# Balanced Resource Allocations Across Multiple Dynamic MapReduce Clusters
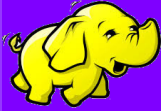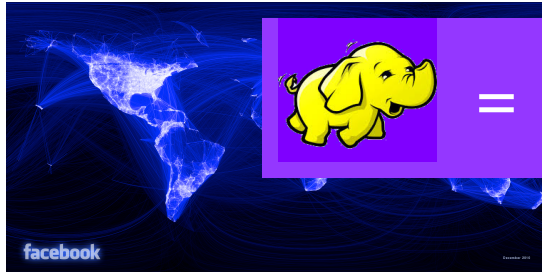
**ACM SIGMETRICS 2014**

Bogdan Ghiţ, Nezih Yigitbasi, Alexandru Iosup, Dick Epema

Parallel and Distributed Systems
Delft University of Technology
Delft, the Netherlands

TUDelft    **COMMIT/**    PDS Group    Challenge the future    DELFT DATA SCIENCE

# The "big data cake" problem

Online Social Networks

Financial Analysts

= Hadoop / MapReduce framework

Universe Explorers

Big Data Enthusiast

Multiple frameworks = Isolation, especially performance

# Our solution, FAWKES

## Two-level scheduling architecture

Job submissions

Frameworks

Resource manager

Infrastructure

**FAWKES**

**NODES** **NODES** **NODES**
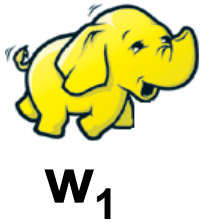
# Why dynamic provisioning?

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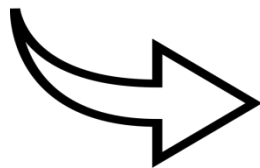  - Poor resource utilization
  - Imbalanced service levels

$$w_1 \quad < \quad w_2 \quad < \quad w_3$$

Framework size:

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

TUDelft

# Roadmap

TUDelft

Challenge the future

# Dynamic MapReduce

## MapReduce framework

- o Distributed file system
- o Execution engine
- o Data locality constraints



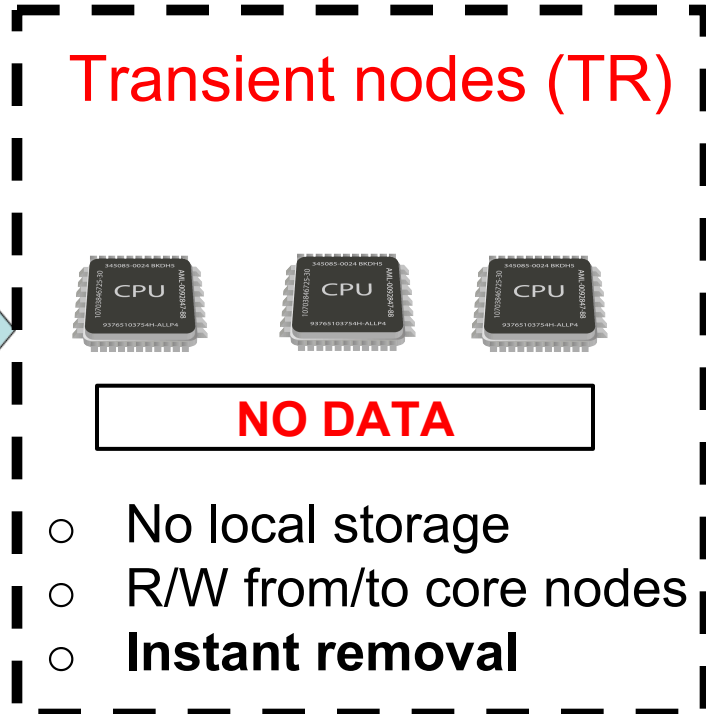Split [k1, v1]   Sort by k1   Merge [k1, [v1, v2, v3 …]]

## Grow and shrink MapReduce

- o Reliable data management
- o Fast reconfiguration
- o Break data locality

**GROW**



**SHRINK**

# No data locality



Core nodes

**IN OUT**

**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**

Performance?

