# Container Networking:
# The Play of BPF &
# Network NS with different Virtual Devices
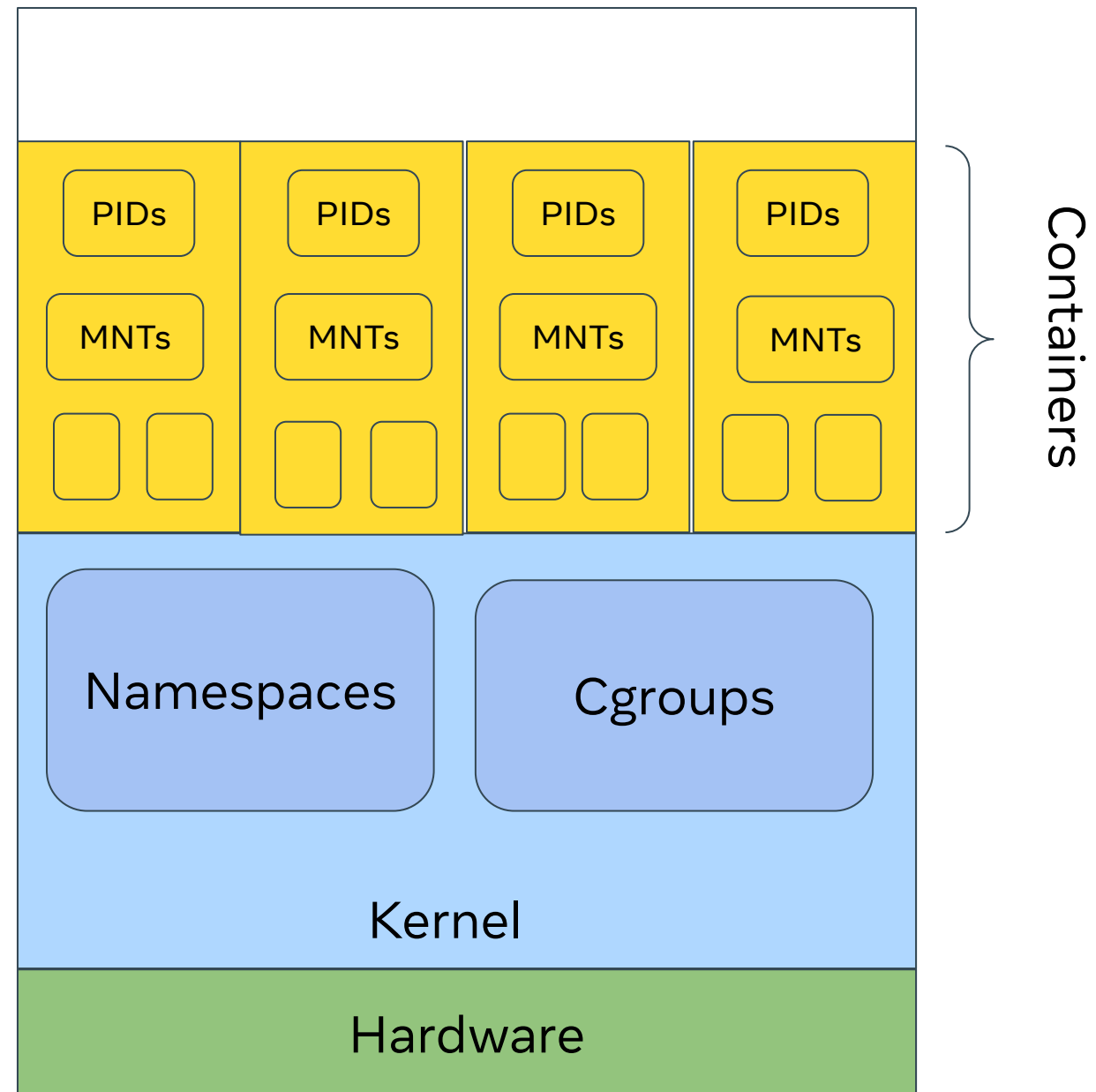
Takshak Chahande   &   Martin KaFai Lau
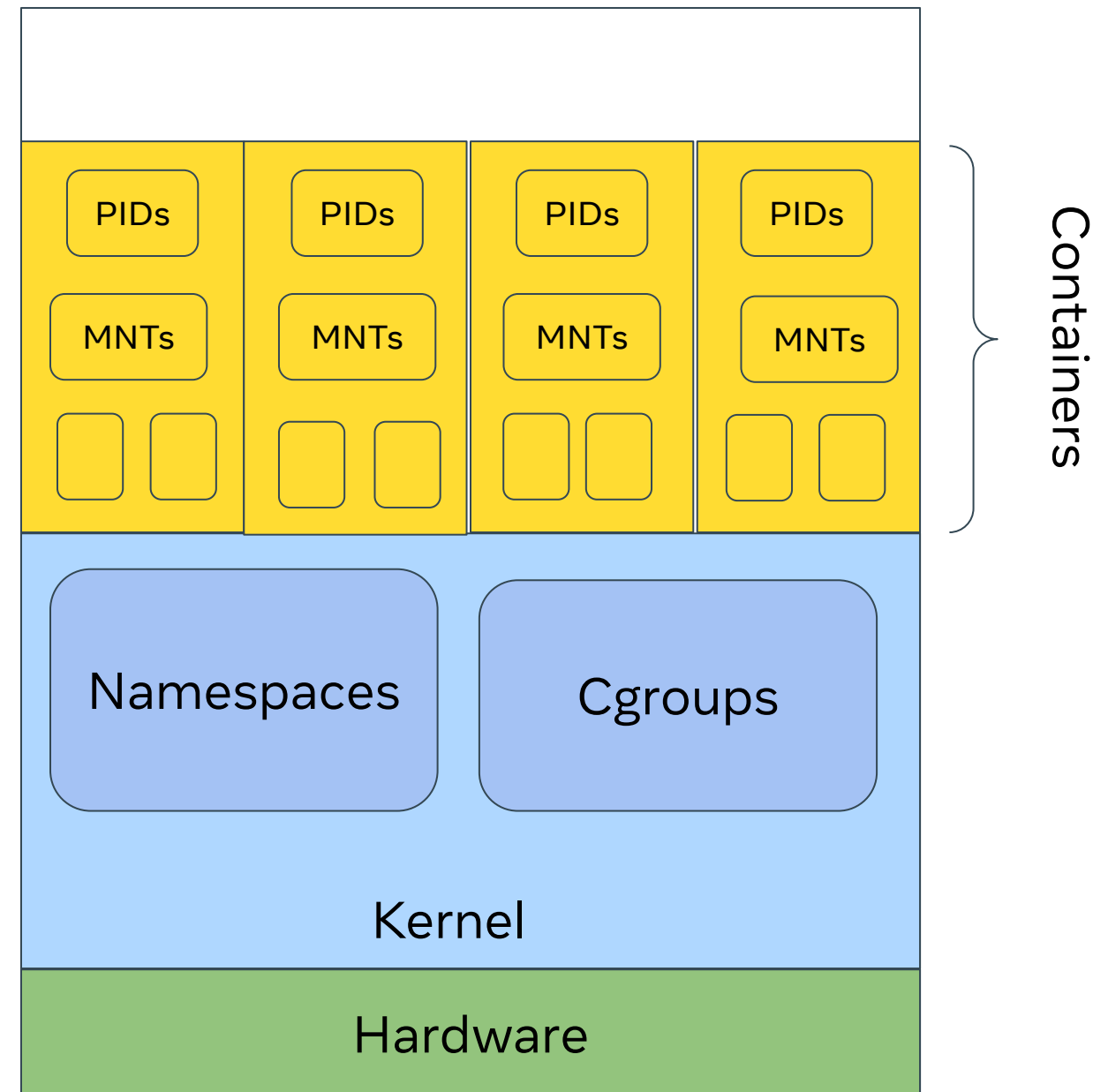
∞ Meta

