# Software-Defined Data Centers

Brighten Godfrey
CS 538 March 6, 2017

# Multi-Tenant Data Centers: The Challenges

# Key Needs

Agility

Strength

Constitution

Dexterity

Charisma

# Key Needs

Agility

Location independent addressing

Performance uniformity

Security
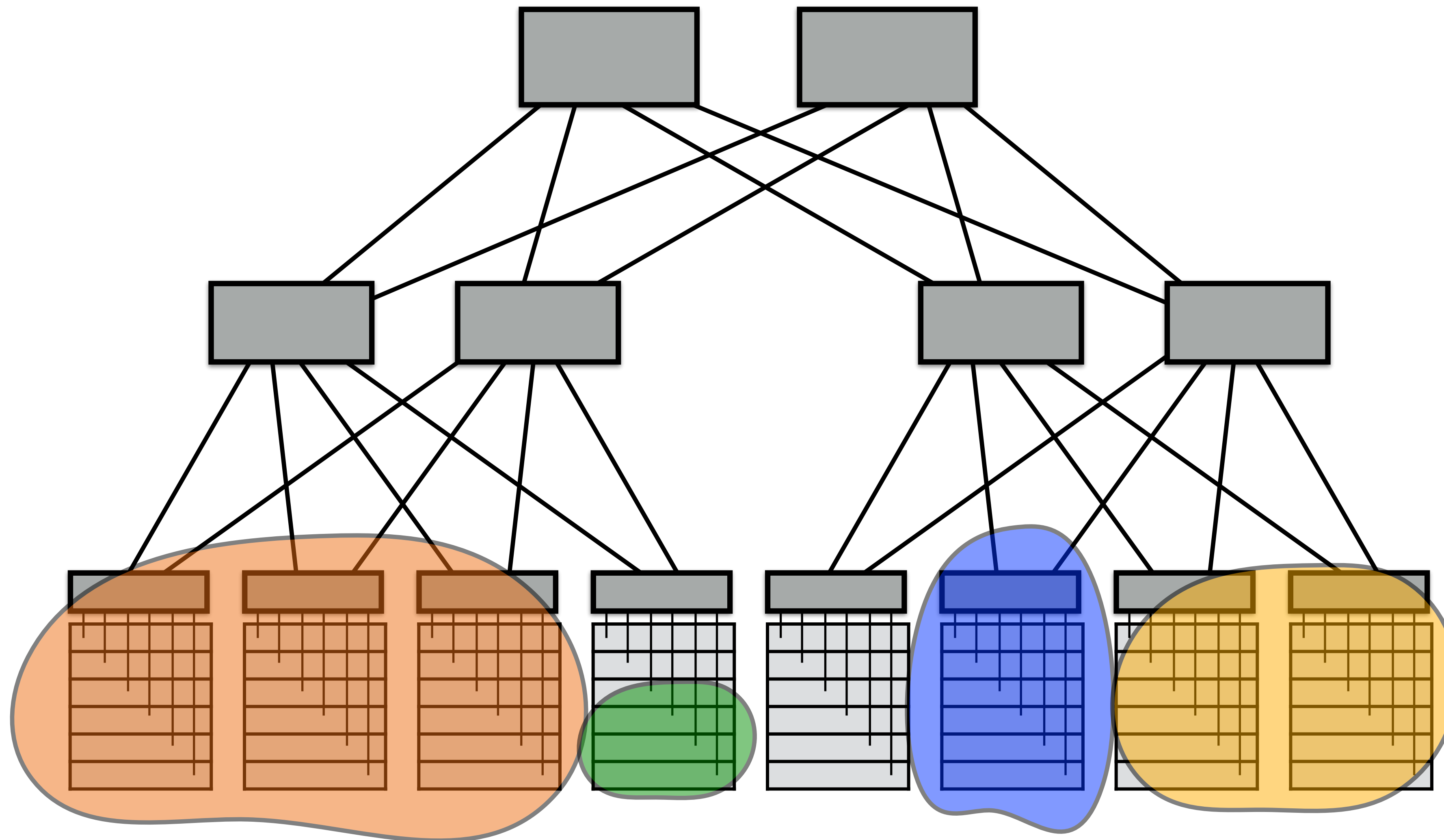
