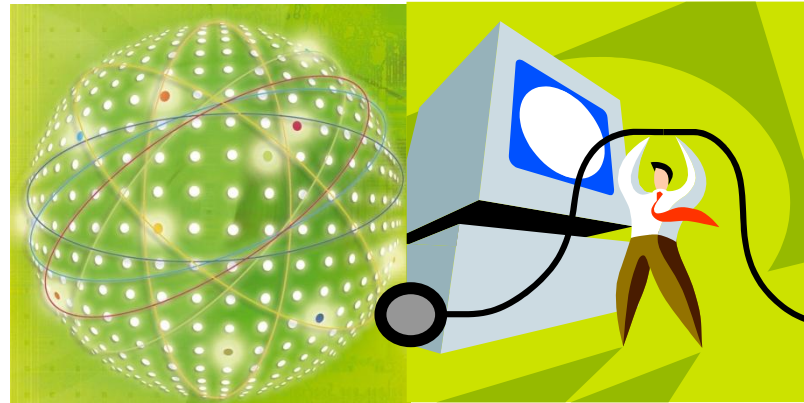# Global-Scale Applications Rely on Datacenters, Datacenters Rely on Scalable Computer Systems
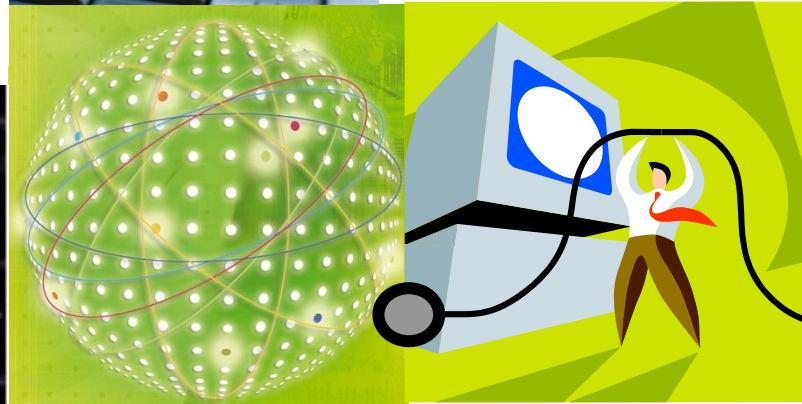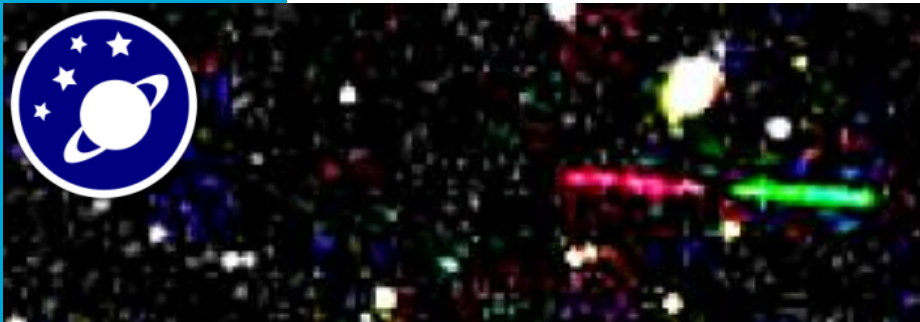
Alexandru Iosup
Parallel and Distributed Systems Group
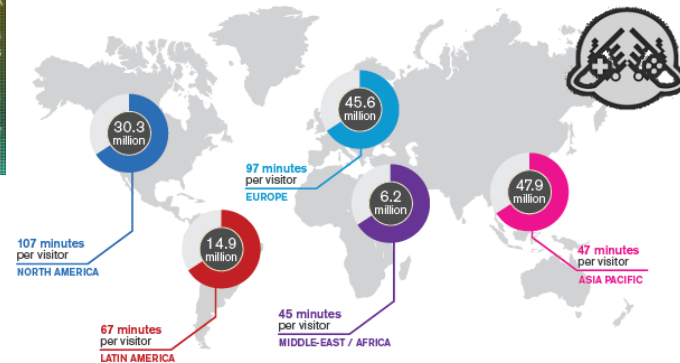
@Alosup

1

PDS Group

# This Is the Golden Age of Scalable Computing

# This Is the Golden Age of Scalable Computing

# This Is the Golden Age of Scalable Computing

# Agenda

1. The Golden Age of scalable computing
2. The core idea of cloud computing

3. Enabling technologies (homework)

4. The scheduling challenge
5. The Ecosystem Navigation challenge
6. The Big Cake challenge
7. Jevon's Effect challenge (IEEE Scale Challenge Award)

8. Take-home message

# Joe Has an Idea ($$$)

6

# Solution #1

## Buy then Maintain

- Big up-front commitment

- Load variability: NOT supported

10%

TUDelft

PDS Group

# Solution #2

## Deploy on IaaS Cloud

- **NO** big up-front commitment

- Load variability: supported

# Inside a Cloud Datacenter: Infrastructure as a Service



Q: So are we just **shifting the ownership problem**, that is, to the cloud owner?

