# The Performance of Big Data Workloads in Cloud Datacenters

**Alexandru Uta**, Alexandru Custura, Harry Obaseki

a.uta@vu.nl

Vrije Universiteit Amsterdam
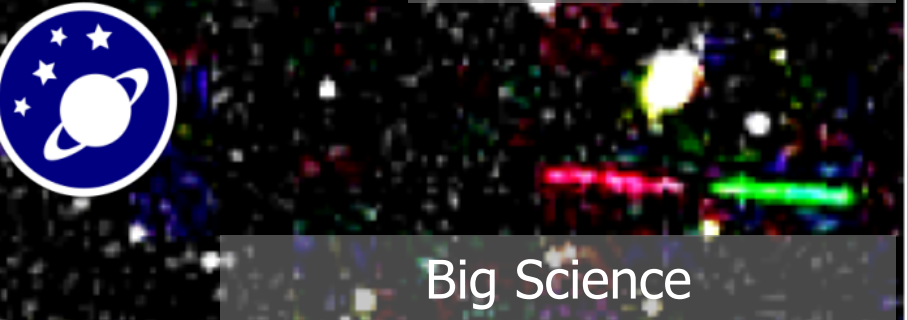
Massivizing Computer Systems
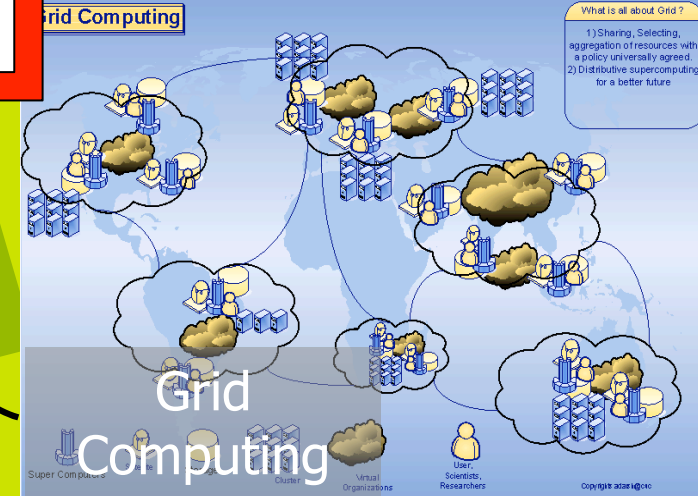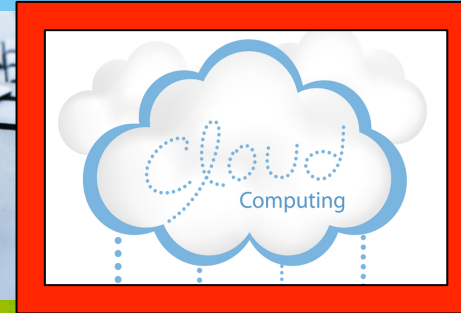
June 11, 2018

VRIJE UNIVERSITEIT AMSTERDAM

# Massivizing Computer Systems
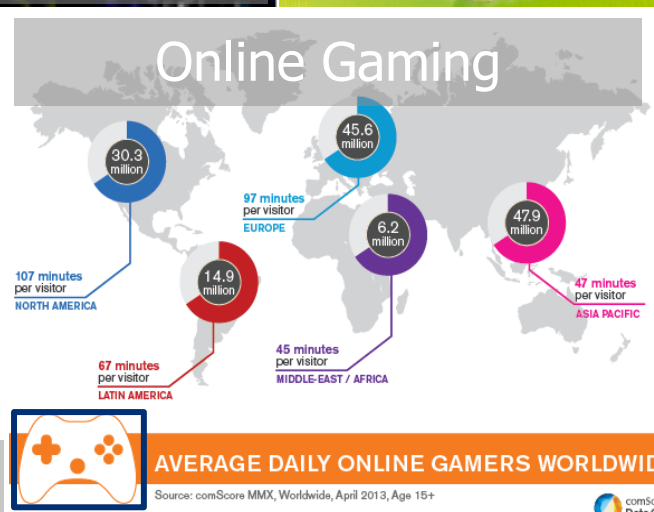


Education for Everyone (Online)

Business Services

Cloud Computing

Grid Computing

Big Science

Grid Computing

Online Gaming

My other computer is a data center

BIG DATA

Datacenters

Daily Life

ABN·AMRO

AVERAGE DAILY ONLINE GAMERS WORLDWIDE
Source: comScore MMX, Worldwide, April 2013, Age 15+
comScore Data Gem

97 minutes per visitor
EUROPE
45.6 million

107 minutes per visitor
NORTH AMERICA
30.3 million

67 minutes per visitor
LATIN AMERICA
14.9 million

45 minutes per visitor
MIDDLE-EAST / AFRICA
6.2 million

47 minutes per visitor
ASIA PACIFIC
47.9 million

What is all about Grid ?
1) Sharing, Selecting, aggregation of resources with a policy universally agreed.
2) Distributive supercomputing for a better future

# Convenient to use big data + cloud

BIG DATA LANDSCAPE 2017

**What happens when everybody runs big data in the cloud?**

Image courtesy of mattturck.com

# Co-location induces (resource) performance variability



How does resource interference affect performance?

Resource contention produces performance variability in clouds!