# Relaxed data locality



**Core nodes**

**INPUT/OUTPUT DATA**

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

**Trans-core nodes (TC)**

**OUTPUT DATA**

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

Better performance?

TUDelft

Challenge the future

# Roadmap

TUDelft

Challenge the future

# FAWKES in a nutshell

1. Updates dynamic weights when:
   - New frameworks arrive
   - Framework states change



2. Shrinks and grows frameworks to:
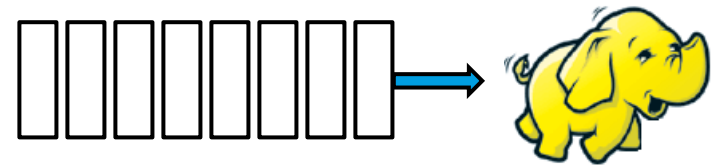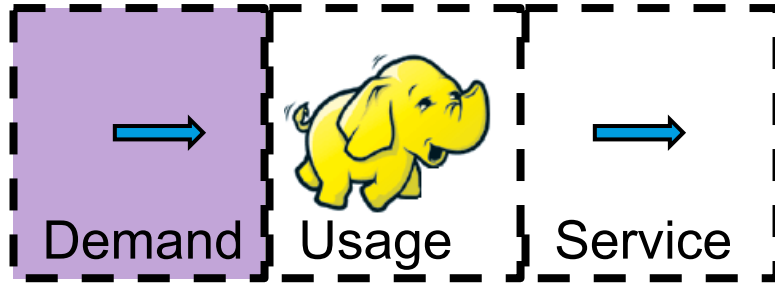   - Allocate new frameworks (min. shares)
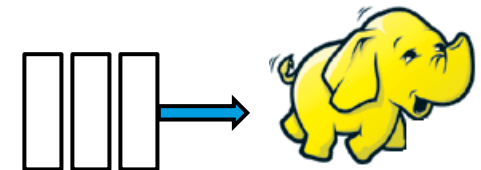   - Give fair shares to existing ones

# How to differentiate frameworks? (1/3)



By demand – 3 policies:
- o Job Demand (JD)
- o Data Demand (DD)
- o Task Demand (TD)

**VS.**

**T**U Delft

# How to differentiate frameworks? (2/3)



By usage – 3 policies:
- o Processor Usage (PU)
- o Disk Usage (DU)
- o Resource Usage (RU)

VS.

**USED**

IDLE

# How to differentiate frameworks? (3/3)



By service – 3 policies:
- o Job Slowdown (JS)
- o Job Throughput (JT)
- o Task Throughput (TT)

**VS.**

TUDelft

# Roadmap

1. Introduction
2. Dynamic MapReduce
3. FAWKES operation
4. **Experimental setup**
5. Results and analysis
6. Conclusions

# Experimental setup



Distributed ASCI Supercomputer - Version 4



## DAS-4 multicluster system:
- o 200 dual-quad-core compute nodes
- o 24 GB memory per node
- o 150 TB total storage
- o 20 Gbps InfiniBand

## Hadoop deployment:
- o Hadoop-1.0 over InfiniBand
- o 6 map + 2 reduce slots per node
- o 128 MB block size

## Overview of experiments:
- o Most experiments on 20 nodes
- o Up to 60 working nodes
- o More than 3 months system time

# MapReduce applications

| Application | Type | Input | Output |
|---|---|---|---|
| Wordcount (WC) | CPU | 200 GB | 5.5 MB |
| Sort (ST) | Disk | 200 GB | 200 GB |
| PageRank (PR) | CPU | 50 GB | 1.5 MB |
| K-Means (KM) | Both | 70 GB | 72 GB |
| TrackerOverTime (TT) | CPU | 100 GB | 3.9 MB |
| ActiveHashes (AH) | Both | 100 GB | 90 KB |
| BTWorld (BT) | Both | 100 GB | 73 GB |

