

Self-* Datacenter Management for Business Critical Workloads



Alexandru Iosup
Delft University of Technology
The Netherlands

Team: **Undergrad** Tim Hegeman, ... **Grad** Vincent van Beek, Lipu Fei, Yong Guo, Mihai Capota, Bogdan Ghit **Researchers** Marcin Biczak, Otto Visser **Staff** Henk Sips, Dick Epema **Collaborators*** Ana Lucia Varbanescu (UvA, Ams), Claudio Martella (VU, Giraph), Intel Research Labs, IBM TJ Watson, Oracle Research Labs SF, ...

1

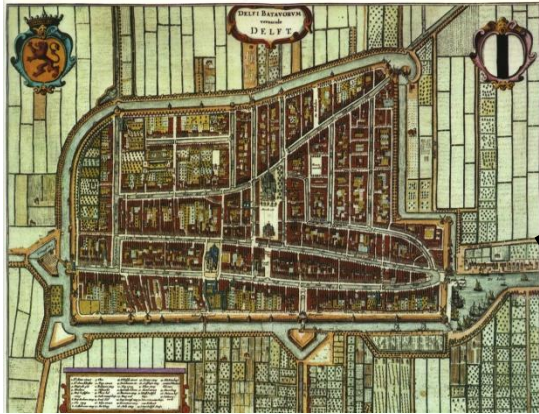
* Not their fault for any mistakes in this presentation. Or so they wish.

January 19, 2015

Dagstuhl seminar Model-driven Algorithms and
Architectures for Self-Aware Computing Systems, Aug 2014



(TU) Delft – the Netherlands – Europe



founded 13th century
pop: 100,000



founded 1842
pop: 13,000



pop: 16.5 M



The Parallel and Distributed Systems Group at TU Delft



VENI

Alexandru Iosup

Grids/Clouds
P2P systems
Big Data
Online gaming

Home page

- www.pds.ewi.tudelft.nl

Publications

- see PDS publication database at publications.st.ewi.tu.nl



Dick Epema

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



VENI

Ana Lucia Varbanescu
(now UvA)
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



August 31, 2011

Winners IEEE TCSC Scale Challenge 2014

Lessons From Grids

From Hypothesis to Data

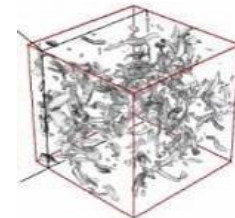


The Fourth Paradigm is suitable for professionals who already know they don't know [enough to formulate good hypotheses], yet need to deliver quickly

ena

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$

- Last few decades:
 - a **computational** branch simulating complex phenomena
- Today (**the Fourth Paradigm**):
 - data exploration**
 - unify theory, experiment, and simulation
 - Data captured by instruments or generated by simulator
 - Processed by software
 - Information/Knowledge stored in computer
 - Scientist analyzes results using data management and statistics



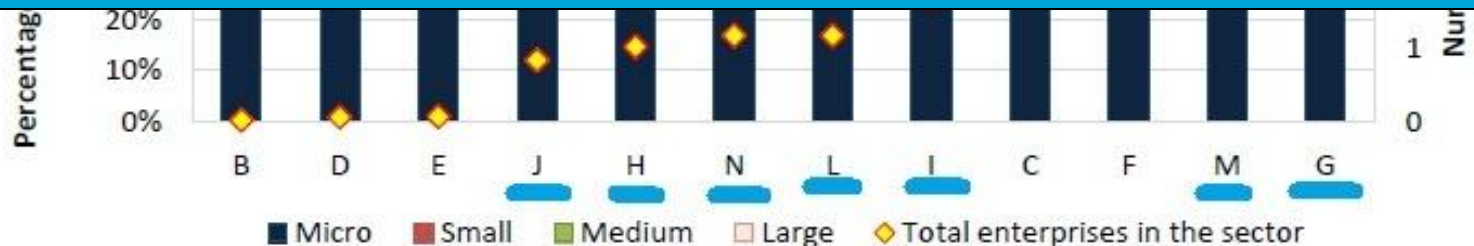
The Vision: Everyone Is a Scientist! (the Fourth Paradigm)



- Data as individual right, enabling high-quality **lifestyle of individuals** and modern **societal services**
- Data as workhorse in creating **commercial services** by SMEs (~60% gross value added, for many years)



Address ICT challenges! (EU)
>500 million people
>85 million employees
>3 trillion euros / year gross value added



What is Cloud Computing?

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

The Energy Ceiling (Can We Afford this Vision?)

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=PLirAqAtl_h2r5g8xGajEwdXd3x1sZh8hC

**PSY Gangnam, this version consumed >300GWh
= more than some countries in a year,
= over 35MW of 24/7/365 diesel, 100M liters of oil,
= 80,000 cars running for a year, ...**

Source: Ian Bitterlin and Jon Summers, UoL, UK, Jul 2013.
Adapted (Sep 2014).

Scheduling in IaaS Clouds

An Overview



Cloud operator:

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

**Need usage and user-aware
scheduling policies**



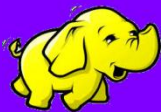
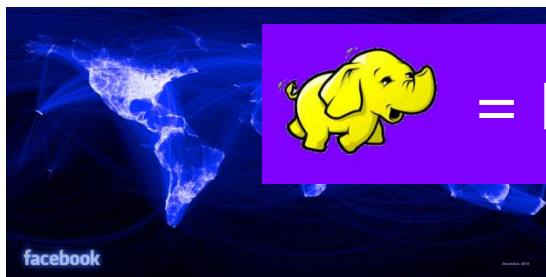
Cloud customer:

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



The "Big Data cake" in the Data Center

Online Social Networks



= Hadoop / MapReduce framework

Financial Analysts



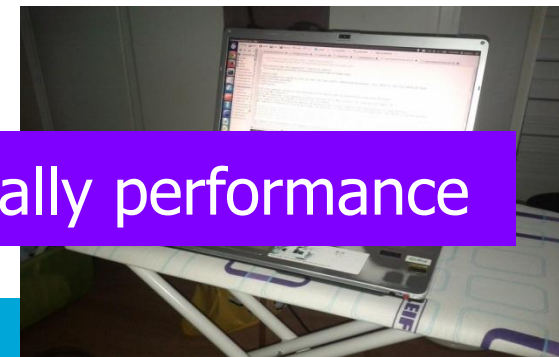
Need multi-tenant, self-metering schedulers and resource managers

Universe Explorers



Multiple frameworks = Isolation, especially performance

Big Data Enthusiast



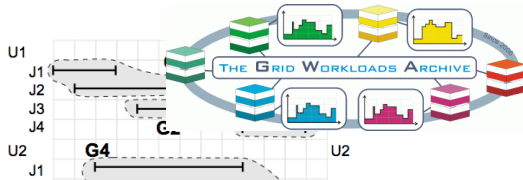
Agenda



Everyone is a Scientist



Can we afford it?



