

# Software Networking and Interfaces

on Linux

SRECon, Singapore | June 2019



Matt Turner

@mt165 | mt165.co.uk

The logo consists of a blue, stylized wave symbol that resembles a lowercase 'n' or 'w' with rounded, flowing lines. It is positioned to the left of the text.

native wave

# Outline

- Networking 101
- Interfaces
- Bridges
- Emergent Systems

# This is not about

- IP, TCP, addressing
- SDN

# Networking 101

# Ethernet and ARP

Ethernet - L2 protocol

MAC Address - Media Access Control Address - Ethernet address, eg  
c0:ff:ee:be:ef:69

ARP - Address Resolution Protocol - DNS for ethernet: IP -> MAC

```
$ arp
Address                HWtype  HWaddress          Flags Mask  Iface
192.168.0.239          ether    48:3b:38:01:6a:23  C           enp2s0
172.28.0.13            ether    02:42:ac:1c:00:0d  C           br-de368312f566
```

# vLANs

Virtual LANs

IEEE 802.1q

Simulates multiple networks using one set of cables and switches

Each vLAN has a short numeric ID

Nested vLANs - IEEE 802.1ad, aka “q in q”

# iptables

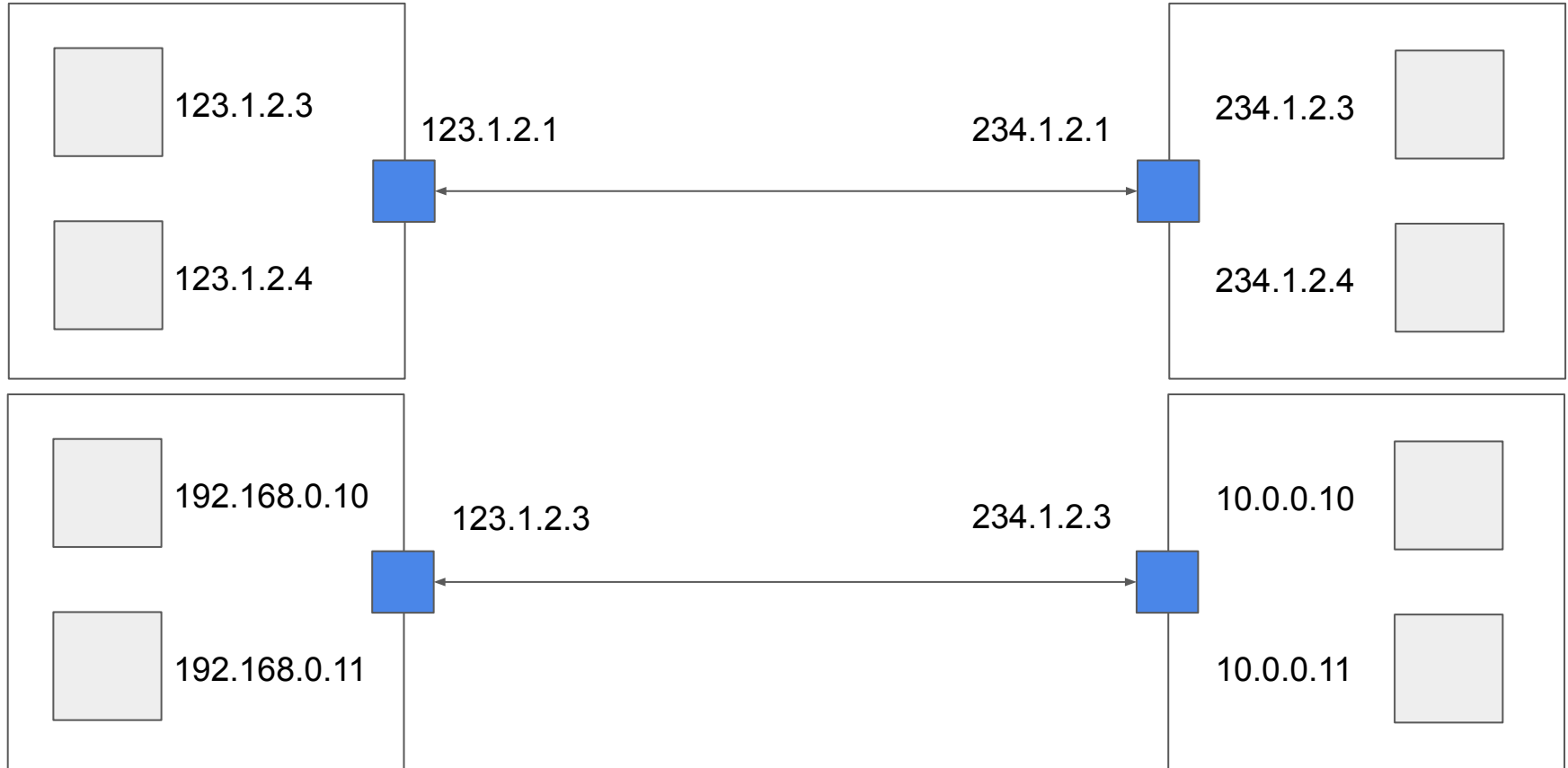
Linux kernel subsystem

Can do lots of things to packets as they pass through a system

Including: packet manipulation



# NAT vs Routing



# Route tables

Gives the next hop for a destination subnet.

```
$ route
Kernel IP routing table
Destination      Gateway          Genmask         Flags Metric Ref    Use Iface
default          moon.lan        0.0.0.0         UG    3     0      0 enp2s0
loopback         localhost       255.0.0.0       UG    0     0      0 lo
172.17.0.0      0.0.0.0         255.255.0.0    U     0     0      0 docker0
```

# Classful Routing

Class A	8 bit	0.0.0.0	127.0.0.1
Class B	16 bit	128.0.0.0	191.255.0.0
Class C	24 bit	192.0.0.0	223.255.255.0
Class D - multicast		224.0.0.0	239.255.255.255
Class E - reserved		240.0.0.0	255.255.255.255

# Private Address Ranges

The “24-bit block”	8 bit prefix	10.0.0.0 – 10.255.255.255	1 class A
The “20-bit block”	12 bit prefix	172.16.0.0 – 172.31.255.255	16 class Bs
The “16-bit block “	16 bit prefix	192.168.0.0 – 192.168.255.255	256 class Cs
Loopback	8 bit prefix	127.0.0.0 - 127.255.255.255	1 class A

# Classless Routing and CIDRs

Classless Inter-Domain Routing

Classful was too rigid and wasteful

Arbitrary ranges of addresses, notated as start address and size (as prefix mask)

CIDR notation: 192.168.42.0/24; 10.0.0.0/8

# Broadcast, Multicast, Anycast

Address all, some, one host

- All hosts in a subnet
- Some hosts, which have opted in
- One arbitrary host from a set