VU VRIJE UNIVERSITEIT AMSTERDAM

# Co-location induces (resource) performance variability



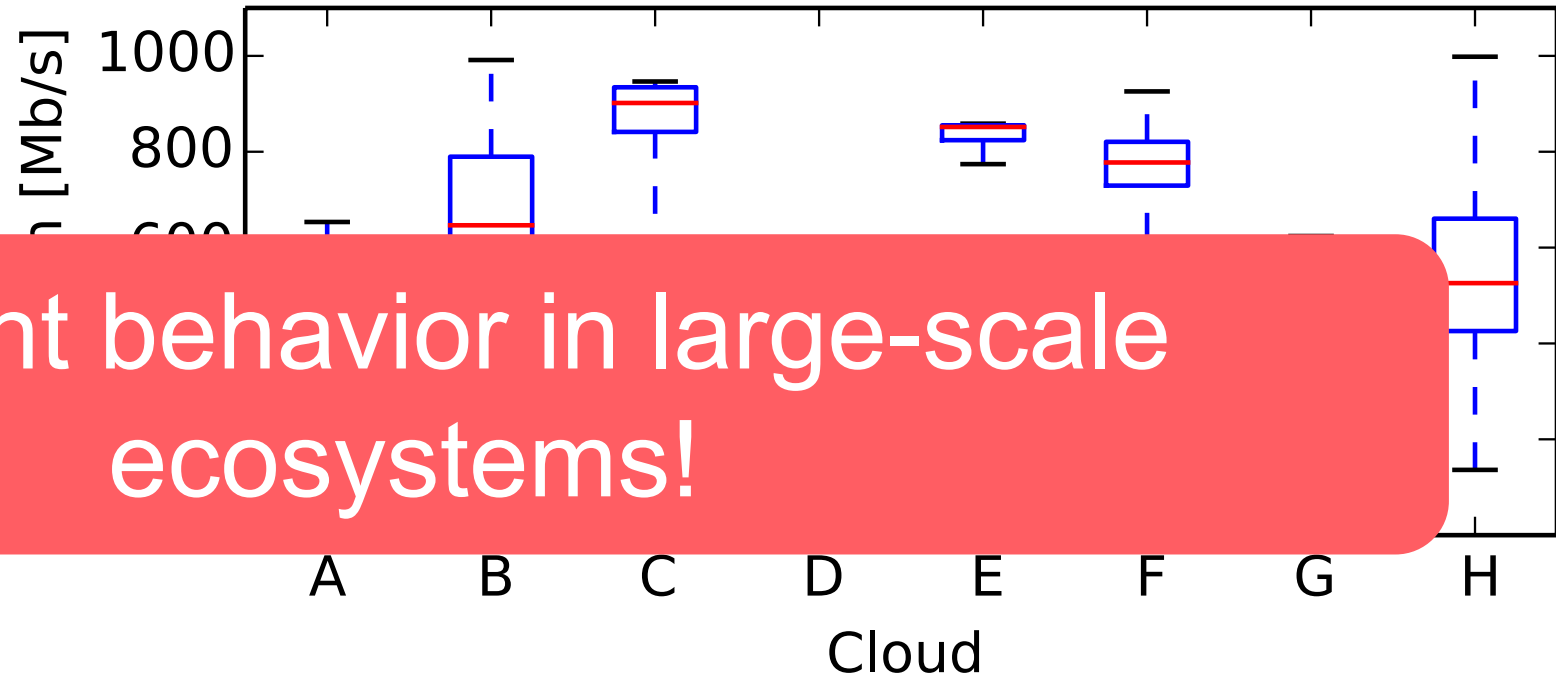How does workload variability affect performance?

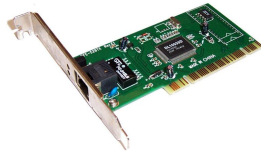Workload variability produces performance variability!

# Cloud (resource) performance is highly variable!

- Due to:
  - Co-location
  - Virtualization
  - Workload variability
  - Network congestion

- Affec



Emergent behavior in large-scale ecosystems!

Ballani et al., SIGCOMM 2011

VRIJE UNIVERSITEIT AMSTERDAM

# Convenient to use big data + cloud, but...

Variability entails:

- Poor performance predictions
- Over-provisioning

- Poor scheduling decisions
- Extra costs

## How to study performance variability? How to control the variability?

VRIJE
UNIVERSITEIT
AMSTERDAM

# How to study performance variability?
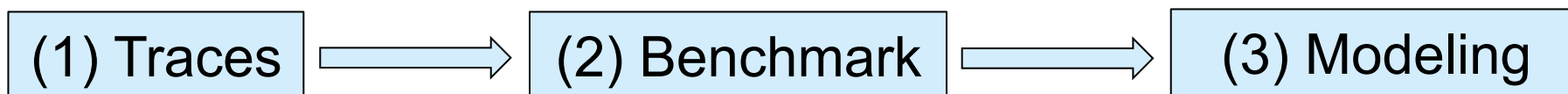
Traditional performance analysis:

- (1) Trace analysis

- (2) Benchmarking

- (3) Performance modeling

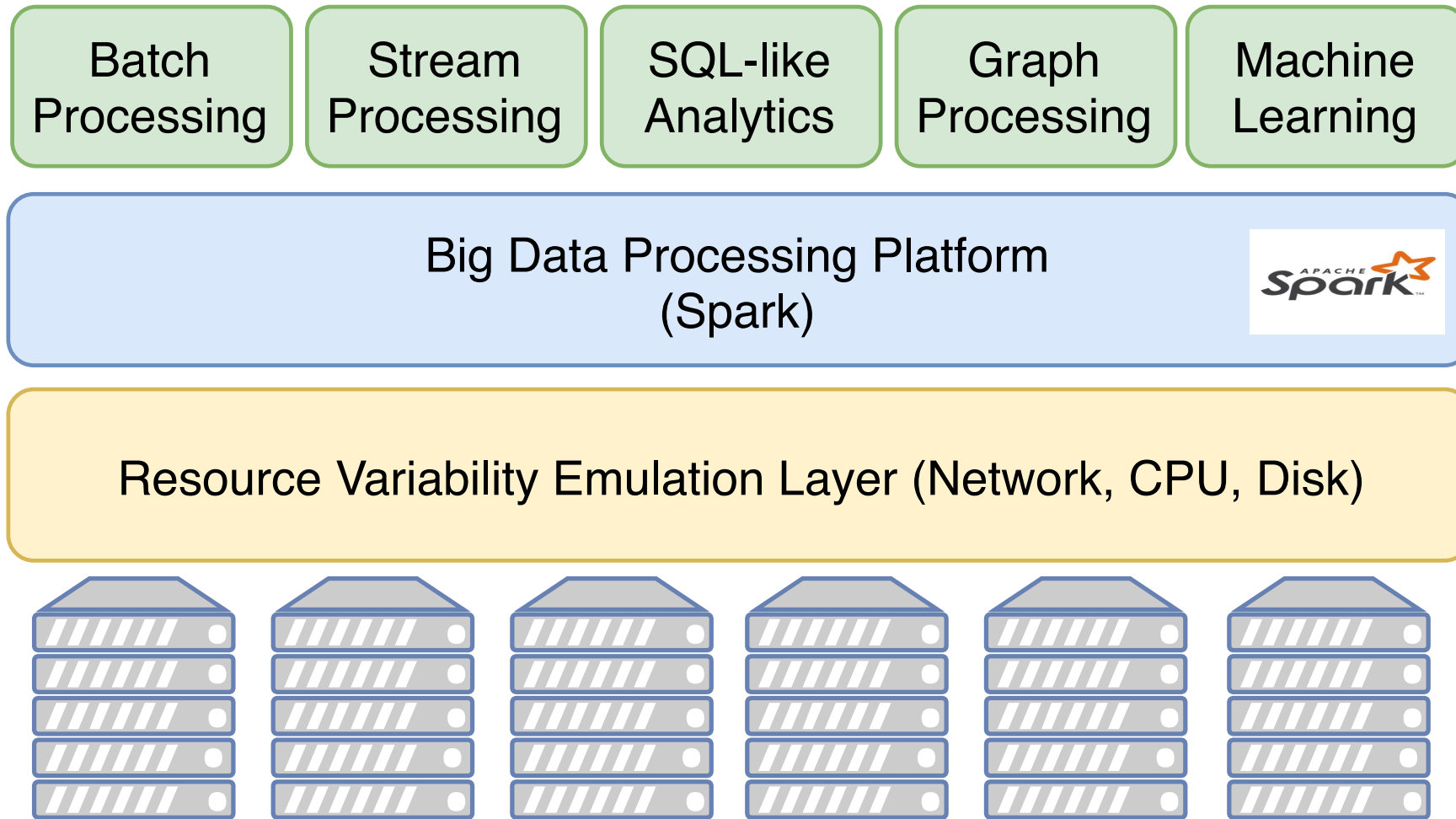Current models, benchmarks do not consider resource variability!

- No study on resource performance variability and big data