Workloads



Scheduling

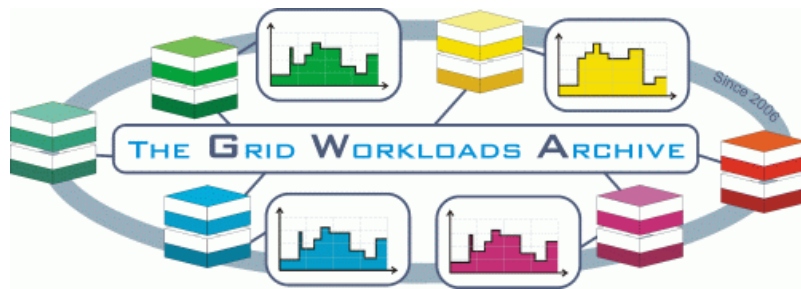
Scheduling

Scheduling



Conclusion

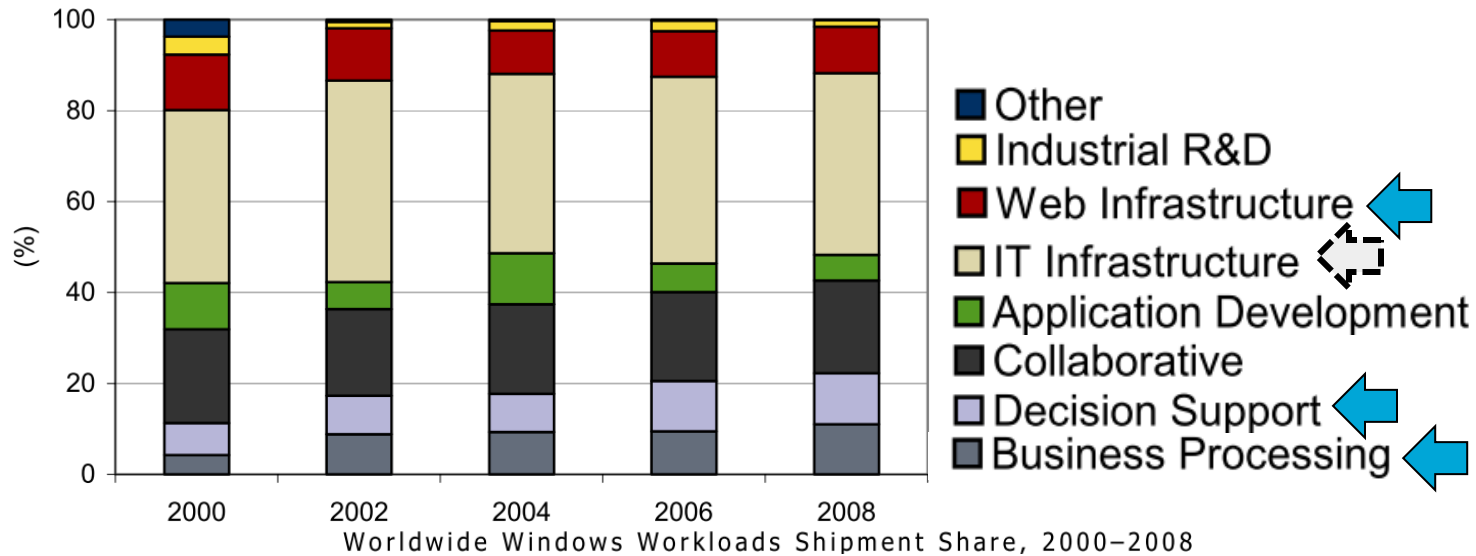
Workloads



Business-Critical Workloads

Business-Critical Workloads

- Growing user base (IDC'10)



- Applications
 - Business intelligence and decision
 - Office back-end
 - Other Business support



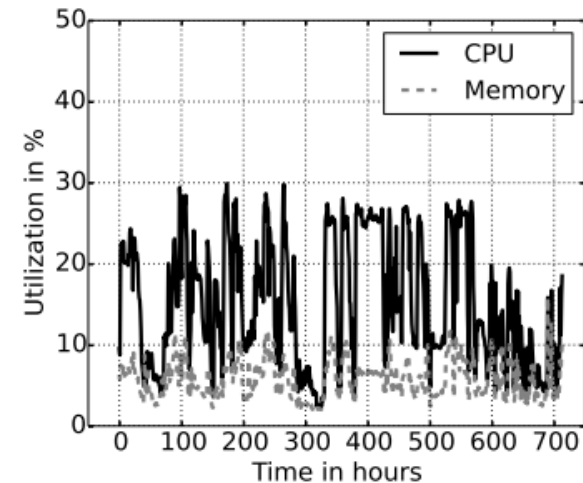
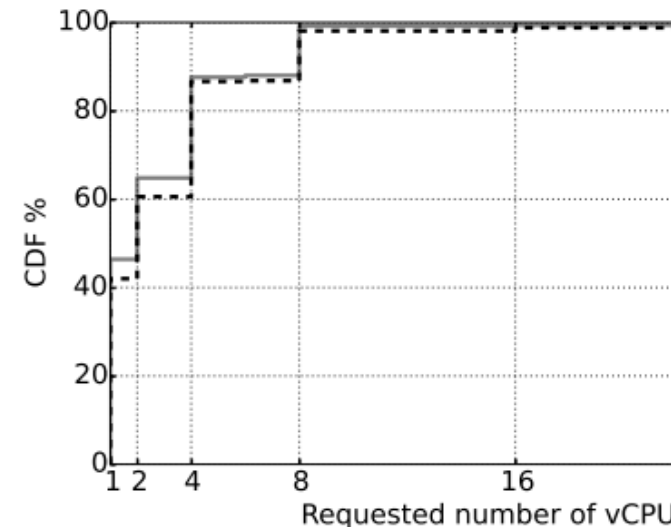
Workload Characterization

Requested resources				Used resources			
CPU	Mem	Disk	Net	CPU	Mem	Disk	Net
—	—	—	—	yes	—	yes	—
yes	yes	—	—	yes	yes	—	—
—	—	—	—	yes	—	yes	—
—	—	—	—	yes	yes	—	—
—	—	—	—	yes	yes	yes	yes
—	yes	—	—	yes	yes	—	—
yes	yes	yes	yes	yes	yes	yes	yes

- Basic statistics
- Correlation analysis
- Time pattern analysis

Workload Characterization Results

1. More than 60% of the VMs use less than 4 cores and 8GB of memory.
2. There is a strong positive correlation between requested CPU and memory.
3. Resource usage is low, under 10% of the requested resources, and the correlation between requested and used resources is also low.
4. Peak workloads can be 10–10,000 times higher than mean workloads, depending on resource type.
5. The CPU and memory resource usage is often predictable over the short-term. Disk and network I/O follow daily patterns.