Synthetic benchmarks:
- o HiBench suite
- o Single applications
- o Random datasets

Real-world applications:
- o BTWorld workflow
- o 14 Pig queries
- o BitTorrent monitoring data

TUDelft

Challenge the future

13

# Roadmap

TUDelft

Challenge the future

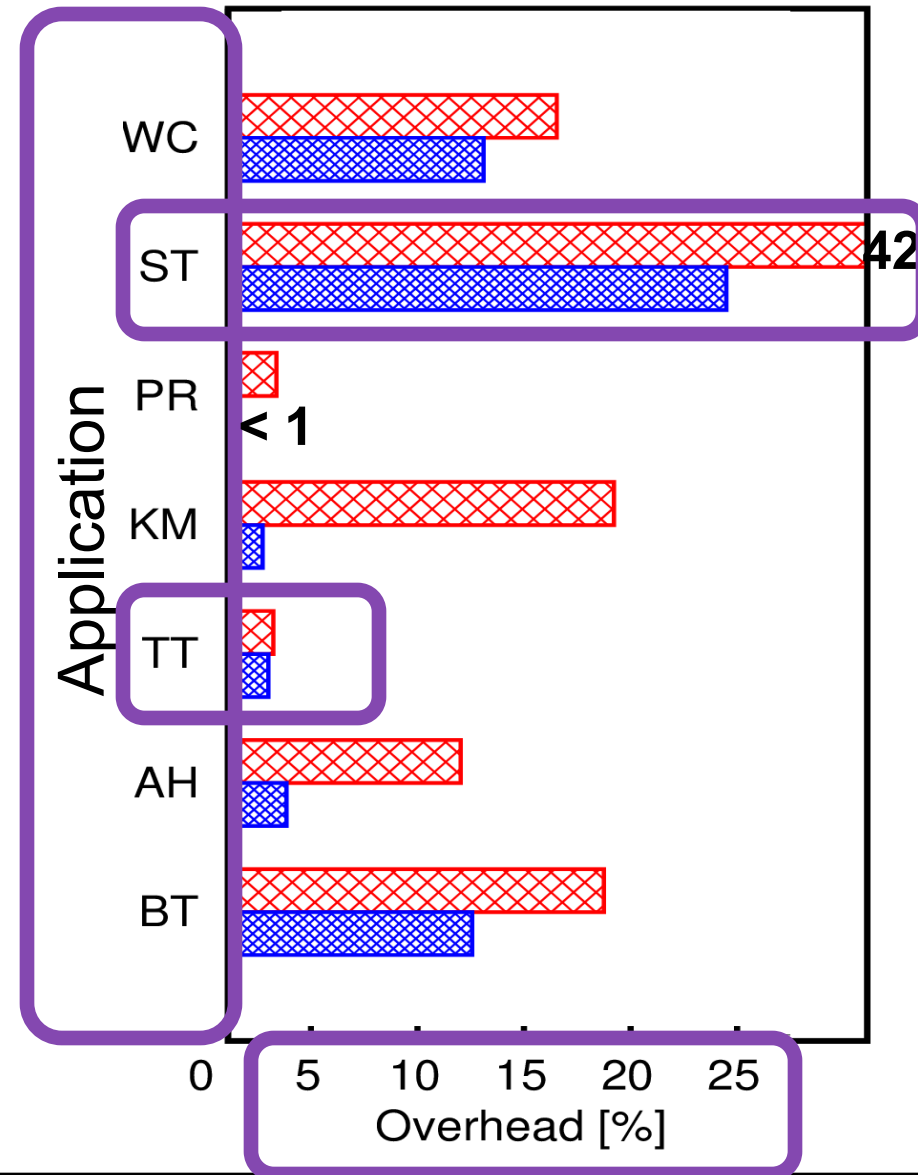# Performance of dynamic MapReduce

10 core + 10xTR ▨▨▨

10 core + 10xTC ▨▨▨

vs.

