

Towards Predictable Datacenter Networks

Hitesh Ballani, Paolo Costa,
Thomas Karagiannis and Ant Rowstron

Microsoft Research, Cambridge

This talk is about ...

Guaranteeing network performance for tenants in multi-tenant datacenters

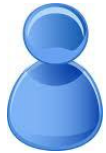
Multi-tenant datacenters

- ▶ Datacenters with multiple (possibly competing) ***tenants***
- ▶ Private datacenters
 - ▶ Run by organizations like Facebook, Intel, etc.
 - ▶ **Tenants:** Product groups and applications
- ▶ Cloud datacenters
 - ▶ Amazon EC2, Microsoft Azure, Rackspace, etc.
 - ▶ **Tenants:** Users renting virtual machines

Cloud datacenters 101

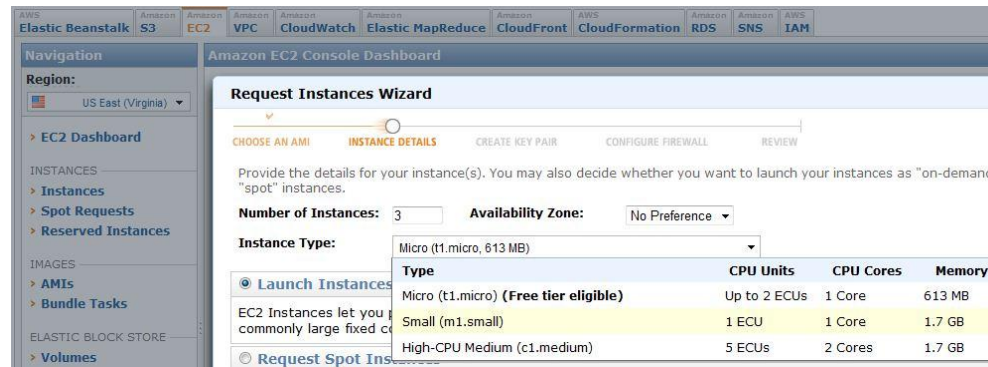
Simple interface: Tenants ask for a set of VMs

Tenant



**Request
VMs**

Amazon EC2 Interface



The screenshot shows the Amazon EC2 Console Dashboard with the 'Request Instances Wizard' active. The wizard has four steps: CHOOSE AN AMI, INSTANCE DETAILS, CREATE KEY PAIR, and CONFIGURE FIREWALL. The 'INSTANCE DETAILS' step is currently selected. It shows 'Number of Instances' set to 3 and 'Availability Zone' set to 'No Preference'. Under 'Instance Type', a dropdown menu is open showing three options: 'Micro (t1.micro, 613 MB)', 'Small (m1.small)', and 'High-CPU Medium (c1.medium)'. A table below the dropdown lists the specifications for these instance types.

Type	CPU Units	CPU Cores	Memory
Micro (t1.micro) (Free tier eligible)	Up to 2 ECUs	1 Core	613 MB
Small (m1.small)	1 ECU	1 Core	1.7 GB
High-CPU Medium (c1.medium)	5 ECUs	2 Cores	1.7 GB



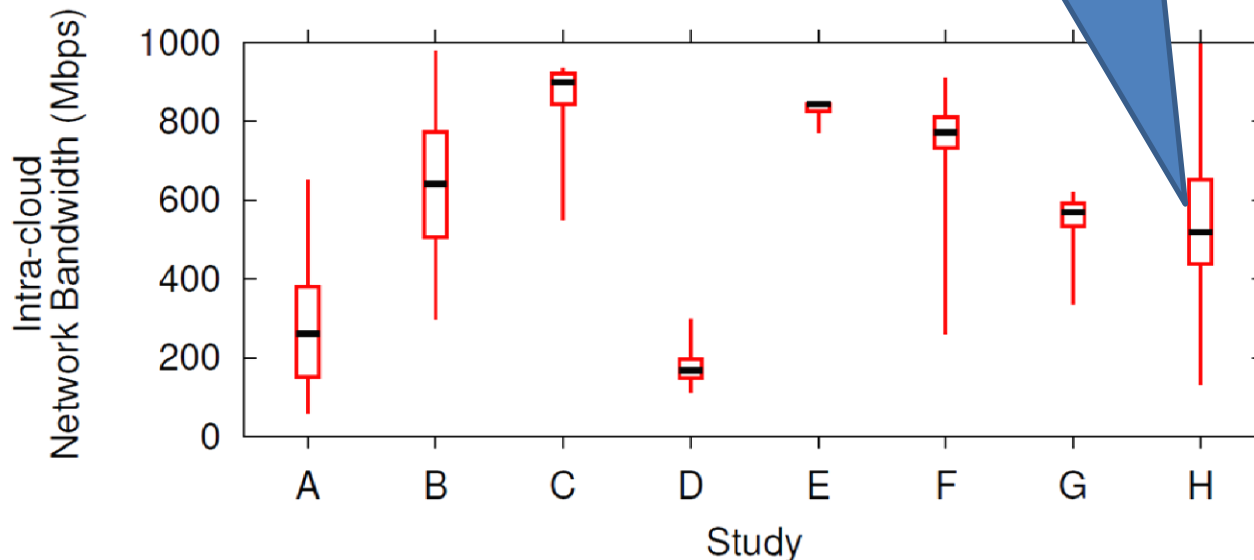
- ▶ Charging is per-VM, per-hour
 - ▶ Amazon EC2 small instances: \$0.085/hour
 - ▶ No (intra-cloud) network cost

Network performance is not guaranteed

Bandwidth between a tenant's VMs depends on their placement, network load, protocols used, etc.

