# MIZAR

*Futurewei Technologies*

Current state and work in progress

# THE PROBLEM

Support provisioning of large number endpoints (10M)

Achieve high network routing throughput and low latency

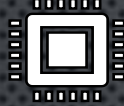Create an extensible plugin framework for cloud network

Unify the network data-plane for VMs, Containers, Serverless and others

# THE PROBLEM WITH CURRENT SOLUTIONS

**Program every host every time a user provisions an endpoint:**

**Approaching cloud-networking with a conventional programming model and network devices**

e.g. OpenFlow programming in OVS

Virtual Switches and Routers are essentially softwarization of hardware switches and routers, but not necessarily programmable to support rapid network changes.

**Existing solutions bring up software network devices, that are primarily created for Telecom, ISP, or datacenter networking and run them in virtual machines to support cloud networking.**
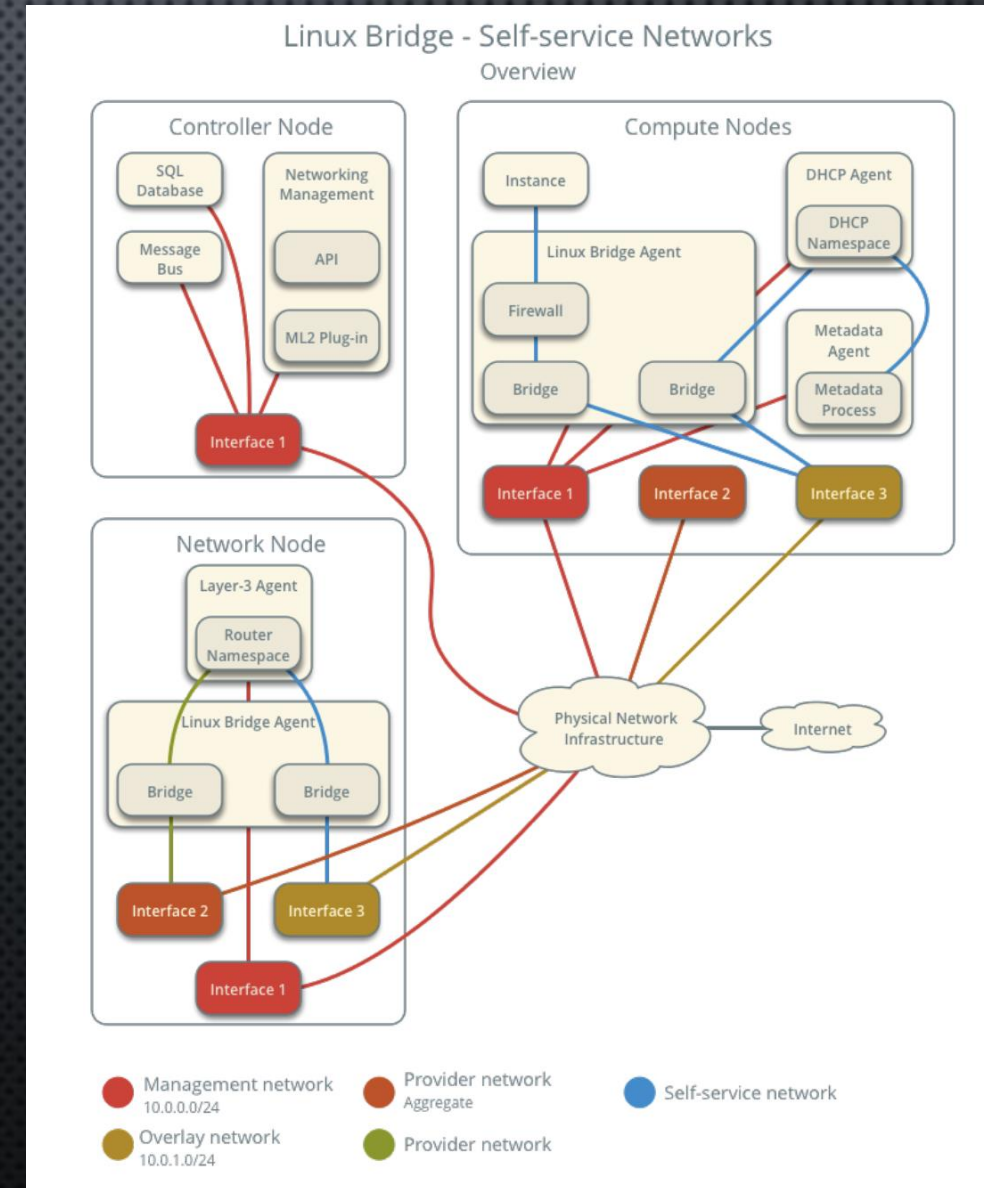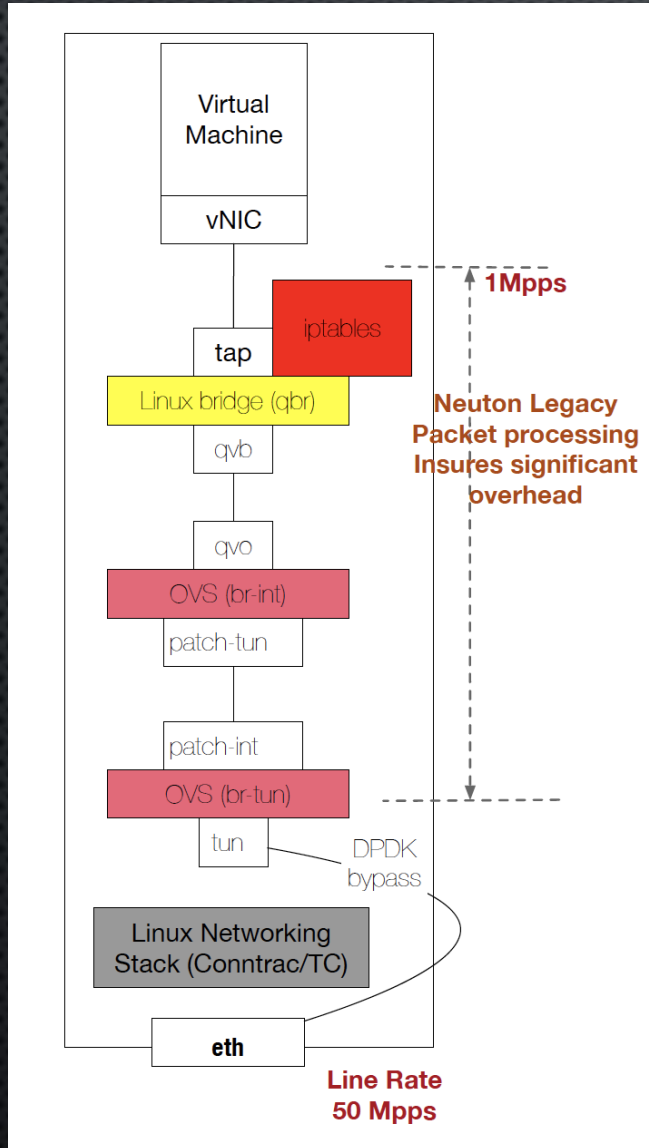
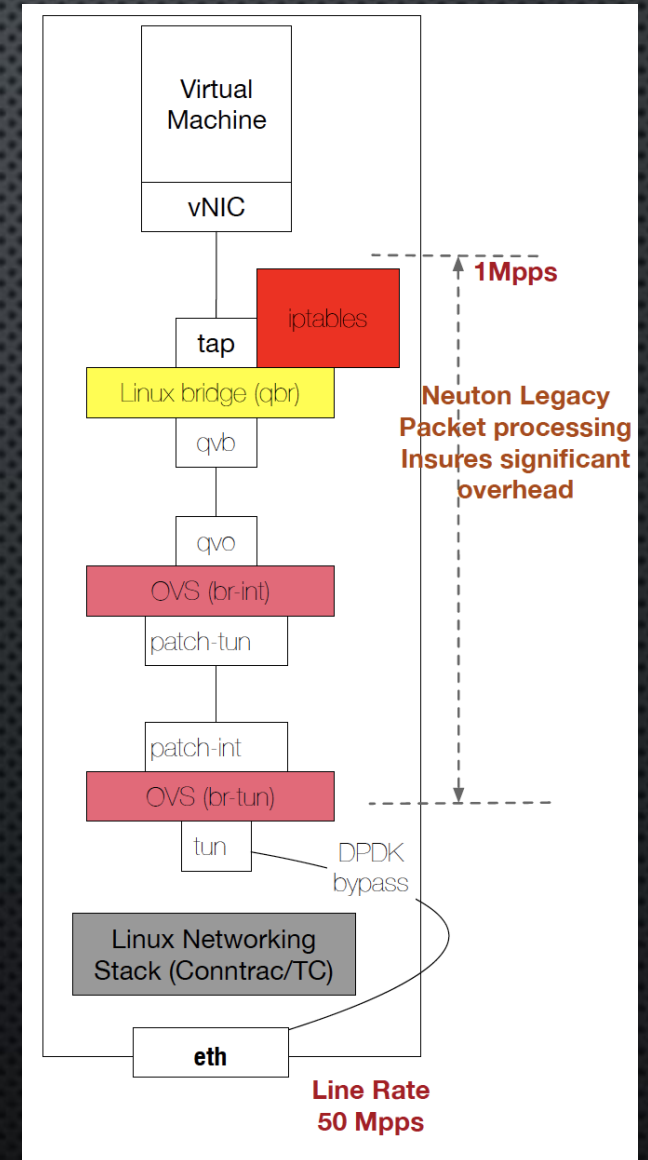**Packets traverse multiple network stacks on the same host**

