

# IaaS Cloud Benchmarking: Approaches, Challenges, and Experience



**Alexandru Iosup**

**Parallel and Distributed Systems Group  
Delft University of Technology  
The Netherlands**

**Our team:** Undergrad Nassos Antoniou, Thomas de Rooter, Ruben Verboon, ...  
**Grad** Siqi Shen, Nezih Yigitbasi, Ozan Sonmez **Staff** Henk Sips, Dick Epema,  
Alexandru Iosup **Collaborators** Ion Stoica and the Mesos team (UC Berkeley),  
Thomas Fahringer, Radu Prodan (U. Innsbruck), Nicolae Tapus, Mihaela Balint, Vlad  
Posea (UPB), Derrick Kondo, Emmanuel Jeannot (INRIA), Assaf Schuster, Orna Ben-  
Yehuda (Technion), Ted Willke (Intel), Claudio Martella (Giraph), ...

# The Parallel and Distributed Systems Group at TU Delft



**VENI**

Alexandru Iosup

Grids/Clouds  
P2P systems  
Big Data  
Online gaming



Dick Epema

Grids/Clouds  
P2P systems  
Video-on-demand  
e-Science



**VENI**

Ana Lucia Varbanescu

HPC systems  
Multi-cores  
Big Data  
e-Science



Henk Sips

HPC systems  
Multi-cores  
P2P systems



**VENI**

Johan Pouwelse

P2P systems  
File-sharing  
Video-on-demand

## Home page

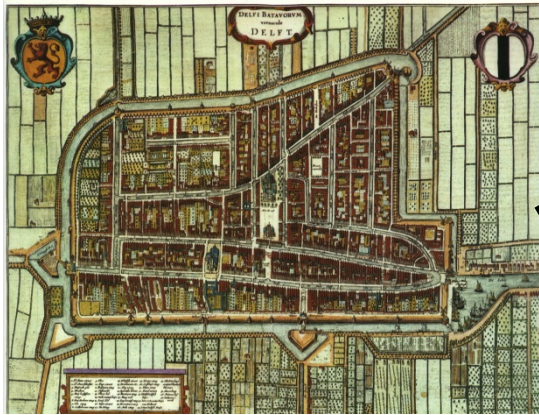
- [www.pds.ewi.tudelft.nl](http://www.pds.ewi.tudelft.nl)

## Publications

- see PDS publication database at [publications.st.ewi.tudelft.nl](http://publications.st.ewi.tudelft.nl)



# (TU) Delft – the Netherlands – Europe



founded 13<sup>th</sup> century  
pop: 100,000



founded 1842  
pop: 13,000



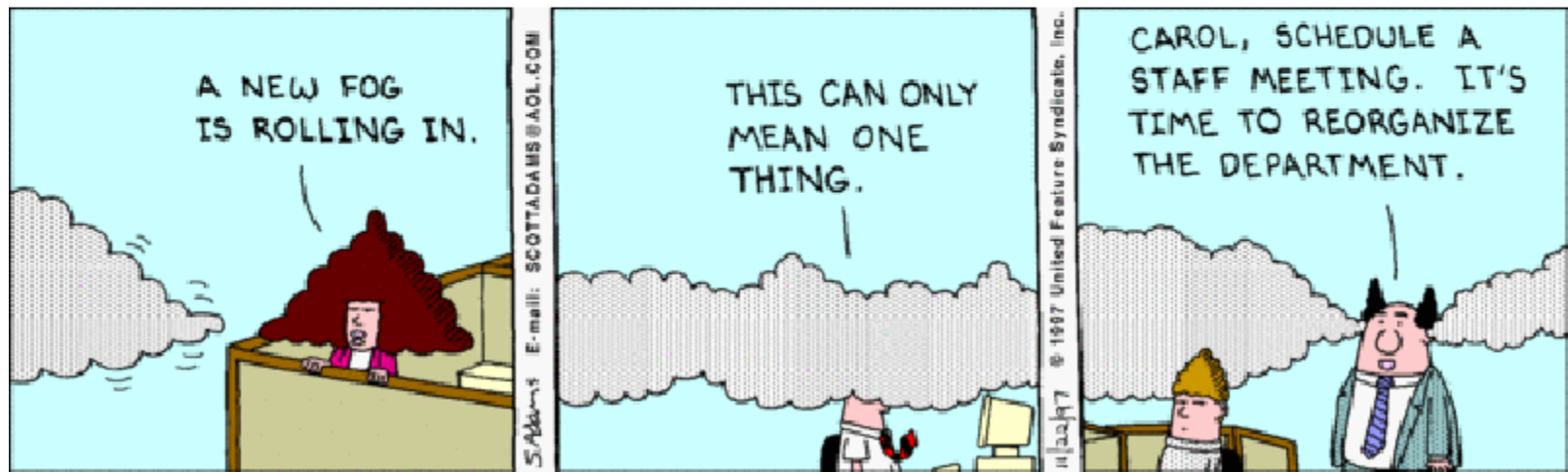
pop: 16.5 M



# What is Cloud Computing?

## 1. A Cloudy Buzzword

- 18 definitions in computer science (ECIS'10). NIST has one. Cal has one. We have one.
- "We have redefined cloud computing to include **everything that we already do.**" Larry Ellison, Oracle, 2009



Source: <http://dilbert.com/strips/comic/1997-11-22/>

# What is Cloud Computing?

## 2. A Descendant\* of the Grid Idea

\* Subset.



Source: <http://royal.pingdom.com/2008/04/11/map-of-all-google-data-center-locations/>

“A computational grid is a hardware and software infrastructure that provides dependable, consistent, pervasive, and inexpensive access to high-end computational capabilities [+ for] nontrivial QoS.” I. Foster, 1998 + 1999

Cloud MW Stack

**Cloud**  
~~Grid~~ Applications

**Cloud**  
~~Grid~~ Very High Level MW

**Cloud**  
~~Grid~~ High Level MW

**Cloud**  
~~Grid~~ Low Level MW

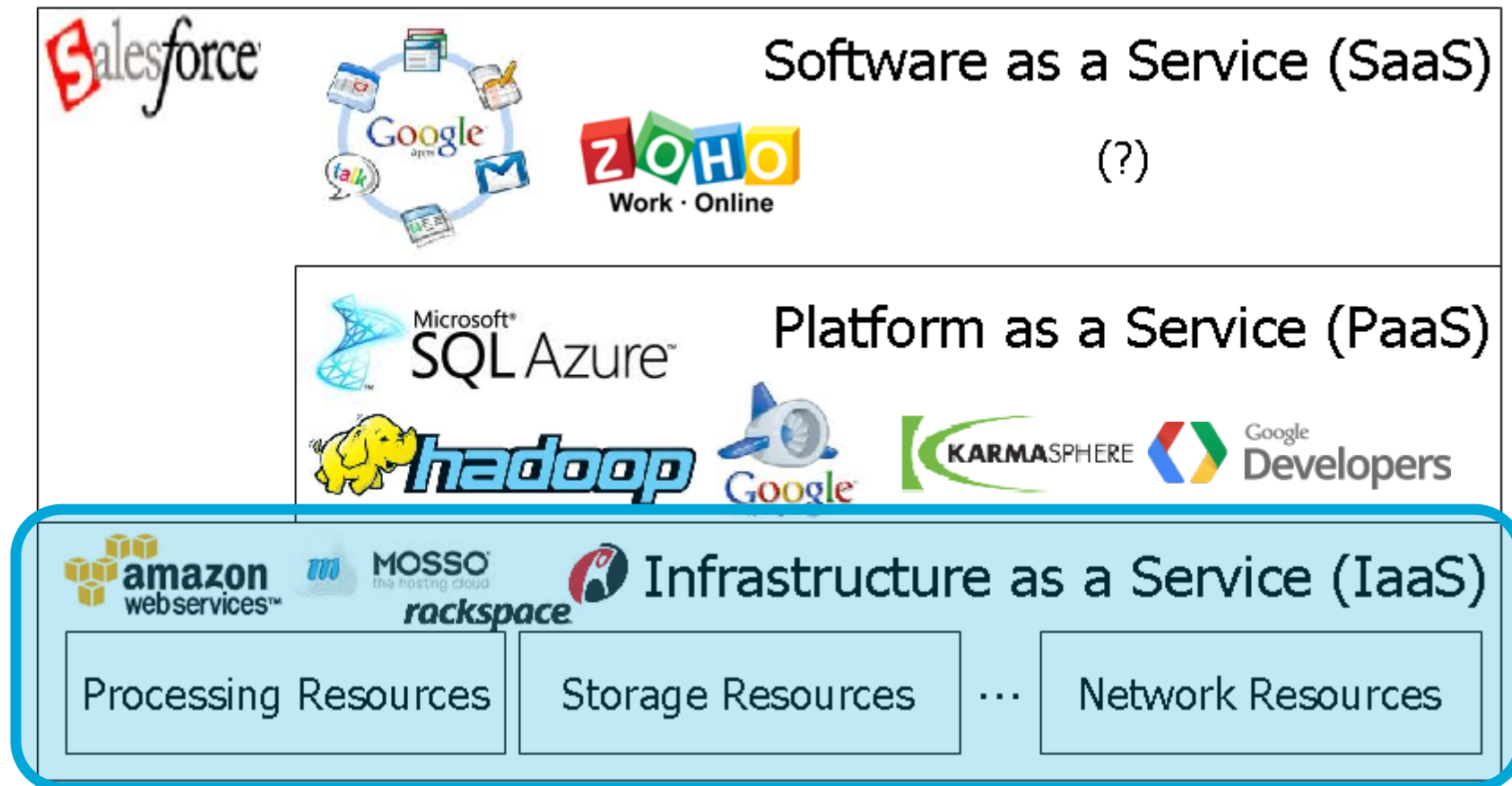
**Virtualized** HW + OS

MW = Middleware

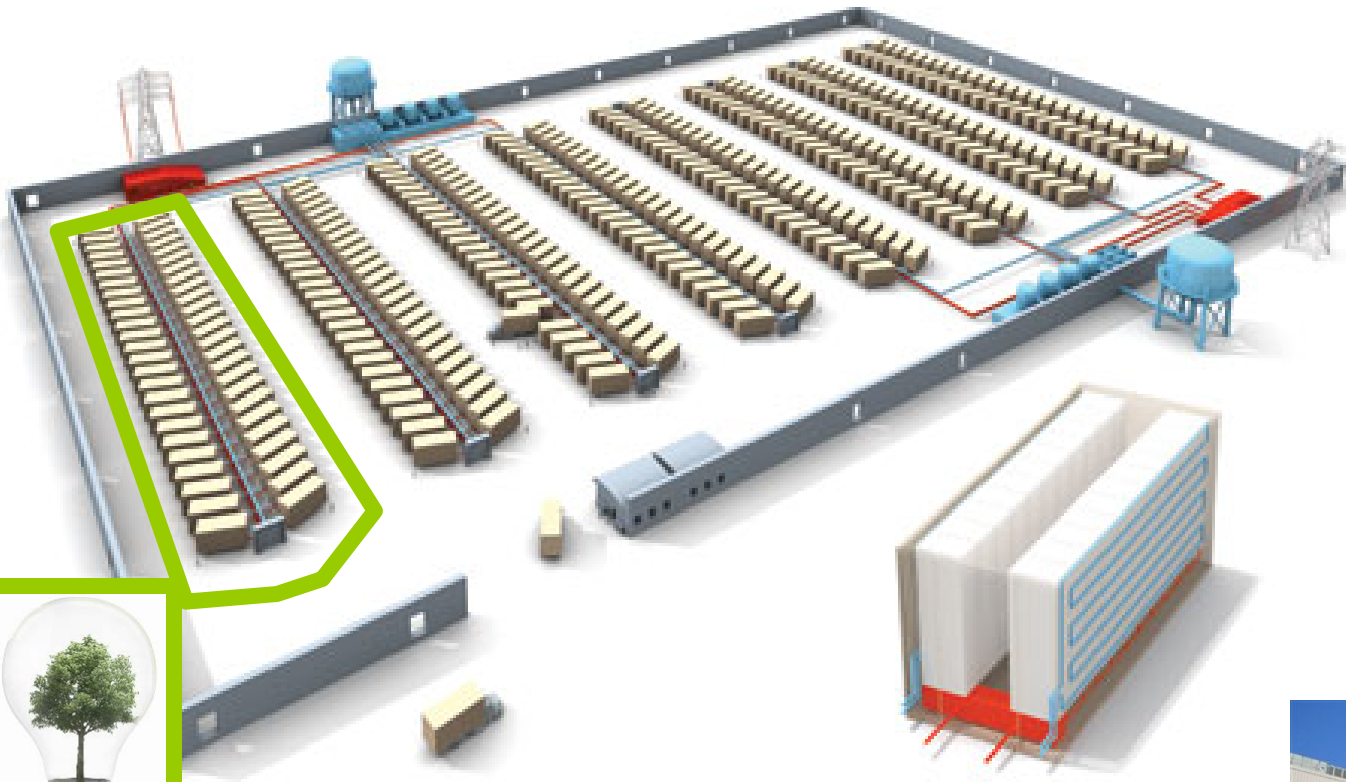
# What is Cloud Computing?

## 3. A Useful IT Service

“Use only when you want! Pay only for what you use!”



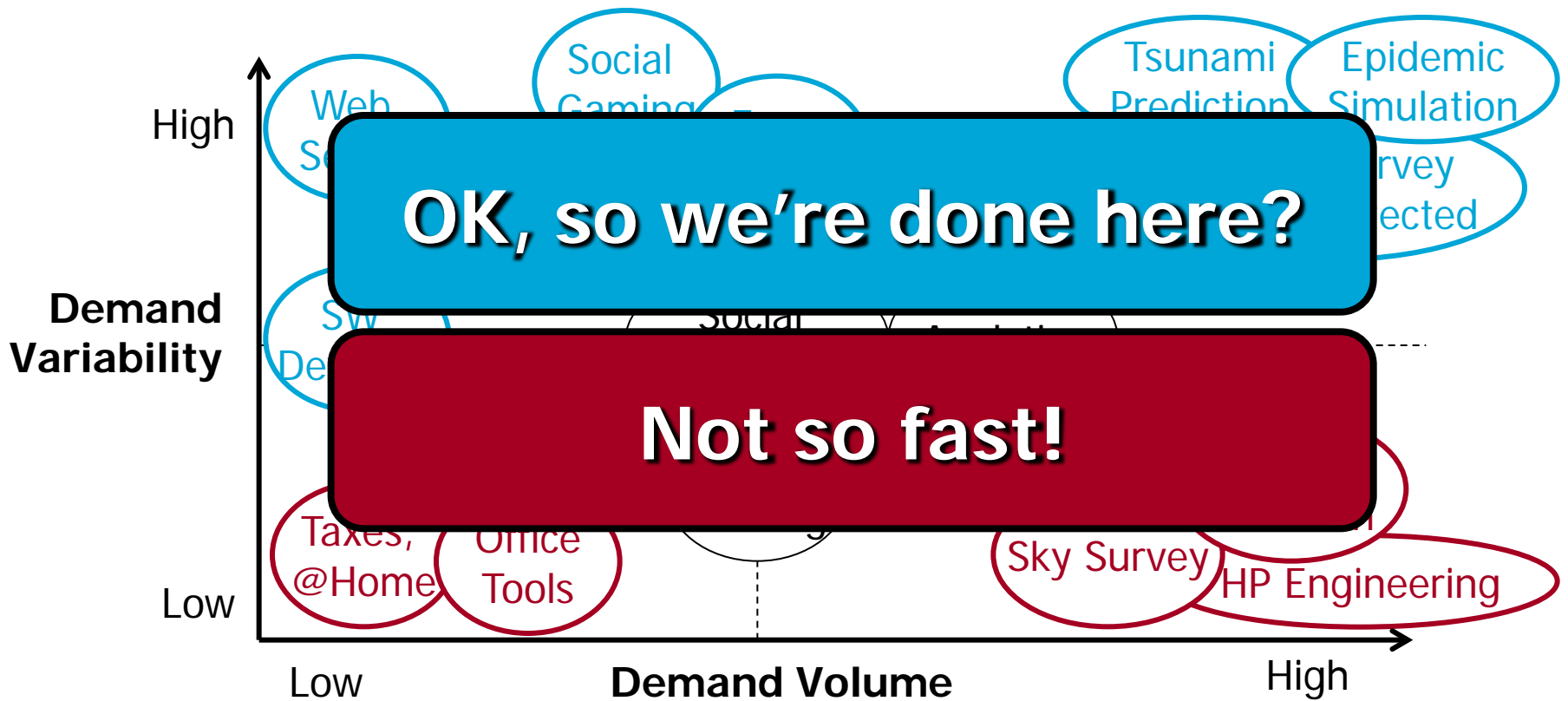
# IaaS Cloud Computing



Many tasks



# Which Applications Need Cloud Computing? A Simplistic View...





# What I Learned From Grids

\* The past

- Average job size is 1 (that is, there are **no [!] tightly-coupled, only conveniently parallel jobs**)

## From Parallel to Many-Task Computing

0 20 40 60 80 100 120 140  
No.CPUs

A. Iosup, C. Dumitrescu, D. H. J. Epema, H. Li, L. Wolters, How are Real Grids Used? The Analysis of Four Grid Traces and Its Implications, Grid 2006.

