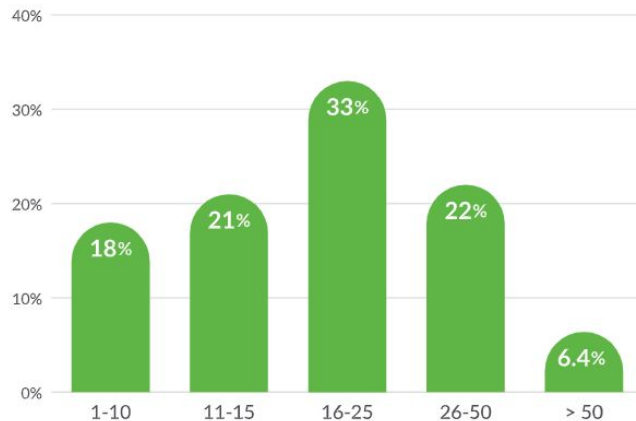# Problem Statement

➜ Increasing Pod density per node

➜ Competition for node resources e.g. CPU and memory

➜ Optimization problem for operators: Resource allocation and efficient use, achieving SLOs, etc

```
apiVersion: v1
kind: Pod
metadata:
  name: frontend
spec:
  containers:
  - name: app
    image: images.my-company.example/app:v4
    resources:
      requests:
        memory: "64Mi"
        cpu: "250m"
      limits:
        memory: "128Mi"
        cpu: "500m"
```
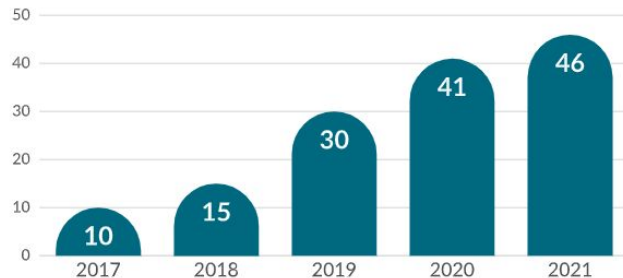
What the container is guaranteed to get, e.g. kubelet will only schedule the Pod on a node which can provide this resource.

Hard upper limit, ensures that container never goes above this threshold.

**Pods per Node**

40%
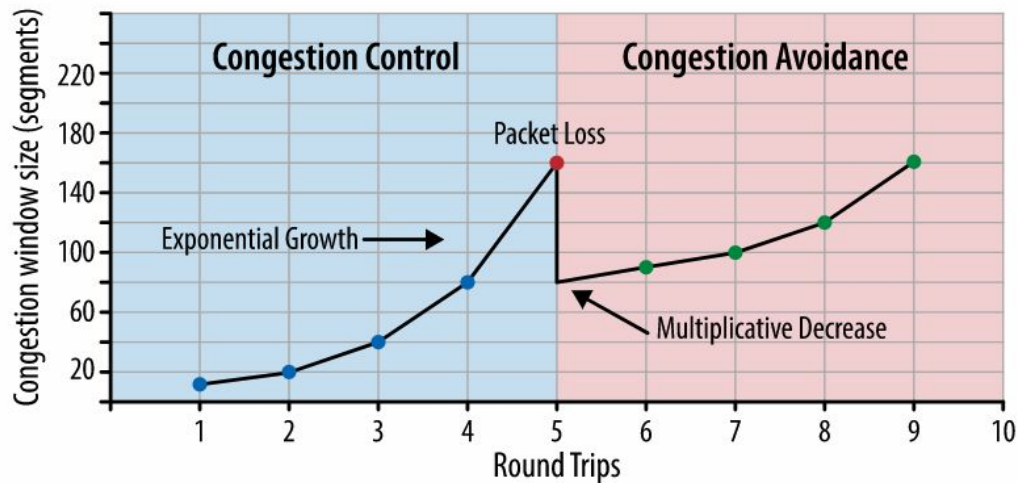
30%  33%

20%  18%  21%  22%

10%

0%  1-10  11-15  16-25  26-50  > 50  6.4%

**Median Containers per Host**

50

40  46

30  41

30

20  15

10  10

0  2017  2018  2019  2020  2021

Source: Sysdig 2022 Cloud Native Security and Usage Report

2

# Problem Statement

➔ But what about networking?
➔ TCP sends AFAP
   (as fast as possible)

➔ AFAP output contract, shaping typically implemented by device output queues
➔ Queue length limit & receive window determines in-flight rate
➔ "How fast" implicit in queue drain rate

➔ Who limits a Pod's network usage in Kubernetes?



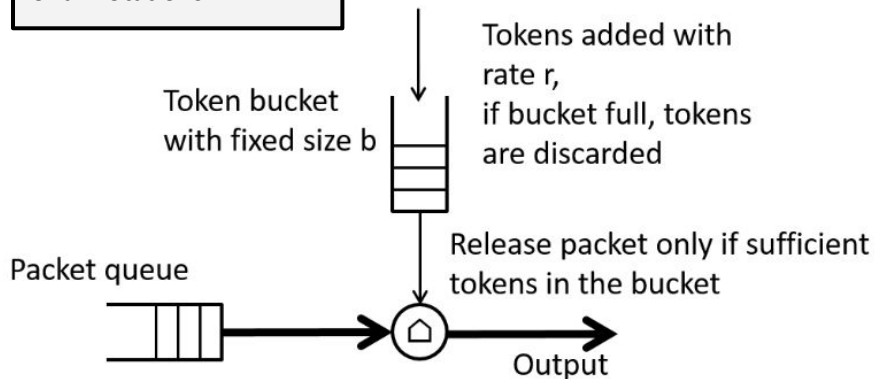Source: https://hpbn.co/building-blocks-of-tcp/

3

# Problem Statement

➔   Kubernetes bandwidth enforcement has only been <u>experimental</u> so far :-(
➔   Support for Pod annotations with 'outsourced' <u>bandwidth meta plugin</u>

```
apiVersion: v1
kind: Pod
metadata:
  annotations:
    kubernetes.io/ingress-bandwidth: 1M
    kubernetes.io/egress-bandwidth: 1M
...
```
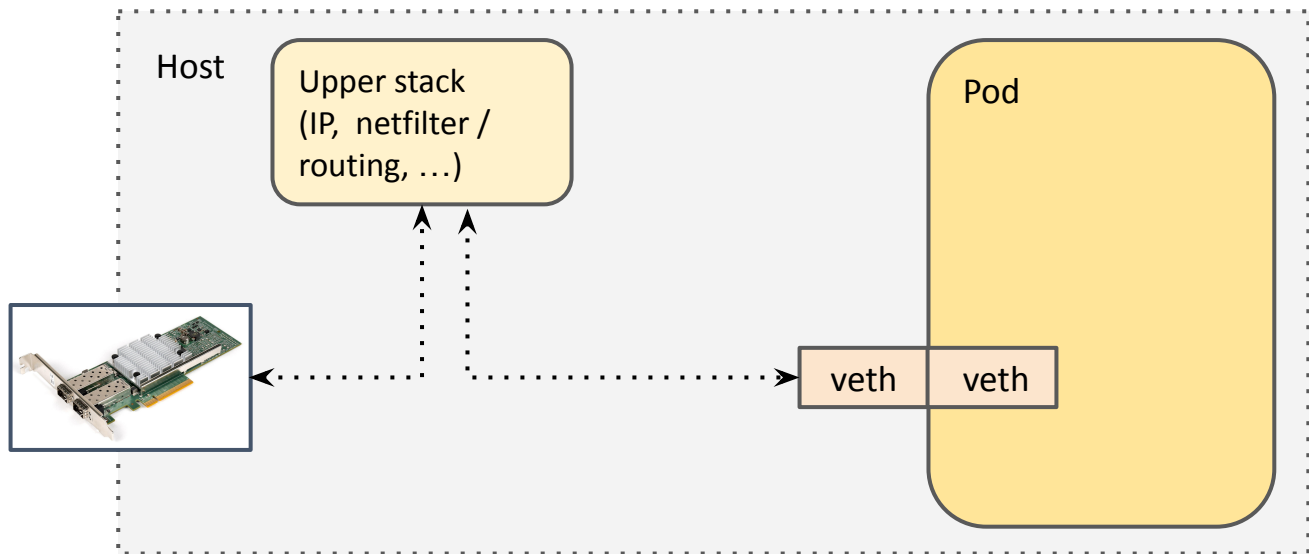
Adds rudimentary **token bucket filters (TBF)** to implement enforcement of annotations.

Token bucket with fixed size b

Tokens added with rate r, if bucket full, tokens are discarded

Release packet only if sufficient tokens in the bucket
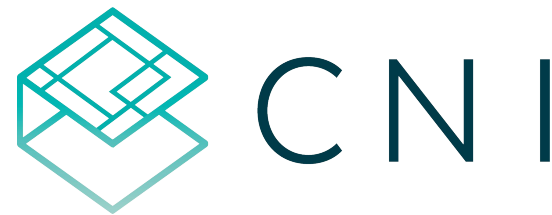
Packet queue

