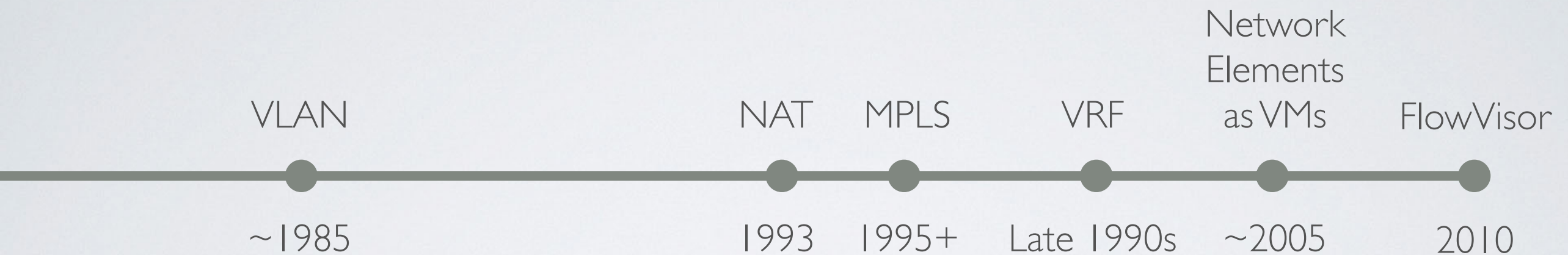# NETWORK VIRTUALIZATION IN MULTI-TENANT DATACENTERS

## Teemu Koponen

with

Keith Amidon, Peter Balland, Martín Casado, Anupam Chanda, Bryan Fulton, Igor Ganichev, Jesse Gross, Natasha Gude, Paul Ingram, Ethan Jackson, Andrew Lambeth, Romain Lenglet, Shih-Hao Li, Amar Padmanabhan, Justin Pettit, Ben Pfaff, Rajiv Ramanathan, Scott Shenker, Alan Shieh, Jeremy Stribling, Pankaj Thakkar, Dan Wendlandt, Alexander Yip, and Ronghua Zhang
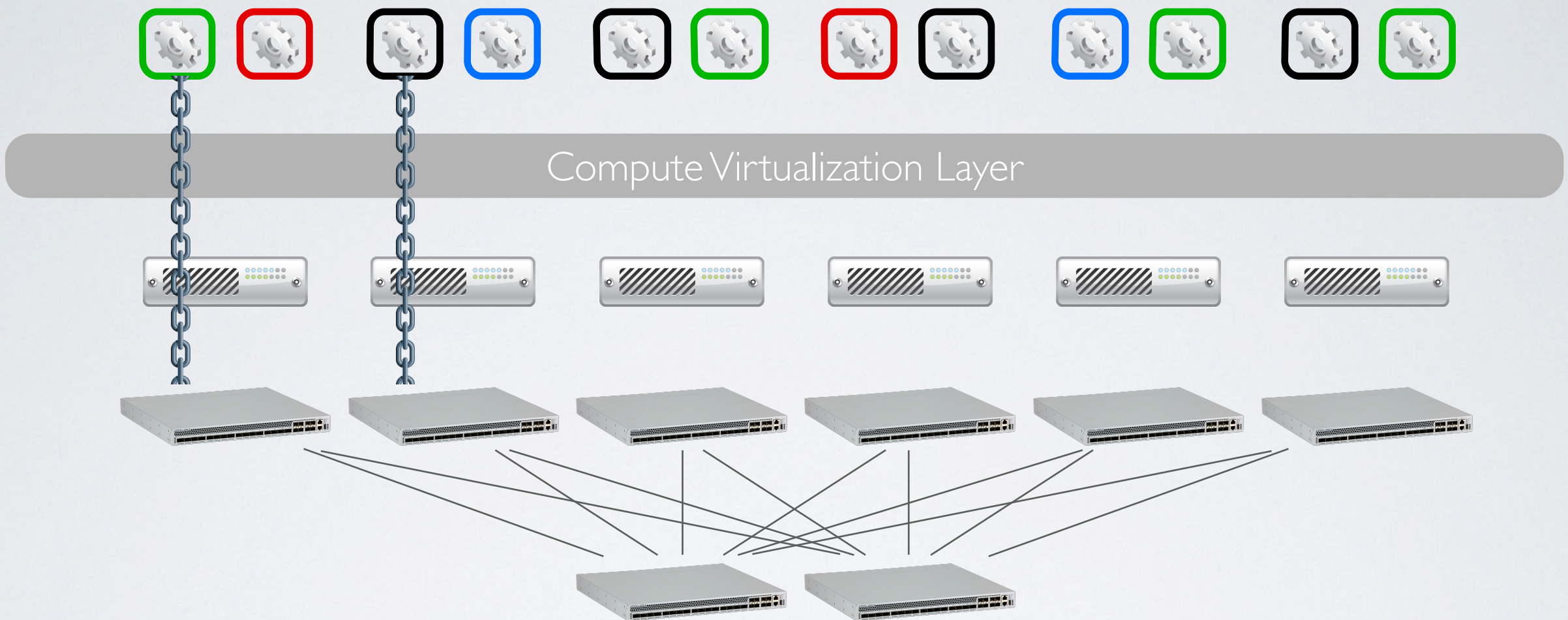
# NETWORK VIRTUALIZATION?