Network semantics

# Agility

Agility: Use any server for any service at any time

- Better economy of scale through increased utilization
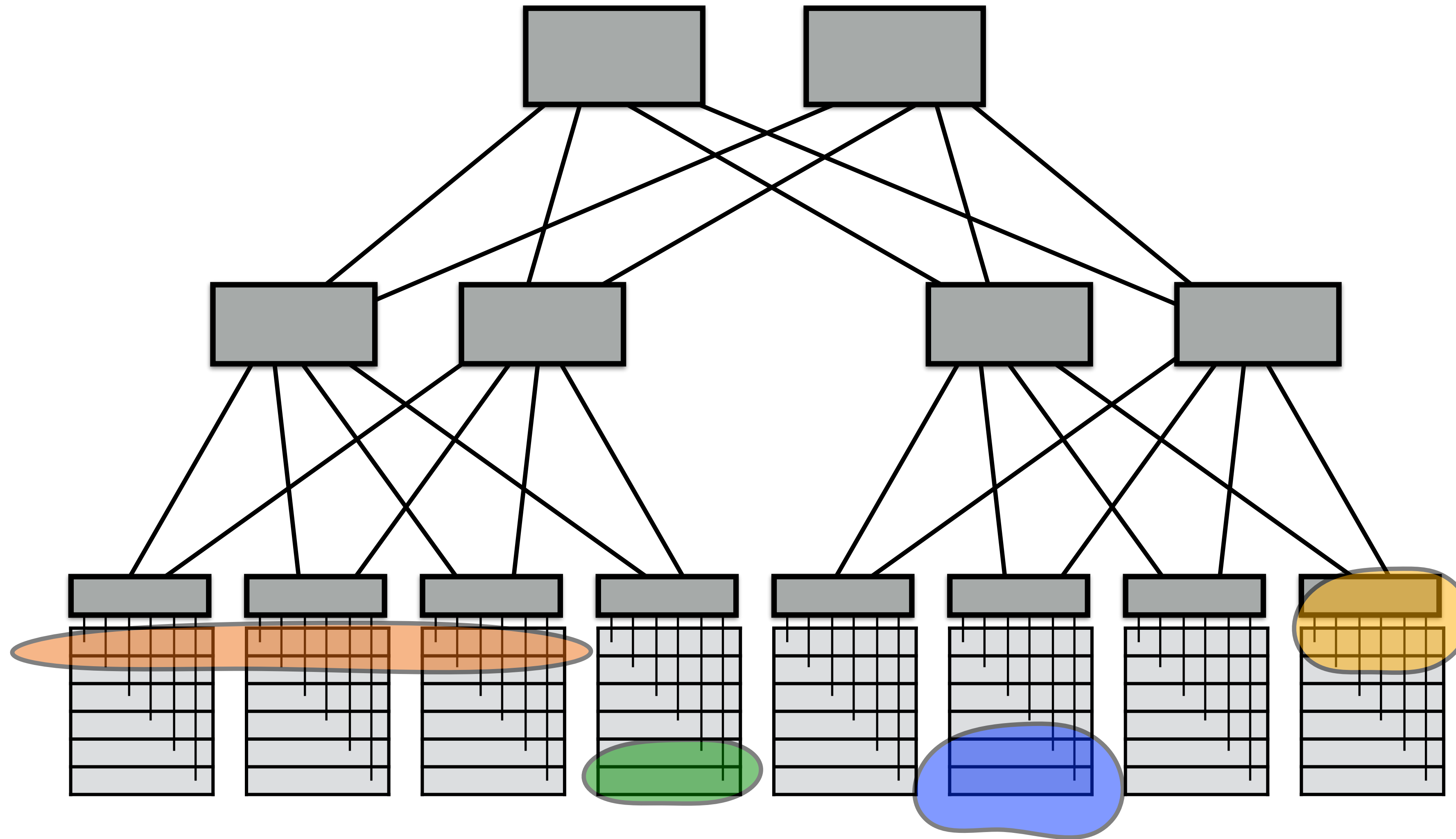- Improved reliability

Service / tenant

- Customer renting space in a public cloud
- Application or service in a private cloud (internal customer)