Output

4

# Problem Statement

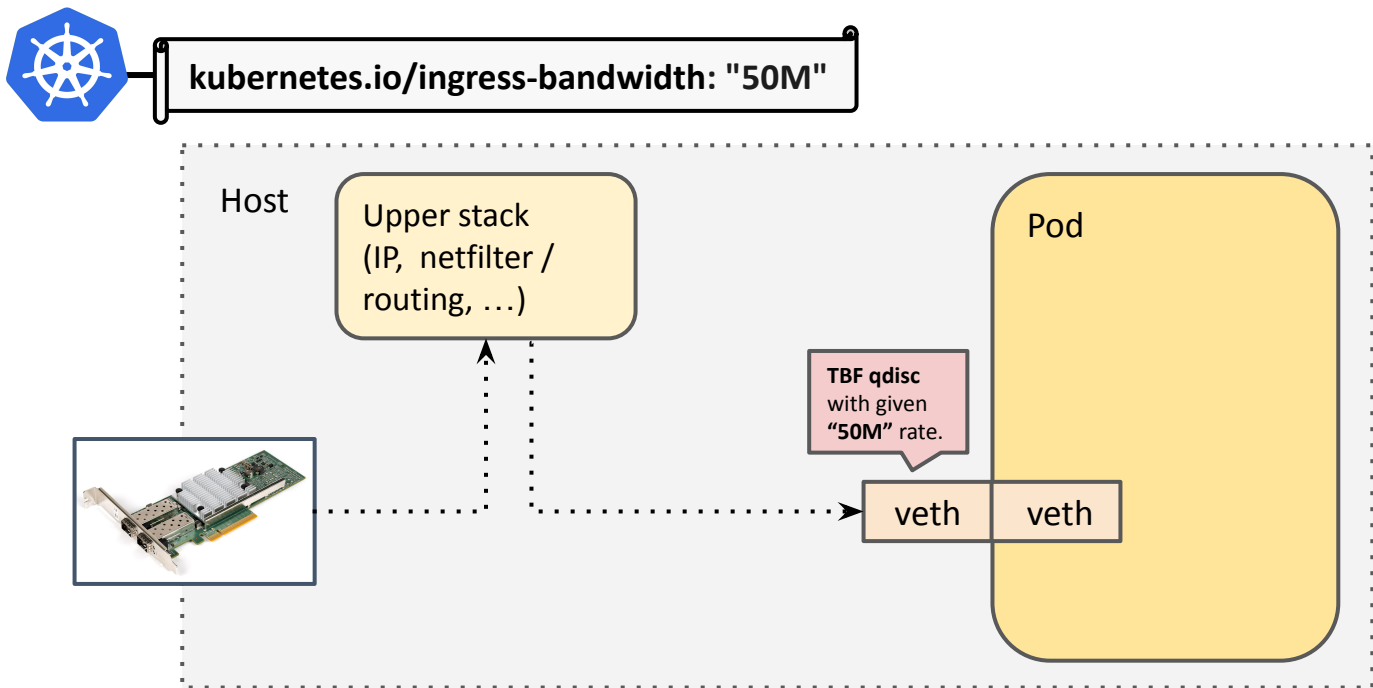➔ bandwidth meta plugin not scalable for production use
➔ TBFs are attached to the Pod's veth devices

# Problem Statement

➔ bandwidth meta plugin ingress example:

# Problem Statement

➔ bandwidth meta plugin ingress example:

**kubernetes.io/ingress-bandwidth:** "50M"

Host

Upper stack
(IP, netfilter /
routing, …)

Pod

**TBF qdisc**
with given
**"50M"** rate.

veth    veth

Client

Client

…

Client

TBF

# Problem Statement



➜ <u>bandwidth meta plugin</u> ingress example:

**kubernetes.io/ingress-bandwidth: "50M"**

Host

**Design issues:**

- Single lock contention point across all CPUs (qdisc root lock)

- Completely defeats phys NIC's multiqueue capability

- Queueing on RX is no-go, hogs up host's and also network's resources (bufferbloat)

TBF qdisc with given **"50M"** rate.

Pod

veth    veth

Client

Client

...

Client

TBF

# Problem Statement

➔ bandwidth meta plugin egress example:

# Problem Statement

➔ <u>bandwidth meta plugin</u> egress example:

**kubernetes.io/egress-bandwidth: "50M"**

Host

Upper stack
(IP, netfilter /
routing, …)

**TBF qdisc**
with given
**"50M"** rate.

All Pod traffic
**redirected** to ifb
device, so that
ifb can throttle
on egress hook.

Pod

ifb

veth

veth

TBF

App

App

…

App

# Problem Statement

➔ <u>bandwidth meta plugin</u> egress example:

**kubernetes.io/egress-bandwidth: "50M"**

**Design issues:**

- Now double queueing for veth *and* phys device (bufferbloat)

- Single lock contention point across all CPUs (qdisc root lock)

- Fooling TCP Small Queues (TSQ) by dropping socket association too early (bufferbloat)

- 2→3 netdevices per Pod needed

**TBF qdisc** with given **"50M"** rate.

All Pod traffic **redirected** to ifb device, so that ifb can throttle on egress.

Pod

ifb ← veth | veth

TBF

App

App
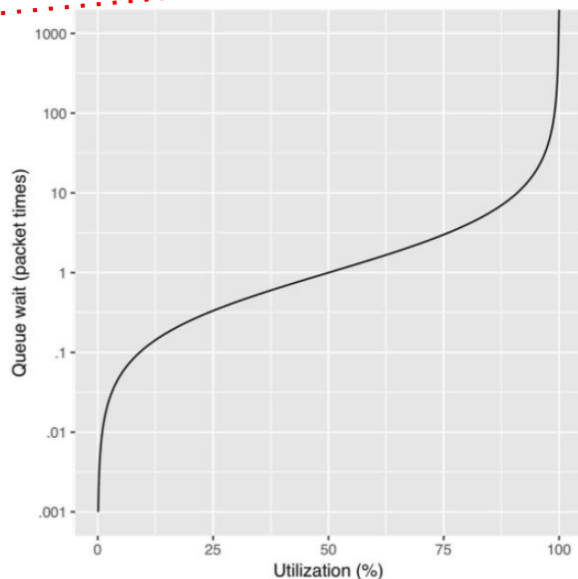
…
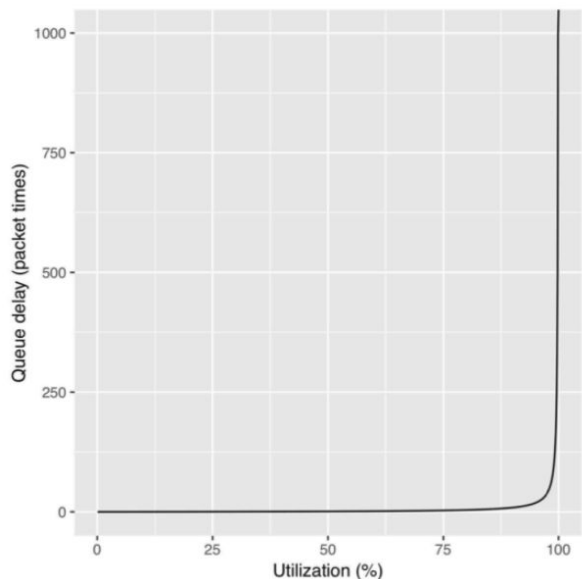
App

11

# Problem Statement

➔ bandwidth meta plugin aka "latency killer"

➔ tl;dr summary:

# Problem Statement

## ... but AFAP makes bottleneck run at 100%

Queuing theory says this is fragile. E.g., for M/D/1:
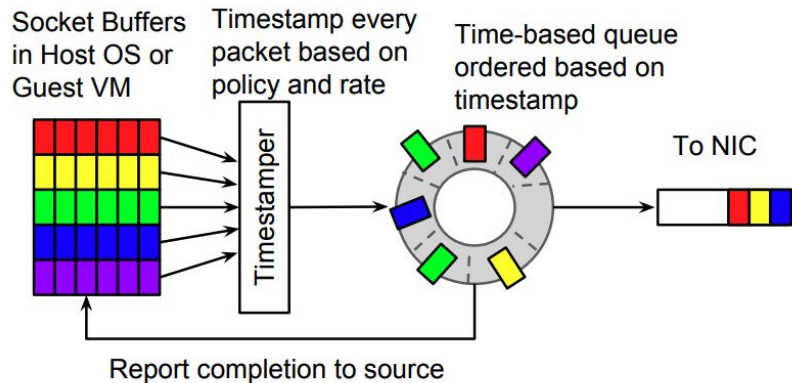
**Packet wait-time in queue skyrockets when bottleneck link gets utilized close to 100%.**

# From Queues to EDT model

Core Idea to replace queues with two simple pieces:

➜ Earliest Departure Time (EDT) time stamp in every packet
➜ Timing-wheel scheduler which replaces the queue



Source: A. Saeed et al., "Carousel: Scalable Traffic Shaping at End Hosts" (SIGCOMM'17)

# How can the EDT model be applied to Kubernetes?



**Programmable** and **performant** in-kernel "virtual machine" that **safely** executes native code on certain events/hooks (aka "JavaScript for the kernel").

# How can the EDT model be applied to Kubernetes?

**eBPF-based** CNI / platform which provides Pod connectivity, service load-balancing, network policies, **bandwidth management**, transparent encryption and more.

# Cilium's Bandwidth Manager

➔ Implements lock-less, EDT-based Pod rate-limiting with eBPF
➔ Enforcement points on phys devices instead of veths to avoid bufferbloat
   and improve TCP TSQ feedback

Host

Upper stack
(IP, netfilter /
routing, …)

Pod

**MQ + FQ** qdiscs attached here;
packet's socket association from
Pod still retained until here

FIB lookup and redirection
directly out of BPF.

veth | veth

eBPF

eBPF

Pod

kubernetes.io/egress-bandwidth: "50M"

eth0

eBPF

MQ

FQ    FQ    FQ    FQ

Q1    Q4

Packet (skb) departure
timestamp management

Multi-queue aware
Packet scheduler

NIC queues

18

Pod

**kubernetes.io/egress-bandwidth: "50M"**

Cilium agent reads Pod's bandwidth annotations, and pushes them into Cilium's BPF datapath.

eth0

eBPF

Packet (skb) departure timestamp management

MQ

Cilium's BPF datapath classifies Pod aggregate traffic and sets packet timestamps based on bandwidth policy.