| | | | | Network Elements as VMs | |
|---|---|---|---|---|---|
| VLAN | | NAT | MPLS | VRF | FlowVisor |
| ~1985 | | 1993 | 1995+ | Late 1990s | ~2005 | 2010 |

| VLAN | NAT | MPLS | VRF | Elements as VMs | FlowVisor |
|---|---|---|---|---|---|
| Subnet | IP address space | Path | L3 FIB | Elements | ASIC |

Plenty of primitives but **no** network virtualization per se.
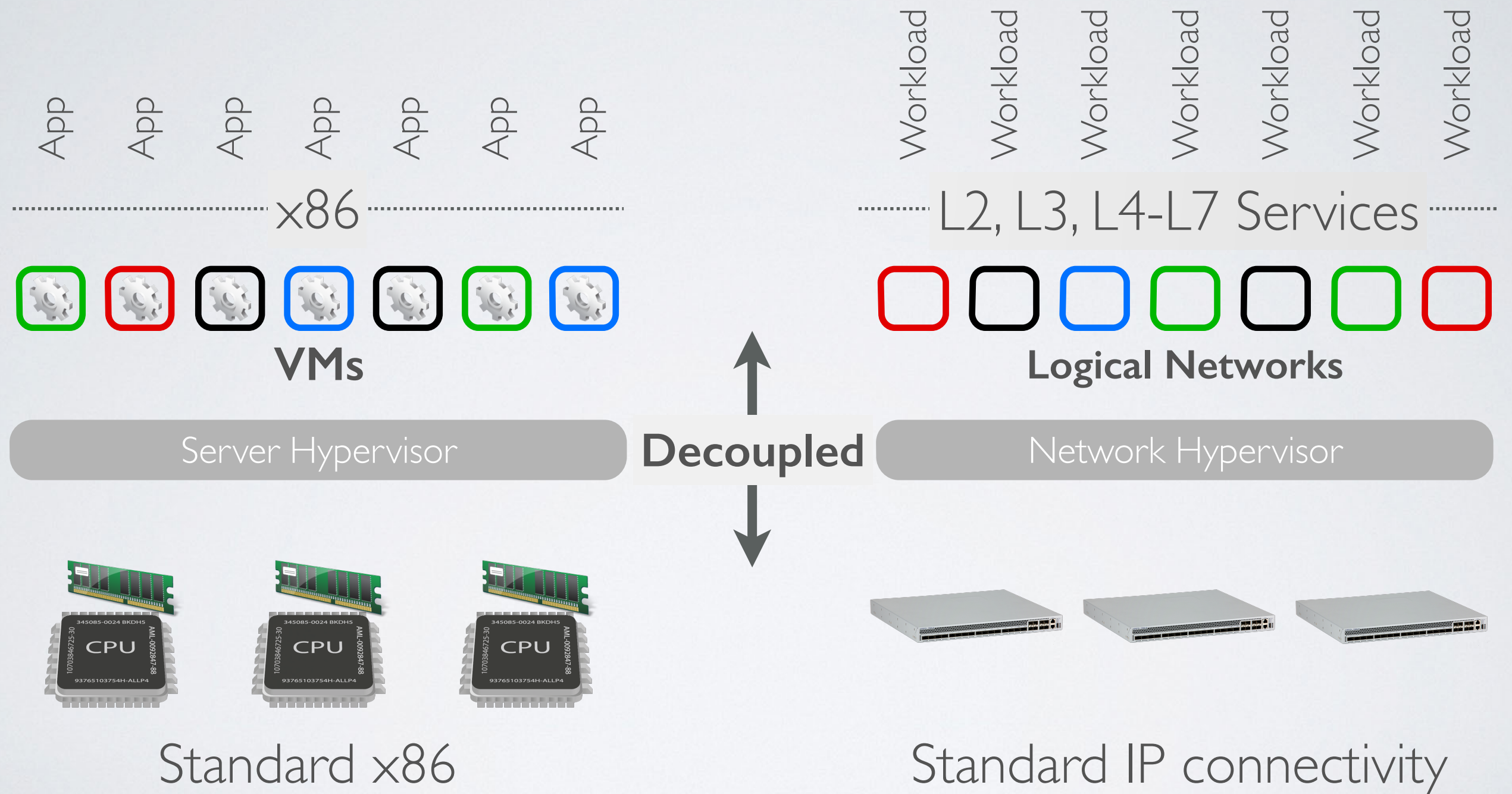
# MULTI-TENANT DATACENTERS

Compute Virtualization Layer

Result with the aforementioned primitives:

- Slow provisioning
- Mobility is limited
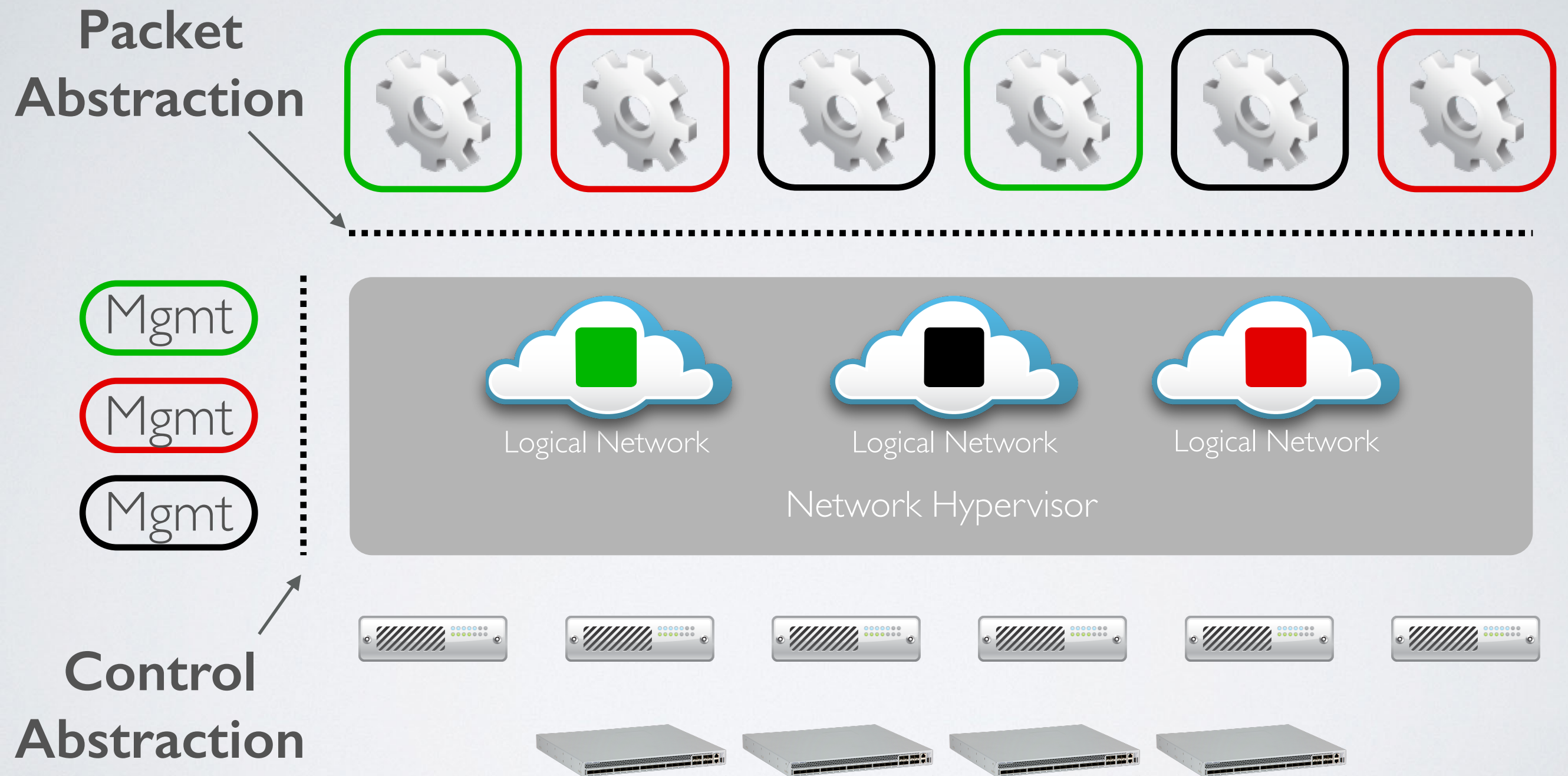- Operationally intensive

- Limited VM placement
- Hardware dependent
- …

# AGENDA

- Overall design of NVP network hypervisor.

- Design challenges.

- Hard lessons learnt.

- What's next in network virtualization?

# WHAT IS A NETWORK HYPERVISOR?

**Packet Abstraction**

**Mgmt**

**Mgmt**

**Mgmt**

Logical Network

Logical Network

Logical Network

Network Hypervisor

**Control Abstraction**

**Packet Abstraction + Control Abstraction = Network Hypervisor**
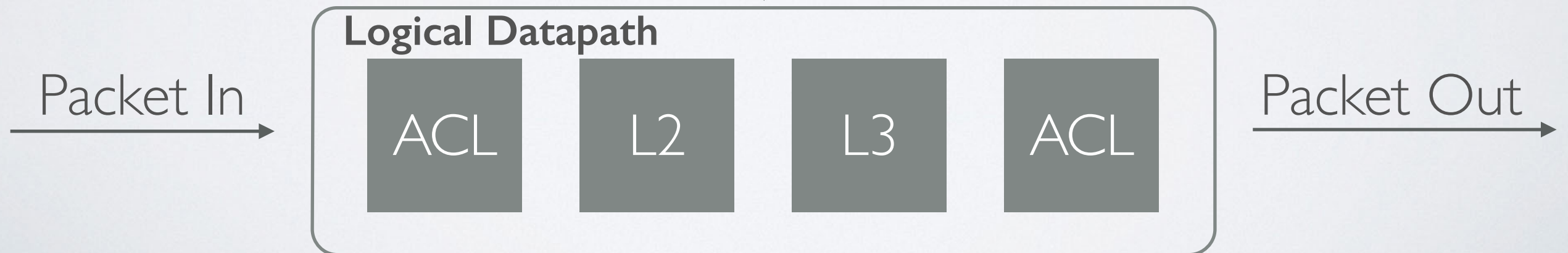
# WHAT ARE THE ABSTRACTIONS?