Scheduling



Portfolio Scheduling



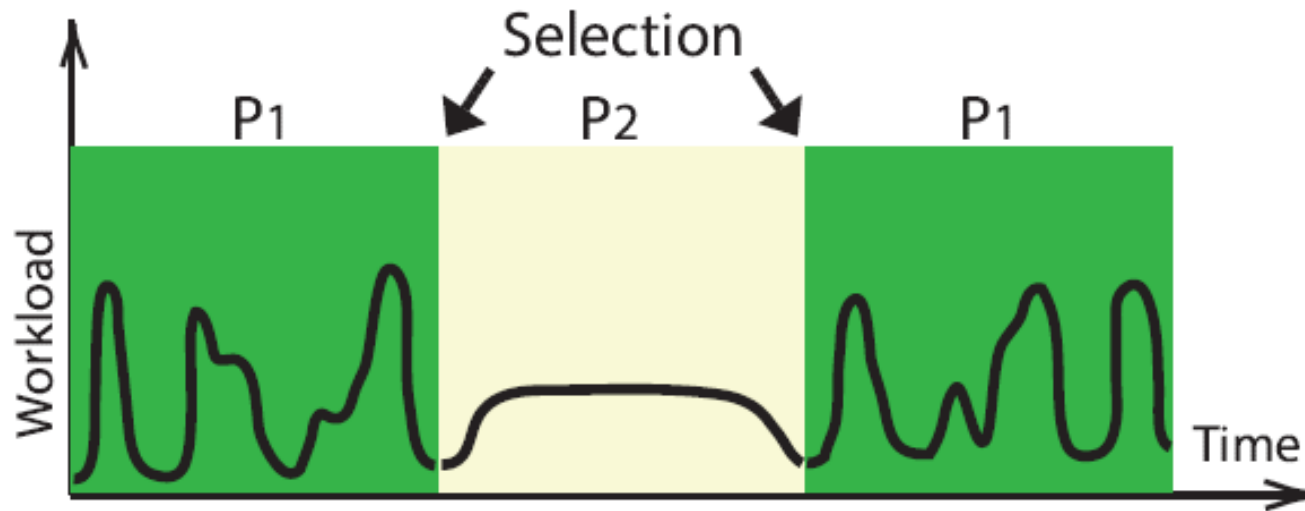
Why Portfolio Scheduling?

- **Data centers increasingly popular**
 - Constant deployment since mid-1990s
 - Users moving their computation to IaaS clouds
 - Consolidation efforts in mid- and large-scale companies
- **Old scheduling aspects**
 - Hundreds of approaches, each targeting specific conditions— which?
 - No one-size-fits-all policy
- **New scheduling aspects**
 - New workloads
 - New data center architectures
 - New cost models
- **Developing a scheduling policy is risky and ephemeral**
- **Selecting a scheduling policy for your data center is difficult**



What is Portfolio Scheduling?

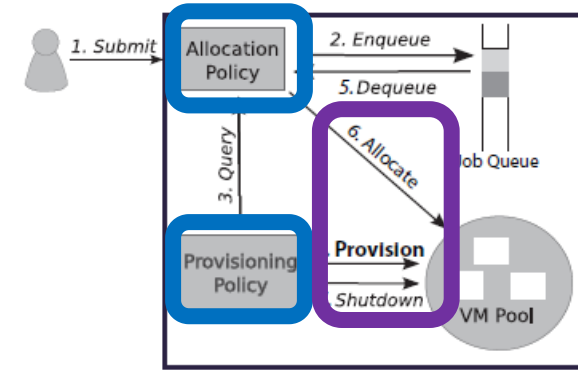
In a Nutshell, for Data Centers



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



Portfolio Scheduling The Process



Which policies to include?

Creation

Reflection

Which policy to activate?

Selection

Application

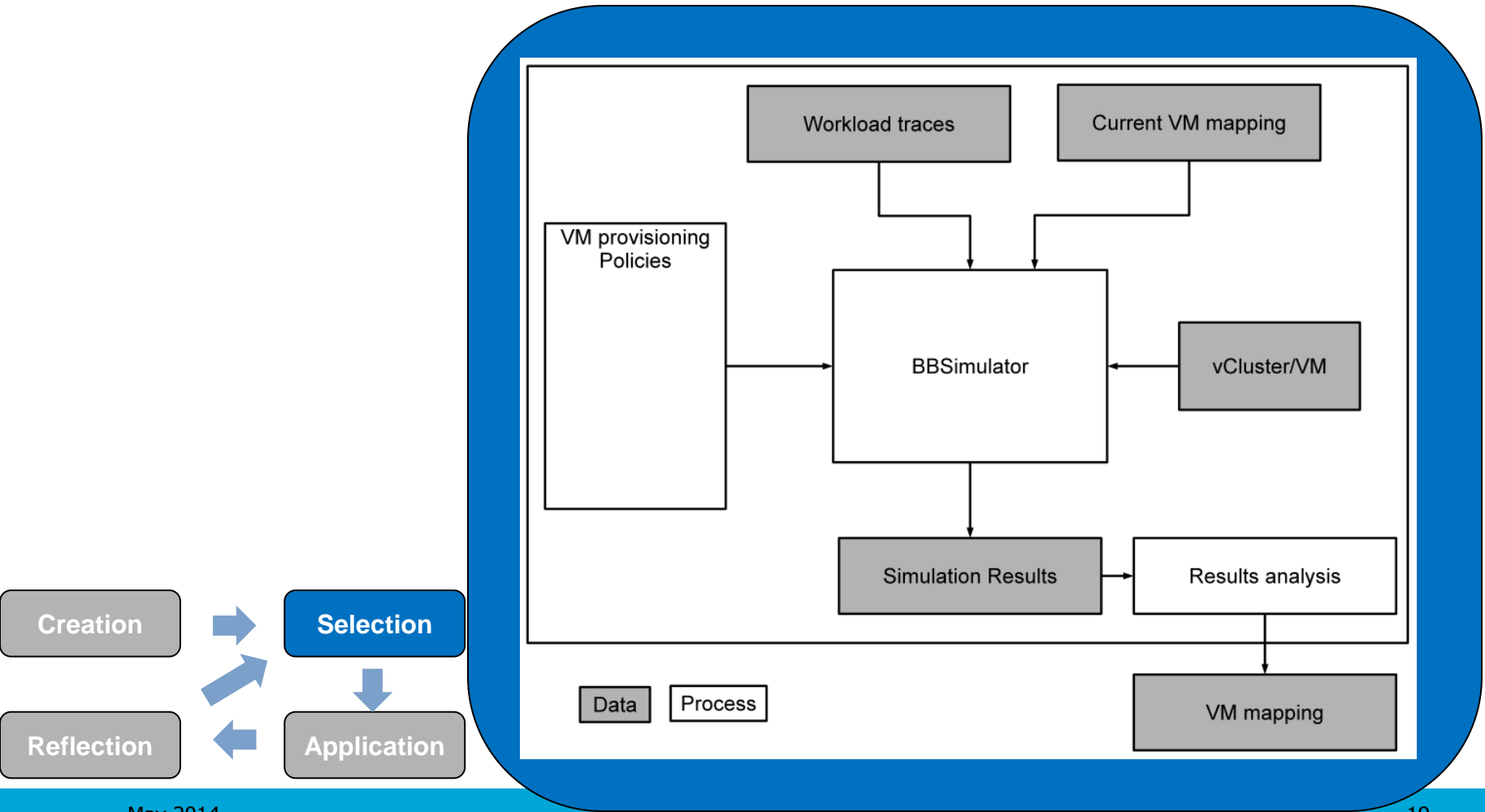
Which changes to the portfolio?

