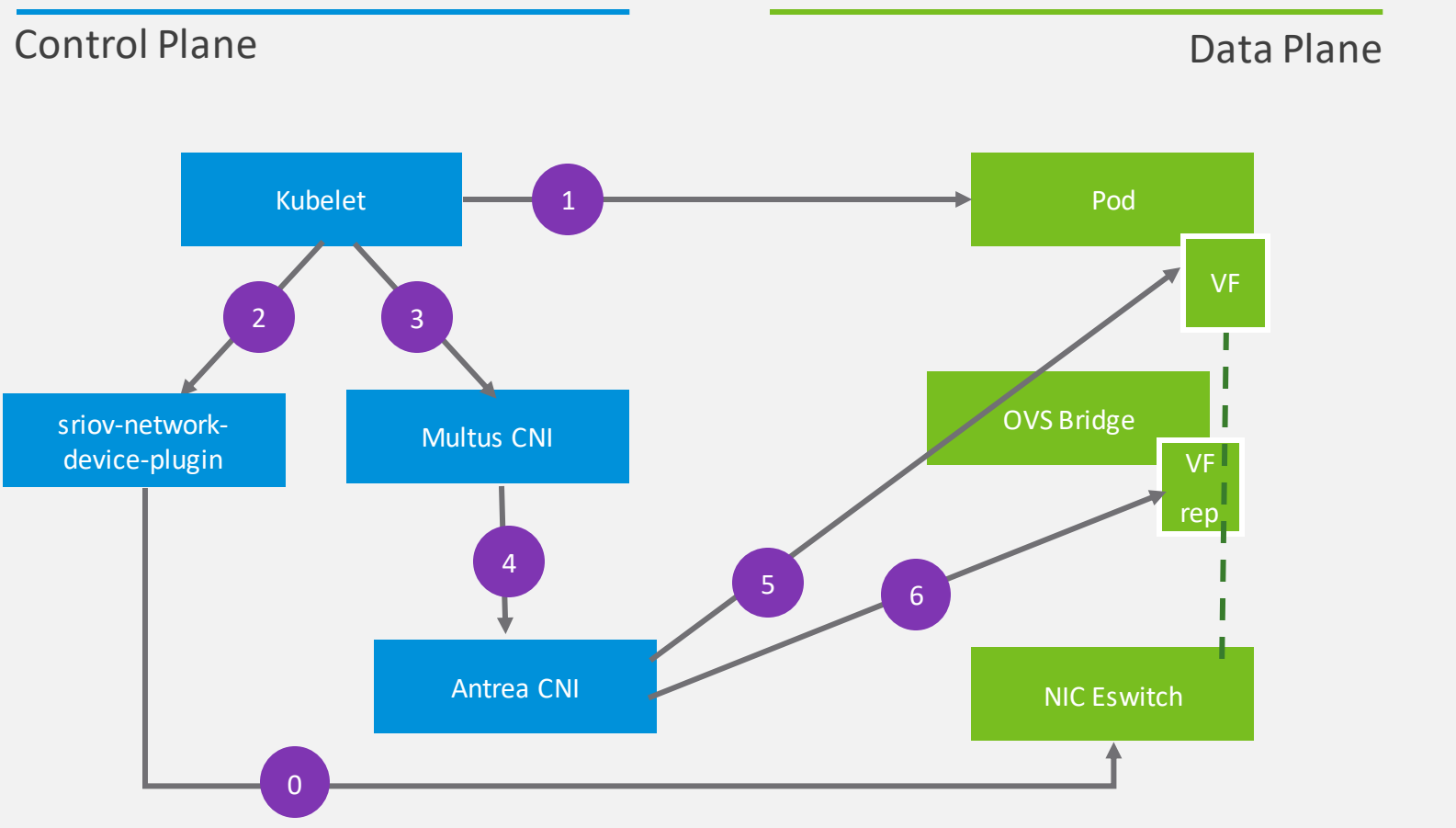4

5

6

Antrea CNI

NIC Eswitch

0

0. VF pool initialization

1. Kubelet creates pod

2. SR-IOV Device Plugin allocates VF PCI address from VF pool to satisfy resource request on pod creation (exposed as environment variable)