Performance variability in

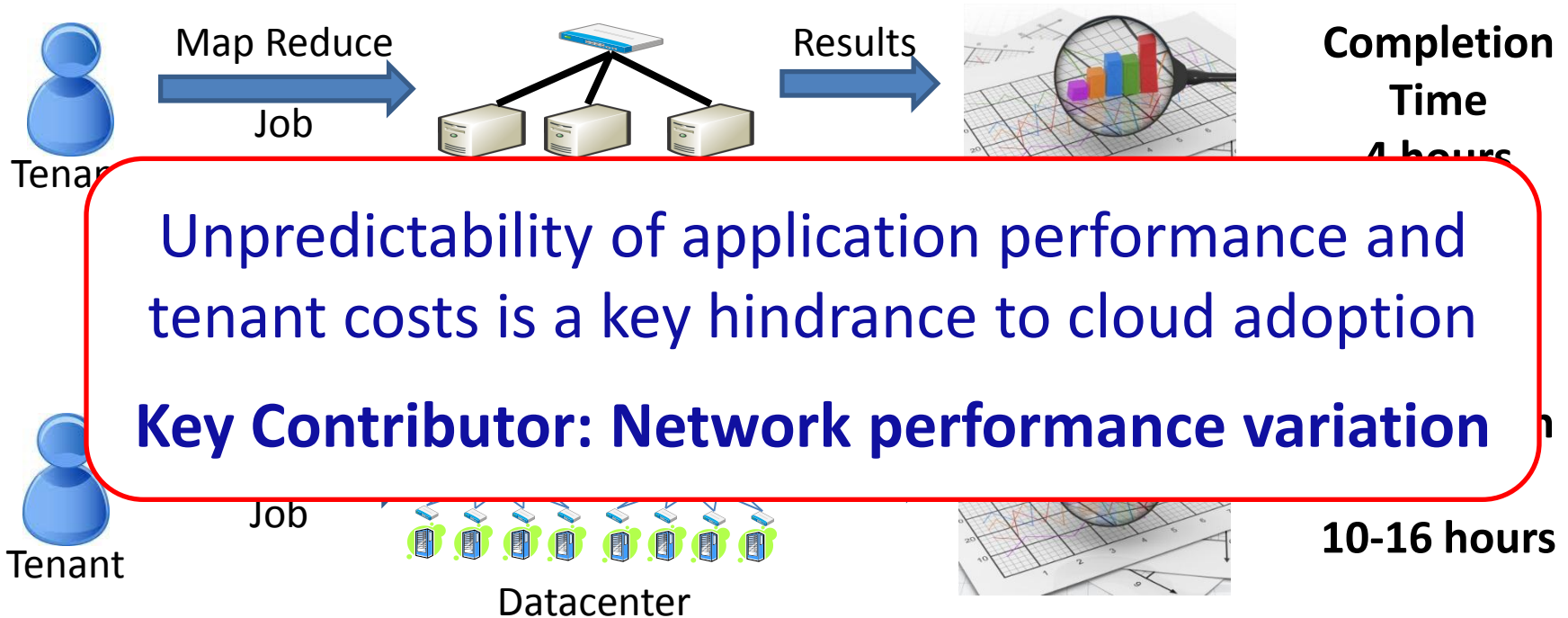
Up to 5x variability



Study	Study	Provider	Duration
A	[Giurgui'10]	Amazon EC2	n/a
B	[Schad'10]	Amazon EC2	31 days
C/D/E	[Li'10]	(Azure, EC2, Rackspace)	1 day
F/G	[Yu'10]	Amazon EC2	1 day
H	[Mangot'09]	Amazon EC2	1 day

Network performance can vary ... so what?

Data analytics on an isolated cluster



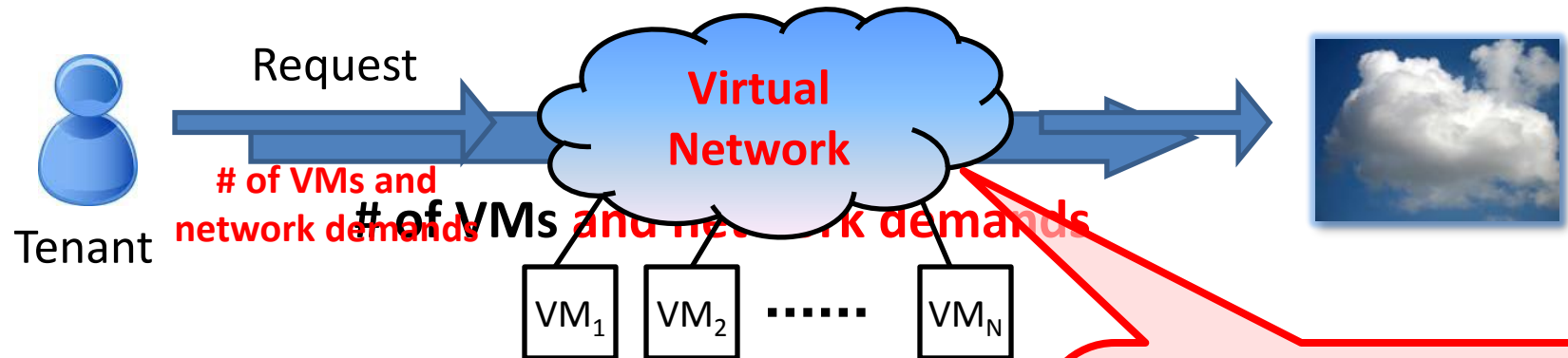
Variable tenant costs

Expected cost (based on 4 hour completion time) = \$100

Actual cost = \$250-400

Predictable datacenter networks

Extend the tenant-provider interface to account for the network



Contributions-

Virtual network abstractions

- ▶ To capture tenant network demands

Oktopus: Proof of concept system

- ▶ Implements virtual networks in multi-tenant datacenters
- ▶ Can be incrementally deployed **today!**