Which resources? What to log?

Bonus: The portfolio scheduler can explain each selection, like your engineers would.

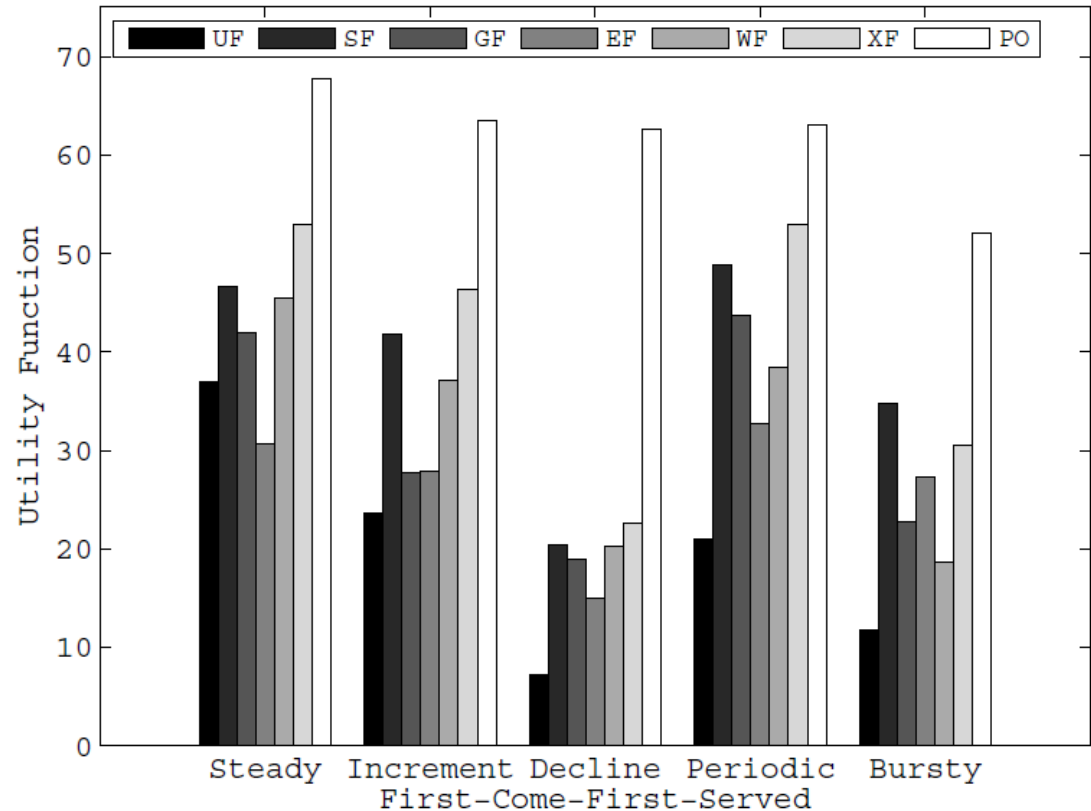
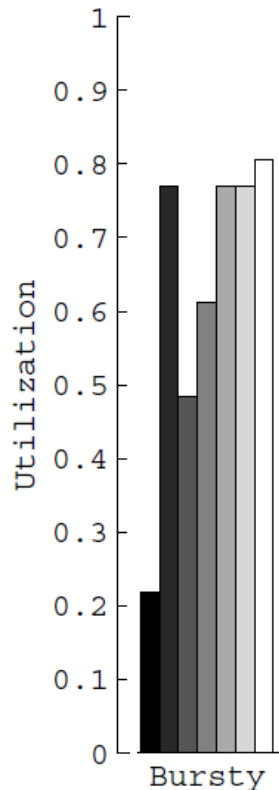
Portfolio Scheduling

An Implementation of the Selection Step



Experimental Results, Synthetic Workloads

Resource Utilization + Workload Utility



- POrtfolio leads to high utilization
- Start-Up leads to poor utilization

- POrtfolio leads to better utility
- Start-Up leads to poor utility

Experimental Results: Real-World Business-Critical Workloads at Bitbrains

- Additional requirements: affinity and anti-affinity, worst-fit/best-fit
- Risk score = Risk(CPU/mem/netw/IO)
 - Considers oversubscription, latency
 - Combines severity and probability
 - Lower is better

What the portfolio optimizes for (same set of policies)

Configurations:	replay	MinScore	MaxScore	MinMem	MaxMem	MinCPU	MaxCPU
Memory	1.68	0.31	2.18	6.16	1.56	0.26	15.09
CPU	0.0	0.0	0.0	0.0	0.0	0.0	0.0
IO write	2.6	1.53	2.99	1.49	1.69	1.35	1.93
IO read	2.31	1.25	2.15	1.2	1.23	2.05	1.41
Network send	1.86	1.18	2.51	1.13	1.32	0.9	1.86
Network receive	1.17	1.13	2.19	1.09	1.1	1.11	1.49
Total:	9.62	5.4	12.01	11.07	6.9	5.66	21.78