- Variability **within** clouds and **between** clouds (performance portability issues)

VU **VRIJE UNIVERSITEIT AMSTERDAM**

# A Framework for Studying Performance Variability

**1** • Fallback to empirical evaluation based on previous observations

**2** • Controlled environment that emulates real-world variability scenarios

• Multiple classes of big data applications

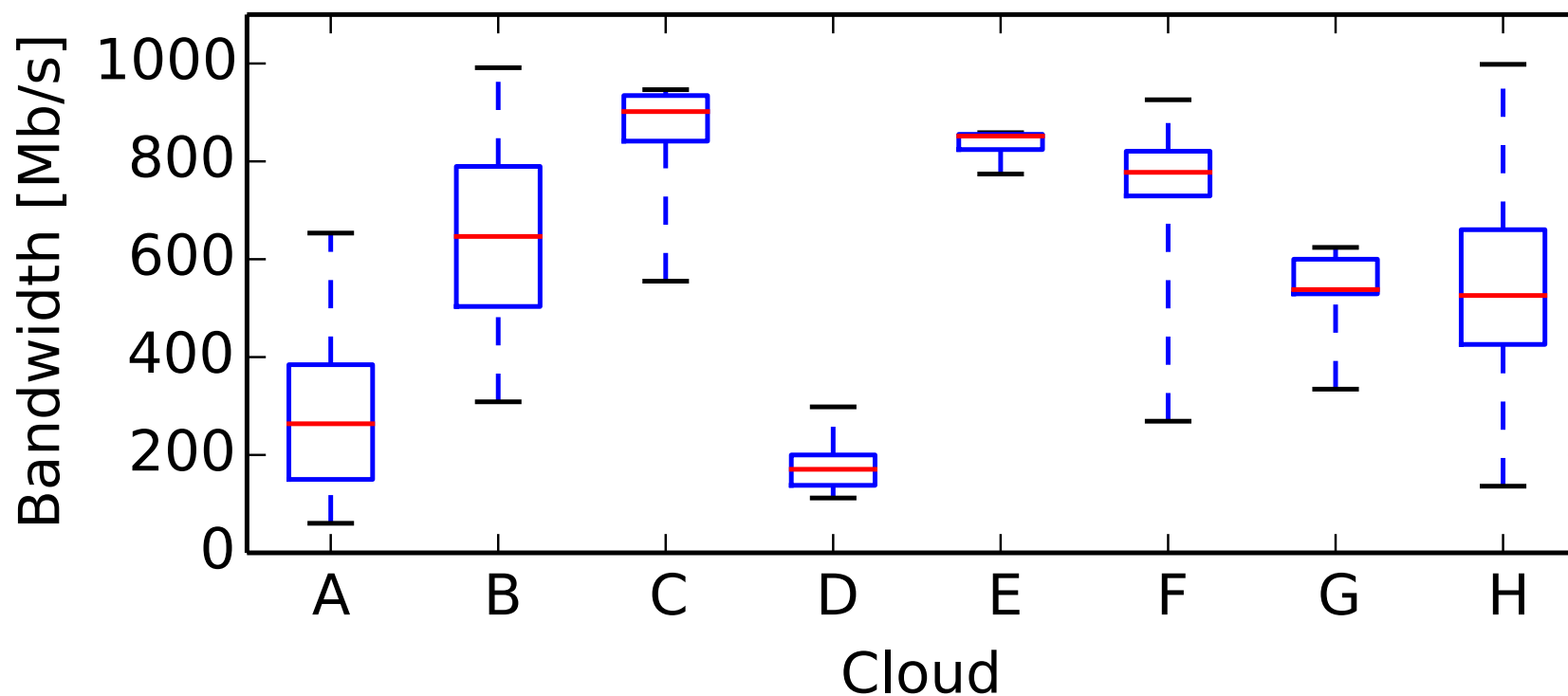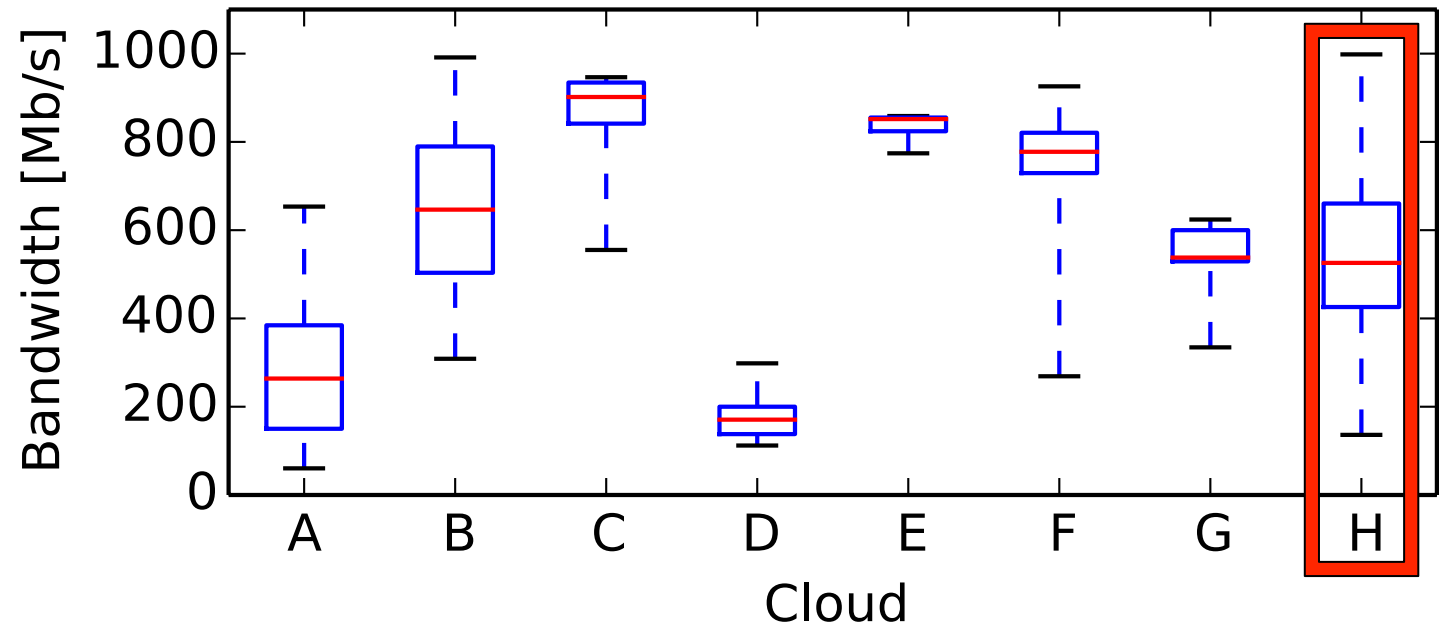**3** • Statistical analysis and performance modeling to understand correlations