Key Idea: Tenants are offered a virtual network with bandwidth guarantees

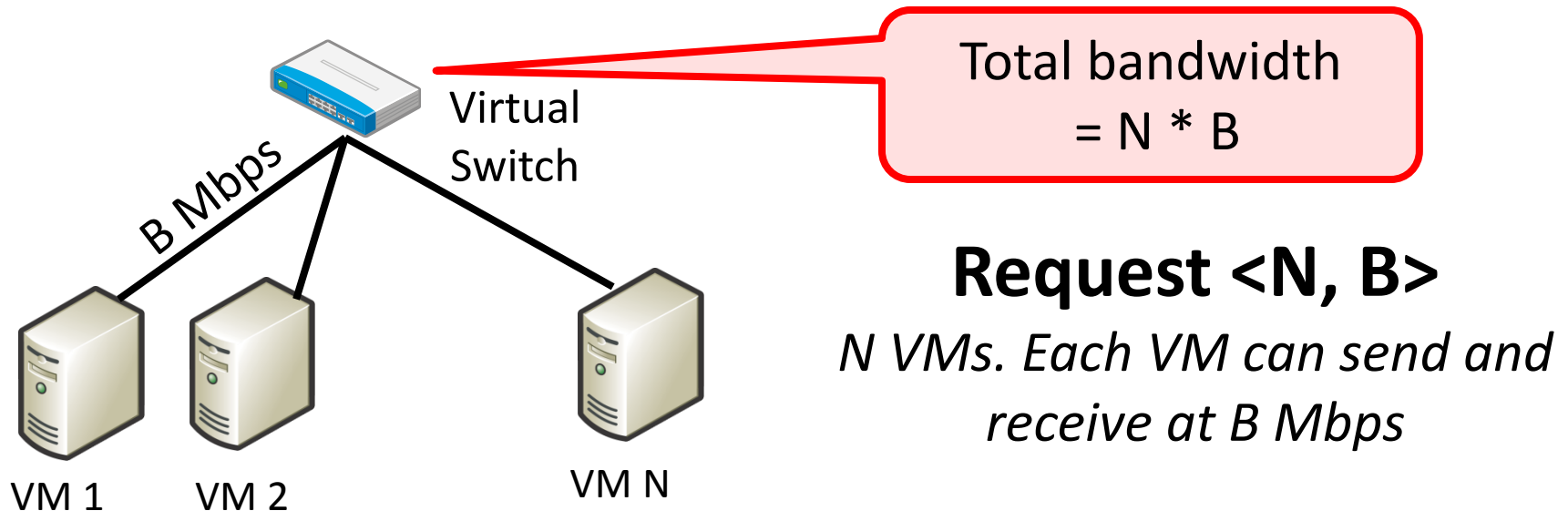
This decouples tenant performance from provider infrastructure

Talk Outline

- ▶ Introduction
- ▶ Virtual network abstractions
- ▶ Oktopus
 - ▶ Allocating virtual networks
 - ▶ Enforcing virtual networks
- ▶ Evaluation

Abstraction 1: Virtual Cluster (**VC**)

Motivation: In enterprises, tenants run applications on dedicated Ethernet clusters



Tenants get a network with no oversubscription

- ✓ Suitable for data-intensive apps. (MapReduce, BLAST)
- ✗ Moderate provider flexibility

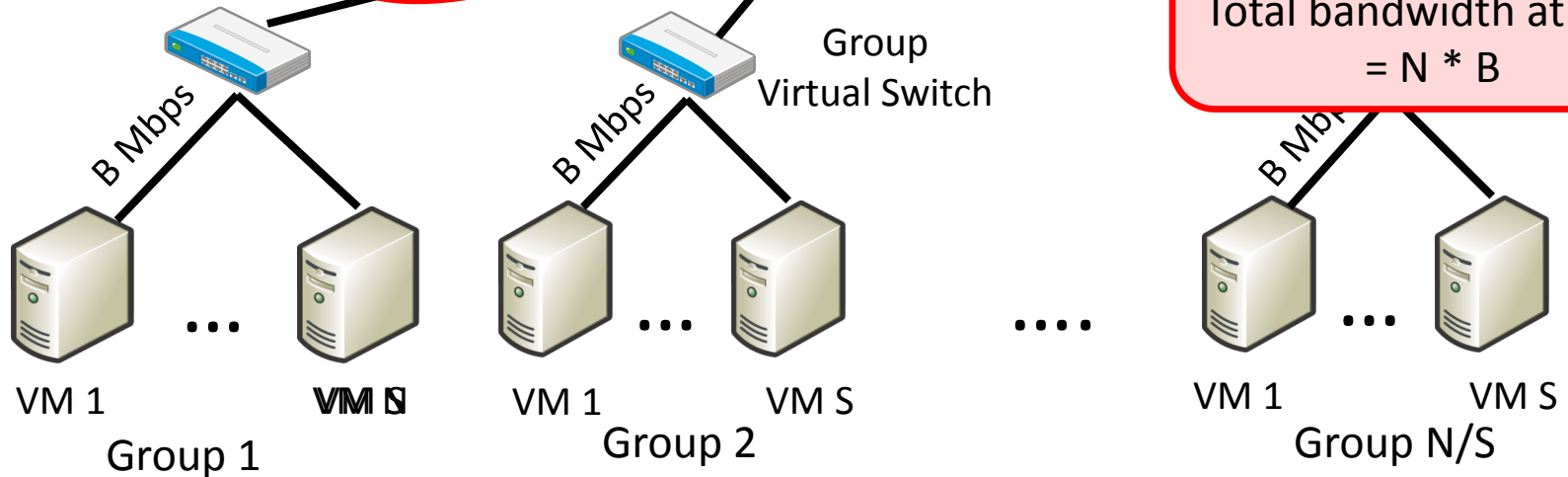
Abstraction 2: Virtual Oversubscribed Cluster (**VOC**)