- Portfolio scheduler can lead to lower risk in the datacenter
- Policy selection very important, otherwise portfolio can perform badly

Scheduling

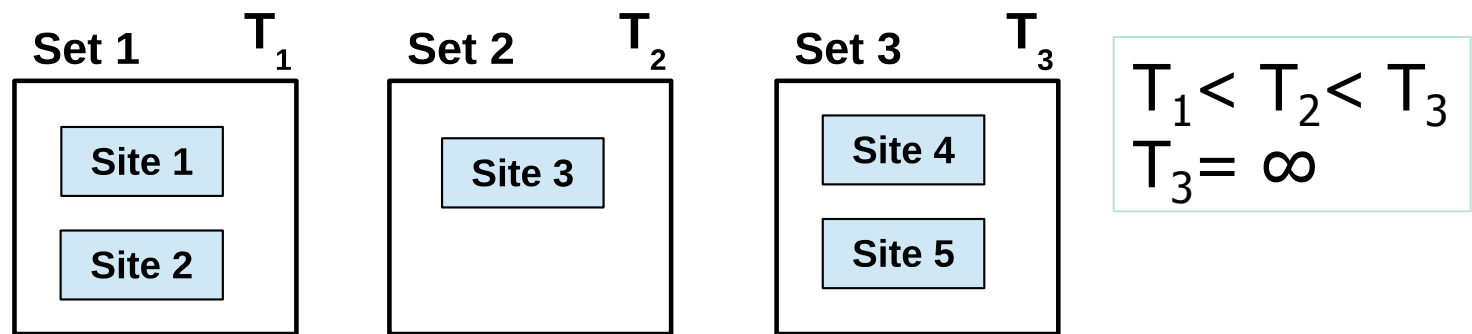


**Multi-Cluster,
Multi-DC TAGS**

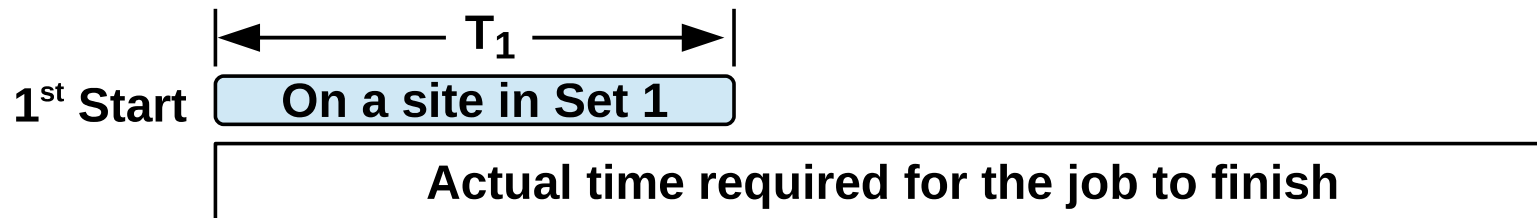


Mor Harchol-Balter's Task Assignment by Guessing Size TAGS-based Policy Design

- Goal: to achieve low slowdown **without prediction**
- Method: to partition the sites into sets to serve jobs of different runtime ranges
- A number of sets of sites
- Set i allows jobs to run for T_i amount of time ($T_i < T_{i+1}$)
- The last set has a T of ∞ (all jobs will finish without being killed)