# Lack of Agility in Traditional DCs



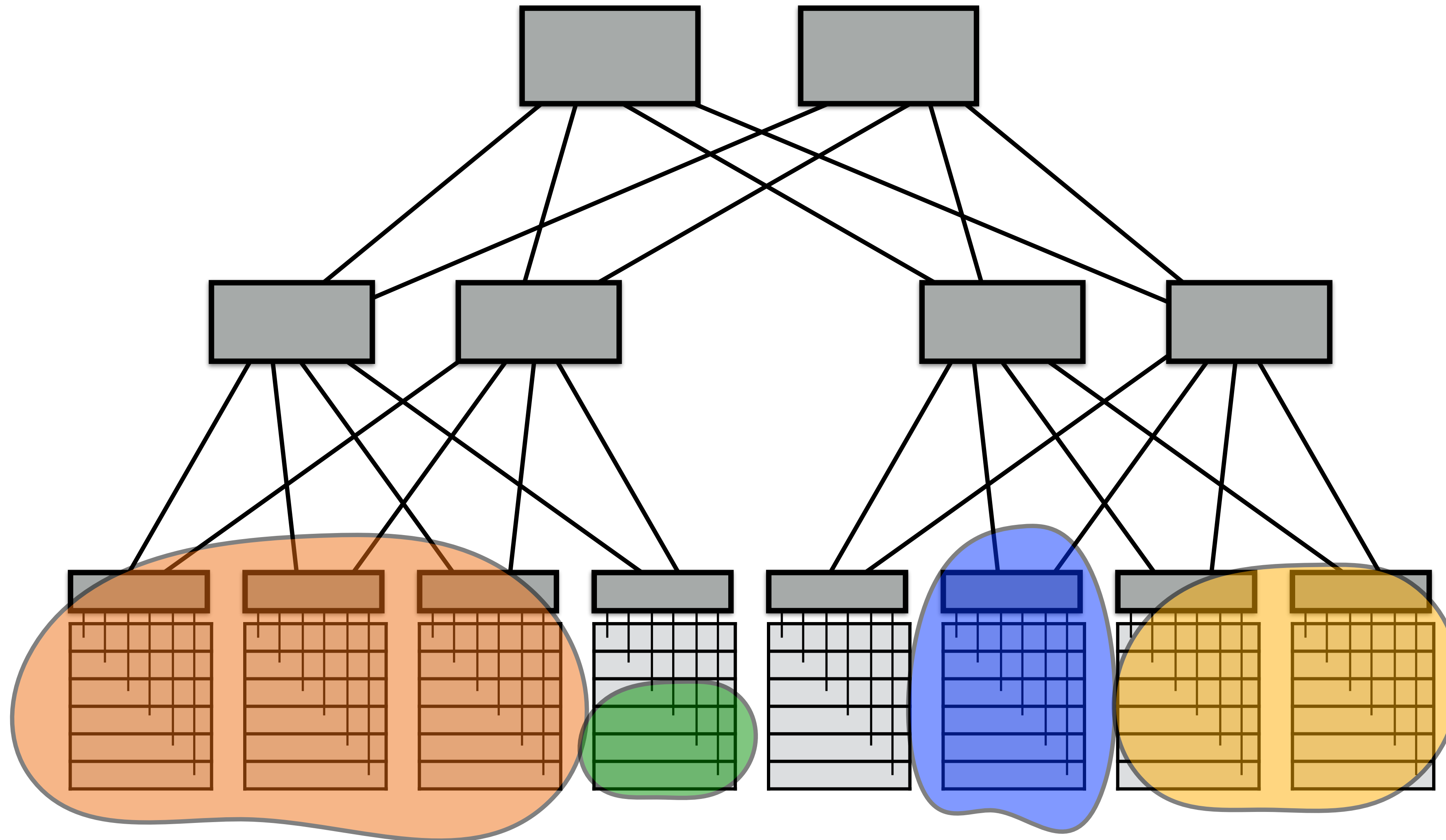Tenants in "silos"

# Lack of Agility in Traditional DCs



Tenants in "silos"

Poor utilization

# Lack of Agility in Traditional DCs



Tenants in "silos"

Poor utilization

Inability to expand

# Lack of Agility in Traditional DCs



10.0.4.0/24

10.0.6.0/24

IP addresses locked to topological location!