20 core nodes

**TR** - **good** for compute-intensive workloads.
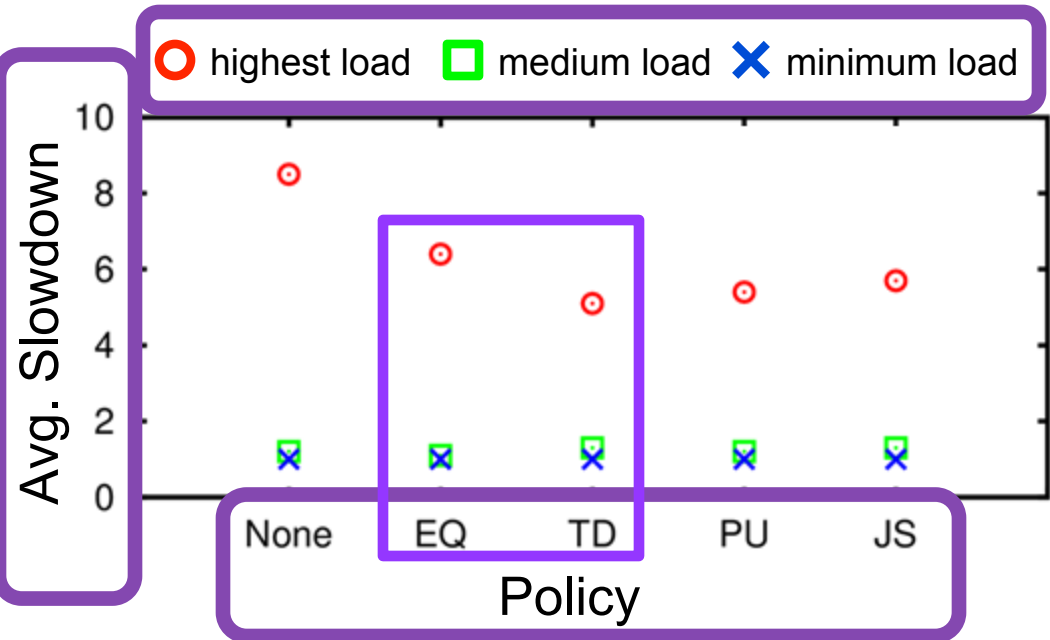
**TC** - **needed** for disk-intensive workloads.

Dynamic MapReduce:
< 25% overhead

**Application**

- WC
- ST — **42**
- PR — **< 1**
- KM
- TT
- AH
- BT

Overhead [%]: 0, 5, 10, 15, 20, 25

**T̃UDelft**

Challenge the future

14

# Performance of FAWKES

| Nodes | 45 |
|---|---|
| Frameworks | 3 |
| Min. shares | 10 |
| Datasets | 300 GB |
| Jobs submitted | 900 |