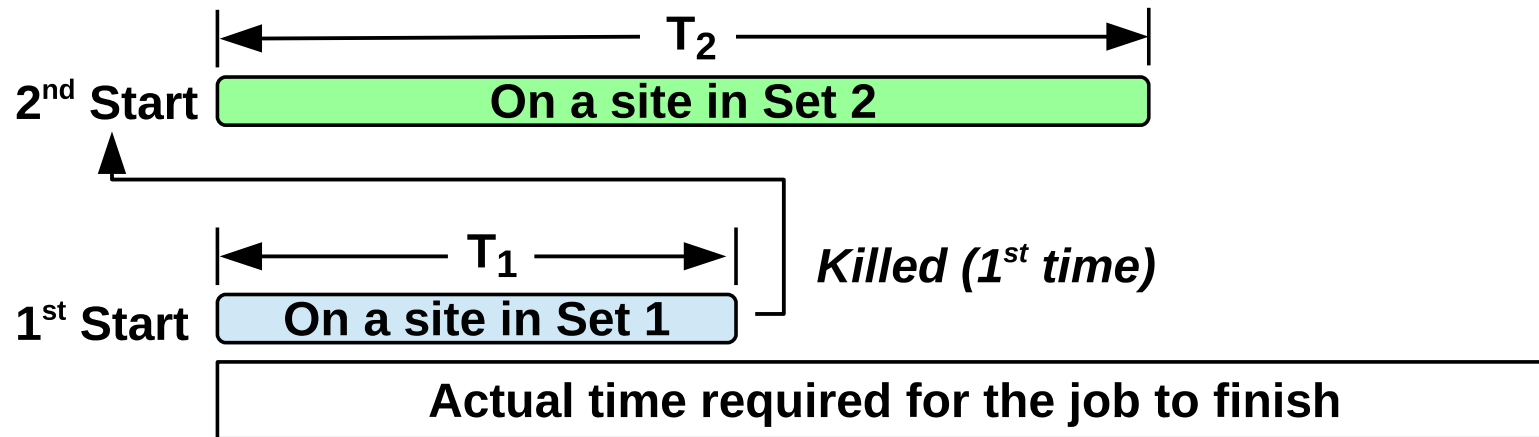


TAGS-based Policy Design



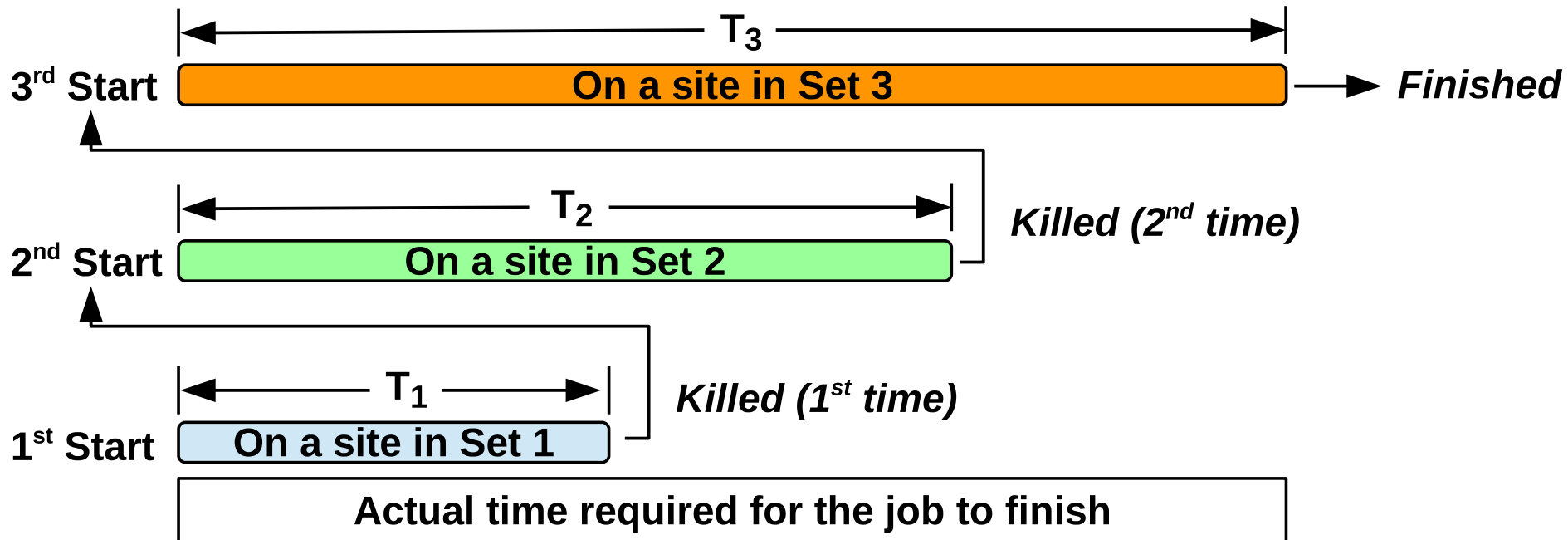
Policy Design

TAGS-based Policies in General



Policy Design

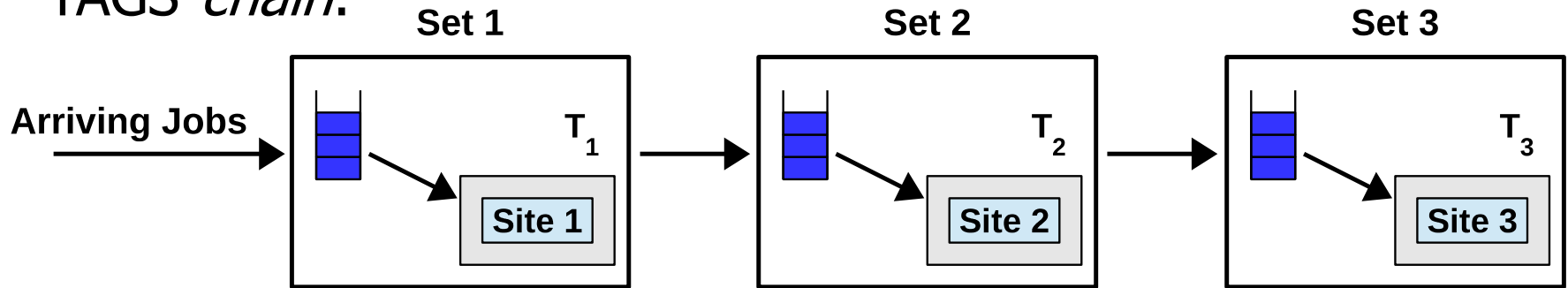
TAGS-based Policies in General



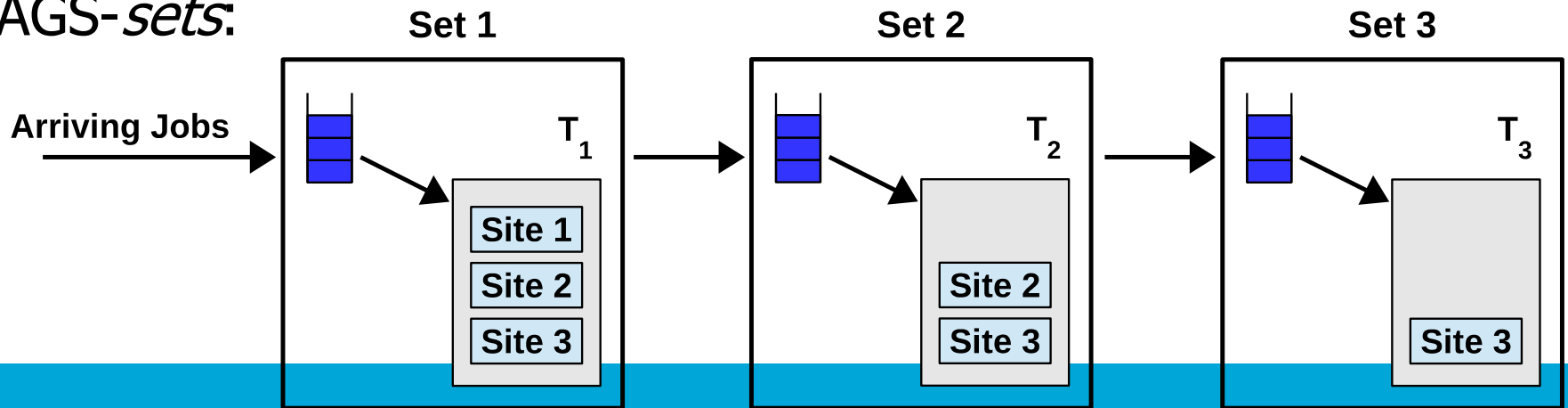
Our Policy Design for Multi-Cluster Multi-DCs

TAGS-*chain* and TAGS-*sets*

TAGS-*chain*:



TAGS-*sets*:

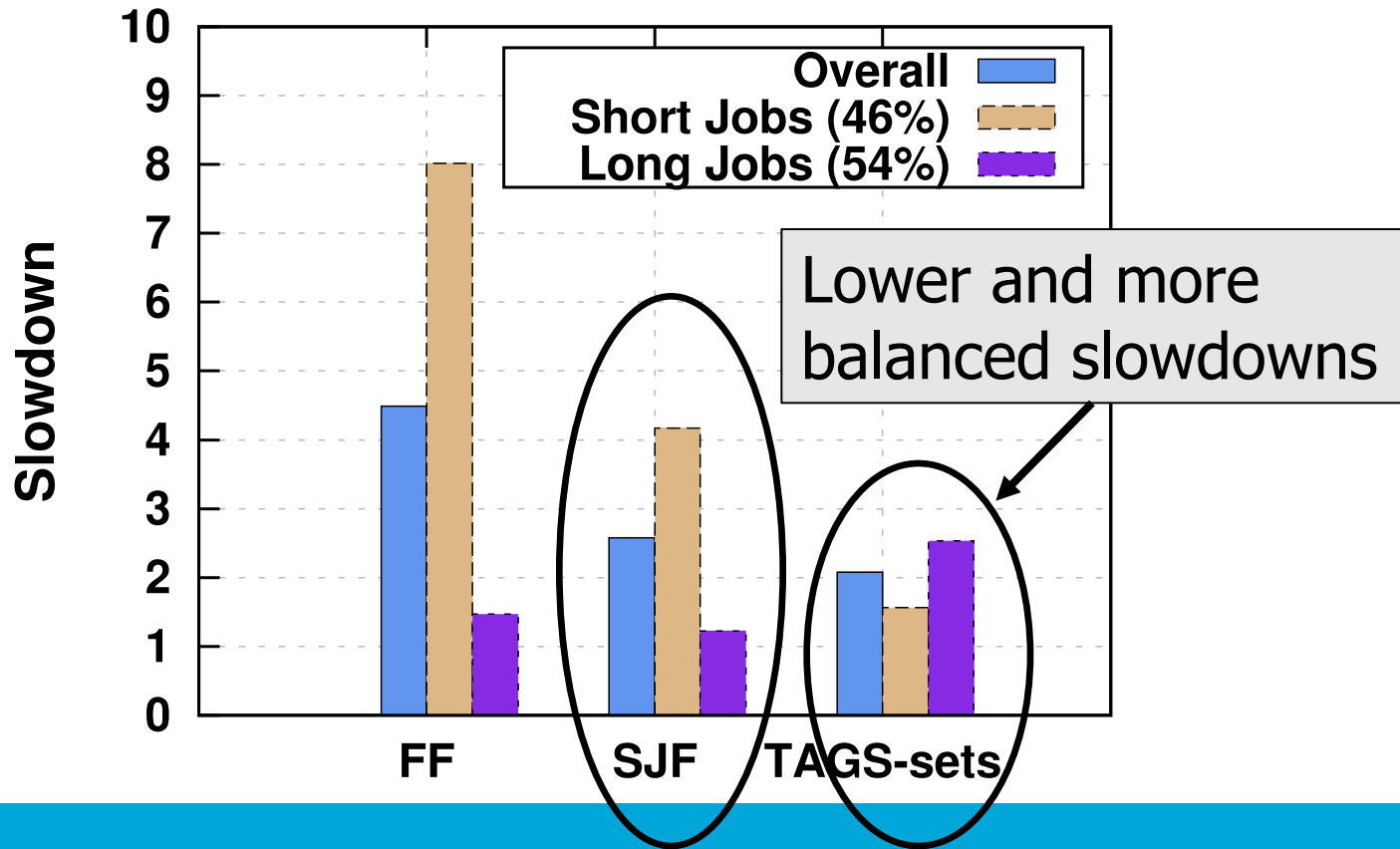


Policy Design Space

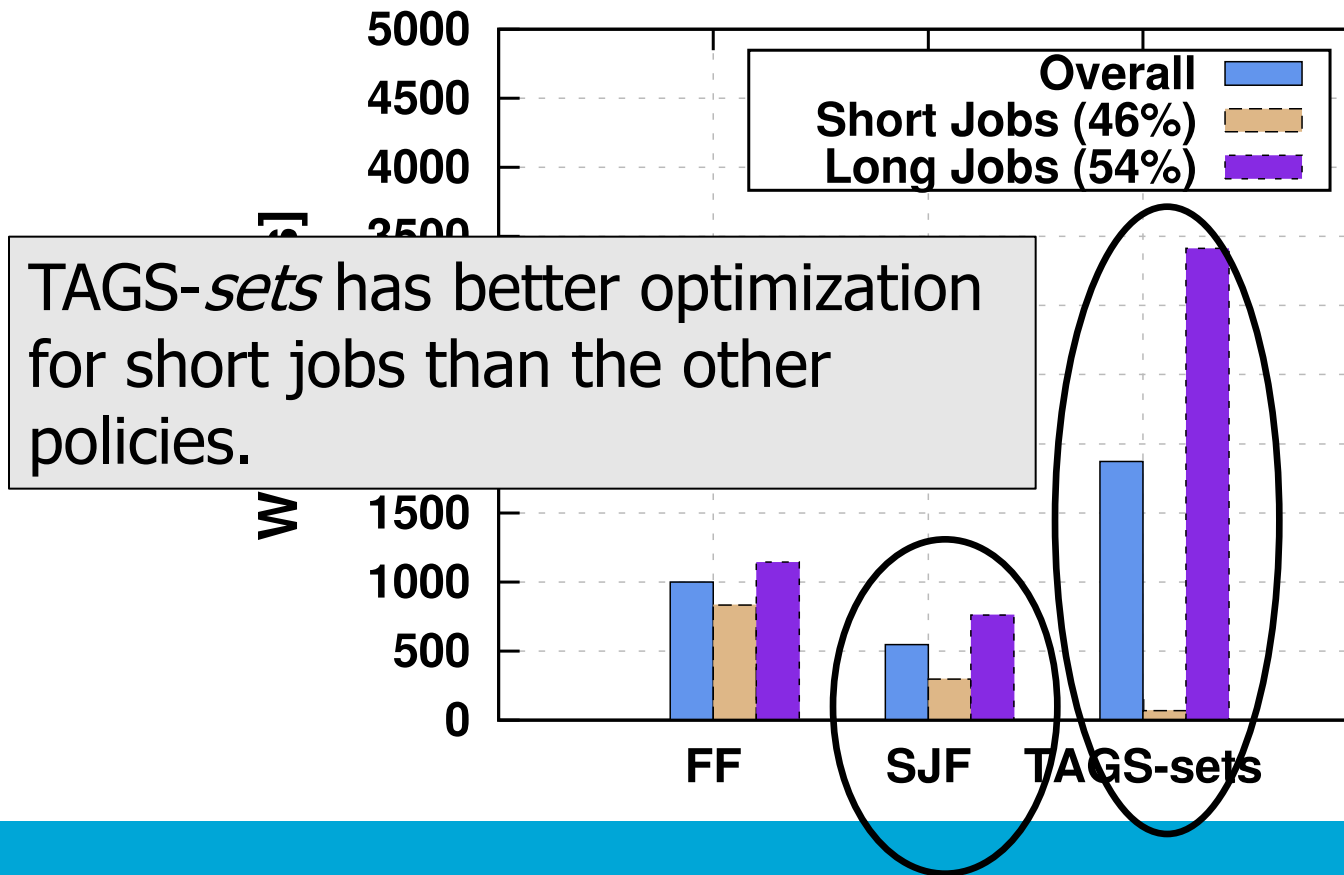
	Prediction	Slowdown	Preemption	MC+MC support
Uninformed (FF, RR)	—	↑	×	×
Informed (SJF, SJF-ideal, HSDF)	+	↓	×	×
Traditional TAGS	—	↓	✓	×
TAGS extensions (TAGS-chain, TAGS-sets) new	—	↓	✓	✓

TAGS-chain and TAGS-sets have nice properties both in simulation and in real-world experimentation.