A. Iosup and D. H. J. Epema, Grid Computing Workloads, IEEE Internet Computing 15(2): 19-26 (2011)

# What I Learned From Grids

\* The past



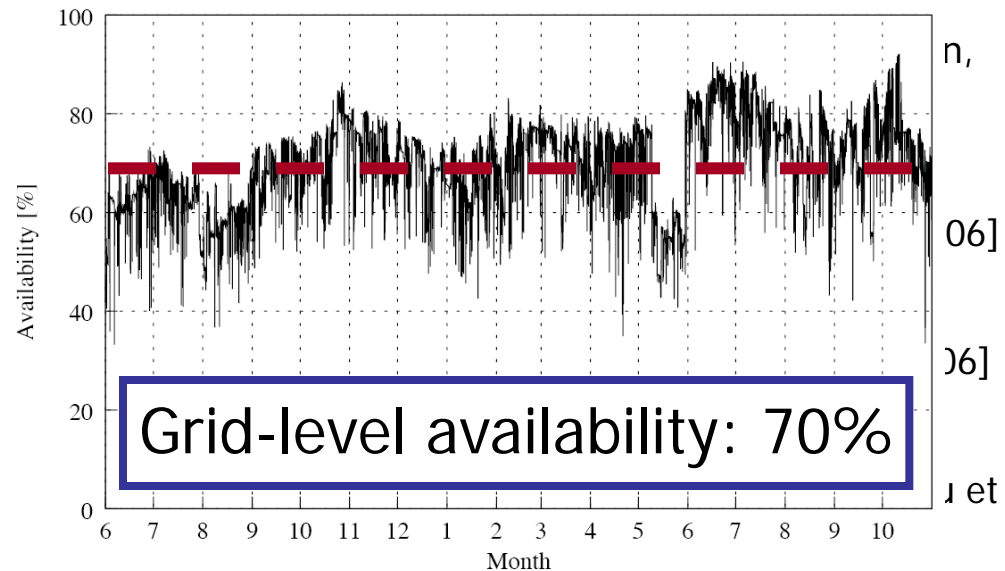
Server

• 99.99999% reliable

## Grids are unreliable infrastructure

BNL-LCG2	0 vs 1 (0.00%)
<b>CERN LCG jobs</b> 74.71% successful <b>25.29% unsuccessful</b>	
INFN-T1	19066 vs 6042 (75.94%)
NIKHEF-ELPROD	5994 vs 22270 (21.21%)
RAL-LCG2	21631 vs 22391 (49.14%)
Taiwan-LCG2	18254 vs 9246 (66.38%)
USCMS-FNAL-WC1	101542 vs 8623 (92.17%)
pic	12851 vs 6627 (65.98%)
<b>TOTAL</b>	<b>495281 vs 167668 (74.71%)</b>

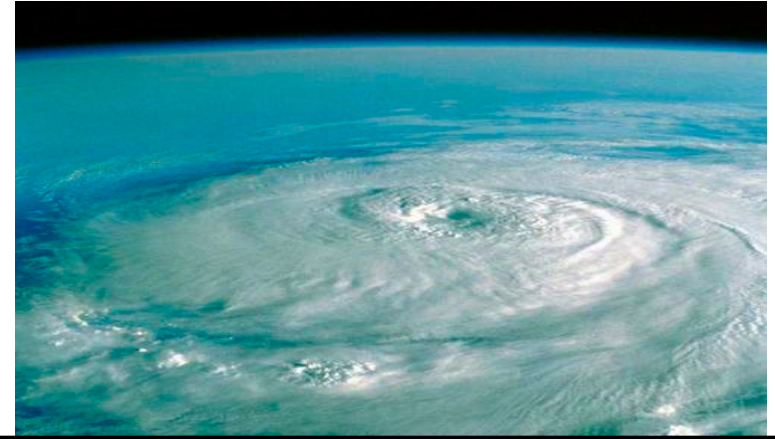
Source: dboard-gr.cern.ch, May'07.



# What I Learned From Grids, Applied to IaaS Clouds



or



**We just don't know!**

- "The path to abundance"
- On-demand capacity
- Cheap for short-term tasks
- Great for web apps (EIP, web crawl, DB ops, I/O)
- "The killer cyclone"
- Performance for scientific applications (compute- or data-intensive)
- Failures, Many-tasks, etc.

# This Presentation: Research Questions

**Q0: What are the workloads of IaaS clouds?**

**Q1: What is the performance of production IaaS cloud services?**

**Q2: How variable is the performance of widely used production cloud services?**

**Q3: How do provisioning and allocation policies affect the performance of IaaS cloud services?**

**Q4: What is the performance of production graph-processing platforms? (ongoing)**

**But ... this is Benchmarking = process of quantifying the performance and other non-functional properties of the system**

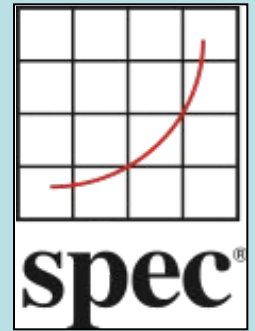
# Why IaaS Cloud Benchmarking?

- Establish and share best-practices in answering important questions about IaaS clouds
- Use in procurement
- Use in system design
- Use in system tuning and operation
- Use in performance management
- **Use in training**

# SPEC Research Group (RG)

*The Research Group of the  
Standard Performance Evaluation Corporation*

\* The present



## Mission Statement

- ▶ Provide a **platform for** collaborative research efforts in the areas of computer benchmarking and quantitative system analysis
- ▶ Provide metrics, tools and benchmarks for evaluating early prototypes and research results as well as full-blown implementations
- ▶ Foster interactions and collaborations btw. industry and academia



# Current Members (Dec 2012)

\* The present



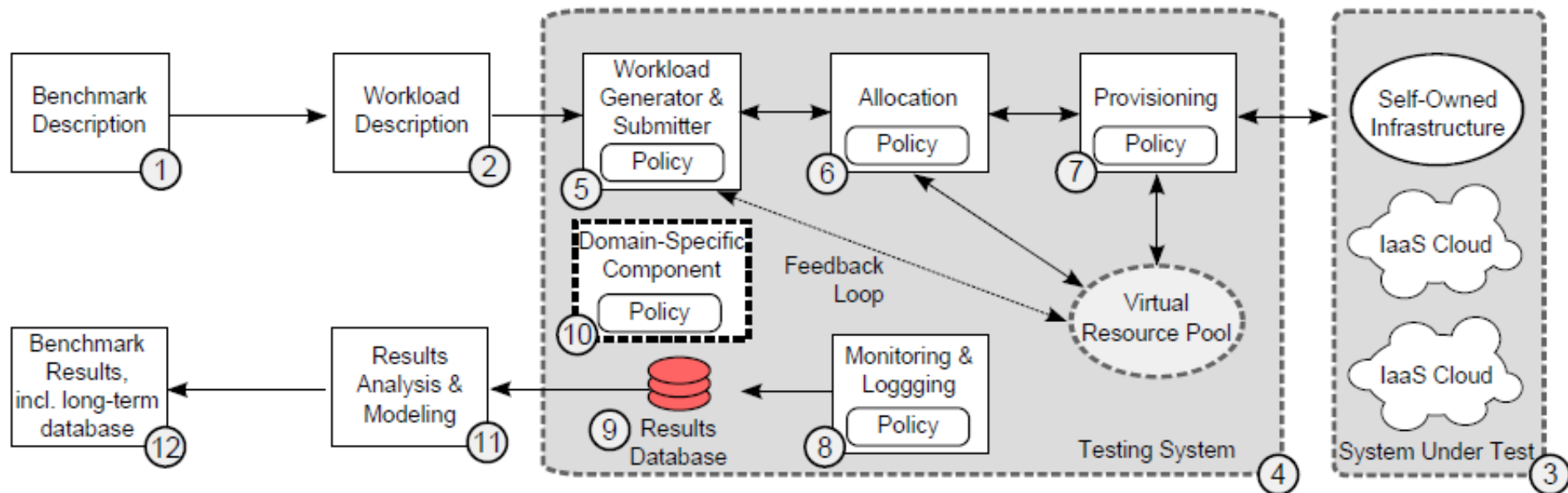
# Agenda

1. An Introduction to IaaS Cloud Computing
2. Research Questions or Why We Need Benchmarking?
- 3. A General Approach and Its Main Challenges**
4. IaaS Cloud Workloads (Q0)
5. IaaS Cloud Performance (Q1) and Perf. Variability (Q2)
6. Provisioning and Allocation Policies for IaaS Clouds (Q3)
7. Big Data: Large-Scale Graph Processing (Q4)
8. Conclusion



# A General Approach for IaaS Cloud Benchmarking

\* The present



# Approach: Real Traces, Models, and Tools + Real-World Experimentation (+ Simulation)

\* The present

- Formalize real-world scenarios
- Exchange real traces
- Model relevant operational elements
- Develop usable tools for meaningful and repeatable experiments
- Conduct comparative studies
  - Simulation only when needed (long-term scenarios, etc.)

**Rule of thumb:  
Put 10-15% project effort  
into benchmarking**

# 10 Main Challenges in 4 Categories\*

\* The future

\* List not exhaustive

## • Methodological

1. Experiment compression
2. Beyond black-box testing through testing short-term dynamics and long-term evolution
3. Impact of middleware

## • System-Related

1. Reliability, availability, and system-related properties
2. Massive-scale, multi-site benchmarking
3. Performance isolation, multi-tenancy models

## • Workload-related

1. Statistical workload models
2. Benchmarking performance isolation under various multi-tenancy workloads

## • Metric-Related

1. Beyond traditional performance: variability, elasticity, etc.
2. Closer integration with cost models

[Read our article](#)

Iosup, Prodan, and Epema, IaaS Cloud Benchmarking: Approaches, Challenges, and Experience, MTAGS 2012. (invited paper)

# Agenda

1. An Introduction to IaaS Cloud Computing
2. Research Questions or Why We Need B
3. A General Approach and Its Main Challenges
4. **IaaS Cloud Workloads (Q0)**
5. **IaaS Cloud Performance (Q1) & Perf. Variability (Q2)**
6. **Provisioning & Allocation Policies for IaaS Clouds (Q3)**
7. **Big Data: Large-Scale Graph Processing (Q4)**
8. [Conclusion](#)



**Workloads**

**Performance**

**Variability**

**Policies**

**Big Data:  
Graphs**

# IaaS Cloud Workloads: Our Team



Alexandru Iosup  
TU Delft

BoTs  
Workflows  
Big Data  
Statistical modeling



Dick Epema  
TU Delft

BoTs  
Grids



Mathieu Jan  
TU Delft/INRIA

BoTs  
Statistical modeling



Ozan Sonmez  
TU Delft

BoTs



Thomas de Ruiter  
TU Delft

MapReduce  
Big Data  
Statistical modeling



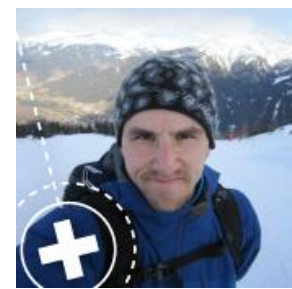
Radu Prodan  
U. Isbk.

Workflows



Thomas Fahringer  
U. Isbk.

Workflows



Simon Ostermann  
U. Isbk.

Workflows

# What I'll Talk About

## IaaS Cloud Workloads (Q0)

1. BoTs
2. Workflows
3. Big Data Programming Models
4. MapReduce workloads

# What is a Bag of Tasks (BoT)? A System View

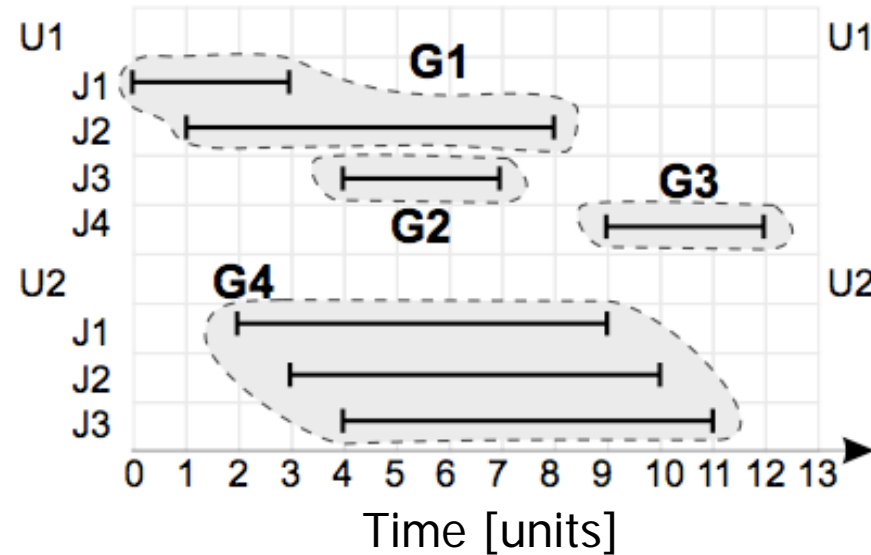
BoT = set of jobs sent by a user...

$$W_u = \{J_i | user(J_i) = u\}$$

...that is submitted at most  $\Delta$ s after the first job

$$ST(J') \leq ST(J) + \Delta$$

- Why **Bag of Tasks**? From the perspective of the user, jobs in set are just **tasks of a larger job**
- A single useful result from the complete BoT
- Result can be combination of all tasks, or a selection of the results of most or even a single task



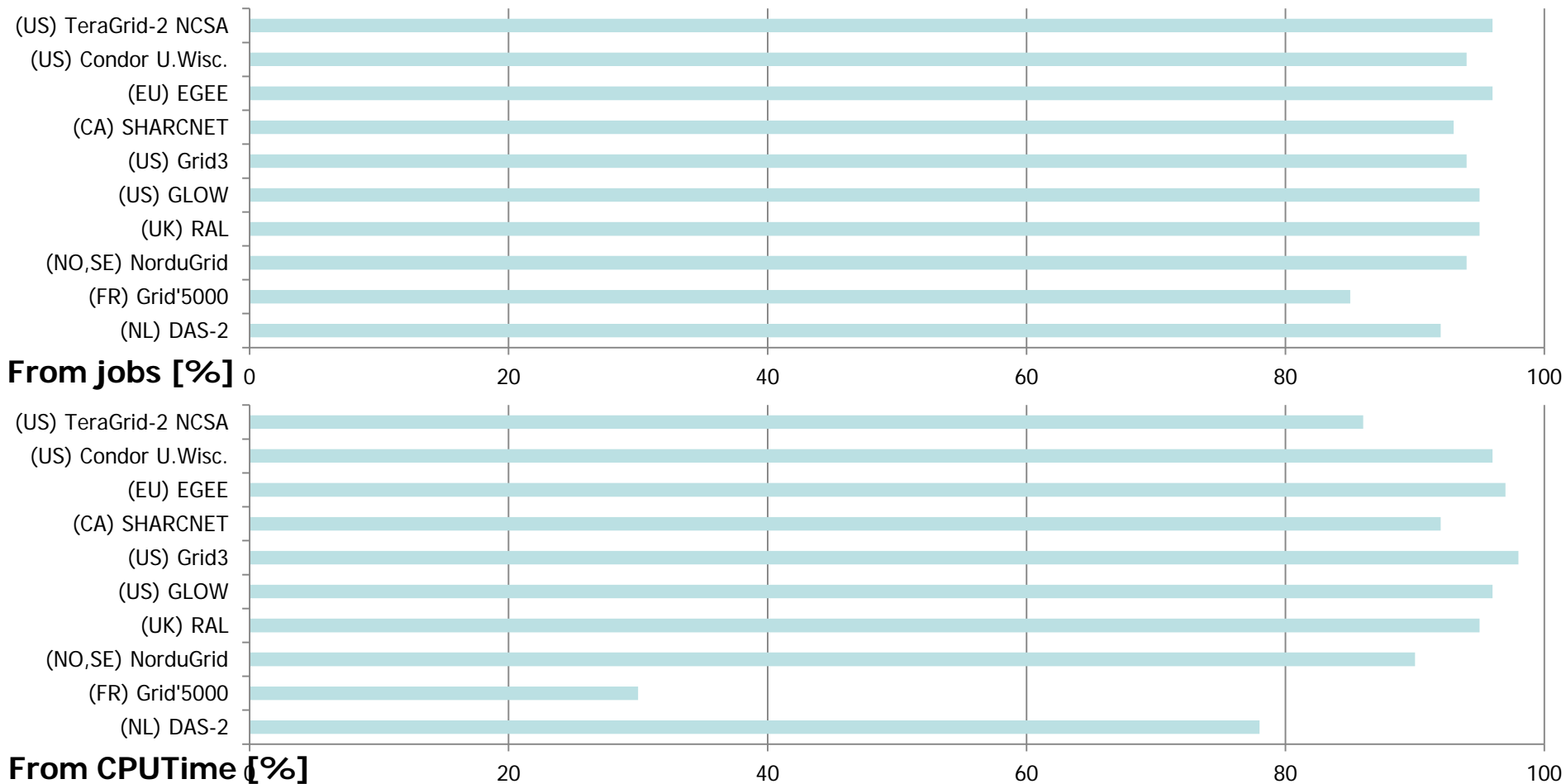
Iosup et al., The Characteristics and Performance of Groups of Jobs in Grids, Euro-Par, LNCS, vol. 4641, pp. 382-393, 2007.