# Agenda

∞ Meta

# 01 Building solution without Network Namespaces

**Meta**

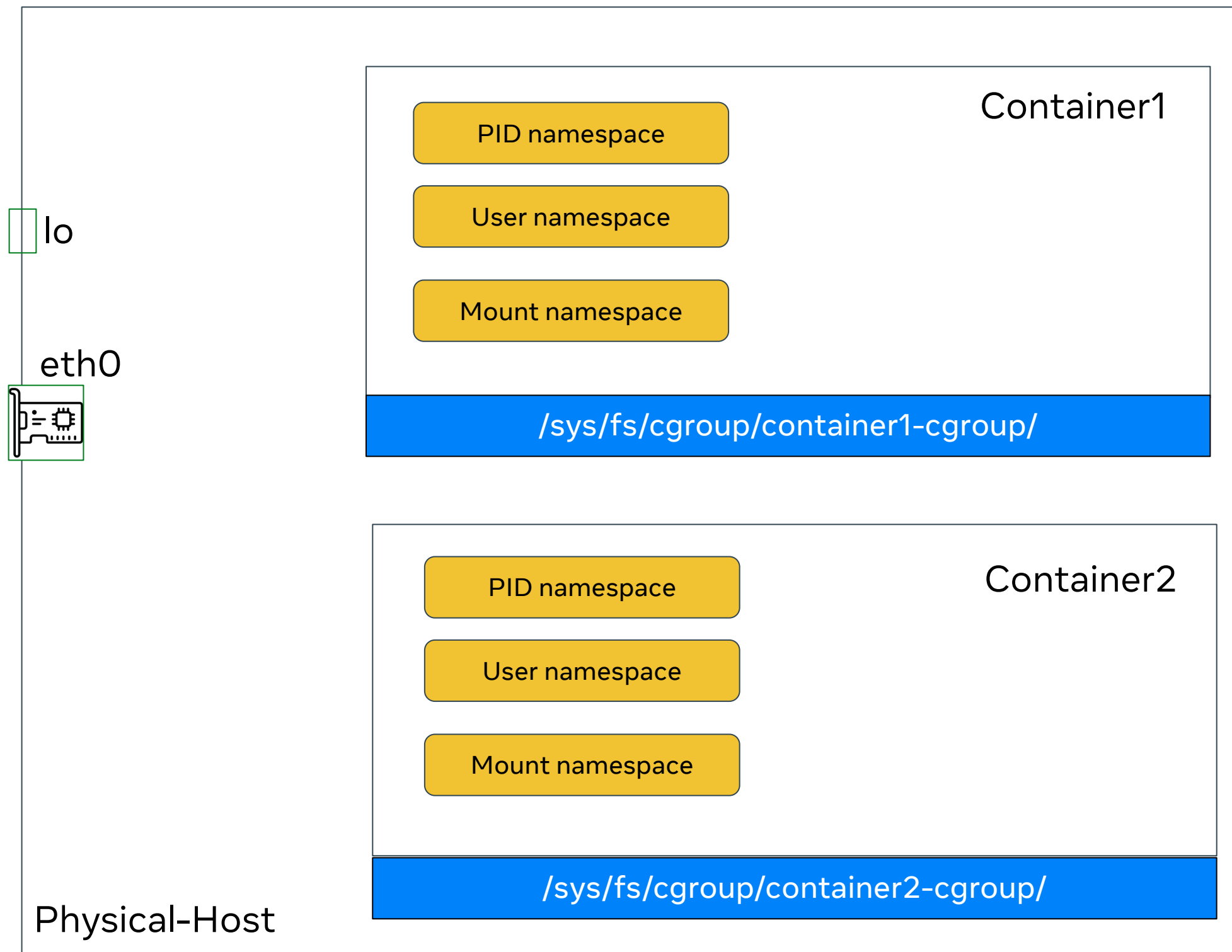# Linux Containers ?

# Linux Containers



- Container shapes
  - Square Shapes
  - L-Shape
- Host Accessibility
  - Single Tenant
  - Multi-Tenant Host
- Container Isolation
  - Resource Isolations
  - Shared Resources

# Linux Containers : Shared-Network Resource

lo

eth0

**Container1**

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container1-cgroup/

**Container2**

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container2-cgroup/

Physical-Host

- Container1 & Container2 both shares the host-network namespace
- No extra-network configuration setup

∞ Meta

# Shared-Network Resource: Port mgmt

lo

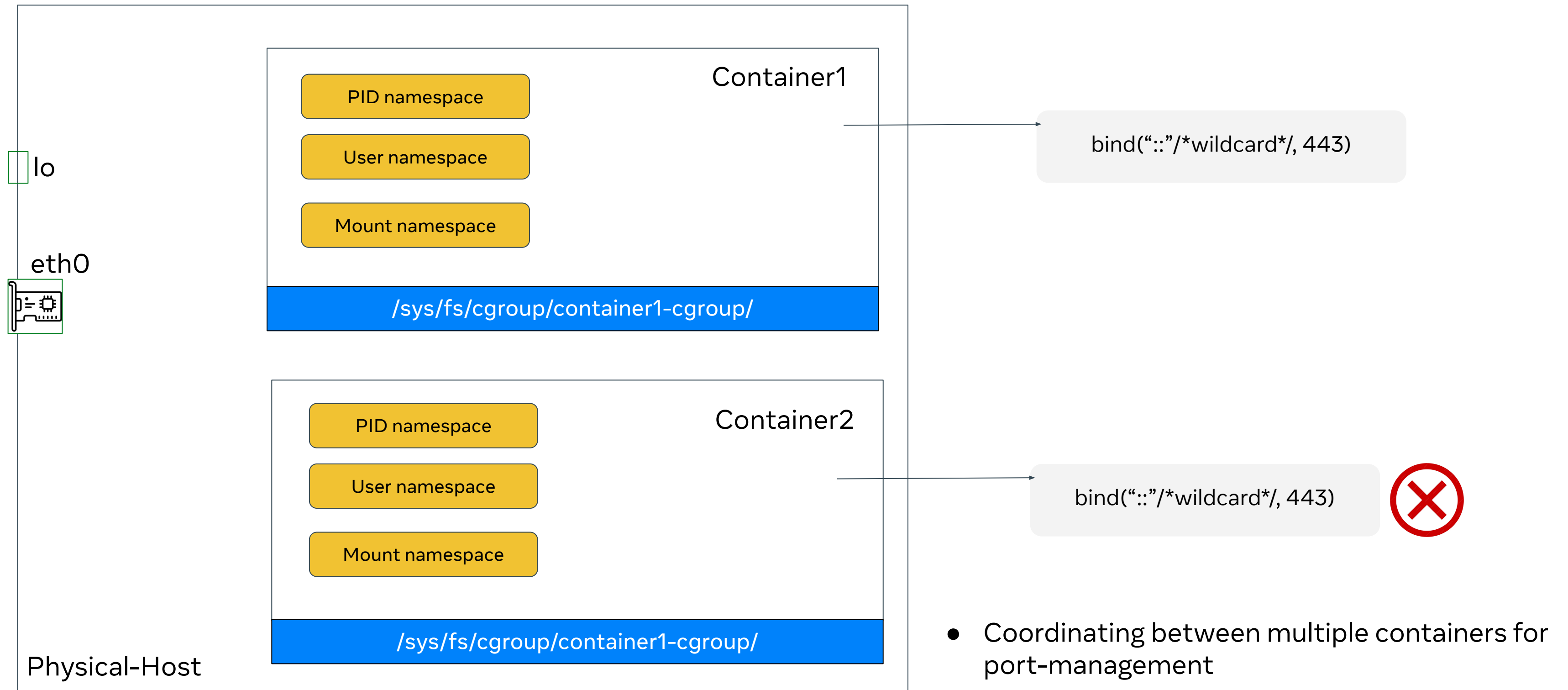eth0

Physical-Host

**Container1**

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container1-cgroup/

bind("::"/*wildcard*/, 443)

