# CNI-agnostic network performance accelerator with eBPF

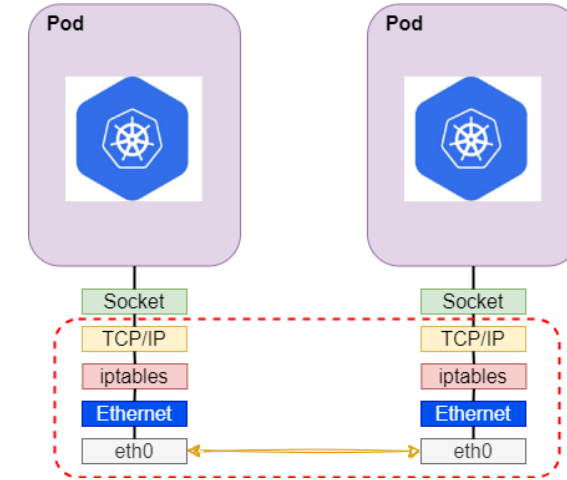*Yizhou Xu, Intel*

*Mengxin Liu, Alauda*

# Agenda

- TCP/IP stack overhead
- eBPF background knowledge
- How to bypass Tcp/Ip with eBPF
- Performance Analysis
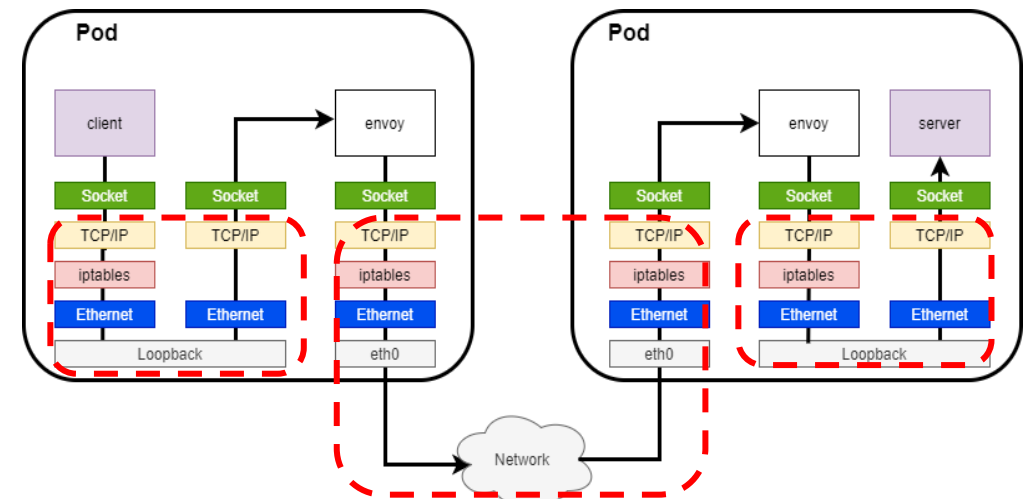- Practice on KubeVirt acceleration

# TCP/IP stack overhead

In Kubernetes, each pod has its own network stack, packet from one pod to another traverse whole stack multiple times