VMs can send traffic to group members at B Mbps

$$B * S / O \text{ Mbps}$$

Root
Virtual Switch

Total bandwidth at root
 $= N * B / O$
Total bandwidth at VMs
 $= N * B$



Request $\langle N, B, S, O \rangle$

Motivation: Many applications moving to the cloud have

VOC capitalizes on tenant communication patterns

- ✓ Suitable for typical applications (*though not all*)
- ✓ Improved provider flexibility

Talk Outline

- ▶ Introduction
- ▶ Virtual network abstractions
- ▶ **Oktopus**
 - ▶ **Allocating virtual networks**
 - ▶ Enforcing virtual networks
- ▶ Evaluation

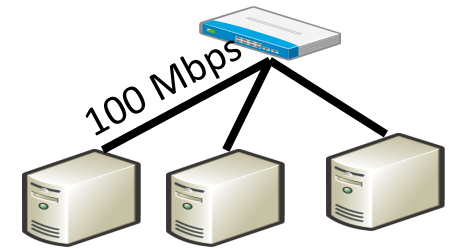
Oktopus

Offers virtual networks to tenants in datacenters

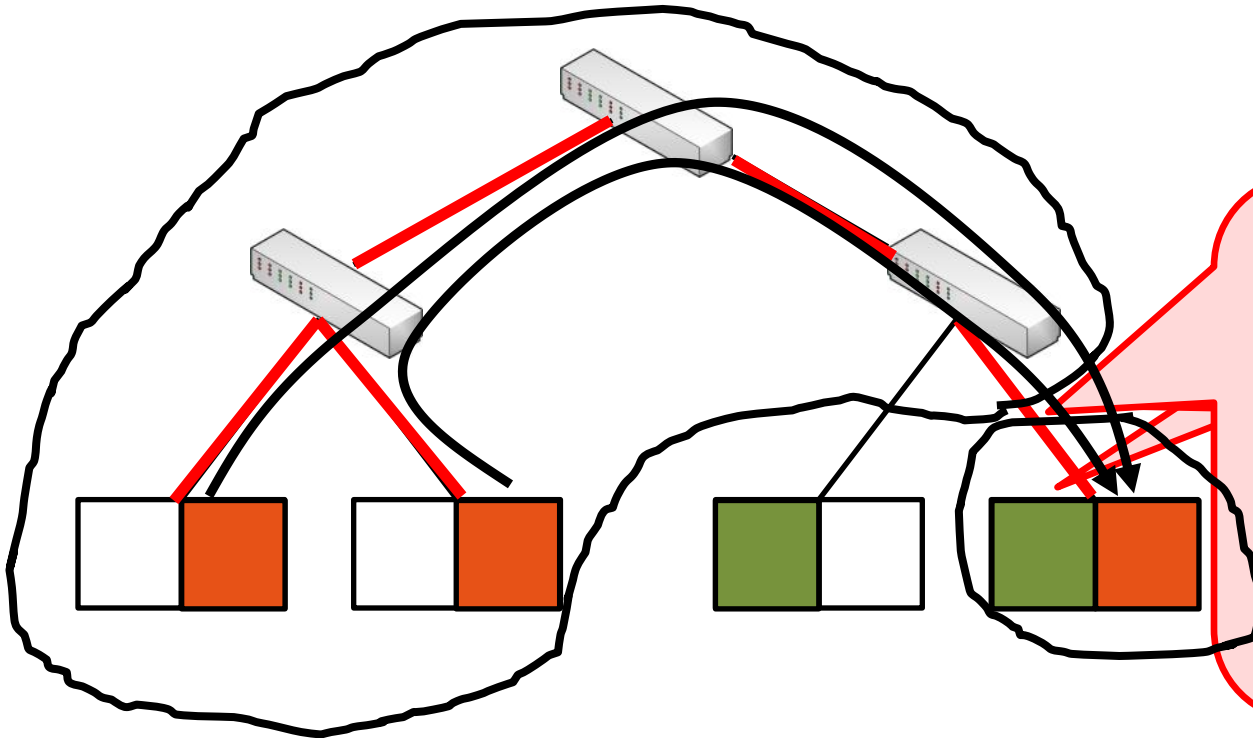
Two main components

- ▶ Management plane: ***Allocation of tenant requests***
 - ▶ Allocates tenant requests to physical infrastructure
 - ▶ Accounts for tenant network bandwidth requirements
- ▶ Data plane: ***Enforcement of virtual networks***
 - ▶ Enforces tenant bandwidth requirements
 - ▶ Achieved through rate limiting at end hosts