**Packets traverse multiple network devices (as if we are operating a data-center), while these functions could be consolidated during design.**

# OBSERVATION: (NOT REALLY A NEW ONE)

- In a cloud network (overlay), most functions can be reduced to
  1. Encapsulate/decapsulate a packet
  2. Modify the outer packet header and forward it
  3. Modify the inner packet header and forward it
  4. Drop unwanted packets

- Several network functions can be thought of in a similar way:
  1. Responding to ARP
  2. DHCP
  3. NAT
  4. Passthrough load-balancing
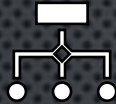
# MIZAR OVERALL ARCHITECTURE!

**Natural Partitioning domains of Cloud Network**

Virtual Private Cloud VPC

Networks within a VPC

Endpoints within a network

**Bouncer**

In-network hash tables

Holds the configuration of endpoints within a network

Determines flow modifications, and bounce the packet for TX

Implements all middleboxes within a network

**Divider**

In-network hash tables

Holds the configuration of Bouncers within a VPC

Divides the VPCs endpoint's configuration into clusters of Bouncers

Implements all middleboxes at the VPC level

# MIZAR OVERALL ARCHITECTURE!

**Natural Partitioning domains of Cloud Network**

Virtual Private Cloud VPC

Networks within a VPC

Endpoints within a network

**Bouncer**

In-network hash tables

Holds the configuration of endpoints within a network

Determines flow modifications, and bounce the packet for TX

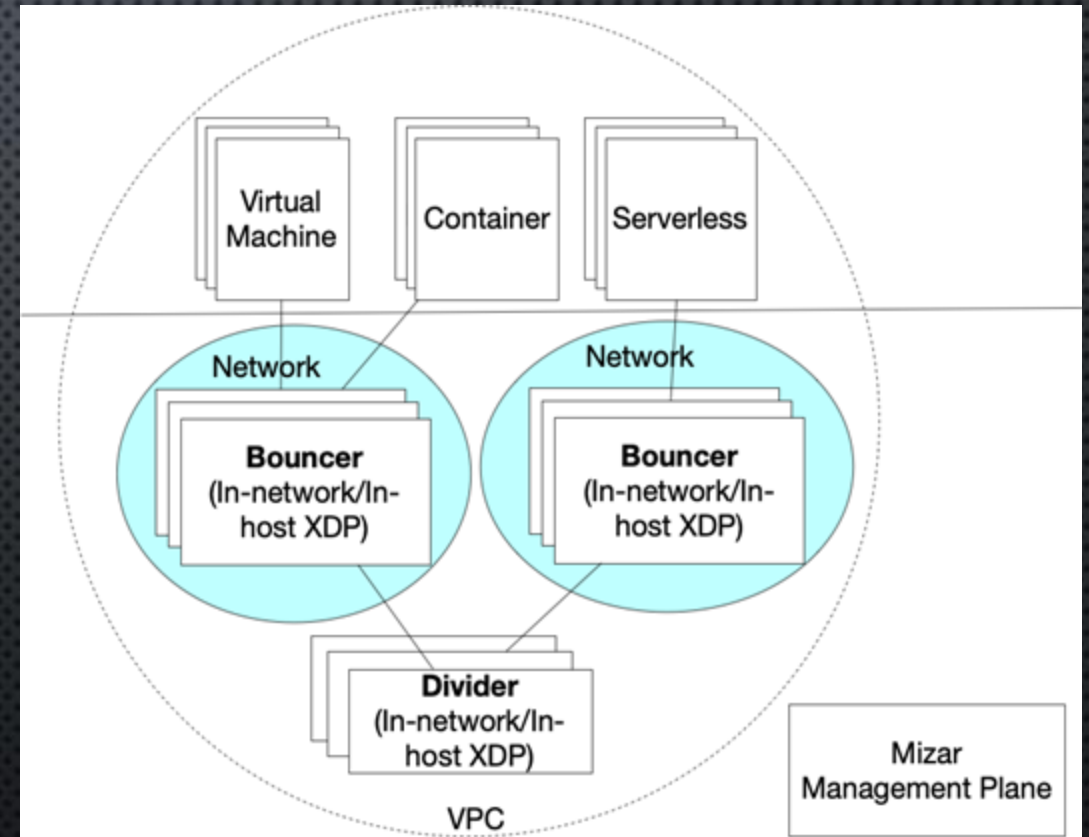Implements all middleboxes within a network

**Divider**

In-network hash tables

Holds the configuration of Bouncers within a VPC

Divides the VPCs endpoint's configuration into clusters of Bouncers

Implements all middleboxes at the VPC level
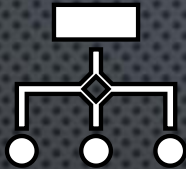
# MANAGEMENT PLANE ARCHITECTURE

## Kubernetes CRDS
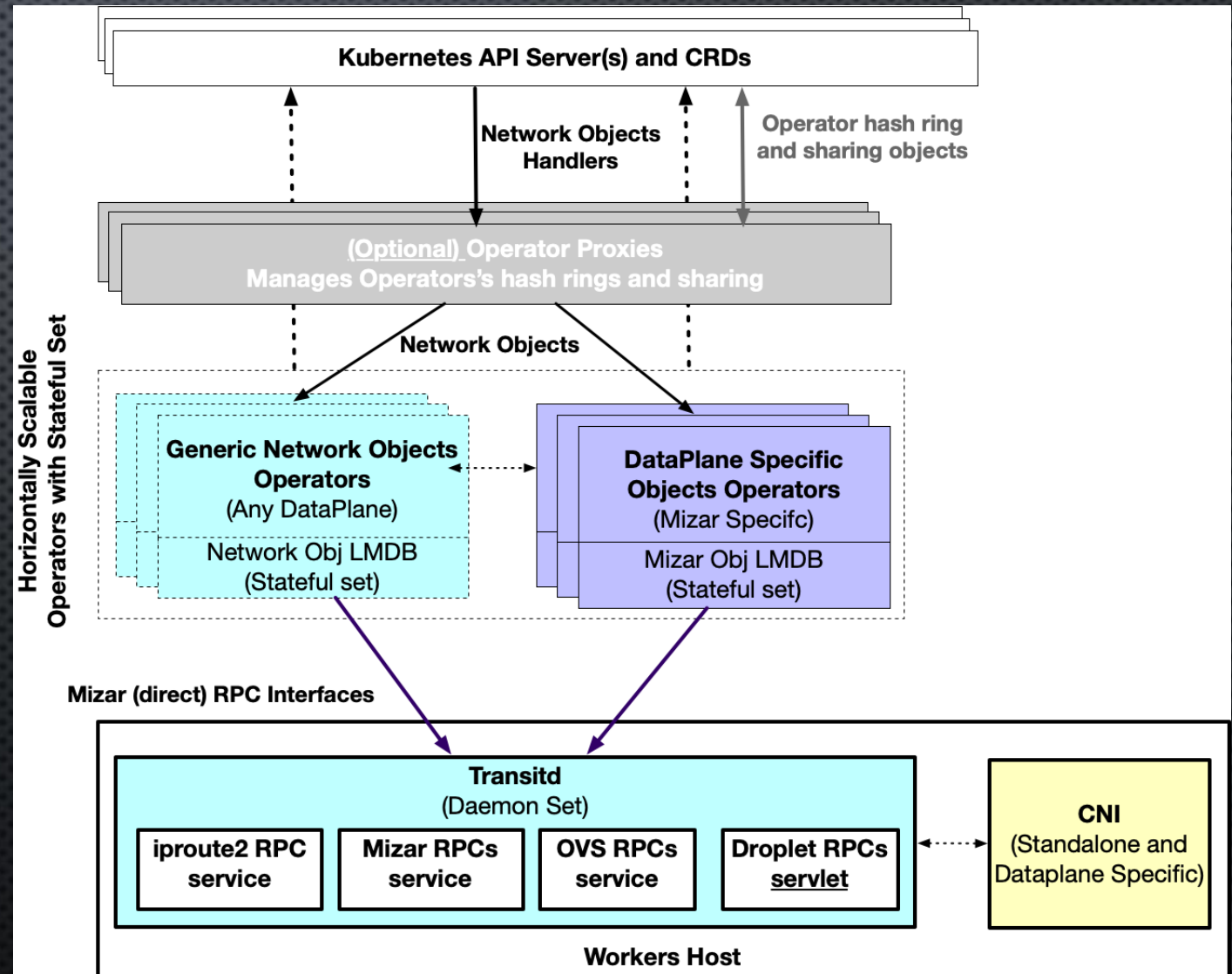
Allows us to extend the K8s API server with networking objects. Some of these objects are generic to any networking solution and some of them are specific to Mizar.