User C

User B

MusicWave

TUDelft

PDS Group

# Agenda

1. The Golden Age of scalable computing
2. The core idea of cloud computing

3. Enabling technologies (homework)

4. The scheduling challenge
5. The Ecosystem Navigation challenge
6. The Big Cake challenge
7. Jevon's Effect challenge (IEEE Scale Challenge Award)

8. Take-home message

**TU**Delft

PDS Group

# THE PIZZA-BOX STACK
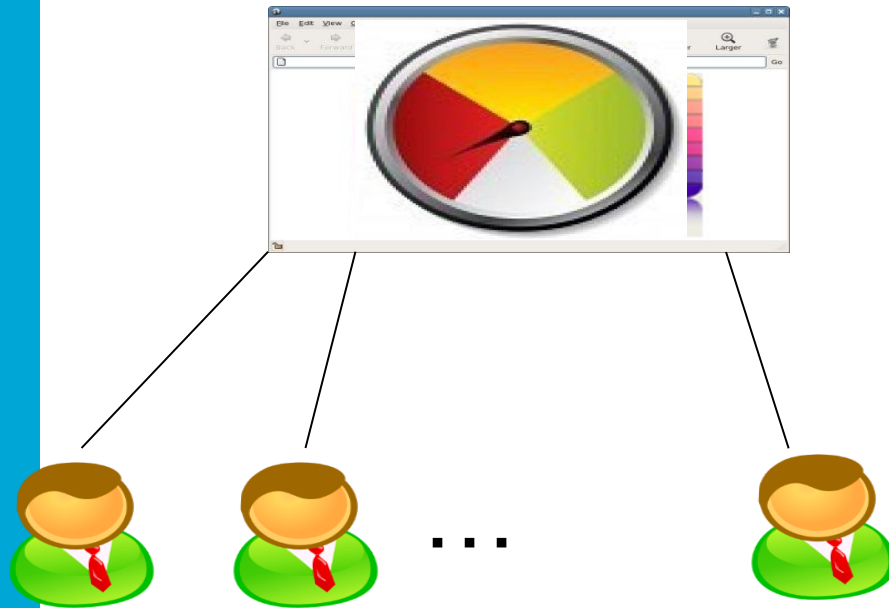
- The 1U server

Image source: http://www.avadirect.com/images/showroom/646739_1.jpg

# THE PIZZA-BOX STACK

- The 1U server

# THE PIZZA-BOX STACK

- The 1U server

13

# THE PIZZA-BOX STACK

- The 1U server

- The 19" server rack (42U is now standard)

14

# THE DATA CENTER NETWORK

- Network bandwidth per rack

  - 1 x 48-port GigE switch = 40 UP-, 8 DOWN-links



Image source: http://www.supermicro.com/a_images/products/Accessories/SSE-X3348T.gif

- Network bandwidth per socket

  - (fast) 1 Gbps for 10 GigE rack switch

  - (slow) 100 Mbps for 1 GigE rack switch

  - (exorbitant) 10 GBps for ncHT3 (supercomputing class)

Source: Dennis Abts (Google, Inc.) and John Kim (KAIST), High Performance Data Center Networks, 2011

# Resource Sharing Models

**Grids**
**Space-Sharing**

MusicWave

Q: Which one is better?

**IaaS Clouds**
**Time-Sharing**

MusicWave · OtherApp

OtherApp

MusicWave · OtherApp

Host OS

Host OS

TUDelft

PDS Group

# Virtualization

**Applications**

**Guest OS**

**Virtual Resources**

App

**Q:** What is the problem?

at to d

**VM Instance**

**Applications**

**Guest OS**

**Virtual Resources**

**VM Instance**

**Virtualization**

**Host OS**

June 4, 2015

PDS Group

# Agenda

1. The Golden Age of scalable computing
2. The core idea of cloud computing

3. Enabling technologies (homework)

4. The scheduling challenge
5. The Ecosystem Navigation challenge
6. The Big Cake challenge
7. Jevon's Effect challenge (IEEE Scale Challenge Award)

8. Take-home message

**TU**Delft

PDS Group

# The Scheduling Challenge

**Cloud operator:**

**Which resources to lease?**
**Where to place? Penalty v reward?**

**Need scheduling policies for both**
**the cloud user and the cloud operator**

**Cloud customer:**

**Which resources to lease?**
**When? How many? When stop?**
**Utility functions?**

# Portfolio Scheduling, In A Nutshell



- Create a set of scheduling policies
  - Resource provisioning and allocation policies, in this work
- Online selection of the active policy, at important moments
- Same principle for other changes: pricing model, system, …

# Portfolio Scheduling: Process



Which policies to include?