Allocating Virtual Clusters



Request : <3 VMs, 100 Mbps>



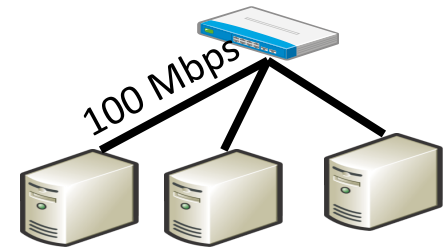
Max Sending Rate =
 $2 * 100 = 200$

Max Receive Rate =
 $1 * 100 = 100$

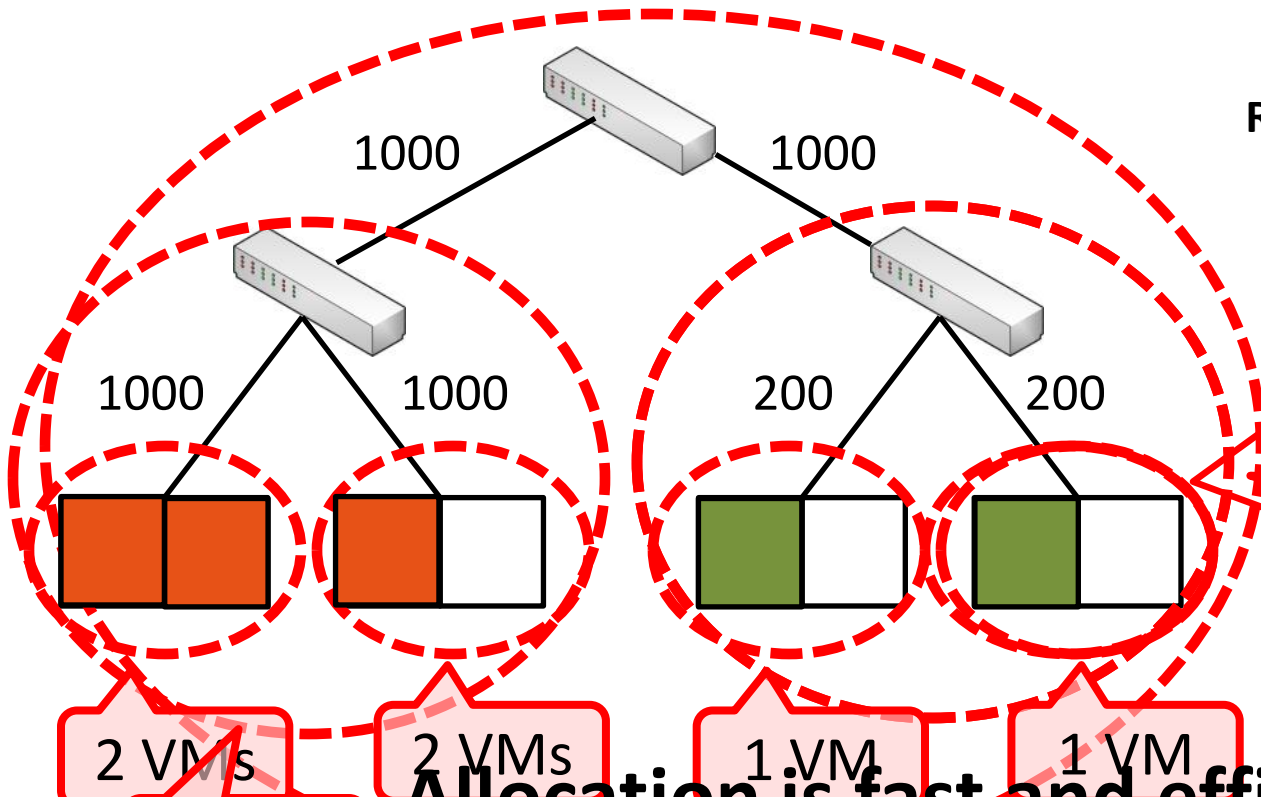
B/W needed on link =
 $\text{Min}(200, 100) =$
100Mbps

For a virtual cluster $\langle N, B \rangle$, bandwidth needed on a link that connects an upstream tenant $\langle M, B \rangle$ to N VMs is $\text{Min}(m, N) * B$.
 An allocation of tenants $\langle M, B \rangle$ into two parts, $\langle m, B \rangle$ and $\langle M-m, B \rangle$, is a **valid allocation** if traffic from the left to right part cluster with 100 Mbps each, i.e. $\langle 3 \text{ VMs}, 100 \text{ Mbps} \rangle$.
 Bandwidth needed \leq Link's Residual Bandwidth

Allocation Algorithm



Request : <3 VMs, 100 Mbps>



Solution

At most 1 VM for this tenant can be allocated here

Allocation is fast and efficient

Greedy Allocation algorithm

- Constraints for # of VMs (key information allocated to the machine-
3 VMs)
1. Packing VMs together motivated by the fact that VMs only need to be placed on one machine.
 2. Traverse up the hierarchy and determine the lowest level at which all 3 VMs can be allocated.
 3. Allocation can be extended for $(m \times 100) \leq 200$, etc.
- number of VMs that can be allocated to any level of the datacenter, machines, racks and so on

Talk Outline

- ▶ Introduction
- ▶ Virtual network abstractions
- ▶ **Oktopus**
 - ▶ Allocating virtual networks
 - ▶ **Enforcing virtual networks**
- ▶ Evaluation

Enforcement in Oktopus: Key highlights

Oktopus enforces virtual networks at end hosts

- ▶ Use egress rate limiters at end hosts
 - ▶ Implement on hypervisor/VMM

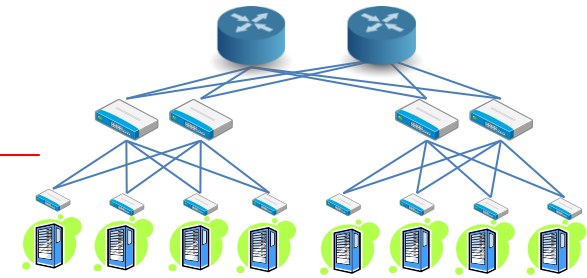
Oktopus can be deployed *today*

- ▶ No changes to tenant applications
- ▶ No network support
- ▶ Tenants without virtual networks can be supported
 - ▶ Good for incremental roll out

Talk Outline

- ▶ Introduction
- ▶ Virtual network abstractions
- ▶ Oktopus
 - ▶ Allocating virtual networks
 - ▶ Enforcing virtual networks
- ▶ Evaluation