## Operator Framework

Extending Kubernetes with domain-specific operators that *act as a controller* for CRDs. Operator framework allow us to write custom lightweight operators that derive the networking objects life-cycle.

# OPERATORS, CRDS, AND WORKFLOWS

VPC　　　　Divider　　　　Network　　　　Bouncer　　　　Endpoint　　　　Droplet

All objects are defined through Kubernetes CRDs under the mizar.com resource group. An operator exposes interfaces to the management-plane workflows.

The lifecycle of an object is governed by three workflows: create, update, and delete.
Each of these workflows are triggered by state changes in the respective object.
For example, the droplet object which represents a physical interface on a node
Will trigger the management-plane delete workflow if the physical interface itself were to be removed.

# THE CRDS

**VPC**
- Contains information about the VPC such as CIDR range, and list of dividers

**Bouncer**
- Contains information about the bouncer such as parent Network, and host information (IP, MAC, etc)

**Divider**
- Contains information about the divider such as parent VPC, and host information (IP, MAC, etc)

**Endpoint**
- Contains information about the bouncer such as parent Network, and host information (IP, MAC, etc)

**Network**
- Contains information about the Network such as parent VPC, CIDR range, and bouncers.

**Droplet**
- Contains information about the current host interface (IP, MAC, etc_

# INSIDE A MIZAR HOST

# BACKGROUND XDP: SIMPLIFIED AND EXTENSIBLE PACKET PROCESSING NEAR LINE RATE

Skip unnecessary stages of network stack whenever possible and transit packet processing to smart NICs.

Makes the best use of kernel packet processing constructs without being locked in to a specific processor architecture.

Packet processing is entirely in-kernel.

Very small programs < 4KB

# IN-HOST PACKET FLOW: BYPASS NETWORK STACK

Packets traverses only the container stack

Namespace A

Namespace B

veth0
(VETH pair in container)

TX    RX

veth0
(VETH pair in container)

TX    RX

Root Namespace

RX

**Transit Agent XDP**

RX

**Transit Agent XDP**

TX

TX

veth_peer
(VETH pair in root namespace)

veth_peer
(VETH pair in root namespace)

**Mizar User-space Network Function**

XDP_REDIRECT (SKB)

XDP_REDIRECT (Frame)

eth0
(physical interface)

TX

**Transit XDP**

RX

AF_XDP

- Implements essential logical networking function within the same XDP program that provides multitenant cloud networking solutions through new Bouncer and Divider concepts

- Mizar autonomously adapts to various traffic demands in immense scale cloud environments. Allowing Mizar to serve various cloud workloads in a multi-tenant environment optimally.

- Extensible support of native networking features through custom chains of optimized XDP programs hooks and Geneve protocol options. Future possible Features including: Security, Loadbalancing, Connectivity, Traffic Shaping Control

One Efficient XDP Program with Extensible Functions

Divider — Tail-Call → Divider XDP Hook — FORWARD →
e.g. Cross-VPC traffic

Bouncer — Tail-Call → Bouncer XDP Hook — FORWARD →

Endpoint lookup — Tail-Call → Endpoint XDP Hook — REDIRECT → Endpoint veth pair
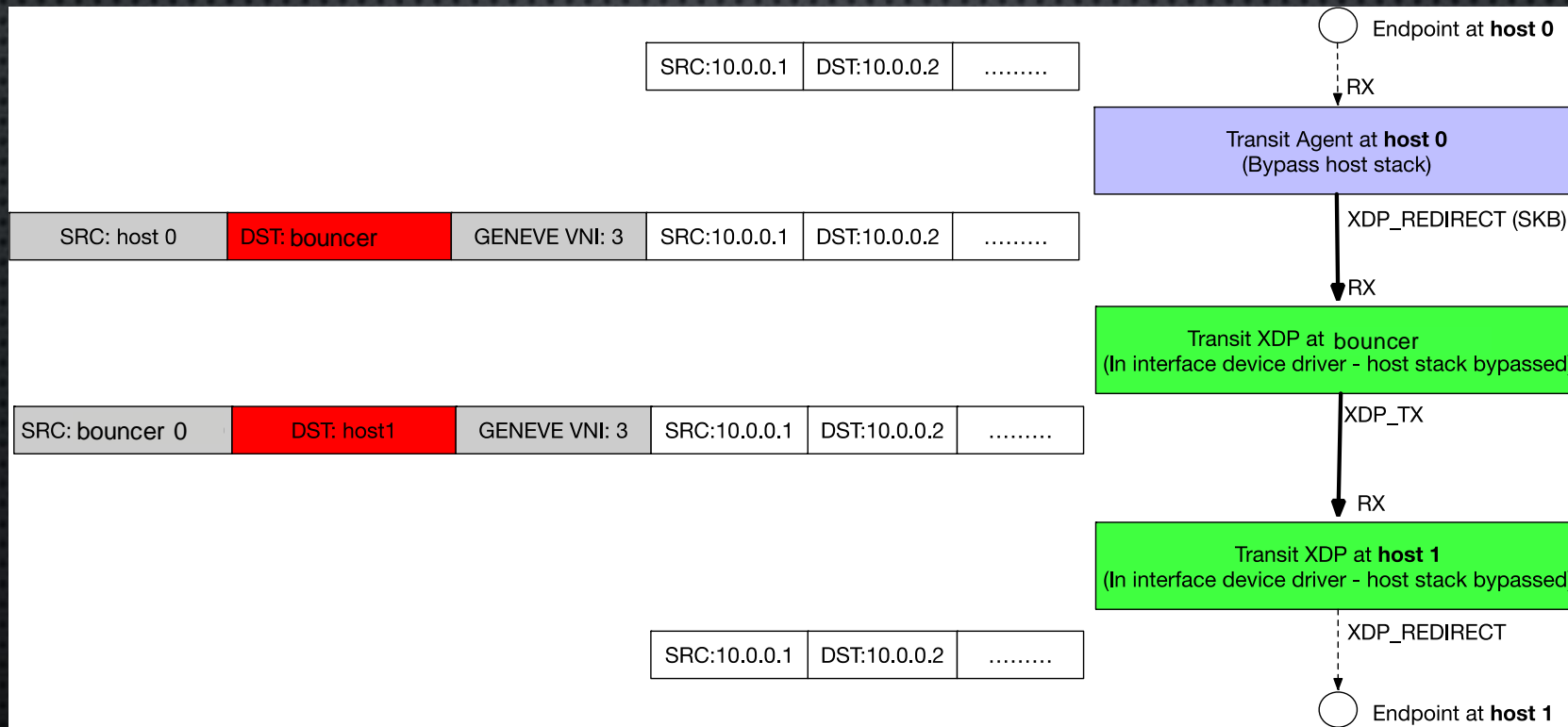e.g. Security Groups