# Key Needs

## Agility

Location independent addressing

- Tenant's IP addresses can be taken anywhere

Performance uniformity

Security

Network semantics

# Lack of Agility in Traditional DCs



1:100 or worse oversubscription

Nonuniform performance

Full line rate

# Key Needs

Agility

Location independent addressing

- Tenant's IP addresses can be taken anywhere

Performance uniformity

- VMs receive same throughput regardless of placement

Security

Network semantics

# Lack of Agility in Traditional DCs



Untrusted environment

# Key Needs

Agility

Location independent addressing

- Tenant's IP addresses can be taken anywhere

Performance uniformity

- VMs receive same throughput regardless of placement

Security

- Micro-segmentation: isolation at tenant granularity

Network semantics

# Lack of Agility in Traditional DCs



x 1000s of legacy apps in a large enterprise

# Key Needs

## Agility

## Location independent addressing

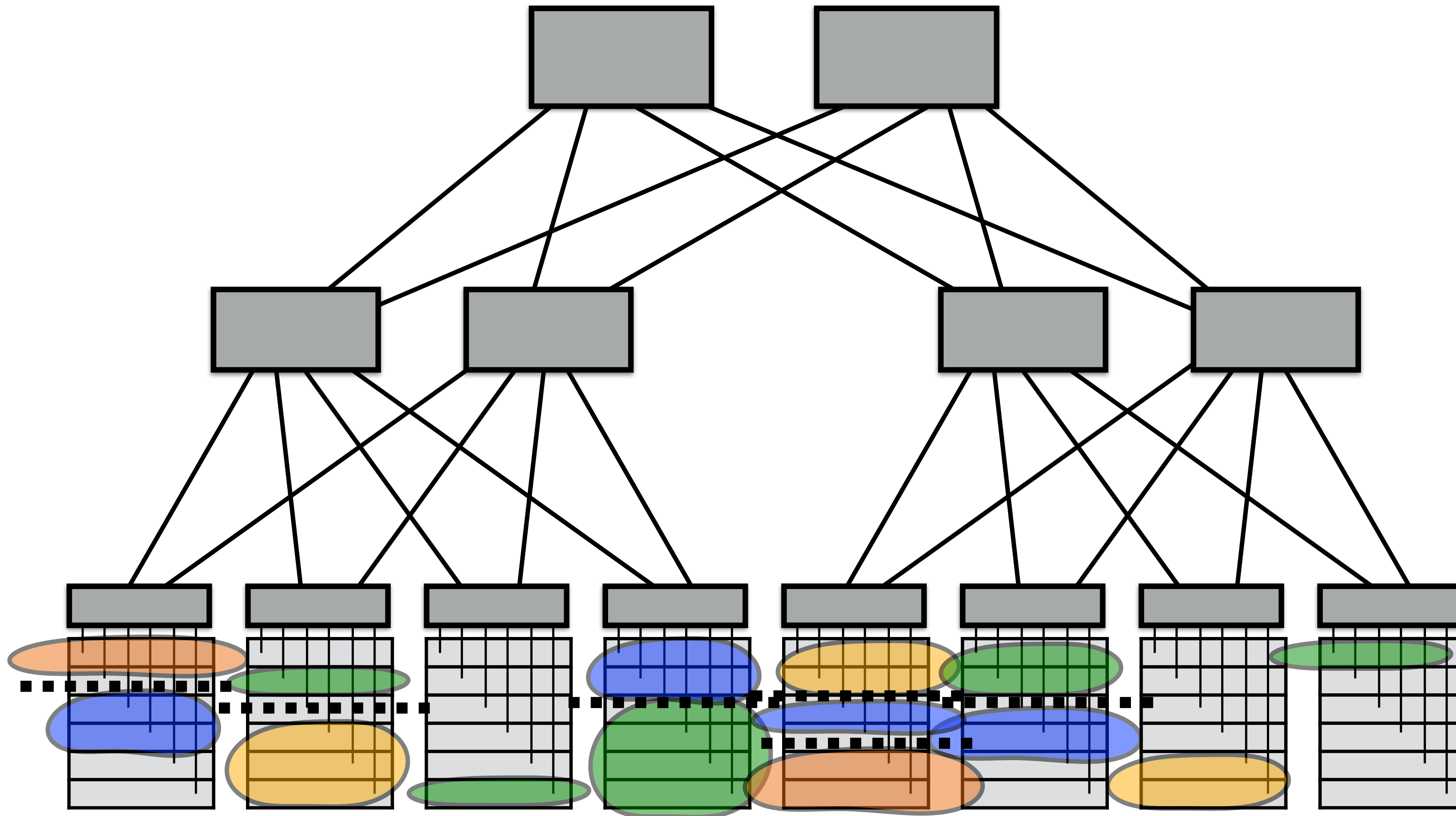- Tenant's IP addresses can be taken anywhere

## Performance uniformity

- VMs receive same throughput regardless of placement

## Security

- Micro-segmentation: isolation at tenant granularity

## Network semantics

- Layer 2 service discovery, multicast, broadcast, …

# Network Virtualization
# Case Study: VL2

# Case Study

## VL2: A Scalable and Flexible Data Center Network

Albert Greenberg      James R. Hamilton      Navendu Jain
Srikanth Kandula      Changhoon Kim      Parantap Lahiri
David A. Maltz      Parveen Patel      Sudipta Sengupta

Microsoft Research

### [ACM SIGCOMM 2009]

Influenced architecture of
Microsoft Azure



VL2 → Azure Clos Fabrics with 40G NICs

Scale-out, active-active

Outcome of >10 years of history, with major revisions every six months

[From Albert Greenberg keynote at SIGCOMM 2015:
http://conferences.sigcomm.org/sigcomm/2015/pdf/papers/keynote.pdf]

# Virtualization

"All problems in computer science can be solved by another level of indirection."