**Container2**

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container2-cgroup/

∞ Meta

# Shared-Network Resource: Port mgmt

Container1

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container1-cgroup/

bind(“::”/*wildcard*/, 443)

Container2

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container1-cgroup/

bind(“::”/*wildcard*/, 443)

lo

eth0

Physical-Host

- Coordinating between multiple containers for port-management

∞ Meta

# Shared-Network Resource : Port mgmt

**Meta**

lo

eth0

**Container1**

| PID namespace |
| User namespace |
| Mount namespace |

**/sys/fs/cgroup/container1-cgroup/**

bind("localhost", 443)

**Container2**

| PID namespace |
| User namespace |
| Mount namespace |

**/sys/fs/cgroup/container2-cgroup/**

bind("localhost", 443) ❌

Physical-Host

- Coordinating between multiple containers for port-management
- Binding to localhost is exposed to other containers & host

# Shared-Network Resource : Service traffic mgmt

DC Network

Fixed Uplink Capacity

Rack Switch (RSW)

Physical-Host

lo

eth0

Marking traffic at source (QoS)

## Container1

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container1-cgroup/

## Container2

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container2-cgroup/

∞ Meta

# Shared-Network Resource : Service traffic mgmt



**Physical-Host**

lo

Fixed Uplink Capacity

eth0

Rack Switch (RSW)

DC Network

eBPF

Marking traffic at source (QoS)

Identify Service and apply network Policies ?

**Container1**

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container1-cgroup/

**Container2**

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container2-cgroup/

∞ Meta

# Shared-Network Resource : Service traffic mgmt

Secure Network : Which service is allowed to pass/deny ?

Fixed Uplink Capacity

Rack Switch (RSW)

Troubleshoot the traffic : Which container originated the traffic from {Host IPv6} ?

Physical-Host

lo

eth0

Host IPv6
Host IPv6

Container1

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container1-cgroup/

Container2

PID namespace

User namespace

Mount namespace

/sys/fs/cgroup/container2-cgroup/

∞ Meta

Need network identifier to each container and some level of isolation ?

We decided to give unique IPv6 identity to each container

Meta

We decided to give unique IPv6 identity to each container

How to tie this IPv6 identity to the container ?

Meta

We decided to give unique IPv6 identity to each container

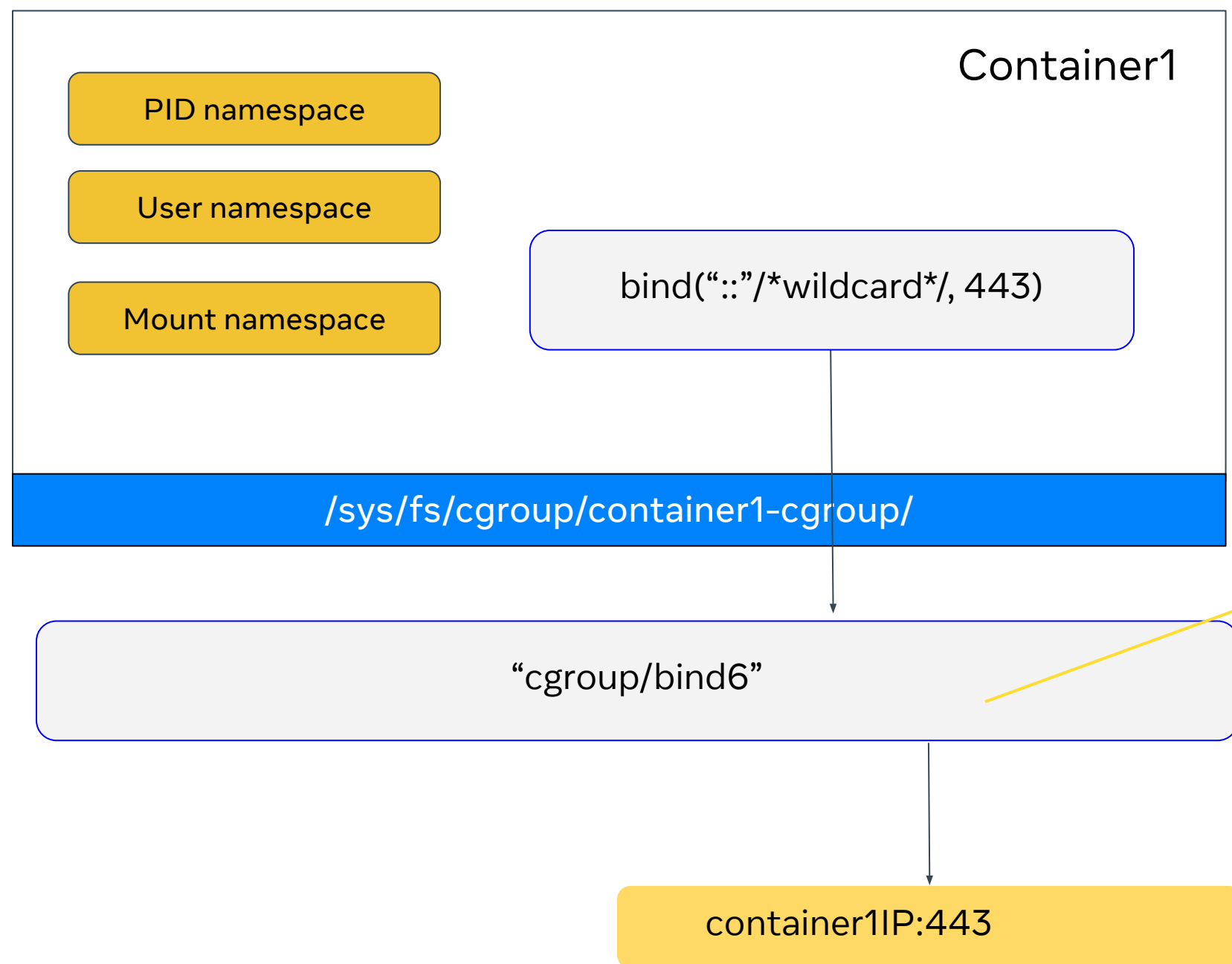How to tie this IPv6 identity to the container ?
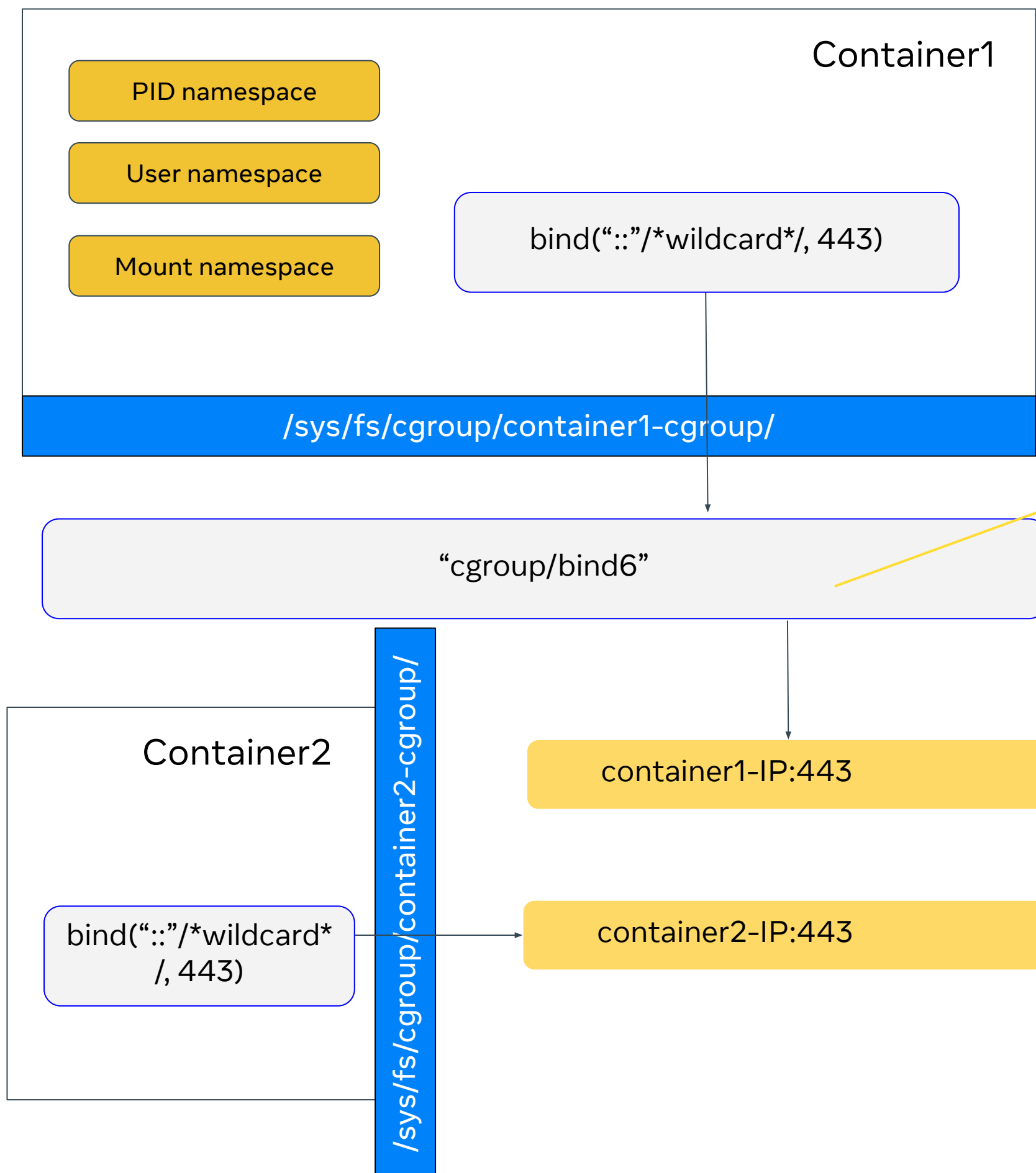
🐝eBPF cgroup-bpf : bind/connect/sendmsg