**Creation**

Which policy to activate?

**Selection**

**Reflection**

Which changes to the portfolio?

**Application**

Which resources? What to log?

# Good Results for Scientific Computing, Business-Critical, and Online Gaming Workloads



**Not performance-related, but: A portfolio scheduler can explain each decision with data.**

Q: Can our sysadmin do this? Can we? (Rhetorical)

- No single dominant policy

# Agenda

1. The Golden Age of scalable computing
2. The core idea of cloud computing

3. Enabling technologies (homework)

4. The scheduling challenge
5. The Ecosystem Navigation challenge
6. The Big Cake challenge
7. Jevon's Effect challenge (IEEE Scale Challenge Award)

8. Take-home message

**TU**Delft

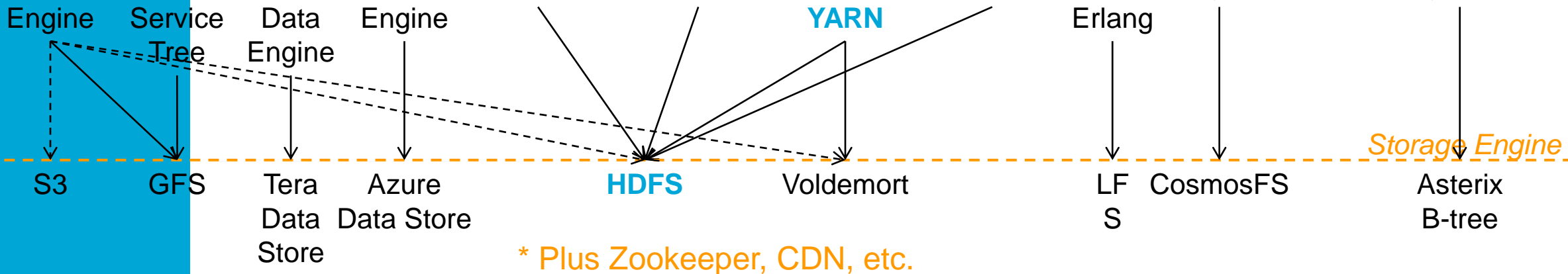PDS Group

# Big Data Processing =
## System of Systems (Ecosystem)

Flume    BigQuery    SQL    Meteor    JAQL    **Hive**    **Pig**    Sawzall    Scope    DryadLINQ    AQL

**Need to support real users who choose their tools:
batch, workflows, stream, transactions, …**

Engine    Service    Data    Engine    **YARN**    Erlang
         Tree       Engine

S3    GFS    Tera    Azure    **HDFS**    Voldemort    LF    CosmosFS    Asterix
             Data   Data Store                          S                B-tree
             Store

*Storage Engine*

* Plus Zookeeper, CDN, etc.

**TU**Delft

Adapted from: Dagstuhl Seminar on Information Management in the Cloud,
http://www.dagstuhl.de/program/calendar/partlist/?semnr=11321&SUOG

25

PDS Group
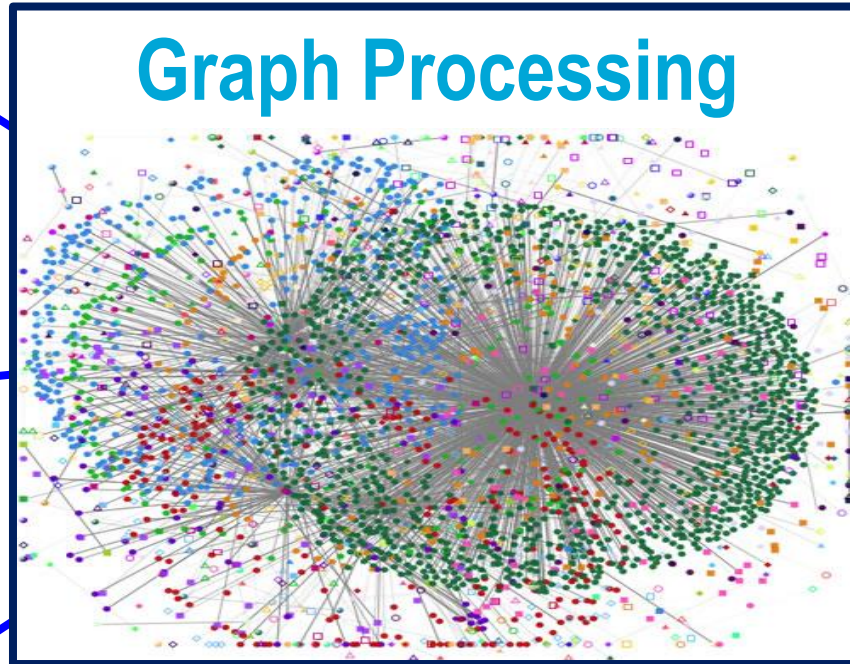
# THE DATA DELUGE: LARGE-SCALE GRAPHS
# TENS OF BILLIONS OF EDGES



**Graph Processing**

# PLATFORM DIVERSITY

Neo4j

Oracle Labs
**PGX**

YARN

APACHE GIRAPH

Spark

hadoop

GraphX

StratoSphere
Above the Clouds

GraphLab

# P-A-D TRIANGLE

Algorithm

Different algorithms
for different datasets

Performance enabled,
portability disabled

Deployment?