*– David Wheeler*

**App / Tenant layer**

- Application Addresses (AAs): Location independent
- Illusion of a single big Layer 2 switch connecting the app

**Virtualization layer**

- Directory server: Maintain AA to LA mapping
- Server agent: Query server, wrap AAs in outer LA header

**Physical network layer**

- Locator Addresses (LAs): Tied to topology, used to route
- Layer 3 routing via OSPF

# End-to-end example



Link-state network with LAs (10/8)

Int (10.1.1.1) · · · Int (10.1.1.1) · · · Int (10.1.1.1)

| H(ft) | 10.1.1.1 |
|---|---|
| H(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| H(ft) | 10.0.0.6 |
|---|---|
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

(10.0.0.4)
ToR
(20.0.0.1)

(10.0.0.6)
ToR
(20.0.0.1)

| H(ft) | 10.1.1.1 |
|---|---|
| H(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| 20.0.0.55 | 20.0.0.66 |
|---|---|
| Payload | |

S (20.0.0.55)

D (20.0.0.56)

IP subnet with AAs (20/8)

IP subnet with AAs (20/8)

[Greenberg et al.]

Application sends
to AA 20.0.0.56

# End-to-end example



Directory servers

Q: Where is AA 20.0.0.56?

A: LA 10.0.06

Link-state network with LAs (10/8)

**Int** (10.1.1.1) . . . **Int** (10.1.1.1) . . . **Int** (10.1.1.1)

| $H$(ft) | 10.1.1.1 |
|---|---|
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| $H$(ft) | 10.0.0.6 |
|---|---|
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

(10.0.0.4) **ToR** (20.0.0.1)

(10.0.0.6) **ToR** (20.0.0.1)

Host agent encapsulates

| $H$(ft) | 10.1.1.1 |
|---|---|
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| 20.0.0.55 | 20.0.0.66 |
|---|---|
| Payload | |

S (20.0.0.55)

D (20.0.0.56)

**IP subnet with AAs (20/8)**

**IP subnet with AAs (20/8)**

[Greenberg et al.]

Application sends to AA 20.0.0.56

# End-to-end example



Directory servers

Q: Where is AA 20.0.0.56?

A: LA 10.0.06

Link-state network with LAs (10/8)

| Int (10.1.1.1) | . . . | Int (10.1.1.1) | . . . | Int (10.1.1.1) |

| H(ft) | 10.1.1.1 |
| H(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| H(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

(10.0.0.4)
**ToR**
(20.0.0.1)

(10.0.0.6)
**ToR**
(20.0.0.1)

Host agent encapsulates

| H(ft) | 10.1.1.1 |
| H(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| 20.0.0.55 | 20.0.0.66 |
| Payload | |

S (20.0.0.55)

D (20.0.0.56)

**IP subnet with AAs (20/8)**

**IP subnet with AAs (20/8)**

[Greenberg et al.]

Application sends to AA 20.0.0.56

# End-to-end example

Directory servers

Q: Where is AA 20.0.0.56?

A: LA 10.0.06

Intermediate switch decapsulates

### Link-state network with LAs (10/8)

**Int** (10.1.1.1)  . . .  **Int** (10.1.1.1)  . . .  **Int** (10.1.1.1)

| $H$(ft) | 10.1.1.1 |
|---|---|
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| $H$(ft) | 10.0.0.6 |
|---|---|
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

(10.0.0.4) **ToR** (20.0.0.1)

(10.0.0.6) **ToR** (20.0.0.1)

Host agent encapsulates

| $H$(ft) | 10.1.1.1 |
|---|---|
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| 20.0.0.55 | 20.0.0.66 |
|---|---|
| Payload | |

S (20.0.0.55)

D (20.0.0.56)

**IP subnet with AAs (20/8)**

**IP subnet with AAs (20/8)**

[Greenberg et al.]

Application sends to AA 20.0.0.56

# End-to-end example

Directory servers

Q: Where is AA 20.0.0.56?