**Avg. Slowdown** vs **Policy**

○ highest load  ▢ medium load  ✕ minimum load

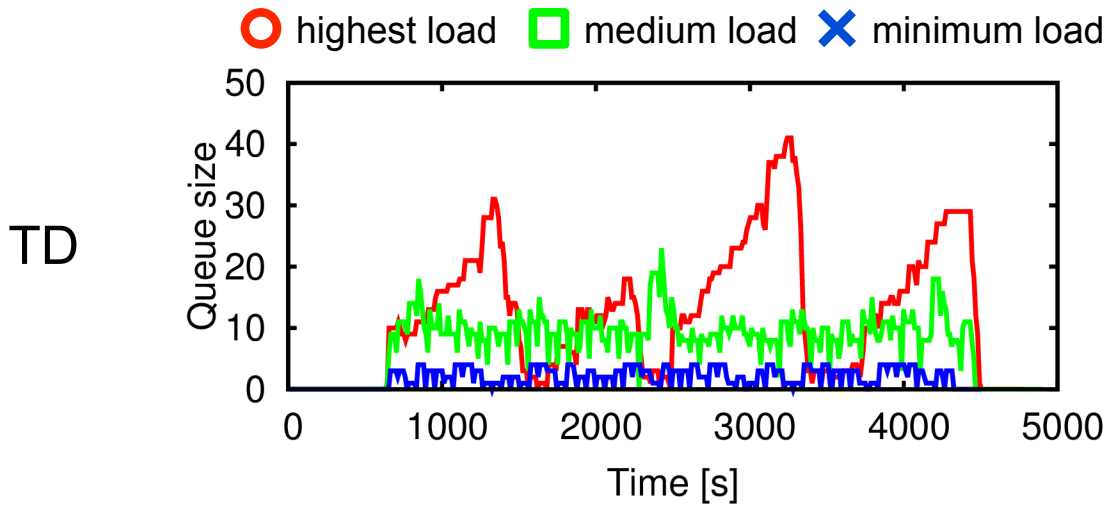Policies on x-axis: None, EQ, TD, PU, JS

Up to 20% lower slowdown

**None** – Minimum shares
**EQ** – EQual shares
**TD** – Task Demand
**PU** – Processor Usage
**JS** – Job Slowdown

# FAWKES: behind the scenes

○ highest load  ▢ medium load  ✕ minimum load

**EQ**



Utilizations:  60% / 23% / 5%

**Imbalanced**

○ highest load  ▢ medium load  ✕ minimum load

**TD**
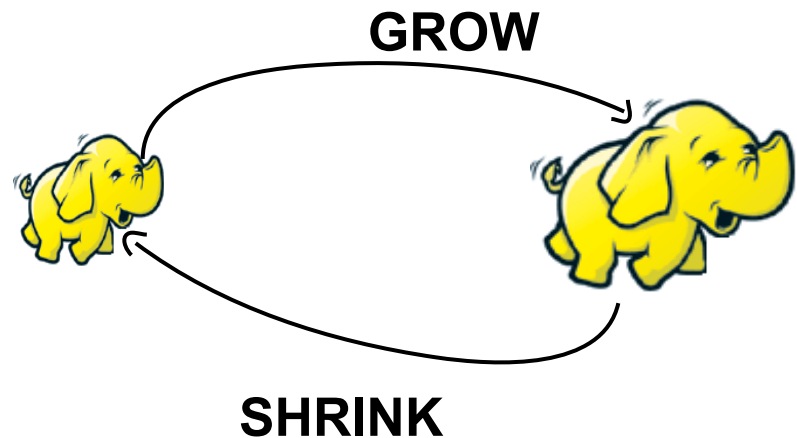


Utilizations: 50% / 30% / 8%

**More balanced**

# Roadmap

1. Introduction
2. Dynamic MapReduce
3. FAWKES operation
4. Experimental setup
5. Results and analysis
6. **Conclusions**

# Take-home message

1. Dynamic MapReduce relaxes data locality

2. FAWKES reduces the imbalance between frameworks

3. More aggressive policies?

# Our PDS group

*Scheduling and resource management research:*
- <u>Systems</u>: multi-cluster systems and clouds
- <u>Applications</u>: workflows, bags-of-tasks, data-intensive, etc.



THE Koala GRID SCHEDULER

MapReduce Resource Cloud Graph Big computing Data Bags-of-tasks Scheduling Workflows Experimentation processing Simulation Management

Contact: B.I.Ghit@tudelft.nl

http://www.pds.ewi.tudelft.nl/ghit/
http://www.pds.ewi.tudelft.nl/research-publications/publications/