# Applications of the BoT Programming Model

- Parameter sweeps
  - Comprehensive, possibly exhaustive investigation of a model
  - Very useful in engineering and simulation-based science
- Monte Carlo simulations
  - Simulation with random elements: fixed time yet limited inaccuracy
  - Very useful in engineering and simulation-based science
- Many other types of batch processing
  - Periodic computation, Cycle scavenging
  - Very useful to automate operations and reduce waste

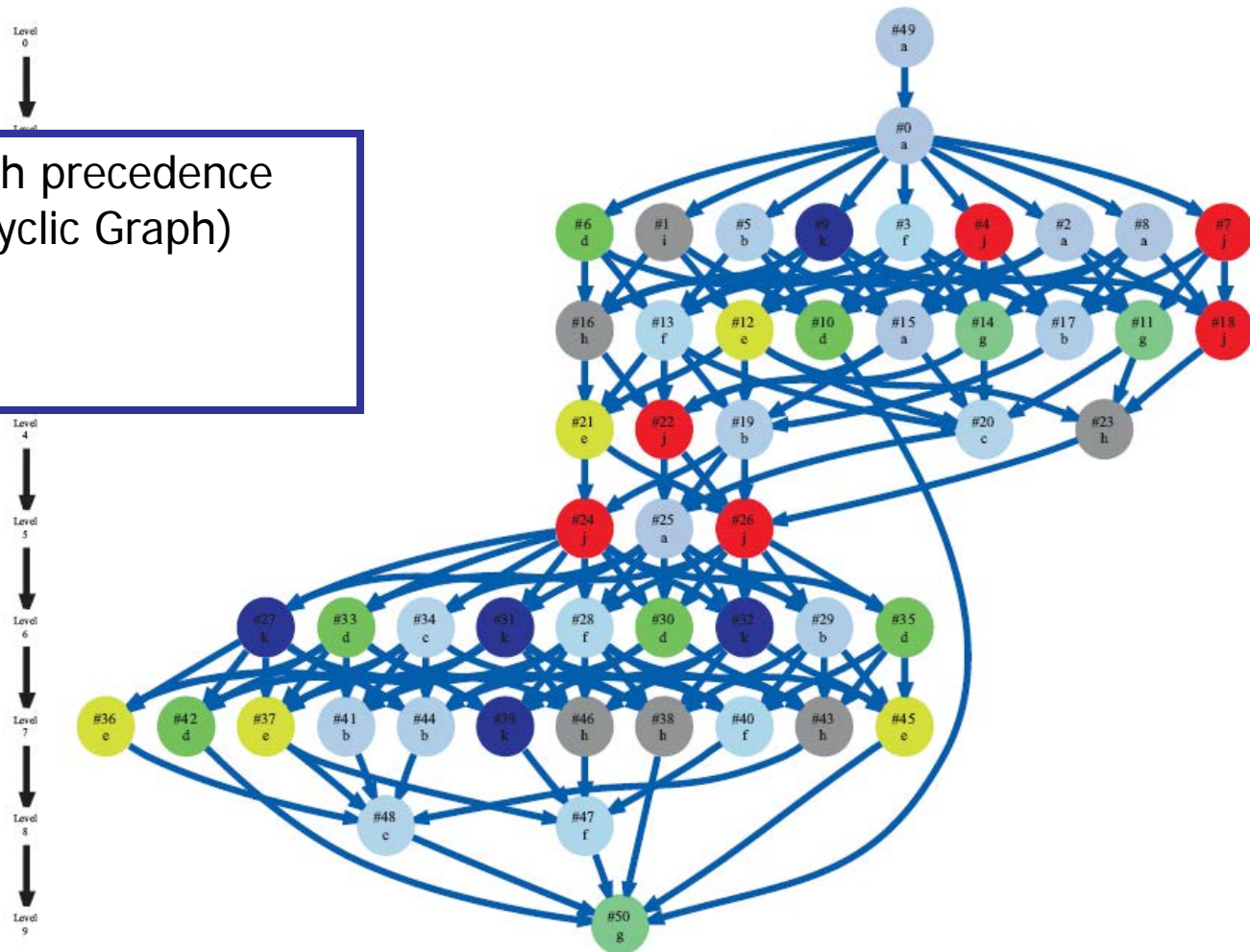


# BoTs Are the Dominant Programming Model for Grid Computing (Many Tasks)



# What is a Workflow?

WF = set of jobs with precedence  
(think Direct Acyclic Graph)



# Applications of the Workflow Programming Model

- Complex applications
  - Complex filtering of data
  - Complex analysis of instrument measurements
- Applications created by non-CS scientists\*
  - Workflows have a natural correspondence in the real-world, as descriptions of a scientific procedure
  - Visual model of a graph sometimes easier to program
- Precursor of the MapReduce Programming Model (next slides)



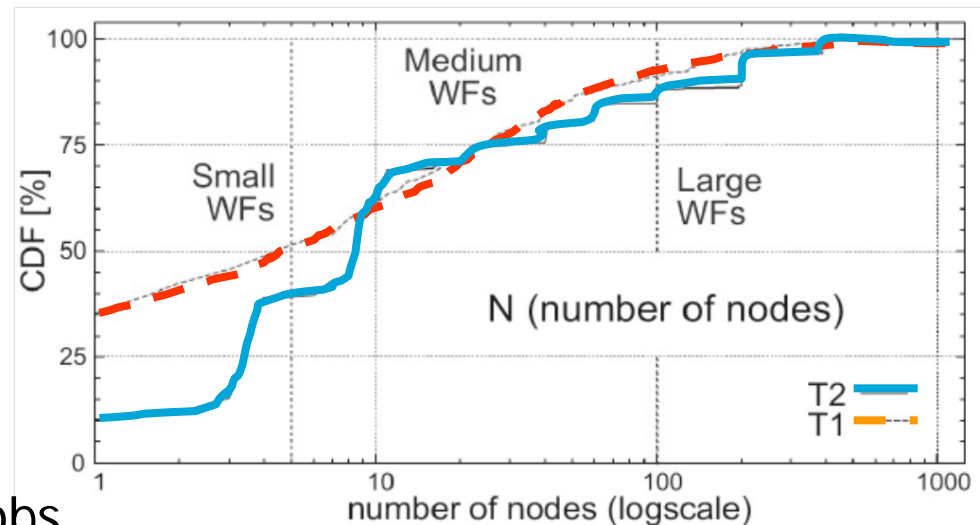
# Workflows Exist in Grids, but Did No Evidence of a Dominant Programming Model

- Traces

Trace	Source	Duration	Number of WFs	Number of Tasks	CPUdays
T1	DEE	09/06-10/07	4,113	122k	152
T2	EE2	05/07-11/07	1,030	46k	41

- Selected Findings

- Loose coupling
- Graph with 3-4 levels
- Average WF size is 30/44 jobs
- 75%+ WFs are sized 40 jobs or less, 95% are sized 200 jobs or less



Ostermann et al., On the Characteristics of Grid Workflows, CoreGRID Integrated Research in Grid Computing (CGIW), 2008.

# What is “Big Data”?

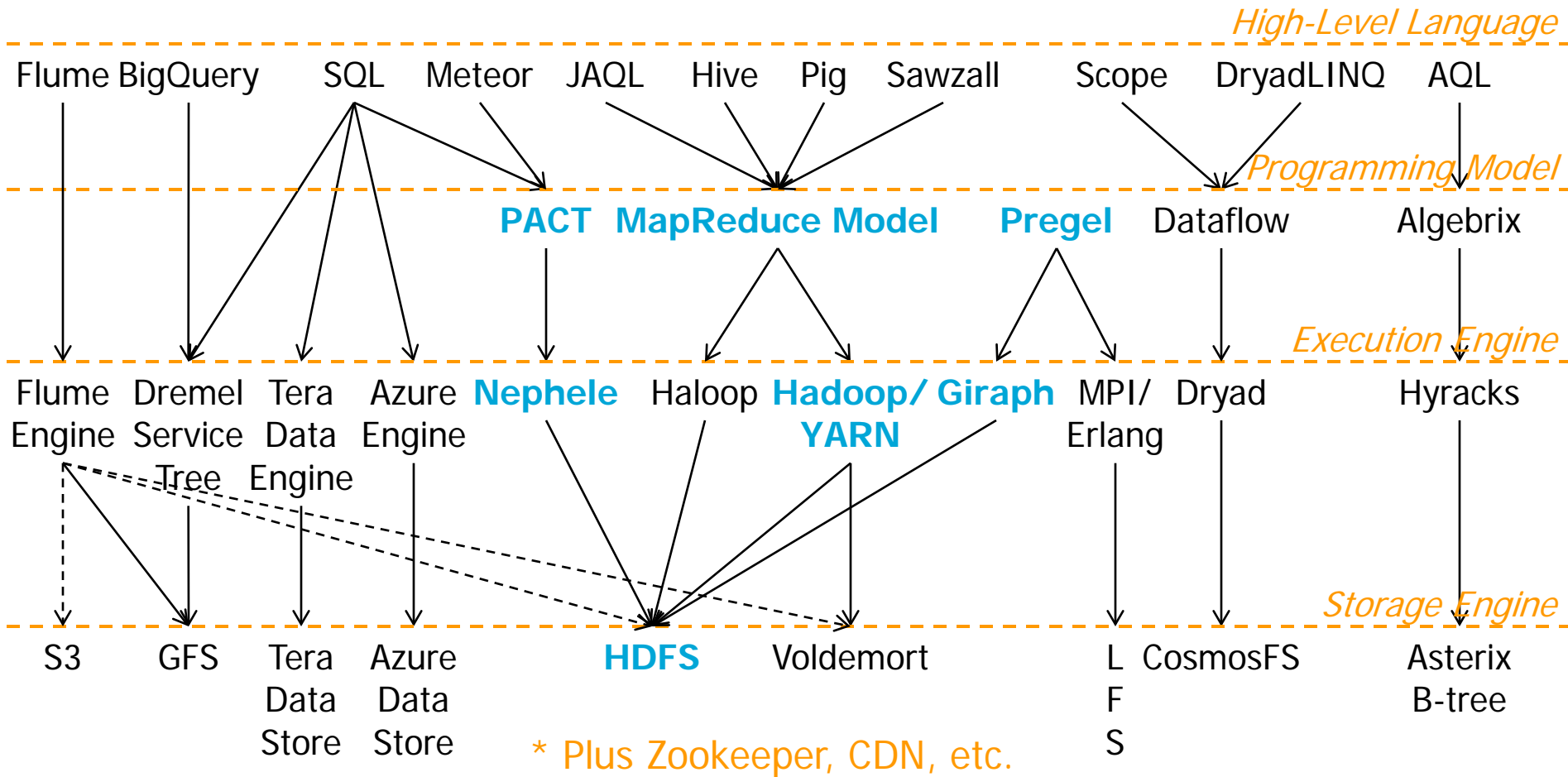
- Very large, distributed aggregations of loosely structured data, often incomplete and inaccessible
- Easily exceeds the processing capacity of conventional database systems
- Principle of Big Data: *“When you can, keep everything!”*
- **Too** big, **too** fast, and **doesn't** comply with the traditional database architectures

# The Three “V”s of Big Data

- Volume
  - More data vs. better models
  - Data grows exponentially
  - Analysis in near-real time to extract value
  - Scalable storage and distributed queries
- Velocity
  - Speed of the feedback loop
  - Gain competitive advantage: fast recommendations
  - Identify fraud, predict customer churn faster
- Variety
  - The data can become messy: text, video, audio, etc.
  - Difficult to integrate into applications

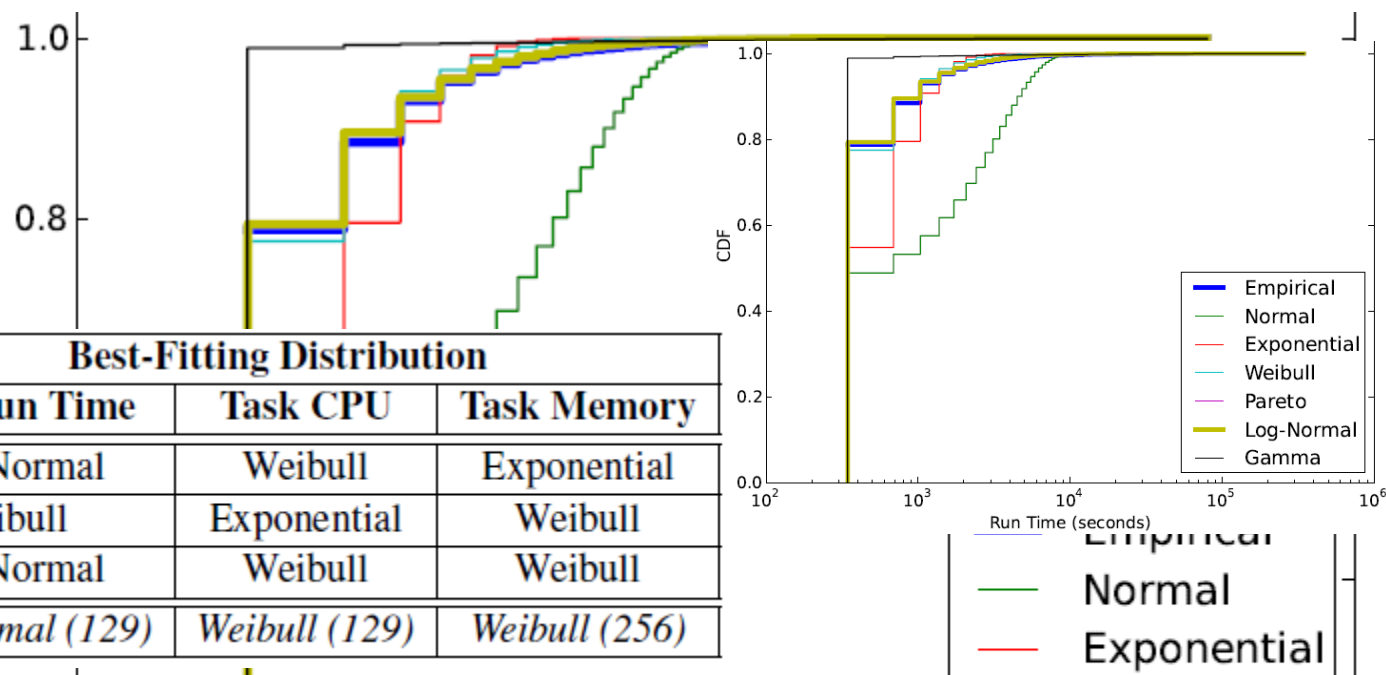


# Ecosystems of Big-Data Programming Models



# Our Statistical MapReduce Models

- Real traces
  - Yahoo
  - Google
  - 2 x Social N



Model	Tasks	Correlation	Map/Reduce Modeled	Sign. Level	Indirect Distr. Sel.
Complex Model	Indirect	Run time – Disk	Separately	0.05	Best fits
Relaxed Complex Model	Indirect	Run time – Disk	Separately	0.02	All fits
Safe Complex Model	Direct	Run time – Disk	Separately	0.05	–
Simple Model	Direct	–	Together	0.05	–



# Agenda

1. An Introduction to IaaS Cloud Comput
2. Research Questions or Why We Need B
3. A General Approach and Its Main Chall
4. **IaaS Cloud Workloads (Q0)**
5. **IaaS Cloud Performance (Q1) & Perf. Variability (Q2)**
6. **Provisioning & Allocation Policies for IaaS Clouds (Q3)**
7. **Big Data: Large-Scale Graph Processing (Q4)**
8. [Conclusion](#)



**Workloads**

**Performance**

**Variability**

**Policies**

**Big Data:  
Graphs**

# IaaS Cloud Performance: Our Team



Alexandru Iosup  
TU Delft

Performance  
Variability  
Isolation  
Multi-tenancy  
Benchmarking



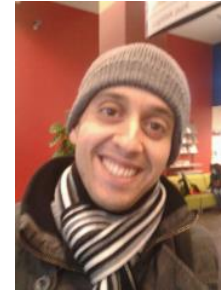
Dick Epema  
TU Delft

Performance  
IaaS clouds



Nezhir Yigitbasi  
TU Delft

Performance  
Variability



Athanasios Antoniou  
TU Delft

Performance  
Isolation



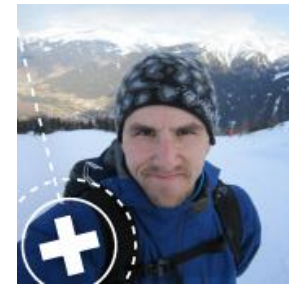
Radu Prodan  
U. Isbk.

Benchmarking



Thomas Fahringer  
U. Isbk.

Benchmarking



Simon Ostermann  
U. Isbk.

Benchmarking

# What I'll Talk About

## IaaS Cloud Performance (Q1)

1. Previous work
2. Experimental setup
3. Experimental results
4. Implications on real-world workloads

# Some Previous Work

## (>50 important references across our studies)

### Virtualization Overhead

- Loss below 5% for computation [Barham03] [Clark04]
- Loss below 15% for networking [Barham03] [Menon05]
- Loss below 30% for parallel I/O [Vetter08]
- Negligible for compute-intensive HPC kernels [You06] [Panda06]

### Cloud Performance Evaluation

- Performance and cost of executing a sci. workflows [Dee08]
- Study of Amazon S3 [Palankar08]
- Amazon EC2 for the NPB benchmark suite [Walker08] or selected HPC benchmarks [Hill08]
- CloudCmp [Li10]
- Kosmann et al.