# EXAMPLE PACKET WITHIN A NETWORK
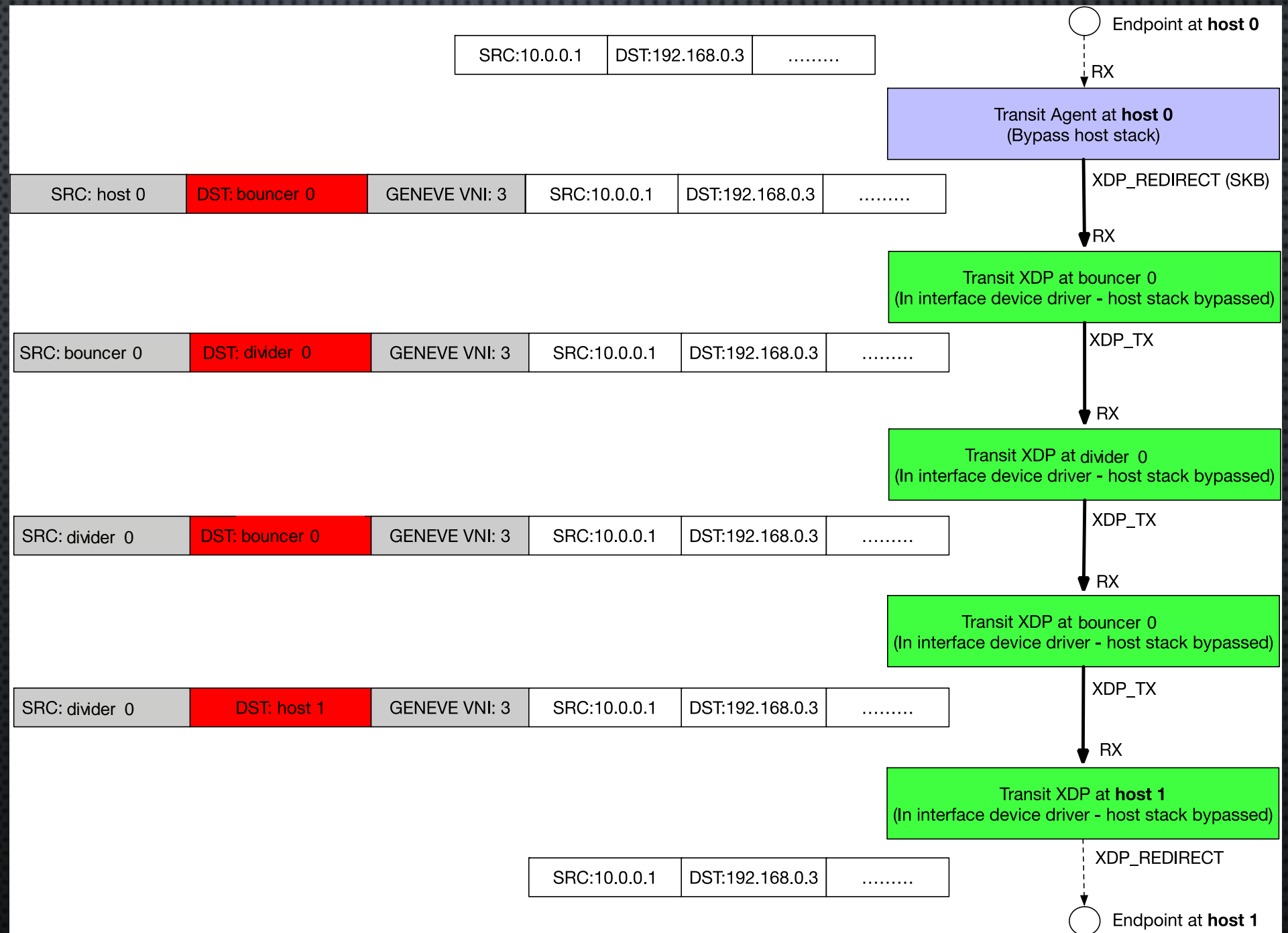
**Three steps to provision an endpoint**

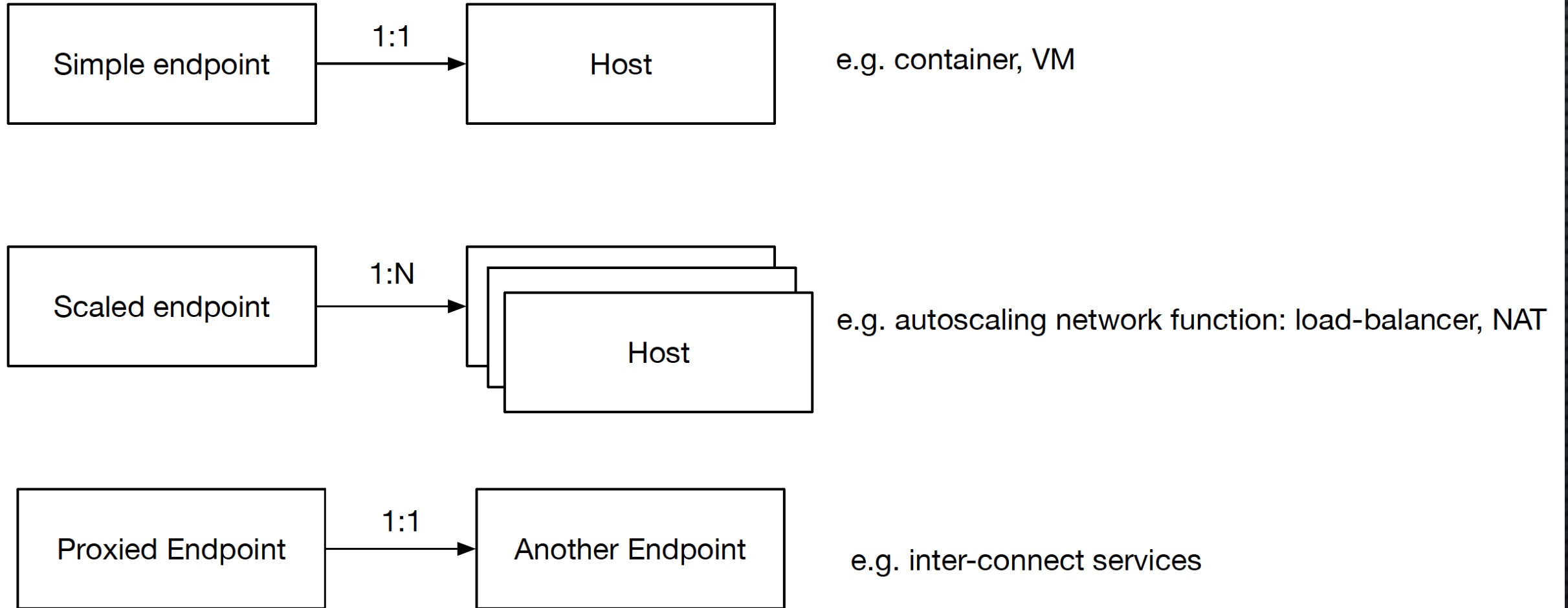| Add the endpoint to N Bouncer table | Provision the endpoint on the host | Configure the host transit agent to tunnel the endpoint traffic to the Bouncer |
|---|---|---|

CROSS NETWORK

Endpoint at **host 0**

| SRC:10.0.0.1 | DST:192.168.0.3 | ........ |

RX

Transit Agent at **host 0**
(Bypass host stack)

XDP_REDIRECT (SKB)

| SRC: host 0 | DST: bouncer 0 | GENEVE VNI: 3 | SRC:10.0.0.1 | DST:192.168.0.3 | ........ |

RX

Transit XDP at bouncer 0
(In interface device driver - host stack bypassed)

XDP_TX

| SRC: bouncer 0 | DST: divider 0 | GENEVE VNI: 3 | SRC:10.0.0.1 | DST:192.168.0.3 | ........ |

RX

Transit XDP at divider 0
(In interface device driver - host stack bypassed)

XDP_TX

| SRC: divider 0 | DST: bouncer 0 | GENEVE VNI: 3 | SRC:10.0.0.1 | DST:192.168.0.3 | ........ |

RX

Transit XDP at bouncer 0
(In interface device driver - host stack bypassed)

XDP_TX

| SRC: divider 0 | DST: host 1 | GENEVE VNI: 3 | SRC:10.0.0.1 | DST:192.168.0.3 | ........ |

RX

Transit XDP at **host 1**
(In interface device driver - host stack bypassed)