# Antrea CNI Plumbing With Offload



**Control Plane**

**Data Plane**

Kubelet — 1 → Pod

Kubelet — 2 → sriov-network-device-plugin

Kubelet — 3 → Multus CNI

Multus CNI — 4 → Antrea CNI

Antrea CNI — 5 → VF

Antrea CNI — 6 → VF rep

OVS Bridge

NIC Eswitch

VF

VF rep

0

3. Kubelet calls CNI (Multus) to add pod to network

4. Multus CNI looks up the allocated SR-IOV VF PCI Address and passes it as extra CNI args to Antrea CNI

5. Antrea CNI moves the VF netdevice to the pod network namespace and renames to eth0

6. Antrea CNI plugs the VF representor intoto the OVS br-int bridge

# Demo - Setup Details

- 3 servers – 1 master and 2 workers

- Linux CentOS 7.7

- Kubernetes 1.18

- Linux 5.7 kernel

- Antrea v0.8.0 with offload patches

- NVIDIA Mellanox ConnectX-5 SmartNICs

# Demo – Flow

- Deploy SR-IOV network device plugin

- Deploy Multus CNI

- Deploy Antrea

- Create veth Pod

- Create offload Pod

- Run iperf3 between 2 veth pods

- Run iperf3 between 2 offload pods

# Demo

# Antrea Roadmap

# Features Available Through v0.8.0

## Overlay Modes

Geneve, VXLAN,
STT, GRE

Policy-only
(CNI chaining)

No-encap

Hybrid

## Clouds

Private Cloud:
bare metal, vSphere, other
VM, kind

Public Cloud:
Azure – AKS Engine
AWS – EC2, EKS (beta)
Google – GKE (alpha)

## Service Load Balancing

kube-proxy support in IPVS
and IPtables modes

OVS based kube-proxy
implementation

# Features Available Through v0.8.0

## Network Policy

networking.k8s.io
NetworkPolicy v1
(upstream)

Native Policy:
ClusterNetworkPolicy

## Security

Server certificate verification
for Controller APIs (user
provided or generated)

Spoof Guard

IPsec over GRE

## Visibility

Prometheus Metrics
& Monitoring CRDs

Traceflow

Support bundle
generation

antctl CLI &
Octant UI Plugin

# Traceflow

## Request 1: traffic is allowed

# Traceflow

Request 2: traffic is denied

# Features Available Through v0.8.0

Operating Systems

---

Linux

Windows Server 2019 (alpha)

# Planned Features This Year

IPFIX flow data export

Advanced traffic matching and pod binding

Tiering to support multi-tenancy and delegation.

IPv6 dual-stack support

IPsec Offload

Expand support for KaaS and Cluster API providers

Enhanced data path including:
DPDK, SR-IOV, AF_XDP, VPP, and XDP

DNS egress filtering

Advanced IP Address Management

Named external endpoints with metadata

Extension mechanisms

# Flow information export and visualization

Track all cluster traffic
- Number of connections
- Bandwidth for each connection
- Inter-Node bandwidth
- Aggregated Service bandwidth

Complements Prometheus metrics

IPFIX records with K8s context (Namespace, Name, Labels, …)

Visualization using Elastic Stack

# Flow information export

## IPFIX Records

# Flow information visualization

## With Elastic Stack

# Get Involved

![ANTREA logo]

# Come help us continually improve Kubernetes Networking!

**@**
projectantrea-announce
projectantrea
projectantrea-dev
(Google Groups)

@ProjectAntrea

Kubernetes Slack
#antrea

Community Meeting, Mondays @ 9PM PT
Zoom ID: 823-654-111

https://github.com/vmware-tanzu/antrea

- Good first issues
- Help us improve our documentation
- Propose new features
- File Bugs

https://antrea.io

- Documentation
- Blogs

# Thank You

# Backup Slides

# Antrea in Public Cloud

The Antrea CNI provides both pod connectivity and network policy enforcement and is flexible to use in either cloud native or overlay IP addressing schemes.