## Packet abstraction

- Compliance with standard TCP/IP stack is a necessity:
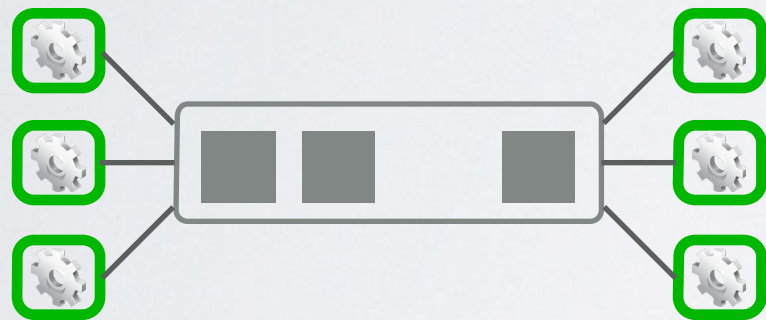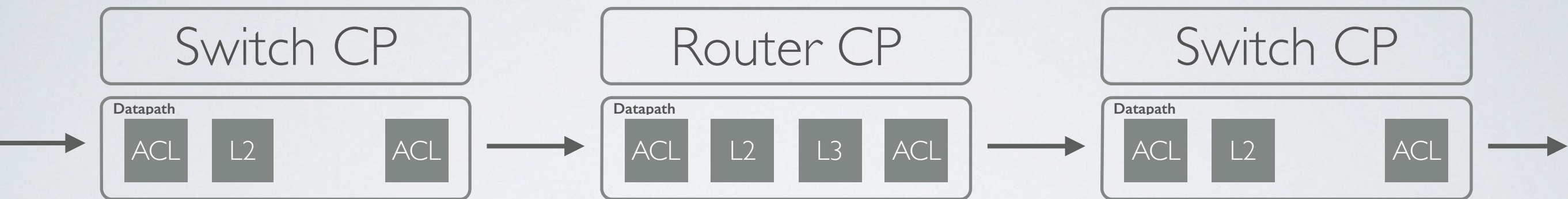
    - L2, L3 semantics (unicast, ARP, …)

## Control abstraction

- Networking has no single high level control interface.
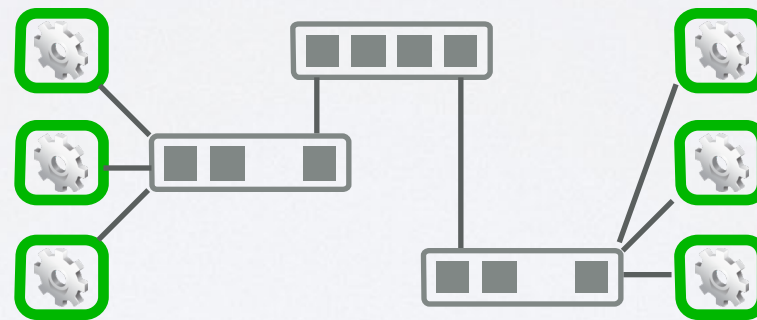
    - There's a low-level one though!