XDP_REDIRECT

| SRC:10.0.0.1 | DST:192.168.0.3 | ........ |

Endpoint at **host 1**

# NEW ENDPOINT TYPES

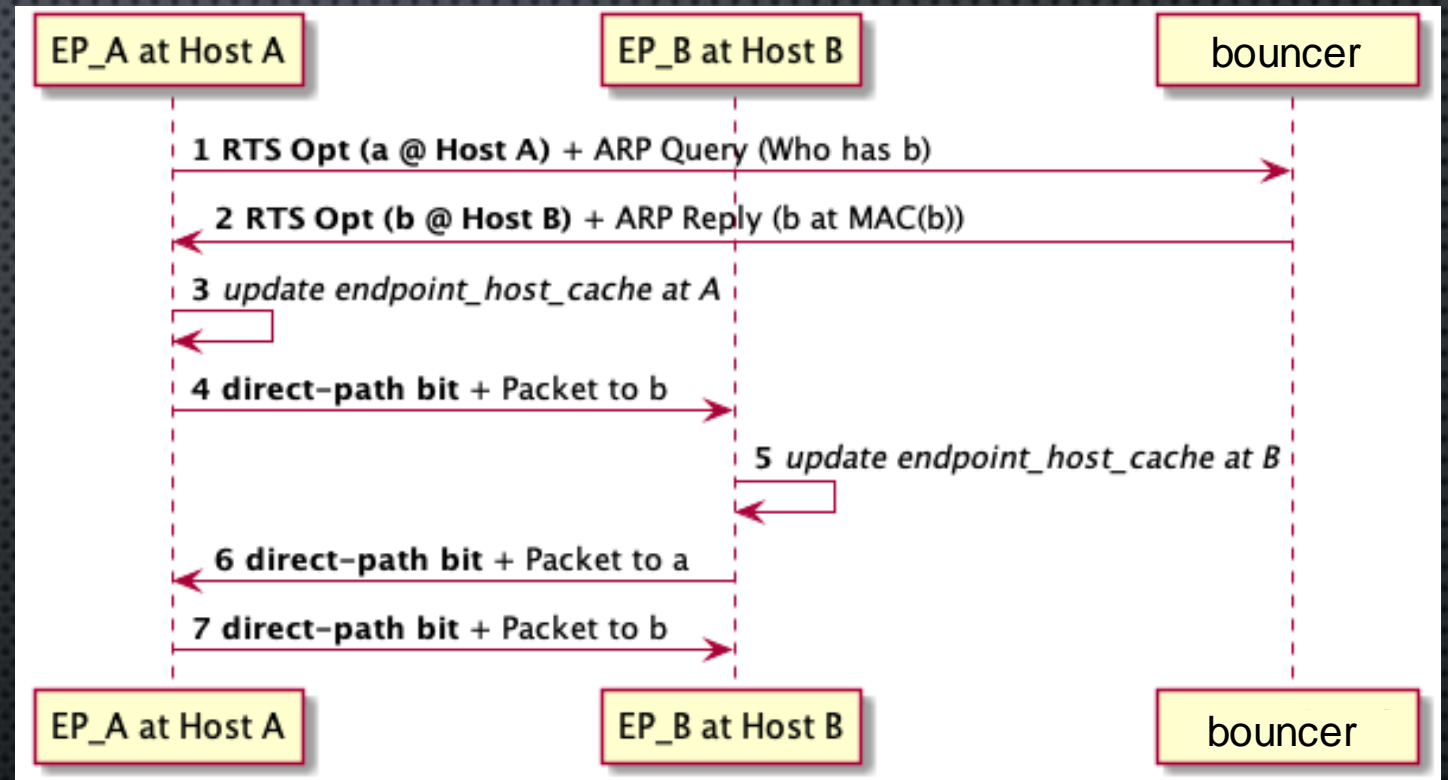| | | |
|---|---|---|
| Simple endpoint | —1:1→ Host | e.g. container, VM |
| Scaled endpoint | —1:N→ Host | e.g. autoscaling network function: load-balancer, NAT |
| Proxied Endpoint | —1:1→ Another Endpoint | e.g. inter-connect services |

# DIRECT PATH

The direct path feature allows a flexible data-plane fast path that avoids multiple hops of flows through bouncers. Endpoints shall be able to communicate directly between their hosting nodes without the need to go through multiple intermediaries.

Endpoint                          Endpoint

# DIRECT PATH: INTRA-NETWORK

To support the direct path feature, the transit agent may add a Return To Sender (RTS) option to the outer packet header. The option carries information about the endpoint's host IP and MAC addresses.
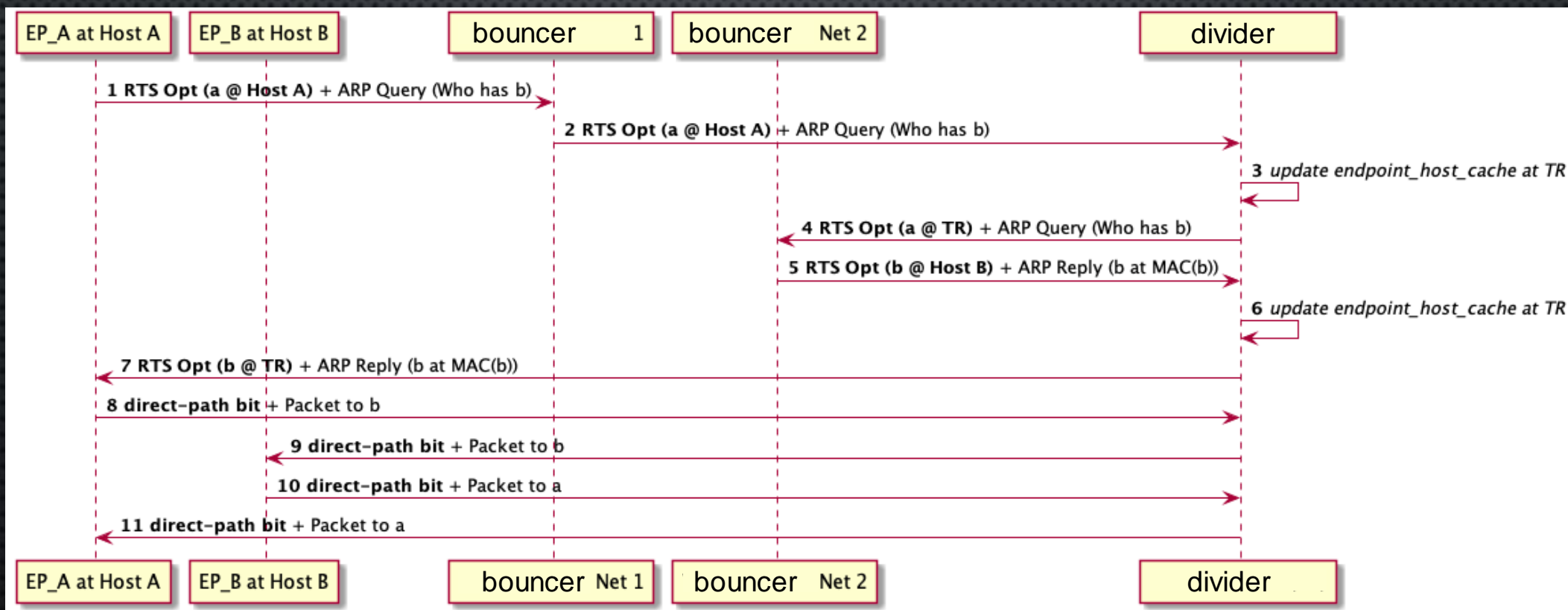
With direct path, all packets sent proceeding the first initial packet is sent directly between hosts
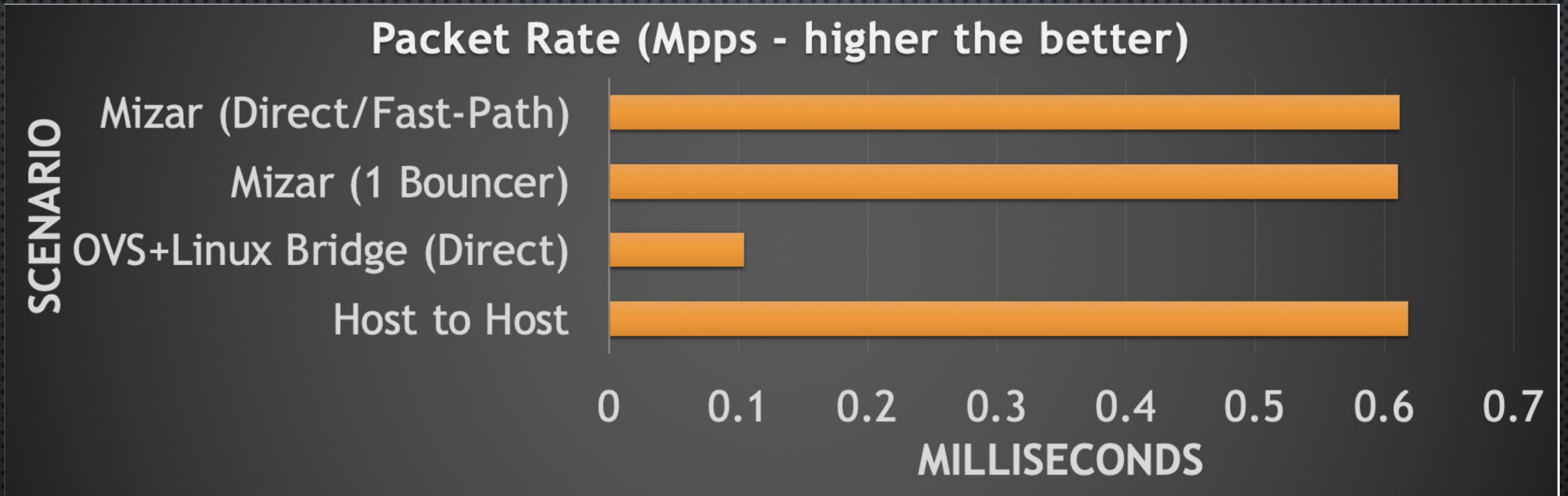
# DIRECT PATH: INTER-NETWORK

As the divider will be evaluating network access control (ACL), it is important that all inter-network and inter-VPC packets be processed by the dividers to take effect.