# DHCP

Asks a central server to allocate you an IP address

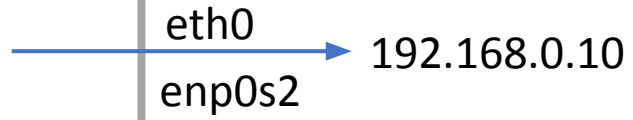
Based on an **ethernet** broadcast

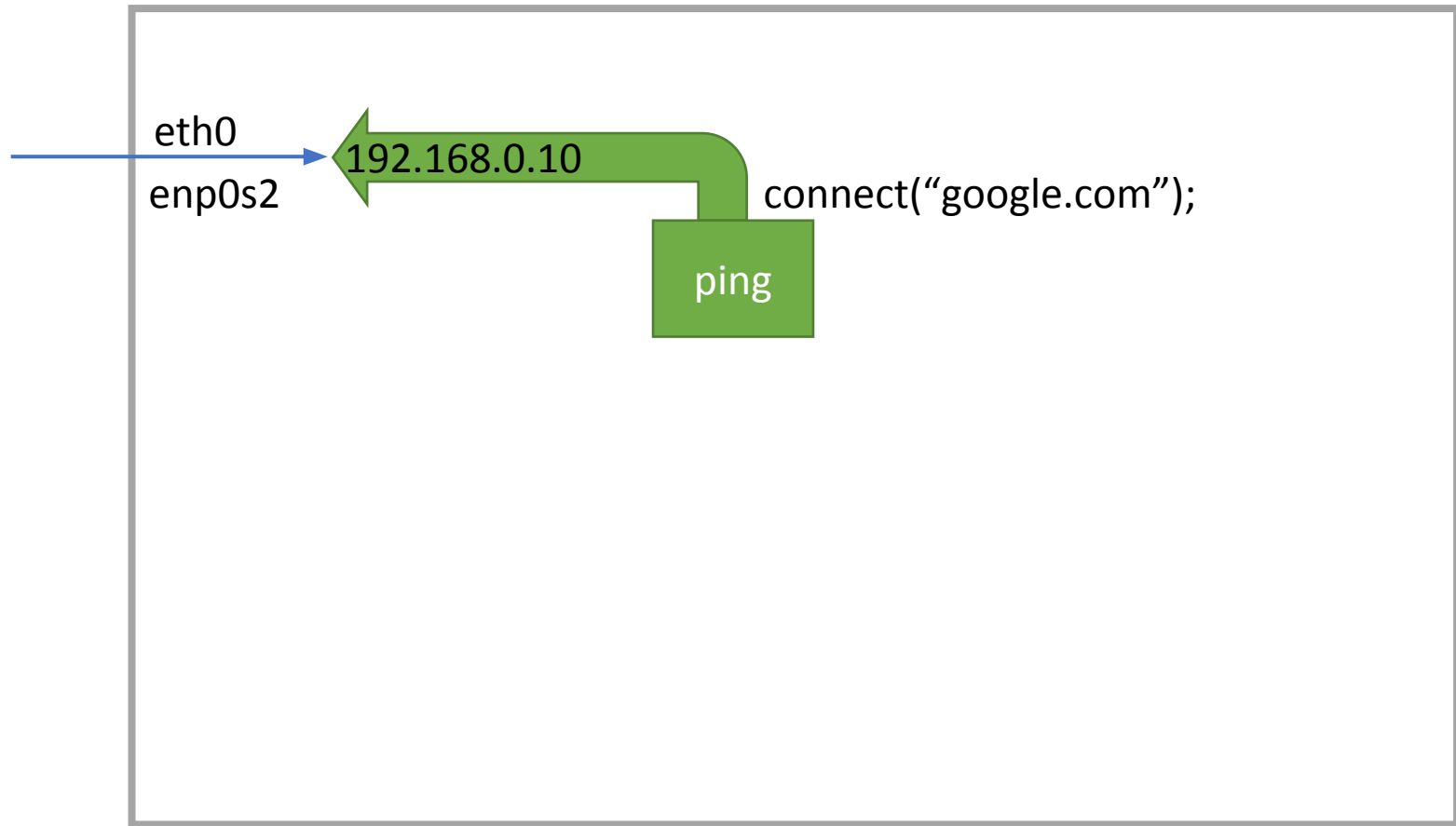
# Interfaces and Bridges

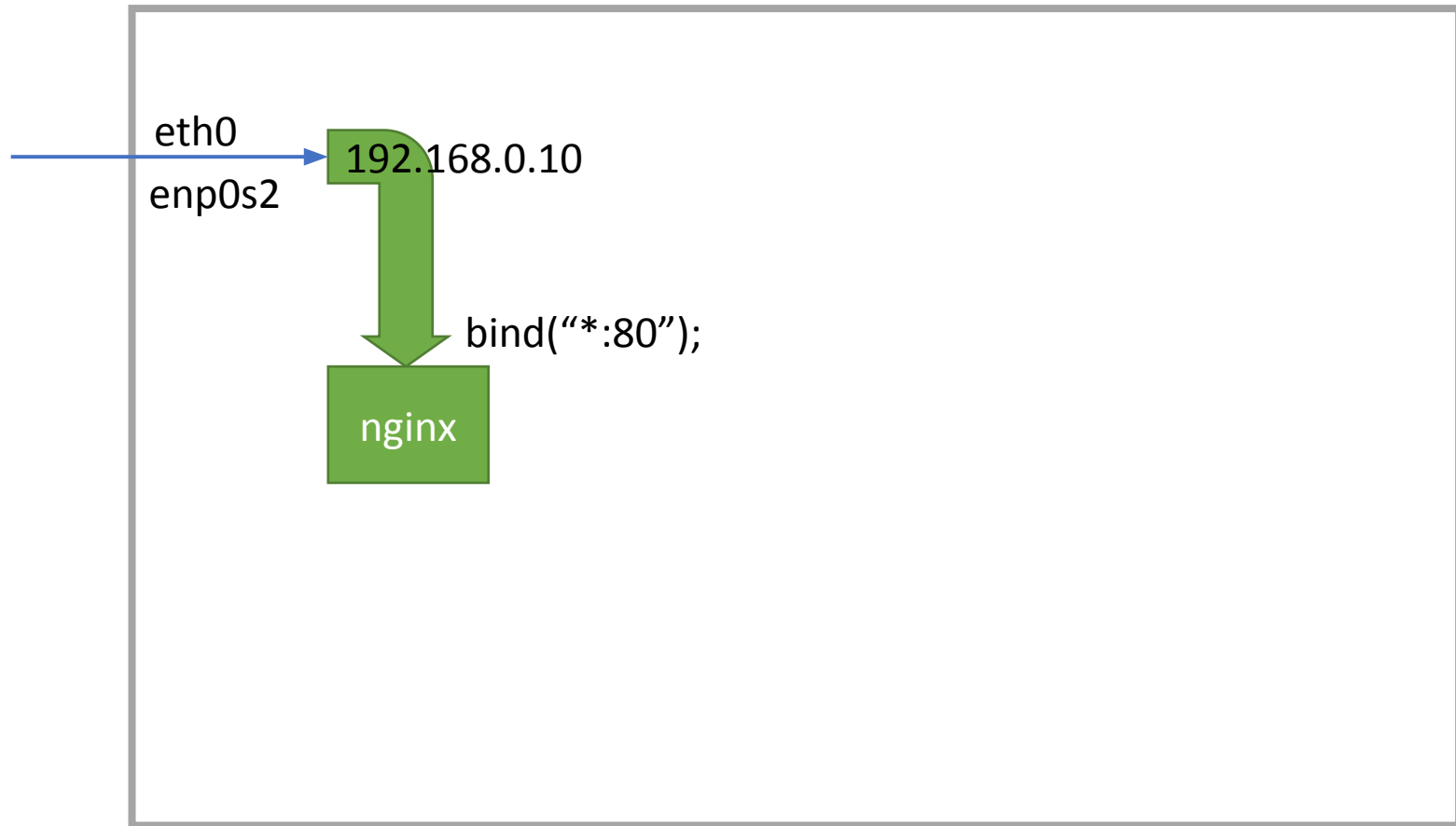


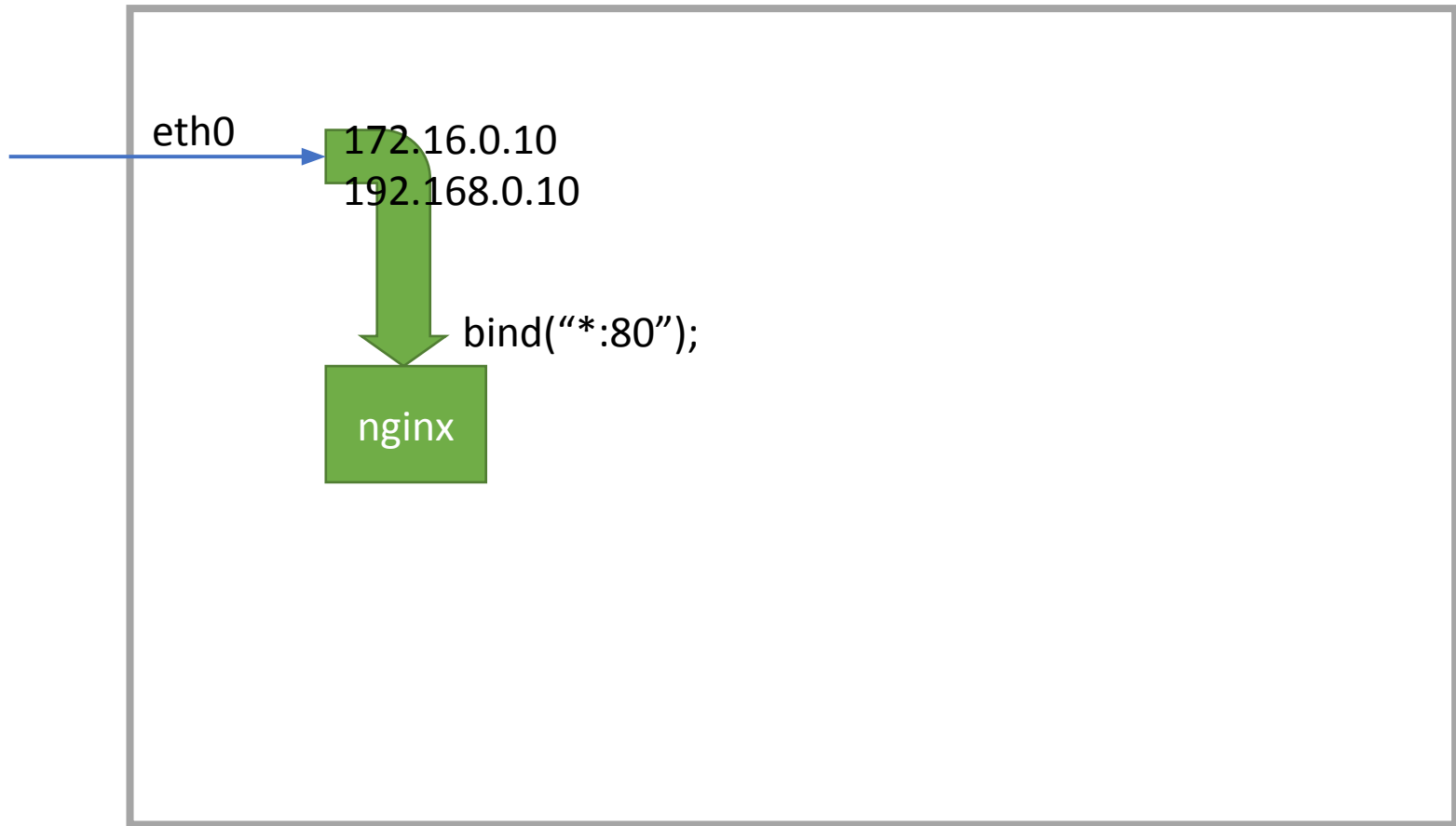


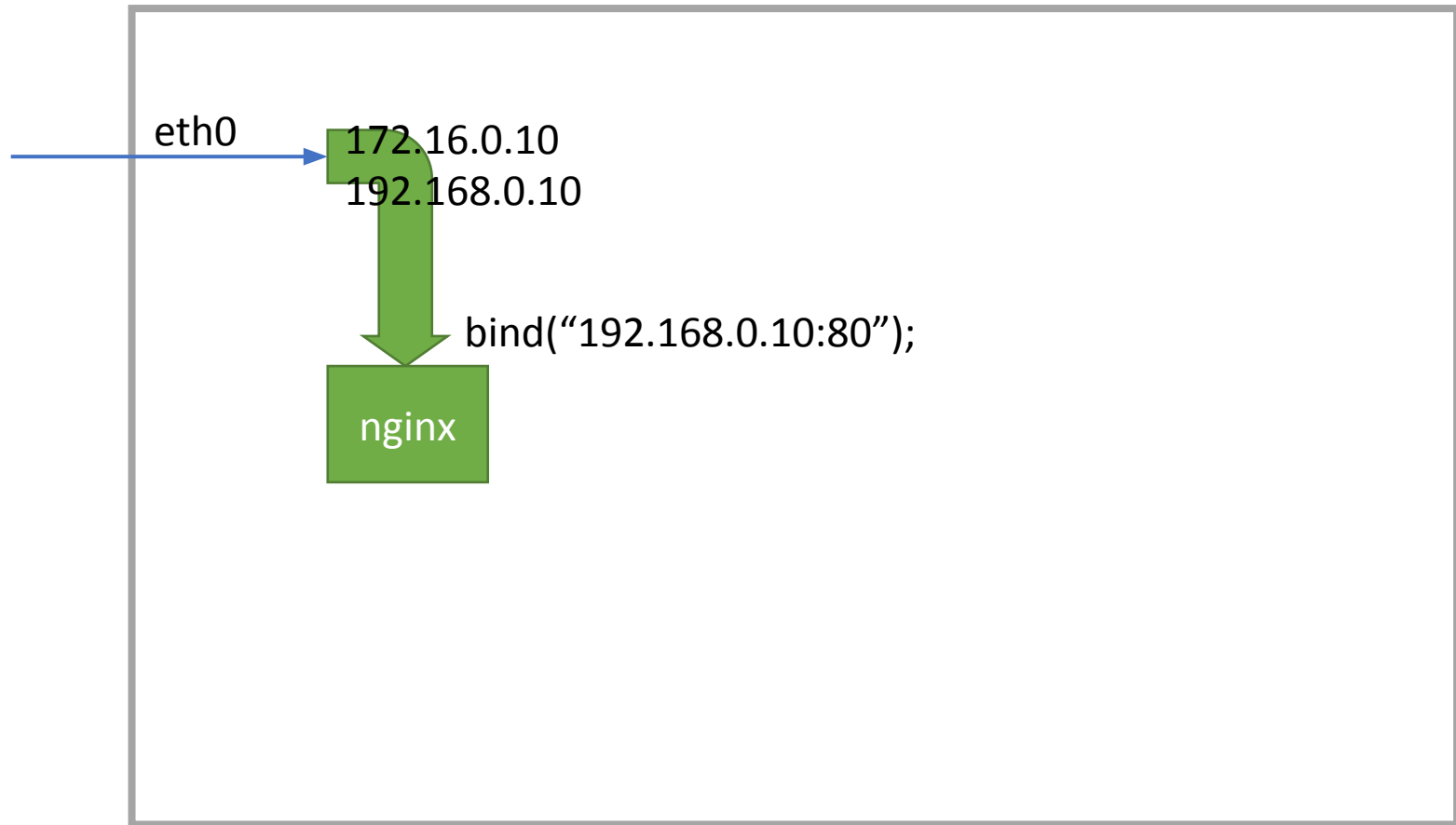


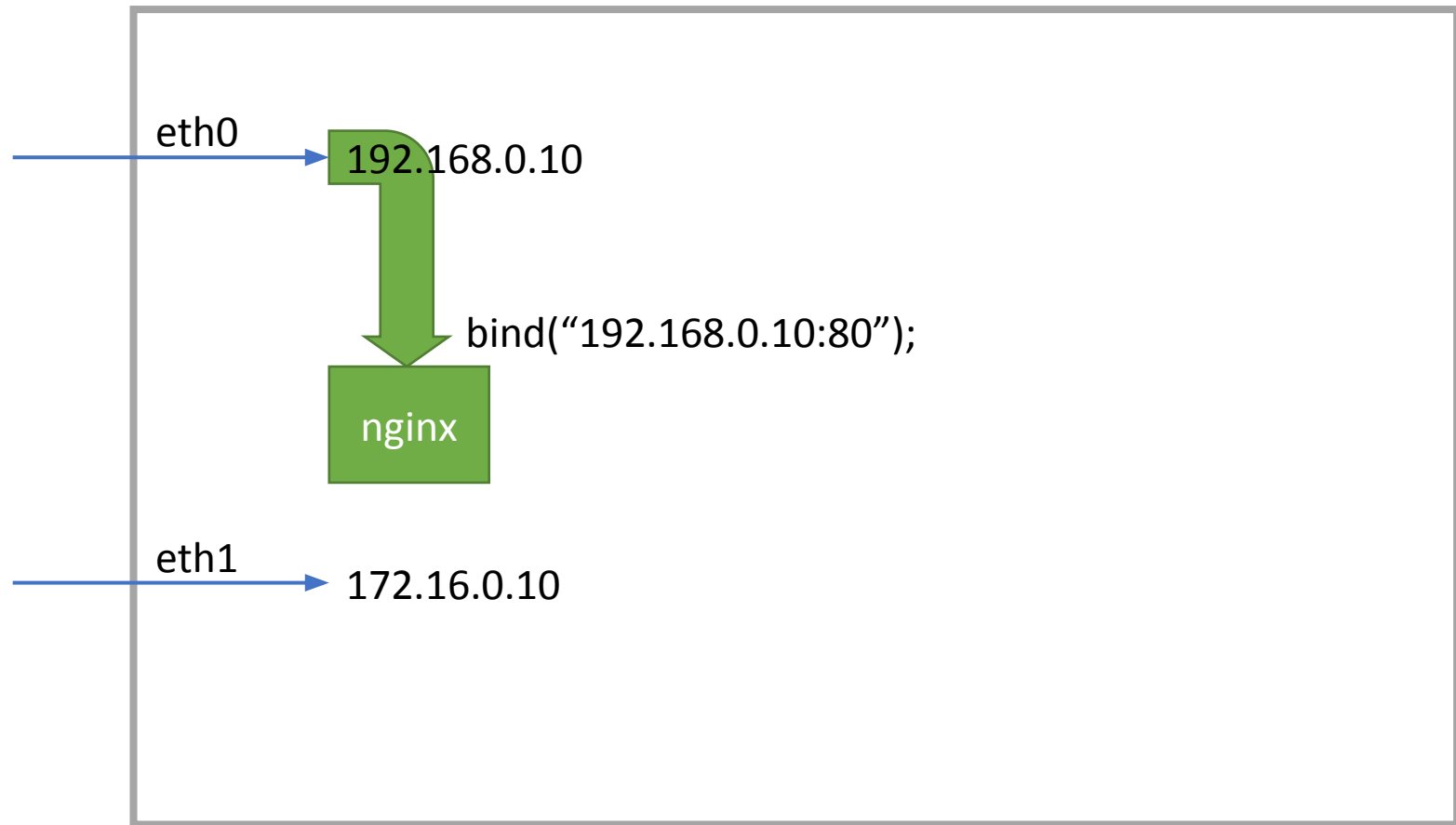




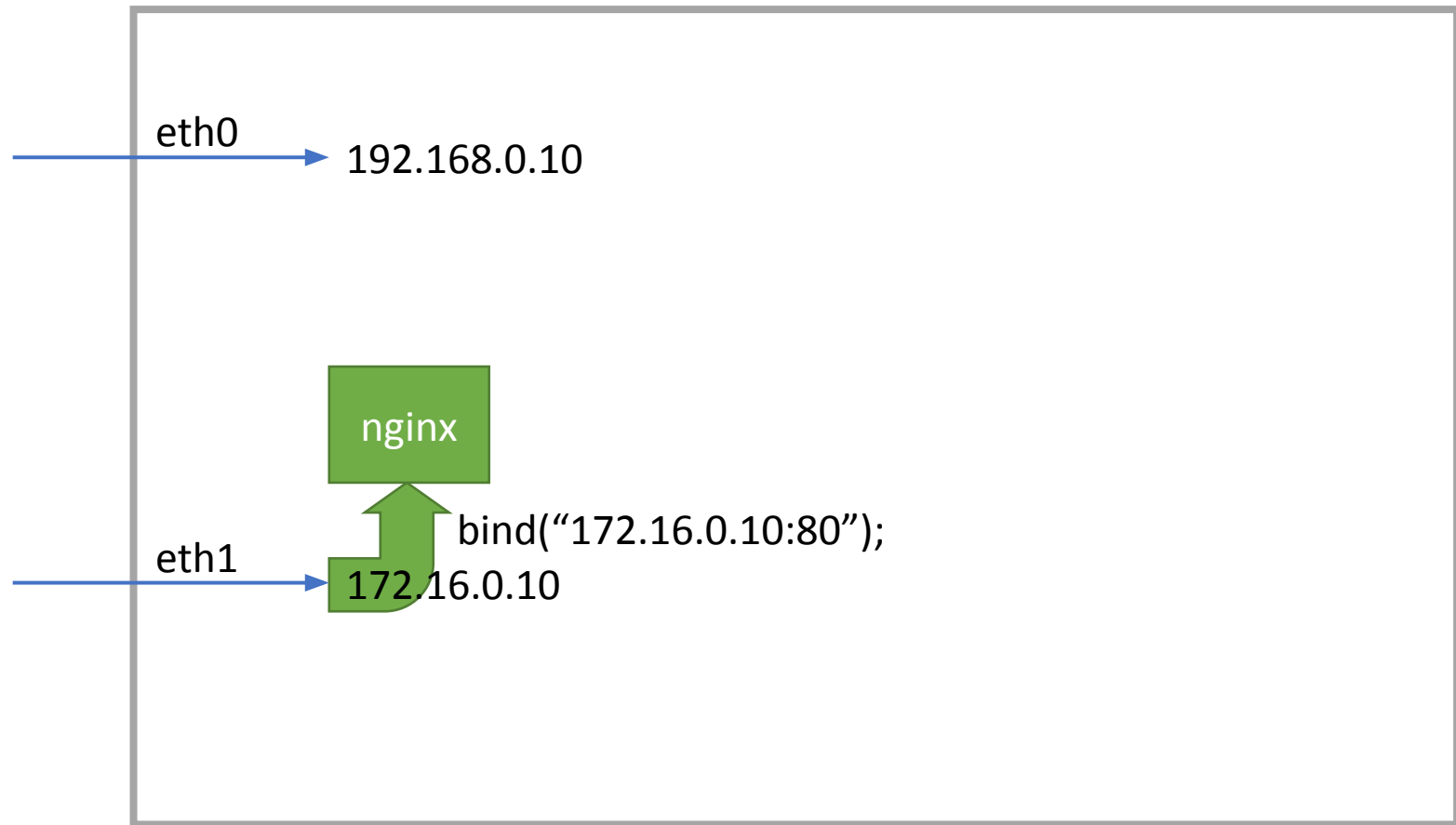