FQ     FQ     FQ     FQ

Multi-queue aware Packet scheduler

Q1     Q4

NIC queues

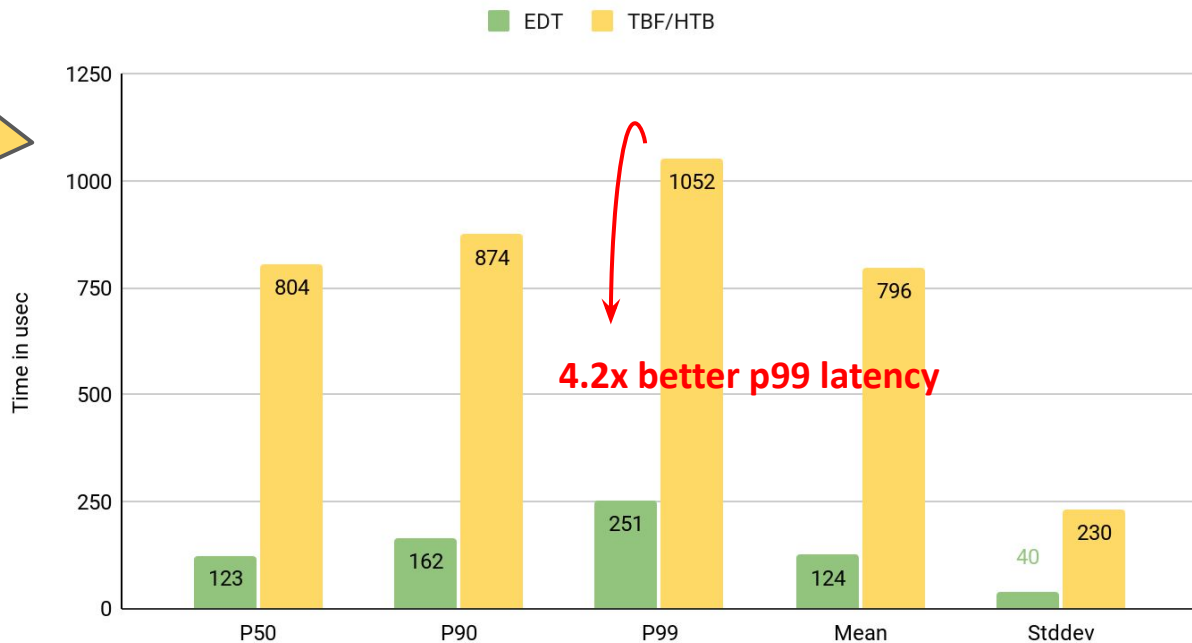Cilium agent sets up MQ/FQ Qdisc combination for all physical devices.

# Comparison of Cilium's EDT implementation vs TBF

Single flow latency for EDT and HTB/TBF model (lower is better)

EDT    TBF/HTB

Env: 256 concurrent request/response type flows (TCP_RR), 100M rate per flow

4.2x better p99 latency

Time in usec

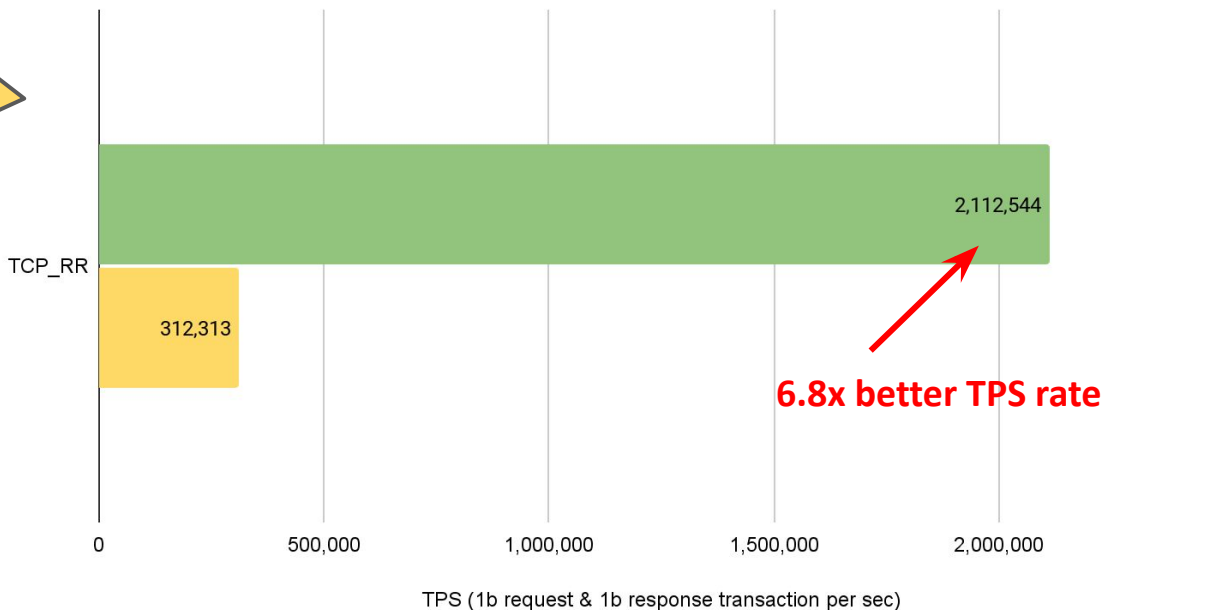| | P50 | P90 | P99 | Mean | Stddev |
|---|---|---|---|---|---|
| EDT | 123 | 162 | 251 | 124 | 40 |
| TBF/HTB | 804 | 874 | 1052 | 796 | 230 |

# Comparison of Cilium's EDT implementation vs TBF

Total transaction rate between EDT and HTB/TBF model (higher is better)

■ EDT    ■ TBF/HTB

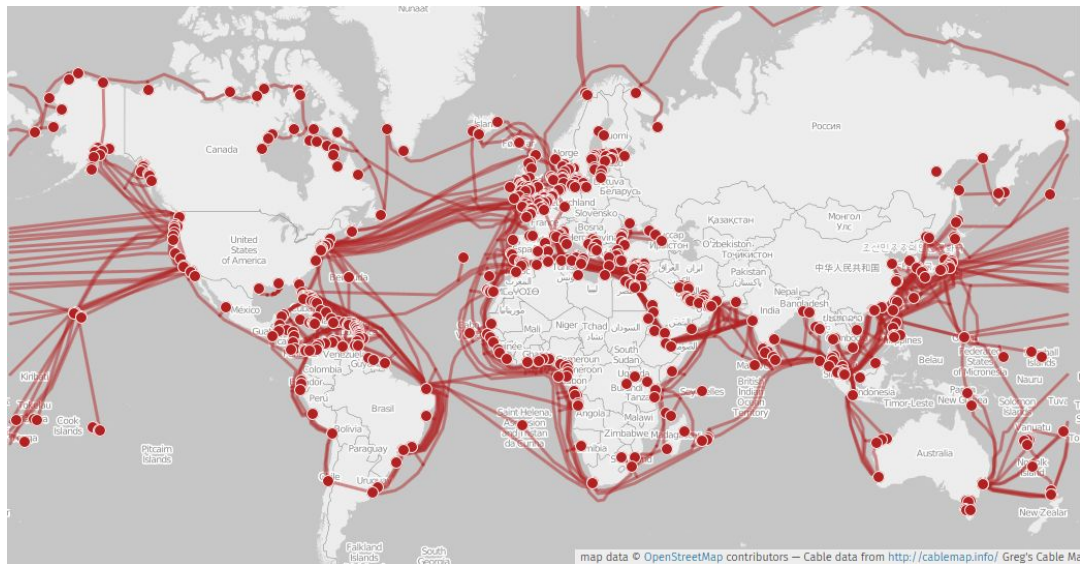**Env: 256 concurrent request/response type flows (TCP_RR), 100M rate per flow**

TCP_RR

2,112,544

312,313

**6.8x better TPS rate**

0    500,000    1,000,000    1,500,000    2,000,000

TPS (1b request & 1b response transaction per sec)

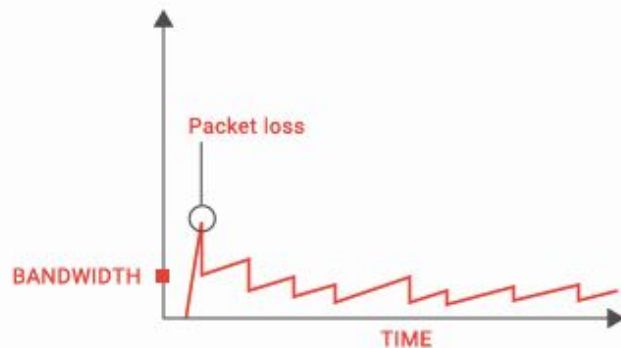# Cilium's Bandwidth Manager: recap for now

➔ So far: Cilium's EDT approach allows for scalable bandwidth enforcement

➔ What about more broadly Internet-level bandwidth management?

# What else does EDT model enable? Enter: BBR

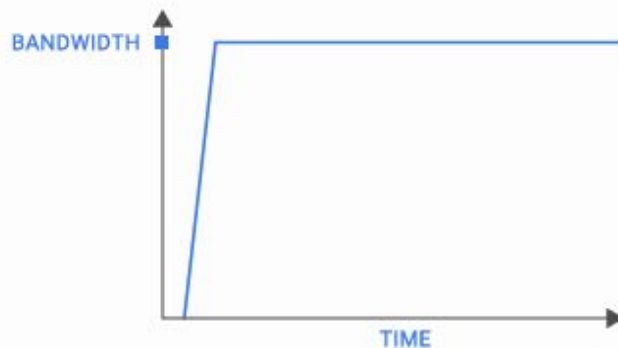# What else does EDT model enable? Enter: BBR



## TCP before BBR

Today's Internet is not moving data as well as it should. TCP sends data at lower bandwidth because the 1980s-era algorithm assumes that packet loss means network congestion.

Packet loss as signal to slow down.

## TCP BBR

BBR models the network to send as fast as the available bandwidth and is 2700x faster than previous TCPs on a 10Gb, 100ms link with 1% loss. BBR powers google.com, youtube.com, and apps using Google Cloud Platform services.
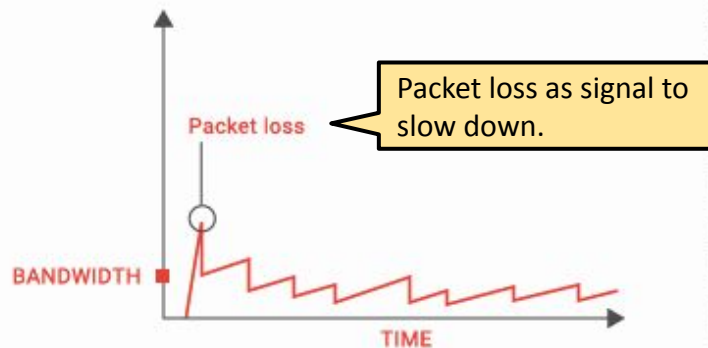
Models bandwidth based on measurements around delivery rate and round trip time.