# Backup slides

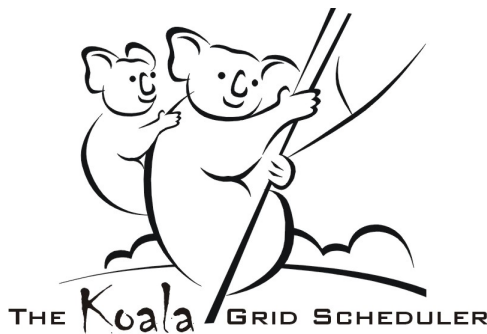# Related work


**YARN**

o   Resource requests from applications
o   Capacity and Fair schedulers

FAWKES uses feedback from system operation


MESOS

o   Resource offers to frameworks
o   Optimizes for data locality

FAWKES schedules frameworks automatically
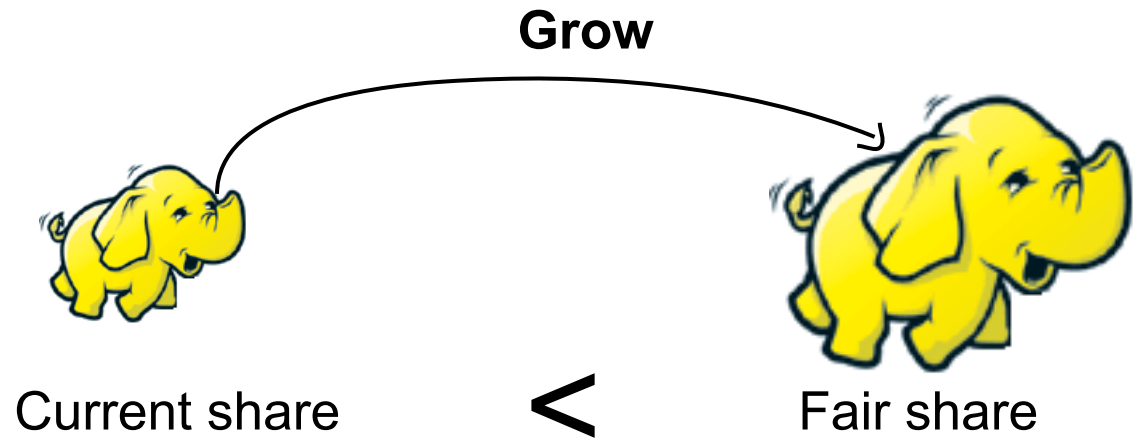