(1) Traces → (2) Benchmark → (3) Modeling

VRIJE
UNIVERSITEIT
AMSTERDAM

# Benchmarking Performance Variability

| Batch Processing | Stream Processing | SQL-like Analytics | Graph Processing | Machine Learning |
|---|---|---|---|---|

Big Data Processing Platform
(Spark)

Resource Variability Emulation Layer (Network, CPU, Disk)

- Systematic study using A-H cloud bandwidth distributions
- Run a series of big data applications

# Cloud network bandwidth emulation
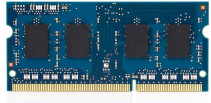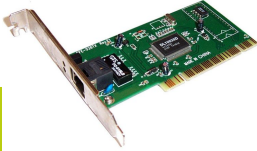
- For each distribution:



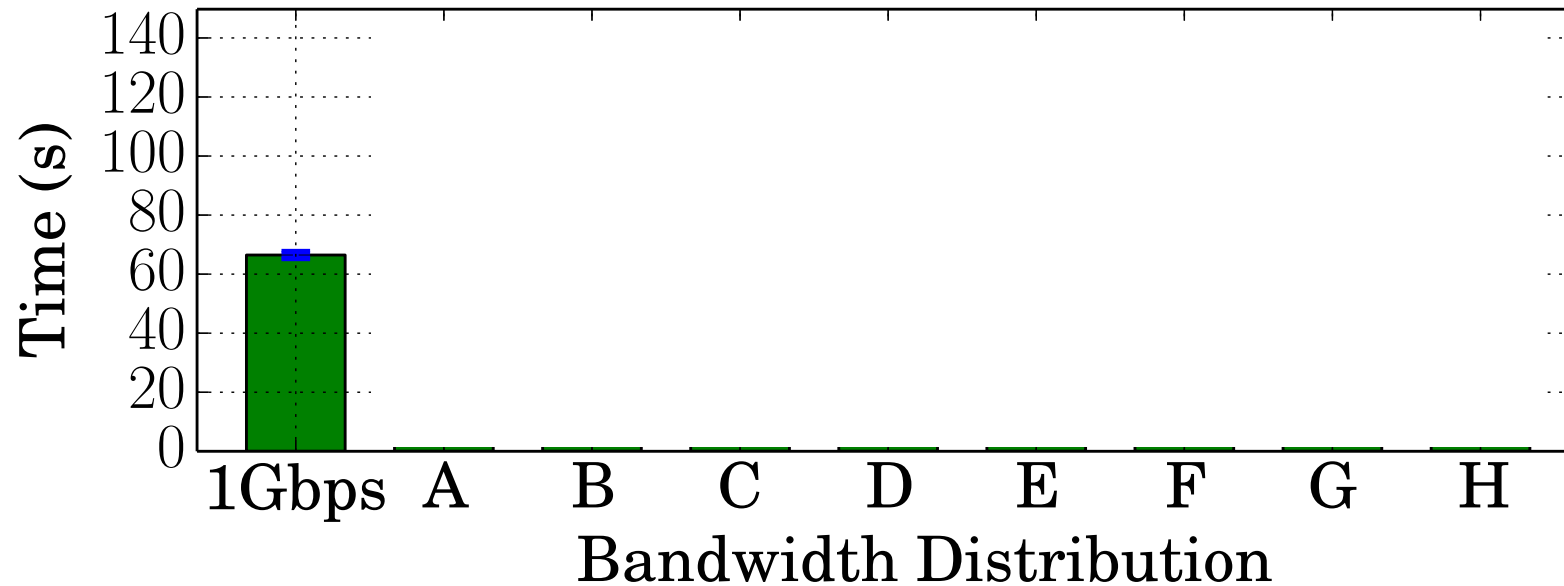Vary bandwidth

Cluster

VRIJE
UNIVERSITEIT
AMSTERDAM

# Big Data Workloads
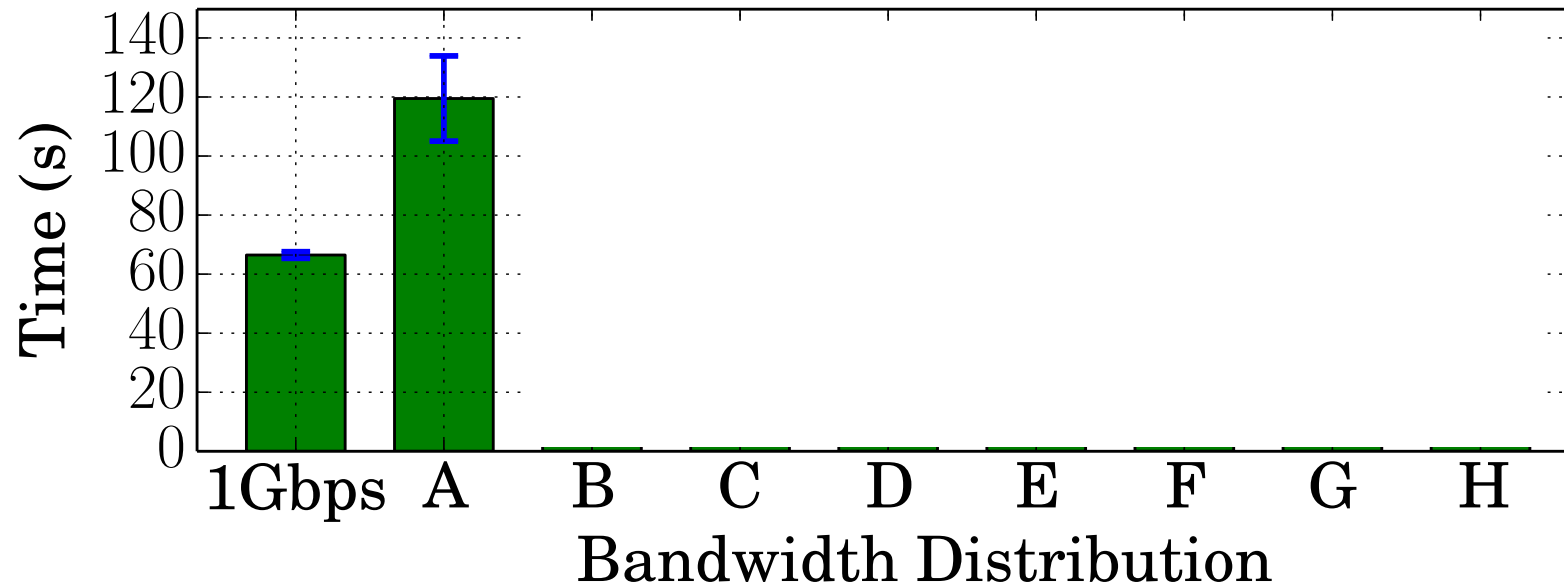
- HiBench suite, MapReduce-style apps
- 6 real-world applications from various domains
- Each app having different resource usage

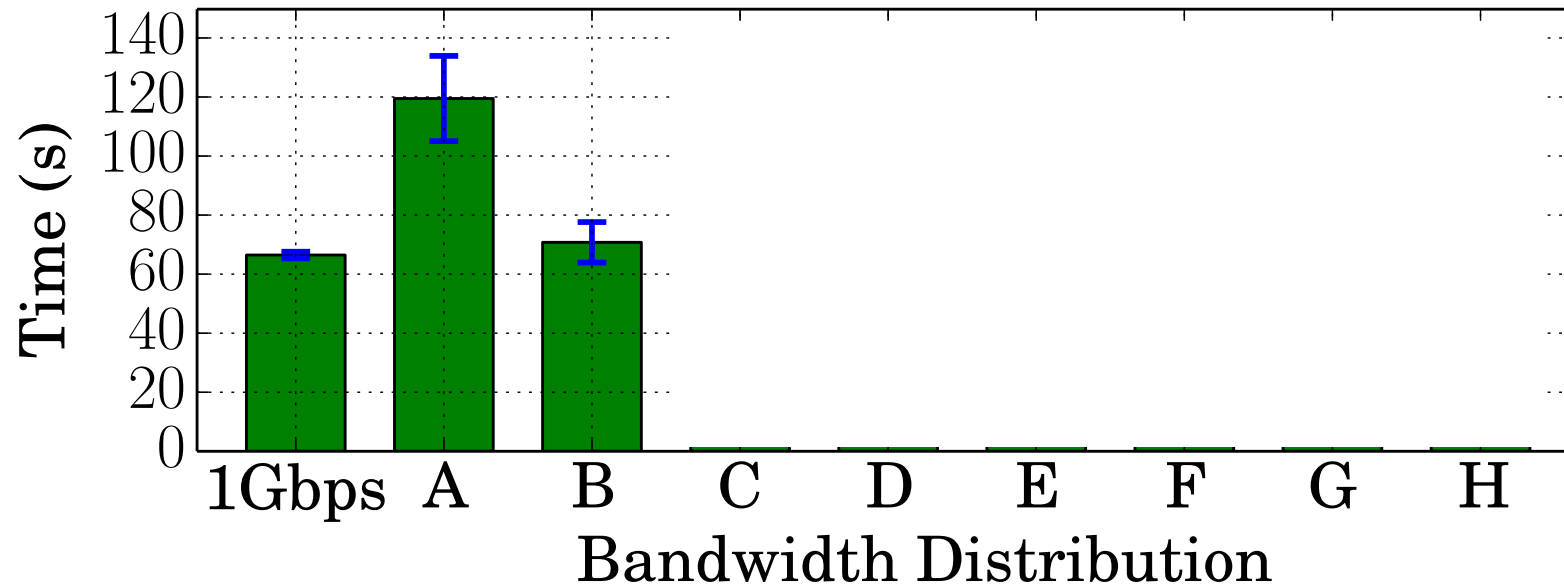| Application |  | |  | |  | |  |
|---|---|---|---|---|---|---|---|
| Wordcount | ++ | | -- | | 0 | | 0 |
| Sort | -- | | ++ | | 0 | | ++ |
| Terasort | ++ | | 0 | | ++ | | ++ |
| Naïve Bayes | 0 | | 0 | | ++ | | -- |
| K-means | ++ | | -- | | 0 | | -- |
| PageRank | 0 | | -- | | 0 | | -- |