**Native CNI**

IP addresses are assigned from cloud native private network (VPC)

**ANTREA**

pods are assigned addresses from IP space opaque to the cloud

**ANTREA**

Enforces network policy and filters traffic to/from pods

**ANTREA**

can provide overlay encapsulation and encryption when connecting pods

**Native CNI**

pods are routable on cloud fabric

IPAM

Native CNI and/or Antrea CNI

**optional chaining**

Policy Enforcement

Connectivity

# Network Policy Resources

Antrea supports both upstream K8S and native policy primitives

Upstream K8S

Antrea

Antrea Roadmap

Unordered
Simple Matching

Ordered
Advanced Matching

Groups policy together for setting global precedence and managing access via RBAC

security.antrea.tanzu.
vmware.com/v1alpha1

NetworkPolicy

security.antrea.tanzu.
vmware.com/v1alpha1

Tier

networking.k8s.io/v1

NetworkPolicy

security.antrea.tanzu.
vmware.com/v1alpha1

ClusterNetworkPolicy

Policy Scope

Namespace

Namespace

Namespace

Cluster

# NetworkPolicy Implementation



Policy = "Pods with label 'app=server' can only receive traffic from Pods with label 'app=client', and only on port 80.

kube-apiserver

Network Policy definitions
Pod definitions -> namespace, labels, IP address
Namespace definitions -> labels

Antrea Controller

K8s apiserver library

AppliedToGroup =
    Name: "foo"
    Pods: {Pod3A(10.1.3.2)}
AddressGroup =
    Name: "bar"
    Pods: {Pod3B(10.1.3.3)}
NetworkPolicy =
    Rule:
        Direction: Ingress
        From: {"bar"}
        Ports: {80}
        AppliedToGroups: {"foo"}

span

...

**Node 1**
Antrea Agent
OpenFlow
OVS bridge
Pod1A
10.1.1.2
app=other

**Node 2**
Antrea Agent
OpenFlow
OVS bridge
Pod2A          Pod2B
10.1.2.2       10.1.2.3
app=server     app=server

**Node 3**
Antrea Agent
OpenFlow
OVS bridge
Pod3A          Pod3B
10.1.3.2       10.1.3.3
app=server     app=client

Centralized controller for Network Policy computation

Each Node's Agent receives only the relevant data

Very lightweight for the Node's Agent (simple conversion to flows)

Controller = single source of truth
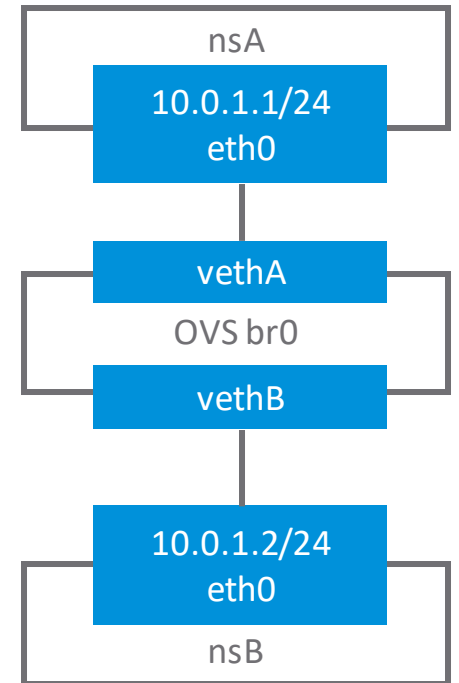- Easier to debug