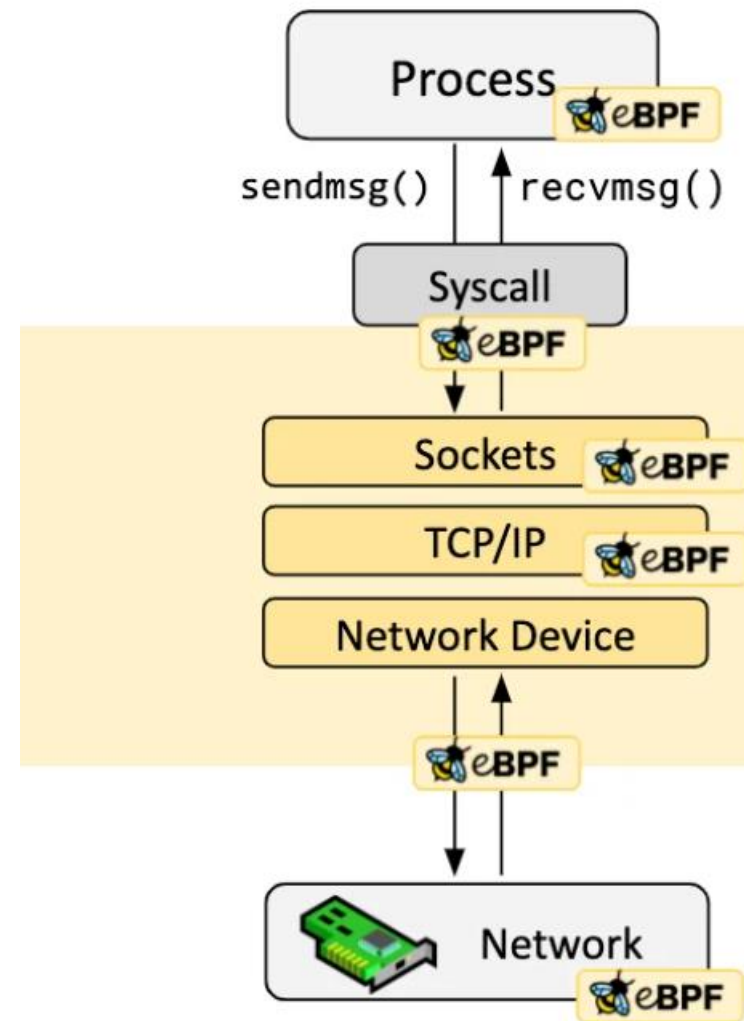In particular, service mesh amplify the overhead by sidecar mode