THE Koala GRID SCHEDULER

o   Grid and cloud scheduler @ TU Delft
o   Single applications and frameworks

FAWKES is a research prototype

# The grow-shrink mechanism

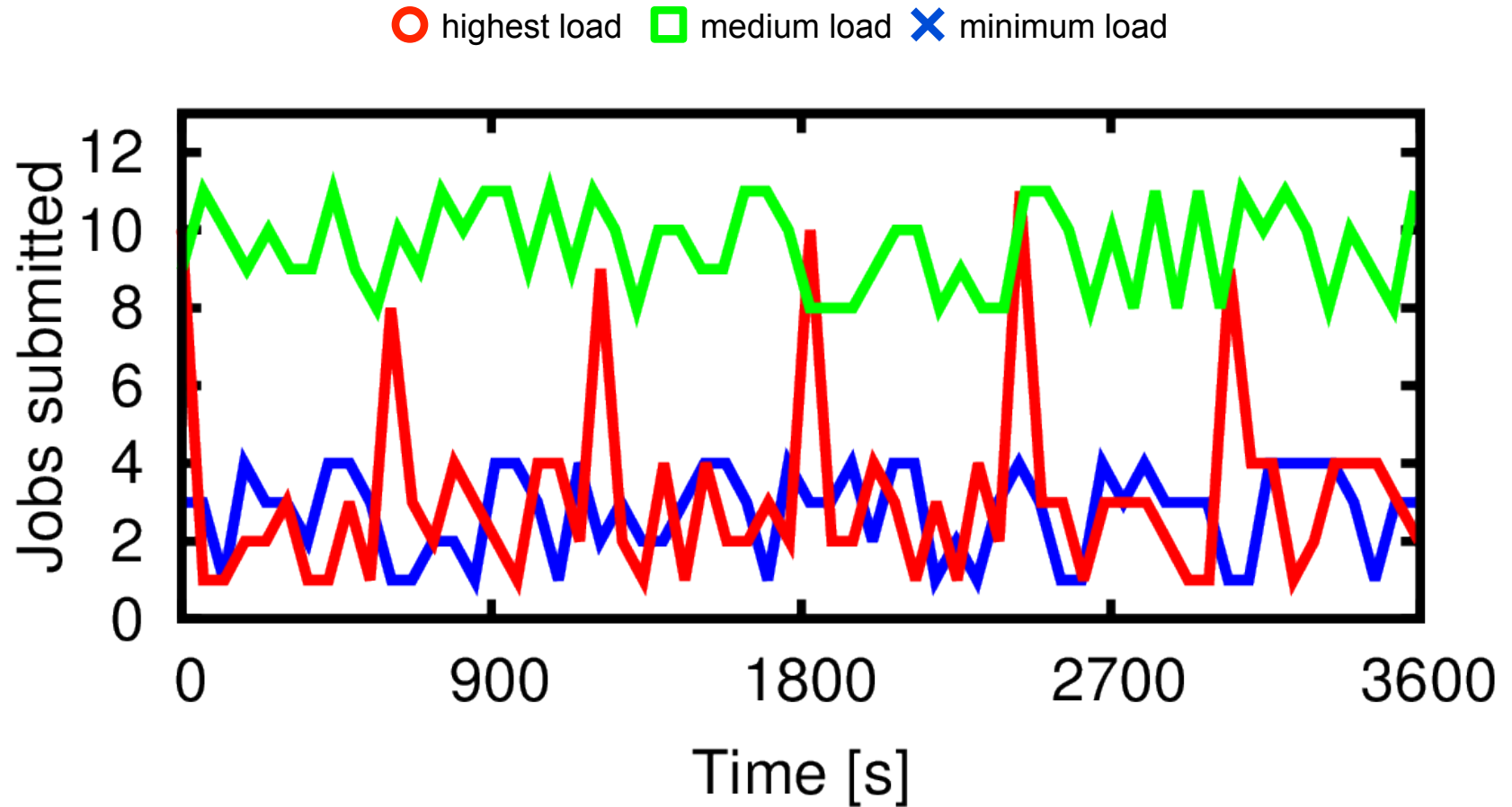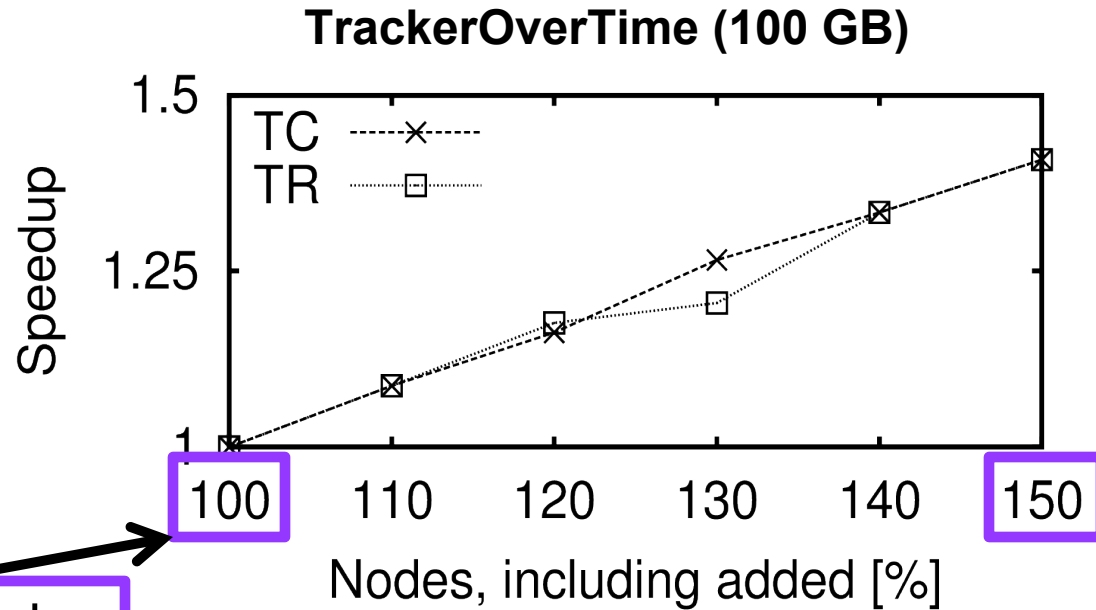**Negative discrimination**