∞ Meta

# eBPF Play : Same Host-Network-namespace

**Container1**

- PID namespace
- User namespace
- Mount namespace

bind("::"/*wildcard*/, 443)

/sys/fs/cgroup/container1-cgroup/

"cgroup/bind6"

container1IP:443

```c
/* User bpf_sock_addr struct to access socket fields and sockaddr struct passed
 * by user and intended to be used by socket (e.g. to bind to, depends on
 * attach type).
 */
struct bpf_sock_addr {
    __u32 user_family;  /* Allows 4-byte read, but no write. */
    __u32 user_ip4;     /* Allows 1,2,4-byte read and 4-byte write.
                         * Stored in network byte order.
                         */
    __u32 user_ip6[4];  /* Allows 1,2,4,8-byte read and 4,8-byte write.
                         * Stored in network byte order.
                         */
    __u32 user_port;    /* Allows 1,2,4-byte read and 4-byte write.
                         * Stored in network byte order
                         */
    __u32 family;       /* Allows 4-byte read, but no write */
    __u32 type;         /* Allows 4-byte read, but no write */
    __u32 protocol;     /* Allows 4-byte read, but no write */
    __u32 msg_src_ip4;  /* Allows 1,2,4-byte read and 4-byte write.
                         * Stored in network byte order.
                         */
    __u32 msg_src_ip6[4]; /* Allows 1,2,4,8-byte read and 4,8-byte write.
                         * Stored in network byte order.
                         */
    __bpf_md_ptr(struct bpf_sock *, sk);
};
```
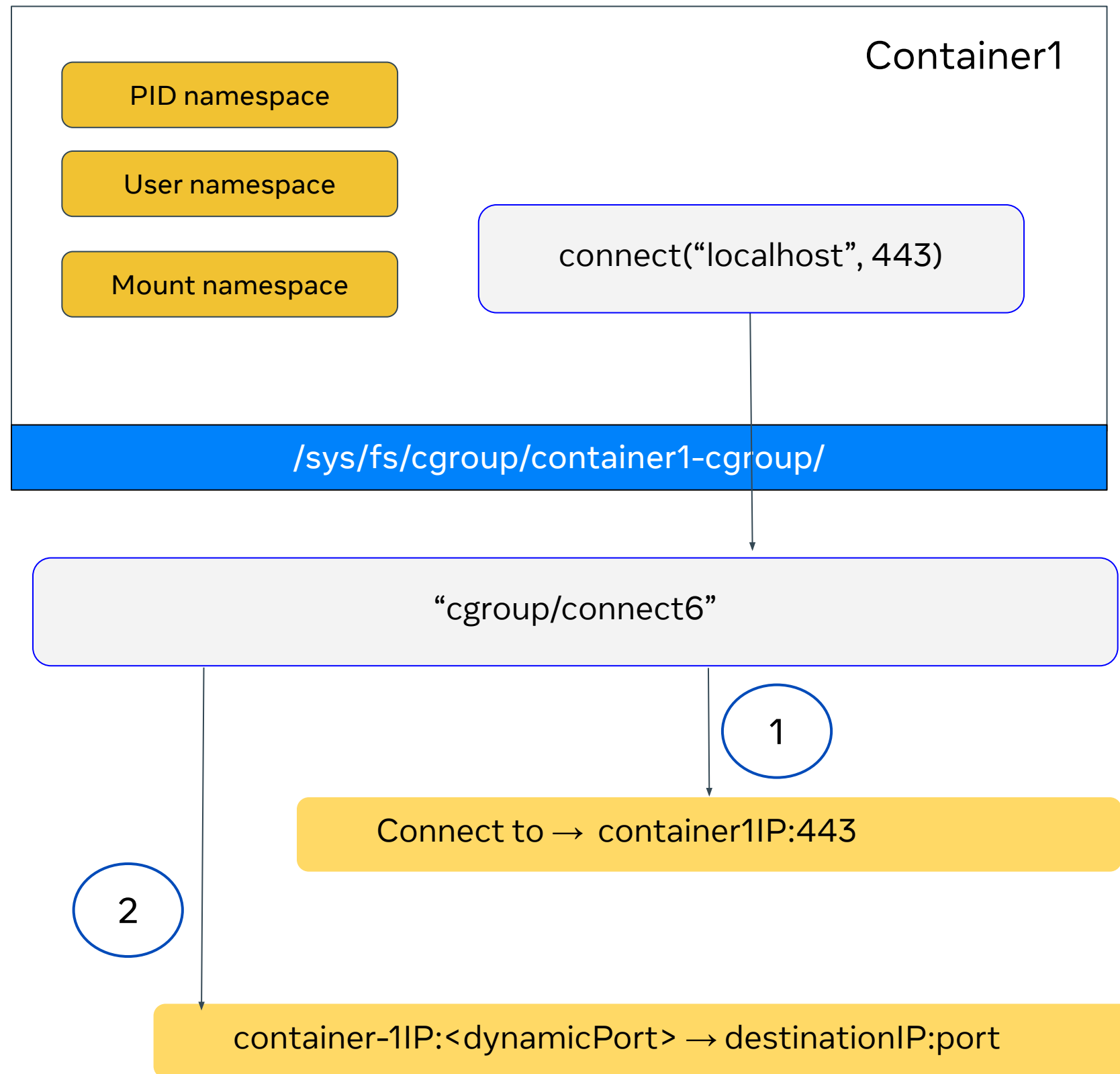
∞ Meta

# eBPF Play : Same Host-Network-namespace

Meta

## Container1

PID namespace

User namespace

Mount namespace

bind("::"/*wildcard*/, 443)

/sys/fs/cgroup/container1-cgroup/

"cgroup/bind6"

container1-IP:443

## Container2

bind("::"/*wildcard*/, 443)