**Tcp/ip overhead in Pod to Pod (same host)**



**Tcp/ip overhead in Service Mesh (same host)**

# Why use eBPF

- Work in Kernel

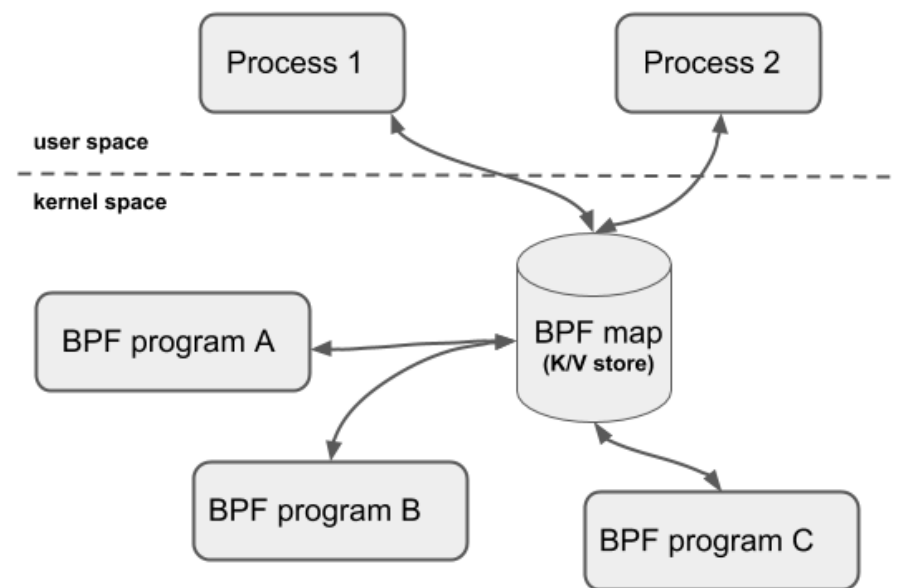- Non-intrusive