A: LA 10.0.0.6

Host agent encapsulates

Application sends to AA 20.0.0.56

Intermediate switch decapsulates

Destination ToR decapsulates again

**Link-state network with LAs (10/8)**

| **Int** (10.1.1.1) | . . . | **Int** (10.1.1.1) | . . . | **Int** (10.1.1.1) |

| $H$(ft) | 10.1.1.1 |
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

(10.0.0.4)
**ToR**
(20.0.0.1)

(10.0.0.6)
**ToR**
(20.0.0.1)

| $H$(ft) | 10.1.1.1 |
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| 20.0.0.55 | 20.0.0.66 |
| Payload | |

S (20.0.0.55)

D (20.0.0.56)

**IP subnet with AAs (20/8)**

**IP subnet with AAs (20/8)**

[Greenberg et al.]

# End-to-end example



Directory servers

Q: Where is AA 20.0.0.56?

A: LA 10.0.06

Link-state network with LAs (10/8)

Intermediate switch decapsulates

| **Int** (10.1.1.1) | . . . | **Int** (10.1.1.1) | . . . | **Int** (10.1.1.1) |

| $H$(ft) | 10.1.1.1 |
|---|---|
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| $H$(ft) | 10.0.0.6 |
|---|---|
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

Destination ToR decapsulates again

(10.0.0.4) **ToR** (20.0.0.1)

(10.0.0.6) **ToR** (20.0.0.1)

Host agent encapsulates

| $H$(ft) | 10.1.1.1 |
|---|---|
| $H$(ft) | 10.0.0.6 |
| 20.0.0.55 | 20.0.0.56 |
| Payload | |

| 20.0.0.55 | 20.0.0.66 |
|---|---|
| Payload | |

Host agent delivers

S (20.0.0.55)

D (20.0.0.56)

**IP subnet with AAs (20/8)**

**IP subnet with AAs (20/8)**

[Greenberg et al.]

Application sends to AA 20.0.0.56

# Did we achieve agility?

Location independent addressing

- AAs are location independent

L2 network semantics

- Agent intercepts and handles L2 broadcast, multicast

- Both of the above require "layer 2.5" shim agent running on host; but, concept transfers to hypervisor-based virtual switch

# Did we achieve agility?

Performance uniformity

- Clos network is nonblocking (non-oversubscribed)
- Uniform capacity everywhere
- ECMP provides good (though not perfect) load balancing
- But, performance isolation among tenants depends on TCP backing off to rate destination can receive
- Leaves open the possibility of fast load balancing

Security

- Directory system can allow/deny connections by choosing whether to resolve an AA to a LA
- But, segmentation not explicitly enforced at hosts

# Where's the SDN?

Directory servers: Logically centralized control

- Orchestrate application locations
- Control communication policy

Host agents: dynamic "programming" of data path

# VL2 Enduring Take-Aways

Scale-out nonblocking Clos network

ECMP for traffic-oblivious routing

Separation of virtual and physical addresses

Centralized control plane

# Network Virtualization
# Case Study: NVP

# Case Study: NVP

Teemu Koponen, 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, and Rajiv Ramanathan, *VMware;* Scott Shenker, *International Computer Science Institute and the University of California, Berkeley;* Alan Shieh, Jeremy Stribling, Pankaj Thakkar, Dan Wendlandt, Alexander Yip, and Ronghua Zhang, *VMware*