Dataset        No systematic findings yet        Platform

# GRAPHALYTICS: THE FIRST COMPREHENSIVE BENCHMARK FOR BIG DATA GRAPH PROCESSING

- Advanced benchmarking harness

- Choke-point analysis

- Realistic graph generator


- Co-sponsored by Oracle

- Supported by LDBC, partially developed through SPEC RG

# RUNTIME

# RUNTIME: THE PLATFORM HAS LARGE IMPACT



2 orders of magnitude difference due to platform

# RUNTIME: THE DATASET HAS LARGE IMPACT



**Neo4j fails**

Neo4j : MapReduce ~ 2:1

Neo4j : MapReduce ~ 1:2

# THROUGHPUT: THE DATASET STRUCTURE MATTERS!

20x difference

# Agenda

1. The Golden Age of scalable computing
2. The core idea of cloud computing

3. Enabling technologies (homework)

4. The scheduling challenge
5. The Ecosystem Navigation challenge
6. The Big Cake challenge
7. Jevon's Effect challenge (IEEE Scale Challenge Award)

8. Take-home message

**TU**Delft

PDS Group

# The "Big Cake" In the Datacenter

Online Social Networks

Financial Analysts

= Hadoop / MapReduce framework

**Need multi-tenant, self-aware schedulers and resource managers**

Universe Explorers

Big Data Enthusiast

Multiple frameworks = Isolation, especially performance

(Source: B. Ghit et al., SIGMETRICS 2014)

# DYNAMIC BIG DATA PROCESSING

Fawkes = Elastic MapReduce

**FAWKES**

Job submissions

Frameworks

**FAWKES/Others**

| NODES | NODES | NODES |
|-------|-------|-------|

Resource manager

Infrastructure

Ghit, Yigitbasi, Iosup, Epema, Iosup. Balanced Resource Allocations Across Multiple Dynamic MapReduce Clusters.
ACM SIGMETRICS 2014.

3

PDS Group

# ELASTICITY FOR MAPREDUCE FRAMEWORRKS

## Core nodes

**INPUT/OUTPUT DATA**

- Classical deployment
- Uniform data distribution
- **No removal**

## Transient nodes (TR)

**NO DATA**

- No local storage
- R/W from/to core nodes
- **Instant removal**

## Trans-core nodes (TC)

**OUTPUT DATA**

- Local storage, no input
- Only R from core nodes
- **Delayed removal**

TUDelft

Because workloads may be time-varying:
- Poor resource utilization
- Imbalanced service levels

**W$_1$** < **W$_2$** < **W$_3$**

1. Fair framework size:

$$s_i = \frac{w_i}{w_1 + w_2 + w_3}, \quad i = 1,2,3$$

**Core**

**TR/TC**

$w > w_{min}$

2. Updates dynamic weights when:
  * New frameworks arrive
  * Framework states change

**FAWKES**

**FAWKES**

3. Shrinks and grows frameworks to:
  * Allocate new frameworks
  * Give fair shares to existing frameworks
  * Eliminate unused frameworks

$w_{min}$

$w=0$

# PERFORMANCE OF DYNAMIC MAPREDUCE

10 core + 10xTR ▨
10 core + 10xTC ▨
vs.
20 core nodes (baseline)

**TR** - **good** for compute-intensive workloads.

**TC** - **needed** for disk-intensive workloads.

Dynamic MapReduce:
< 25% overhead

Fawkes also reduces imbalance



**data-intensive app**

< 1

**CPU-intensive app**

TUDelft

# Agenda

1. The Golden Age of scalable computing
2. The core idea of cloud computing

3. Enabling technologies (homework)

4. The scheduling challenge
5. The Ecosystem Navigation challenge
6. The Big Cake challenge
7. Jevon's Effect challenge (IEEE Scale Challenge Award)

8. Take-home message

**TU**Delft

PDS Group

**Over 500 YouTube videos have at least 100,000,000 viewers each.**

**If you want to help kill the planet:**
https://www.youtube.com/playlist?list=PLirAqAtI_h2r5g8xGajEwdXd3x1sZh8hC

**PSY Gangnam consumed ~500GWh**

**= more than entire countries\* in a year (\*41 countries),**
**= over 50MW of 24/7/365 diesel, 135M liters of oil,**
**= 100,000 cars running for a year, ...**

Source: Ian Bitterlin and Jon Summers, UoL, UK, Jul 2013.
Note: Psy has now >3 billion views (Jun 2015).