# Production IaaS Cloud Services

- **Production IaaS cloud:** lease resources (infrastructure) to users, operate on the market and have active customers

Name	Cores (ECUs)	RAM [GB]	Archi. [bit]	Disk [GB]	Cost [\$/h]
<i>Amazon EC2</i>					
m1.small	1 (1)	1.7	32	160	0.1
m1.large	2 (4)	7.5	64	850	0.4
m1.xlarge	4 (8)	15.0	64	1,690	0.8
c1.medium	2 (5)	1.7	32	350	0.2
c1.xlarge	8 (20)	7.0	64	1,690	0.8
<i>GoGrid (GG)</i>					
GG.small	1	1.0	32	60	0.19
GG.large	1	1.0	64	60	0.19
GG.xlarge	3	4.0	64	240	0.76
<i>Elastic Hosts (EH)</i>					
EH.small	1	1.0	32	30	£0.042
EH.large	1	4.0	64	30	£0.09
<i>Mosso</i>					
Mosso.small	4	1.0	64	40	0.06
Mosso.large	4	4.0	64	160	0.24

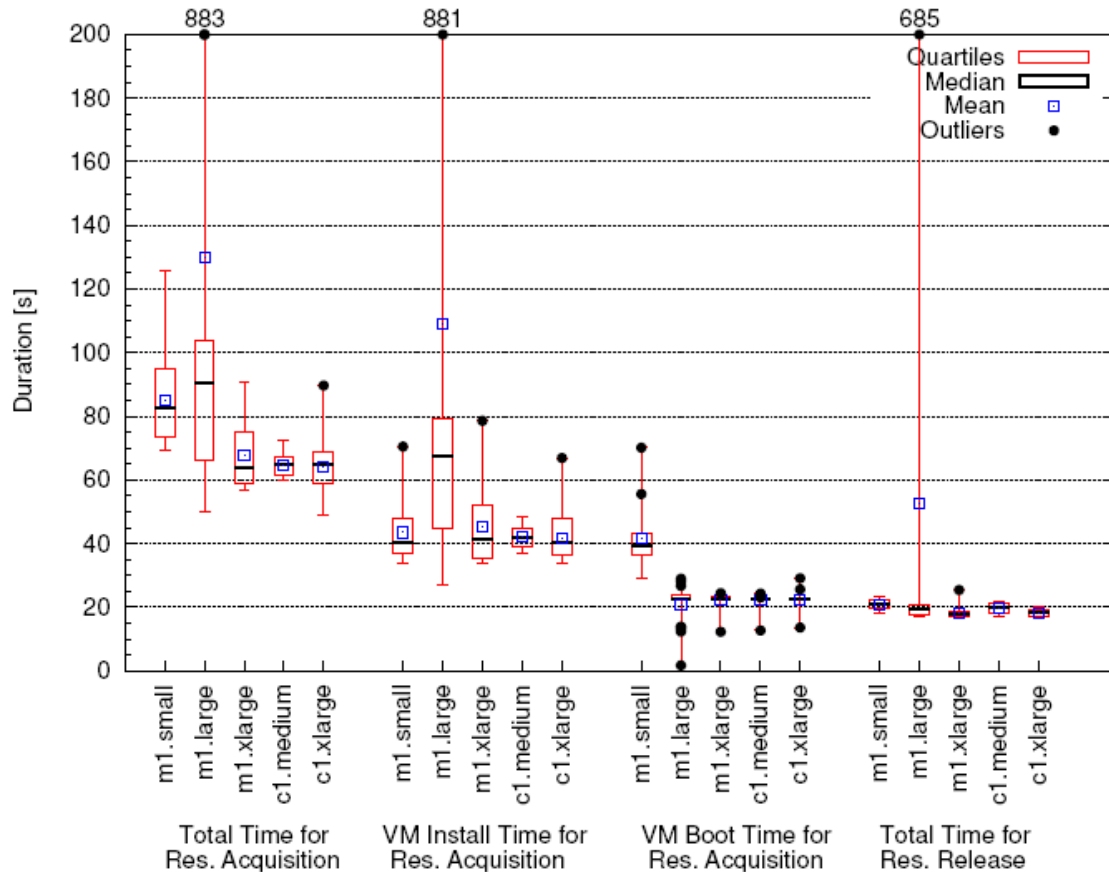
# Our Method

- Based on general performance technique: model performance of individual components; system performance is performance of workload + model [Saavedra and Smith, ACM TOCS'96]
- Adapt to clouds:
  1. Cloud-specific elements: resource provisioning and allocation
  2. Benchmarks for single- and multi-machine jobs
  3. Benchmark CPU, memory, I/O, etc.:

Type	Suite/Benchmark	Resource	Unit
SI	LMbench/all [24]	Many	Many
SI	Bonnie/all [25], [26]	Disk	MBps
SI	CacheBench/all [27]	Memory	MBps
MI	HPCC/HPL [28], [29]	CPU	GFLOPS
MI	HPCC/DGEMM [30]	CPU	GFLOPS
MI	HPCC/STREAM [30]	Memory	GBps
MI	HPCC/RandomAccess [31]	Network	MUPS
MI	HPCC/ $b_{eff}$ (lat,bw.) [32]	Comm.	$\mu s$ , GBps

# Single Resource Provisioning/Release

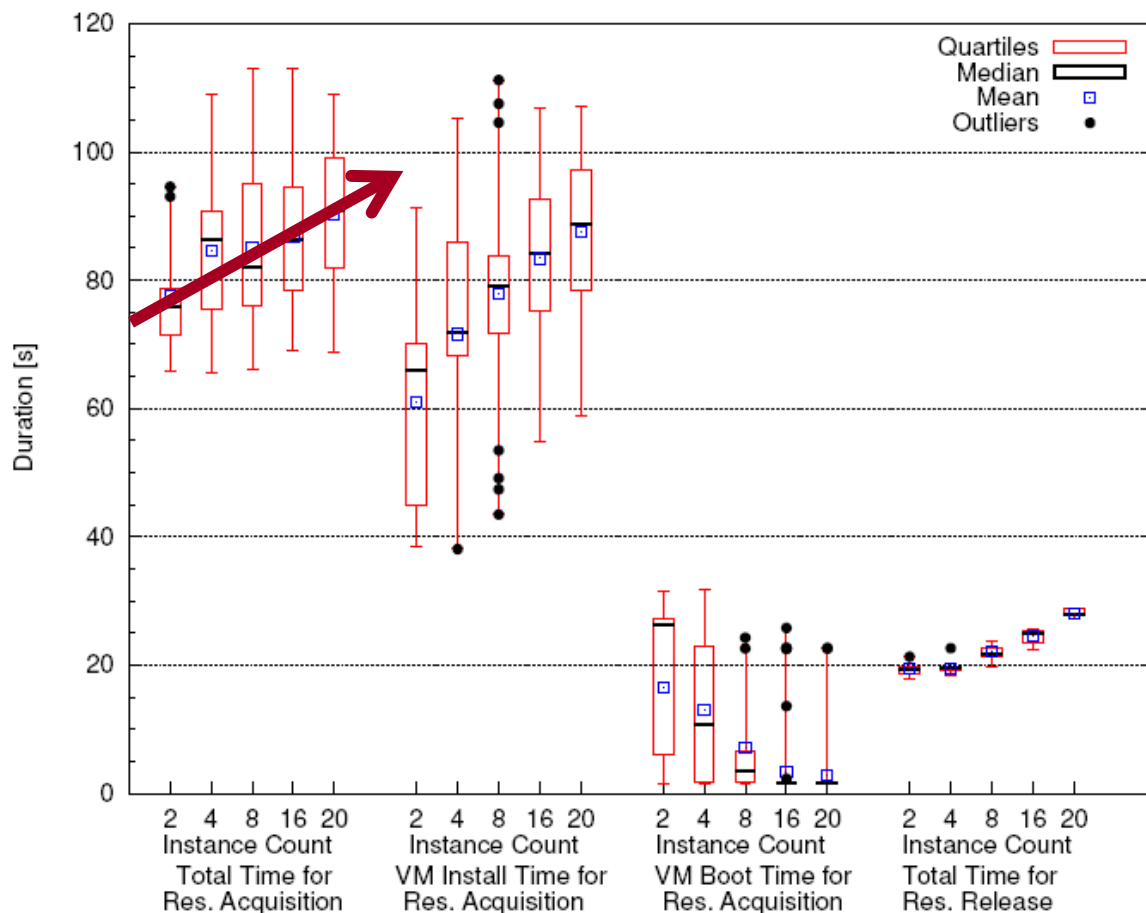
Q1



- Time depends on instance type
- Boot time non-negligible

# Multi-Resource Provisioning/Release

Q1

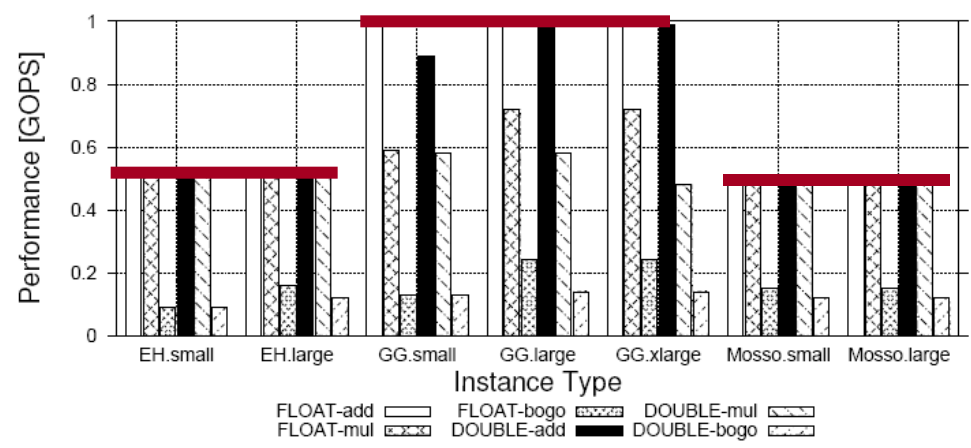
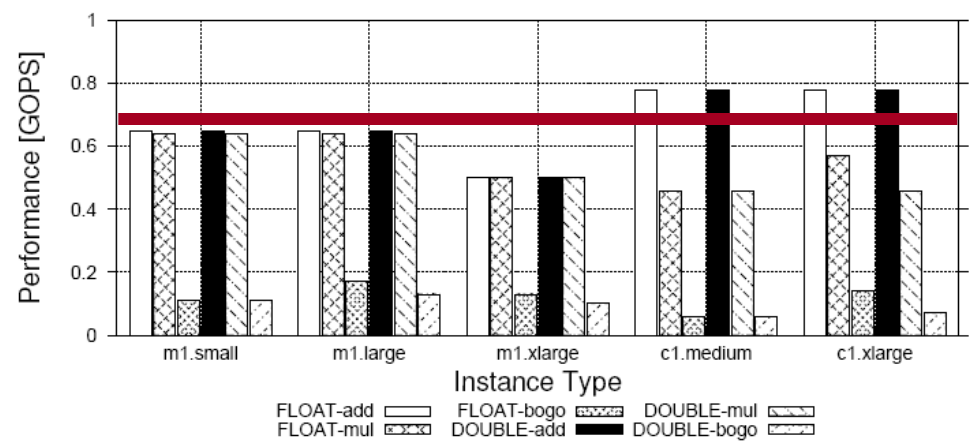


- Time for *multi*-resource increases with number of resources

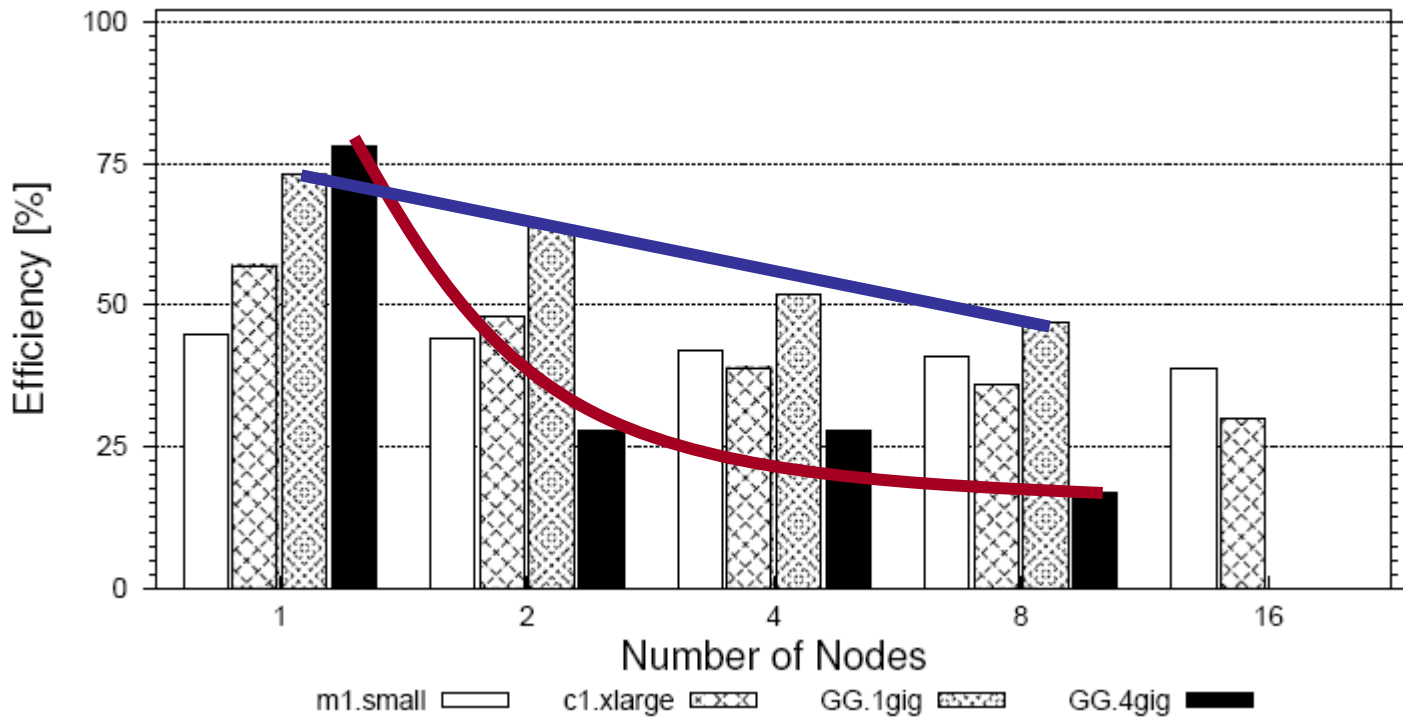


# CPU Performance of Single Resource

- ECU definition: "a 1.1 GHz 2007 Opteron" ~ 4 flops per cycle at full pipeline, which means at peak performance one ECU equals 4.4 gigaflops per second (GFLOPS)
- Real performance 0.6..0.1 GFLOPS = ~1/4..1/7 theoretical peak



# HPLinpack Performance (Parallel)

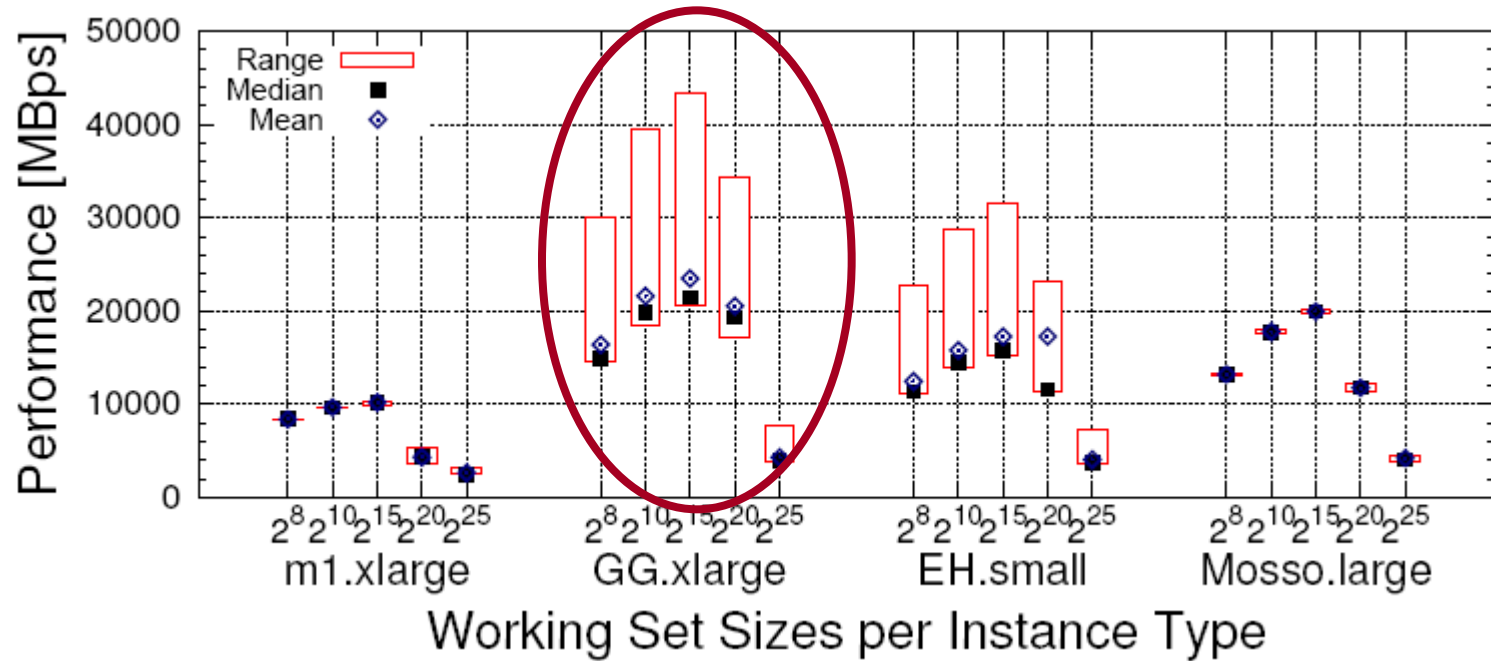


- Low efficiency for parallel compute-intensive applications
- Low performance vs cluster computing and supercomputing

# Performance Stability (Variability)

Q1

Q2



- High performance **variability** for the best-performing instances

# Summary

- Much lower performance than theoretical peak
  - Especially CPU (GFLOPS)
- Performance variability
- Compared results with some of the commercial alternatives (see report)

# Implications: Simulations

- Input: real-world workload traces, grids and PPEs

- Running in

- Original env.
- Cloud with source-like perf.
- Cloud with measured perf.