VRIJE
UNIVERSITEIT
AMSTERDAM

# Variable network = Variable Runtime (Terasort)

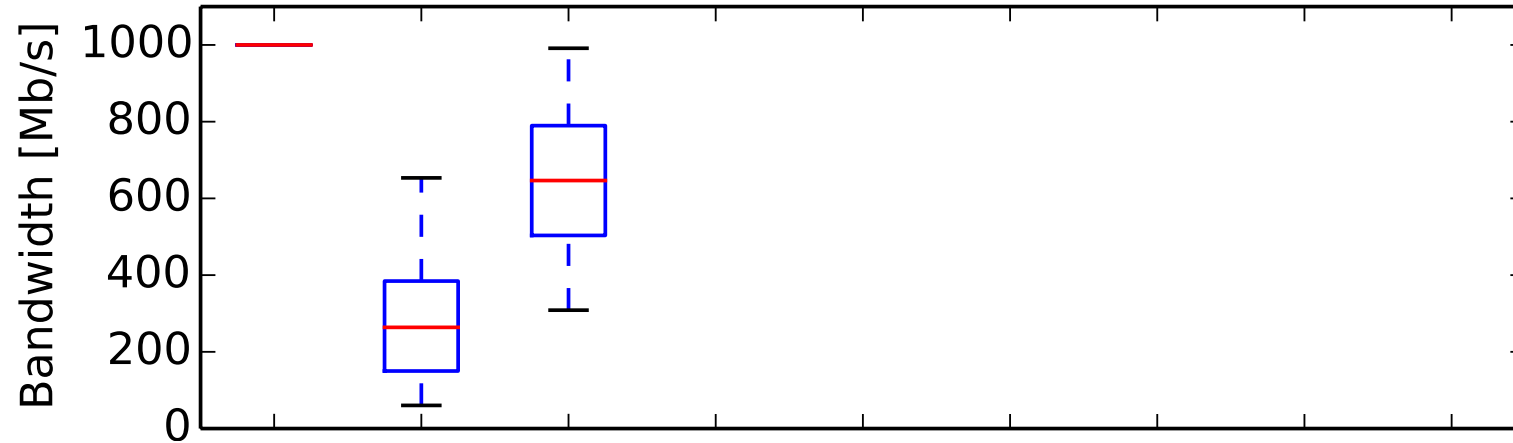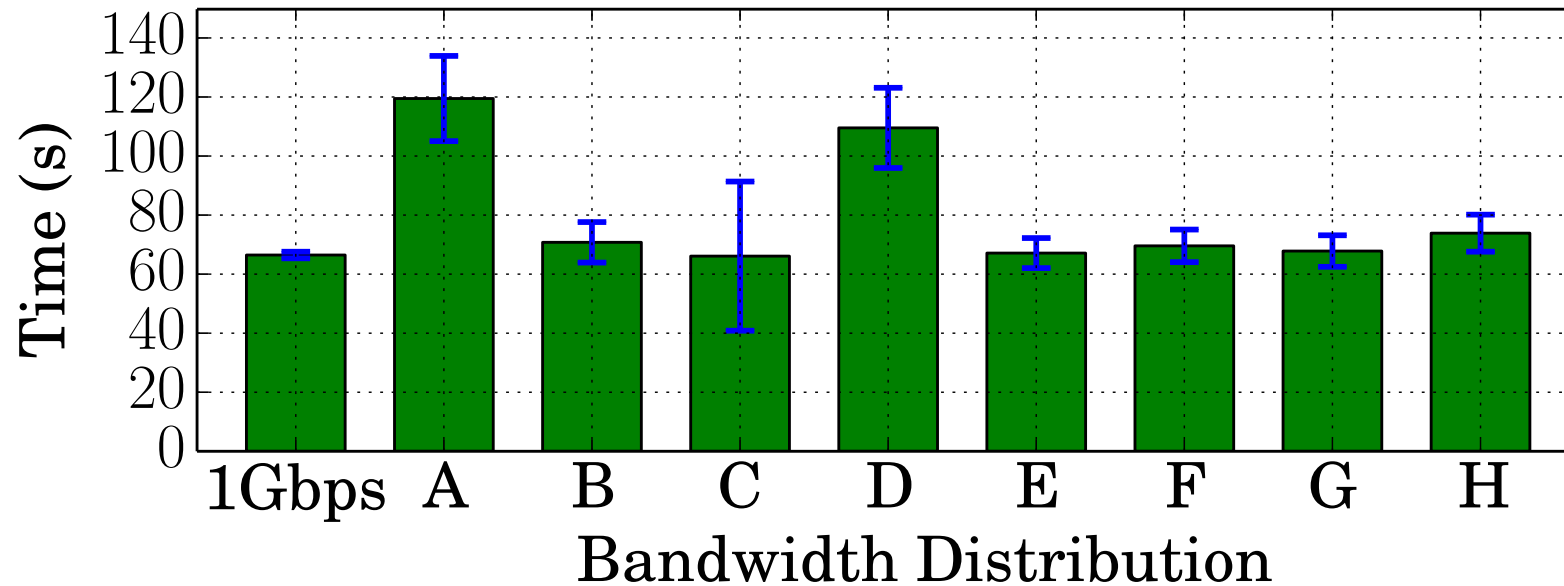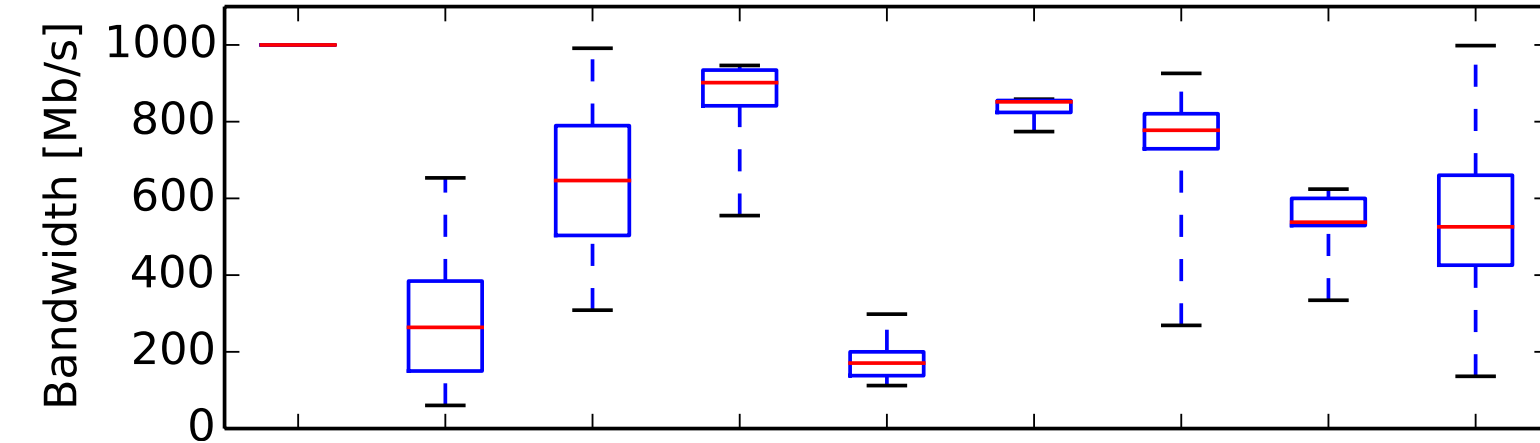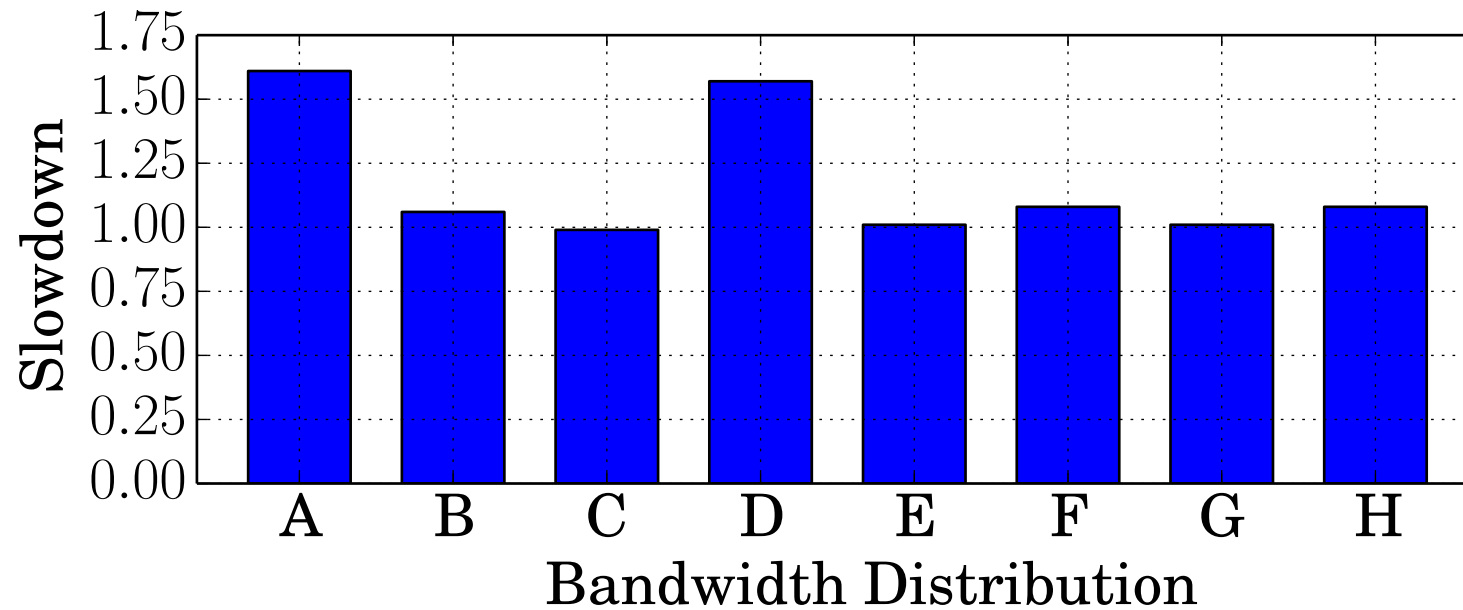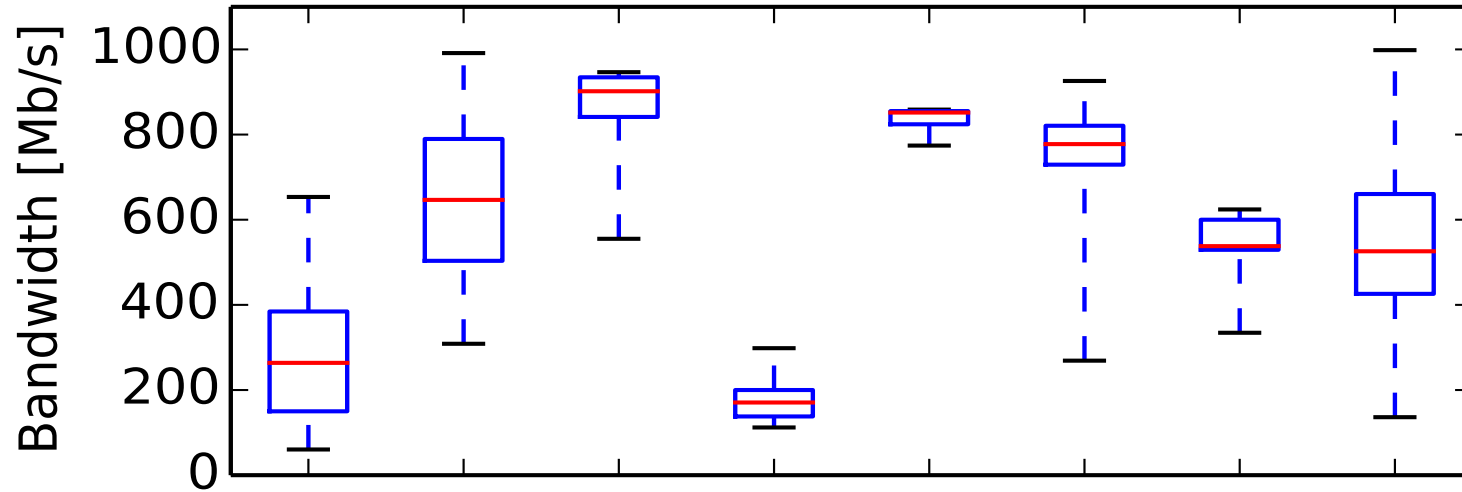# Variable network = Variable Runtime (Terasort)
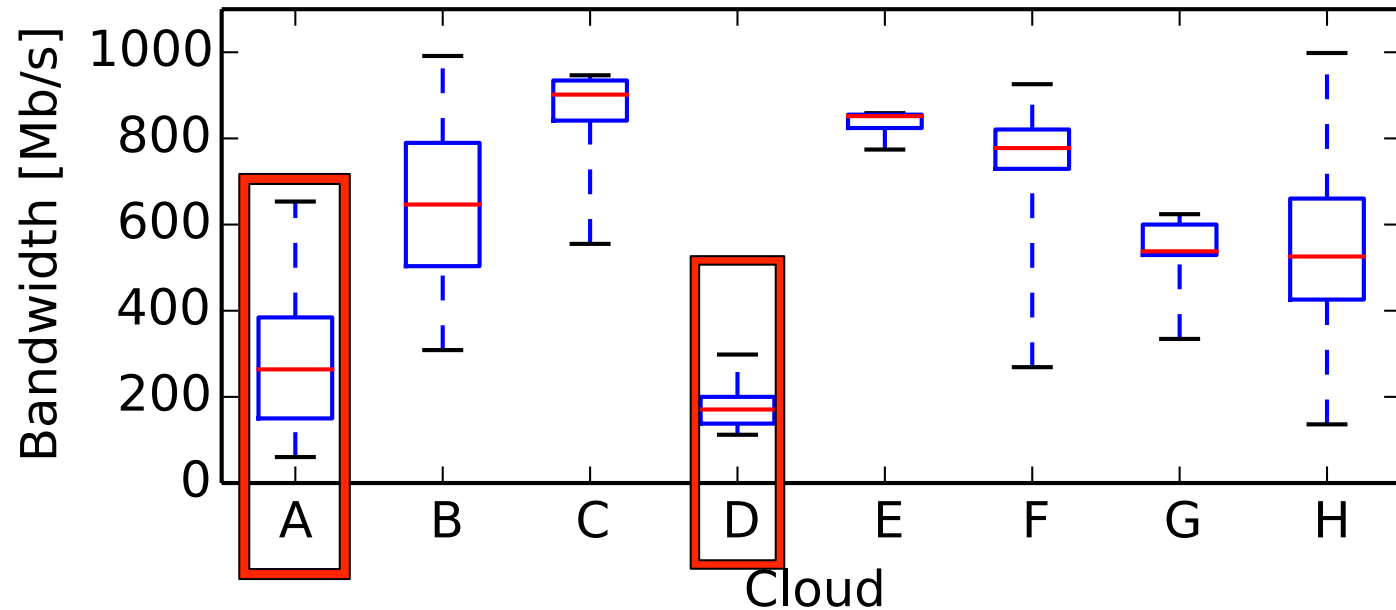
# Variable network = Variable Runtime (Terasort)

# Surprisingly, non-network-intensive Wordcount slowed down

# Most apps are slowed down on real clouds

| Application | Maximum Slowdown | Bandwidth Distribution |
|---|---|---|
| Wordcount | 1.61 | A |
| Sort | 1.51 | D |
| Terasort | 1.79 | A |
| K-Means | 1.48 | D |
| Bayes | 1.14 | A |
| Pagerank | 1.07 | A |



VRIJE
UNIVERSITEIT
AMSTERDAM

# Take-home message

- Network variability leads to high slowdown for big data in the cloud

- Network variability also affects performance portability

- Surprisingly, also apps not network-bound applications slow down

Future work:
  - In-depth statistical analysis

  - Performance modeling tools

  - Control through better scheduling

**Alexandru Uta**

a.uta@vu.nl

Vrije Universiteit Amsterdam

Massivizing Computer Systems