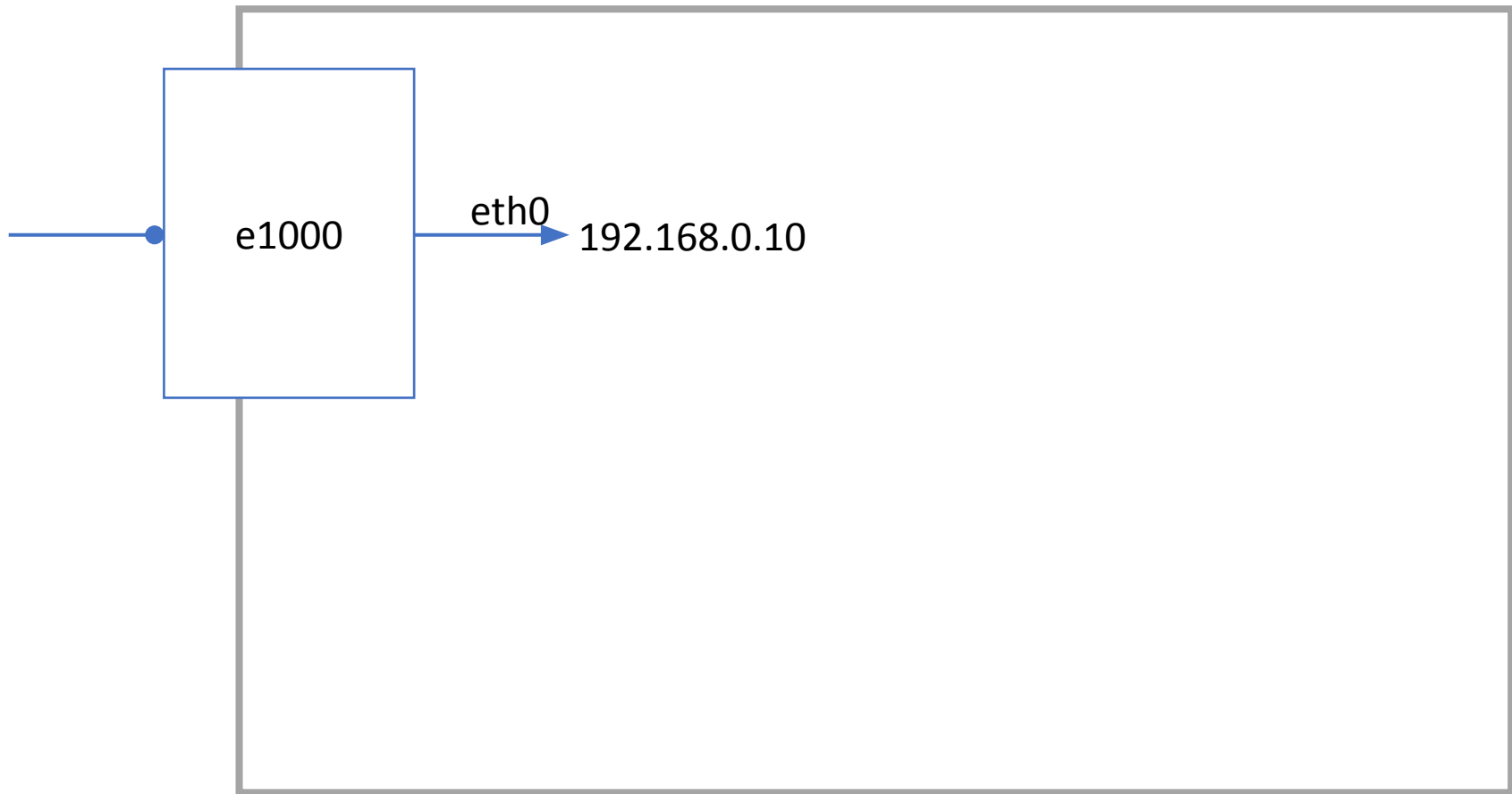


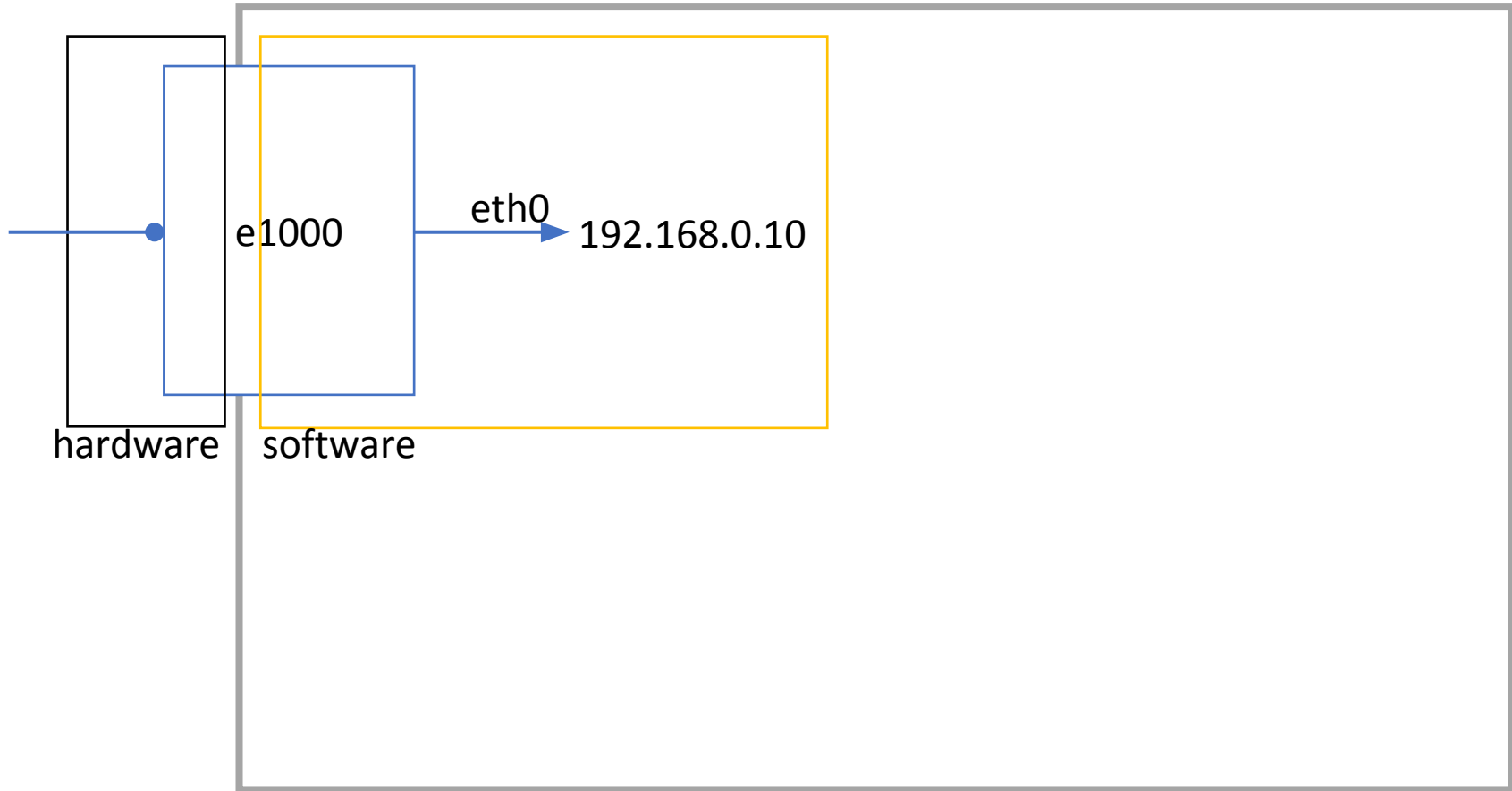


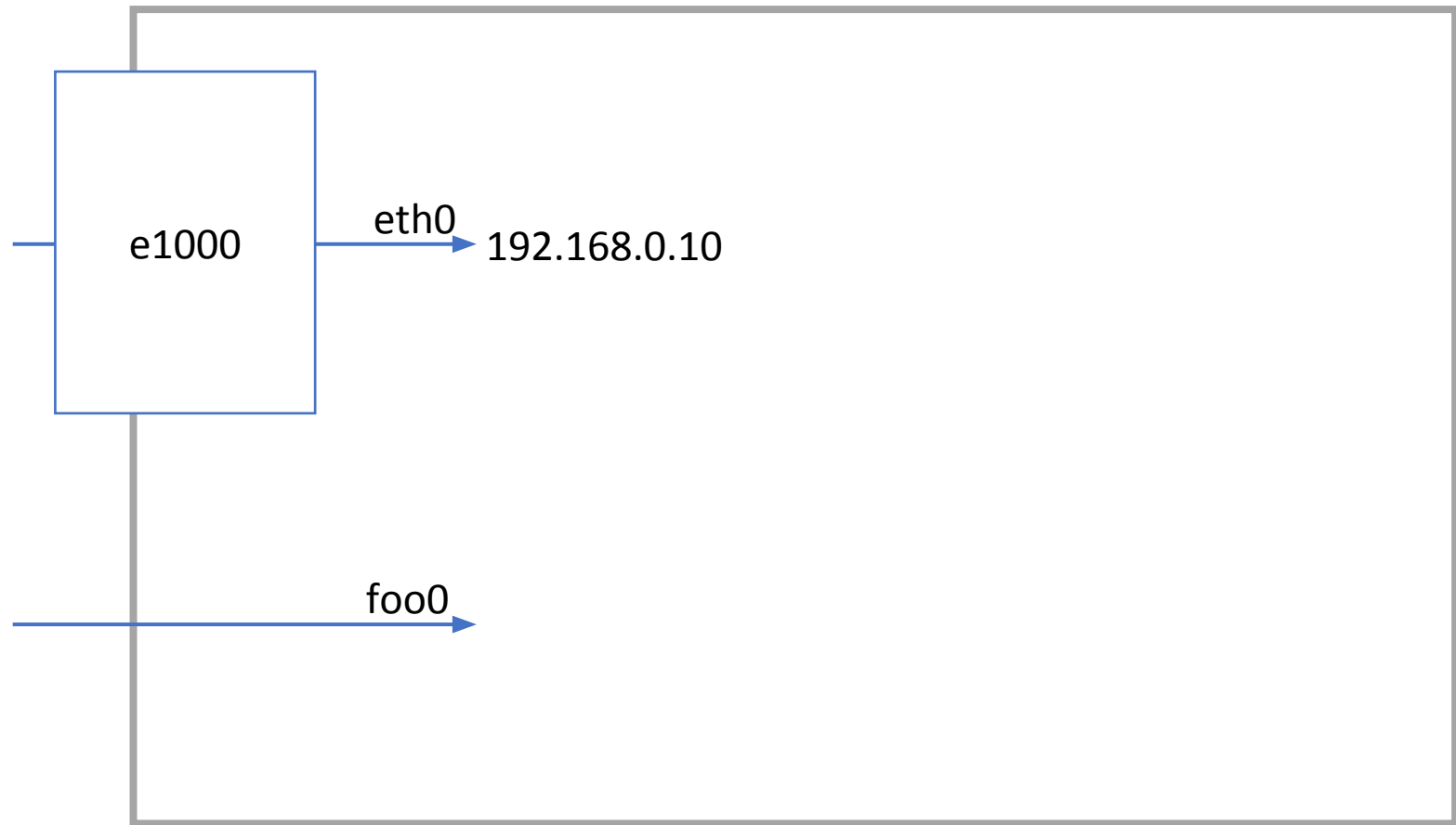


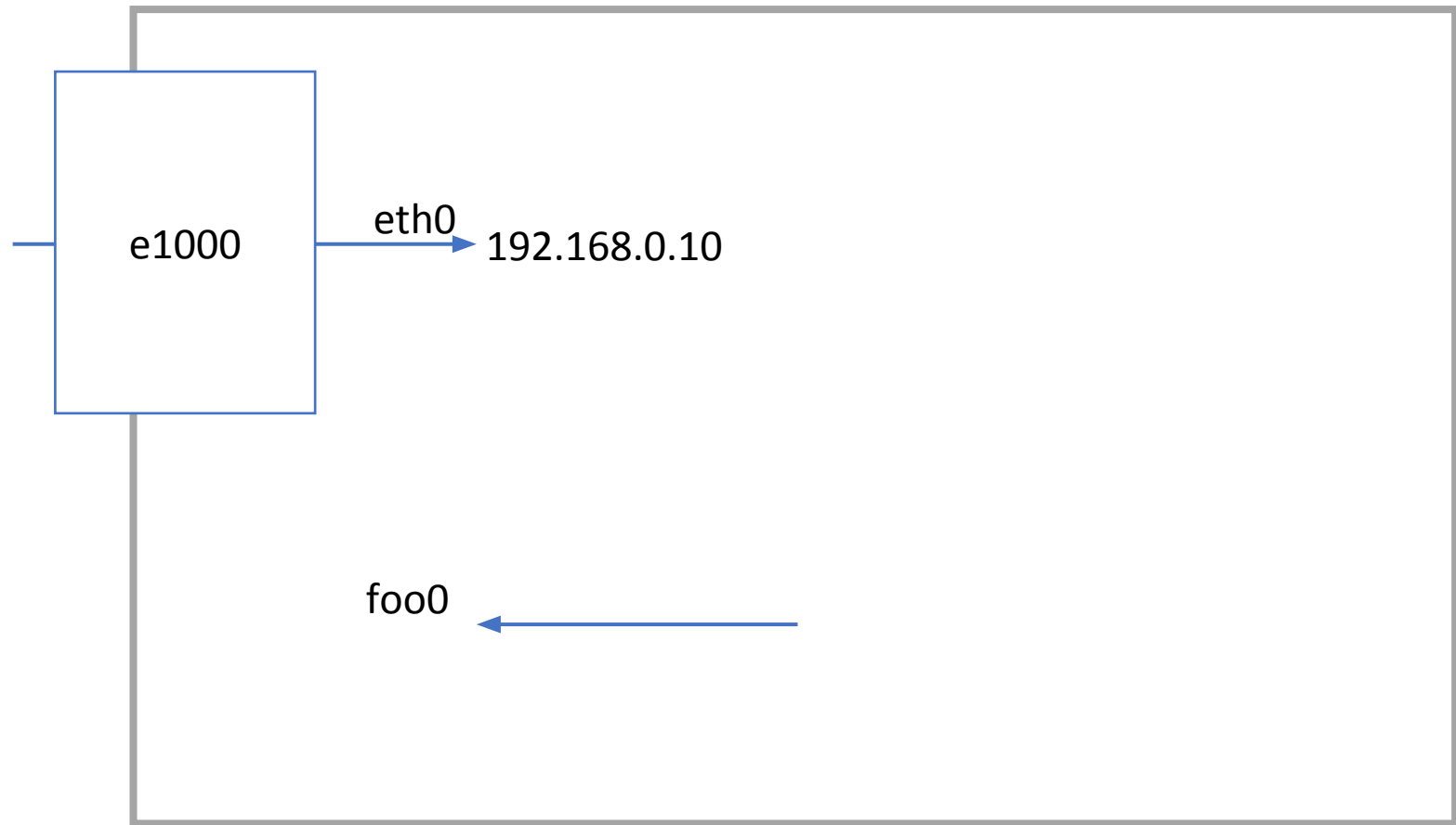


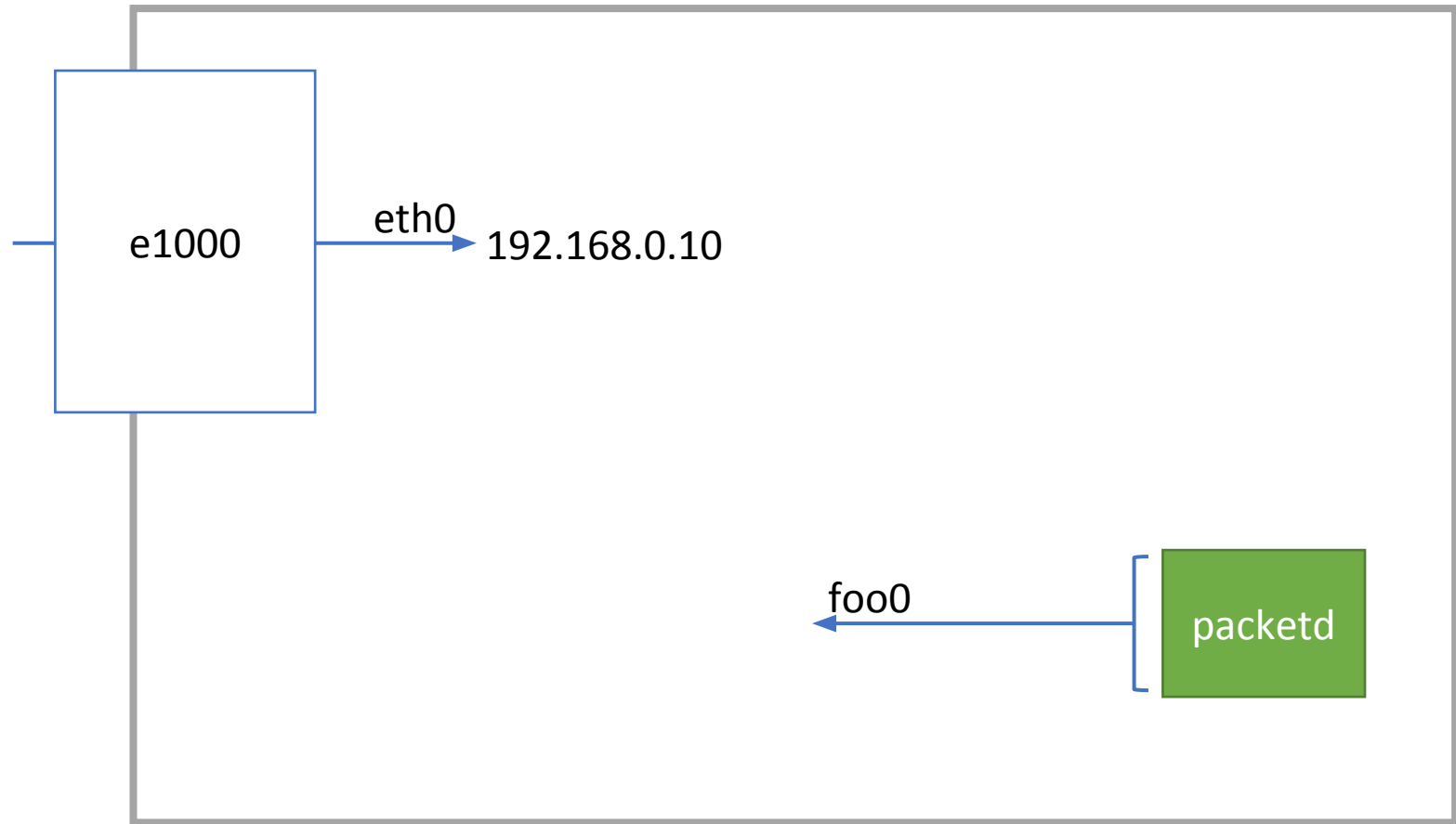


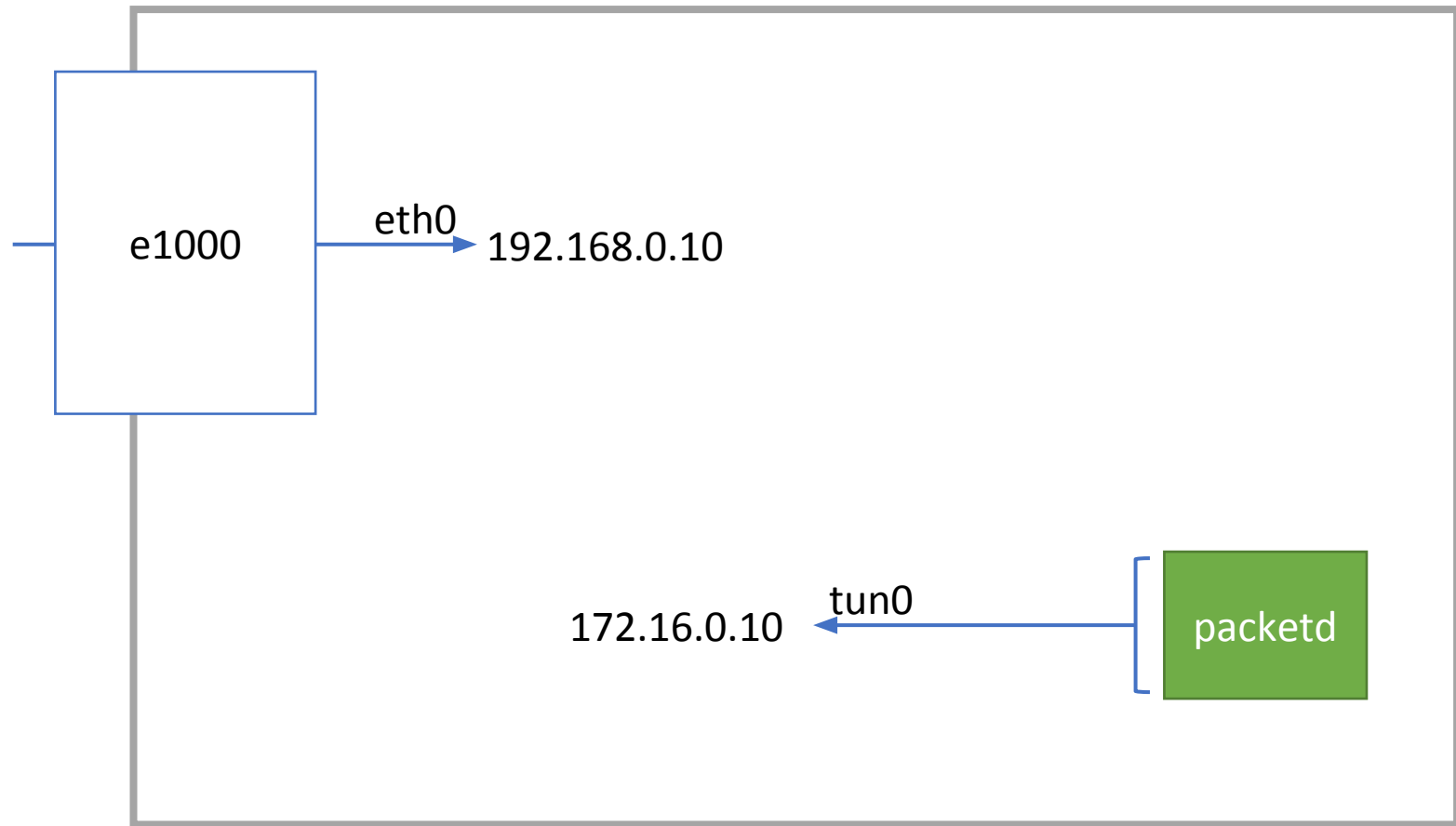


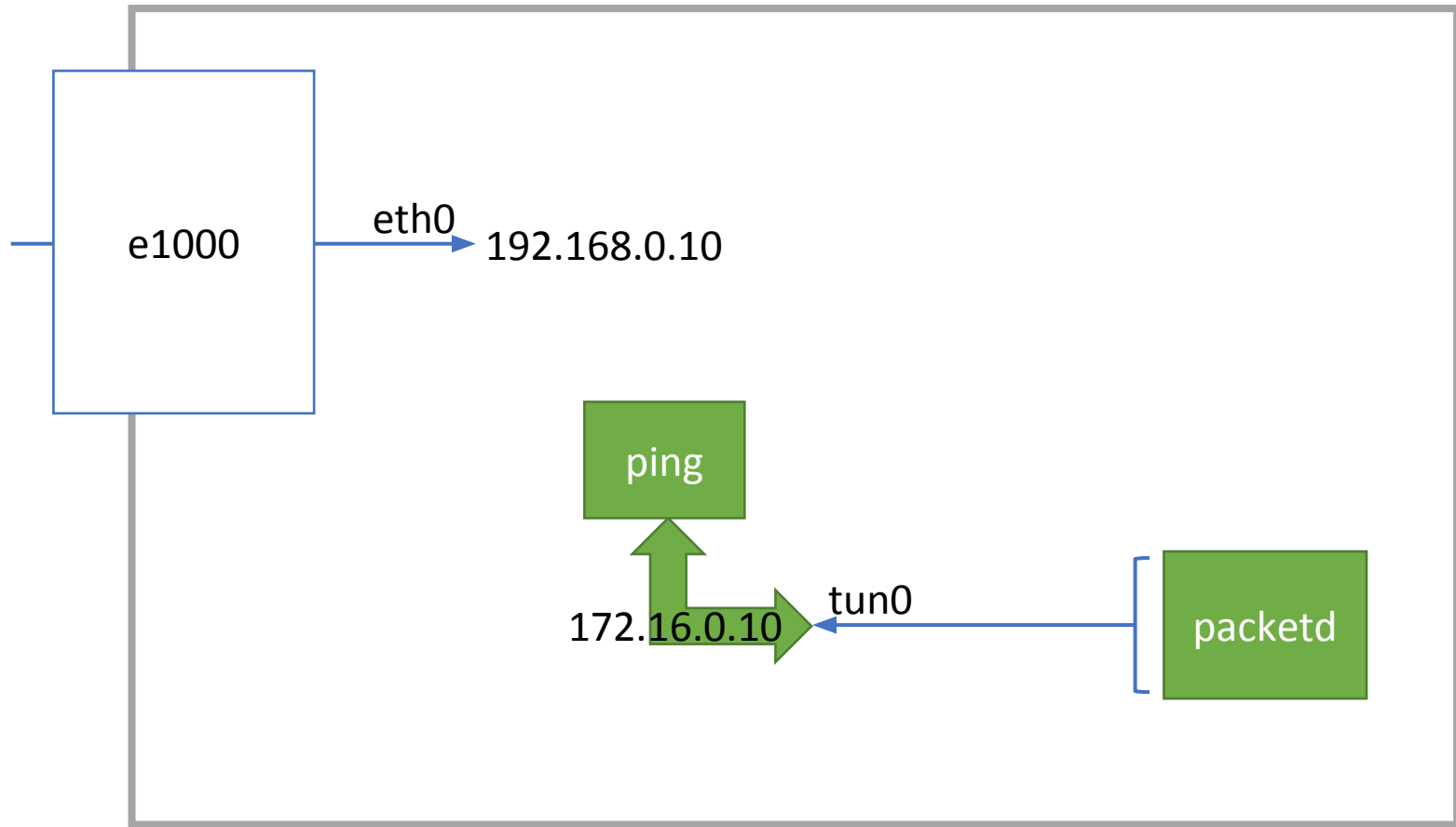




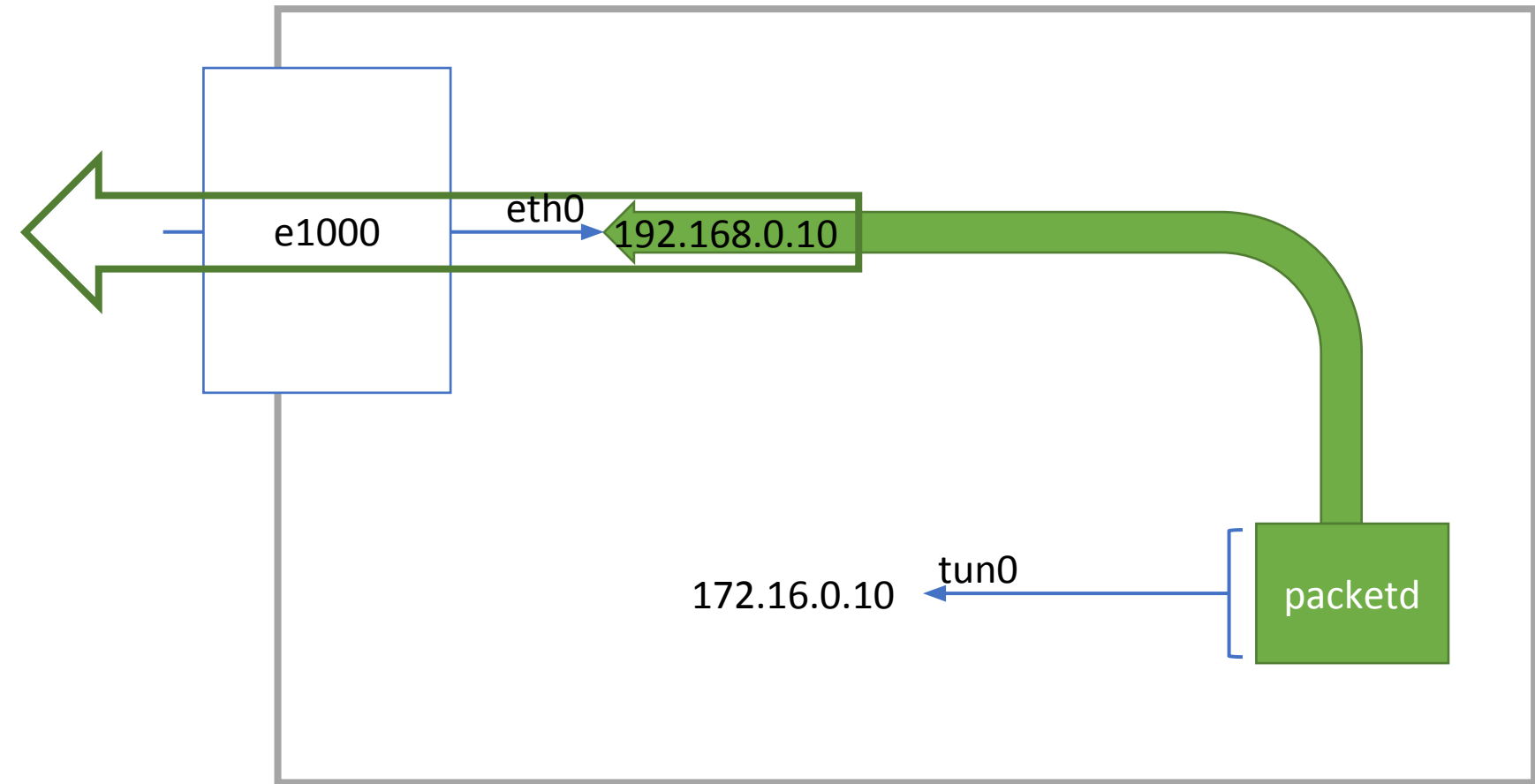


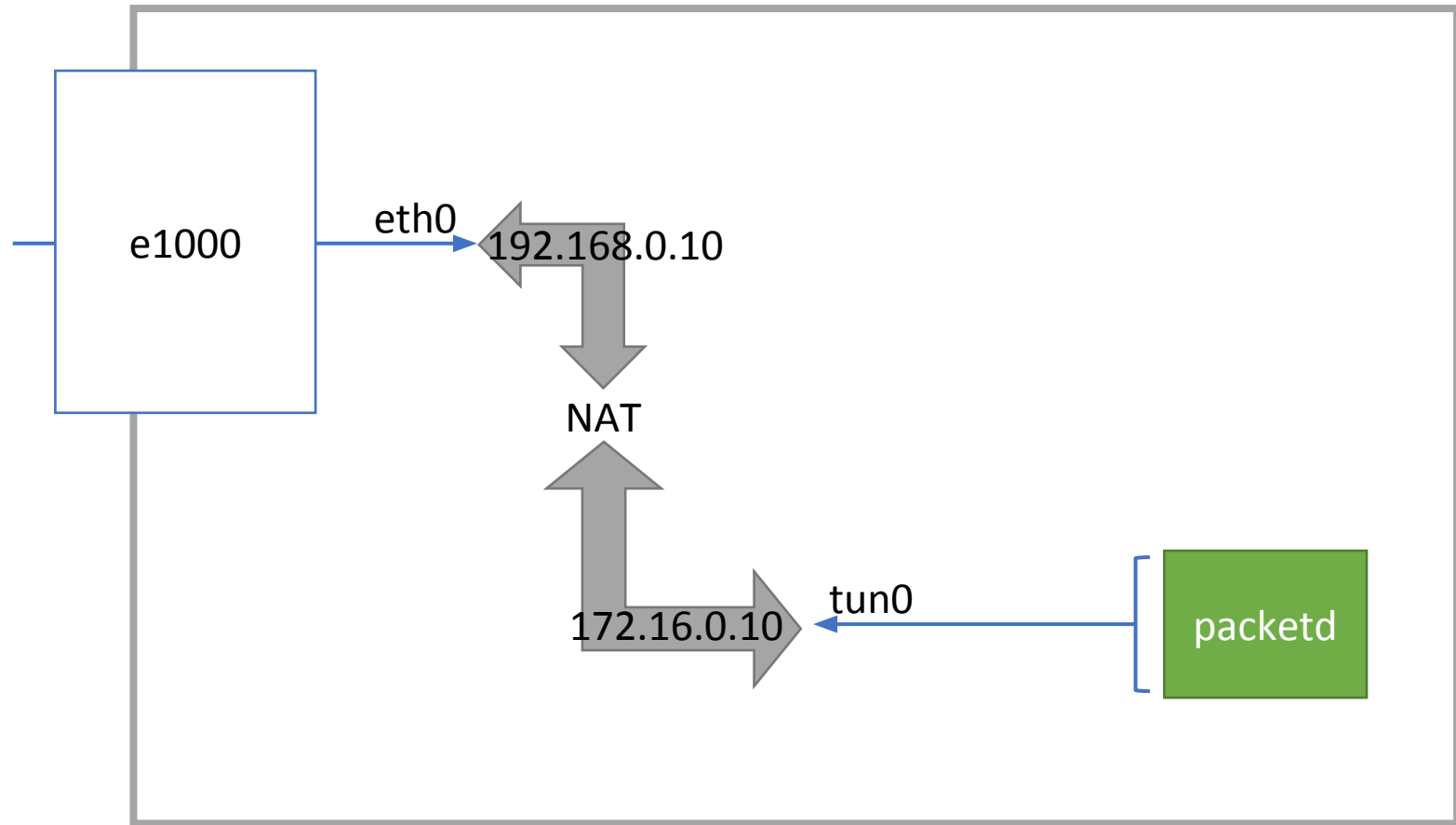


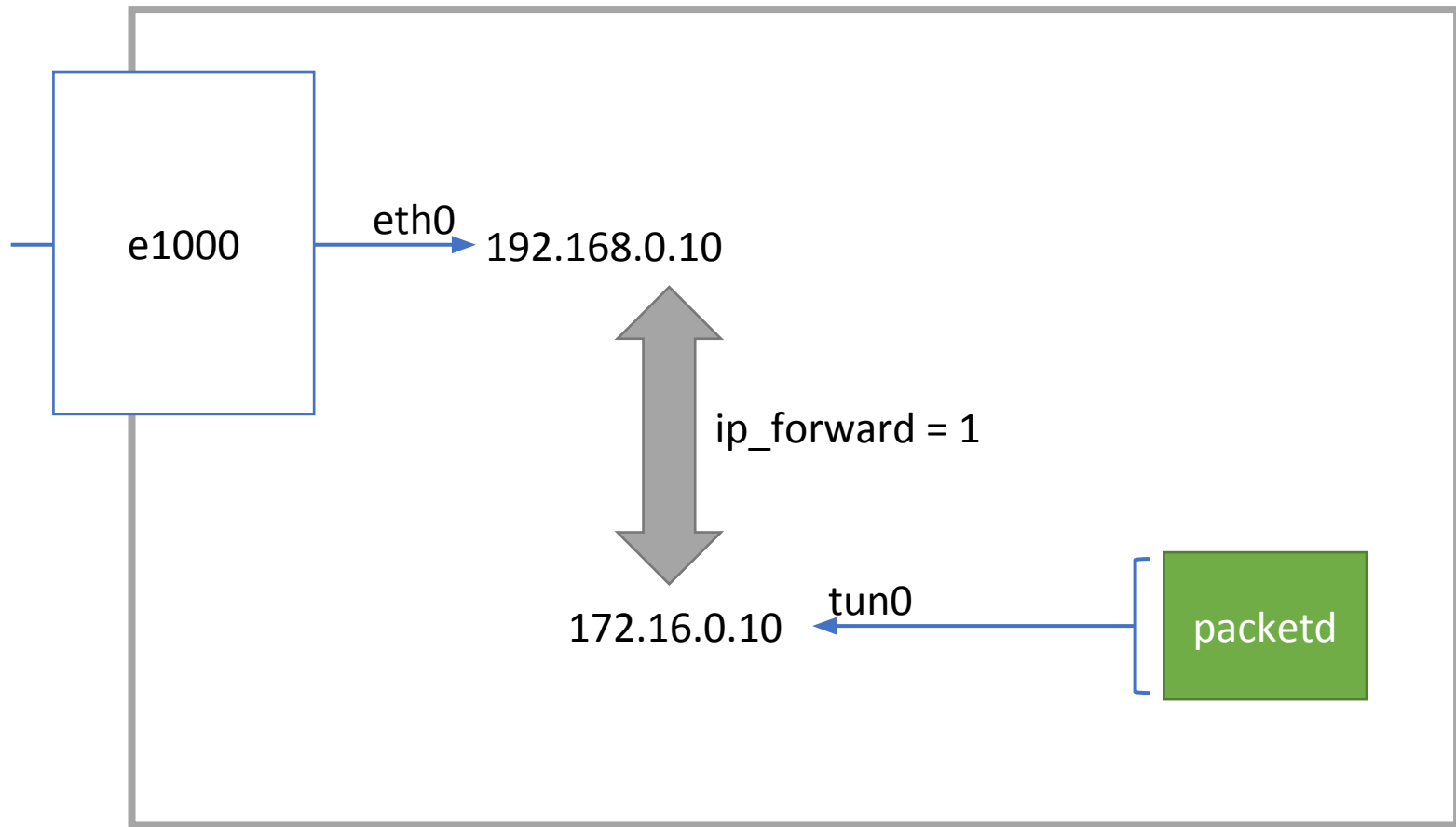