Datacenter Simulator



Flow-based simulator

- ▶ 16,000 servers and 4 VMs/server \Rightarrow 64,000 VMs
- ▶ Three-tier network topology (10:1 oversubscription)

Tenants submit requests for VMs and execute jobs

- ▶ Job: VMs process and shuffle data between each other

Baseline: representative of today's setup

- ▶ Tenants simply ask for VMs
- ▶ VMs are allocated in a locality-aware fashion

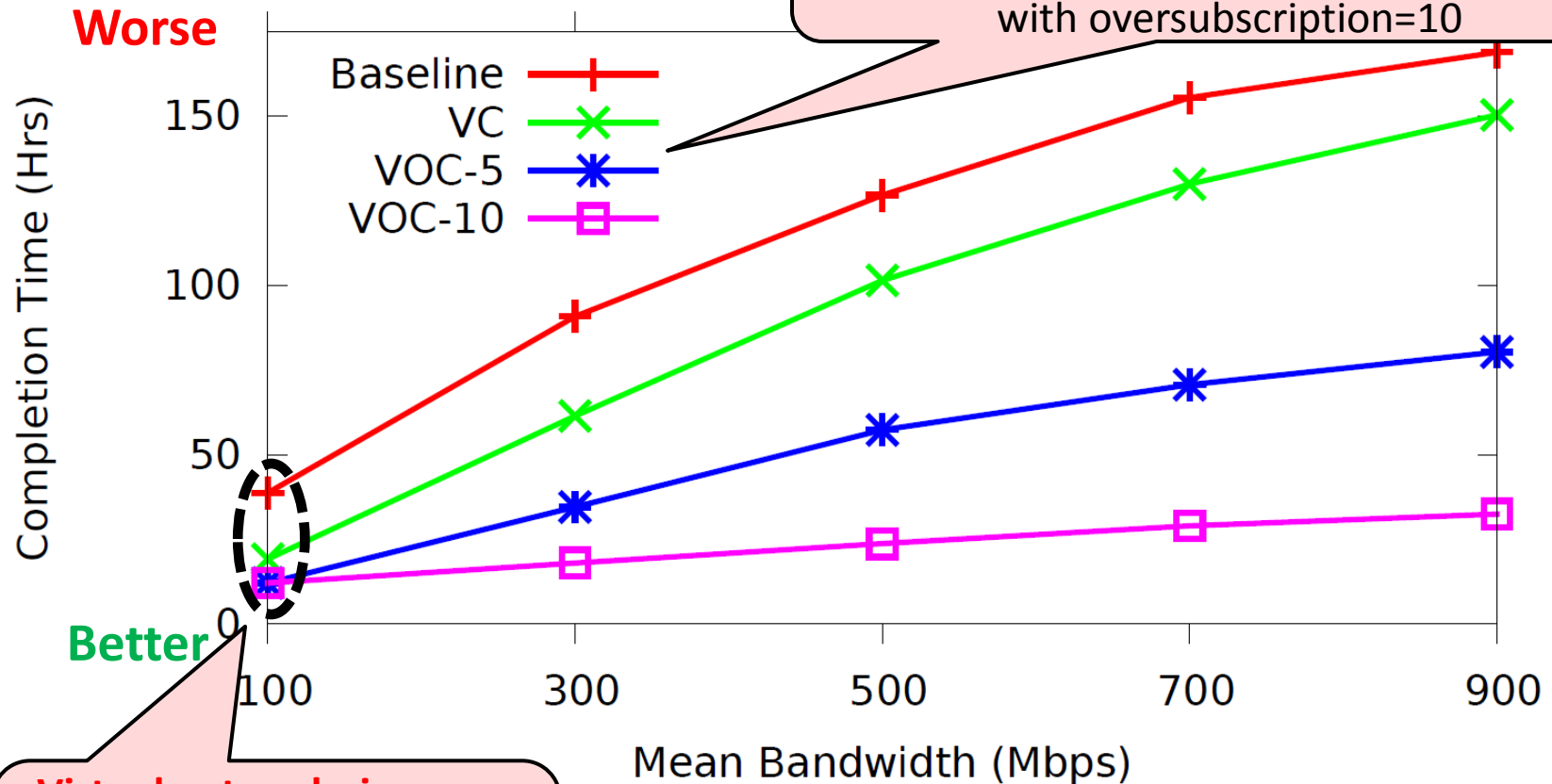
Virtual network request

- ▶ Tenants ask for Virtual Cluster (VC) or Virtual Oversubscribed Cluster (VOC)