**Tenant's Control Plane**

↓

Packet In →  **Logical Datapath** | ACL | L2 | L3 | ACL | → Packet Out

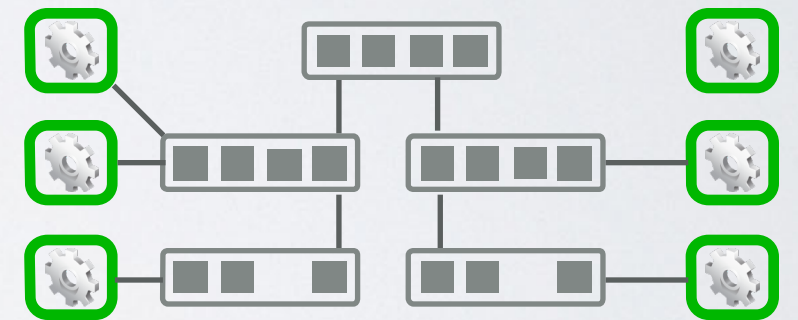# GENERALITY OF DATAPATH



One logical switch

2-tier logical topology

Arbitrary logical topology

Faithful reproduction of physical network service model.
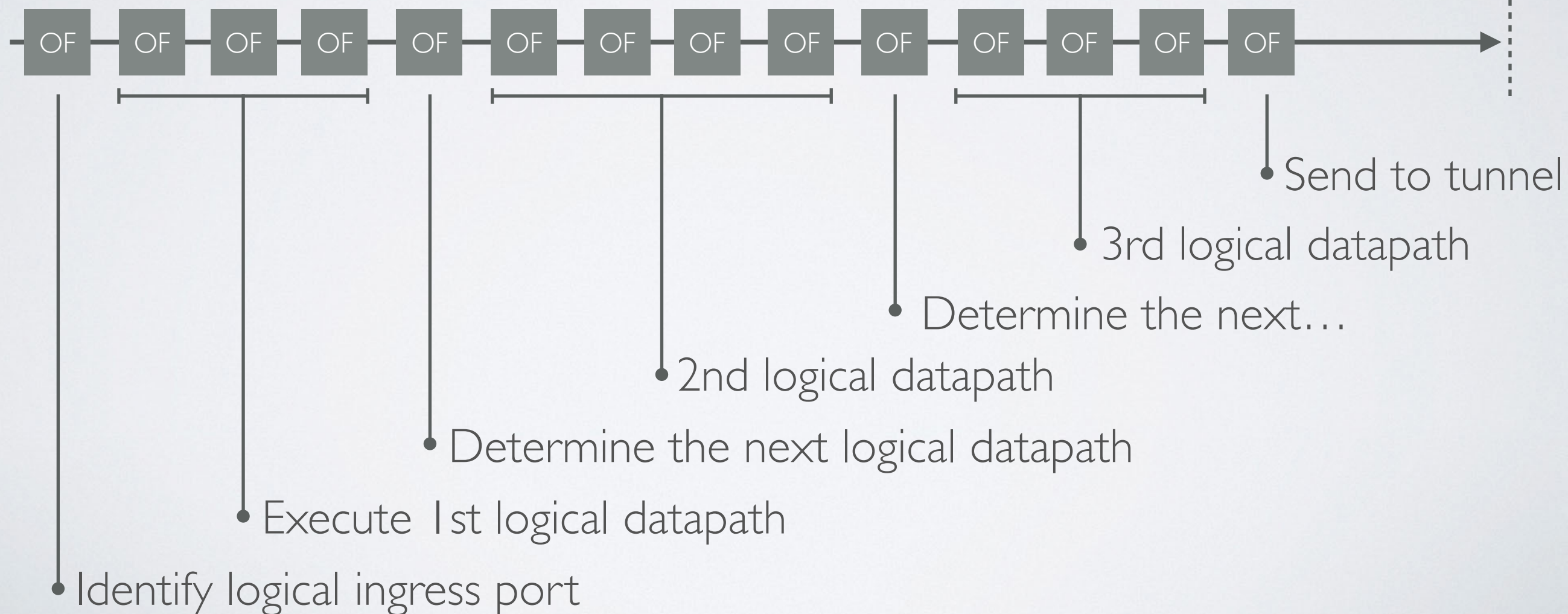
# WHERE TO IMPLEMENT?

Tenant CP

Network Hypervisor

- Independence from physical hardware.
- Programmatic control.
- Operational model of compute virtualization.

No extra x86 hops: just the source and destination hypervisor!