Trace ID, Source (Trace ID in Archive)	Time [mo.]	Trace		System		Load [%]
		Number of Jobs	Users	Sites	Size CPUs	
<i>Grid Workloads Archive [13], 6 traces</i>						
1. DAS-2 (1)	18	1.1M	333	5	0.4K	15+
2. RAL (6)	12	0.2M	208	1	0.8K	85+
3. GLOW (7)	3	0.2M	18	1	1.6K	60+
4. Grid3 (8)	18	1.3M	19	29	3.5K	-
5. SharcNet (10)	13	1.1M	412	10	6.8K	-
6. LCG (11)	1	0.2M	216	200+	24.4K	-
<i>Parallel Workloads Archive [16], 4 traces</i>						
7. CTC SP2 (6)	11	0.1M	679	1	0.4K	66
8. SDSC SP2 (9)	24	0.1M	437	1	0.1K	83
9. LANLO2K (10)	5	0.1M	337	1	2.0K	64
10. SDSC DS (19)	13	0.1M	460	1	1.7K	63

- Metrics

- WT, ReT, BSD(10s)
- Cost [CPU-h]

# Implications: Results

Trace ID	Source env. (Grid/PPI)			Cloud (real performance)			Cloud (source performance)		
	AWT [s]	AReT [s]	ABSD (10s)	AReT [s]	ABSD (10s)	Total Cost [CPU-h,M]	AReT [s]	ABSD (10s)	Total Cost [CPU-h,M]
DAS-2	432	802	11	2,292	2.39	2	450	2	1.19
RAL	13,214	27,807	68	131,300	1	40	18,837	1	6.39
GLOW	9,162	17,643	55	59,448	1	3	8,561	1	0.60
Grid3	-	7,199	-	50,470	3	19	7,279	3	3.60
SharcNet	31,017	61,682	242	219,212	1	73	31,711	1	11.34
LCG	-	9,011	-	63,158	1	3	9,091	1	0.62
CTC SP2	25,748	37,019	78	75,706	1	2	11,351	1	0.30
SDSC SP2	26,705	33,388	389	46,818	2	1	6,763	2	0.16
LANL O2K	4,658	9,594	61	37,786	2	1	5,016	2	0.26
SDSC DS	32,271	33,807	516	57,065	2	2	6,790	2	0.25

- Cost: Clouds, real >> Clouds, source



- Performance:

- AReT: Clouds, real >> Source env. (bad)



- AWT,ABSD: Clouds, real << Source env. (good)



# Agenda

1. An Introduction to IaaS Cloud Computing
2. Research Questions or Why We Need B
3. A General Approach and Its Main Challenges
4. **IaaS Cloud Workloads (Q0)**
5. **IaaS Cloud Performance (Q1) & Perf. Variability (Q2)**
6. **Provisioning & Allocation Policies for IaaS Clouds (Q3)**
7. **Big Data: Large-Scale Graph Processing (Q4)**
8. [Conclusion](#)



**Workloads**

**Performance**

**Variability**

**Policies**

**Big Data:  
Graphs**

# IaaS Cloud Performance: Our Team



Alexandru Iosup  
TU Delft

Performance  
Variability  
Isolation  
Multi-tenancy  
Benchmarking



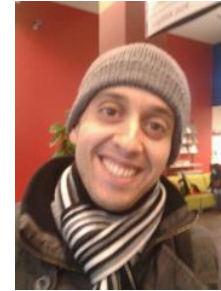
Dick Epema  
TU Delft

Performance  
IaaS clouds



Nezhir Yigitbasi  
TU Delft

Performance  
Variability



Athanasios Antoniou  
TU Delft

Performance  
Isolation



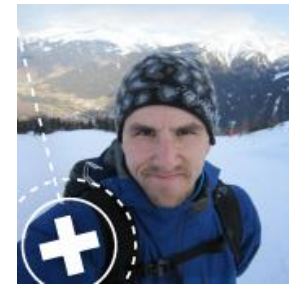
Radu Prodan  
U. Isbk.

Benchmarking



Thomas Fahringer  
U. Isbk.

Benchmarking



Simon Ostermann  
U. Isbk.

Benchmarking



# What I'll Talk About

## IaaS Cloud Performance Variability (Q2)

1. Experimental setup
2. Experimental results
3. Implications on real-world workloads

# Production Cloud Services

- **Production cloud:** operate on the market and have active customers
- **IaaS/PaaS:**
  - **Amazon Web Services (AWS)**
    - EC2 (Elastic Compute Cloud)
    - S3 (Simple Storage Service)
    - SQS (Simple Queueing Service)
    - SDB (Simple Database)
    - FPS (Flexible Payment Service)
  - **PaaS:**
    - **Google App Engine (GAE)**
      - Run (Python/Java runtime)
      - Datastore (Database) ~ SDB
      - Memcache (Caching)
      - URL Fetch (Web crawling)

# Our Method

## Performance Traces

[1/3]

- CloudStatus\*
  - Real-time values and weekly averages for most of the AWS and GAE services
- Periodic performance probes
  - Sampling rate is under 2 minutes

\* [www.cloudstatus.com](http://www.cloudstatus.com)

# Our Method Analysis

[2/3]

Q2

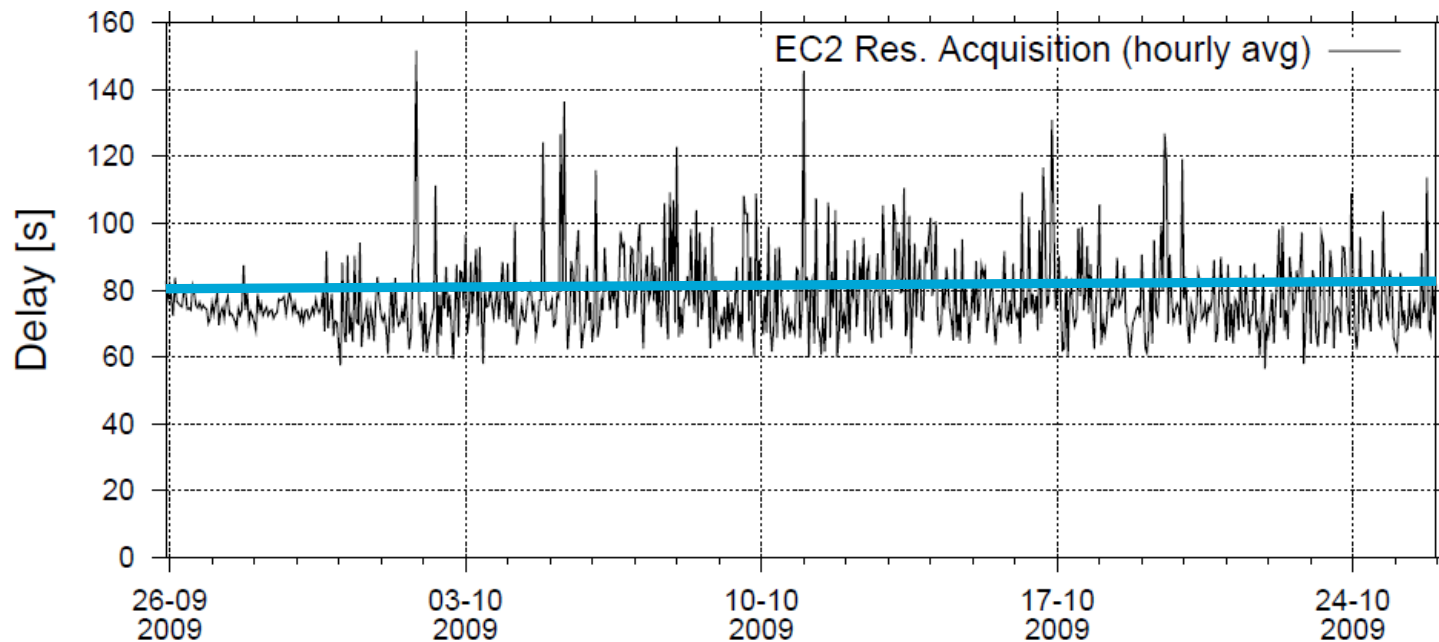
1. Find out whether variability is present
  - Investigate several months whether the performance metric is highly variable
2. Find out the characteristics of variability
  - Basic statistics: the five quartiles ( $Q_0$ - $Q_4$ ) including the median ( $Q_2$ ), the mean, the standard deviation
  - Derivative statistic: the IQR ( $Q_3$ - $Q_1$ )
  - $CoV > 1.1$  indicate high variability
3. Analyze the performance variability time patterns
  - Investigate for each performance metric the presence of daily/monthly/weekly/yearly time patterns
  - E.g., for monthly patterns divide the dataset into twelve subsets and for each subset compute the statistics and plot for visual inspection

# Our Method

## Is Variability Present?

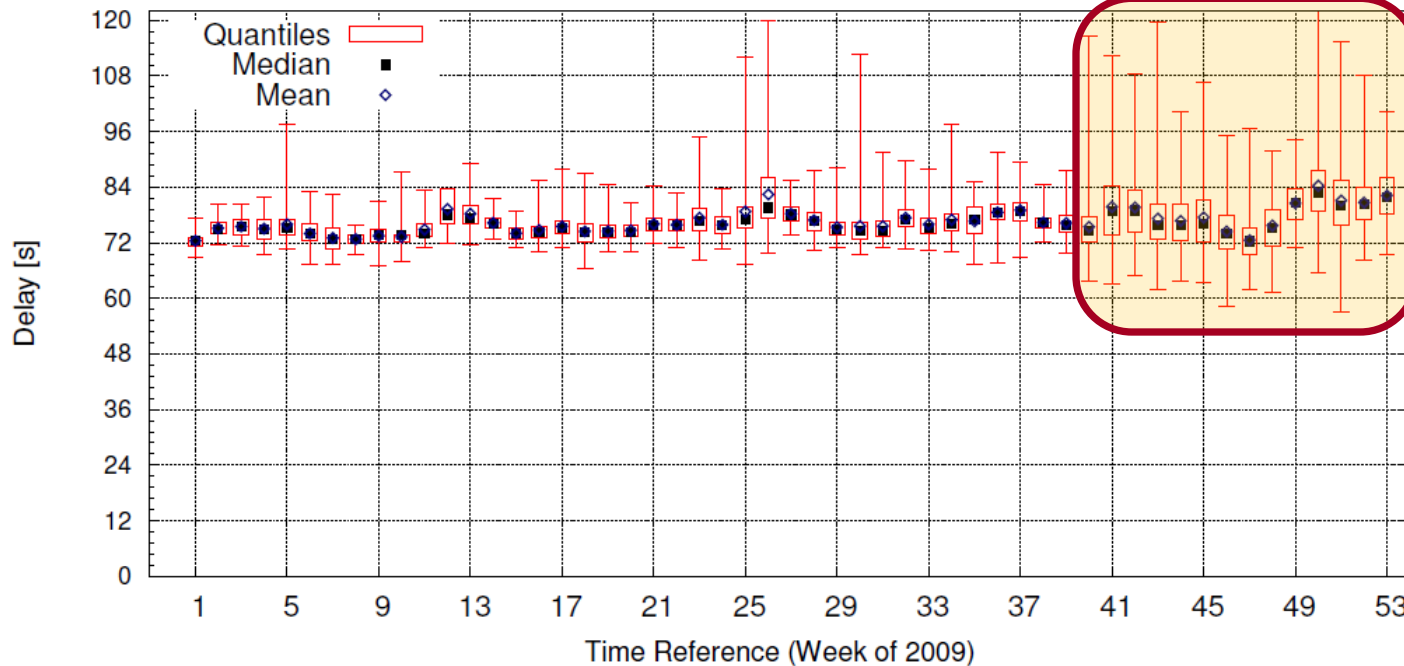
[3/3]

- **Validated Assumption:** The performance delivered by production services is variable.



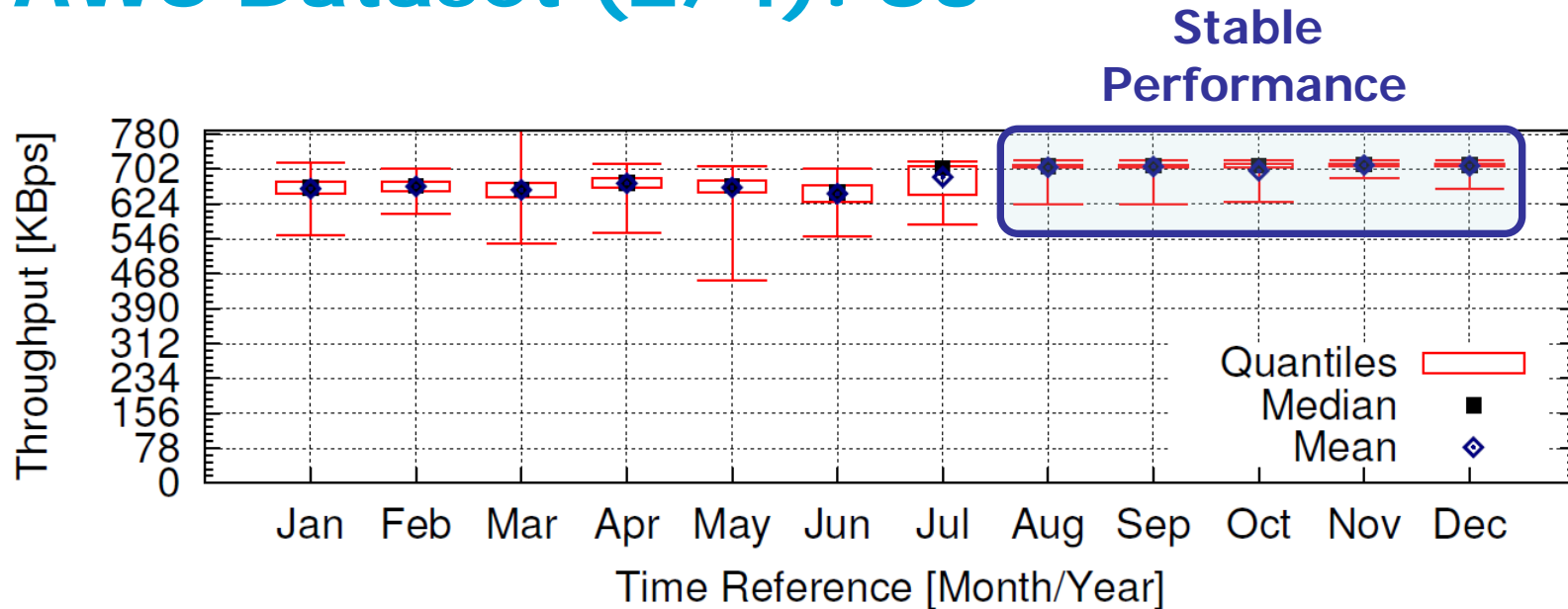
# AWS Dataset (1/4): EC2

## Variable Performance



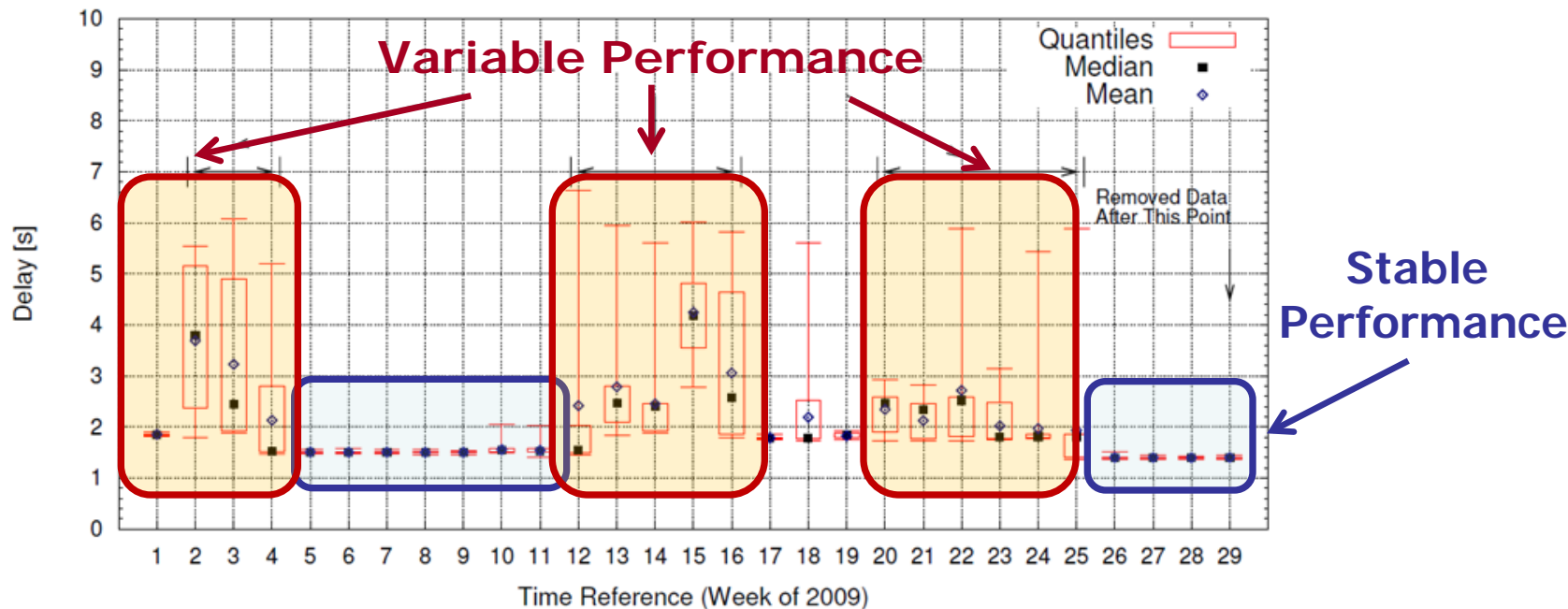
- **Deployment Latency [s]:** Time it takes to start a small instance, from the startup to the time the instance is available
- Higher IQR and range from week 41 to the end of the year; possible reasons:
  - Increasing EC2 user base
  - Impact on applications using EC2 for auto-scaling