Private datacenters

VC is Virtual Cluster

VOC-10 is Virtual Oversubscribed Cluster with oversubscription=10



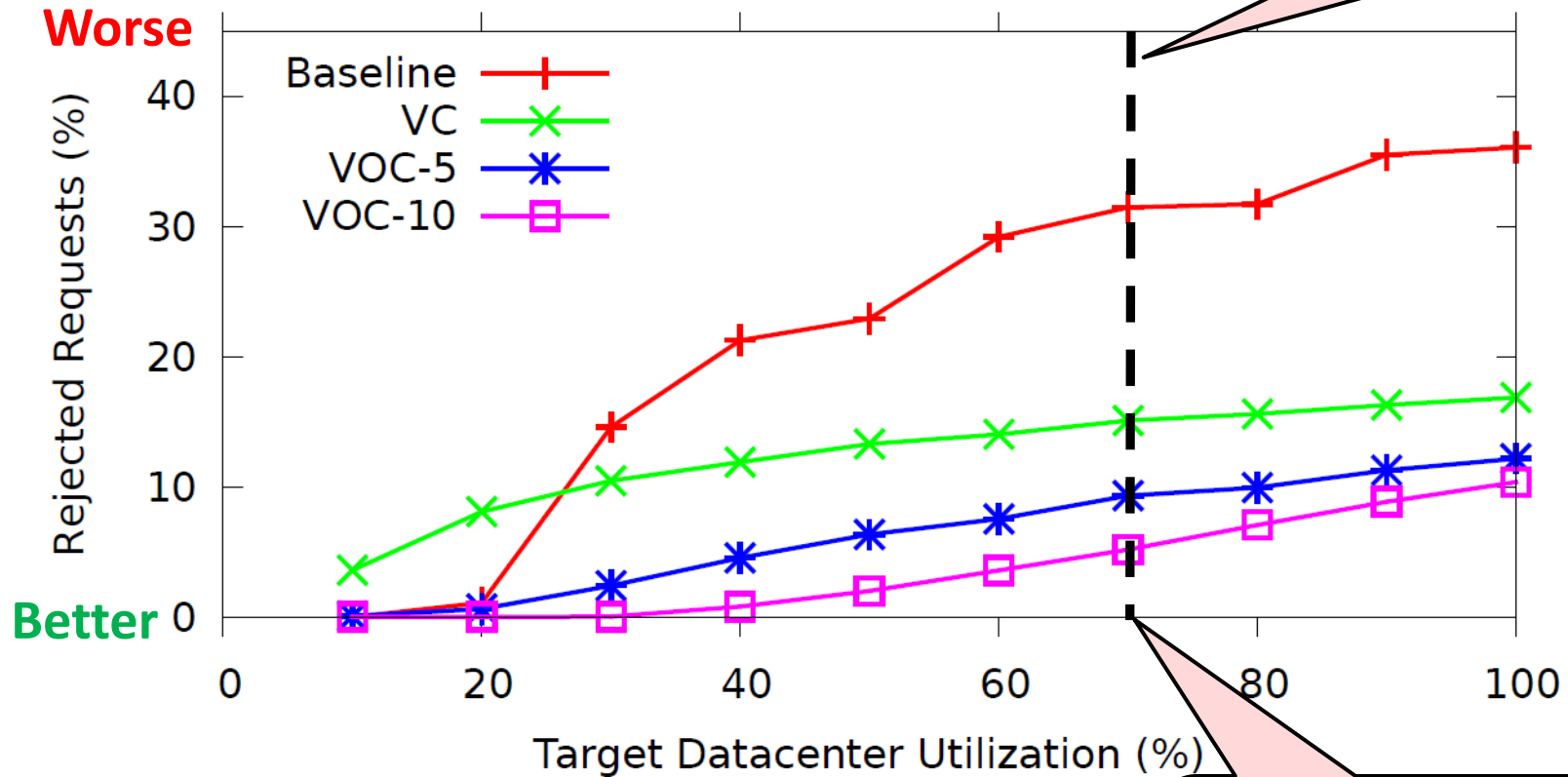
Virtual networks improve completion time

VC: 50% of Baseline
VOC-10: 31% of Baseline

As bandwidth increases, systems become more network intensive

Cloud Datacenters

Amazon EC2's
reported target
utilization



Job requests arrive faster

Rejected Requests

Baseline: 31%

VC: 15%

VOC-10: 5%

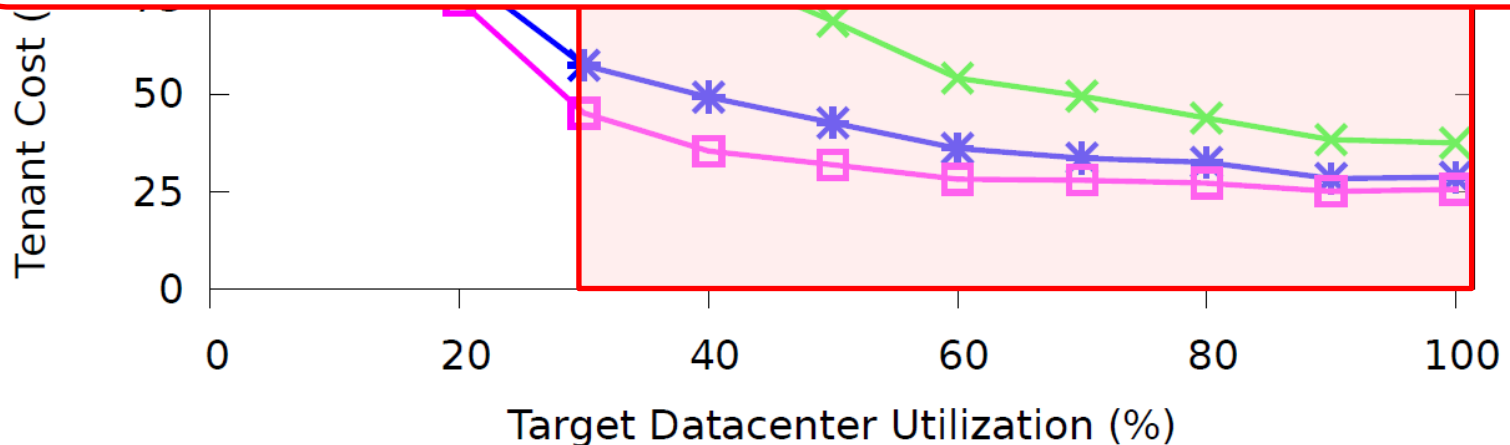
Tenant Costs

What should tenants pay to ensure *provider revenue neutrality*, i.e. provider revenue remains the same with all approaches

Based on today's EC2 prices, i.e. \$0.085/hour for each VM

Provider revenue increases while tenants pay less

At 70% target utilization, provider revenue increases by 20% and median tenant cost reduces by 42%