- CNI agnostic

- Safety and efficient

## MAP

provide generic data structure for user space
and kernel space communication

- o HASHMAP
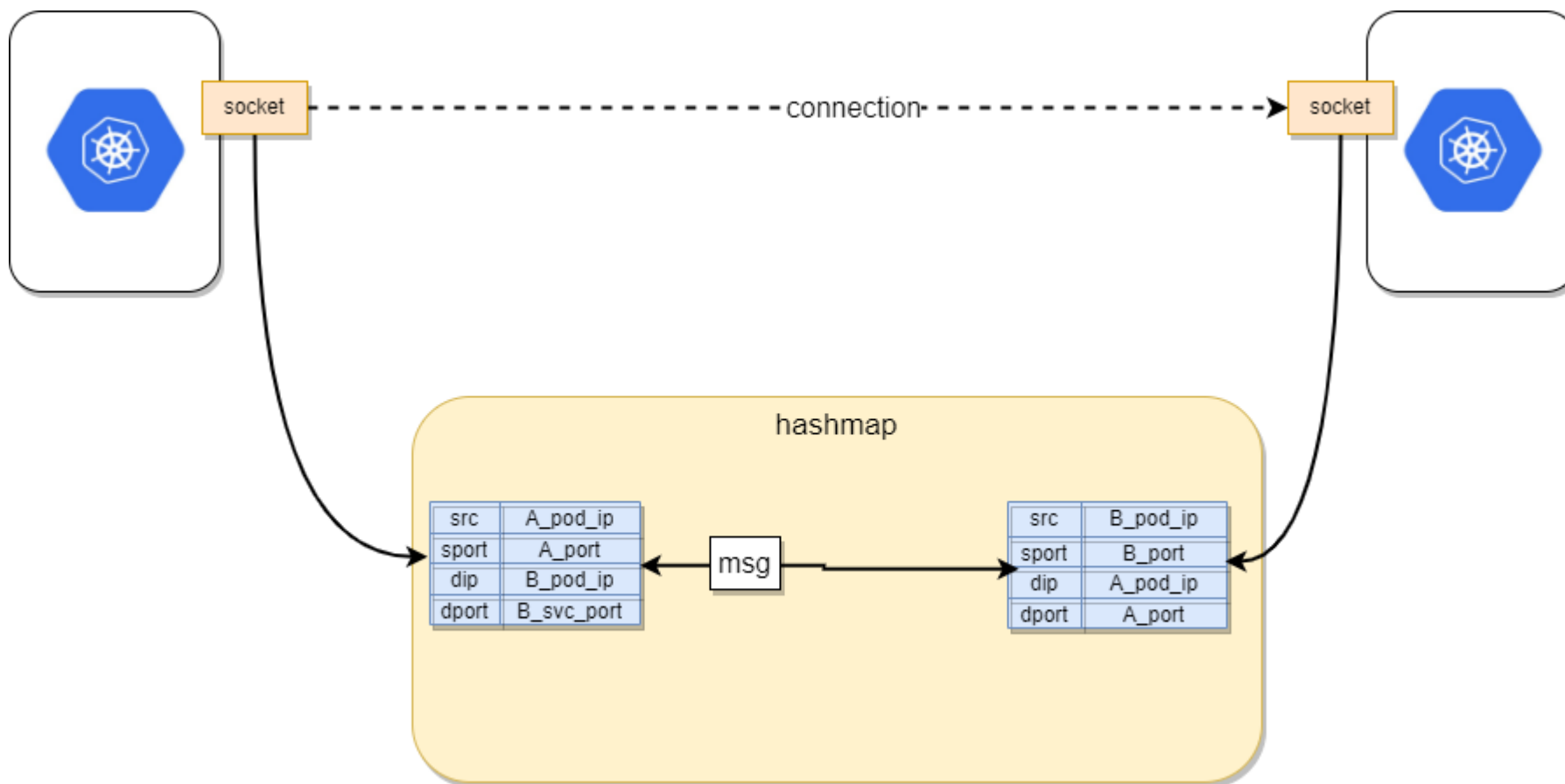- o SOCKHASH