All packets proceeding the initial packet will bypass the bouncers and be sent between the two hosts and the divider.

# PACKET RATE (NON TCP) – SCALING NETWORK SERVICES



Packet Rate (Mpps - higher the better)

- Near line rate packet per second

HIT

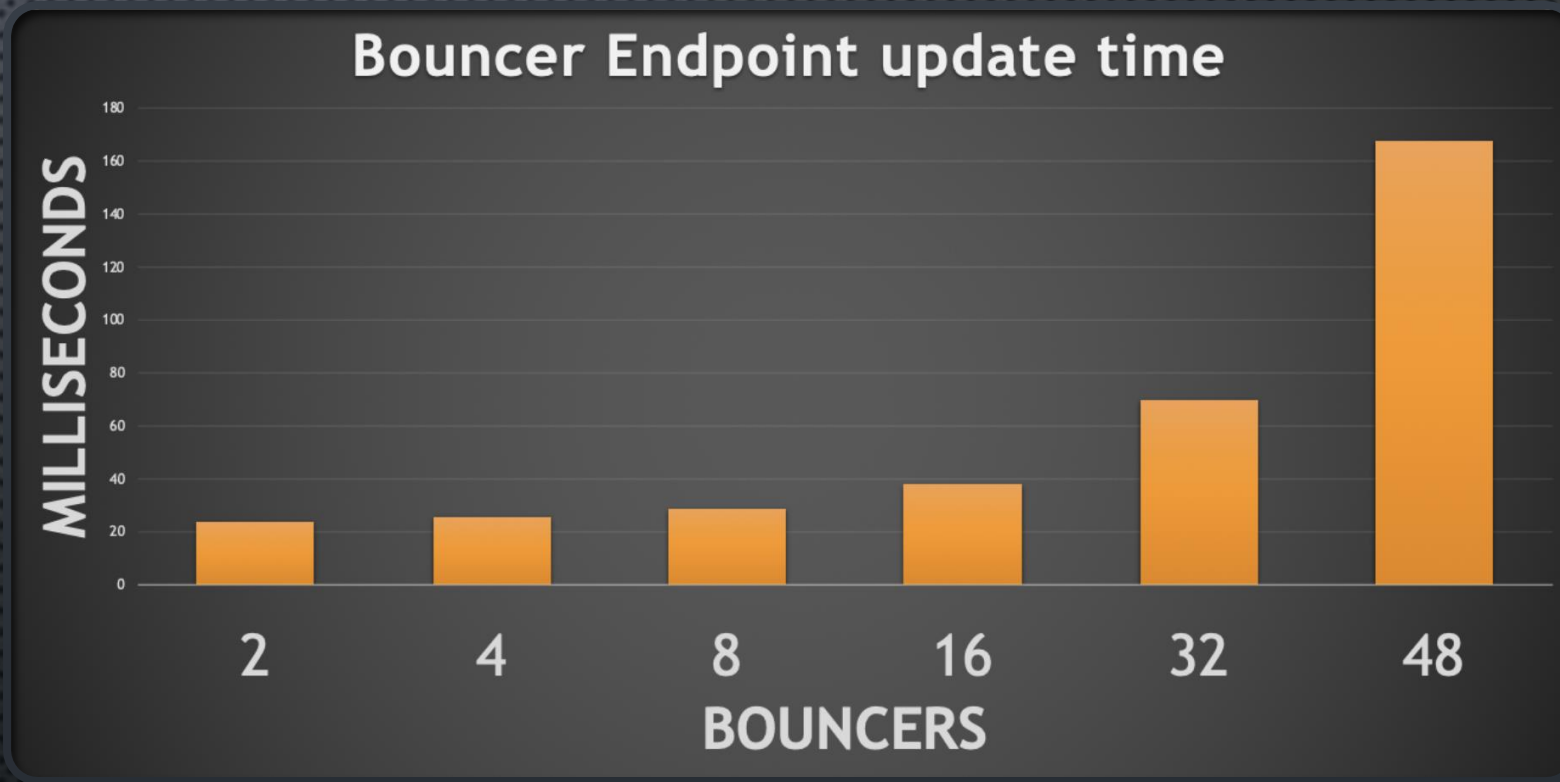# ENDPOINT UPDATE TIME WITH MULTIPLE BOUNCERS

## HIT

- Constant time with parallel updates (20ms) until the Test Controller starts to Hit its re

## Scale

- With a scalable Control-plane (on multiple machines), we foresee maintenance of constant time scaling.

## IMPROVEMENT

- Simplifications in data-plane as we introduce the scaled endpoint. One core required.

# ENDPOINT E2E PROVISIONING TIME MULTIPLE BOUNCERS
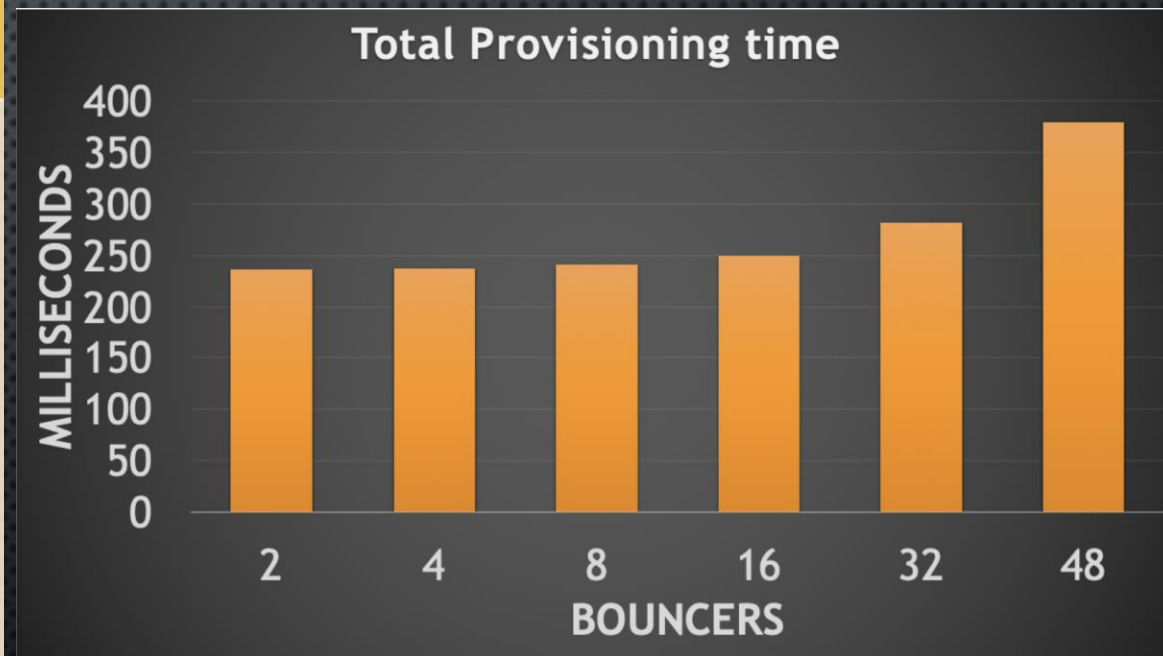
| HIT | Overhead | IMPROVEMENT |
|---|---|---|
| • Scale remains constant (until hitting test controller machine limits) | • Primarily overhead on the host from creating the virtual interfaces by executing shell command (~250 ms). | • Expected to improve with production ready control-plane as it makes use of netlink. |



Total Provisioning time

# ROUND TRIP TIME EFFECT ON END-USER

**HIT**

Mizar direct path is faster than OVS+Linux Bridge. Though, Still has minimal impact on PPS and TCP BW.

**HIT**

Even with an increased latency due to the extra hop, the packet per second processed by endpoints remains close to line rate.

**Benefit**

Primarily benefit of fast-path is latency sensitive applications.