**Grow**



Current share  <  Fair share

**Positive discrimination**

**Shrink**



Current share  <  Fair share

# Submission patterns

# Speedup when growing (1/2)

**TrackerOverTime (100 GB)**



20 core nodes

30 nodes

TR nodes deliver good performance for CPU bound workloads

Challenge the future

*TU*Delft

# Speedup when growing (2/2)
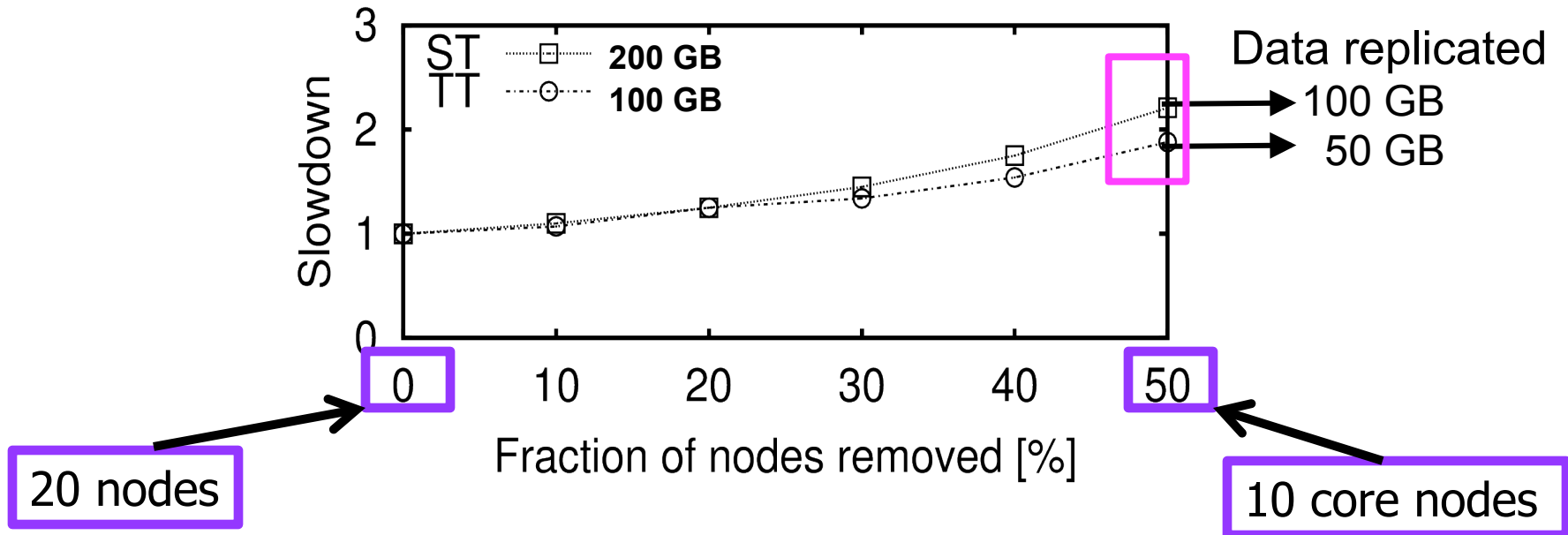
**Sort (200 GB)**



20 core nodes

30 nodes

(Only) TC nodes deliver good performance for disk-bound workloads

# Slowdown when shrinking



Job slowdown increases linearly with the amount of replicated data

TUDelft