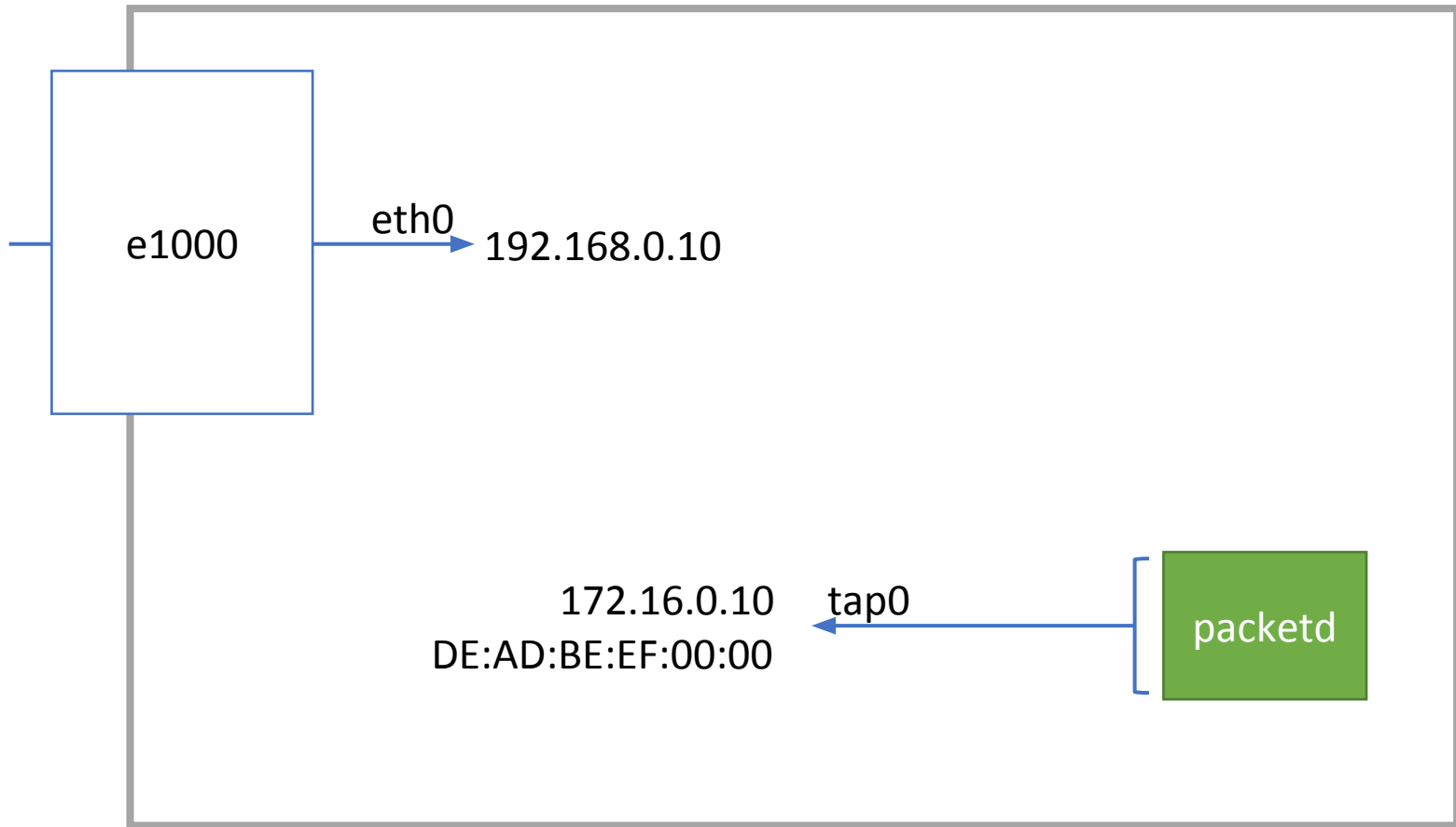


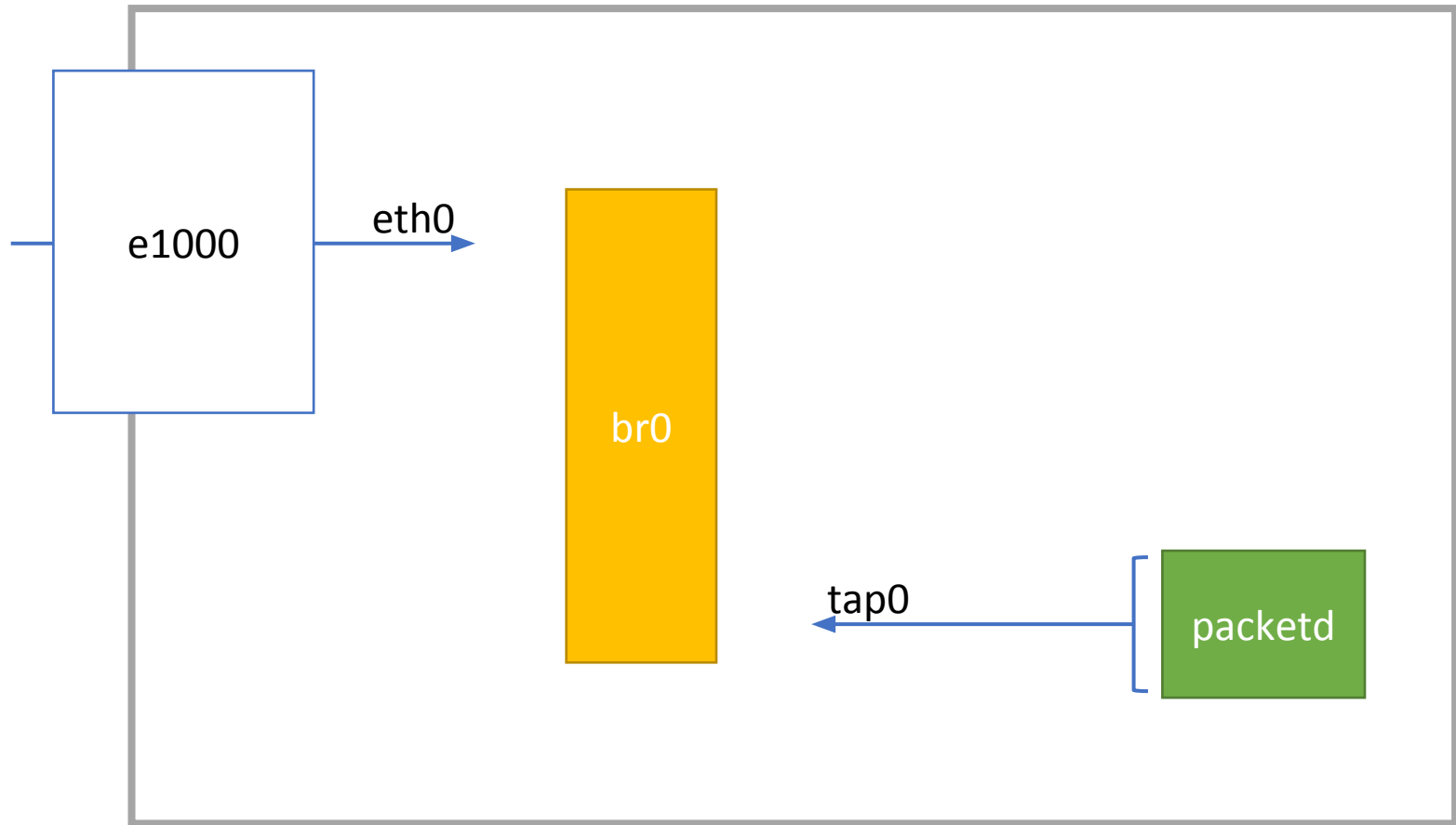


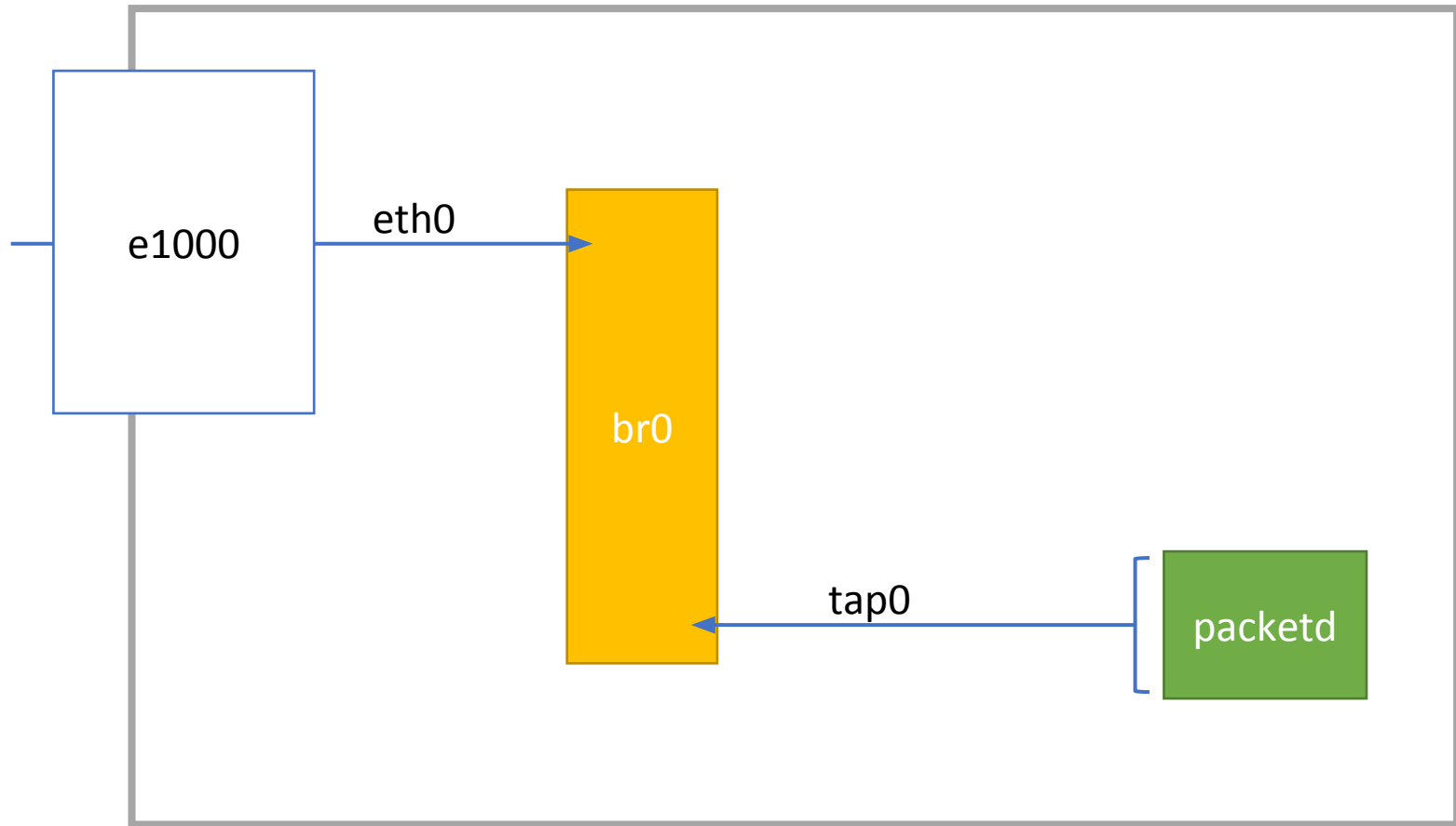


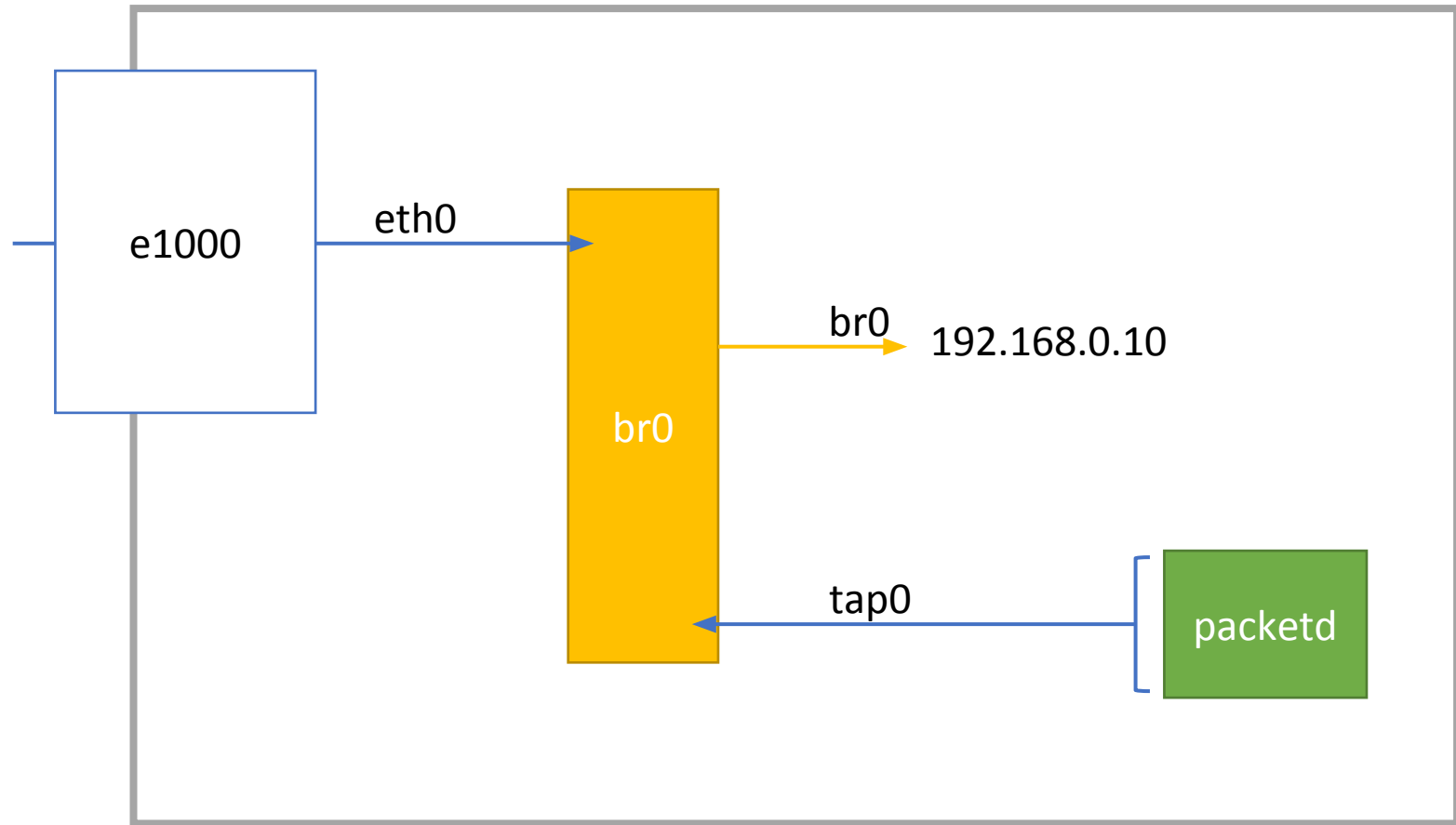


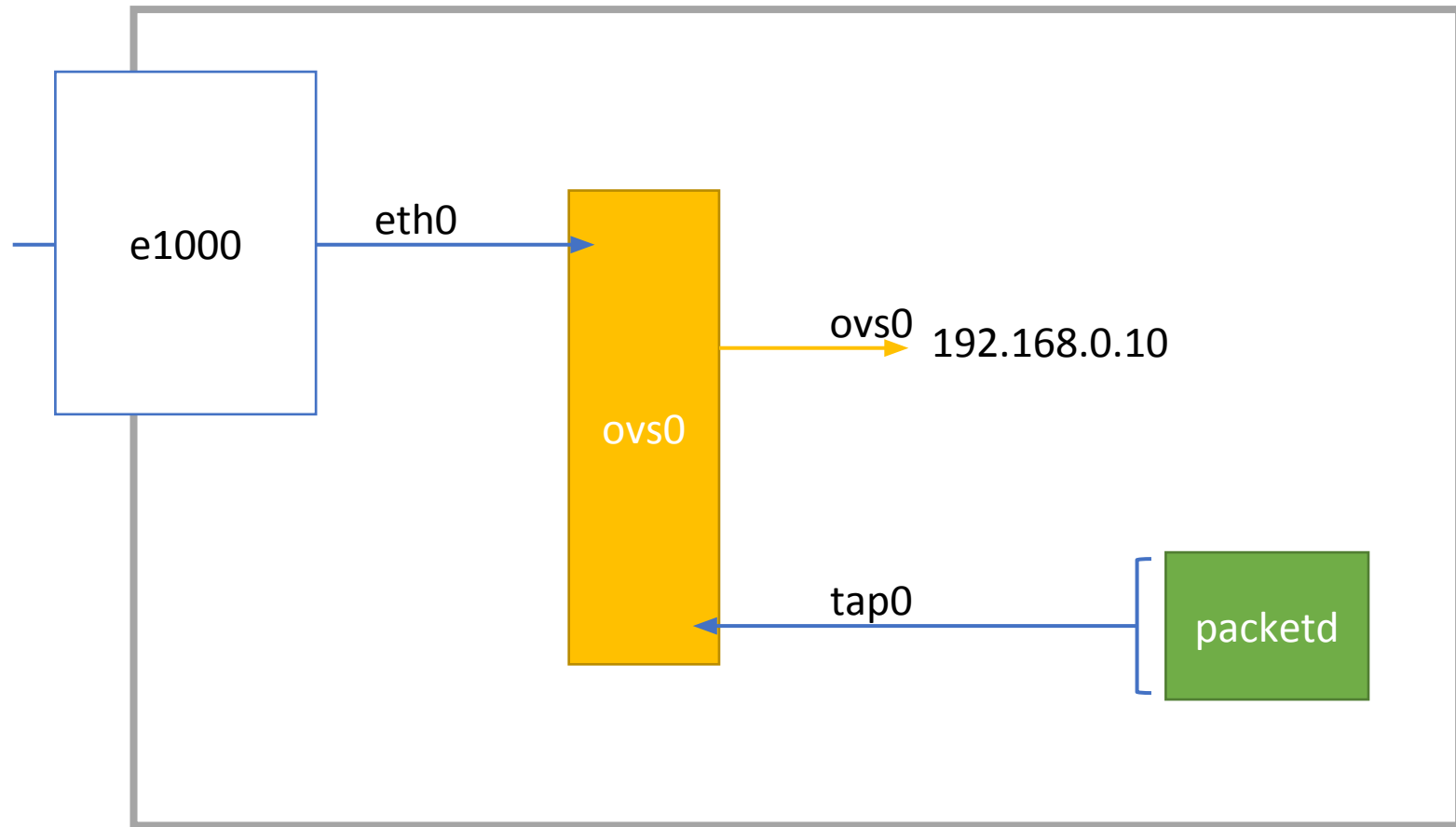




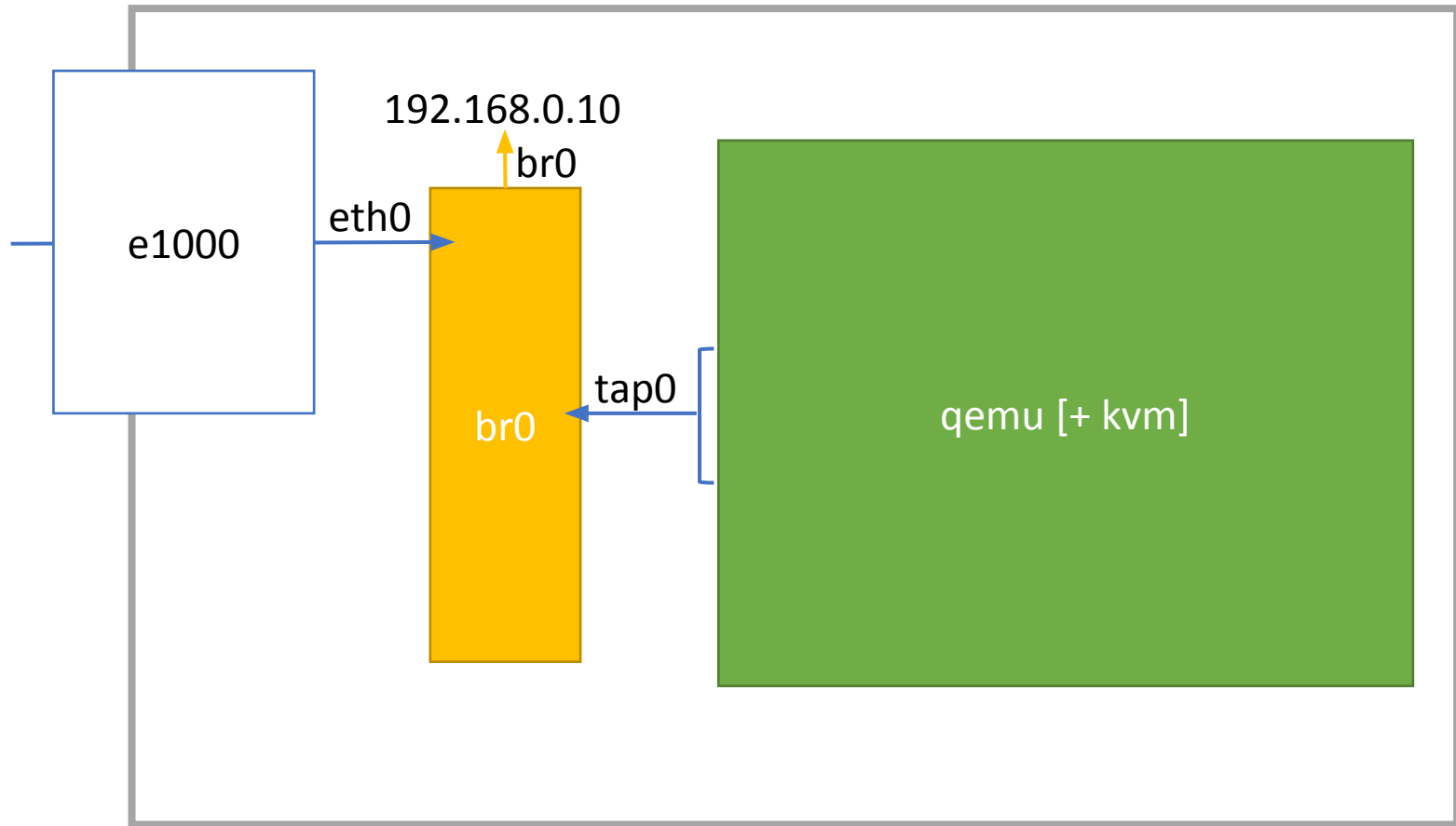


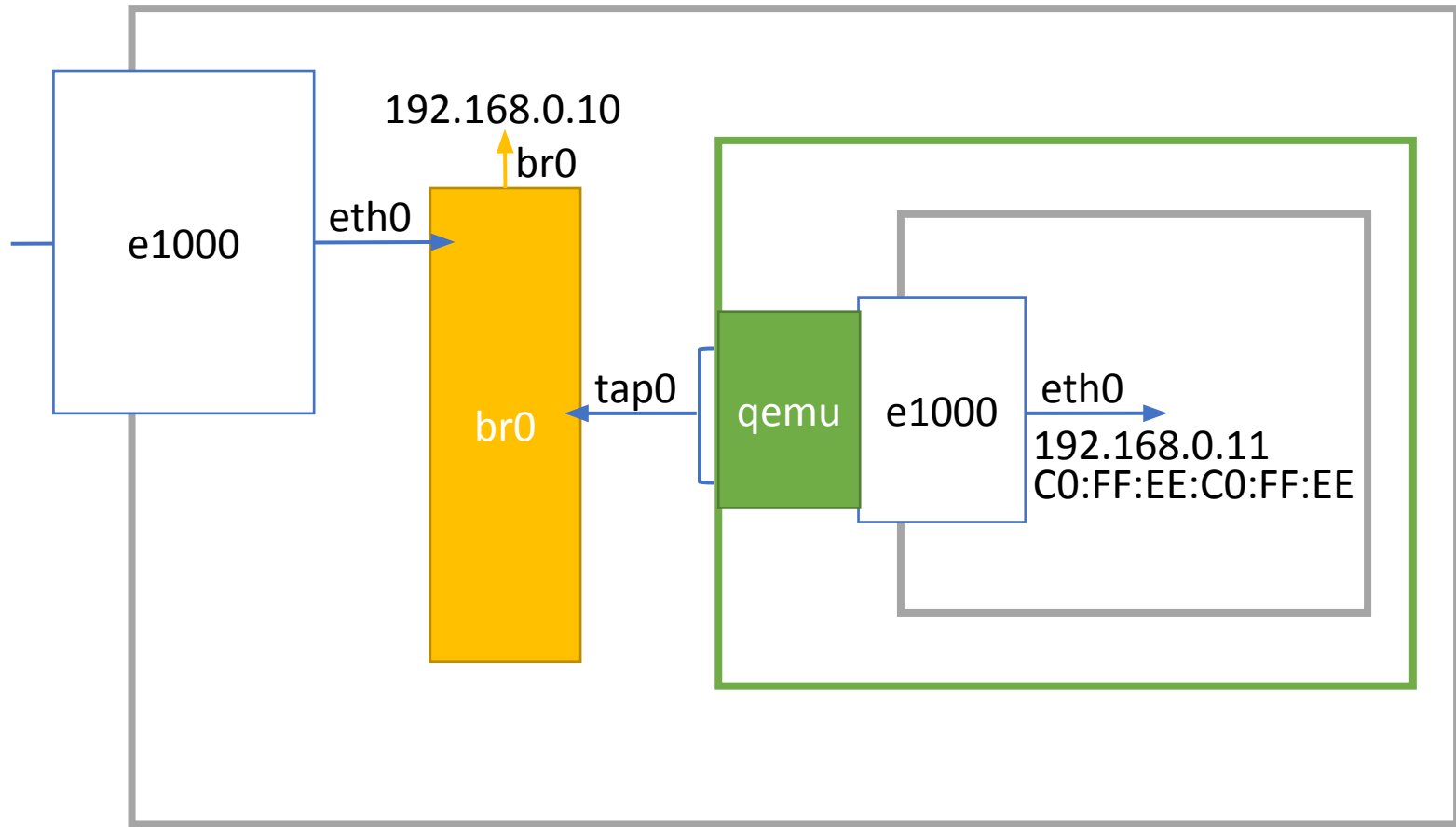


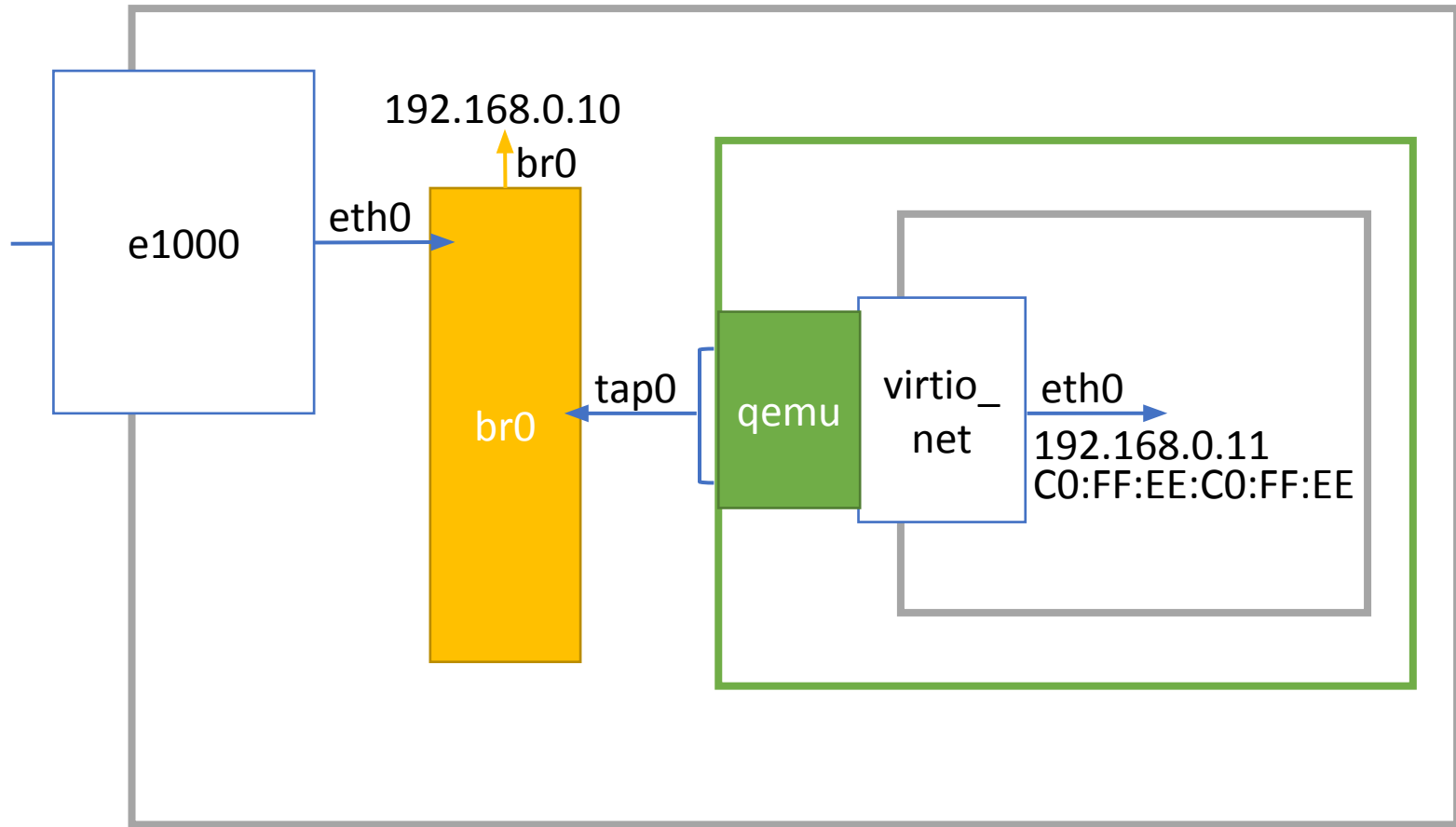


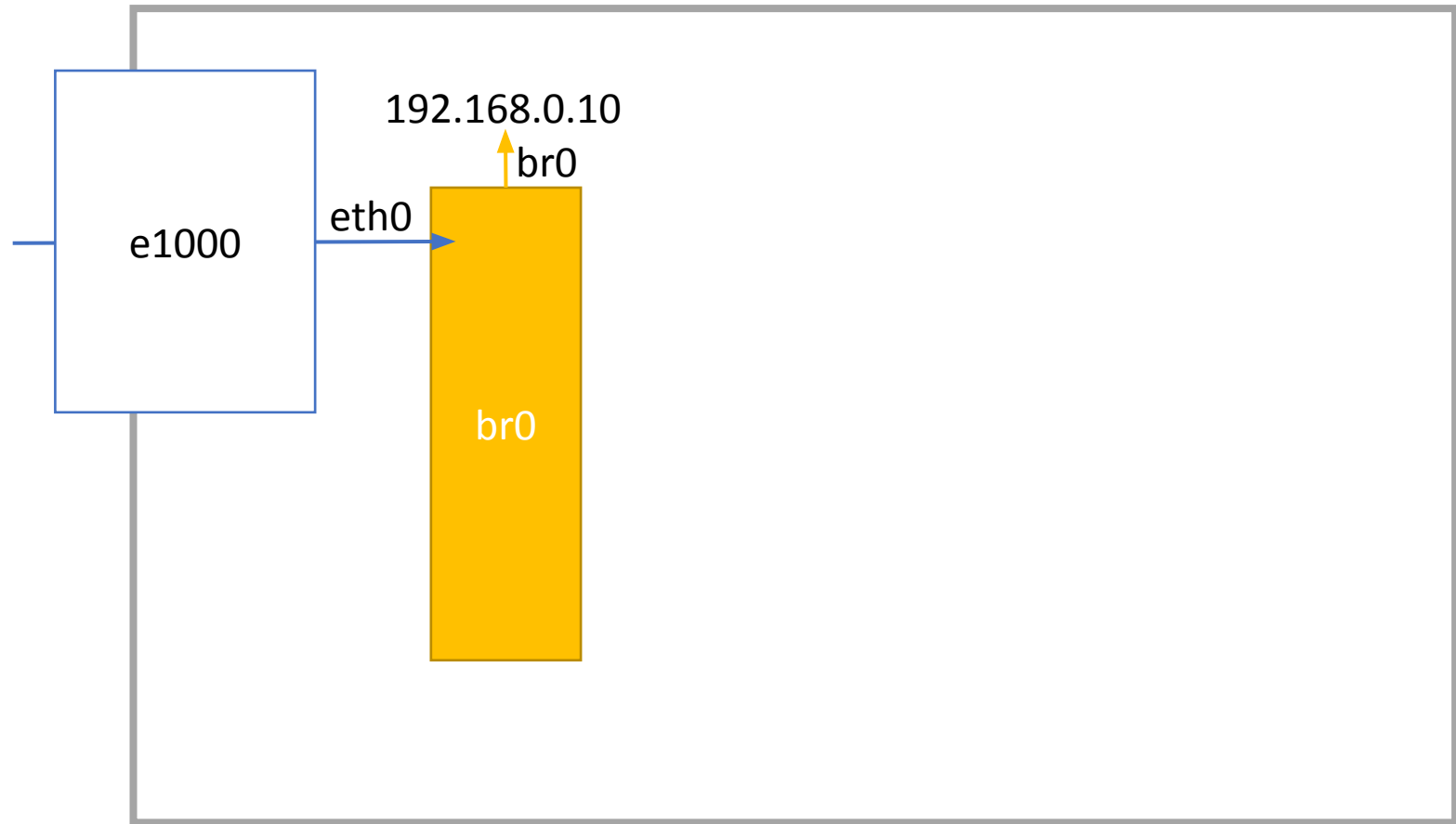


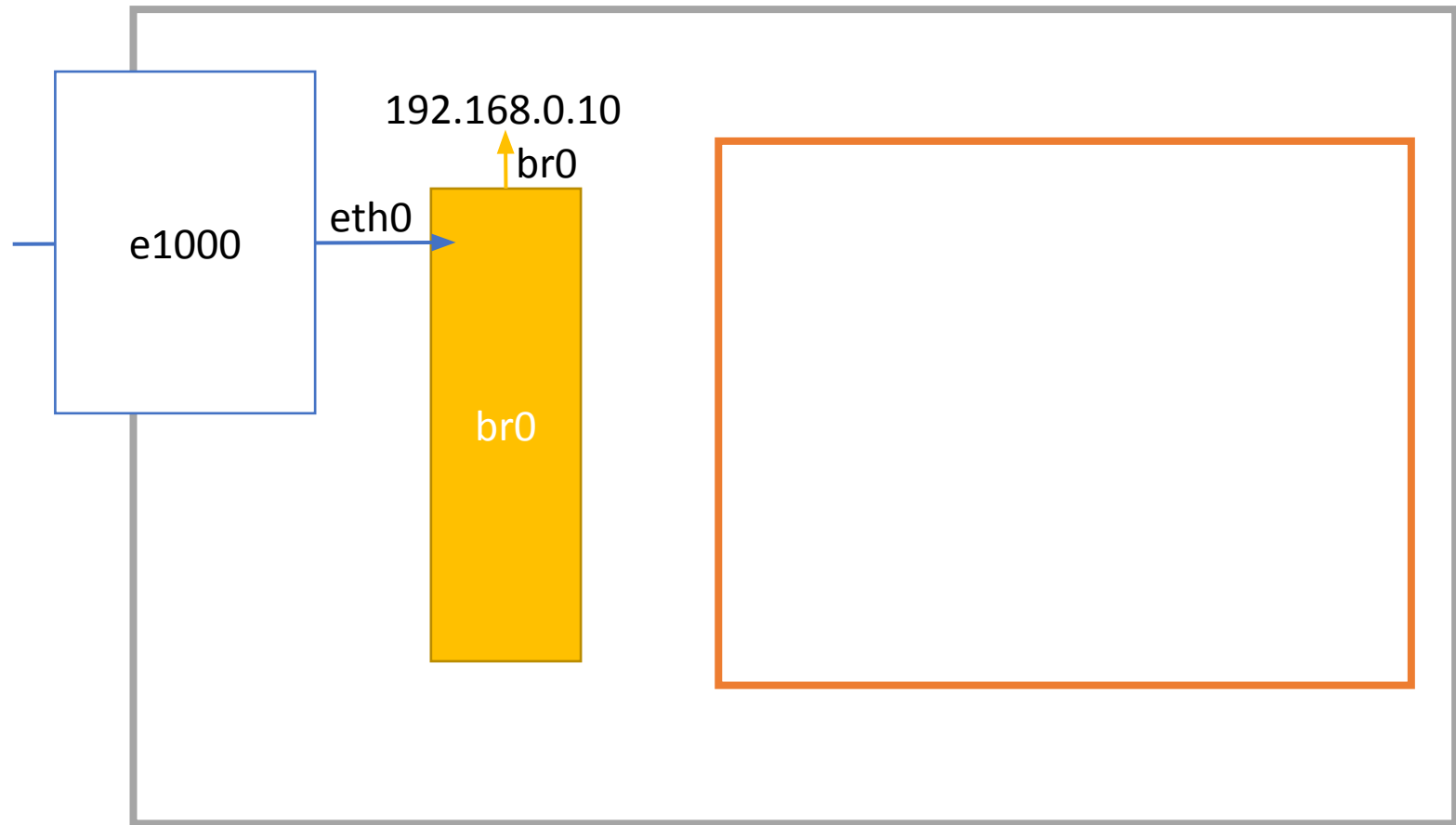


