Microsoft Research





Bridging the Gap Between Applications and Networks in Data Centers Paolo Costa

Microsoft Research Cambridge

joint work with

Hitesh Ballani, Thomas Karagiannis, Ant Rowstron

Motivation



7.5 GB memory
4 EC2 Compute Units (2 virtual cores with 2 EC2 Compute Units each)
850 GB instance storage
64-bit platform
I/O Performance: High
API name: m1.large

Network performance is not guaranteed!







Is this really the case?



Study	Study	Provider	Duration
А	[Giurgui'10]	Amazon EC2	n/a
В	[Schad'10]	Amazon EC2	31 days
C/D/E	[Li'10]	(Azure, EC2, Rackspace)	1 day
F/G	[Yu'10]	Amazon EC2	1 day
Н	[Mangot'o9]	Amazon EC2	1 day

Up to 5x variability

Is this really the case?



Study	Study	Provider	Several Causes
А	[Giurgui'10]	Amazon EC2	· VMs placement
В	[Schad'10]	Amazon EC2	· Network load
C/D/E	[Li'10]	(Azure, EC2, Rackspa	· Protocols used
F/G	[Yu'10]	Amazon EC2	•
Н	[Mangot'o9]	Amazon EC2	ı day

Fair enough...it's free so who cares?



Fair enough...it's free so who cares?

Data analytics on an isolated cluster



Completion Time 4 hours

Data analytics in a multi-tenant data center



Completion Time 10-16 hours



from Mosharaf Chowdhury et al. "Managing data transfers in computer clusters with orchestra", SIGCOMM'11

Paolo Costa

Fair enough....it's free so who cares?

Data analytics on an isolated cluster



Completion Time 4 hours

Data analytics in a multi-tenant data center



Completion Time 10-16 hours

Variable tenant costs

Expected cost (based on 4 hour completion time) = \$100 Actual cost = \$250-400

Paolo Costa

Fair enough...it's free so who cares?

Data analytics on an isolated cluster

Lack of predictability

Results

Unpredictability of application performance and tenant costs is a key hindrance to cloud adoption

Tenant

Job



10-16 hours

omplation

Variable tenant costs Expected cost (based on 4 hour completion time) = \$100 Actual cost = \$250-400

Paolo Costa

Addressing the Problem

• Amazon EC2 cluster compute

✓ Guaranteed 10 Gbps bandwidth

- Very expensive for provider / tenants
 - \odot Requires full-bisection bandwidth and 1VM/server
 - \circ \$ 1.68 / hour \cong 20 small instances
- Very limited tenant flexibility
 - either 10Gbps or nothing

23 GB of memory

33.5 EC2 Compute Units (2 x Intel Xeon X5570, quad-core "Nehalem" architecture) 1690 GB of instance storage

64-bit platform

I/O Performance: Very High (10 Gigabit Ethernet)

API name: cc1.4xlarge

Paolo Costa

Predictable Data Center Networks



Predictable Data Center Networks



Predictable Data Center Networks



- Extend the tenant-provider interface to account for the network
- Easier transition for tenants
 - Tenants should be able to predict the performance of applications running atop the virtual network
- Provider flexibility
 - Providers should be able to multiplex many virtual networks on the physical network

These are competing design goals

Our abstractions strive to strike a balance between them

Abstraction #1: Virtual Cluster



Request <N> NVMs

Abstraction #1: Virtual Cluster



Abstraction #1: Virtual Cluster

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



Request <N, B>

= N * B

N VMs. Each VM can send and receive at B Mbps

Tenants get a network with no oversubscription

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







Motivation: Many applications moving to the cloud have localized communication patterns

Applications are composed of groups with more traffic within groups than across groups

Paolo Costa



No oversubscription for intra-group communication Intra-group communication is the common case! Bridging the Gap Between Applications and Networks in Data Centers

Paolo Costa



Oversubscription factor O for inter-group communication

(captures the sparseness of inter-group communication) Bridging the Gap Between Applications and Networks in Data Centers

Paolo Costa



VOC capitalizes on tenant communication patterns

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

Oktopus [SIGCOMM'11]

• Offers virtual networks to tenants in data centers



Oktopus [SIGCOMM'11]

- Offers virtual networks to tenants in data centers
- Management plane: Allocation of tenant requests
 - Goal: Maximize the ability to allocate future requests
 - \circ Variant of network embedding problem (NP-hard)
 - \odot Existing heuristics don't scale beyond O(10²) nodes
 - Key idea: Restrict the set of virtual and physical topologies
 O Allocation is fast (median time is 0.35ms)
- Data plane: Enforcement of virtual networks
 - Enforcement of virtual networks is cheap

 Implemented on hypervisor only (no network support)
 It can be deployed today

Rejected Requests







• What should tenants pay to ensure *provider revenue neutrality*?

 kN_rT_r

k= \$0.085 /hour N= # of VMs T = completion time

• What should tenants pay to ensure provider revenue neutrality? $-\sum_r (kN_rT_r) = \sum_r (N_rT_r(k_n + k_bB))$

k= \$0.085 /hour N= # of VMs T = completion time

 What should tenants pay to ensure provider revenue neutrality?