# INSIDE THE VIRTUAL SWITCH

**Datapath** ACL L2 ACL

**Datapath** ACL L2 L3 ACL

**Datapath** ACL L2 ACL

**Logical Topology**

**First-hop vSwitch**

OF OF OF OF OF OF OF OF OF OF OF OF OF OF

Send to tunnel

3rd logical datapath

Determine the next…

2nd logical datapath

Determine the next logical datapath

Execute 1st logical datapath

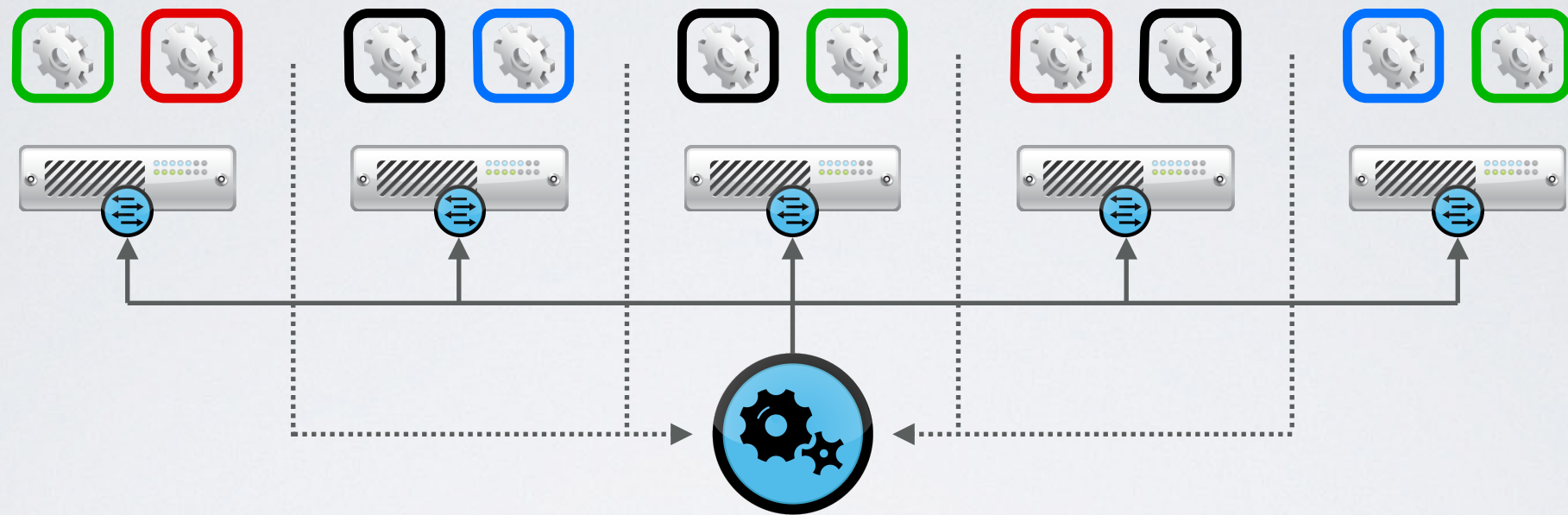Identify logical ingress port

# COMPUTATIONAL CHALLENGE



1. Controllers learn the location of VMs.

2. Controllers proactively compute & push all forwarding state required to connect each VM.

→ Forwarding State = **F**(configuration, VM locations)

Repeat above as logical configuration or physical configuration (VM placement) changes.

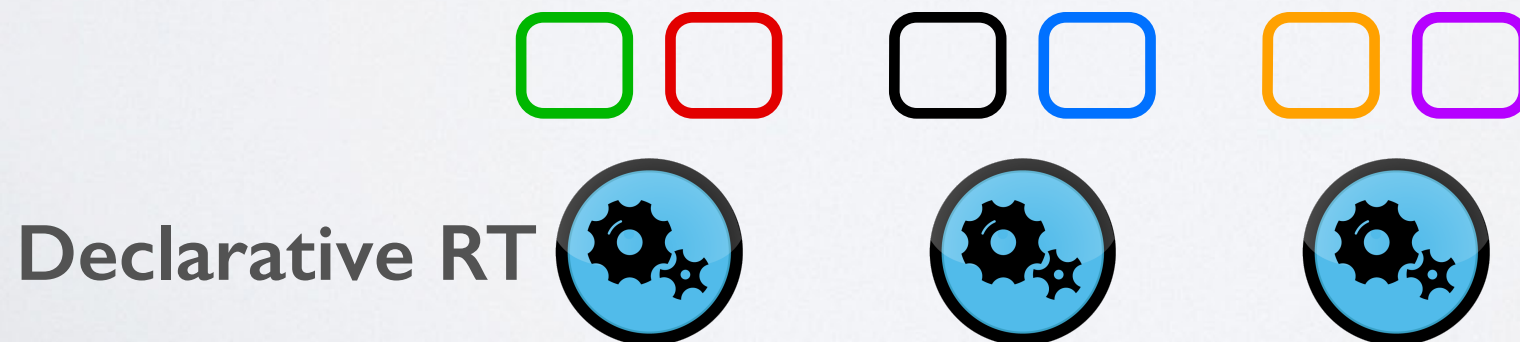**Challenge**: How to compute $O(N^2)$ volume of low-level OpenFlow and OVSDB state, when inputs change all the time.

# STATE COMPUTATION
## *Forwarding State = F(configuration, VM locations)*

### 1. How to Scale Computation

- **Incremental** computation and pushing for quick updates.

### 2. How to Guarantee Correctness

- Avoid all handwritten finite state machines, **machine generated** instead.

→ Datalog based declarative language to program F.

→ Shard the computation across controller cluster.

**Declarative RT**

# LESSONS LEARNT: ABSTRACTIONS
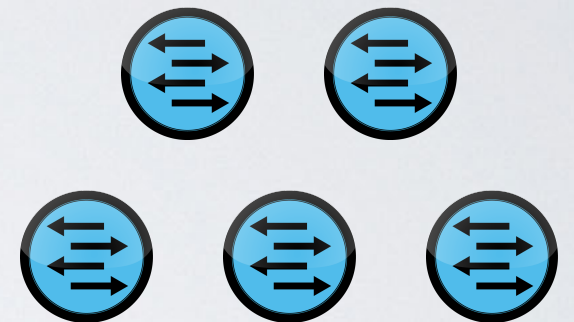
**"Basic Enterprise App"**  **"Modern App"**  **"Bank"**