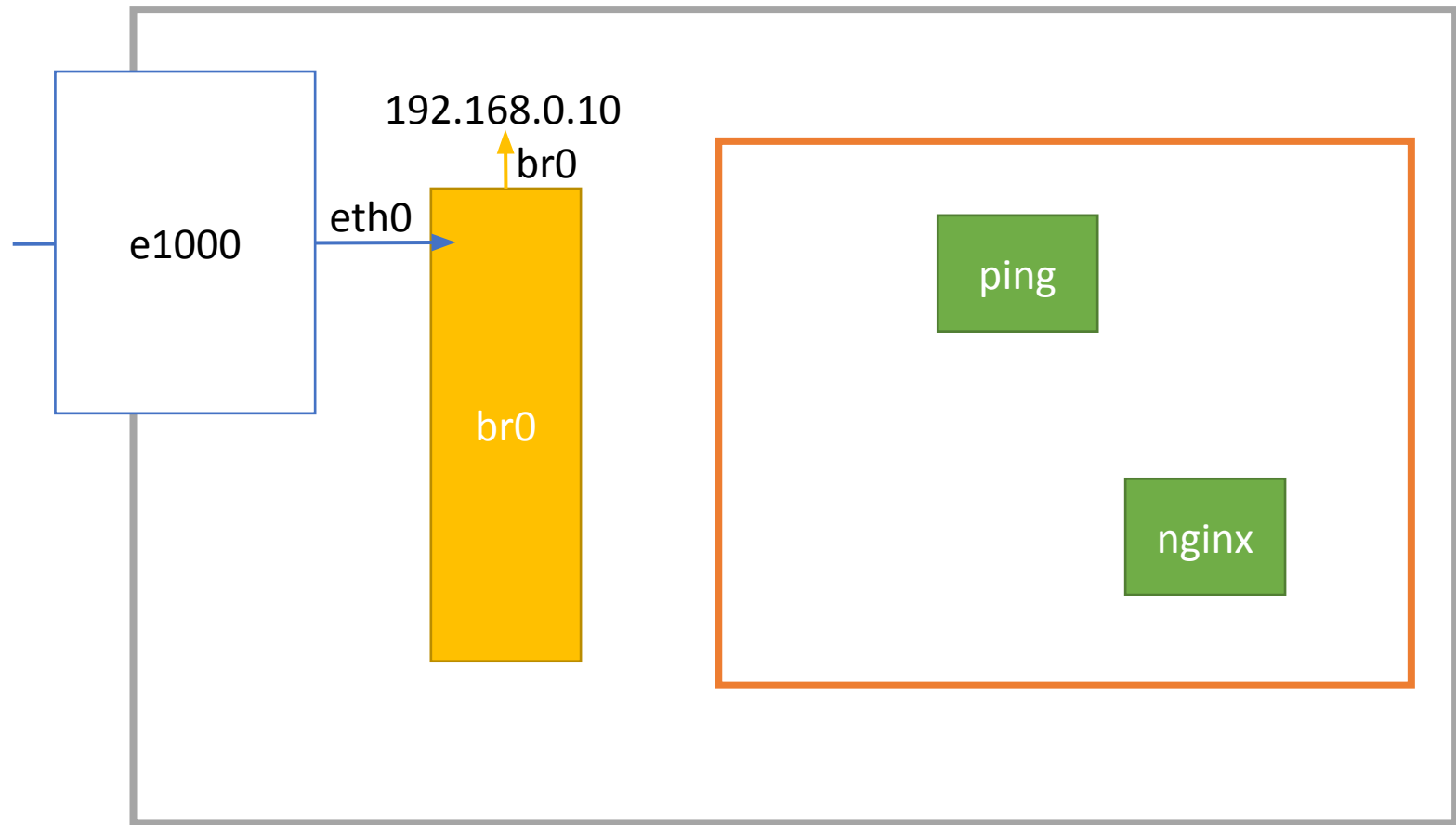


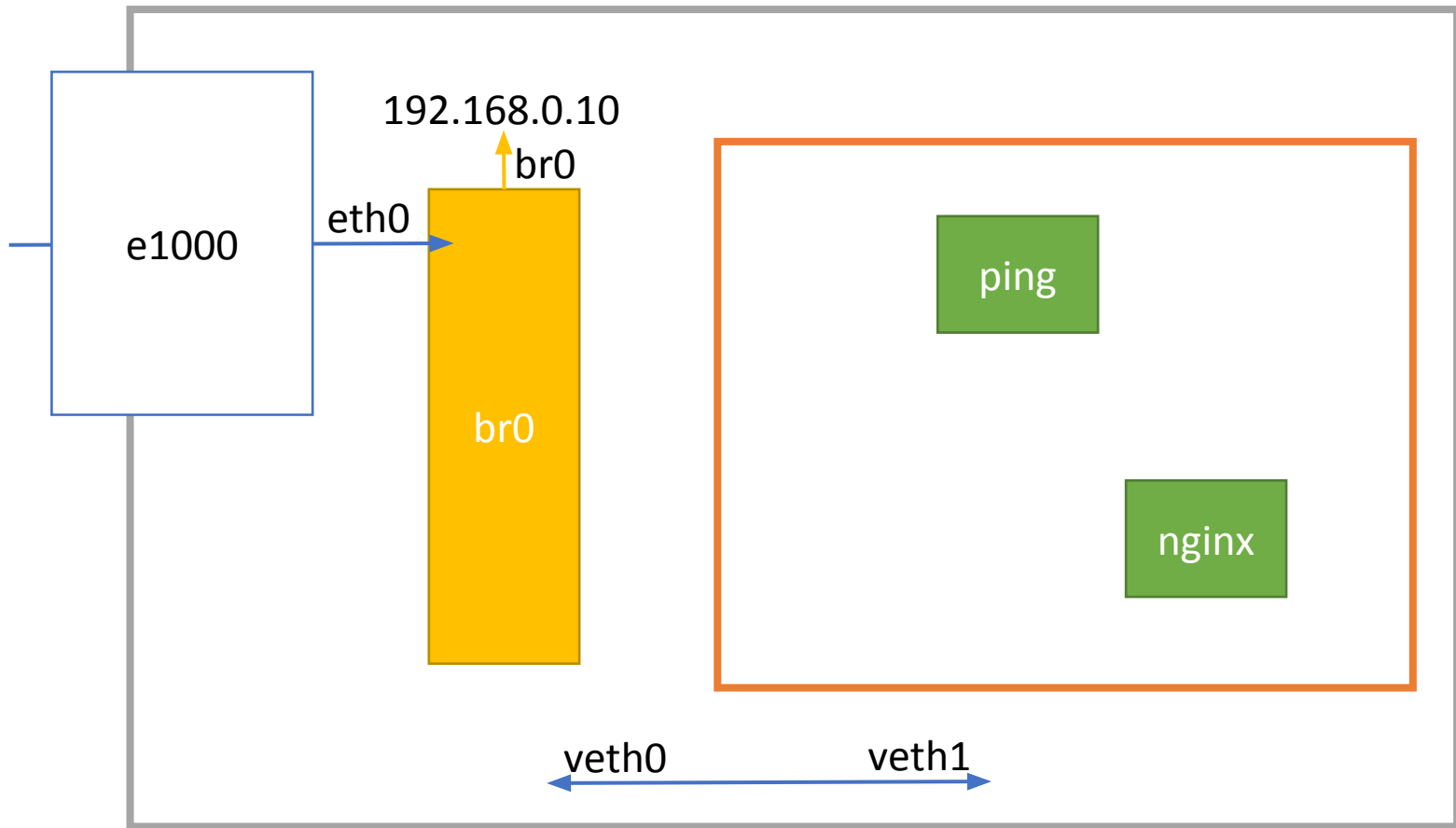


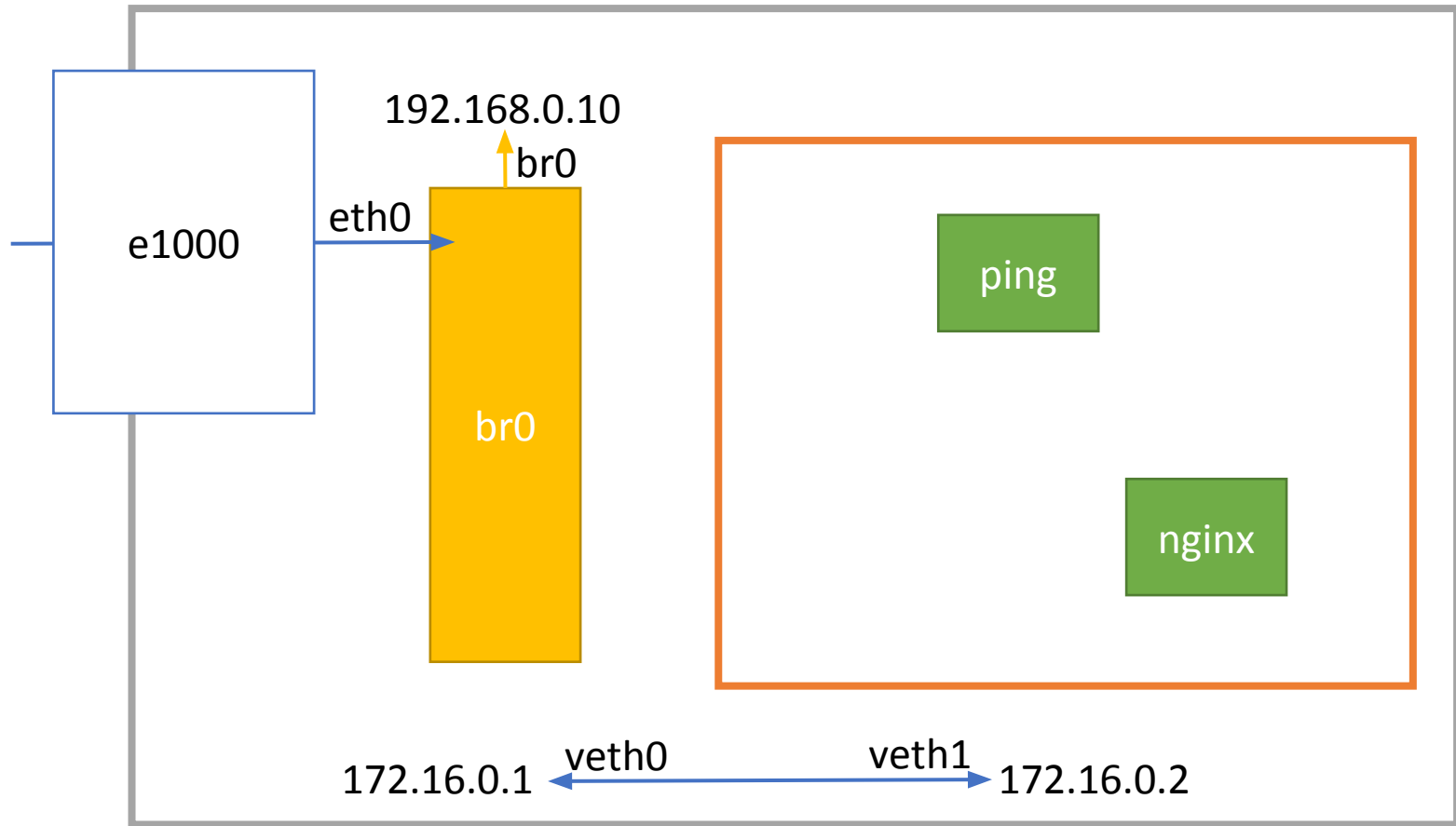




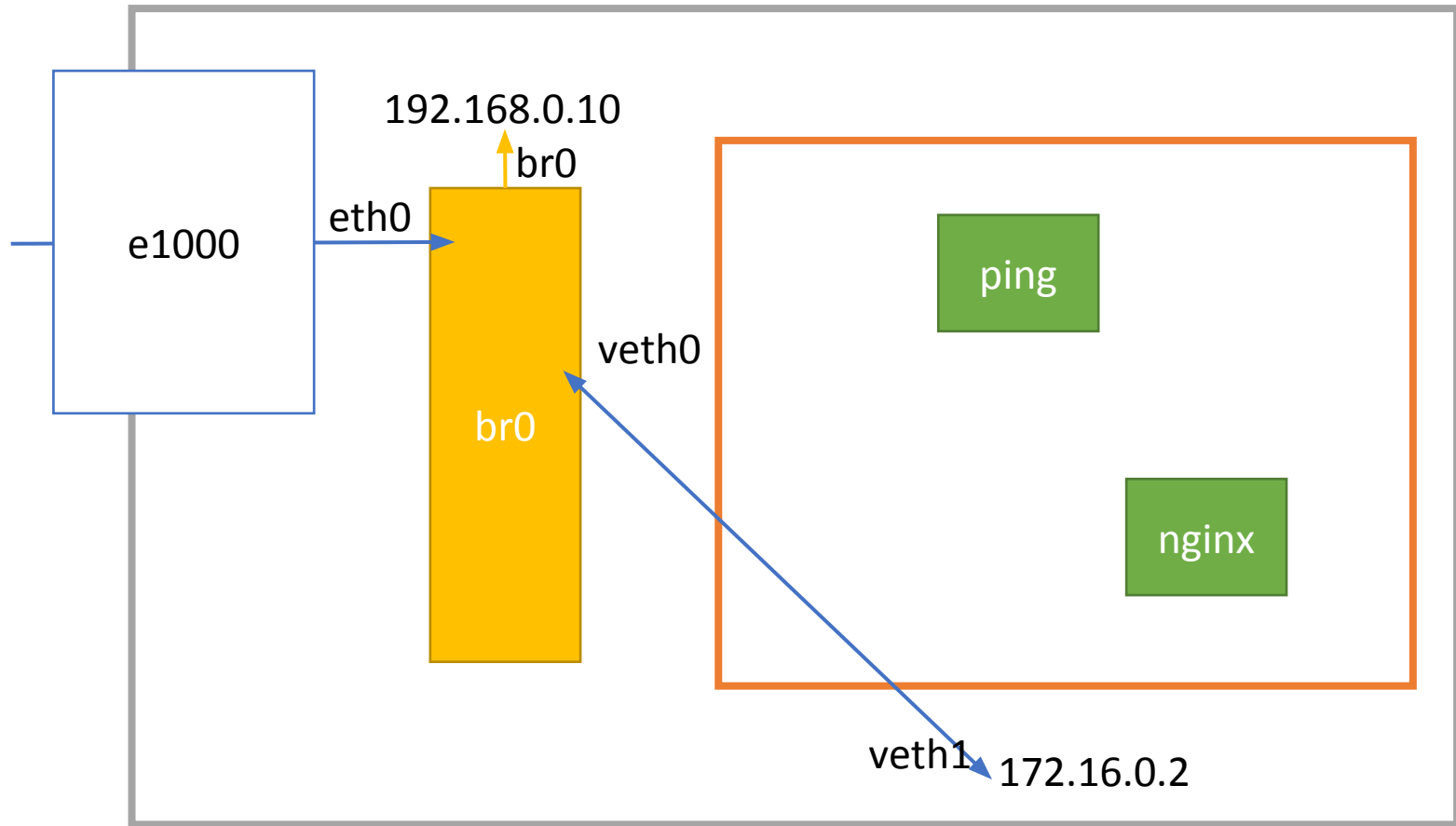


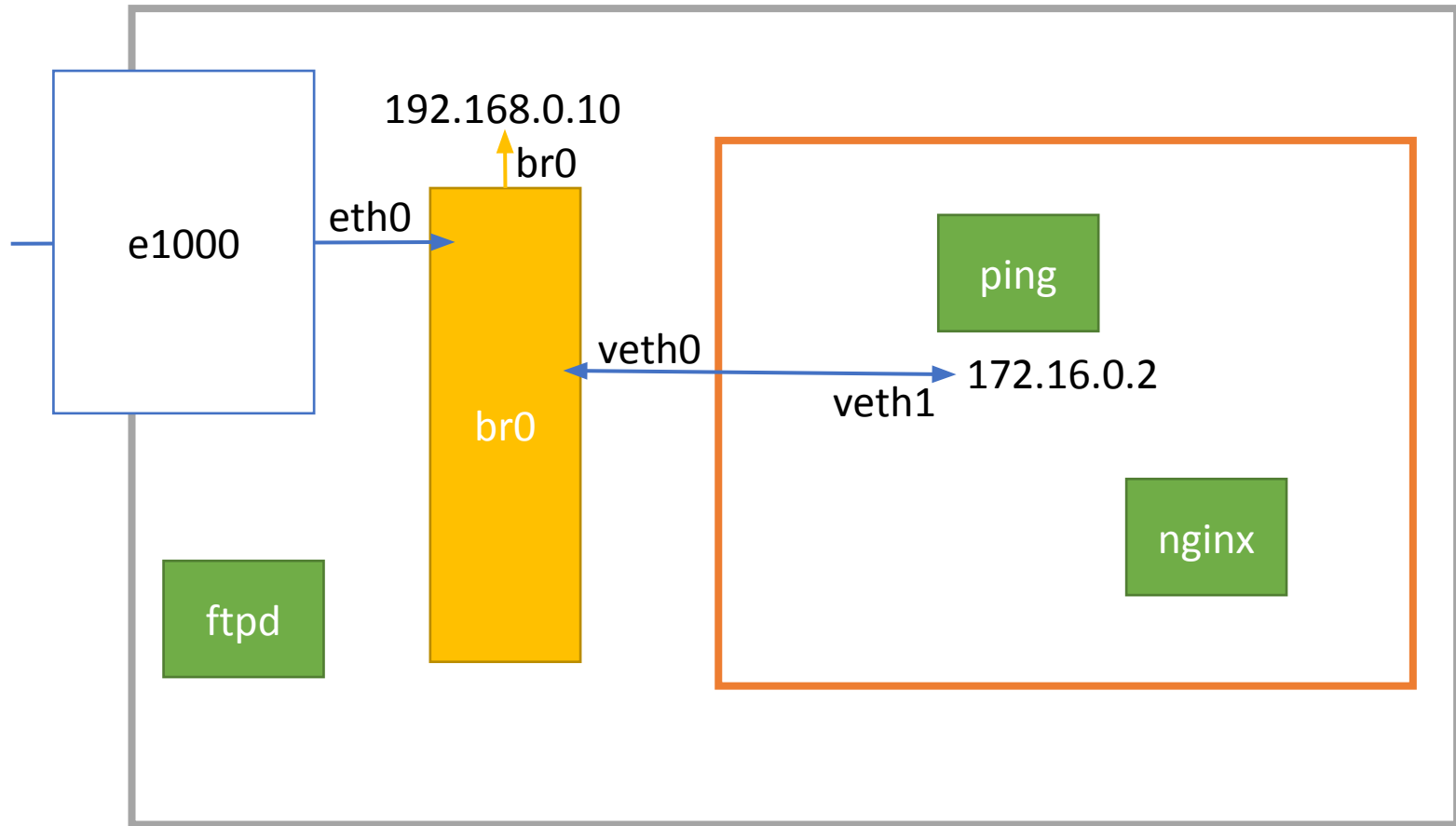


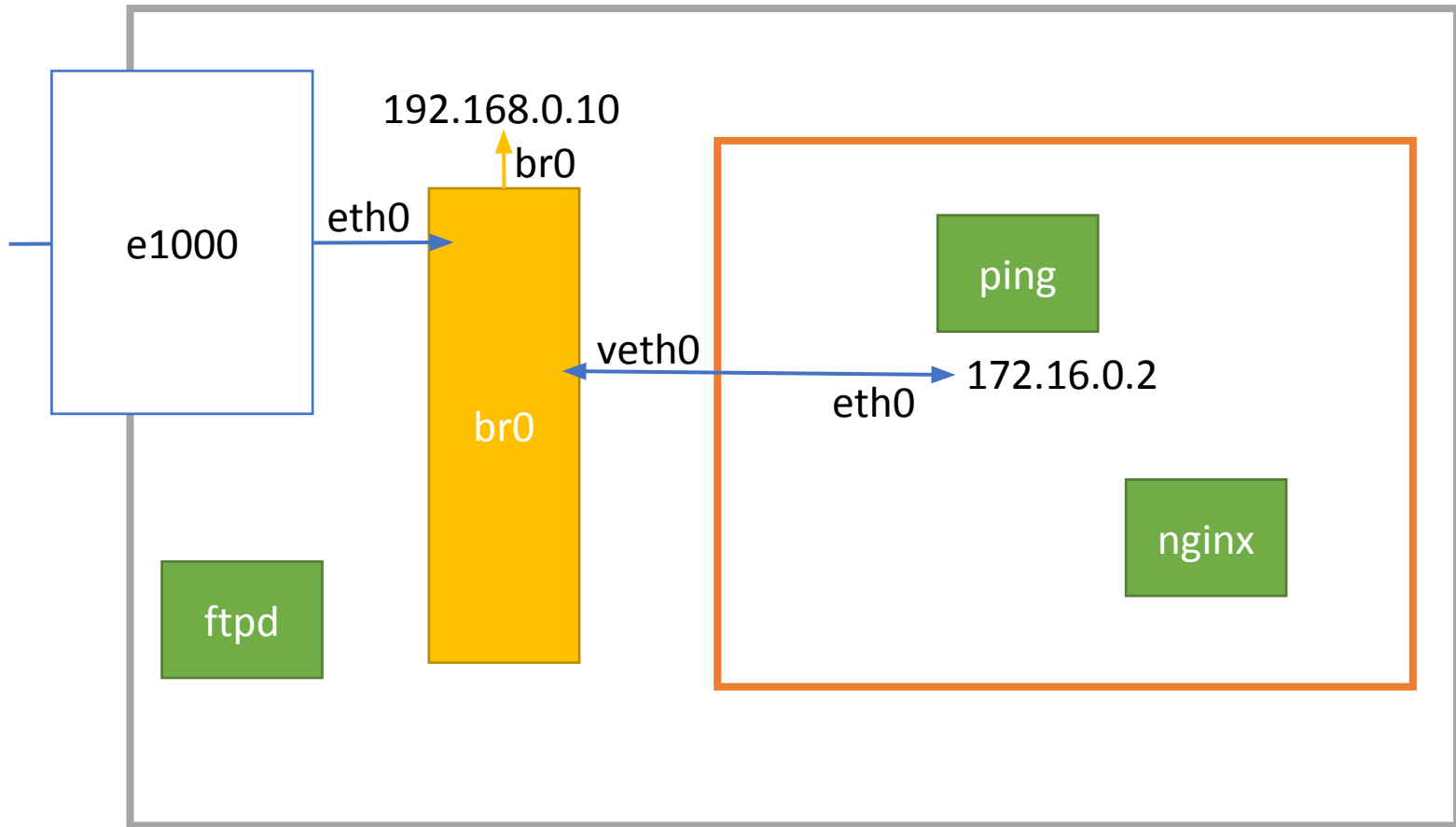


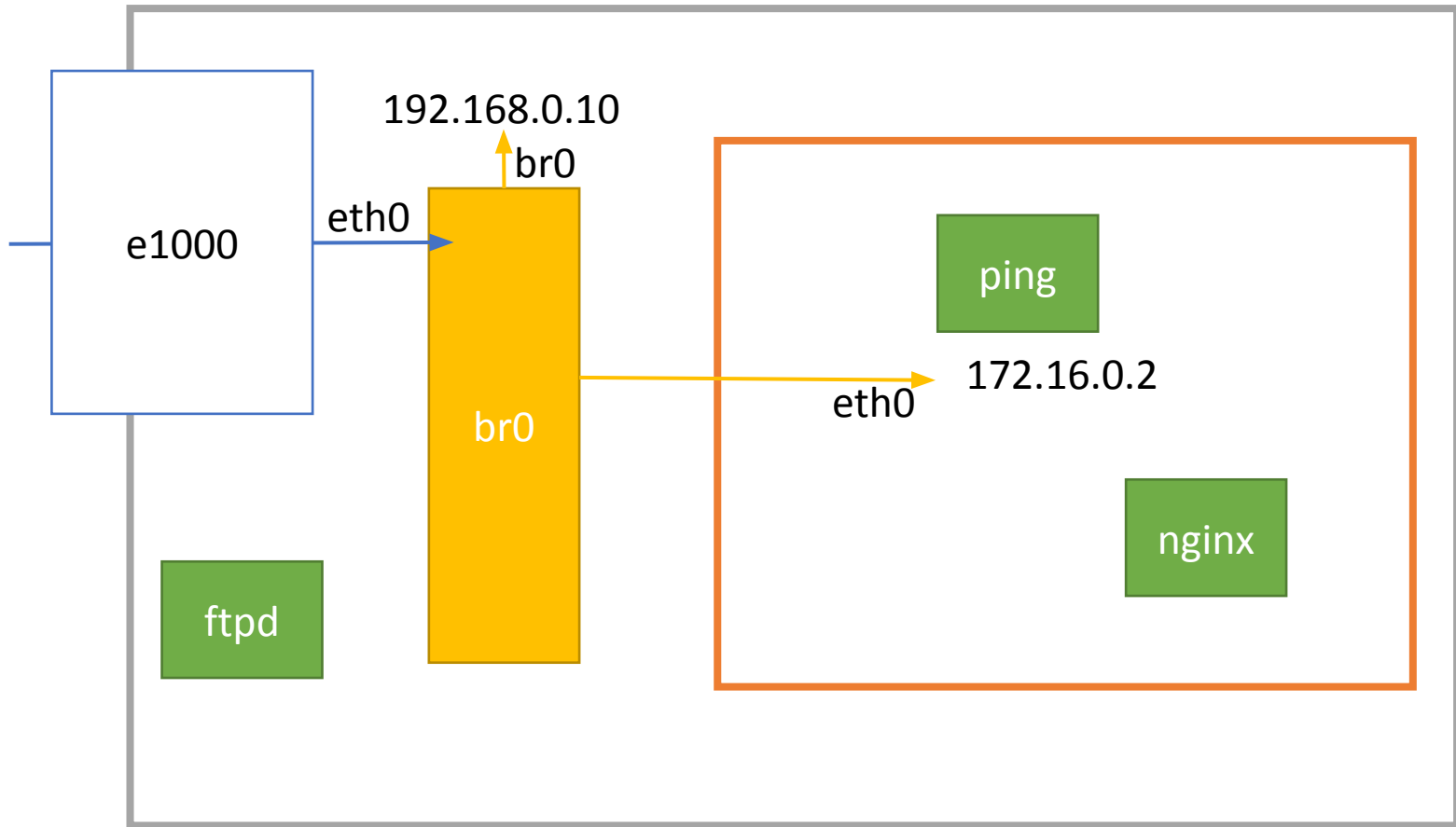


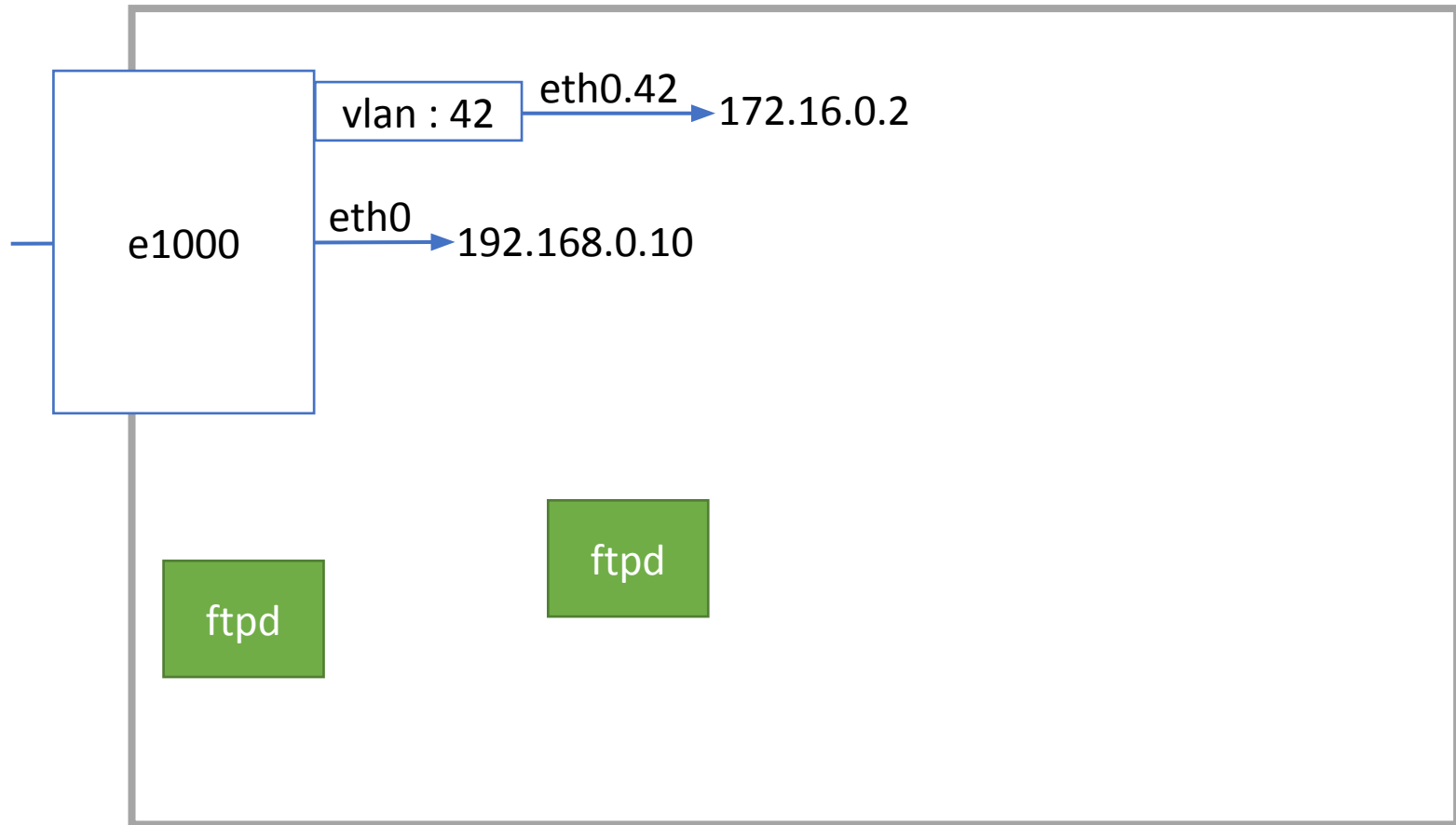


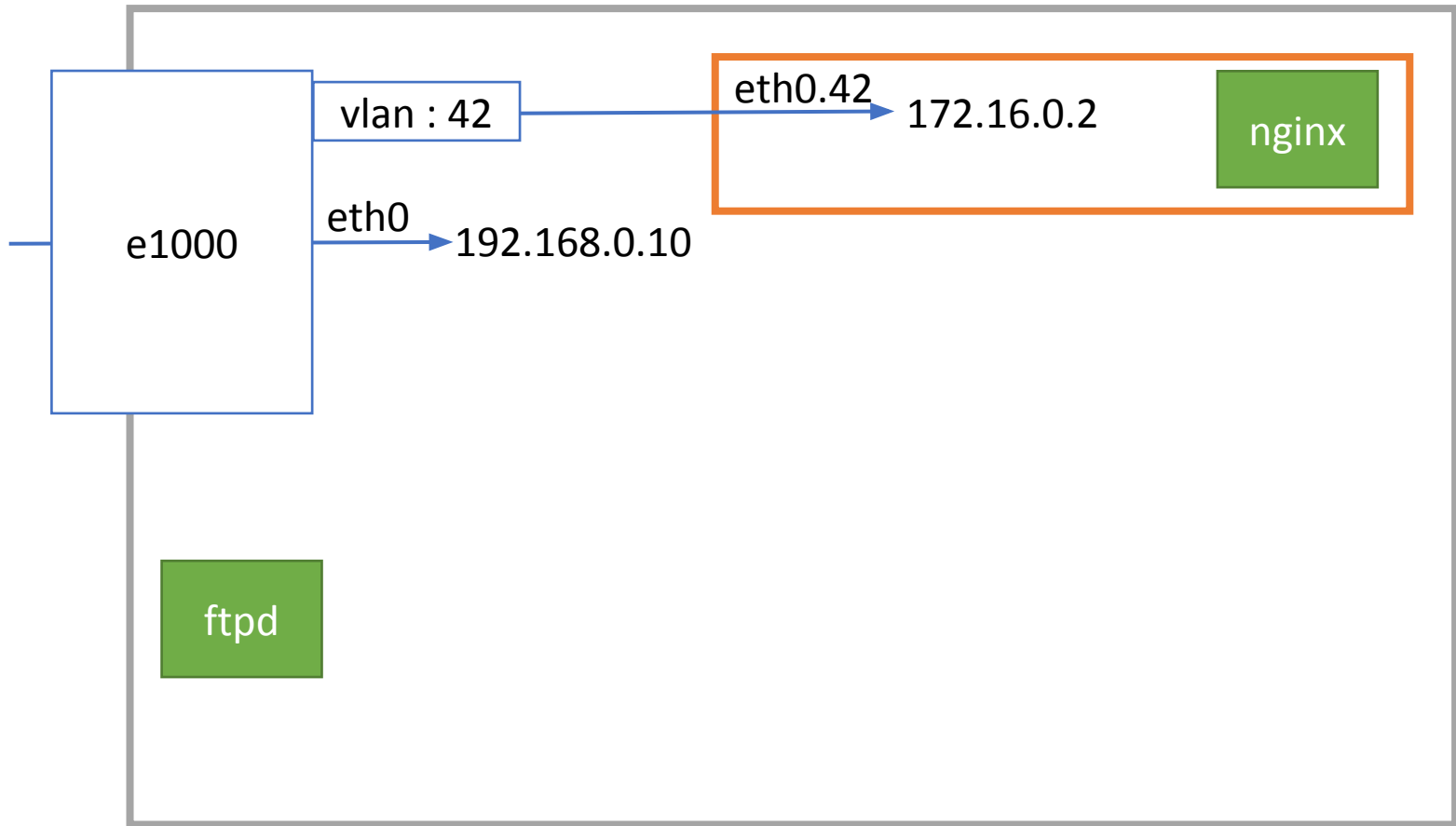