**T̃U**Delft

43

PDS Group

# The New "Jevon's Effect":
## The "Data Deluge"



**44 ZETTABYTES**
2020

2013 TOTAL

GENERATED BY CONSUMERS

NOT TOUCHED BY ENTERPRISES

2.9 ZB

0.6 ZB (15%)

2.3 ZB (85%)

4.4 ZETTABYTES

1.5 ZB

GENERATED BY ENTERPRISES

Source: IDC, 2014

**Data Deluge =**
**data generated by humans and devices (IoT)**

- Interacting
- Understanding
- Deciding
- Creating

**Need to address**
**Volume, Velocity, Variety of Big Data\***

\* New Vs later: ours is "vicissitude"

**T U** Delft

PDS Group

44

# Vs of big data

$V_1,...,V_n$

Data

- Volume – large scale of data

- Variety – different forms of data

- Velocity – timeliness of data

- Veracity – uncertainty of data

- **Vicissitude** – dynamic combination of several big data Vs in processing systems that support the addition of new queries at run-time

Vicissitude

Value
€$£¤

vicissitude *noun* [vĭˈsɪsɪˌtu()d]:
a favorable or unfavorable event or situation that occurs by chance; a fluctuation of state or condition

http://merriam-webster.com/dictionary/vicissitude

PDS Group

# OBSERVING BITTORRENT: MANAGING A TYPICAL GLOBALLY DISTRIBUTED SYSTEM



Most used protocol on Internet, by upload volume [1]
One third (US) to half (EU) of residential upload
Over 100 million users [2]

[1] https://sandvine.com/downloads/general/global-internet-phenomena/
2013/2h-2013-global-internet-phenomena-report.pdf
[2] http://www.bittorrent.com/company/about/ces_2012_150m_users

# BTWORLD: A TYPICAL BIG DATA PROJECT
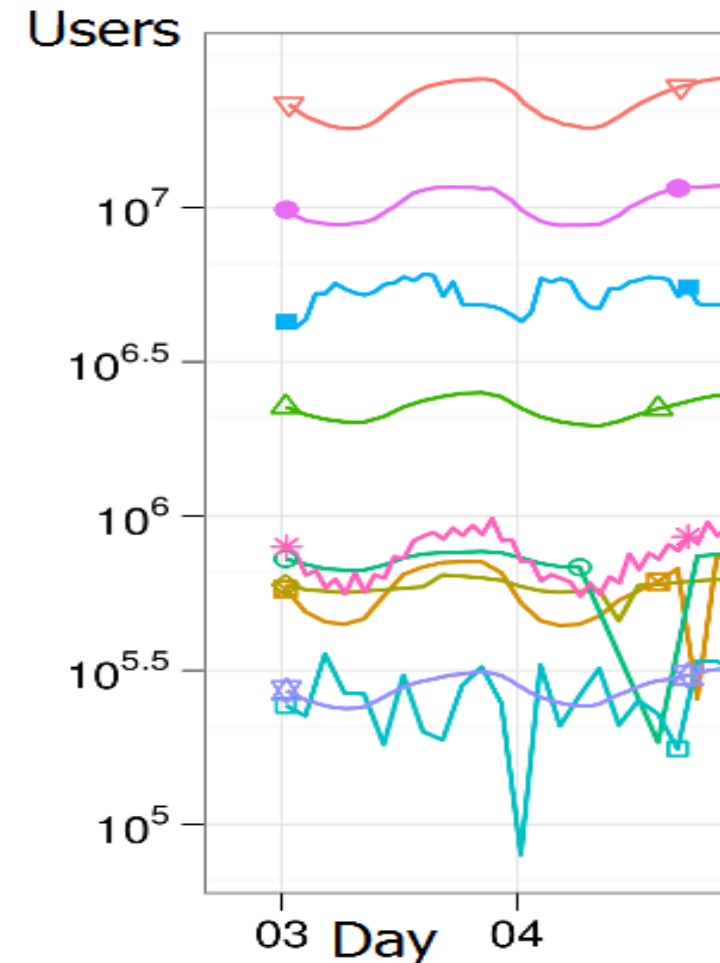
- Ongoing longitudinal study, 5 YEARS

- Data-driven project to understand BitTorrent:
  data first, ask questions later

  - Over 15 TB of structured and semi-structured data
    added during the project

  - Queries added during project, e.g.,
    How does the BitTorrent population vary?
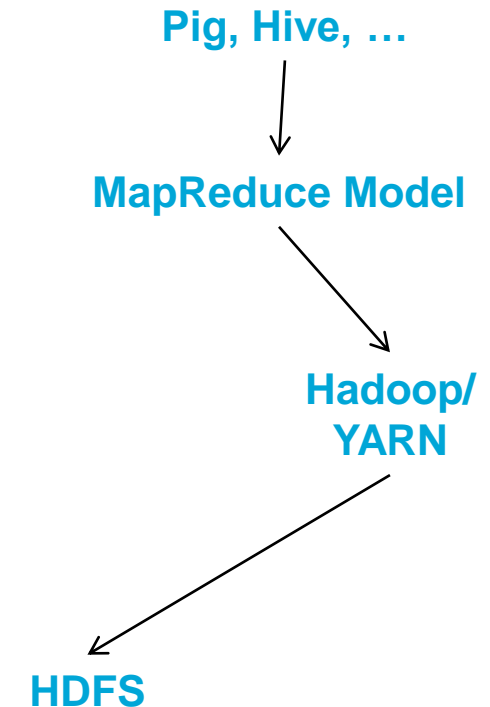    How does BitTorrent change over time?