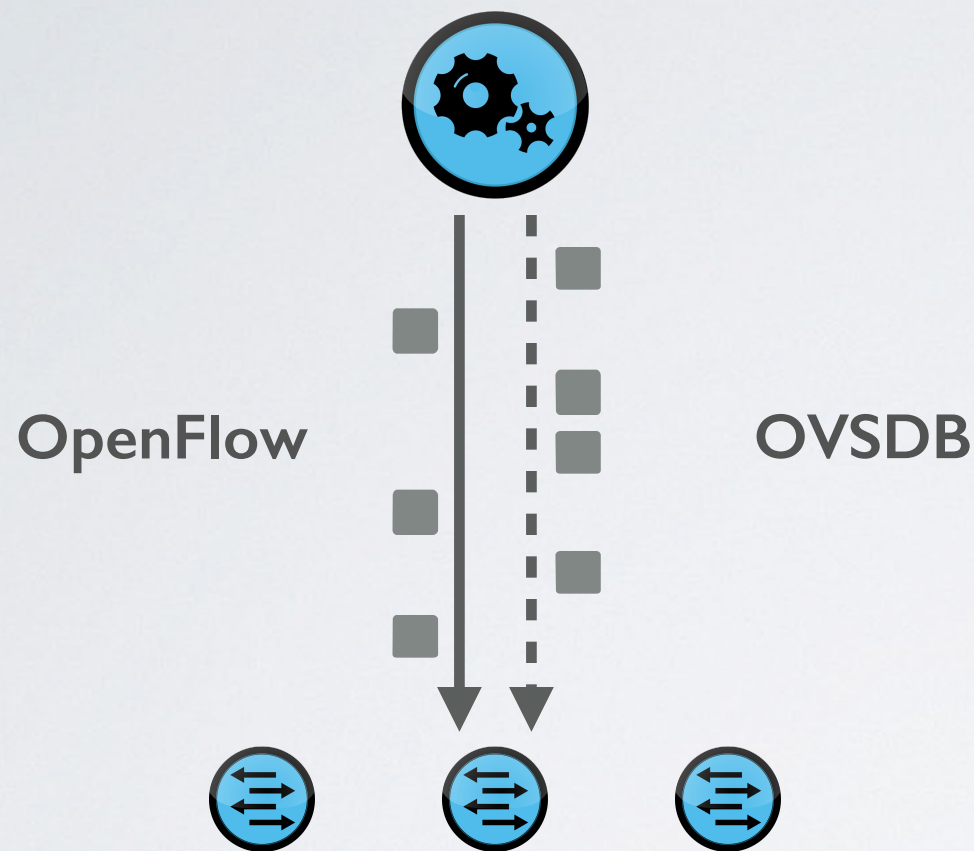**A logical switch**  **Two tier logical network**  **Arbitrary logical network**

- Assumptions about logical network structure often embedded into the workload.

- A single L2 domain sufficient for initial, simple workloads.

- To support more complex workloads without changing them, more complex logical topologies become a necessity.
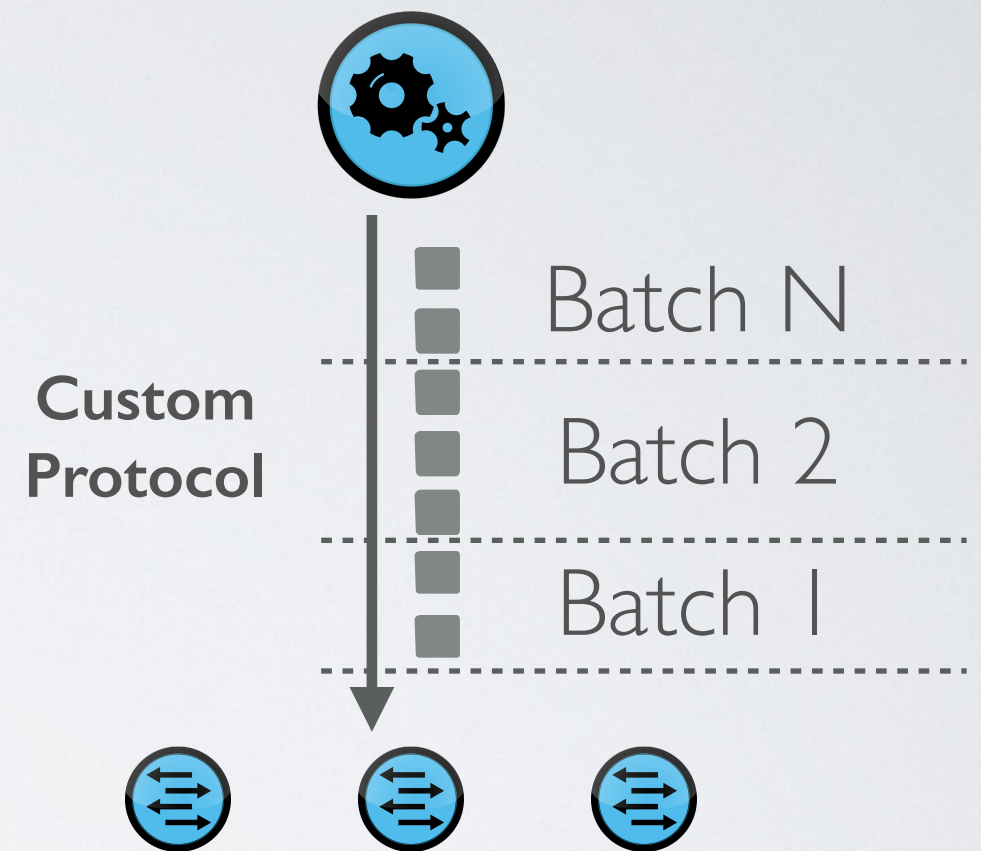
# LESSONS: FAILURE ISOLATION

**OpenFlow** **OVSDB**

**Custom Protocol**

Batch N

Batch 2

Batch 1

**Two Channels, No Atomic Updates**

**One Channel, Atomic Updates**

- Proactive pushing of all state not enough to decouple controllers from data plane.