# AWS Dataset (2/4): S3



- **Get Throughput [bytes/s]:** Estimated rate at which an object in a bucket is read
- The last five months of the year exhibit much lower IQR and range
  - More stable performance for the last five months
  - Probably due to software/infrastructure upgrades

# AWS Dataset (3/4): SQS



- **Average Lag Time [s]:** Time it takes for a posted message to become available to read. Average over multiple queues.
- Long periods of stability (low IQR and range)
- Periods of high performance variability also exist



# AWS Dataset (4/4): Summary

- **All services exhibit time patterns in performance**
- EC2: periods of special behavior
- SDB and S3: daily, monthly and yearly patterns
- SQS and FPS: periods of special behavior

# Experimental Setup (1/2): Simulations

Q2

- Trace based simulations for three applications
- **Input**
  - GWA traces
  - Number of daily unique users
  - Monthly performance variability

Application	Service
Job Execution	GAE Run
Selling Virtual Goods	AWS FPS
Game Status Maintenance	AWS SDB/GAE Datastore

# Experimental Setup (2/2): Metrics

- Average Response Time and Average Bounded Slowdown
- Cost in millions of consumed CPU hours
- **Aggregate Performance Penalty** -- APP(t)

$$\frac{P(t)}{P_{ref}} \times \frac{U(t)}{\max U(t)}$$

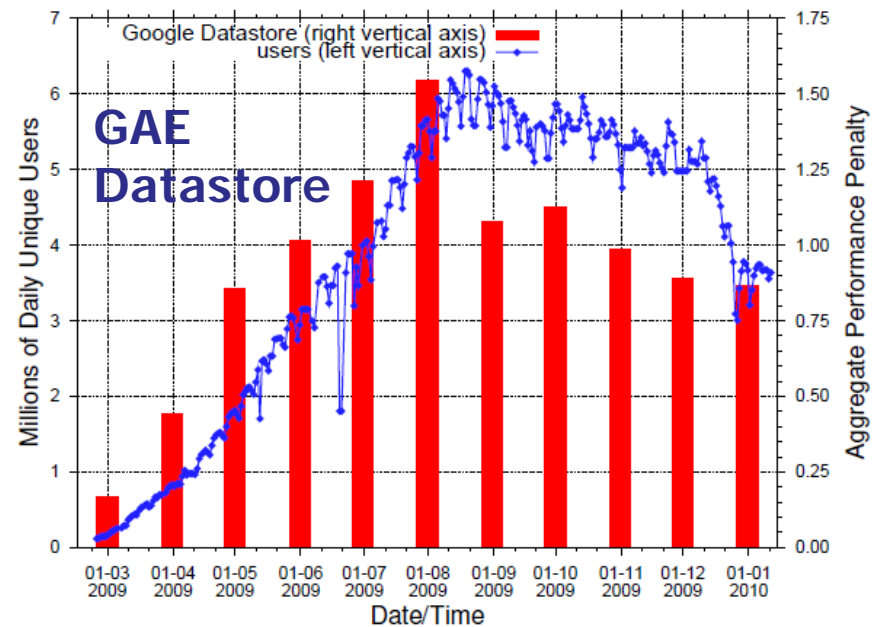
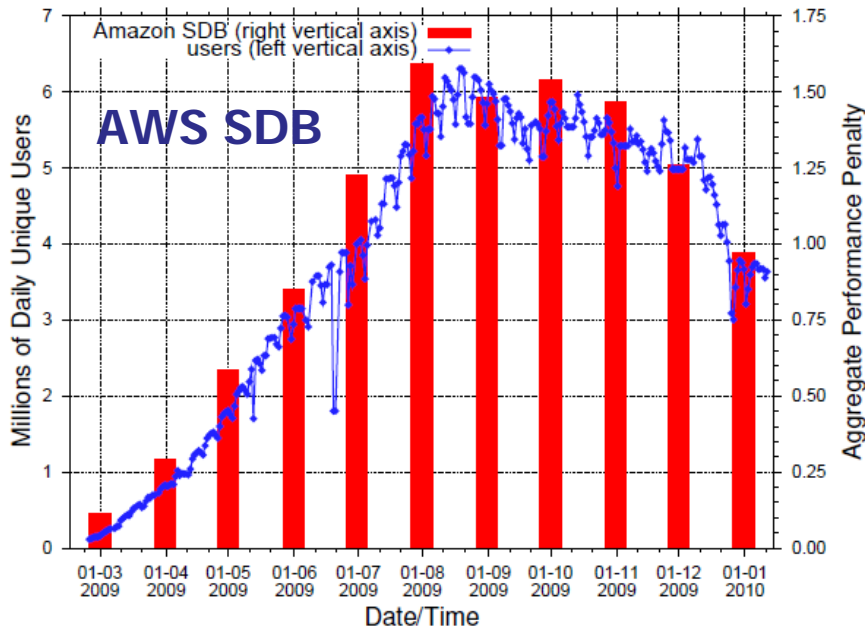
- Pref (Reference Performance): Average of the twelve monthly medians
- P(t): **random** value sampled from the distribution corresponding to the current month at time t (*Performance is like a box of chocolates, you never know what you're gonna get ~ Forrest Gump*)
- max U(t): max number of users over the whole trace
- U(t): number of users at time t
- **APP—the lower the better**

# Game Status Maintenance (1/2): Scenario

- Maintenance of game status for a large-scale social game such as Farm Town or Mafia Wars which have millions of unique users daily
- AWS SDB and GAE Datastore
- We assume that the number of database operations depends linearly on the number of daily unique users

# Game Status Maintenance (2): Results

Q2



- Big discrepancy between SDB and Datastore services
- **Sep'09-Jan'10**: APP of Datastore is well below than that of SDB due to increasing performance of Datastore
- APP of Datastore  $\sim 1$   $\Rightarrow$  no performance penalty
- APP of SDB  $\sim 1.4$   $\Rightarrow$  %40 higher performance penalty than SDB

# Agenda

1. An Introduction to IaaS Cloud Comput
2. Research Questions or Why We Need B
3. A General Approach and Its Main Chall
4. **IaaS Cloud Workloads (Q0)**
5. **IaaS Cloud Performance (Q1) & Perf. Variability (Q2)**
6. **Provisioning & Allocation Policies for IaaS Clouds (Q3)**
7. **Big Data: Large-Scale Graph Processing (Q4)**
8. [Conclusion](#)



**Workloads**

**Performance**

**Variability**

**Policies**

**Big Data:  
Graphs**

# IaaS Cloud Policies: Our Team



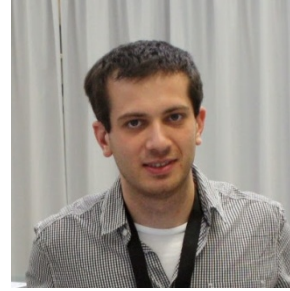
Alexandru Iosup  
TU Delft

Provisioning  
Allocation  
Elasticity  
Utility  
Isolation  
Multi-Tenancy



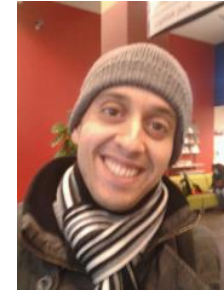
Dick Epema  
TU Delft

Provisioning  
Allocation  
Koala



Bogdan Ghit  
TU Delft

Provisioning  
Allocation  
Koala

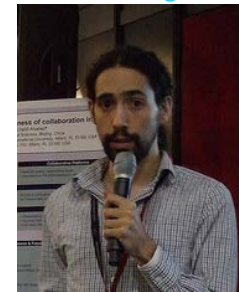


Athanasios Antoniou  
TU Delft

Provisioning  
Allocation  
Isolation  
Utility



Orna Agmon-Ben Yehuda  
Technion  
Elasticity, Utility



David Villegas  
FIU/IBM  
Elasticity, Utility

# What I'll Talk About

## Provisioning and Allocation Policies for IaaS Clouds (Q3)

1. Experimental setup
2. Experimental results



# Provisioning and Allocation Policies\*

\* For User-Level Scheduling

- Provisioning

Policy	Class	Trigger	Adaptive
Startup	Static	–	–
OnDemand	Dynamic	QueueSize	No
ExecTime	Dynamic	Exec.Time	Yes
ExecAvg	Dynamic	Exec.Time	Yes
ExecKN	Dynamic	Exec.Time	Yes
QueueWait	Dynamic	Wait Time	Yes

- Allocation

Policy	Queue-based	Known job durations
FCFS	Yes	No
FCFS-NW	No	No
SJF	Yes	Yes

- Also looked at combined Provisioning + Allocation policies

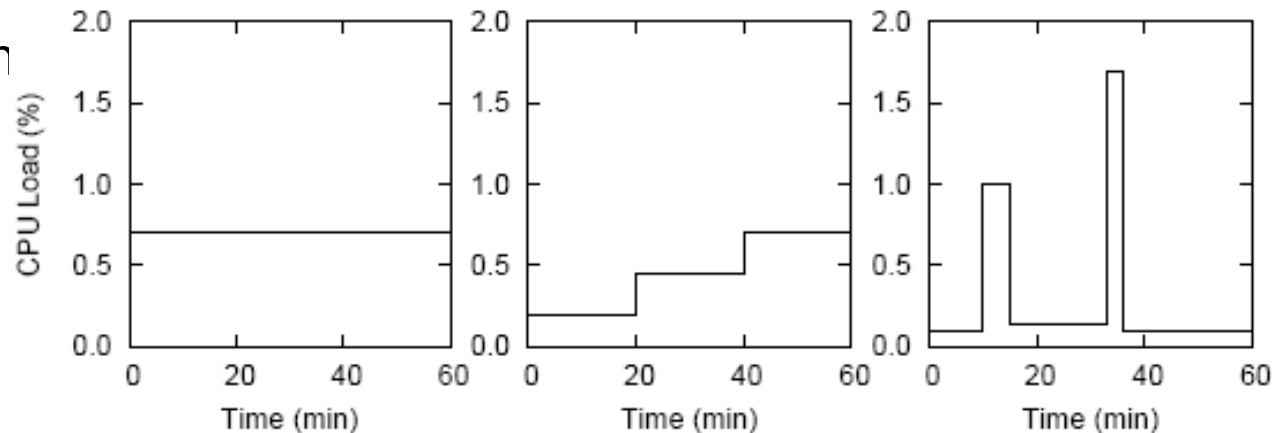
**The SkyMark Tool for IaaS Cloud Benchmarking**

# Experimental Setup (1)

- Environments
  - DAS4, Florida International University (FIU)
  - Amazon EC2

- Workloads
  - Bottleneck
  - Arrival pattern

Workload Unit	CPU	Memory	I/O	Appears in
WU1	X			WL1
WU2		X		WL2, WL4
WU3			X	WL3, WL4



# Experimental Setup (2)

## • Performance Metrics

- Traditional: Makespan, Job Slowdown
- Workload Speedup One (SU1)
- Workload Slowdown Infinite (SUinf)

$$SU_1(W) = \frac{MS(W)}{\sum_{i \in W} t_R(i)}$$

$$SU_\infty(W) = \frac{MS(W)}{\max_{i \in W} t_R(i)}$$

## • Cost Metrics

- Actual Cost (Ca)
- Charged Cost (Cc)

$$C_a(W) = \sum_{i \in \text{leased VMs}} t_{\text{stop}}(i) - t_{\text{start}}(i)$$

$$C_c(W) = \sum_{i \in \text{leased VMs}} [t_{\text{stop}}(i) - t_{\text{start}}(i)]$$

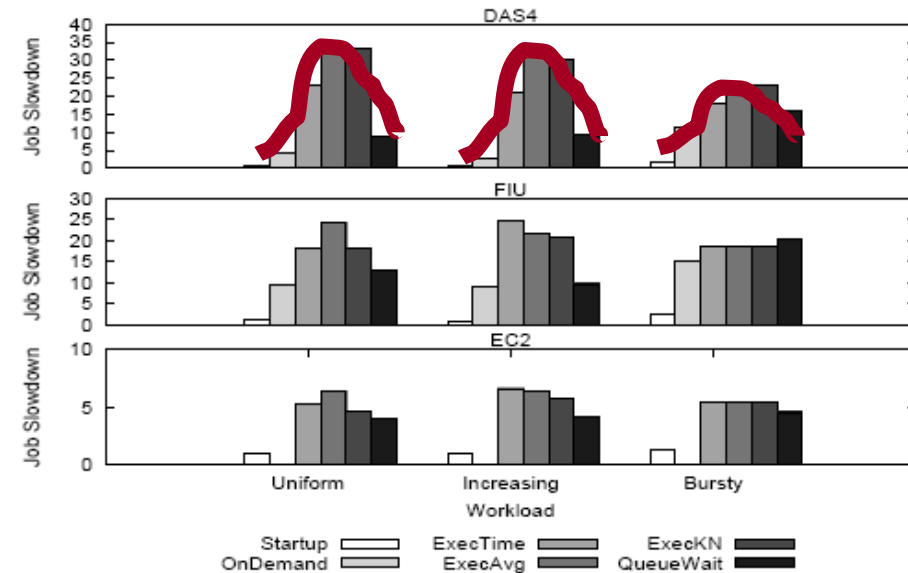
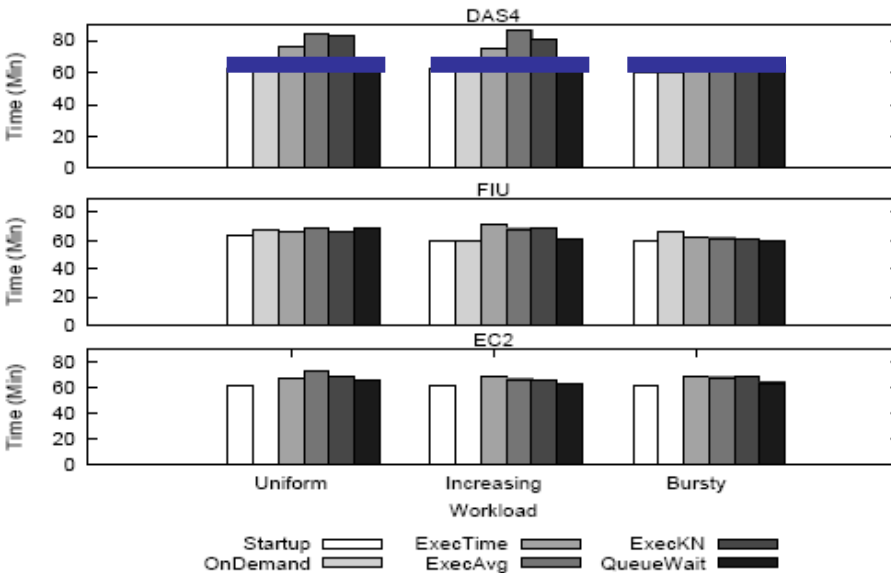
## • Compound Metrics