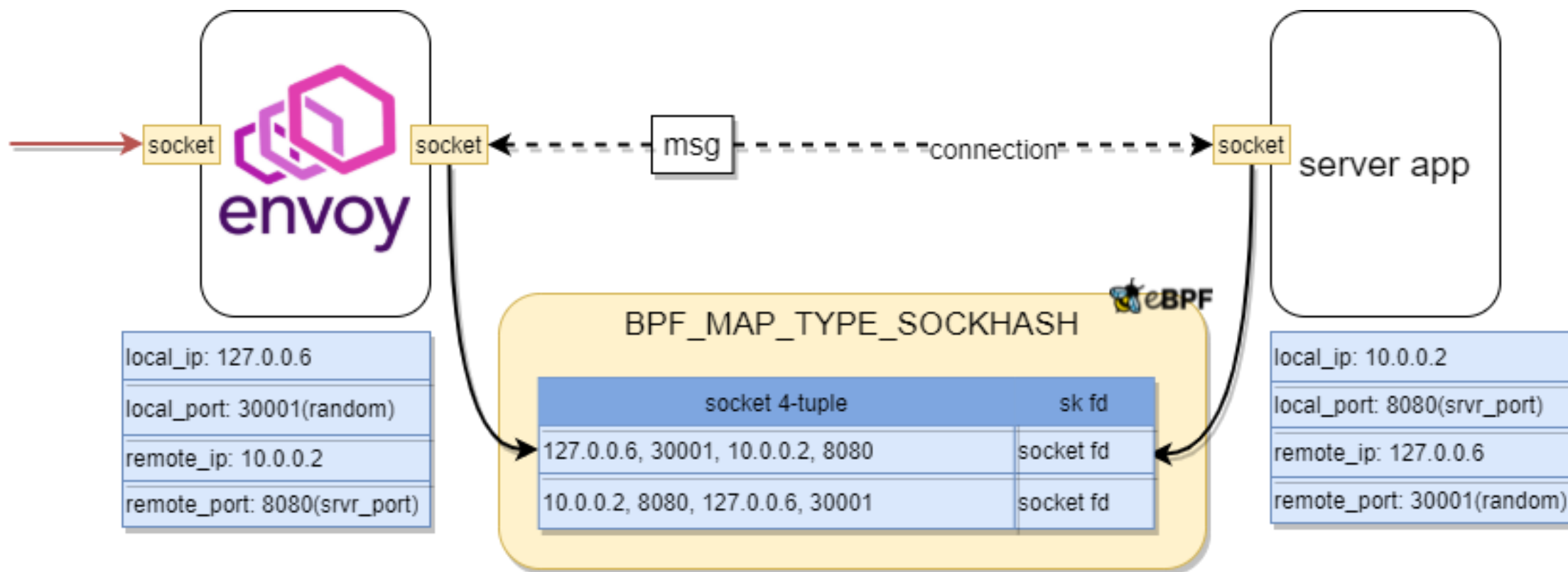## Program type

- o SOCK_OPS

- o SK_MSG

- sock_ops
  - capture socket in given status, populate into map
- sk_msg