 $-\sum_{r}(kN_{r}T_{r}) = \sum_{r}(N_{r}T_{r}(k_{n} + k_{b}B))$





• What should tenants pay to ensure provider revenue neutrality? $-\sum_r (kN_rT_r) = \sum_r (N_rT_r(k_n + k_bB))$

Other pricing models are possible, e.g., Provider revenue increases while tenants pay less

Provider revenue increases by 20% and median tenant cost reduces by 42%"

20

Bridging the Gap Between Applications and Networks in Data Centers

Load (%)

Beyond Oktopus

- Issues with Oktopus model
 - Hard for tenants to derive B (and even N)...
 - Unfortunate <N,B> choices can hurt provider's throughput
- Vision: Job-centric interface
 - Tenants specify what they want
 - Performance: complete the job in 1 hour
 - O Cost: complete the job for less than \$ 100
 - Providers decide how to achieve this
 - \odot E.g., job needs 10 VMs and 200 Mbps of network

Why bother?

- Tenants
 - Simpler and more intuitive interface
 - Especially for non-IT experts
- Provider
 - Exploit multi-resource tradeoff



Why bother?

Intuition

Provider can choose the *most convenient* combination of <N,B> that satisfies tenants' goals

LinkGraph in 300s





Tenants expose constraints

Performance: Complete job in less than 3 hours *Cost*: Complete job for less than \$400

System Model



Bazaar determines resources required

Example, 10 VMs with guaranteed 200Mbps network between them will ensure the job finishes in less than 3 hours

Paolo Costa

System Model



Multiple resource combinations may satisfy tenant goal

Completion time = 3 hrs ⇒ 1). 25 VMs with 500 Mbps, or 2). 40 VMs with 400 Mbps

System Model



Provider can use this flexibility

{25 VMs with 500 Mbps} vs. {40 VMs with 400 Mbps} Choose the resource combination to minimize internal impact

Paolo Costa

Bazaar [SoCC'12]





Performance prediction

- Translating performance goal is hard
 - For general applications...possibly impossible
 - For popular cloud-based applications...maybe!
- Early evidence
 - Amazon Dynamo DB (e.g., 100K reqs/s)
 - We focus on MapReduce: (simple) analytic model

Paolo Costa

Bazaar [SoCC'12]



Resource selection

- Multiple resource combination may satisfy user goals
 E.g., 25 VMs with 500 Mbps or 40 VMs with 400 Mbps
- Similar to Multi-dimensional Bin Packing
- Heuristic: Minimize Resource Imbalance Metric
 - Choose the combination that minimizes resource imbalance

Paolo Costa

Prediction Accuracy

• Setup: Hadoop on 35-node Emulab cluster



Evaluation







Pricing Implications

- Today's resource-based pricing
 - Results in mismatched tenant-provider interests

Tenants want cheapest resource tuple <N, B>

- \circ Provider wants resource tuple with least impact
- No incentive to reduce completion time!
- Bazaar enables job-based pricing
 - E.g.: Finish Sort over 200GB in 4hrs costs \$100
 - Price based on job characteristics and completion time
- Finish jobs "before time"
 - Intuition: if spare resources, finish the job earlier so more resources will be available in future

Better alignment of tenant-provider interests

Data Center Network Research

Many proposals:

- **Oversubscription**: Fat-tree[SIGCOMM'08], VL2[SIGCOMM'09], ...
- Path collision: Hedera[NSDI'10], MPTCP[SIGCOMM'11], SPAIN[NSDI'10], ...
- TCP Incast: DCTCP [SIGCOMM'10], ICTCP[CoNEXT'10], FDS[OSDI'12],...
- Traffic prioritization: Orchestra [SIGCOMM'11], D2TCP[SIGCOMM'11], ...
- Fair sharing: Seawall [NSDI'11], FairCloud [SIGCOMM'12], ...

Paolo Costa

Data Center Network Research

Many proposals:

• **Oversubscription**: Fat-tree[SIGCOMM'08], VL2[SIGCOMM'09], ...

The network is a black box for applications (and vice versa) Applications need to reverse-engineer the network The network has to infer application patterns

Paolo Costa

Internet & Data Centers

• This is due to how the Internet was designed...

Internet

- Multiple administration domains
- Heterogeneous HW and network
- Topology not known
- Malicious software



Internet & Data Centers

This is due to how the Internet was designed...
 - ...but data centers are not mini-Internets

Internet

- Multiple administration domains
- Heterogeneous HW and network
- Topology not known
- Malicious software

Data Centers

- Single administration domain
- Homogenous HW and network
 x86 and Ethernet
- Topology known
 - and can be customised
- Trusted components
 - e.g., using virtualization

Takeaway

- Data centers present both unique challenges and opportunities to network designers
 - Good time to revisit previous assumptions and rethink application and protocol design
- Bridging the gap
 - Enable applications to "control" the network
 - Expose application logic to the network



CamCube [SIGCOMM'10,NSDI'12, HPDC'13]



Predictable Data Centers [SIGCOMM'11,HotNets'11, SoCC'12]

Takeaway

- Data centers present both unique challenges and opportunities to network designers
 - Good time to revisit previous assumptions and rethink application and protocol design
- Bridging the gap
 - Enable applications to "control" the network
 - Expose application logic to the network
- Next steps: Network-as-a-Service (NaaS) [Hotlce'12]
 - Fully programmable cloud
 - Remove the distinction between network and computation