- Cost Efficiency (Ceff)
- Utility

$$C_{\text{eff}}(W) = \frac{C_c(W)}{C_a(W)}$$

$$U(W) = \frac{SU_1(W)}{C_c(W)}$$

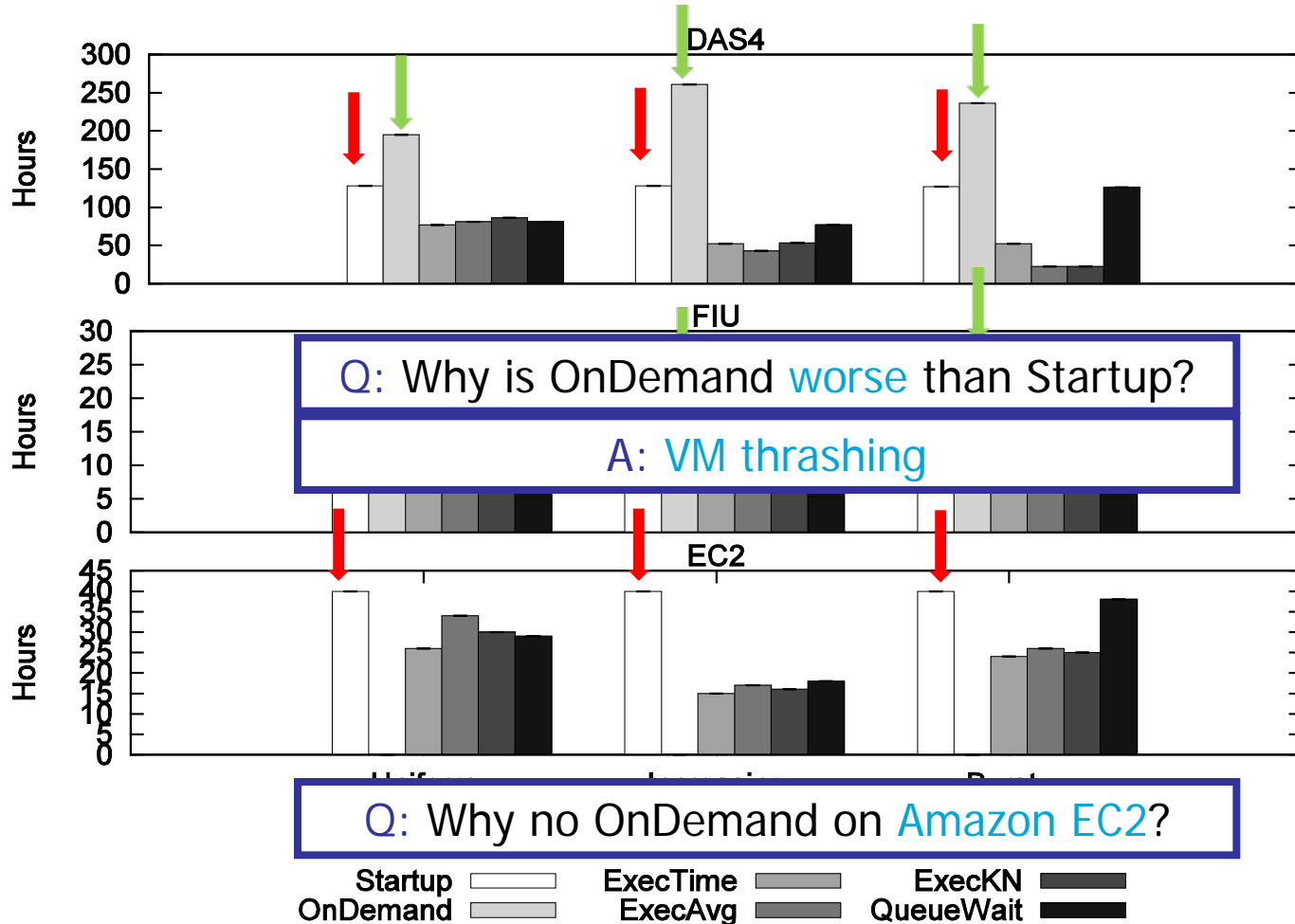
# Performance Metrics



- Makespan very similar
- Very different job slowdown

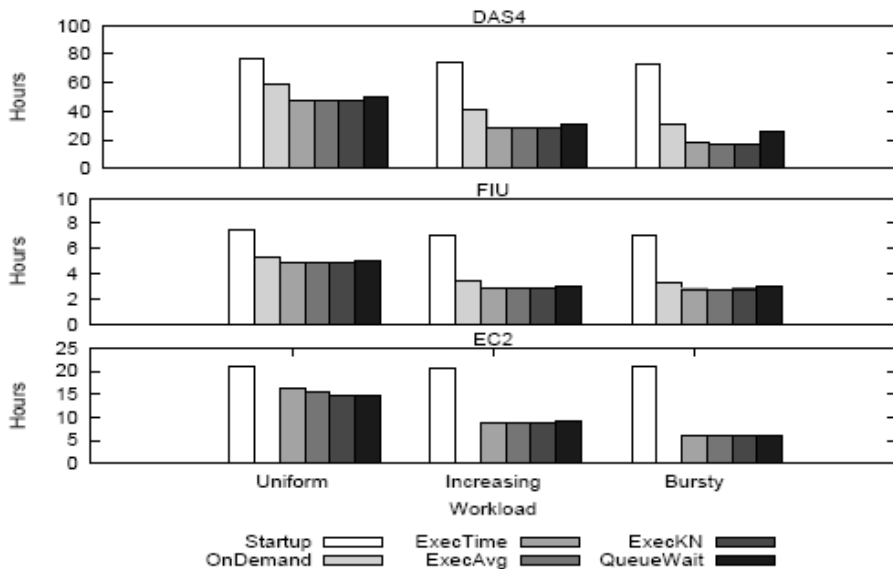
# Cost Metrics

## Charged Cost ( $C_c$ )

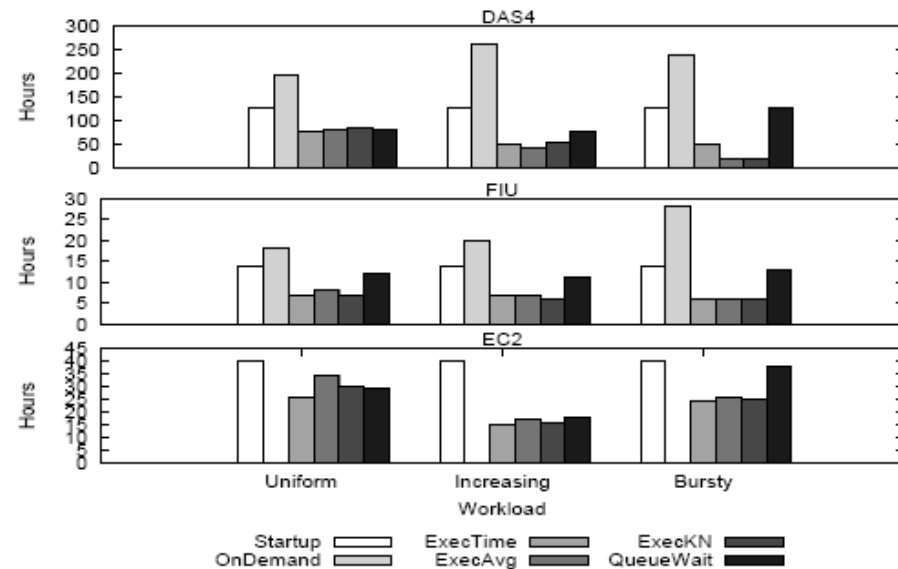


# Cost Metrics

## Actual Cost



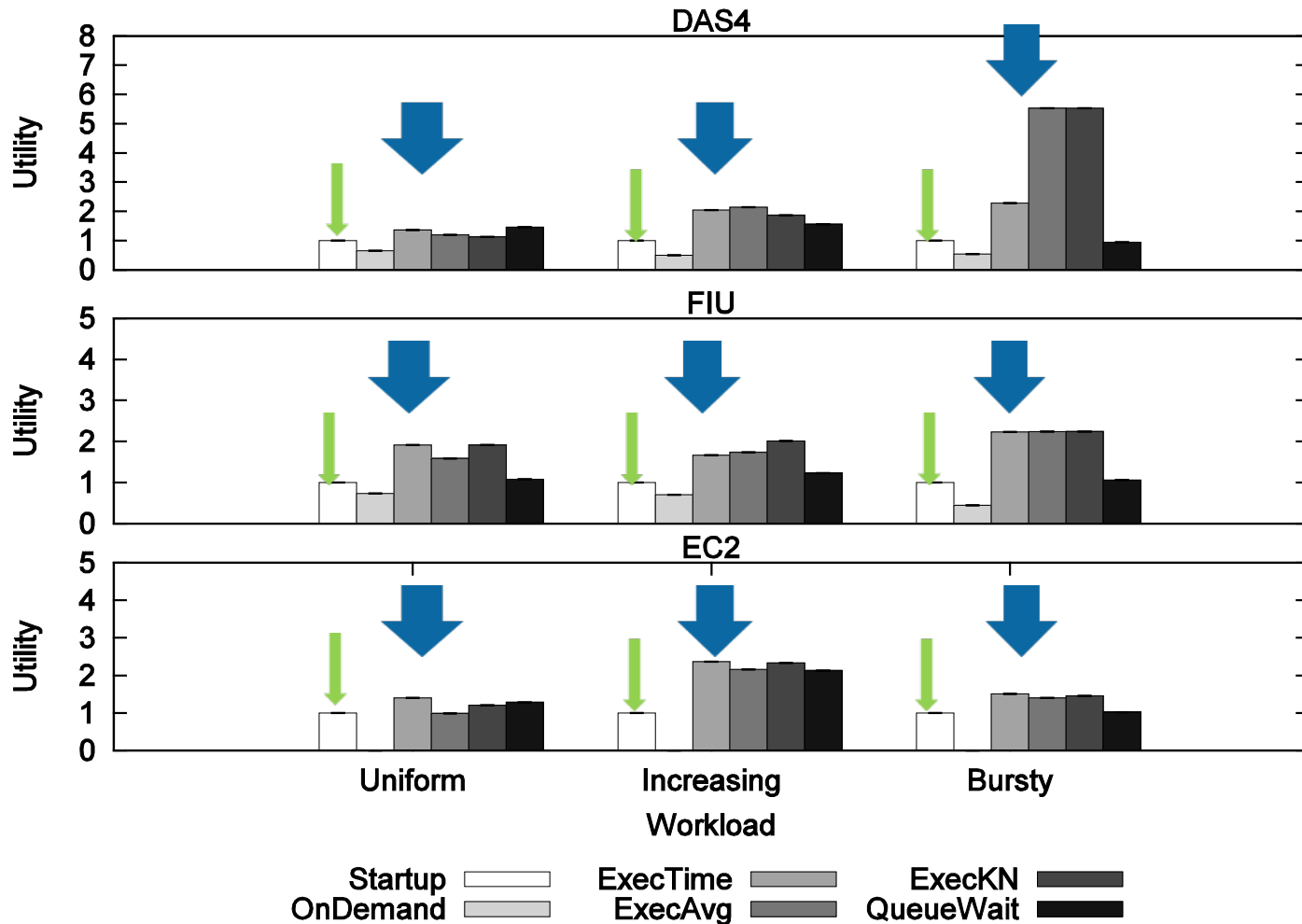
## Charged Cost



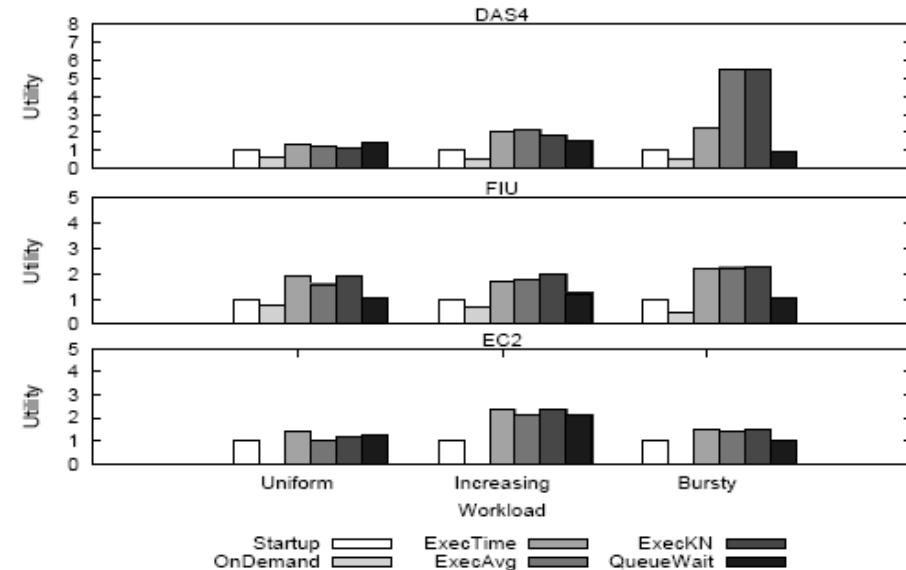
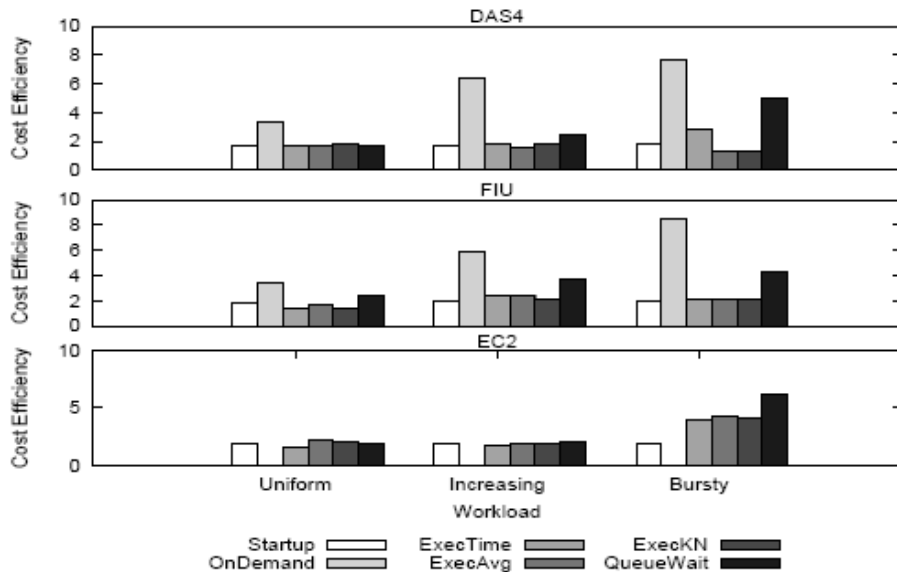
- Very different results between actual and charged
  - Cloud charging function an important selection criterion
- All policies better than Startup in actual cost
- Policies much better/worse than Startup in charged cost

# Compound Metrics (Utilities)

Utility ( $U$ )



# Compound Metrics



- Trade-off Utility-Cost still needs investigation
- **Performance or Cost, not both:**  
the policies we have studied improve one, but not both



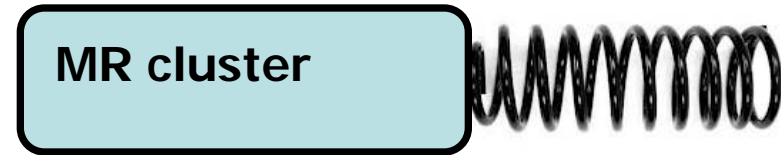
# Ad: Resizing MapReduce Clusters

- **Motivation:**

- Performance and data isolation
- Deployment version and user isolation
- Capacity planning : efficiency—accuracy trade-off

- **Constraints:**

- Data is big and difficult to move
- Resources need to be released fast



- **Approach:**

- Grow / shrink at processing layer
- Resize based on resource utilization
- Policies for provisioning and allocation

# Agenda

1. An Introduction to IaaS Cloud Computing
2. Research Questions or Why We Need B...
3. A General Approach and Its Main Challenges
4. **IaaS Cloud Workloads (Q0)**
5. **IaaS Cloud Performance (Q1) & Perf. Variability (Q2)**
6. **Provisioning & Allocation Policies for IaaS Clouds (Q3)**
7. **Big Data: Large-Scale Graph Processing (Q4)**
8. [Conclusion](#)



**Workloads**

**Performance**

**Variability**

**Policies**

**Big Data:  
Graphs**

# Big Data/Graph Processing: Our Team



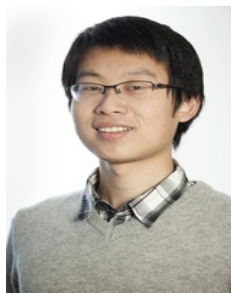
Alexandru Iosup  
TU Delft

Cloud Computing  
Gaming Analytics  
Performance Eval.  
Benchmarking  
Variability



Ana Lucia Varbanescu  
UvA

Parallel Computing  
Multi-cores/GPUs  
Performance Eval.  
Benchmarking  
Prediction



Yong Guo  
TU Delft

Cloud Computing  
Gaming Analytics  
Performance Eval.  
Benchmarking



Marcin Biczak  
TU Delft

Cloud Computing  
Performance Eval.  
Development



<http://www.pds.ewi.tudelft.nl/graphitti/>

Consultant for the project.

Not responsible for issues related  
to this work. Not representing

official products and/or company views.