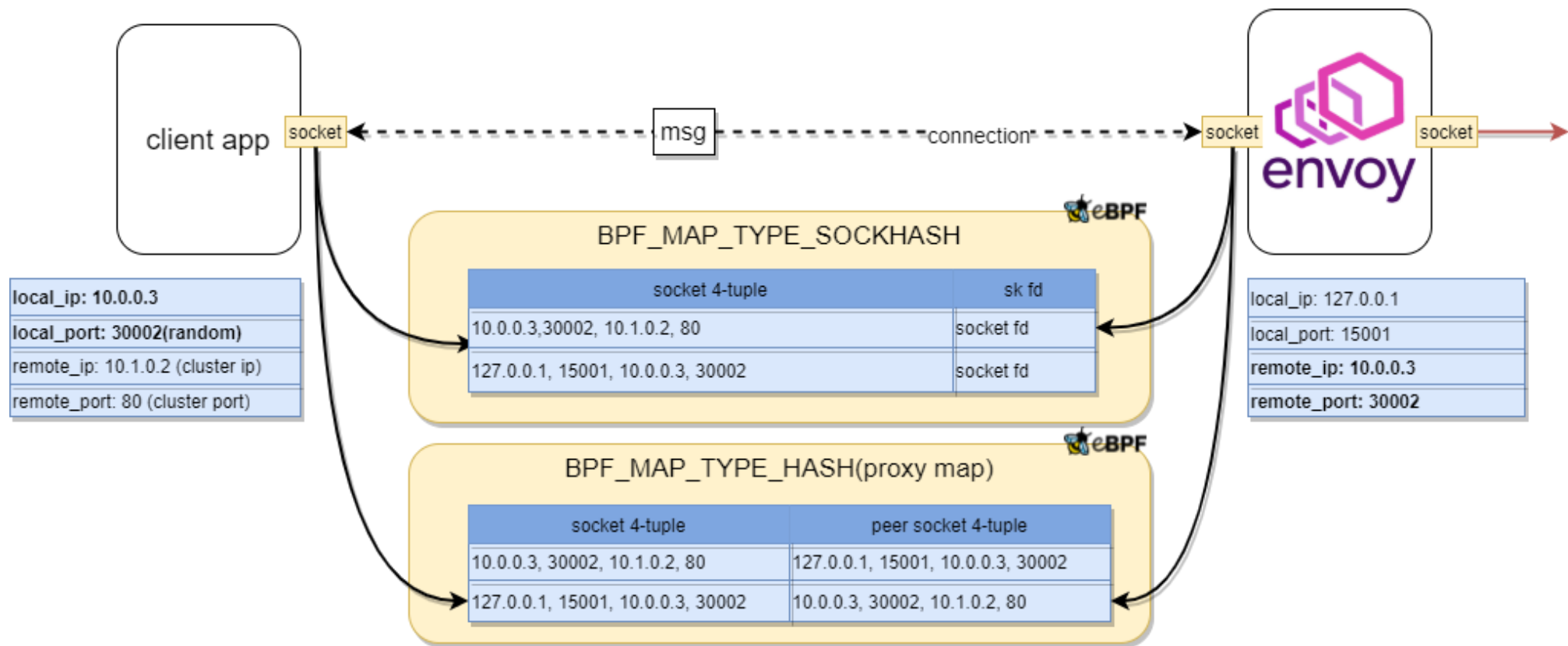  - when there is data in socket, lookup peer socket
  - transfer data to peer
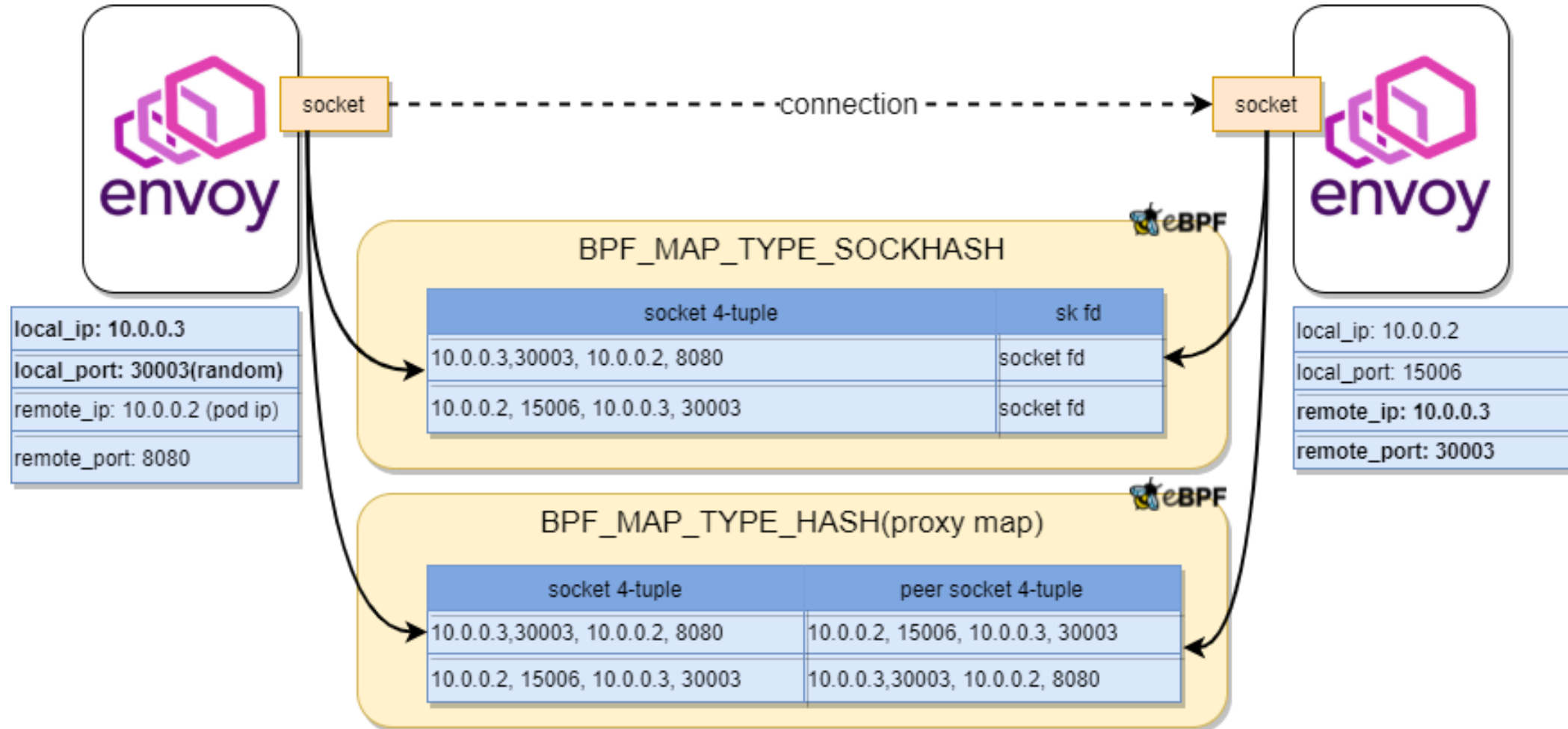
# Pod to Pod acceleration(Kubernetes)