# NVP Approach to Virtualization

## 1. Service: Arbitrary network topology

# NVP Approach to Virtualization

## 1. Service: Arbitrary network topology

# NVP Approach to Virtualization

Service: Arbitrary network topology
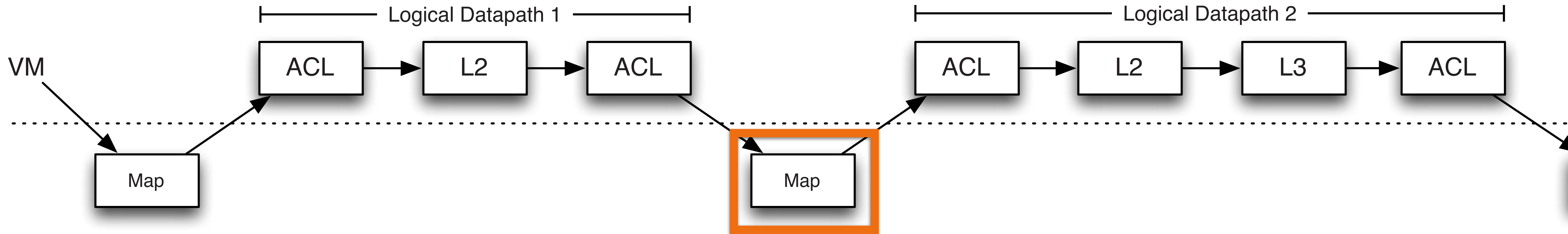


Network Hypervisor

Physical Network:
Any standard layer 3 network
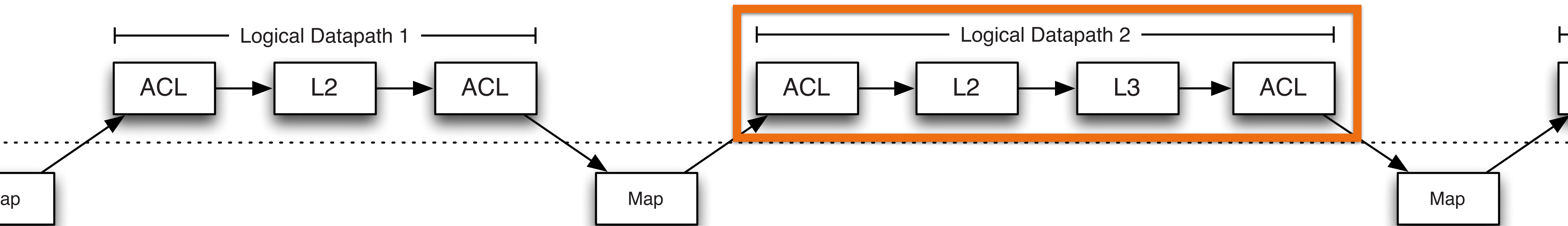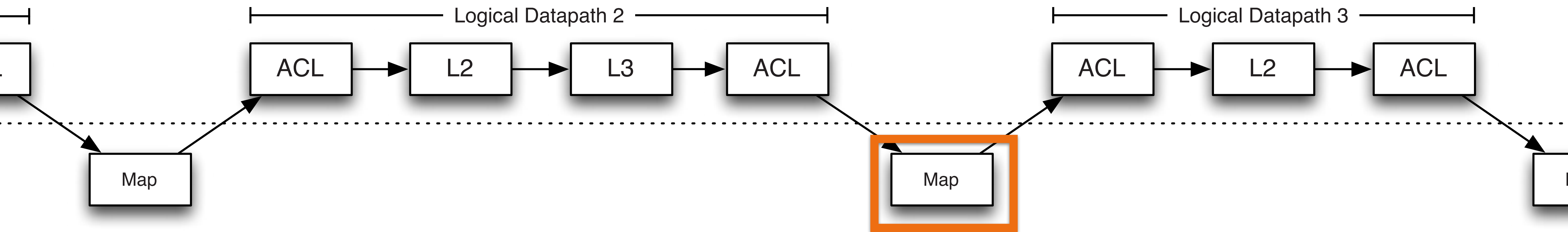
# Virtual network service

# Virtual network service



[Figure: Koponen et al.]