Claudio Martella  
VU Amsterdam  
All things Giraph



Ted Willke  
Intel Corp.  
All things graph-processing

# What I'll Talk About

How well do graph-processing platforms perform?  
(Q4)

1. Motivation
2. Previous work
3. Method / Benchmarking suite
4. Experimental setup
5. Selected experimental results
6. Conclusion and ongoing work

# Why “How Well do Graph-Processing Platforms Perform?”

- Large-scale graphs exists in a wide range of areas: social networks, website links, online games, etc.
- Large number of **platforms** available to developers
  - Desktop: Neo4J, SNAP, etc.
  - Distributed: Giraph, GraphLab, etc.
  - Parallel: too many to mention

# Some Previous Work

Graph500.org: BFS on synthetic graphs

Performance evaluation in graph-processing (limited algorithms and graphs)

- Hadoop does not perform well [Warneke09]
- Graph partitioning improves the performance of Hadoop [Kambatla12]
- Trinity outperforms Giraph in BFS [Shao12]
- Comparison of graph databases [Dominguez-Sal10]

Performance comparison in other applications

- Hadoop vs parallel DBMSs: grep, selection, aggregation, and join [Pavlo09]
- Hadoop vs High Performance Computing Cluster (HPCC): queries [Ouaknine12]
- Neo4j vs MySQL: queries [Vicknair10]

**Problem:** Large differences in performance profiles across different graph-processing **algorithms** and **data sets**

# Our Method

A benchmark suite for  
performance evaluation of graph-processing platforms

1. Multiple Metrics, e.g.,
  - Execution time
  - Normalized: EPS, VPS
  - Utilization
  
2. Representative graphs with various characteristics, e.g.,
  - Size
  - Directivity
  - Density
  
3. Typical graph algorithms, e.g.,
  - BFS
  - Connected components

# Benchmarking suite

## Data sets

Graphs	# V	# E	$d (\times 10^{-5})$	$\bar{D}$	Size	Directivity
Amazon	262.1 K	1.2 M	1.8	4.7	18 MB	directed
WikiTalk	2.4 M	5.0 M	0.1	2.1	87 MB	directed
KGS	293.3 K	16.6 M	38.5	112.9	210 MB	undirected
Citation	3.8 M	16.5 M	0.1	4.4	297 MB	directed
DotaLeague	61.2 K	50.9 M	2,719.0	1,663.2	655 MB	undirected
Synth	2.4 M	64.2 M	2.2	53.6	964 MB	undirected
Friendster	65.6 M	1.8 B	0.1	55.1	31 GB	undirected



Graph500

<http://www.graph500.org/>



The Game Trace Archive

<http://gta.st.ewi.tudelft.nl/>

Guo, Biczak, Varbanescu, Iosup, Martella, Willke.  
 How Well do Graph-Processing Platforms Perform?  
 An Empirical Performance Evaluation and Analysis

Graphitti



# Benchmarking Suite

## Algorithm classes

1. General Statistics (STATS: # vertices and edges, LCC)
2. Breadth First Search (BFS)
3. Connected Component (CONN)
4. Community Detection (COMM)
5. Graph Evolution (EVO)

# Benchmarking suite

## Platforms and Process

- Platforms



Giraph

- Process

- Evaluate baseline (out of the box) and tuned performance
- Evaluate performance on fixed-size system
- Future: evaluate performance on elastic-size system
- Evaluate scalability

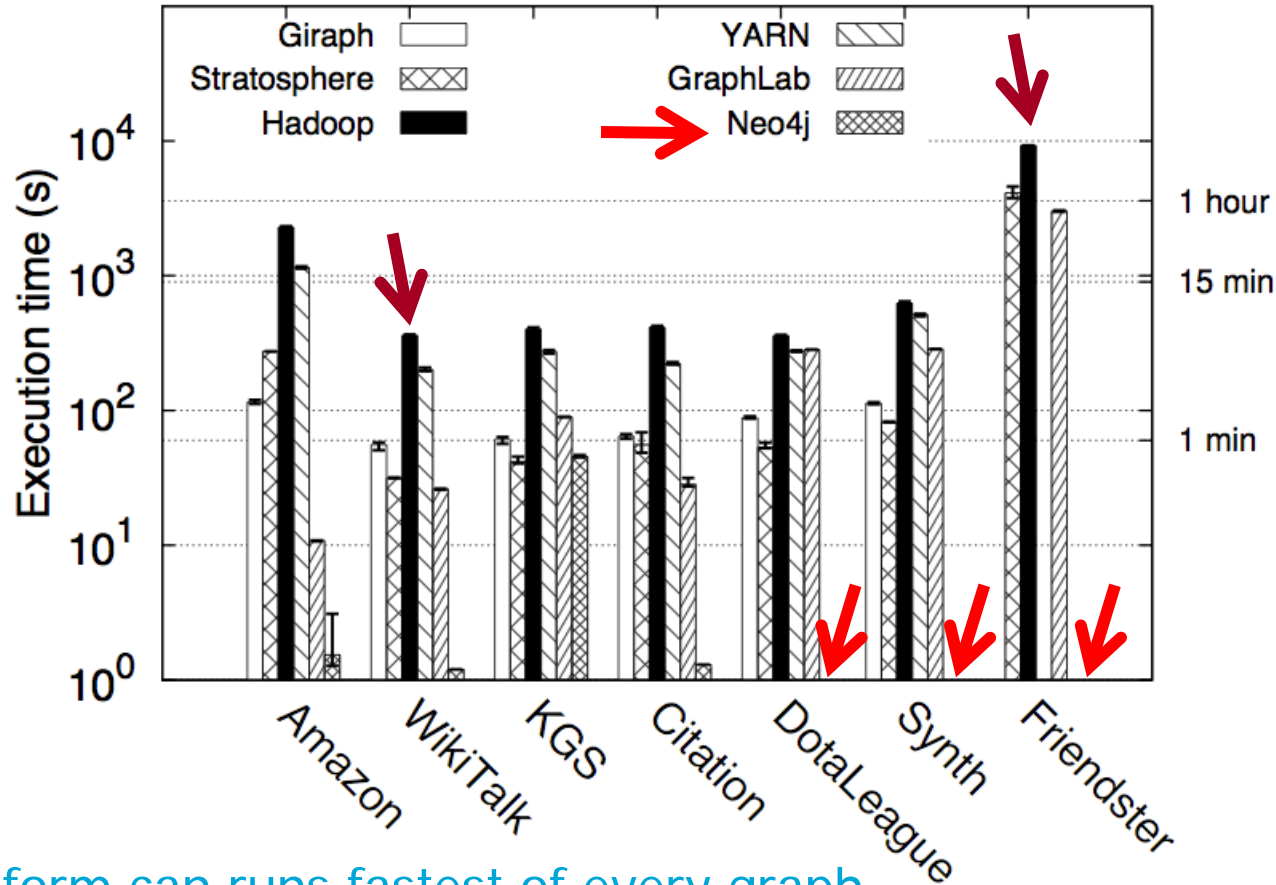
# Experimental setup

- Size
  - Most experiments take 20 working nodes
  - Up to 50 working nodes
- DAS4: a multi-cluster Dutch grid/cloud
  - Intel Xeon 2.4 GHz CPU (dual quad-core, 12 MB cache)
  - Memory 24 GB
  - 10 Gbit/s Infiniband network and 1 Gbit/s Ethernet network
  - Utilization monitoring: Ganglia
- HDFS used here as distributed file systems



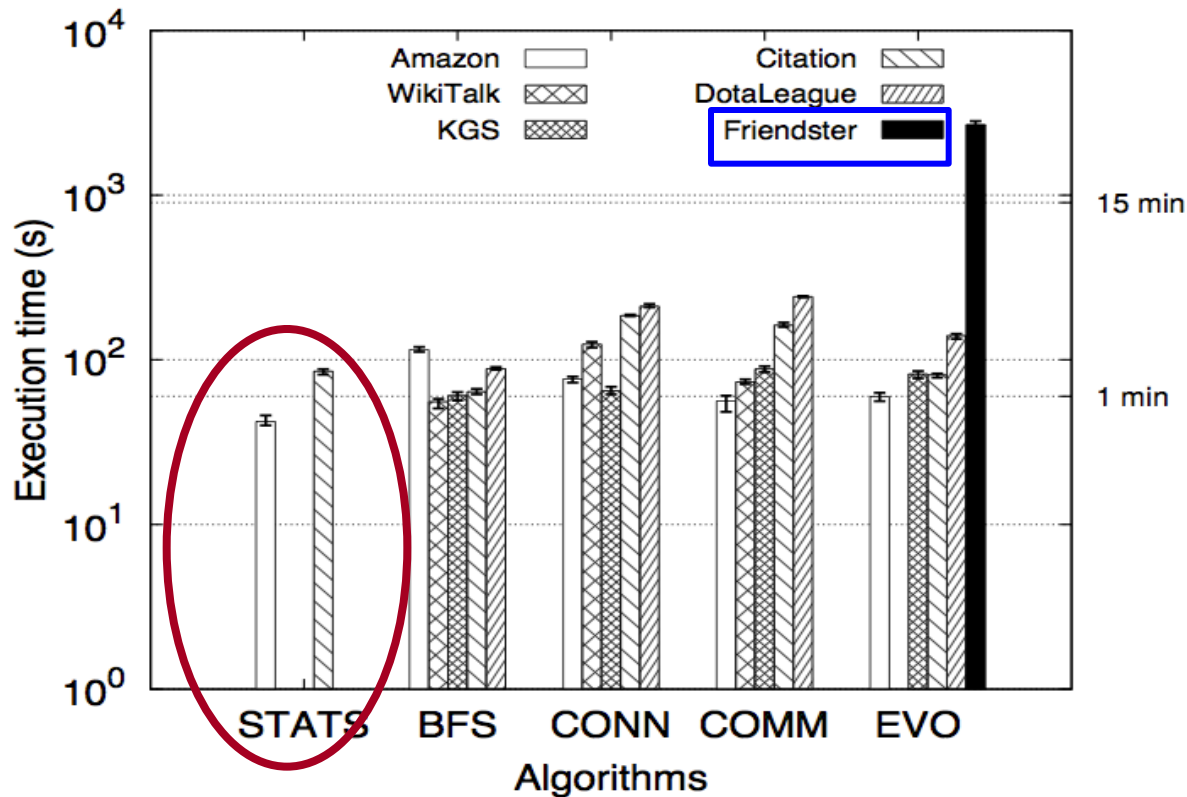
# BFS: results for all platforms, all data sets

Q4



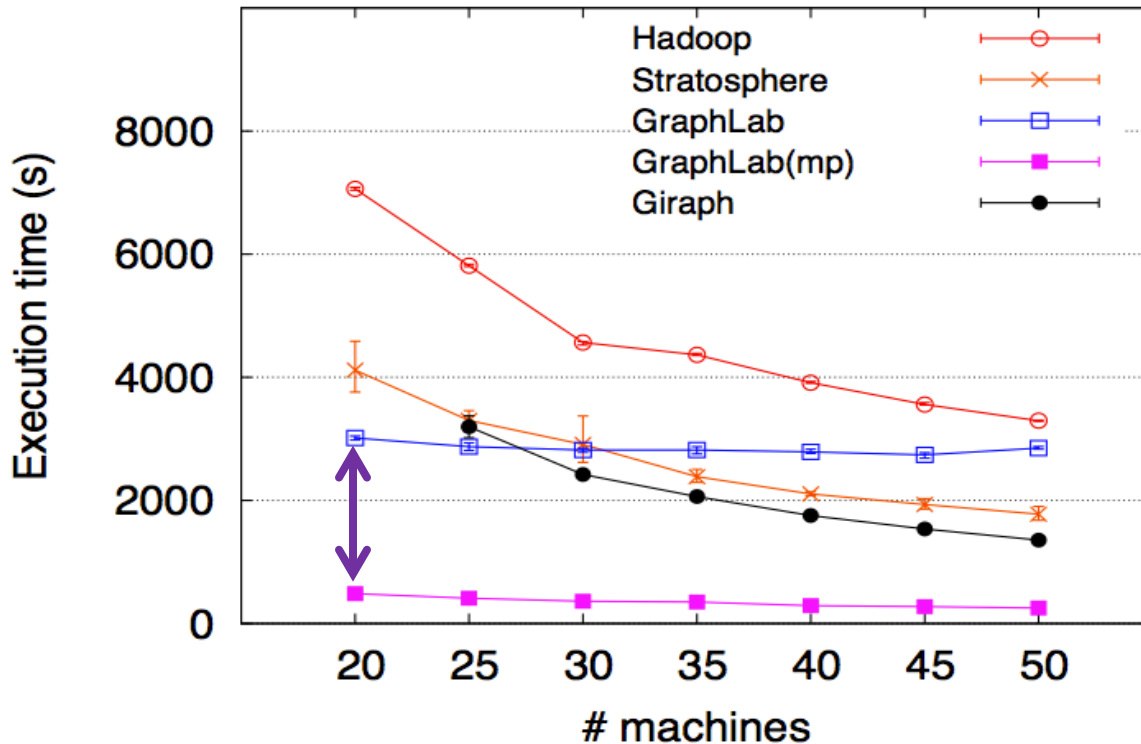
- No platform can run fastest of every graph
- Not all platforms can process all graphs
- Hadoop is the worst performer

# Giraph: results for all algorithms, all data sets



- Storing the whole graph in memory helps Giraph perform well
- Giraph may crash when **graphs** or **messages** become larger

# Horizontal scalability: BFS on Friendster (31 GB)



- Using more computing machines can reduce execution time
- Tuning needed for horizontal scalability, e.g., for GraphLab, split large input files into number of chunks equal to the number of machines

# Additional Overheads

## Data ingestion time

- Data ingestion
  - Batch system: one ingestion, multiple processing
  - Transactional system: one ingestion, one processing
- Data ingestion matters even for batch systems

	Amazon	DotaLeague	Friendster
HDFS	1 second	7 seconds	5 minutes
Neo4J	4 hours	6 <b>days</b>	<b>n/a</b>

# Conclusion and ongoing work

- Performance is  $f(\text{Data set, Algorithm, Platform, Deployment})$
- Cannot tell yet which of (Data set, Algorithm, Platform) the most important (depends on Platform and Deployment)
- Platforms have their own drawbacks
- Some platforms can scale up reasonably with cluster size (horizontally) or number of cores (vertically)
- Ongoing work
  - *Benchmarking* suite
  - Build a performance boundary model
  - Explore performance variability



# Agenda

1. An Introduction to IaaS Cloud Computing
2. Research Questions or Why We Need B...
3. A General Approach and Its Main Chall...
4. **IaaS Cloud Workloads (Q0)**
5. **IaaS Cloud Performance (Q1) & Perf. Variability (Q2)**
6. **Provisioning & Allocation Policies for IaaS Clouds (Q3)**
7. **Big Data: Large-Scale Graph Processing (Q4)**
8. [Conclusion](#)



**Workloads**

**Performance**

**Variability**

**Policies**

**Big Data:  
Graphs**

# Agenda

1. An Introduction to IaaS Cloud Computing
2. Research Questions or Why We Need Benchmarking?
3. A General Approach and Its Main Challenges
4. IaaS Cloud Workloads (Q0)
5. IaaS Cloud Performance (Q1) and Perf. Variability (Q2)
6. Provisioning and Allocation Policies for IaaS Clouds (Q3)
- 7. Conclusion**

# ~~Conclusion~~ Take-Home Message

- **IaaS cloud benchmarking: approach + 10 challenges**
- **Put 10-15% project effort in benchmarking = understanding how IaaS clouds really work**
  - Q0: Statistical workload models
  - Q1/Q2: Performance/variability
  - Q3: Provisioning and allocation
  - Q4: Big Data, Graph processing
- **Tools and Workload Models**
  - SkyMark
  - MapReduce
  - Graph processing benchmarking suite



<http://www.flickr.com/photos/dimitrisotiropoulos/4204766418/>

# Thank you for your attention!

## Questions? Suggestions? Observations?

More Info:

- <http://www.st.ewi.tudelft.nl/~iosup/research.html>
- [http://www.st.ewi.tudelft.nl/~iosup/research\\_cloud.html](http://www.st.ewi.tudelft.nl/~iosup/research_cloud.html)
- <http://www.pds.ewi.tudelft.nl/>



## Alexandru Iosup

[A.iosup@tudelft.nl](mailto:A.iosup@tudelft.nl)

<http://www.pds.ewi.tudelft.nl/~iosup/> (or google "iosup")

Parallel and Distributed Systems Group  
Delft University of Technology

Do not hesitate to  
contact me...



# WARNING: Ads

# CCGrid 2013

MAY 13-16, 2013 • DELFT, THE NETHERLANDS

The 13<sup>th</sup> IEEE/ACM International Symposium on  
Cluster, Cloud and Grid Computing



[www.pds.ewi.tudelft.nl/ccgrid2013](http://www.pds.ewi.tudelft.nl/ccgrid2013)

**Dick Epema, General Chair**  
Delft University of Technology Delft

**Delft, the Netherlands**  
**May 13-16, 2013**

**Thomas Fahringer, PC Chair**  
University of Innsbruck

**Paper submission deadline:**  
**November 22, 2012**



# If you have an interest in novel aspects of performance, you should join the SPEC RG

## ▶ Find a new venue to discuss your work

- ▶ Exchange with experts on how the performance of systems can be measured and engineered
- ▶ Find out about novel methods and current trends in performance engineering
- ▶ Get in contact with leading organizations in the field of performance evaluation

## ▶ Find a new group of potential employees

## ▶ Join a SPEC standardization process

## ▶ Performance in a broad sense:

- ▶ *Classical performance metrics:* Response time, throughput, scalability, resource/cost/energy, efficiency, elasticity
- ▶ *Plus dependability in general:* Availability, reliability, and security