# Inbound Acceleraion(Service Mesh)

# Outbound Acceleration(Service Mesh)

# Envoy to Envoy Acceleration (same host)

# It's opensource!



**REPO:** https://github.com/intel/istio-tcpip-bypass

# Benchmark

➢ Deploy two test Pods in the same Node

➢ Use qperf to benchmark TCP latency/bandwidth with packet sizes from 1byte to 16KB

➢ qperf -t 60 100.64.0.3 -ub -oo msg_size:1:16K:*4 -vu tcp_lat tcp_bw

➢ Compare with the benchmark when the optimization is disabled



TCP Latency — eBPF(us)

TCP Bandwidth Ratio — eBPF/Default

# Benchmark Analysis

➢ TCP latency will decrease by 40% ~ 60%

➢ Throughput will increase by 40% ~ 80% when the packet size is greater than 1024 bytes

➢ When the packet size is less than 512 bytes, the throughput decrease due to per packet processing overhead

# KubeVirt network performance

- ➢ Native Kubevirt Bridge networks have a much larger performance gap compared to Pods.

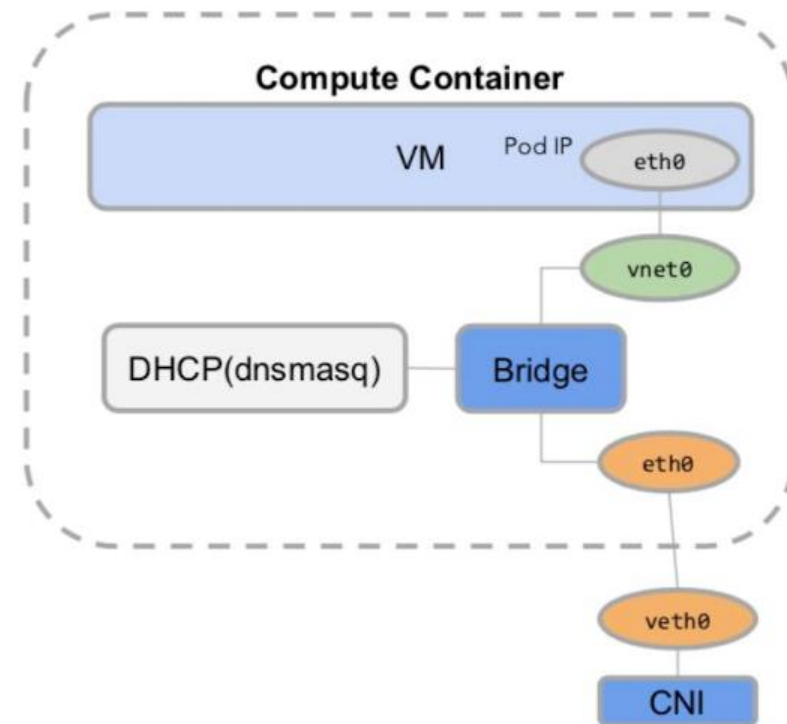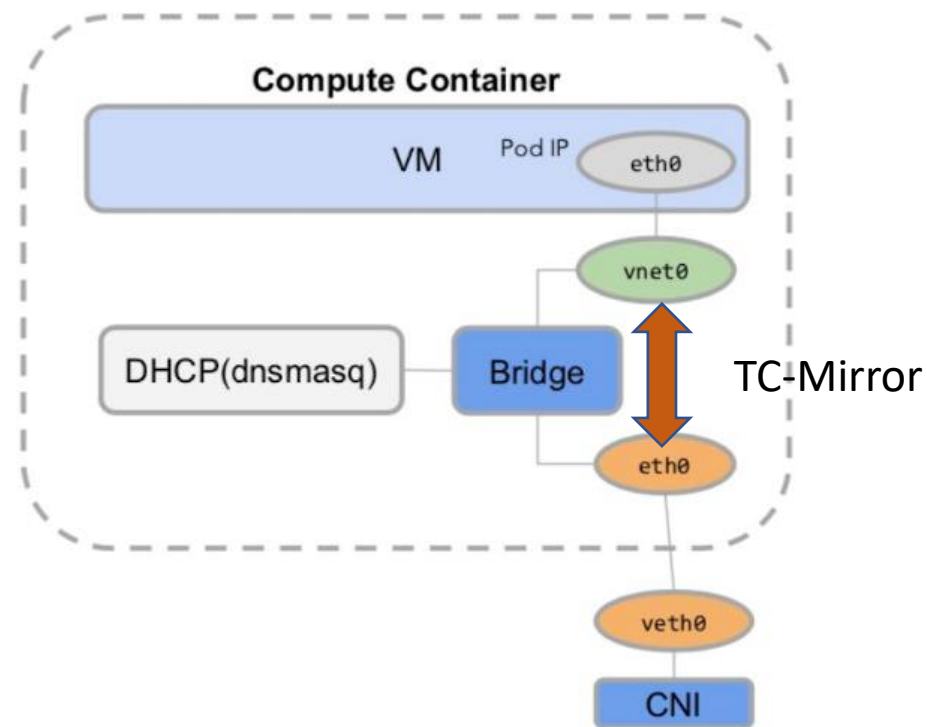  - ➢ Latency increased by 60%

  - ➢ PPS decreased by 50%.

- ➢ Possible causes:

  - ➢ The consumption brought by the virtualized Linux network stack itself

  - ➢ Additional losses caused by the multiple layers of KubeVirt Bridge network.

- ➢ Kata has already utilized tc-mirror to bypass the bridge
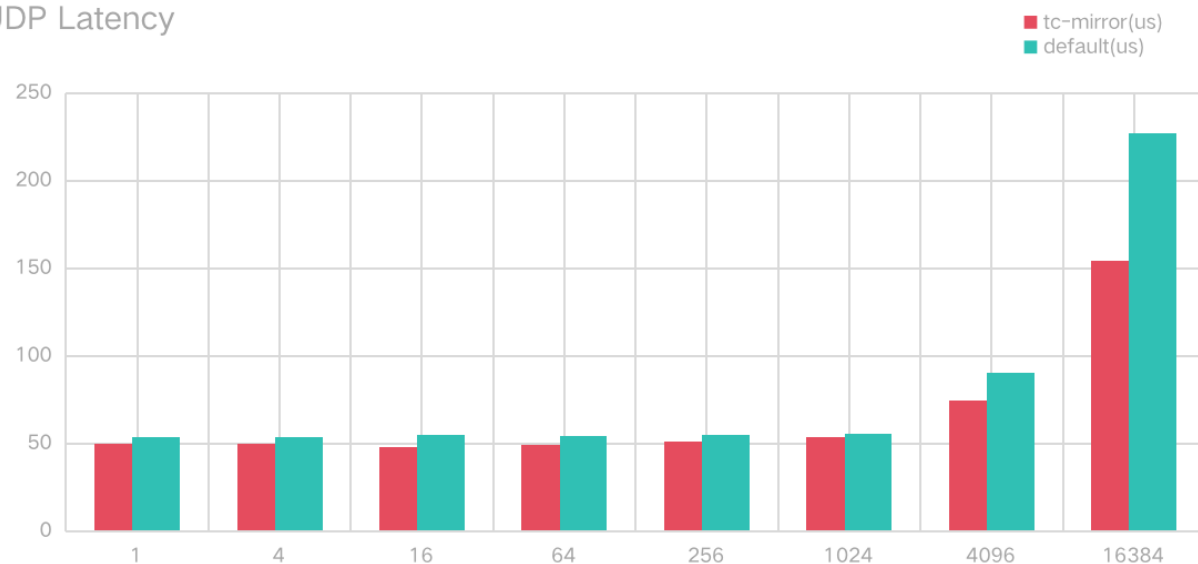
# KubeVirt network acceleration