Source: https://cloud.google.com/blog/products/networking/tcp-bbr-congestion-control-comes-to-gcp-your-internet-just-got-faster
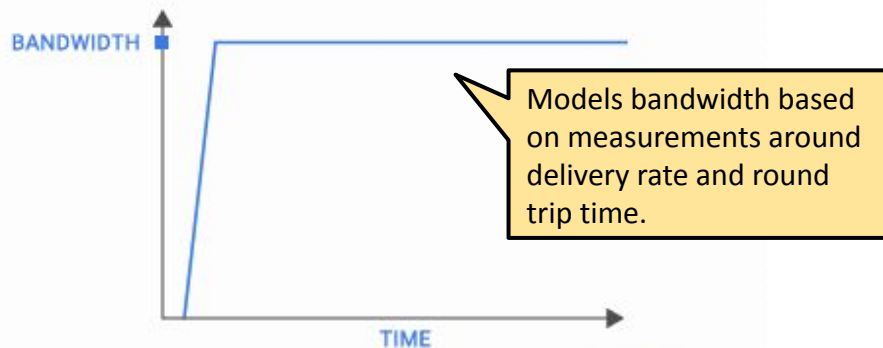
# When is it useful to consider BBR?

➜ Kubernetes cluster exposing services to clients over the Internet
   ◆ Significant latency improvements for low-end last-mile networks
   ◆ Significant throughput improvements for high-speed long-haul links

| | CUBIC (default) | BBR (v1) |
|---|---|---|
| Model parameters to the state machine | N/A | Throughput, RTT |
| Loss | Reduce cwnd by 30% on window with any loss | N/A |
| ECN | RFC3168 (classic ECN) | N/A |
| Startup | Slow-start until RTT rises (Hystart) or any loss | Slow-start until throughput plateaus |

Source: Alexey Ivanov, "Evaluating BBRv2 On the Edge"

# Example: **New York** (packet.net) -> **Zurich**

```
darkstar@linux:~/trees/bpf-...  ×    root@zh-lab-node-1: ~    ×    darkstar@linux:~/trees/bpf  ×    root@ny-c3-small-x86-01: ~  ×    root

root@zh-lab-node-1:~# while [ 1 ]; do iperf3 -c 147.75.66.15 -t 55 -R -i 5 -O 5 ; done
Connecting to host 147.75.66.15, port 5201
Reverse mode, remote host 147.75.66.15 is sending
[  5] local 192.168.178.91 port 52148 connected to 147.75.66.15 port 5201
[ ID] Interval           Transfer     Bitrate
[  5]   0.00-5.00   sec  78.4 MBytes   132 Mbits/sec                  (omitted)
[  5]   0.00-5.00   sec  98.2 MBytes   165 Mbits/sec
[  5]   5.00-10.00  sec  99.1 MBytes   166 Mbits/sec
[  5]  10.00-15.00  sec   113 MBytes   189 Mbits/sec
[  5]  15.00-20.00  sec   159 MBytes   267 Mbits/sec
[  5]  20.00-25.00  sec   257 MBytes   431 Mbits/sec
[  5]  25.00-30.00  sec   153 MBytes   256 Mbits/sec
[  5]  30.00-35.00  sec   146 MBytes   245 Mbits/sec
[  5]  35.00-40.00  sec   148 MBytes   248 Mbits/sec
[  5]  40.00-45.00  sec   157 MBytes   264 Mbits/sec
[  5]  45.00-50.00  sec   193 MBytes   323 Mbits/sec
[  5]  50.00-55.00  sec   272 MBytes   457 Mbits/sec
- - - - - - - - - - - - - - - - - - - - - - - - -
[ ID] Interval           Transfer     Bitrate         Retr
[  5]   0.00-55.10  sec  1.76 GBytes   274 Mbits/sec   1501            sender
[  5]   0.00-55.00  sec  1.75 GBytes   274 Mbits/sec                   receiver
```

Bandwidth probing, overreaction to loss!
(Sawtooth pattern nicely visible)

**Default**, server runs:

- TCP CUBIC
- fq_codel Qdisc

26

# Example: **New York** (packet.net) -> **Zurich**

```
darkstar@linux:~/trees/bpf-...   ×       root@zh-lab-node-1: ~     ×      darkstar@linux:~/trees/bpf  ×        root@ny-c3-small-x86-01: ~   ×      root

root@zh-lab-node-1:~# while [ 1 ]; do iperf3 -c 147.75.66.15 -t 55 -R -i 5 -O 5 ; done
Connecting to host 147.75.66.15, port 5201
Reverse mode, remote host 147.75.66.15 is sending
[  5] local 192.168.178.91 port 52254 connected to 147.75.66.15 port 5201
[ ID] Interval           Transfer     Bitrate
[  5]   0.00-5.00   sec   152 MBytes   254 Mbits/sec                  (omitted)
[  5]   0.00-5.00   sec   258 MBytes   433 Mbits/sec
[  5]   5.00-10.00  sec   240 MBytes   403 Mbits/sec
[  5]  10.00-15.00  sec   255 MBytes   427 Mbits/sec
[  5]  15.00-20.00  sec   247 MBytes   414 Mbits/sec
[  5]  20.00-25.00  sec   255 MBytes   428 Mbits/sec
[  5]  25.00-30.00  sec   255 MBytes   428 Mbits/sec
[  5]  30.00-35.00  sec   238 MBytes   400 Mbits/sec
[  5]  35.00-40.00  sec   255 MBytes   428 Mbits/sec
[  5]  40.00-45.00  sec   239 MBytes   401 Mbits/sec
[  5]  45.00-50.00  sec   253 MBytes   425 Mbits/sec
[  5]  50.00-55.00  sec   242 MBytes   407 Mbits/sec
- - - - - - - - - - - - - - - - - - - - - - - - -
[ ID] Interval           Transfer     Bitrate        Retr
[  5]   0.00-55.10  sec  2.68 GBytes   418 Mbits/sec  58812          sender
[  5]   0.00-55.00  sec  2.67 GBytes   418 Mbits/sec                 receiver
```

No overreaction, stable.

**Updated**, server runs:

- TCP BBR
- FQ Qdisc (for EDT)

(results also reproduce with netperf)

# Example: **New York** (packet.net) -> **Zurich**

```
darkstar@linux:~/trees/bpf-...   ×     root@zh-lab-node-1: ~   ×     darkstar@linux:~/trees/bpf   ×     root@ny-c3-small-x86-01: ~   ×     root

root@zh-lab-node-1:~# while [ 1 ]; do iperf3 -c 147.75.66.15 -t 55 -R -i 5 -O 5 ; done
Connecting to host 147.75.66.15, port 5201
Reverse mode, remote host 147.75.66.15 is sending
[  5] local 192.168.178.91 port 52254 connected to 147.75.66.15 port 5201
[ ID] Interval           Transfer     Bitrate
[  5]   0.00-5.00   sec   152 MBytes   254 Mbits/sec                    (omitted)
[  5]   0.00-5.00   sec   258 MBytes   433 Mbits/sec
[  5]   5.00-10.00  sec   240 MBytes   403 Mbits/sec
[  5]  10.00-15.00  sec   255 MBytes   427 Mbits/sec
[  5]  15.00-20.00  sec   247 MBytes   414 Mbits/sec
[  5]  20.00-25.00  sec   255 MBytes   428 Mbits/sec
[  5]  25.00-30.00  sec   255 MBytes   428 Mbits/sec
[  5]  30.00-35.00  sec   238 MBytes   400 Mbits/sec
[  5]  35.00-40.00  sec   255 MBytes   428 Mbits/sec
[  5]  40.00-45.00  sec   239 MBytes   401 Mbits/sec
[  5]  45.00-50.00  sec   253 MBytes   425 Mbits/sec
[  5]  50.00-55.00  sec   242 MBytes   407 Mbits/sec
- - - - - - - - - - - - - - - - - - - - - - - - - - -
[ ID] Interval           Transfer     Bitrate         Retr
[  5]   0.00-55.10  sec   2.68 GBytes  418 Mbits/sec   58812
[  5]   0.00-55.00  sec   2.67 GBytes  418 Mbits/sec
```

No overreaction, stable.

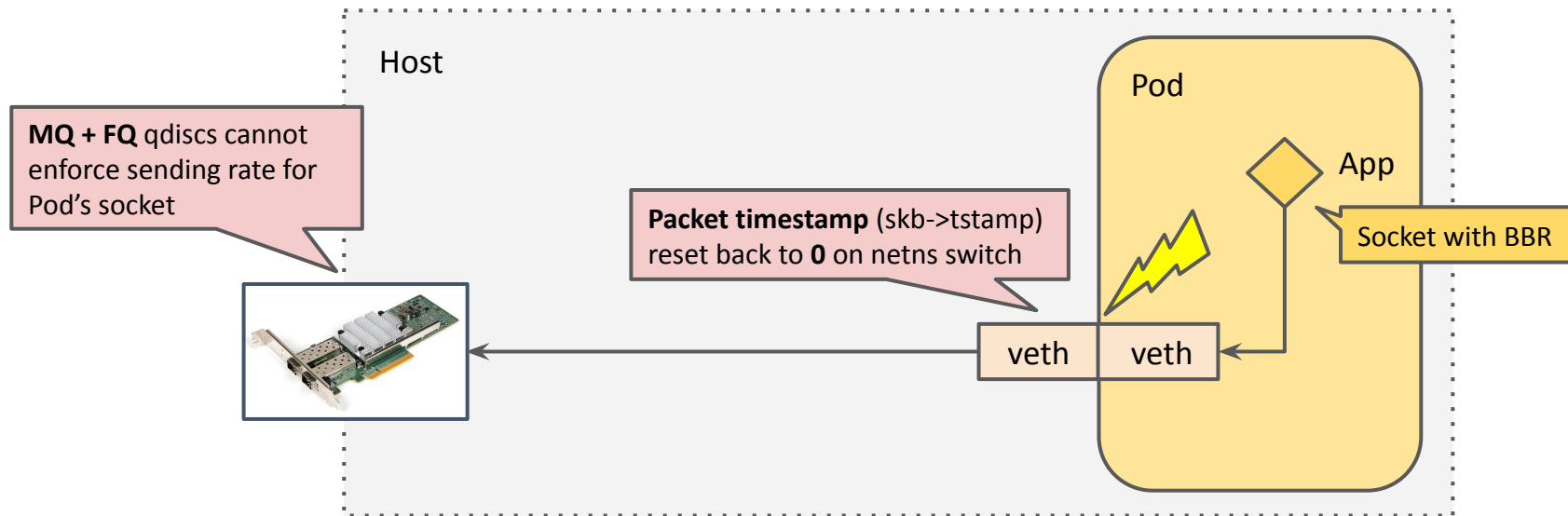**Updated**, server runs:

- TCP BBR
- FQ Qdisc (for EDT)

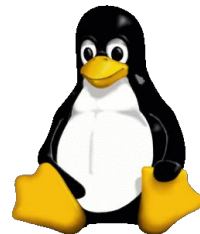**CUBIC → BBR** bumps:

**274 → 418 Mbit/s**

28

# Can BBR be used for Kubernetes Pods?

➔ BBR works in conjunction with FQ and sets packet delivery timestamps

➔ Kernel clears timestamp for packets leaving Pods (== netns)

  ◆ Usage of BBR for Pods not possible/broken in general today



Host

Pod

App

**MQ + FQ** qdiscs cannot enforce sending rate for Pod's socket

**Packet timestamp** (skb->tstamp) reset back to **0** on netns switch

Socket with BBR

veth | veth
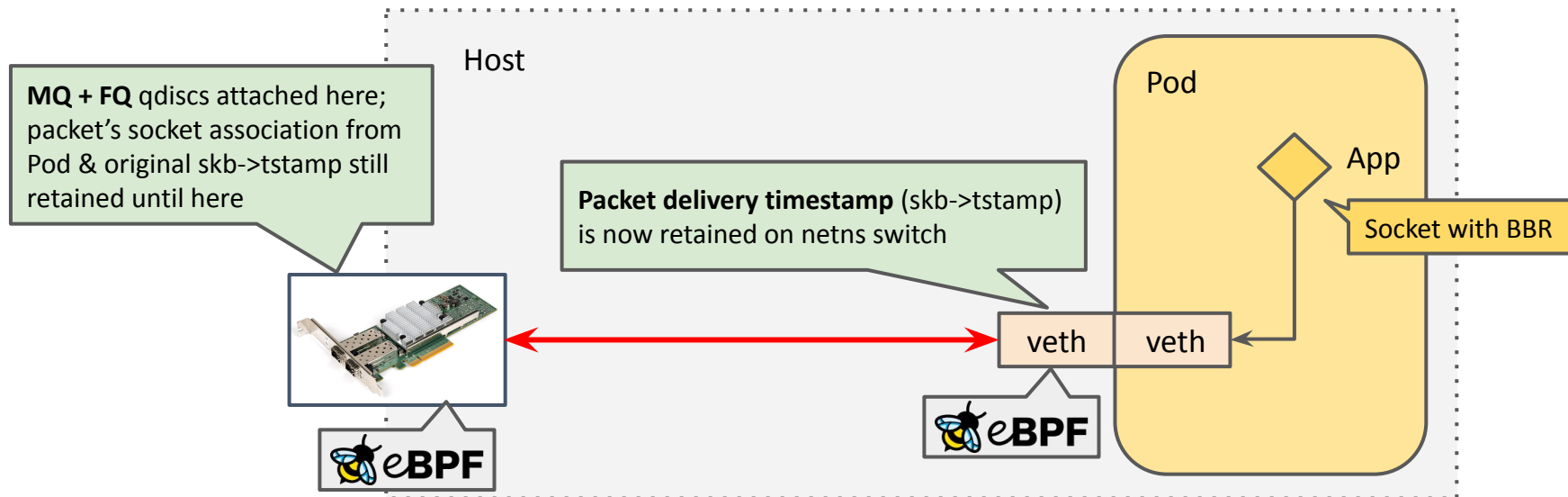
# Rationale on today's timestamp reset ([see our LPC talk](#))

Kernel uses different clock bases for skb->tstamp:

➔ Ingress is CLOCK_TAI, egress is CLOCK_MONOTONIC (as is FQ)
➔ Forwarding from RX to TX would cause drop in FQ due to overreaching FQ's drop horizon given different clock's offsets
➔ No means to figure out clock base from skb->tstamp, hence reset