Round Trip Time (milliseconds - lower the better)

| SCENARIO | |
|---|---|
| Mizar (Direct/Fast-Path) | |
| Mizar (1 Bouncer) | |
| OVS+Linux Bridge (Direct) | |
| Host to Host | |

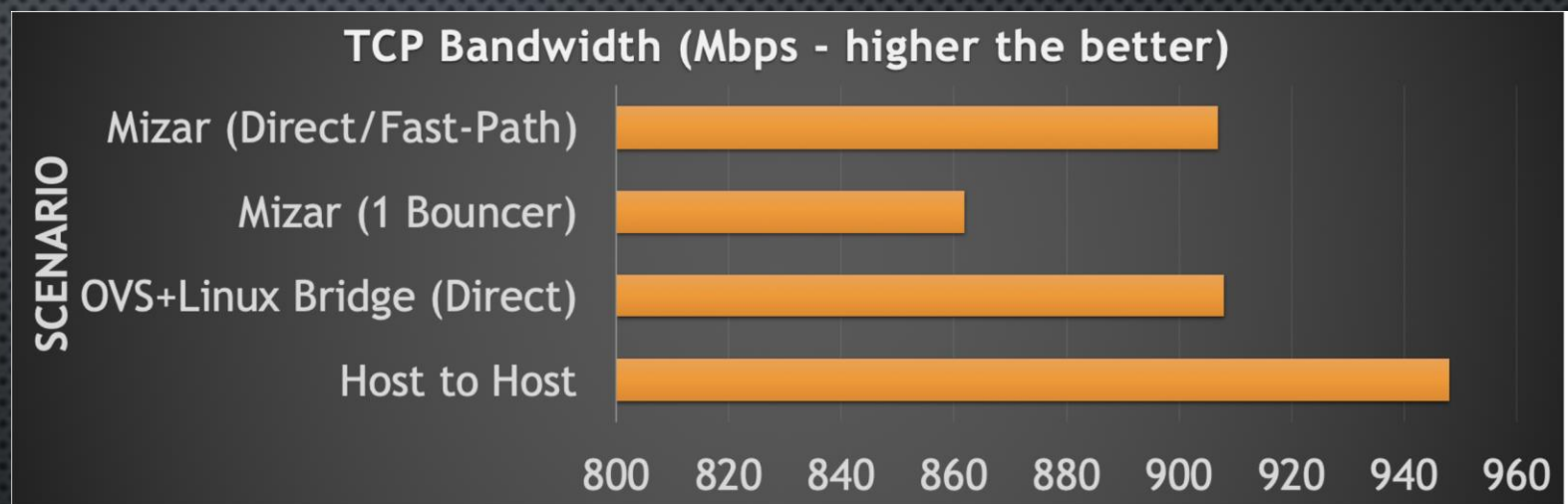0    0.05    0.1    0.15    0.2    0.25    0.3
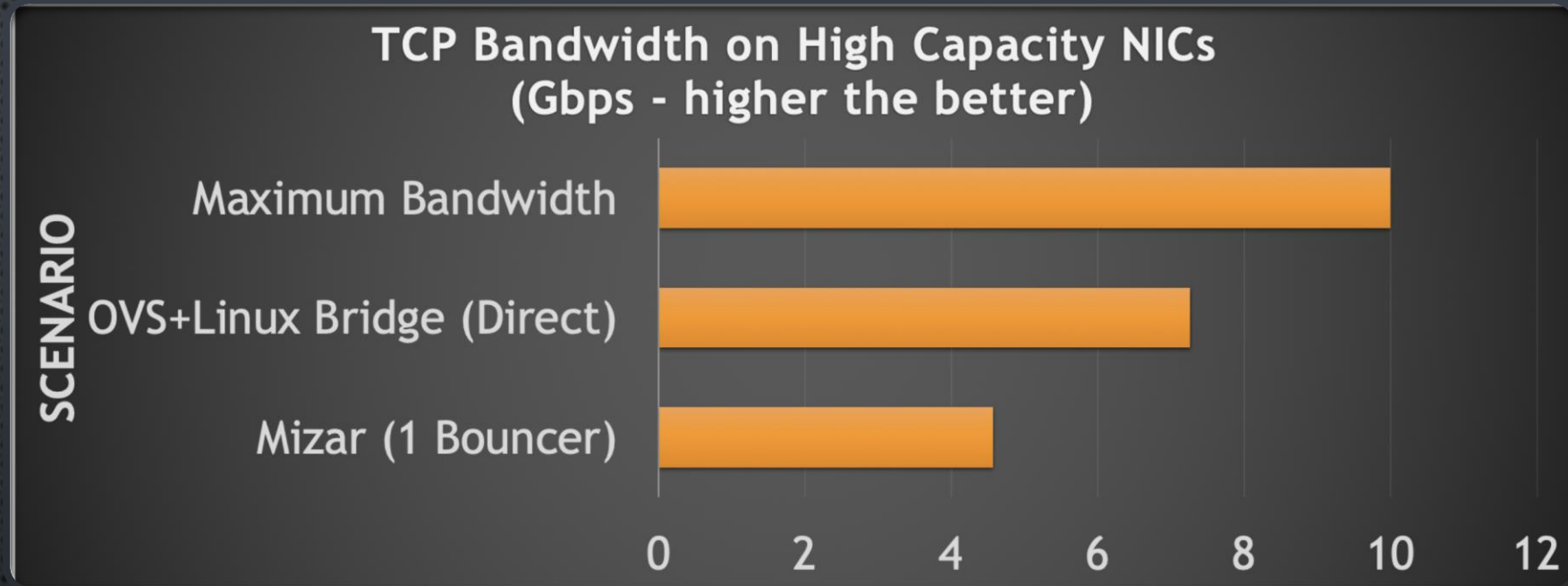
# TCP BANDWIDTH (ON A SLOW NIC 1GBPS)

**HIT**

Comparable throughput to OVS+Bridge (even though we don't use XDP driver mode). *This is applicable for NICS < 4Gbps*

**Hops:**

The Bouncer hop accounts only for 5% less TCP throughput, which shall be negligible for very high bandwidth NICs. This is despite that RTT of the extra hop accounts for 45% more latency.



TCP Bandwidth (Mbps - higher the better)

# TCP BANDWIDTH (ON A FASTER NIC 10 GBPS)



TCP Bandwidth on High Capacity NICs
(Gbps - higher the better)

**MISS:**
The TCP bandwidth caps at around 4Gbps.

**IMPROVEMENT:**
Change to Driver mode (require support in NIC)

**IMPROVEMENT:**
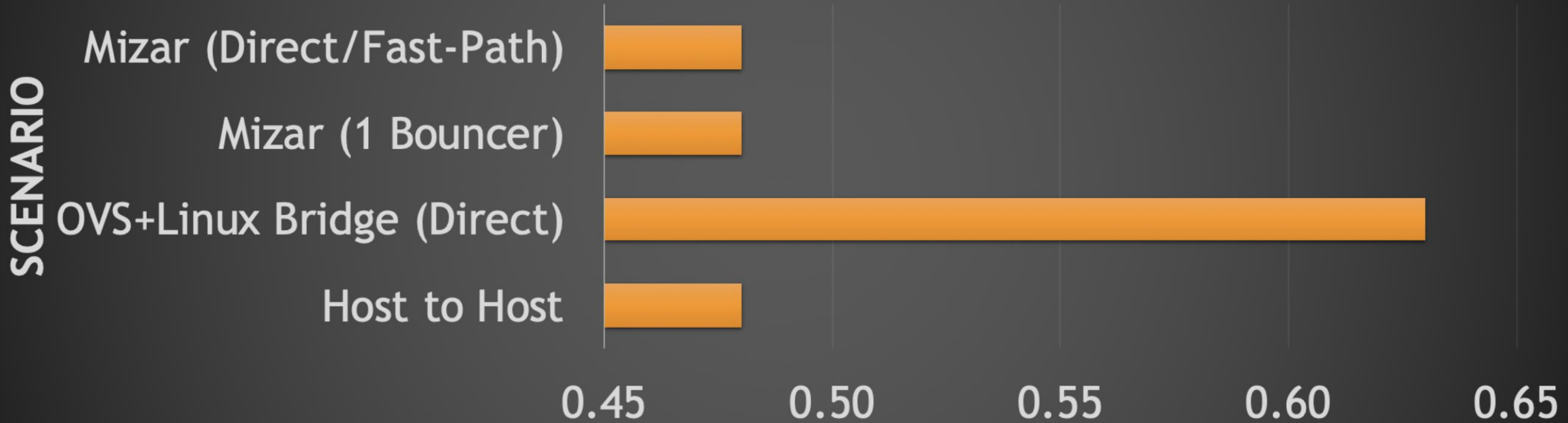Change on-host wiring architecture and reduce reliance on Transit Agent