/sys/fs/cgroup/container2-cgroup/

container2-IP:443

```c
/* User bpf_sock_addr struct to access socket fields and sockaddr
struct passed
 * by user and intended to be used by socket (e.g. to bind to,
depends on
 * attach type).
 */
struct bpf_sock_addr {
    __u32 user_family;  /* Allows 4-byte read, but no write. */
    __u32 user_ip4;     /* Allows 1,2,4-byte read and 4-byte write.
            * Stored in network byte order.
            */
    __u32 user_ip6[4];  /* Allows 1,2,4,8-byte read and 4,8-byte
write.
            * Stored in network byte order.
            */
    __u32 user_port;    /* Allows 1,2,4-byte read and 4-byte write.
            * Stored in network byte order
            */
    __u32 family;       /* Allows 4-byte read, but no write */
    __u32 type;         /* Allows 4-byte read, but no write */
    __u32 protocol;     /* Allows 4-byte read, but no write */
    __u32 msg_src_ip4;  /* Allows 1,2,4-byte read and 4-byte write.
            * Stored in network byte order.
            */
    __u32 msg_src_ip6[4]; /* Allows 1,2,4,8-byte read and 4,8-byte
write.
            * Stored in network byte order.
            */
    __bpf_md_ptr(struct bpf_sock *, sk);
};
```

# eBPF Play : Same Host-Network-namespace

**Container1**

- PID namespace
- User namespace
- Mount namespace

connect("localhost", 443)

/sys/fs/cgroup/container1-cgroup/

"cgroup/connect6"

**1**

Connect to → container1IP:443

**2**

container-1IP:<dynamicPort> → destinationIP:port

- Re-write the destination to container IP
- Set the source identity for the outgoing traffic

```
sa.sin6_family = AF_INET6;
sa.sin6_port = bpf_htons(0);

in6cpy(&sa.sin6_addr, task_ip);

/* Rewrite source IP. */
if (bpf_bind(ctx, (struct sockaddr*)&sa, sizeof(sa)) != 0)
    return FAIL_OPEN;
```

∞ Meta

# Shared host-network space : Re-look challenges

➔ Two containers can not start service on same fixed port
- ◆ Unique IP per container helps at certain extent
- ◆ Fails if the same port binds on the wildcard by other host-based services

➔ Container localhost service get exposed to host & whole meta
- ◆ Services binds on container-IP which is routable in Meta fleet
- ◆ eBPF helps but stills adds an additional overhead to handle it

➔ Does not allow wildcard binding inside the container (hacks for additional VIPs)
- ◆ Hard to share same port among container IP and BGP VIPs

∞ Meta

# Traffic Redirection over TLS

∞ Meta



- Using eBPF hooks with socket cookies, it is easy to track TCP connections
- For UDP sockets, where the same source IP:port can be used for multiple destinations; proxy can't track the connections
  - Packet encapsulation helps to solve this but that requires tc-bpf based solution
  - Moving every container's UDP traffic tracking at host-eth0 is again a challenge in multi-tenant host
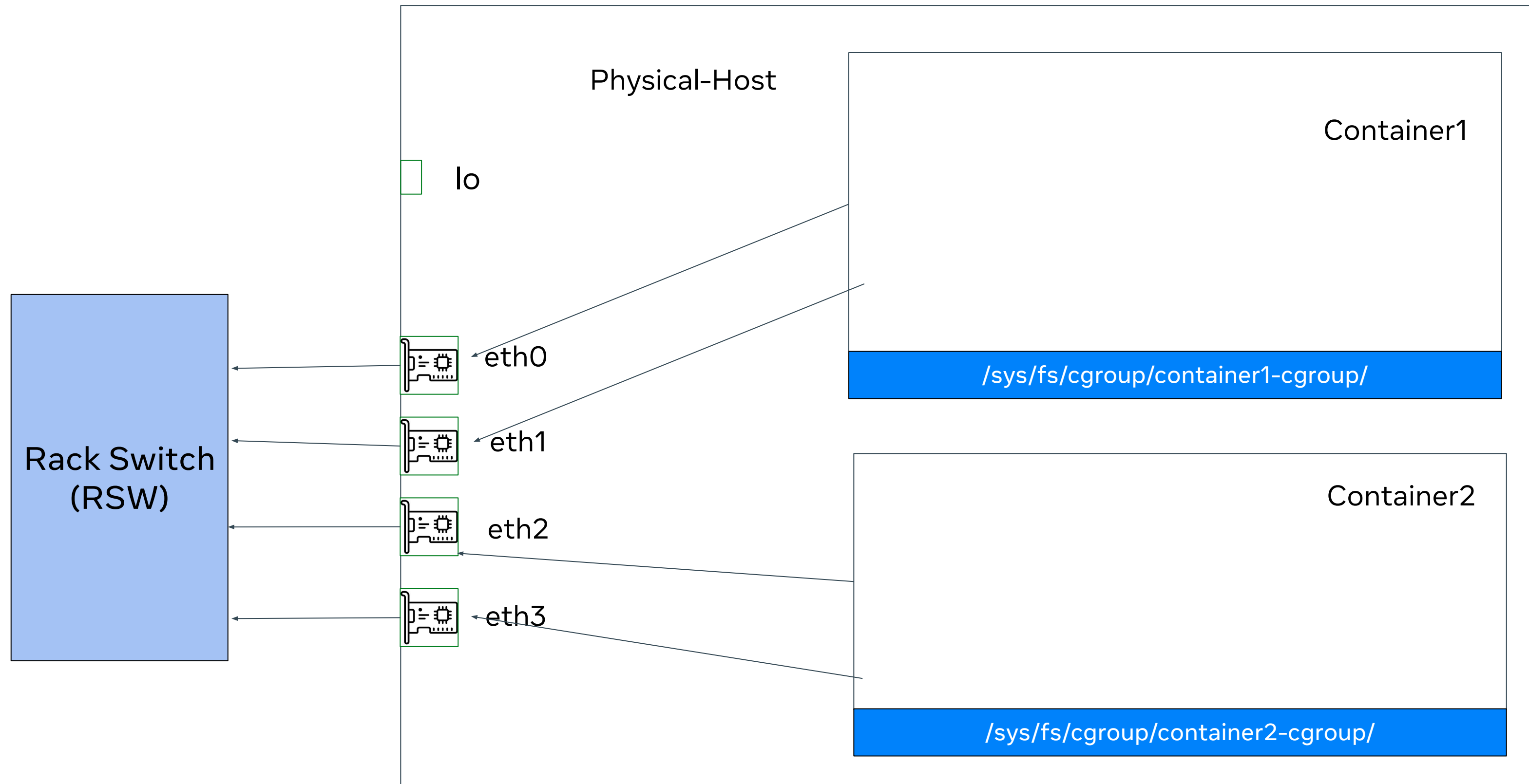
# Process VM inside the container

∞ Meta

Tw container

VM

eth0

Task Cgroup

ttls-foward-proxy

FW egress tc-bpf

TTLS NAT 646 tc-bpf

TTLS egress redirection tc-bpf

HostNS

tap0

FW ingress tc-bpf

TTLS NAT 464 tc-bpf

TTLS ingress redirection tc-bpf

eth0

# Process VM inside the container

∞ Meta

Stacking large number of containers, makes life hard to manage shared host resources.