Oktopus Deployment

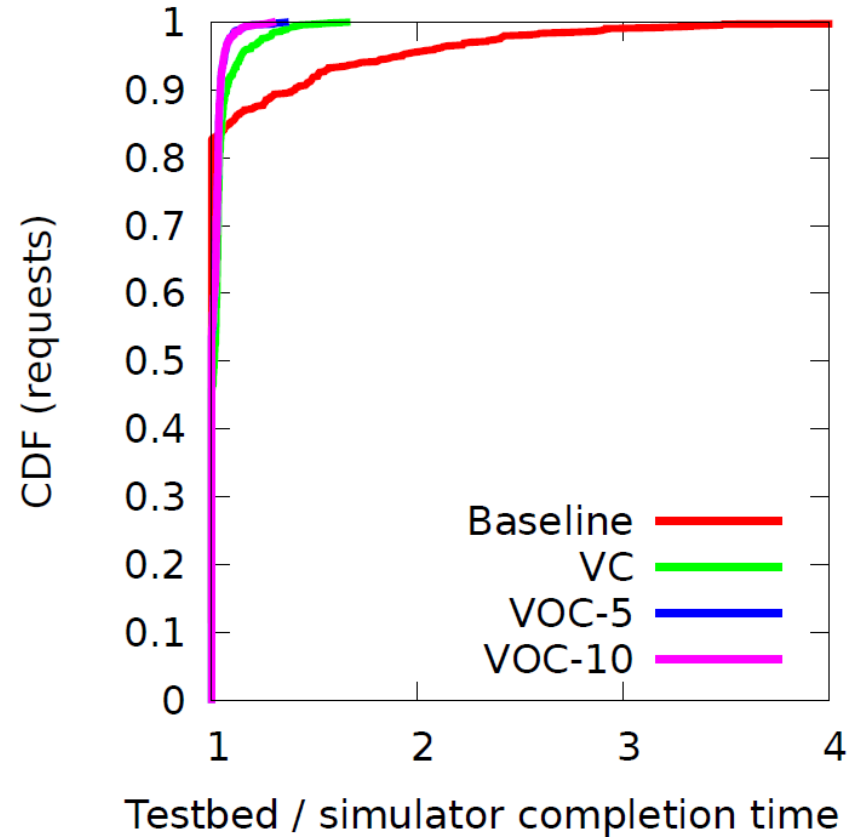
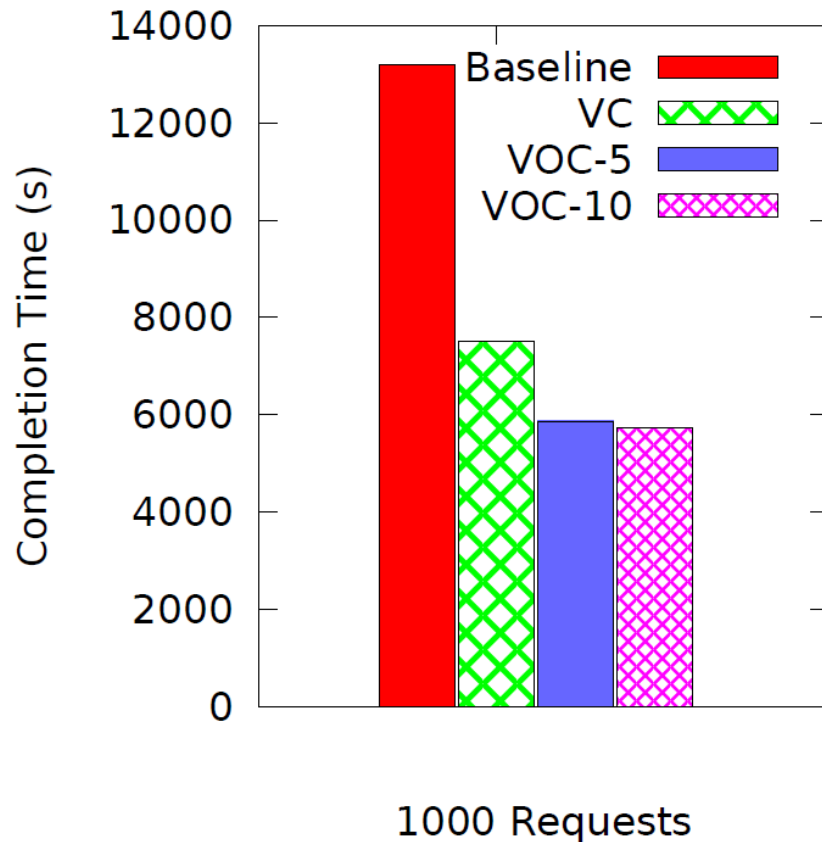
Implementation scales well and imposes low overhead

- ▶ Allocation of virtual networks is fast
 - ▶ In a datacenter with 10^5 machines, median allocation time is 0.35ms
- ▶ Enforcement of virtual networks is cheap
 - ▶ Use Traffic Control API to enforce rate limits at end hosts

Deployment on testbed with 25 end hosts

- ▶ End hosts arranged in five racks

Oktopus Deployment



Cross-validation of simulation results

Completion time for jobs in the simulator matches that on the testbed

Summary

Proposal: Offer virtual networks to tenants

- ▶ Virtual network abstractions
 - ▶ Resemble physical networks in enterprises
 - ▶ Make transition easier for tenants

Proof of concept: Oktopus

- ▶ Tenants get guaranteed network performance
- ▶ Sufficient multiplexing for providers
- ▶ Win-win: *tenants pay less, providers earn more!*

How to determine tenant network demands?

Ongoing work: Map high-level goals (like desired completion time) to Oktopus abstractions

Thank you