- Connection may die while pushing updates.

- Atomically applied, batched updates.

- Connection failure does not result in incomplete state.

Data plane may operate over incomplete state!

At most old state.

# LESSONS: SCALING
## OPENFLOW IS EXPENSIVE

**Too primitive**

**Too tightly coupled**

- Simple operations take several flow entries.

  - For example, tunnel failover, encapsulation header ops.

  - Lots of redundancy.

- Each switch requires some flow customization; can't just blindly replicate flows.

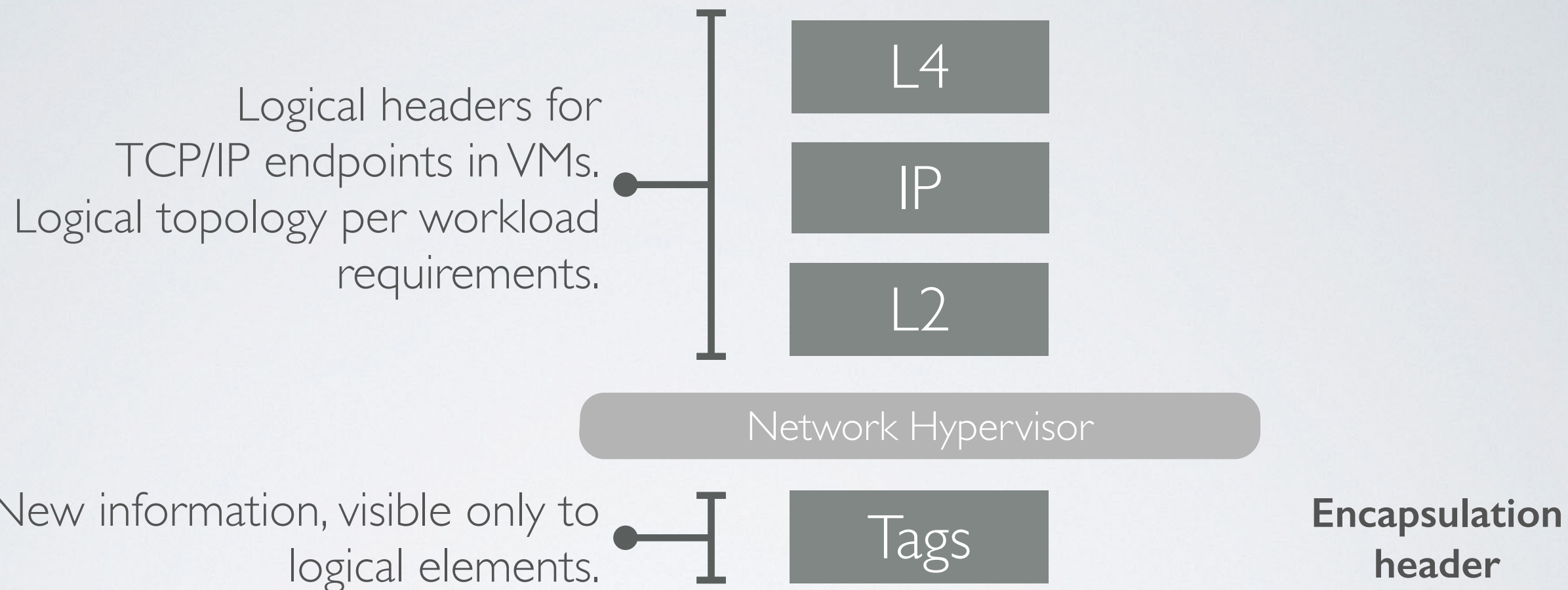- To compute flow entries, may have to wait for responses from the OVS configuration database.

↓

Replace a OF & OVSDB with a network virtualization specific protocol.

↓

OpenFlow becomes a protocol internal to the hypervisor.

# CONCLUSION: WHAT'S NEXT

Logical headers for TCP/IP endpoints in VMs. Logical topology per workload requirements.

| L4 |
|---|
| IP |
| L2 |

Network Hypervisor

New information, visible only to logical elements.

| Tags |
|---|

**Encapsulation header**

## Without Network Virtualization

- Workload may run on a topology where addresses provide little information.

- For instance, firewall rules defined over exact /32 addresses!

## With Network Virtualization

- New "out-of-band" header fields without breaking legacy TCP/IP stacks.

- **Huge** implications to enforcing security policies: groups, users in packet…

# THANK YOU! QUESTIONS?