# Multi-NIC Host : Routing Mgmt (Ingress/Egress)

Physical-Host

lo

Container1

/sys/fs/cgroup/container1-cgroup/

eth0

eth1

Rack Switch
(RSW)

Container2

eth2

eth3

/sys/fs/cgroup/container2-cgroup/

∞ Meta

# Multi-NIC Host : Routing Mgmt

Physical-Host

Io

Container1

BGP multi nexthop route1

ECMP based MultipPath Routing eth0/eth1

eth0

eth1

Rack Switch (RSW)

/sys/fs/cgroup/container1-cgroup/

Container2

eth2

ECMP based MultipPath Routing eth2/eth3

eth3

BGP multi nexthop route2

/sys/fs/cgroup/container2-cgroup/

∞ Meta

# Apart from the gaps, other use-cases ?

➔ L2 level secure isolation to avoid all enforcement at host level

➔ per-container tc/XDP eBPF support again to avoid physical-eth0 a choking point

➔ Some of the emulated services need IPv4 support

➔ Running third-party services/applications with jailed environment

➔ Debug the container level traffic without having access to host

∞ Meta

02    Building solution with Network Namespaces

∞ Meta

# Finally Network NS : Network Connection model

Virtual Ethernet (VETH)



**Container-1**

veth0

eth0

Routing
process

**Container-2**

veth0

Host network space

∞ Meta

# Network Namespaces : Build network connectivity



Upper Stack (IP, Netfilter, Routing..)

eth0

task-netns

task-veth-0

host-veth-0

Host network space

∞ Meta

# Network Namespaces : Build network connectivity

Upper Stack (IP, Netfilter, Routing..)

task-netns

eth0

task-veth-0

host-veth-0

Host network space

➢ Global IP Forwarding enablement
➢ veth(4) is slower due to additional traversal of network stack

∞ Meta

# 03  How bpf helped to work around the issue

Meta

# Network Namespaces : Use of eBPF Kernel extensions

## bpf: Add redirect_neigh helper as redirect drop-in

```
Add a redirect_neigh() helper as redirect() drop-in replacement
for the xmit side. Main idea for the helper is to be very similar
in semantics to the latter just that the skb gets injected into
the neighboring subsystem in order to let the stack do the work
it knows best anyway to populate the L2 addresses of the packet
and then hand over to dev_queue_xmit() as redirect() does.

This solves two bigger items: i) skbs don't need to go up to the
stack on the host facing veth ingress side for traffic egressing
the container to achieve the same for populating L2 which also
has the huge advantage that ii) the skb->sk won't get orphaned in
ip_rcv_core() when entering the IP routing layer on the host stack.
```

∞ Meta
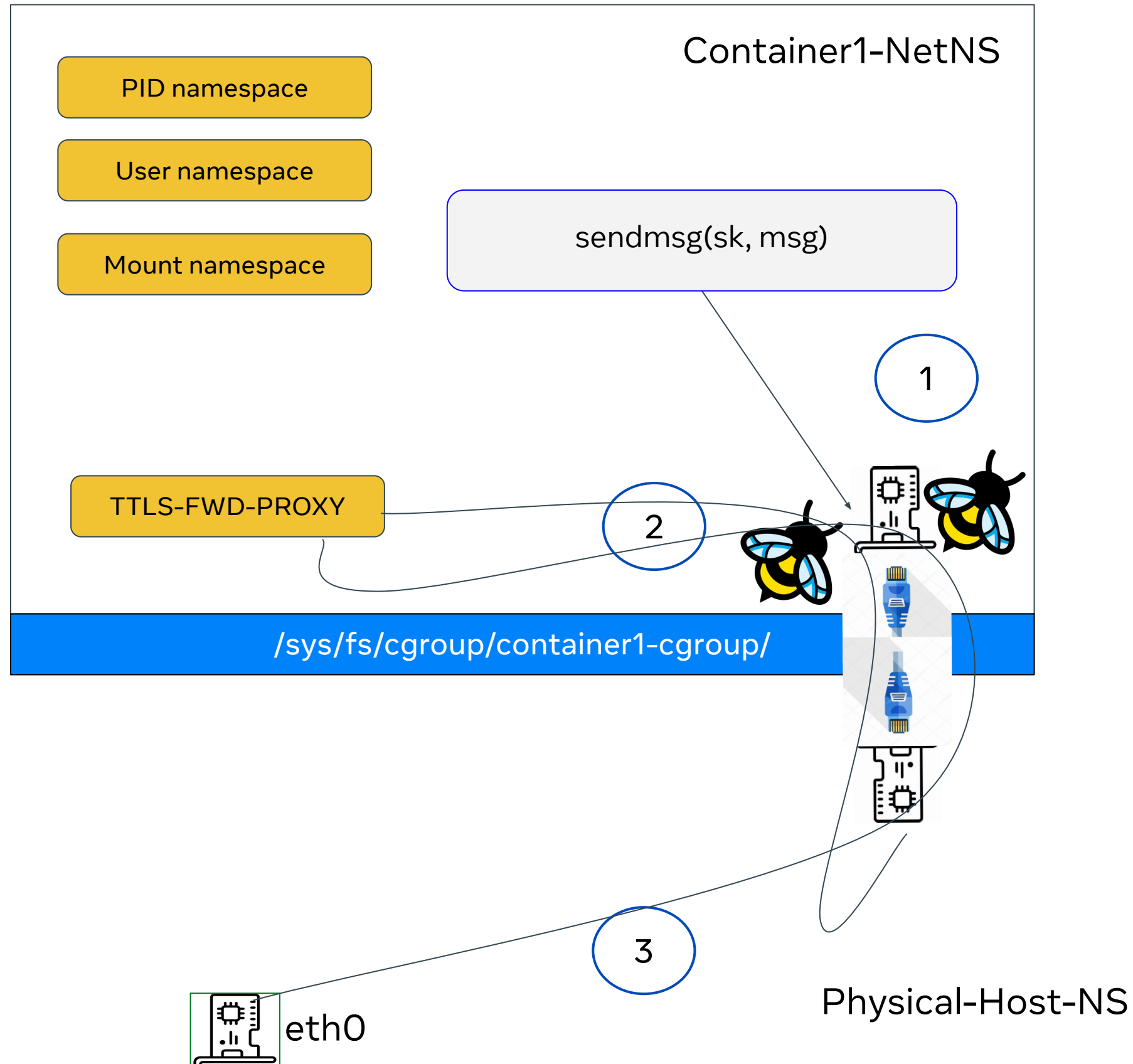
# Network Namespaces : Use of eBPF Kernel extensions

Upper Stack (IP, Netfilter, Routing..) 🚫

task-netns

eth0

🐝 eBPF

🐝 eBPF

task-veth-0

host-veth-0

Host network space

∞ Meta

# Network Namespaces : Build network connectivity



eth0

Upper Stack (IP, Netfilter, Routing..)

task-netns

task-veth-0

host-veth-0

Host network space

❌ ➤ Global IP Forwarding enablement
➤ veth(4) is slower due to additional traversal of network stack

∞ Meta

# UDP Traffic Redirection over TLS

**∞ Meta**

Container1-NetNS

PID namespace

User namespace

Mount namespace

sendmsg(sk, msg)

1

TTLS-FWD-PROXY

2

/sys/fs/cgroup/container1-cgroup/

3

eth0

Physical-Host-NS

- Challenges
  - Packet makes a round trip from veth0-egress to host-end and back to task's netns
  - Ingress program at veth0 is in-effective

# UDP Traffic Redirection over TLS



- Challenges
  - Packet makes a round trip from veth0-egress to host-end and back to task's netns
  - Ingress program at veth0 is in-effective for ttls-fwd-proxy → user client.