# Cilium's Bandwidth Manager: BBR

➔ We helped fixing networking stack in Linux v5.18+ to retain timestamps

➔ Bandwidth Manager plumbs the appropriate underlying infrastructure

◆ Receives new knob for switching whole cluster over to BBR by default

**MQ + FQ** qdiscs attached here; packet's socket association from Pod & original skb->tstamp still retained until here

**Packet delivery timestamp** (skb->tstamp) is now retained on netns switch

Host

Pod

App

Socket with BBR

veth    veth
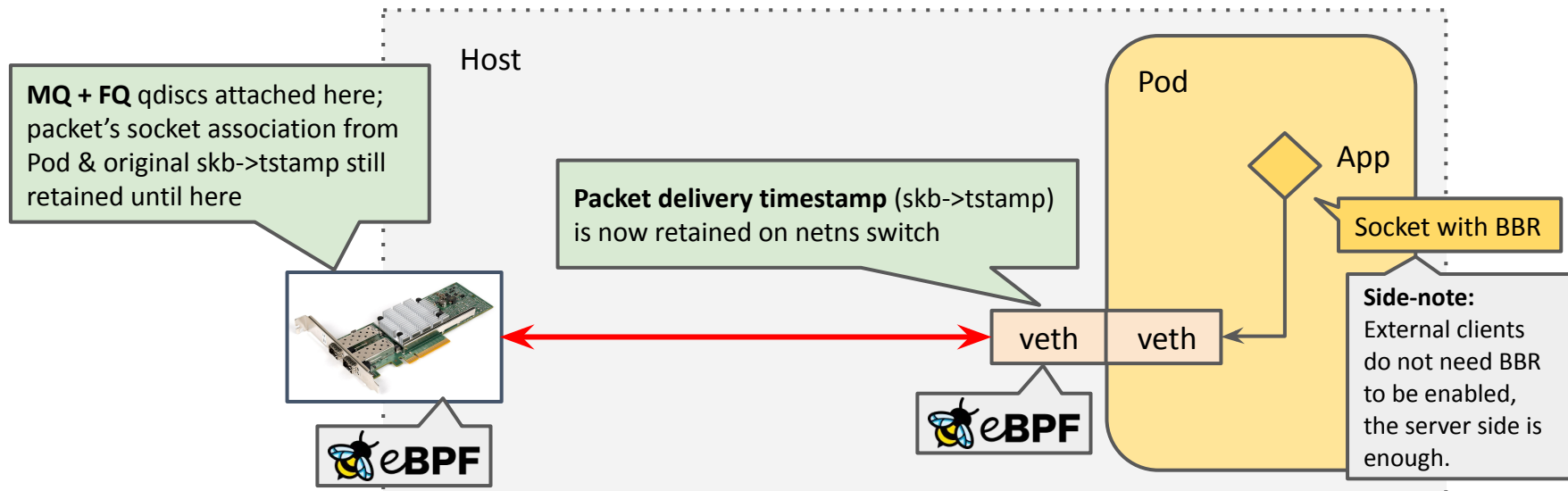
eBPF

eBPF

31

# Cilium's Bandwidth Manager: BBR

➔ We helped fixing networking stack in Linux v5.18+ to retain timestamps

➔ Bandwidth Manager plumbs the appropriate underlying infrastructure

   ◆ Receives new knob for switching whole cluster over to BBR by default



**MQ + FQ** qdiscs attached here; packet's socket association from Pod & original skb->tstamp still retained until here

**Packet delivery timestamp** (skb->tstamp) is now retained on netns switch

Host

Pod

App

Socket with BBR

**Side-note:** External clients do not need BBR to be enabled, the server side is enough.
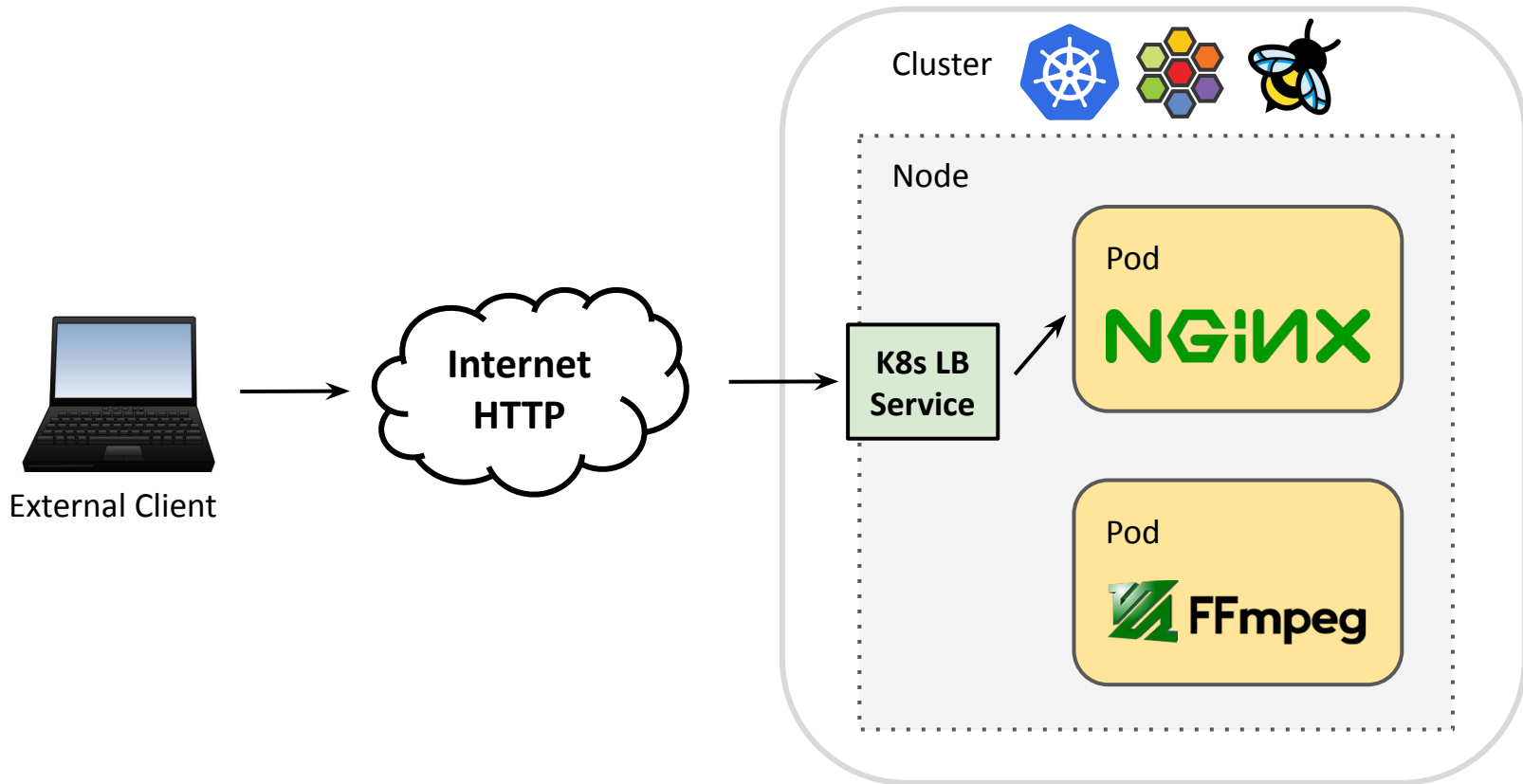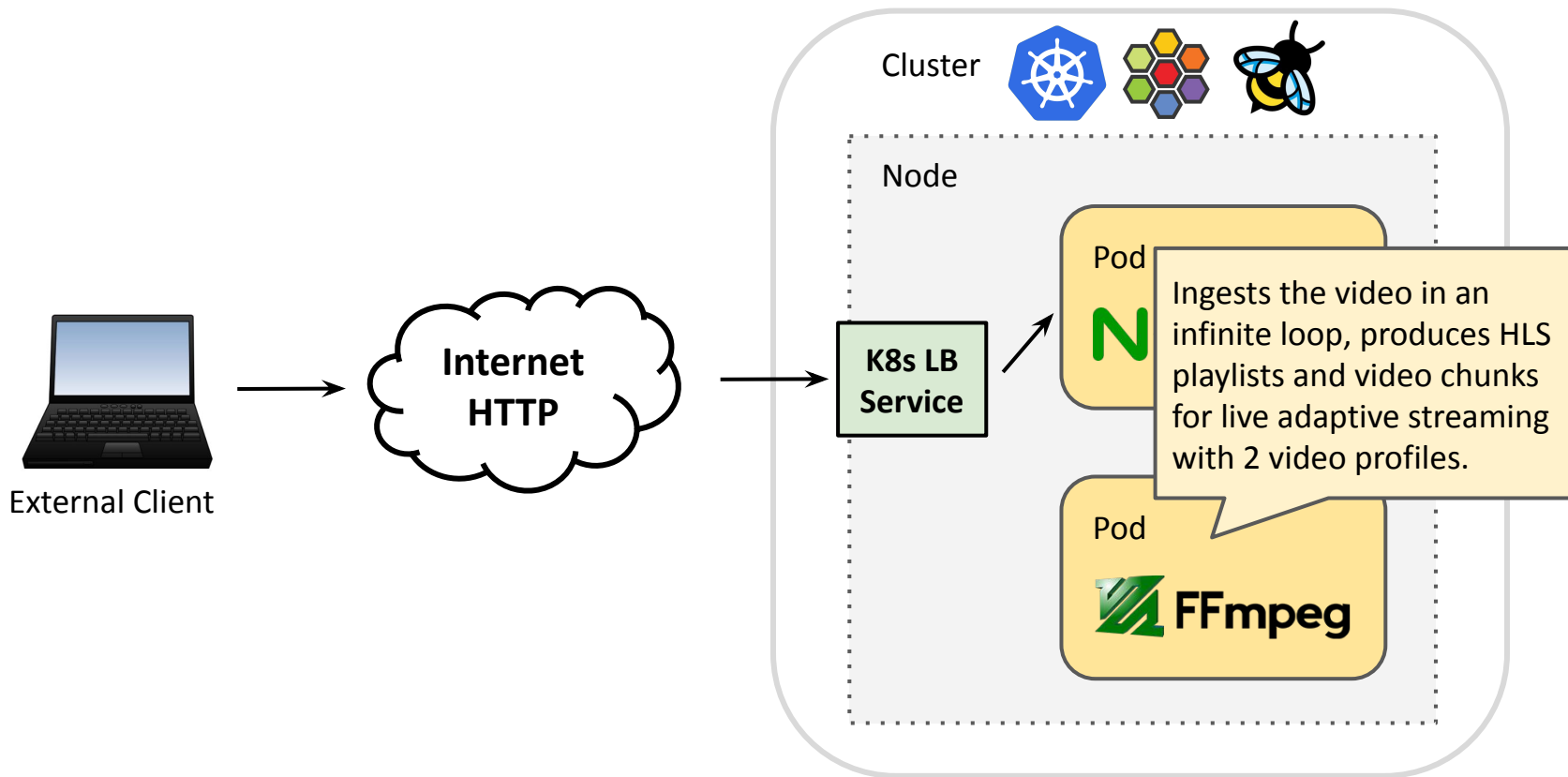
veth | veth

ISOVALENT
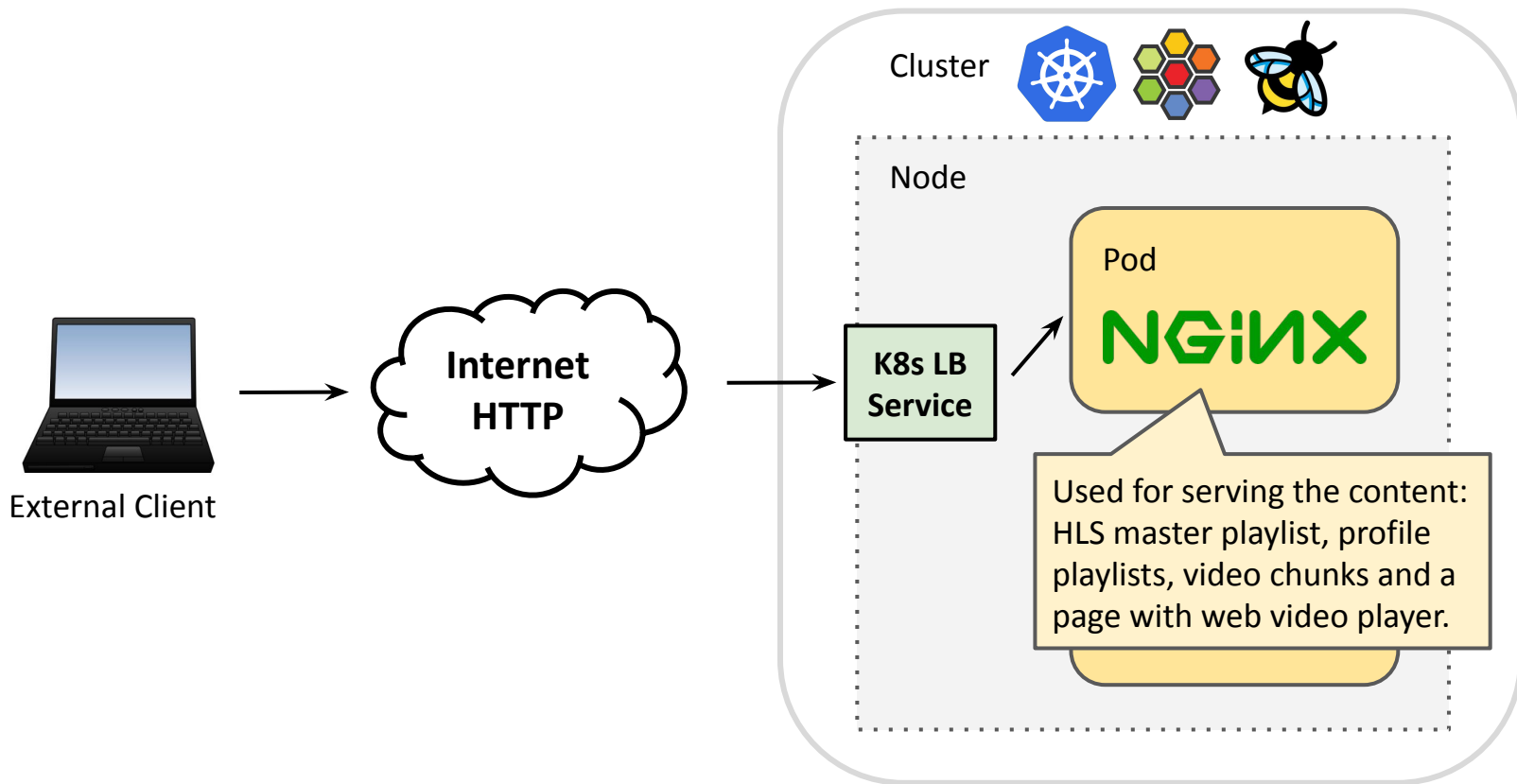
# Demo:  BBR for Pods

(K8s/Cilium-backed video streaming service: CUBIC versus BBR)