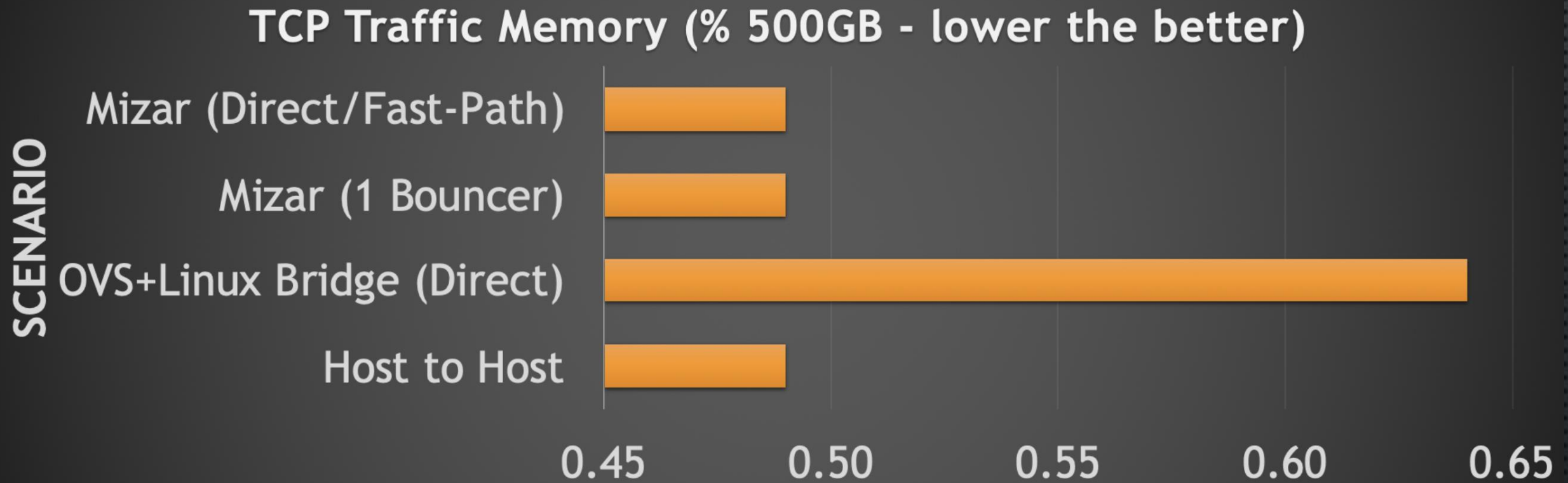**IMPROVEMENT:**
Improved device driver for veth

# MEMORY IDLE CASE
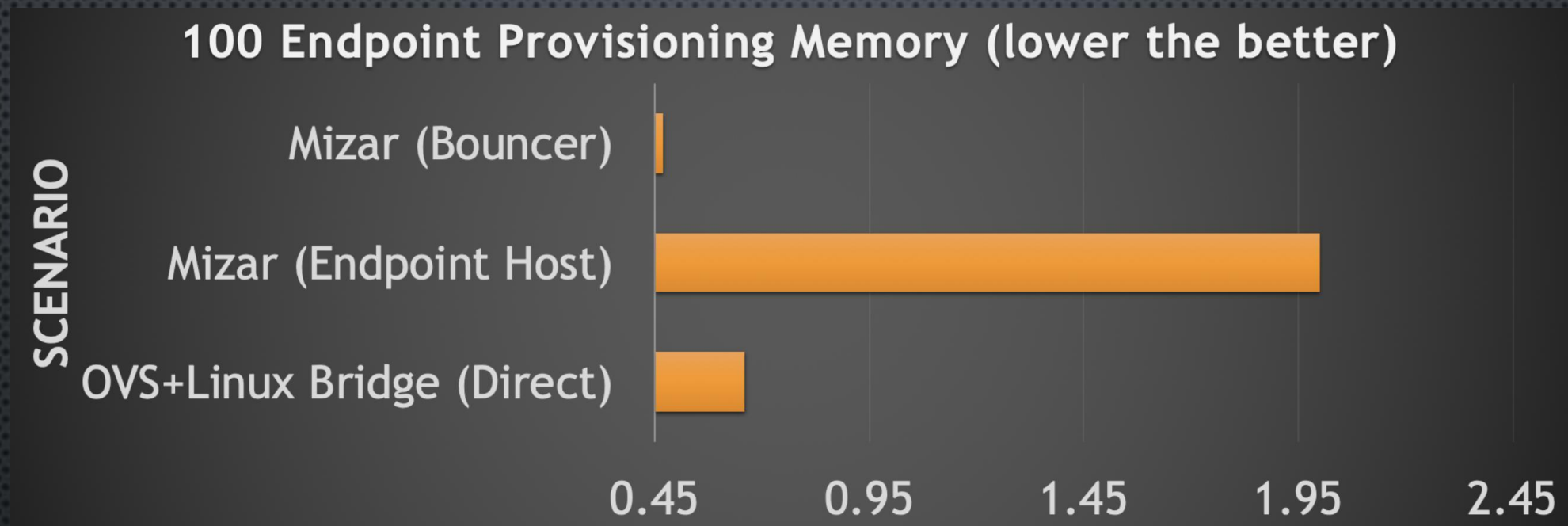


Baseline Memory (% 500GB - lower the better)

SCENARIO

- Mizar (Direct/Fast-Path)
- Mizar (1 Bouncer)
- OVS+Linux Bridge (Direct)
- Host to Host

0.45    0.50    0.55    0.60    0.65

# Memory During TCP Performance Tests



TCP Traffic Memory (% 500GB - lower the better)

SCENARIO

- Mizar (Direct/Fast-Path)
- Mizar (1 Bouncer)
- OVS+Linux Bridge (Direct)
- Host to Host

0.45    0.50    0.55    0.60    0.65

✓
HIT

- Negligible Memory overhead very close to an idle host without networking constructs event with Traffic processing

# MEMORY IDLE CASE (100 ENDPOINTS PER HOST)



**100 Endpoint Provisioning Memory (lower the better)**

SCENARIO

| Scenario | |
|---|---|
| Mizar (Bouncer) | |
| Mizar (Endpoint Host) | |
| OVS+Linux Bridge (Direct) | |

0.45    0.95    1.45    1.95    2.45

## HIT
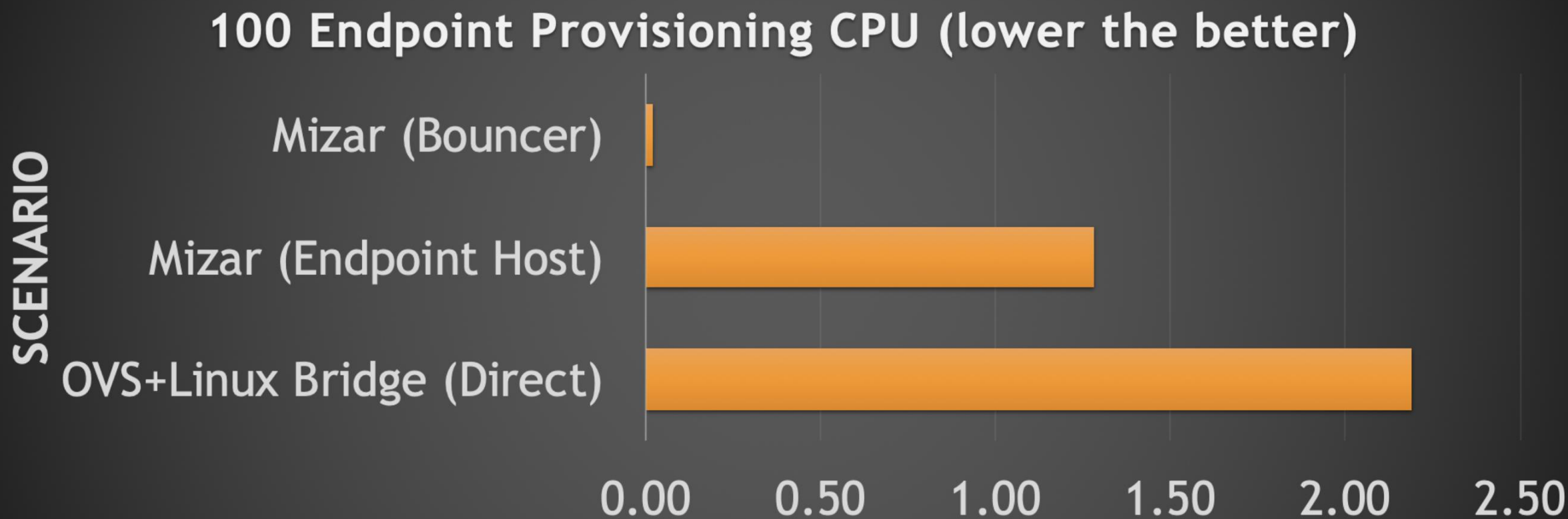Memory overhead on Bouncers remain at baseline level

## MISS
On Host memory increases as we provision more endpoints

## IMPROVEMENT
Share one transit agent across multiple endpoints

# CPU DURING TCP PERFORMANCE TESTS



100 Endpoint Provisioning CPU (lower the better)

- Significantly less CPU overhead during provisioning on both bouncer and host

HIT