- Current Solution:
  - Ingress program needed to attach at "lo" due to kernel optimizing the route.
  - bpf_redirect("eth0"→ifindex, BPF_F_INGRESS) & update MAC
  - Change direction from EGRESS to INGRESS

04  Other Virtual Devices & Performance Improvements

∞ Meta

# ipvlan/veth/netkit/bare-metal
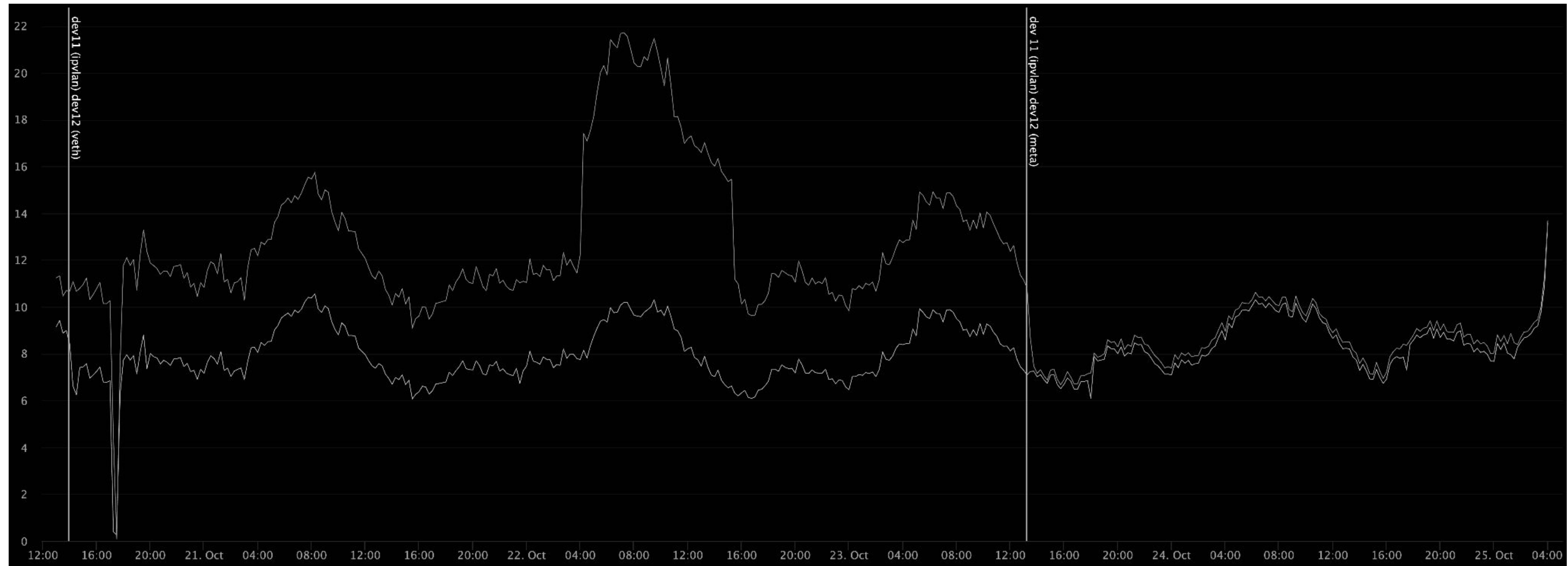
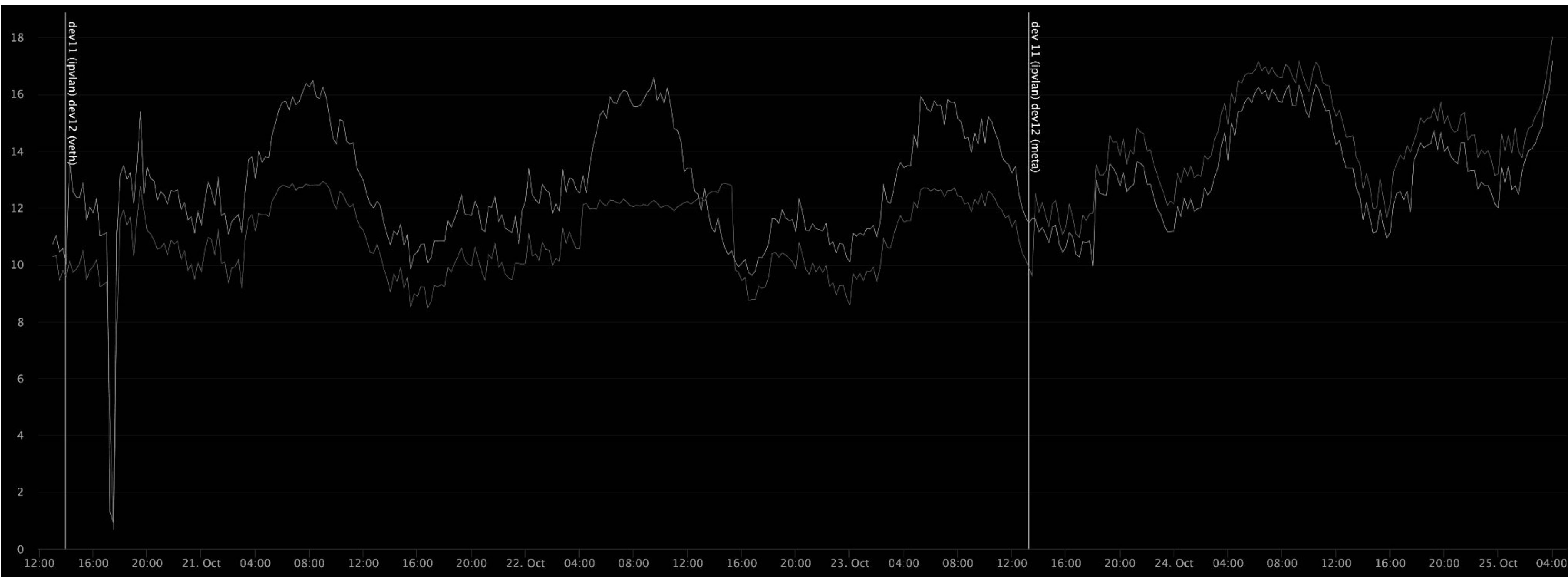# veth -> netkit (cpu-util)