- ➢ Both tc-mirror and bpf_redirect are tried

  - ➢ tc-mirror performed slightly better than bpf_redirect in all tests

  - ➢ tc-mirror do not require high version kernel

  - ➢ We didn't find way to redirect from VM eth0 to veth0 directly

  - ➢ Latency decrease 5% but throughput decrease about 20%

- ➢ Disable checksum

  - ➢ No need checksum for internal packet processing
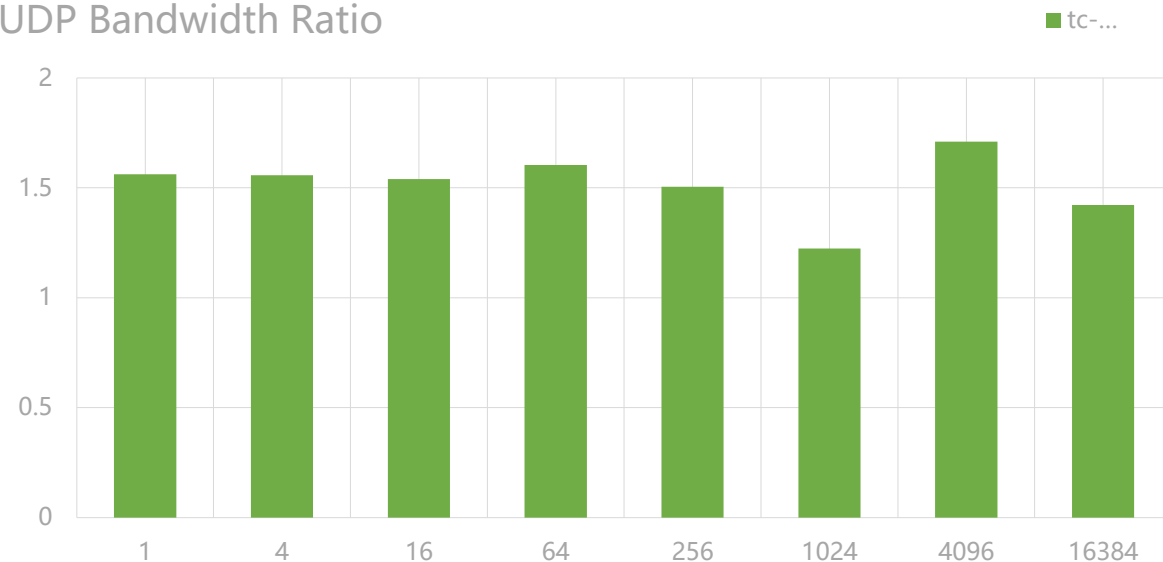
  - ➢ Both rx/tx are dissabled

# Benchmark

➢ Use qperf to benchmark UDP latency/bandwidth with packet sizes from 1byte to 16KB

    ➢ Latency 4%~30% decrease

    ➢ Throughput 20%~70% improvement

# Feature work

➢ Accelerate UDP inter-node traffic

➢ Fine-grained control the switch between eBPF datapath and kernel network datapath

➢ Accelerate Kubevirt VMs traffic within the same Node

# Notice and Disclaimer

**Intel technologies may require enabled hardware, software or service activation.**
**No product or component can be absolutely secure.**
**Your costs and results may vary.**