Wojciechowski et al. Towards observing the global BitTorrent file-sharing network. HPDC 2010
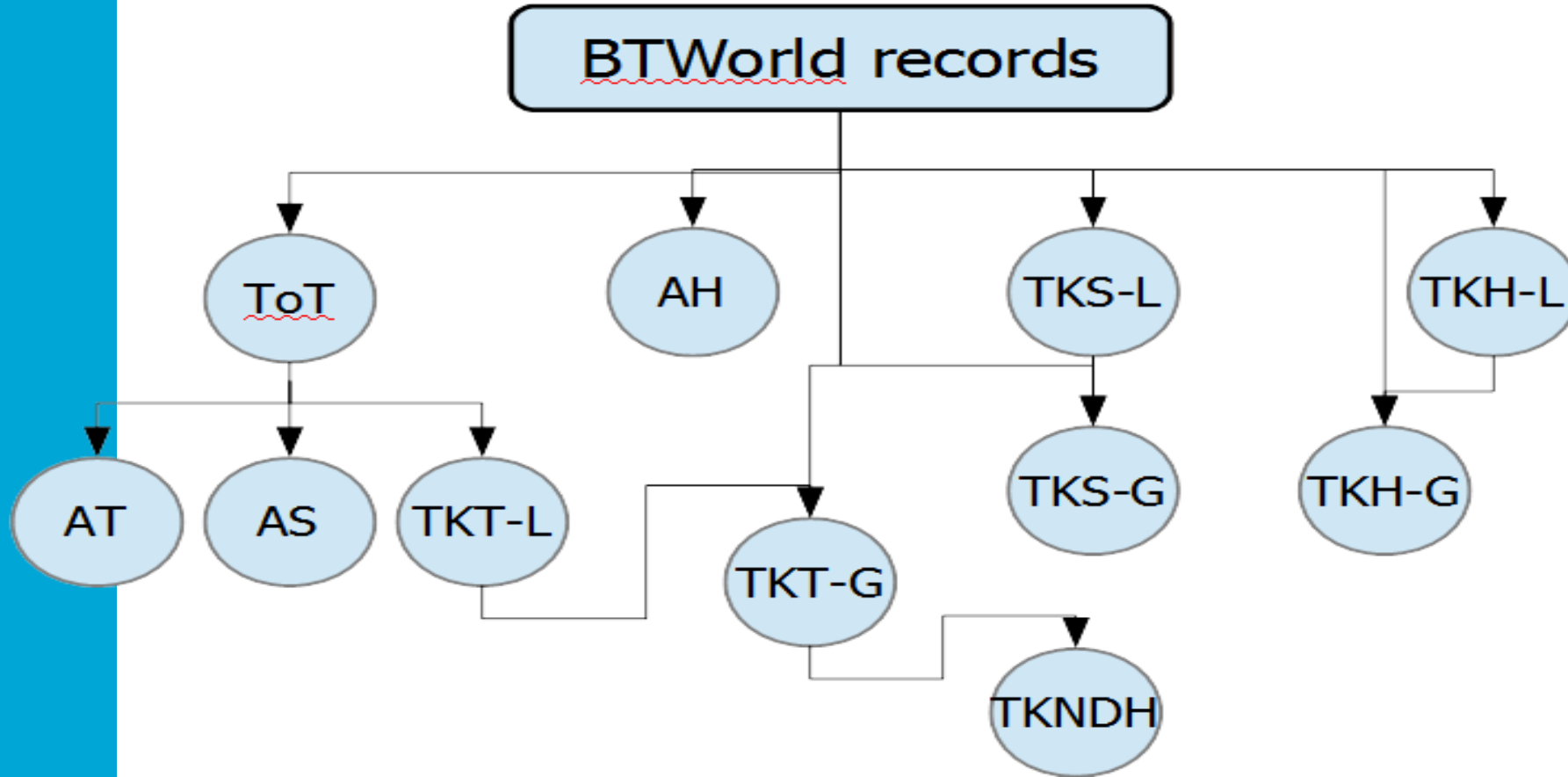
# THE MAPREDUCE ECOSYSTEM
# (A BIG PROBLEM IN BIG DATA)