# Video Streaming Service Demo Setup: CUBIC vs BBR

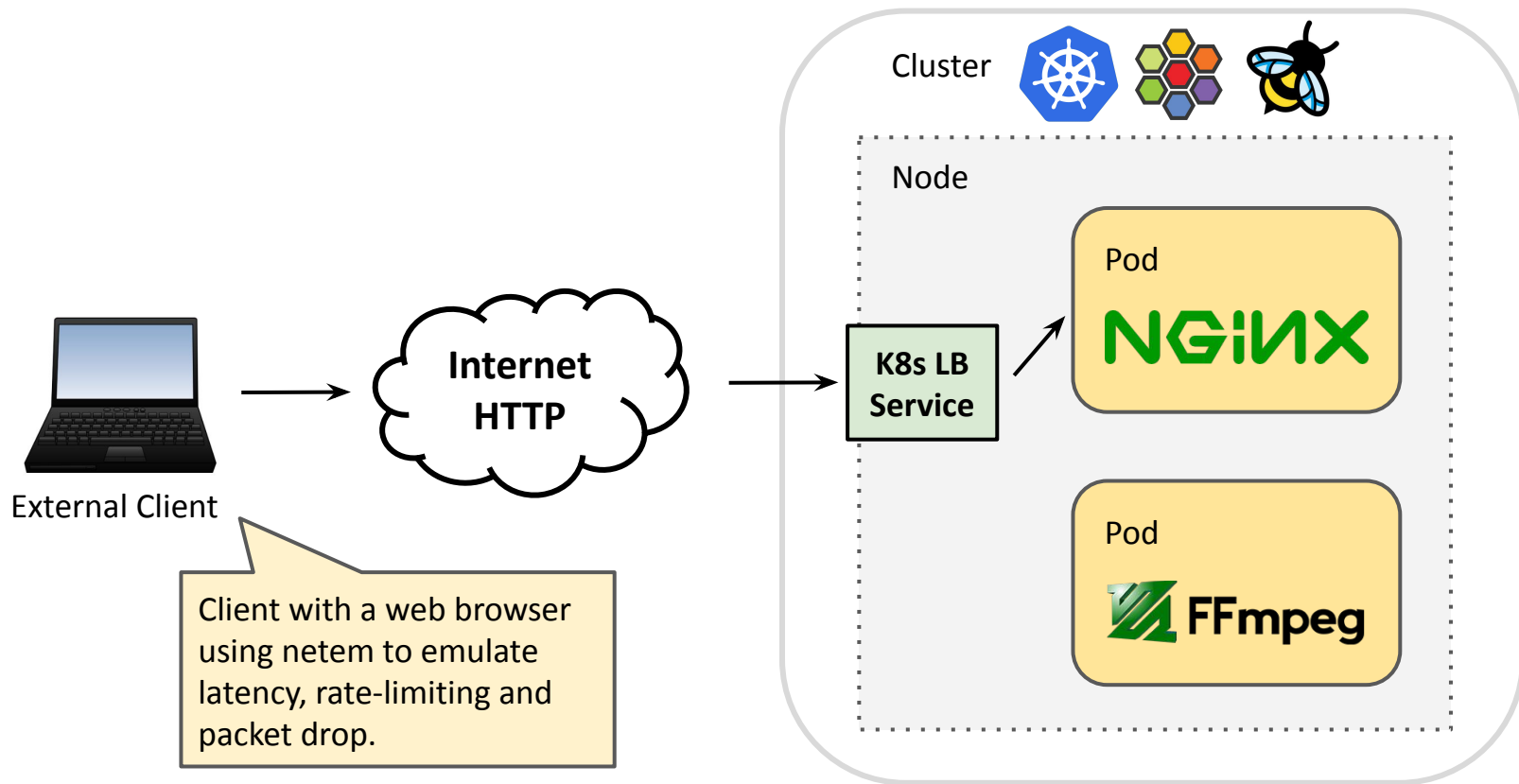# Video Streaming Service Demo Setup: CUBIC vs BBR



Cluster

Node

External Client

Internet HTTP

K8s LB Service

Pod

Pod

Ingests the video in an infinite loop, produces HLS playlists and video chunks for live adaptive streaming with 2 video profiles.

# Video Streaming Service Demo Setup: CUBIC vs BBR



Cluster

Node

Pod

NGINX

External Client

Internet
HTTP

K8s LB
Service

Used for serving the content:
HLS master playlist, profile
playlists, video chunks and a
page with web video player.

# Video Streaming Service Demo Setup: CUBIC vs BBR



Client with a web browser using netem to emulate latency, rate-limiting and packet drop.

# Video Streaming Service Demo Setup: CUBIC vs BBR



Kubernetes Cluster with Cilium in different 2 configs:

```
BandwidthManager:          EDT with BPF [CUBIC] [enp5s0]
BandwidthManager:          EDT with BPF [BBR] [enp5s0]
```

ISOVALENT

Now to the stream ...

# What needs to be considered with use of BBR?

➔ BBR has potential unfairness issues towards CUBIC when env uses both
➔ BBR will trigger a higher TCP retransmission rate (more aggressive probing)
➔ BBRv2 in the works to overcome them



Source: Alexey Ivanov, "Evaluating BBRv2 On the Edge"

# Revisiting earlier Problem Statement

➔ Kubernetes bandwidth enforcement does not need to be in a poor state

➔ Native implementation via Cilium's Bandwidth Manager (GA since v1.12)

   ◆ Efficient, eBPF-based bandwidth enforcement via EDT model

   ◆ First CNI to support BBR (& socket pacing) for Pods

   ◆ Side-note: Realizing such architecture only possible with eBPF

**Getting Started Guide for Bandwidth Manager:**

```
helm upgrade cilium ./cilium \
    --namespace kube-system \
    --reuse-values \
    --set bandwidthManager.enabled=true \
    --set bandwidthManager.bbr=true
kubectl -n kube-system rollout restart ds/cilium
```

(needs Linux kernel v5.1+)

(needs Linux kernel v5.18+)

# Acknowledgements

- ➜   Van Jacobson
- ➜   Eric Dumazet
- ➜   Vytautas Valancius
- ➜   Stanislav Fomichev
- ➜   Martin Lau
- ➜   John Fastabend
- ➜   Cilium, BPF & netdev kernel community

ISOVALENT

# Thank you! Questions?

github.com/cilium/cilium
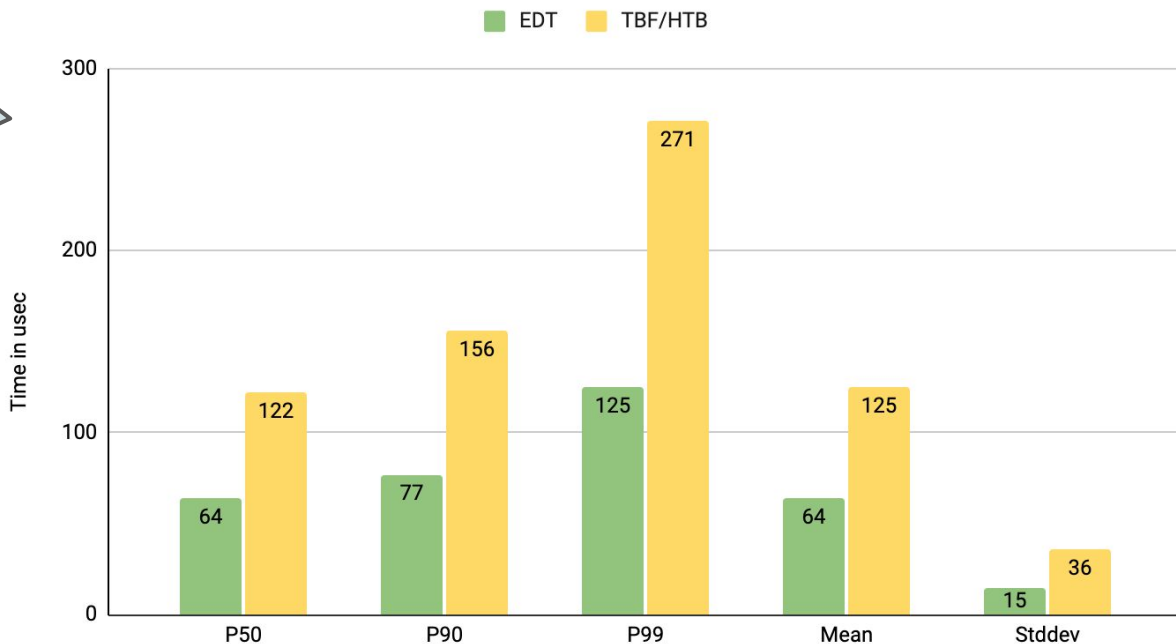
cilium.io

ebpf.io

Isovalent:    booth S21

Cilium:        kiosk 12 (@ project pavilion)

ISOVALENT

# Appendix: Latency & TPS for 64/128/256 flows

# Comparison of Cilium's EDT implementation vs TBF

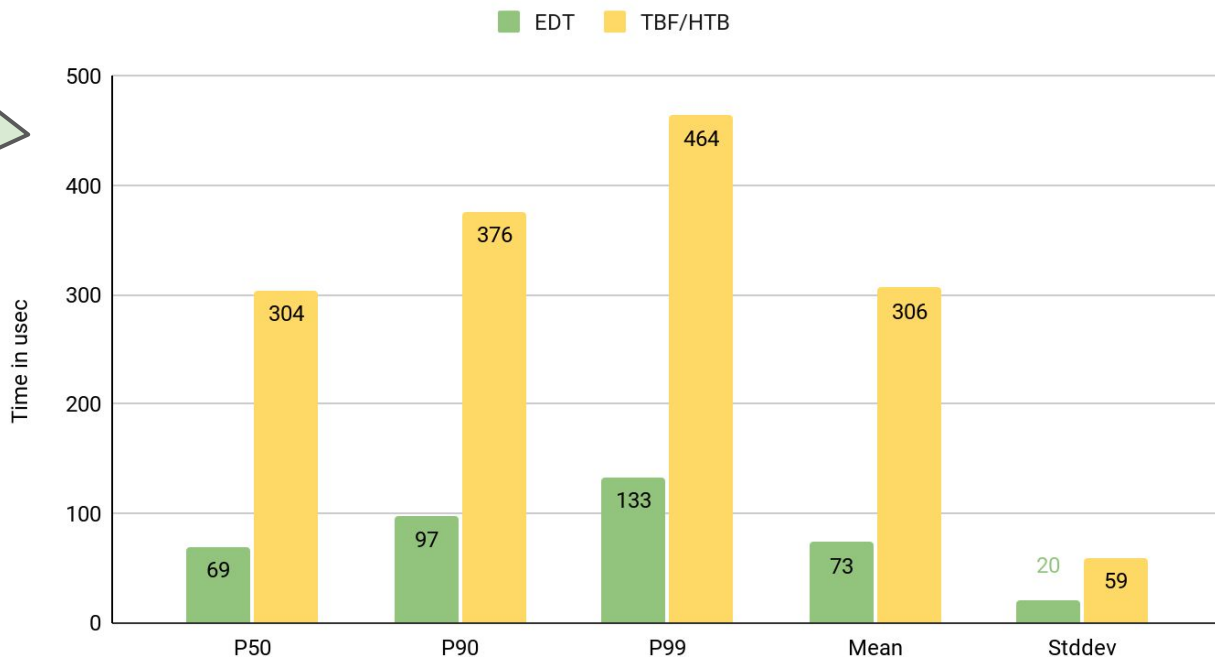Single flow latency for EDT and HTB/TBF model (lower is better)