# Virtual network service

# Virtual network service

VM

Map

Logical Datapath 1
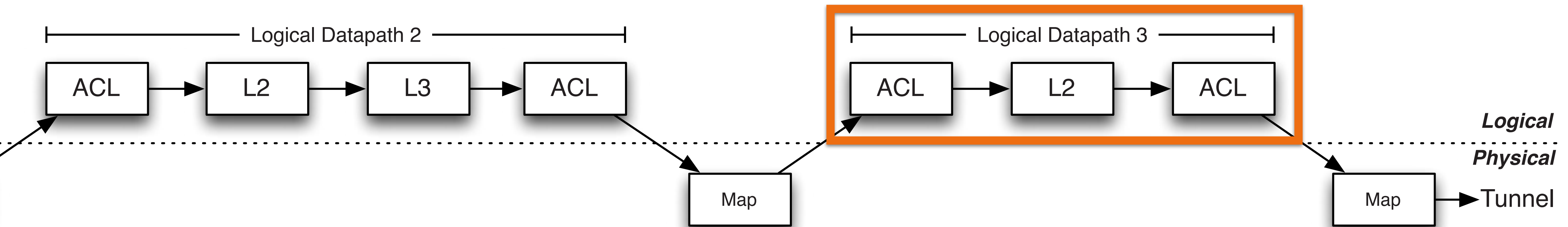
ACL → L2 → ACL

Map

Logical Datapath 2

ACL → L2 → L3

# Virtual network service

# Virtual network service
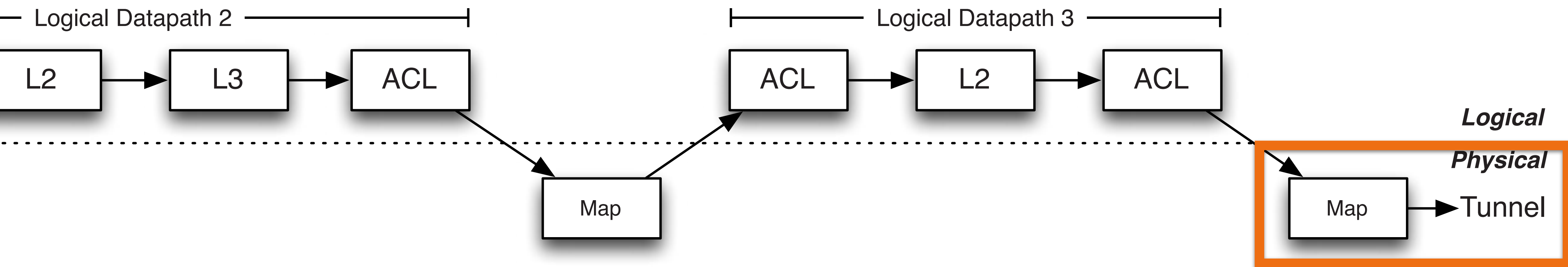
# Virtual network service

# Virtual network service

# Virtual network service
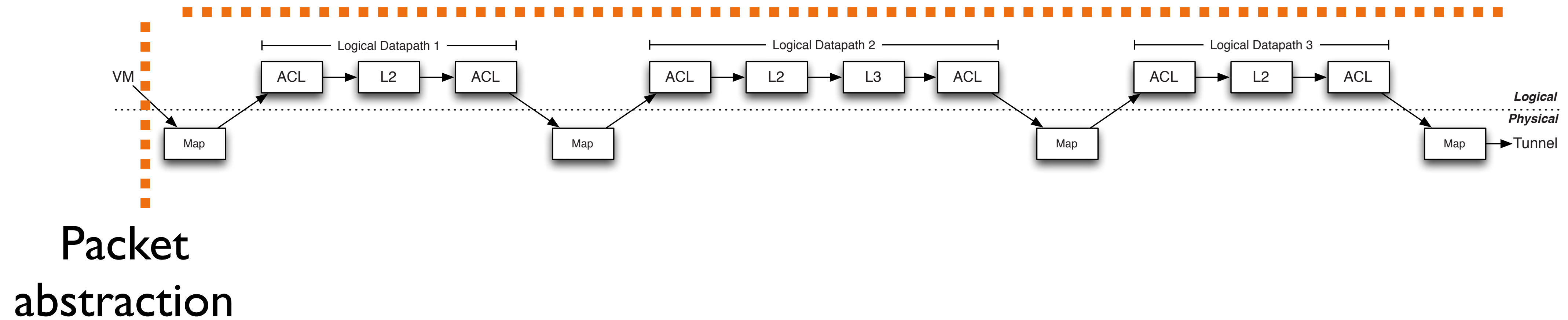
# Virtual network service

# Virtual network service

Control abstraction
(sequence of OpenFlow flow tables)

Tenant control abstraction

L2 — L3 — L2

Network Hypervisor Controllers

Tenant VM

Open vSwitch

server

Tenant VM

Open vSwitch

server

server

Tenant VM

Physical L3 Network

tunnel (GRE, VXLAN)

Open vSwitch

server

server

# Challenge: Performance

Large amount of state to compute

- Full virtual network state at every host with a tenant VM!
- $O(n^2)$ tunnels for tenant with $n$ VMs
- Solution 1: Automated incremental state computation with *nlog* declarative language
- Solution 2: Logical controller computes single set of universal flows for a tenant, translated more locally by "physical controllers"

# Challenge: Performance

Pipeline processing in virtual switch can be slow

- Solution: Send first packet of a flow through the full pipeline; thereafter, put an exact-match packet entry in the kernel

Tunneling interferes with TCP Segmentation Offload (TSO)

- NIC can't see TCP outer header
- Solution: STT tunnels adds "fake" outer TCP header

# Discussion

Where's the SDN?

- API to data plane
- centralized controller
- control abstractions

Why was micro-segmentation a "killer app" for SDN?

- Needed to automate control of a dynamic, virtualized environment, not suited to manual solutions

How does it compare to wide-area control in B4?

# Industry Impact

Multiple vendors with software-defined data center "micro-segmentation" products

- VMware's NSX
- Cisco's ACI
- Startups vArmour, Illumio

- VMware claims more than 2,400 customers, $1B/yr sales

## Next time

- Higher-level programming abstractions for SDN

# Mid-term project presentations

Two key goals

- Demonstrate concrete progress
- Feedback & discussion with your peers

Content

- What problem are you solving?
- Why has past work not addressed the problem?
- What is your approach for solving it?
- What are your preliminary results & progress?