**YARN**

**Pig, Hive, …**

**MapReduce Model**

**Hadoop/ YARN**

**HDFS**

- Widely used in industry and academia

  - Similar to other big data stacks

- Complex software to tune

  - 100s of parameters

  - Non-linear effects common

- Lots of issues cause crashes [1]

- Focus on Small and Medium Enterprises (60% GPD)

  - No resources or even competence to fix issues
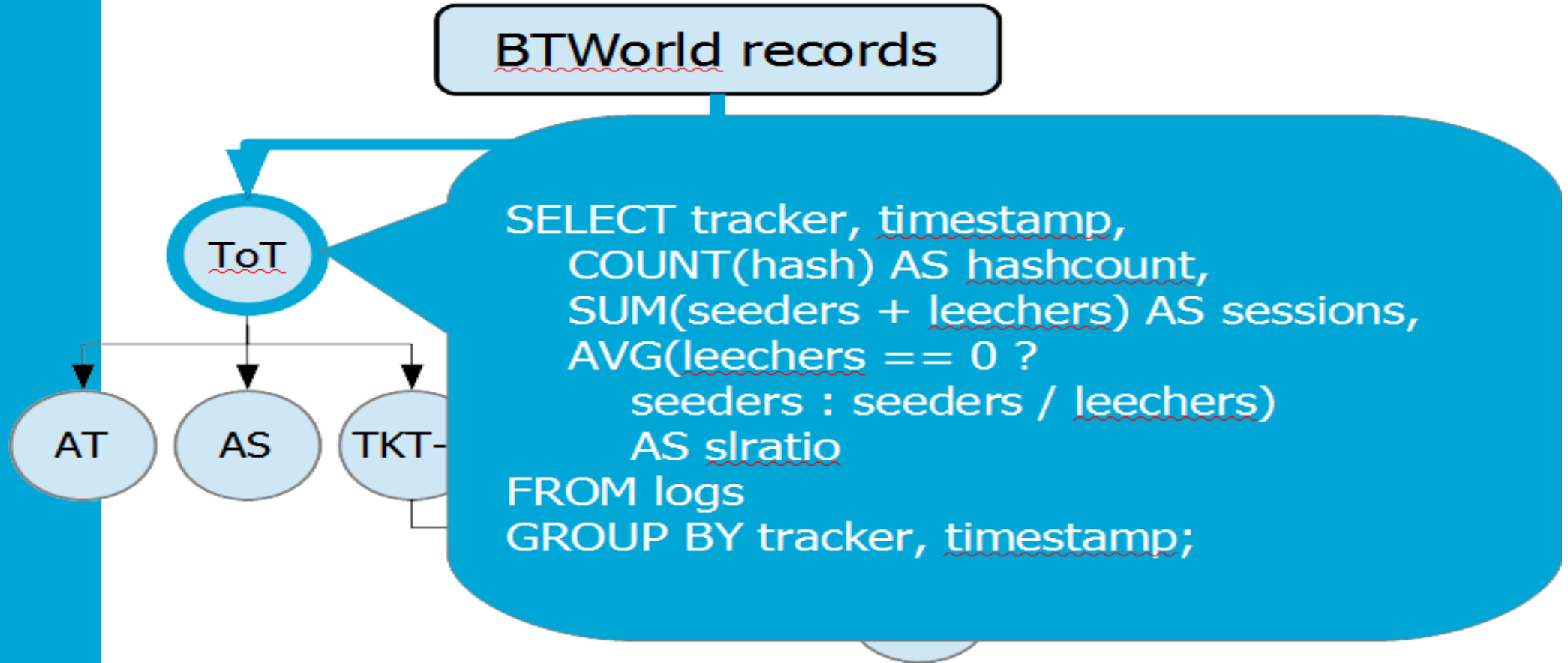
  - Difficult to make stack work for own problems

[1] Ewen et al., "Spinning Fast Iterative Data Flows", PVLDB 2012

PDS Group

# THE BTWORLD WORKFLOW

# THE BTWORLD WORKLOAD



**BTWorld records**

ToT

AT    AS    TKT-

```
SELECT tracker, timestamp,
       COUNT(hash) AS hashcount,
       SUM(seeders + leechers) AS sessions,
       AVG(leechers == 0 ?
               seeders : seeders / leechers)
               AS slratio
FROM logs
GROUP BY tracker, timestamp;
```