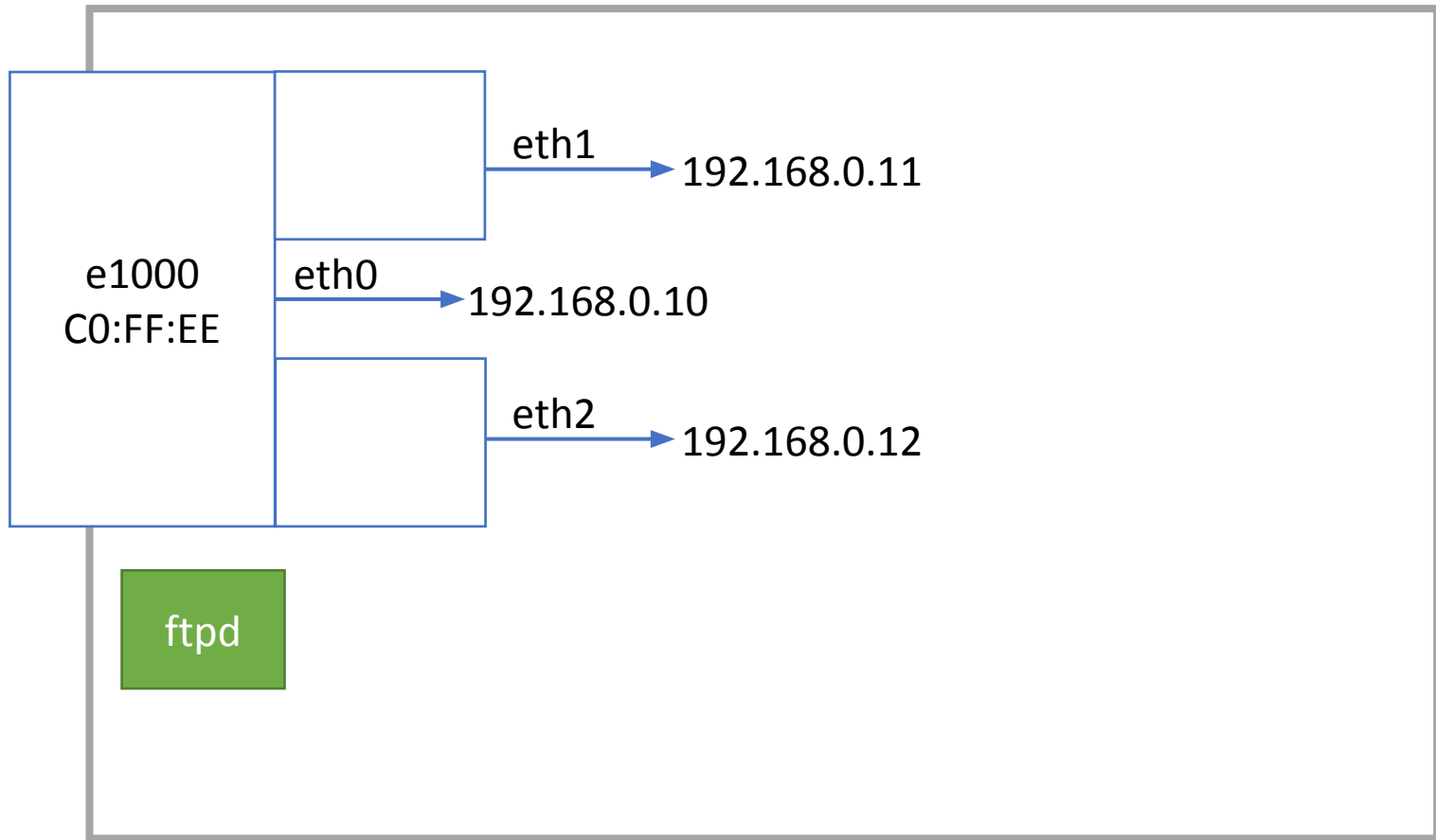


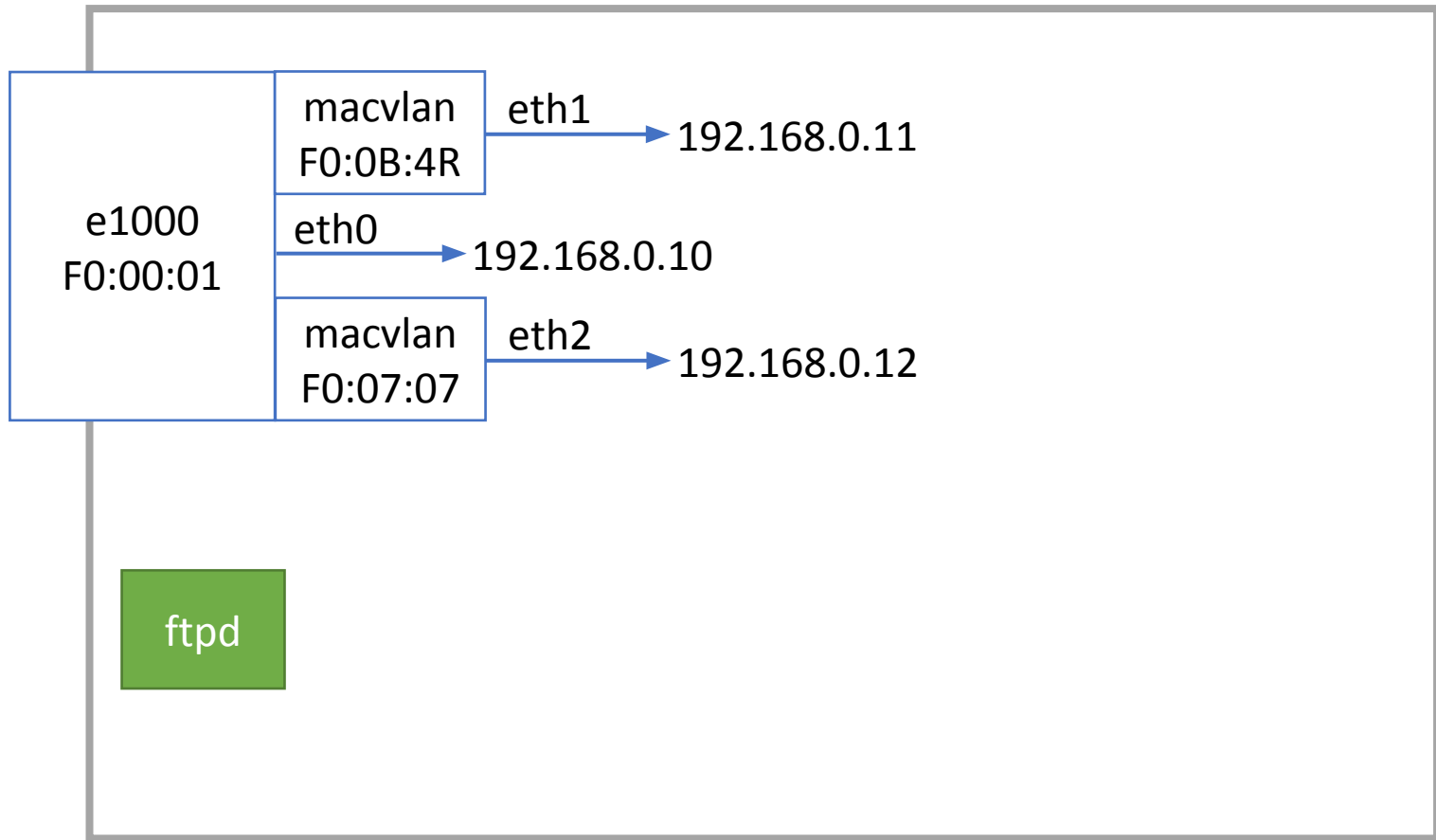




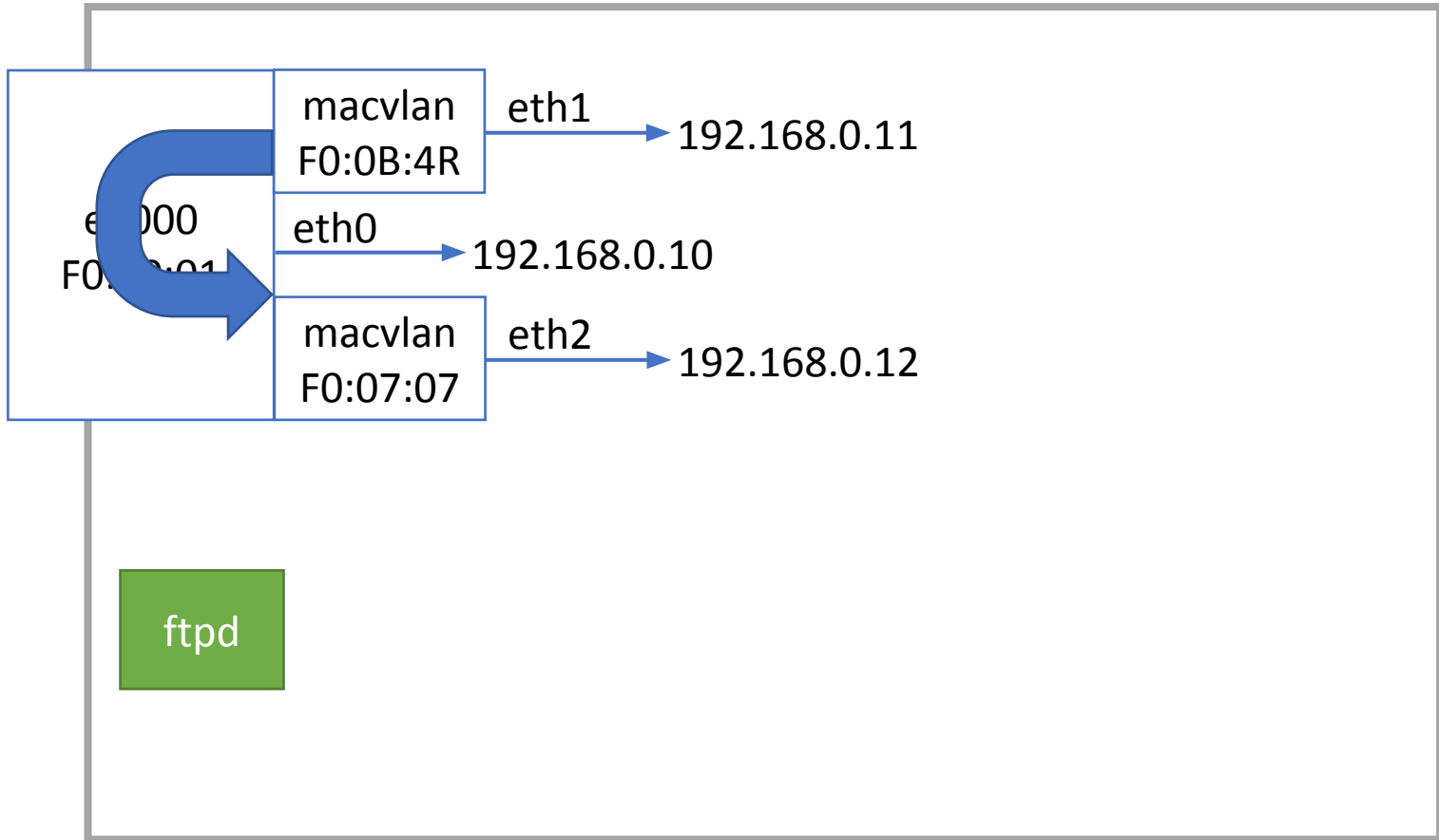


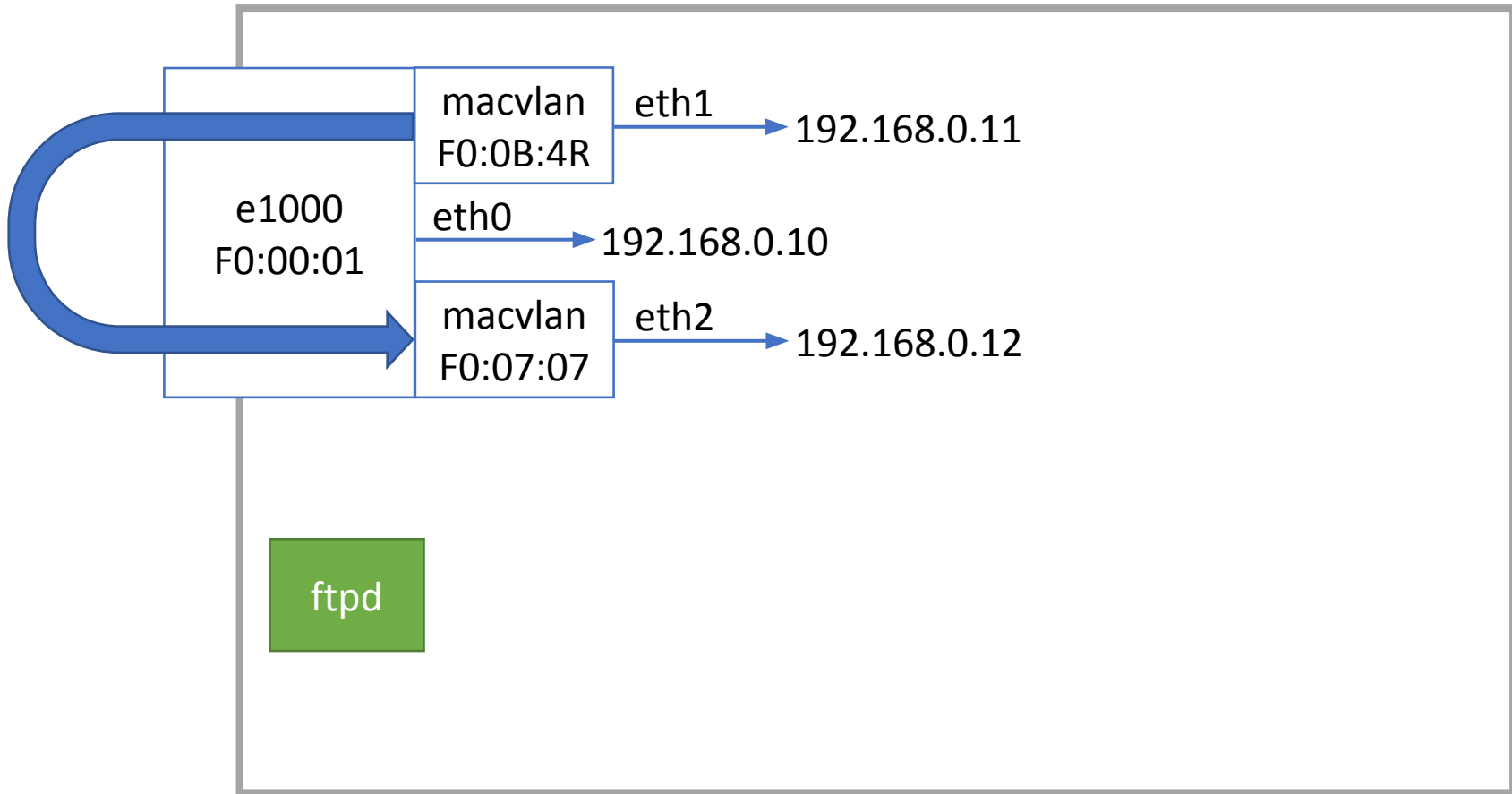


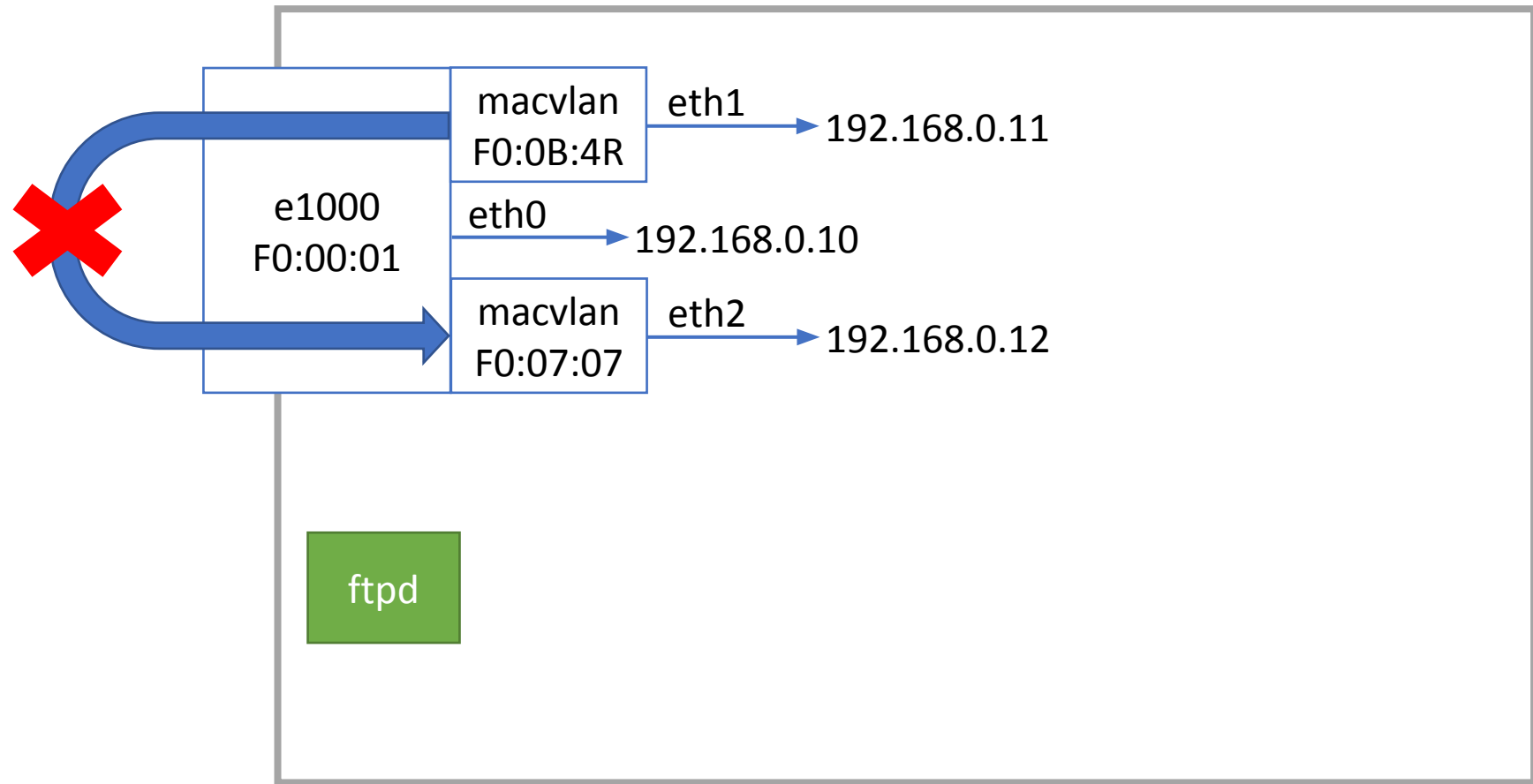


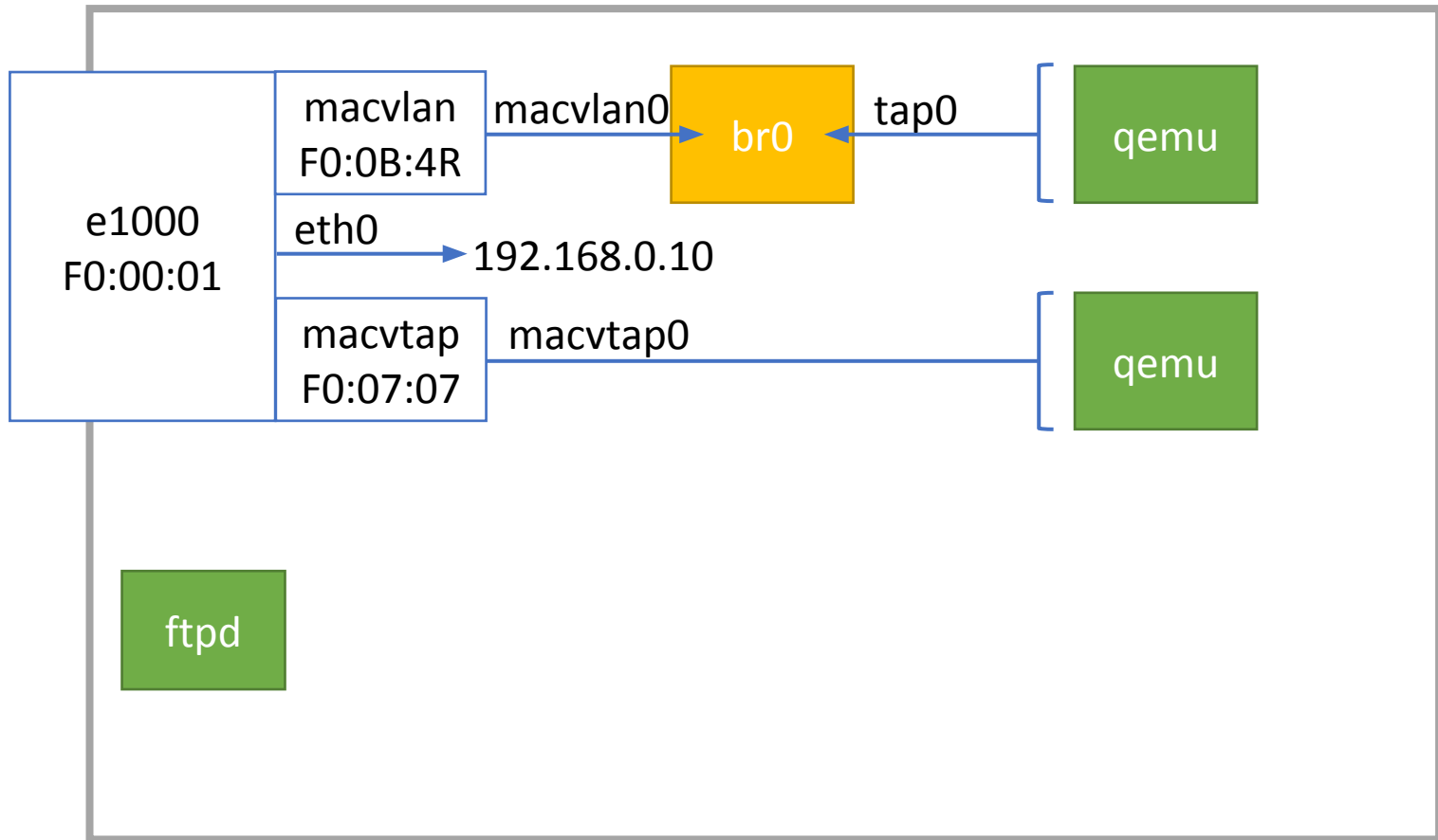


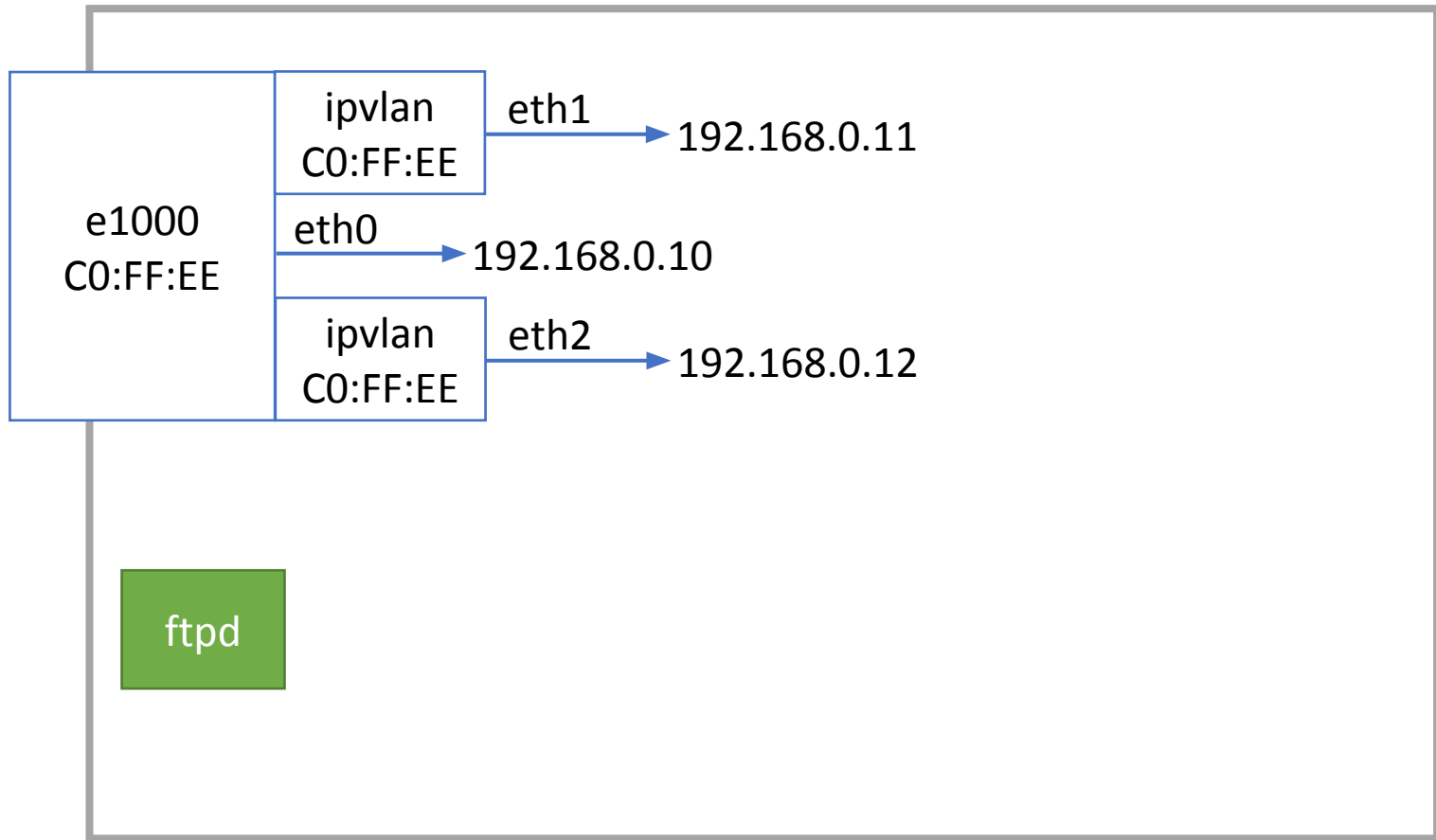


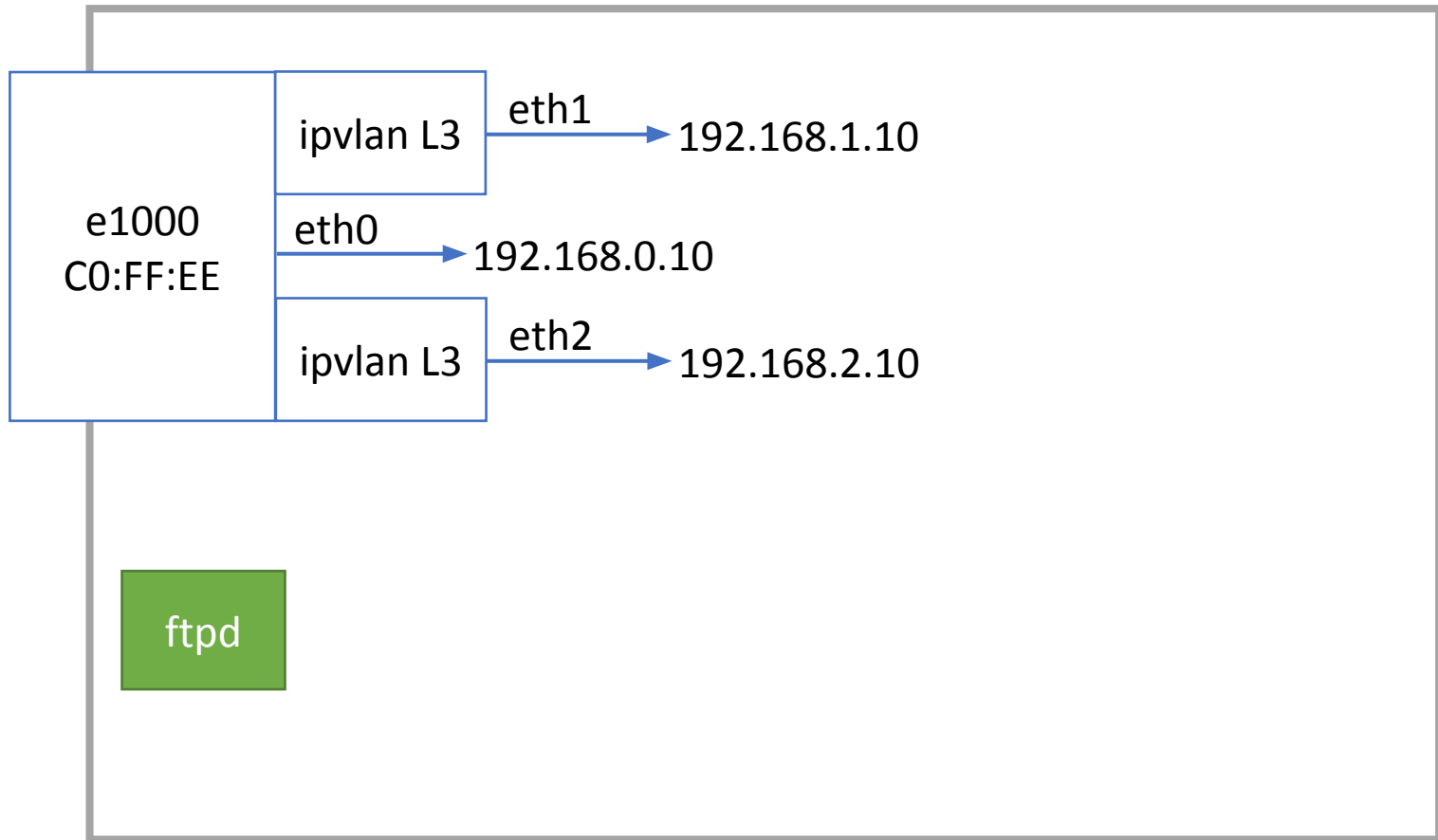


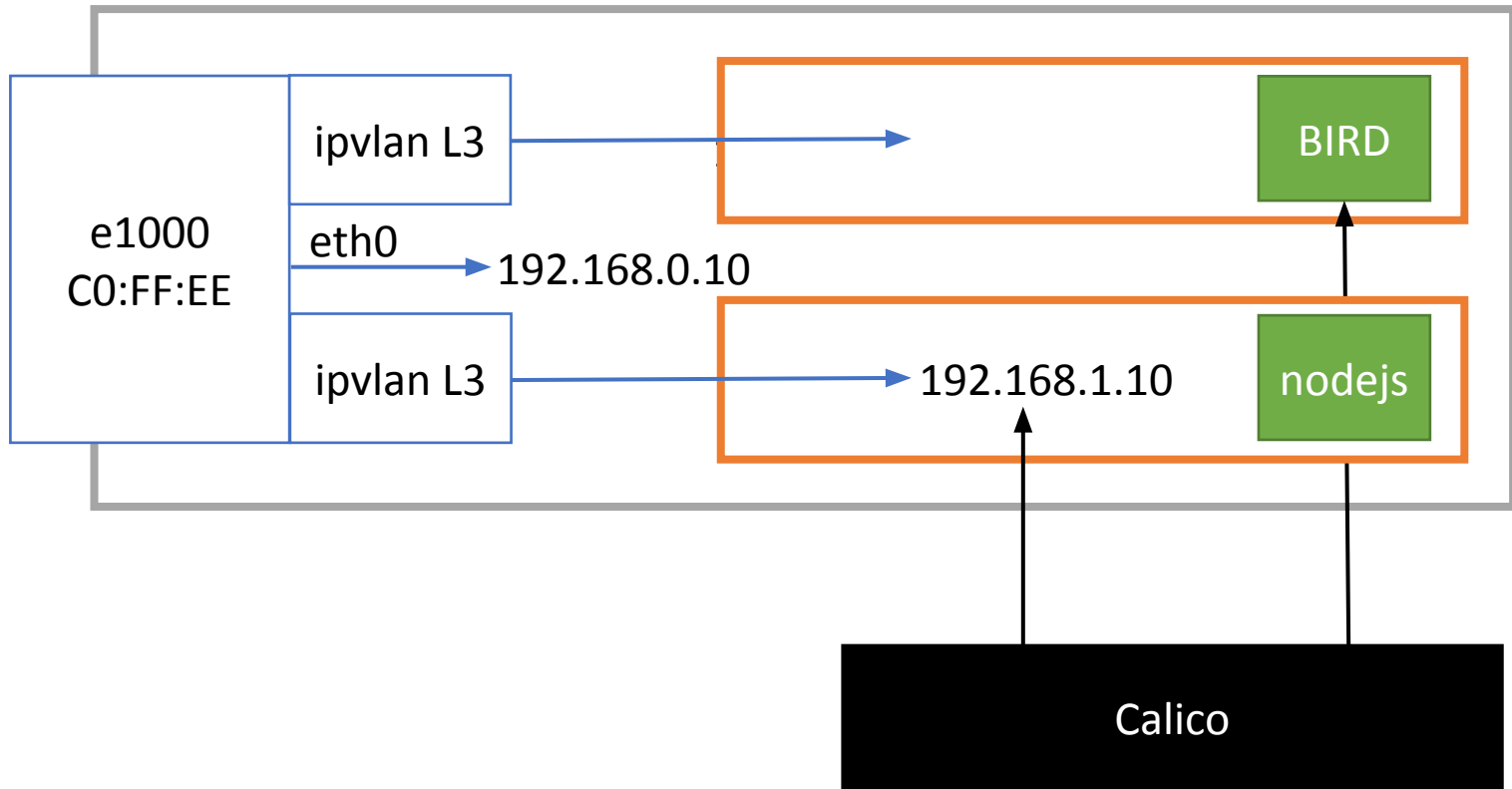


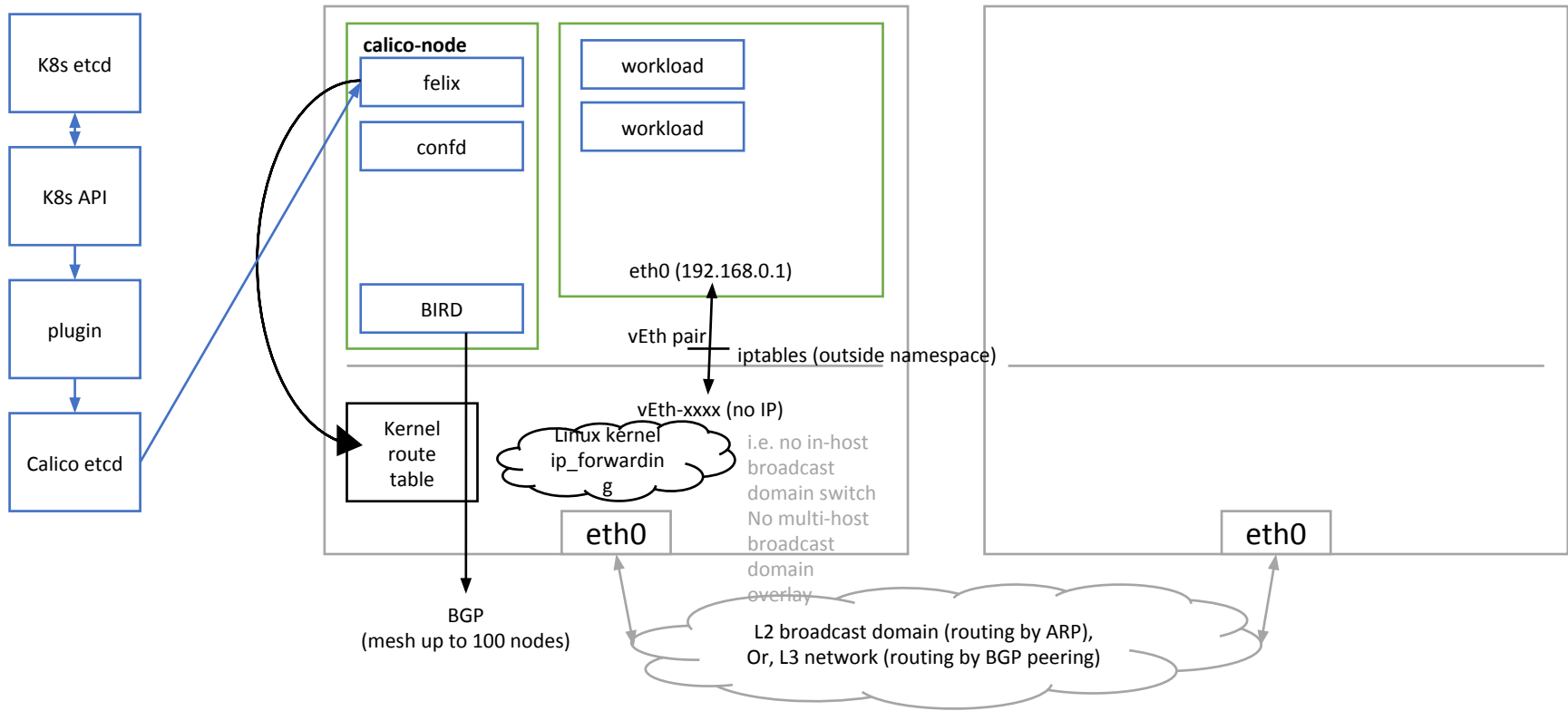












K8s etcd  
K8s API  
plugin  
Calico etcd

**calico-node**  
felix  
confd  
BIRD  
Kernel route table  
BGP (mesh up to 100 nodes)  
Linux kernel ip\_forwarding  
eth0

workload  
workload  
eth0 (192.168.0.1)

vEth pair  
iptables (outside namespace)

vEth-xxxx (no IP)  
i.e. no in-host broadcast domain switch  
No multi-host broadcast domain overlay

L2 broadcast domain (routing by ARP),  
Or, L3 network (routing by BGP peering)

eth0



# Recap - Interface Types

- Loopback
- Dummy
- “real”
- (multiple L3 addresses)
- TAP / TUN
- vEth
- Vlan
- Macvlan / macvtap
- Ipvlan / ipvtap L2
- Ipvlan / ipvtap L3

# Recap - Bridge Types

- Linux Bridge
- Open vSwitch (inc DPDK)
- Macvtap
- (netmap / VALE)
- (snabbswitch)
- (SR-IOV NIC, Cisco vNIC)

# Thanks!

@mt165

Slides |  
Videos | [mt165.co.uk](http://mt165.co.uk)  
Demo code |