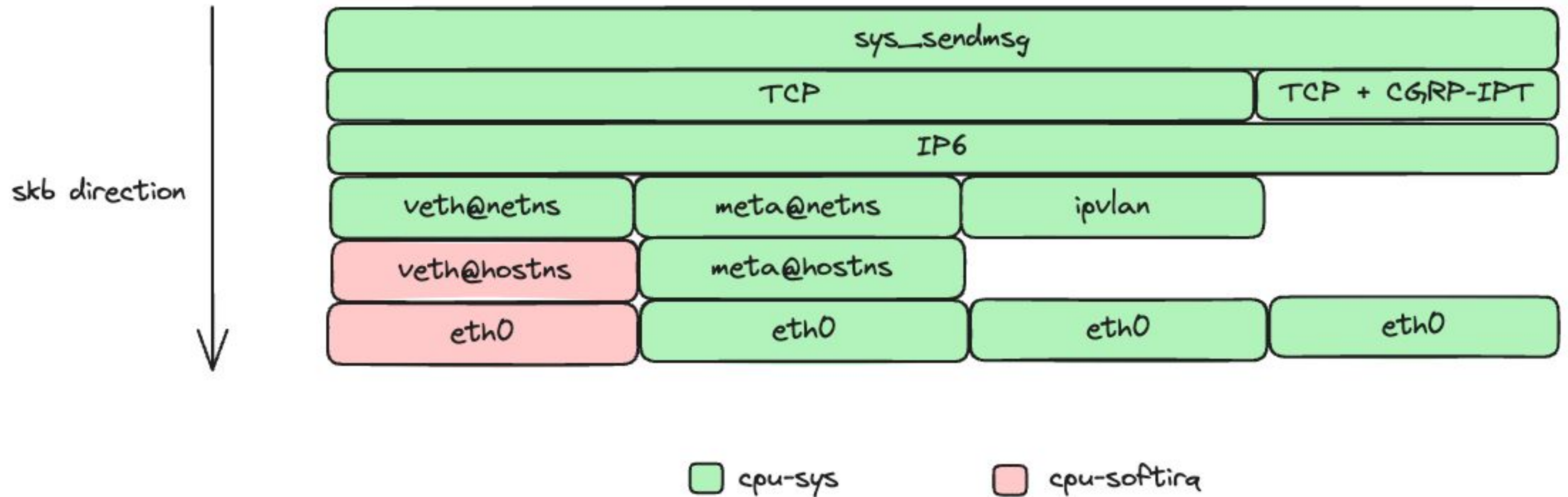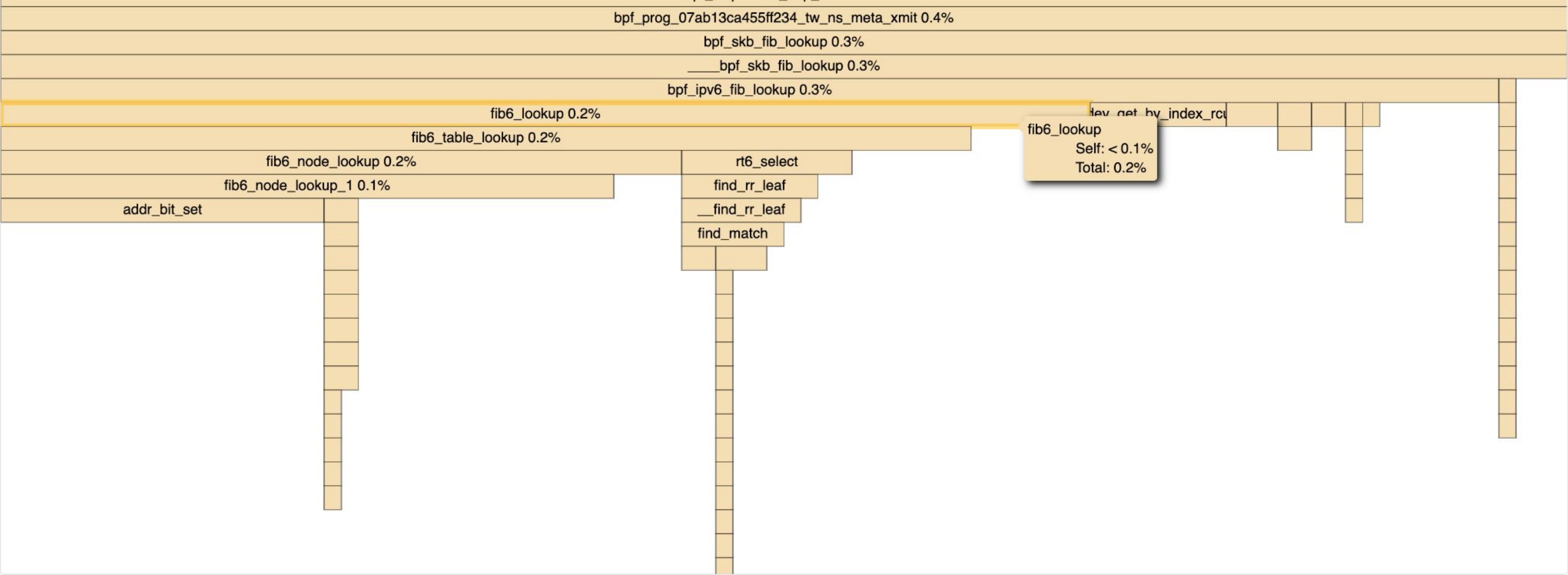# veth -> netkit (cpu-sys + cpu-softirq)

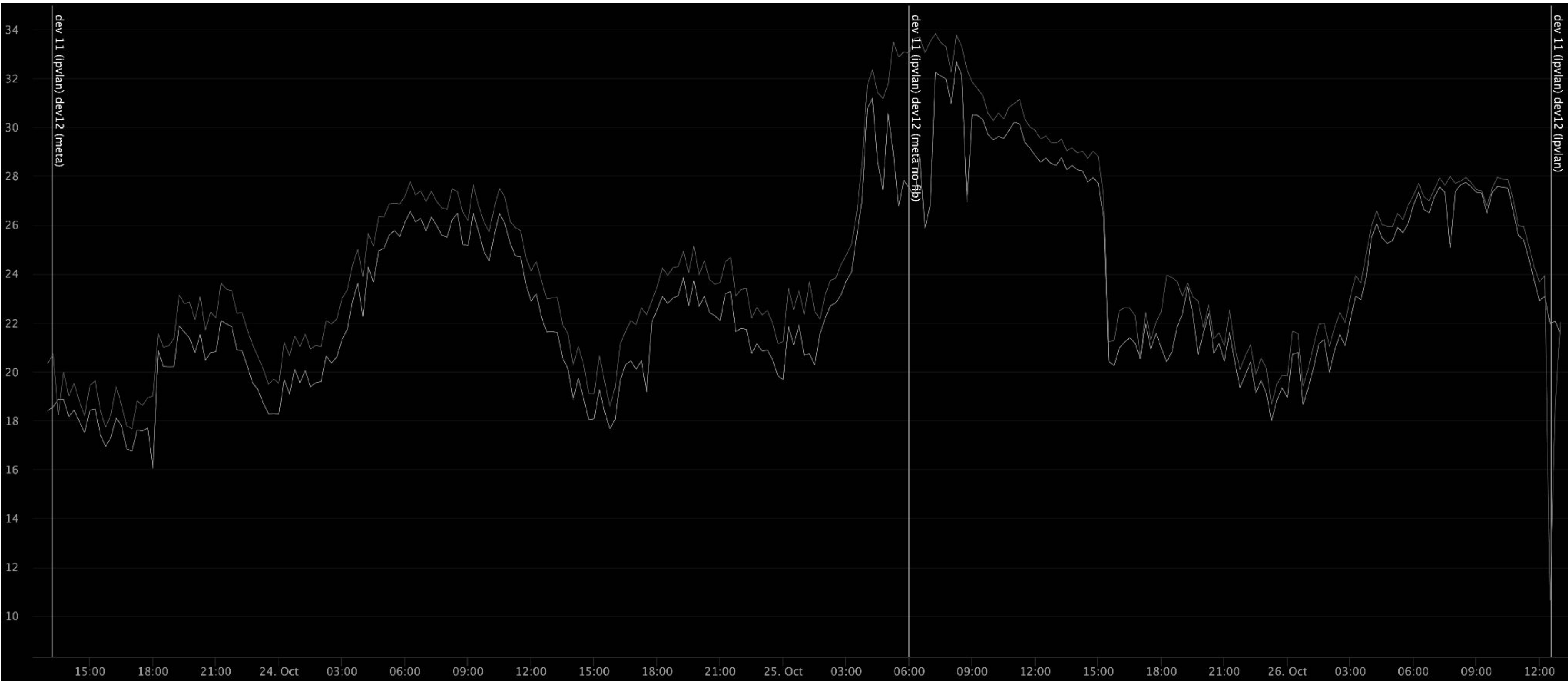# veth -> netkit (cpu-softirq)

# veth -> netkit (cpu-sys)

# egress veth/netkit => phy-eth0

# bpf prog at netkit



bpf_prog_07ab13ca455ff234_tw_ns_meta_xmit 0.4%

bpf_skb_fib_lookup 0.3%

____bpf_skb_fib_lookup 0.3%

bpf_ipv6_fib_lookup 0.3%

fib6_lookup 0.2%

dev_get_by_index_rcu

fib6_table_lookup 0.2%

fib6_lookup
Self: < 0.1%
Total: 0.2%

fib6_node_lookup 0.2%

rt6_select

fib6_node_lookup_1 0.1%

find_rr_leaf

addr_bit_set

__find_rr_leaf

find_match

# netkit at L2 mode (cpu-sys + softirq)

# tcp_rr 500 flows 36 threads 1 byte req/rep

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | 1 x task | no netns (bare metal) | ipvlan L2 | meta L2 (no fib lookup) | veth L2 (no fib lookup) |
| 2 | cpu-util | 84.78 | 85.19 | 85.85 | 87.09 |
| 3 | cpu-softirq | 2.48 | 2.85 | 2.17 | 14.78 |
| 4 | cpu-sys | 63.24 | 63.73 | 64.79 | 53.89 |
| 5 | cpu-user | 18.46 | 18.09 | 18.38 | 17.89 |
| 6 | #Transactions (M) / s | 1.63 | 1.61 | 1.62 | 1.62 |
| 7 | trans_per_s (K)/ cpu% | 19.17 | 18.90 | 18.88 | 18.61 |

∞ Meta

# Ipvlan vs meta (L2)

# ipvlan vs ipvlan (background difference)