# THE BTWORLD WORKLOAD

# OPTIMIZATION CYCLE

Storage
HDFS

↓

Execution
MapReduce

↓

Query
Pig

↓

Application
Workflow

**Crash!**

100 GB

**Crash!**

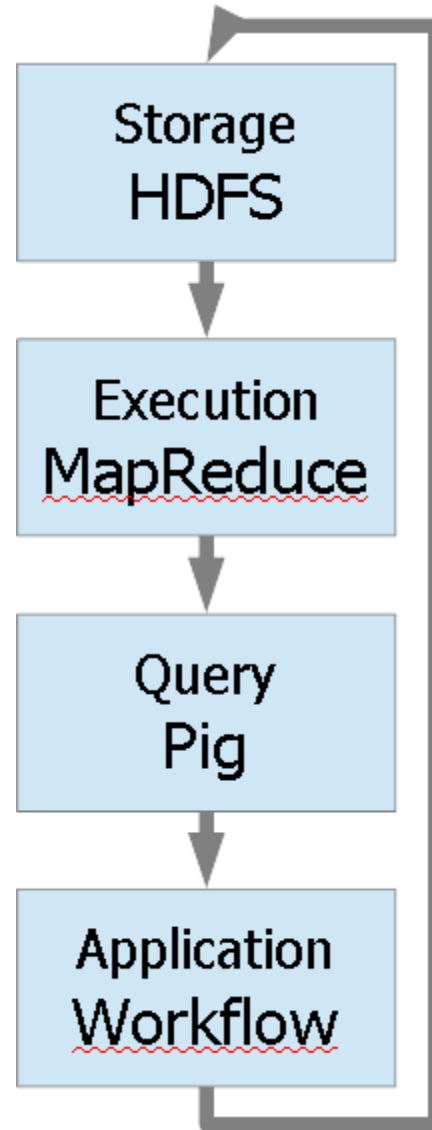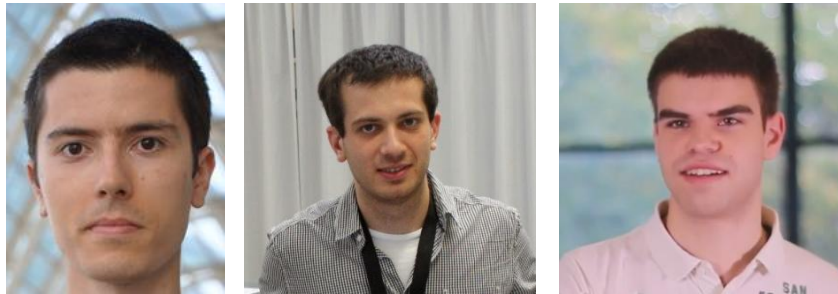1.5 TB

- HDFS: reduced replication, concatenate small files

- MapReduce: memory per task vs number of tasks, mappers then reducers

- Pig: specialized joins, multistage adaptive joins

- Workflow: reuse data between stages, common queries

# GENERAL PROBLEM

| Domain | Data Collection | Entities | Identifiers |
|---|---|---|---|
| BitTorrent | Trackers | Swarms | Hashes |
| Finance | Stock markets | Stock listings | Stocks |
| Tourism | Travel agents | Vacation packages | Venues |

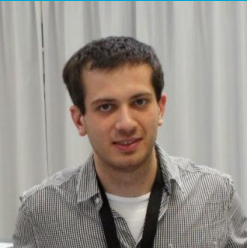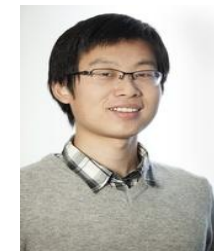Won IEEE Scale Challenge 2014!

# Agenda

**TU**Delft

# Take-Home Message

## The Golden Age of scalable computing

## Cloud computing + Big Data

## Important New Challenges

1. *The scheduling challenge*
2. *The ecosystem navigation challenge*
3. *The big cake challenge*
4. *Jevon's effect challenge*

PDS Group

# Recommended Reading

## Elastic Big Data and Computing

- B. Ghit, N. Yigitbasi (Intel Research Labs, Portland), A. Iosup, and D. Epema. Balanced Resource Allocations Across Multiple Dynamic MapReduce Clusters. SIGMETRICS 2014
- L. Fei, B. Ghit, A. Iosup, D. H. J. Epema: KOALA-C: A task allocator for integrated multicluster and multicloud environments. CLUSTER 2014: 57-65
- K. Deng, J. Song, K. Ren, A. Iosup: Exploring portfolio scheduling for long-term execution of scientific workloads in IaaS clouds. SC 2013: 55

## Time-Based Analytics

- B. Ghit, M. Capota, T. Hegeman, J. Hidders, D. Epema, and A. Iosup. V for Vicissitude: The Challenge of Scaling Complex Big Data Workflows. Winners IEEE Scale Challenge 2014

## Graph Processing

- Y. Guo, M. Biczak, A. L. Varbanescu, A. Iosup, C. Martella, T. L. Willke: How Well Do Graph-Processing Platforms Perform? An Empirical Performance Evaluation and Analysis. IPDPS 2014: 395-404
- A. L. Varbanescu, M. Verstraaten, C. de Laat, A. Penders, A. Iosup, H. J. Sips: Can Portability Improve Performance?: An Empirical Study of Parallel Graph Analytics. ICPE 2015: 277-287

**TU**Delft

PDS Group