Env: 64 concurrent request/response type flows (TCP_RR), 100M rate per flow

# Comparison of Cilium's EDT implementation vs TBF

Single flow latency for EDT and HTB/TBF model (lower is better)

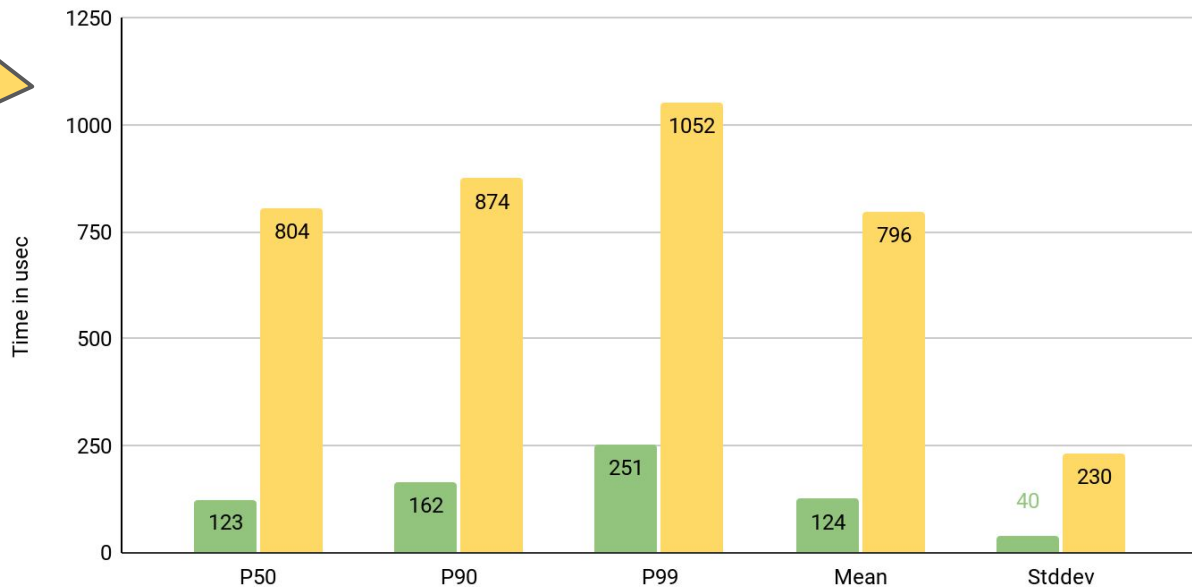Env: 128 concurrent request/response type flows (TCP_RR), 100M rate per flow

Legend: EDT, TBF/HTB

Time in usec

| | P50 | P90 | P99 | Mean | Stddev |
|---|---|---|---|---|---|
| EDT | 69 | 97 | 133 | 73 | 20 |
| TBF/HTB | 304 | 376 | 464 | 306 | 59 |

# Comparison of Cilium's EDT implementation vs TBF

Single flow latency for EDT and HTB/TBF model (lower is better)

**Env: 256 concurrent request/response type flows (TCP_RR), 100M rate per flow**

■ EDT    ■ TBF/HTB

| | P50 | P90 | P99 | Mean | Stddev |
|---|---|---|---|---|---|
| EDT | 123 | 162 | 251 | 124 | 40 |
| TBF/HTB | 804 | 874 | 1052 | 796 | 230 |

Time in usec

# Comparison of Cilium's EDT implementation vs TBF

Total transaction rate between EDT and HTB/TBF model (higher is better)

■ EDT   ■ TBF/HTB

**Env: 64 concurrent request/response type flows (TCP_RR), 100M rate per flow**

TCP_RR

1,012,258

498,633

0       250,000       500,000       750,000       1,000,000

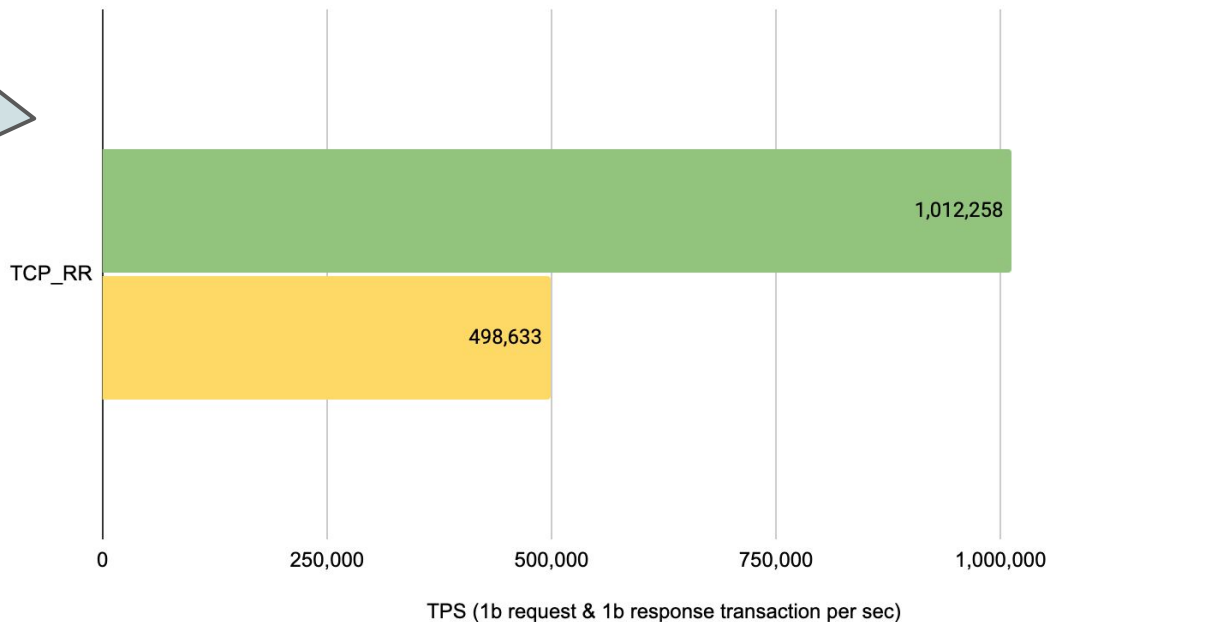TPS (1b request & 1b response transaction per sec)

# Comparison of Cilium's EDT implementation vs TBF

Total transaction rate between EDT and HTB/TBF model (higher is better)

■ EDT    ■ TBF/HTB

**Env: 128 concurrent request/response type flows (TCP_RR), 100M rate per flow**

TCP_RR

1,712,115

412,332

0          500,000        1,000,000       1,500,000

TPS (1b request & 1b response transaction per sec)

# Comparison of Cilium's EDT implementation vs TBF

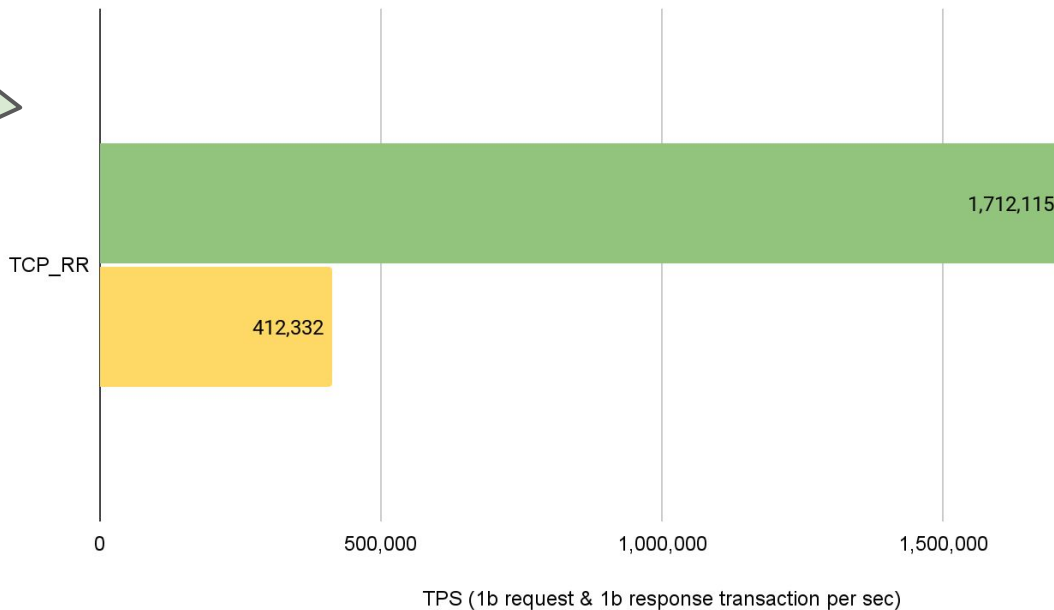Total transaction rate between EDT and HTB/TBF model (higher is better)

■ EDT  ■ TBF/HTB

Env: 256 concurrent request/response type flows (TCP_RR), 100M rate per flow

TCP_RR

2,112,544

312,313

0    500,000    1,000,000    1,500,000    2,000,000

TPS (1b request & 1b response transaction per sec)