Multiple controllers possible
- HA
- Controller scale-out
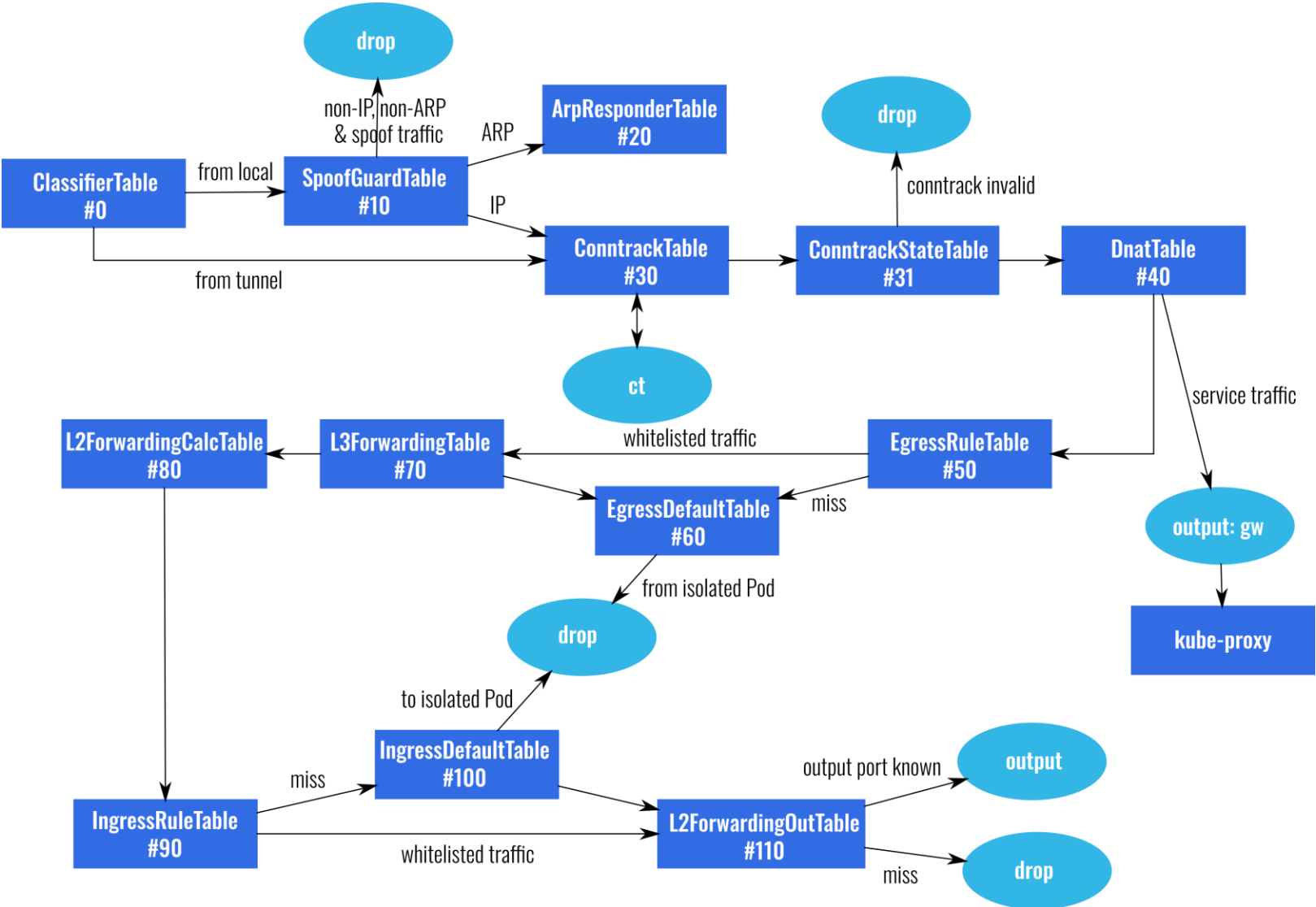
Use OVS flow conjunction
- Reduces number of flows

# OVS Hello World

```
> ovs-vsctl add-br br0
> ovs-vsctl add-port br0 vethA
> ovs-vsctl add-port br1 vethB
> ip netns exec nsA ping -c 1 -W 1 10.0.1.2 && echo "SUCCESS"
SUCCESS
>
> ovs-ofctl add-flow br0 priority=100,icmp,actions=drop
> ip netns exec nsA ping -c 1 -W 1 10.0.1.2 || echo "FAILED"
FAILED
>
> ovs-ofctl dump-flows br0
table=0, n_packets=1, n_bytes=98, priority=100,icmp actions=drop
table=0, n_packets=18, n_bytes=1092, priority=0 actions=NORMAL
```

# OVS Pipeline

# Antrea Packet Walk Across Network Layers