Real-World Experimental Results



Real-world Experimental Results



Scheduling

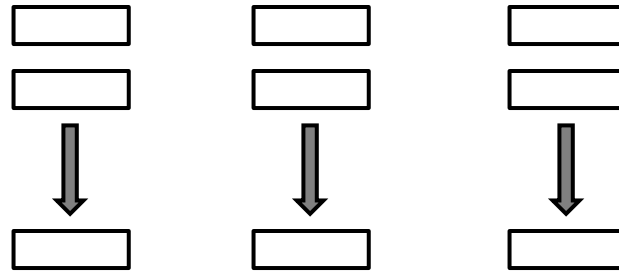


Elastic MapReduce



Dynamic Big Data Processing

Fawkes = Elastic MapReduce via Two-level scheduling architecture

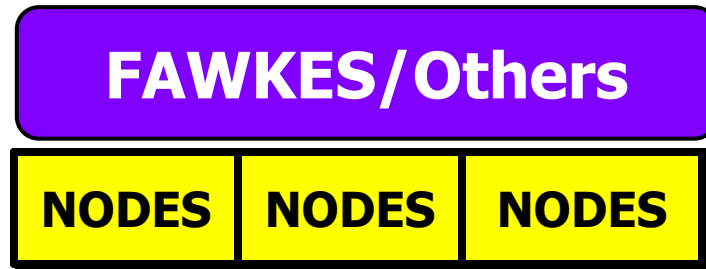


Job submissions

Frameworks

Resource manager

Infrastructure

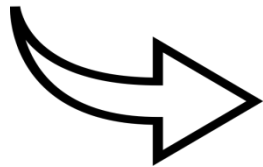
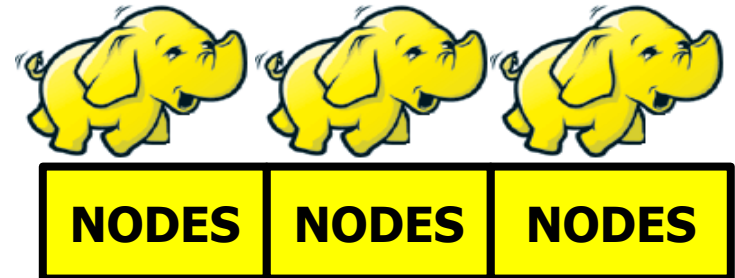


Ghit, Yigitbasi, Iosup, Epema, Iosup. Balanced Resource Allocations Across Multiple Dynamic MapReduce Clusters. ACM SIGMETRICS 2014.

Elastic MapReduce

MapReduce framework

- Distributed file system
- Execution engine
- Data locality constraints



Because workloads may be time-varying:

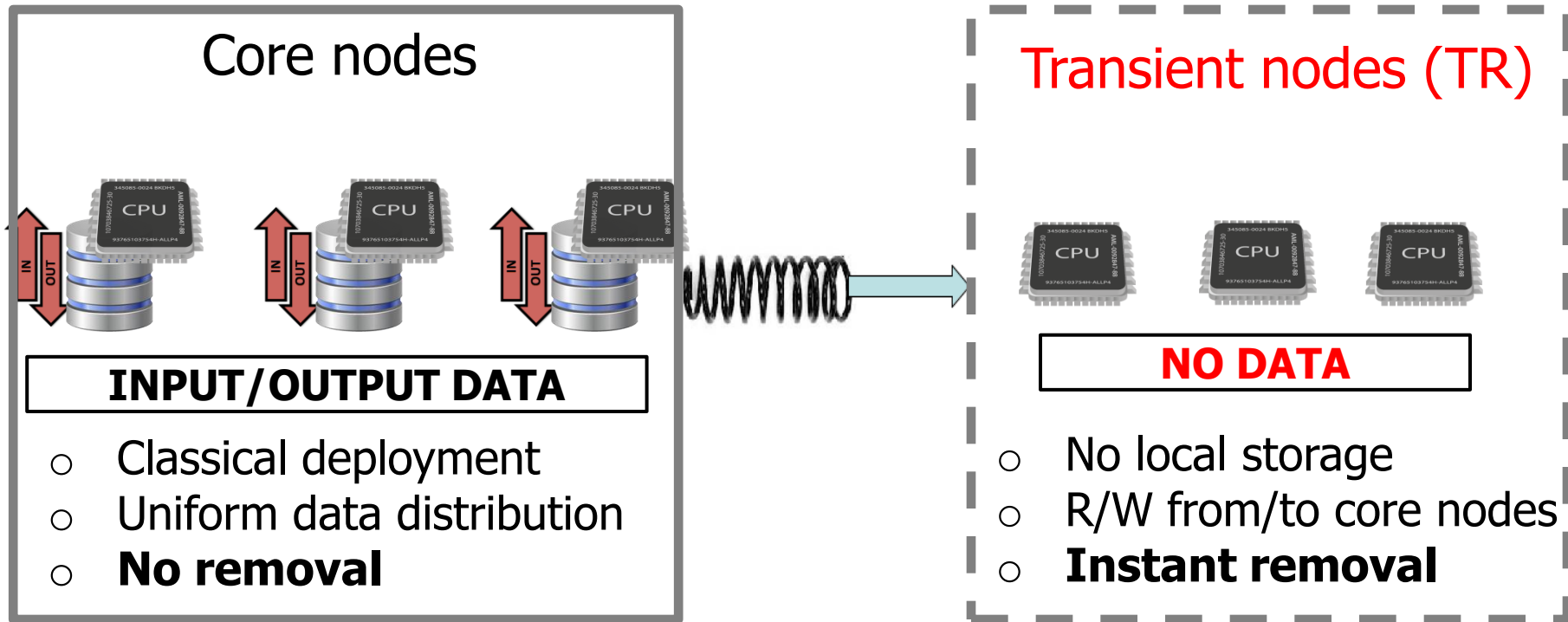
- Poor resource utilization
- Imbalanced service levels

Grow and shrink MapReduce

- High resource utilization
- Reconfiguration for balanced service levels
- Break data locality

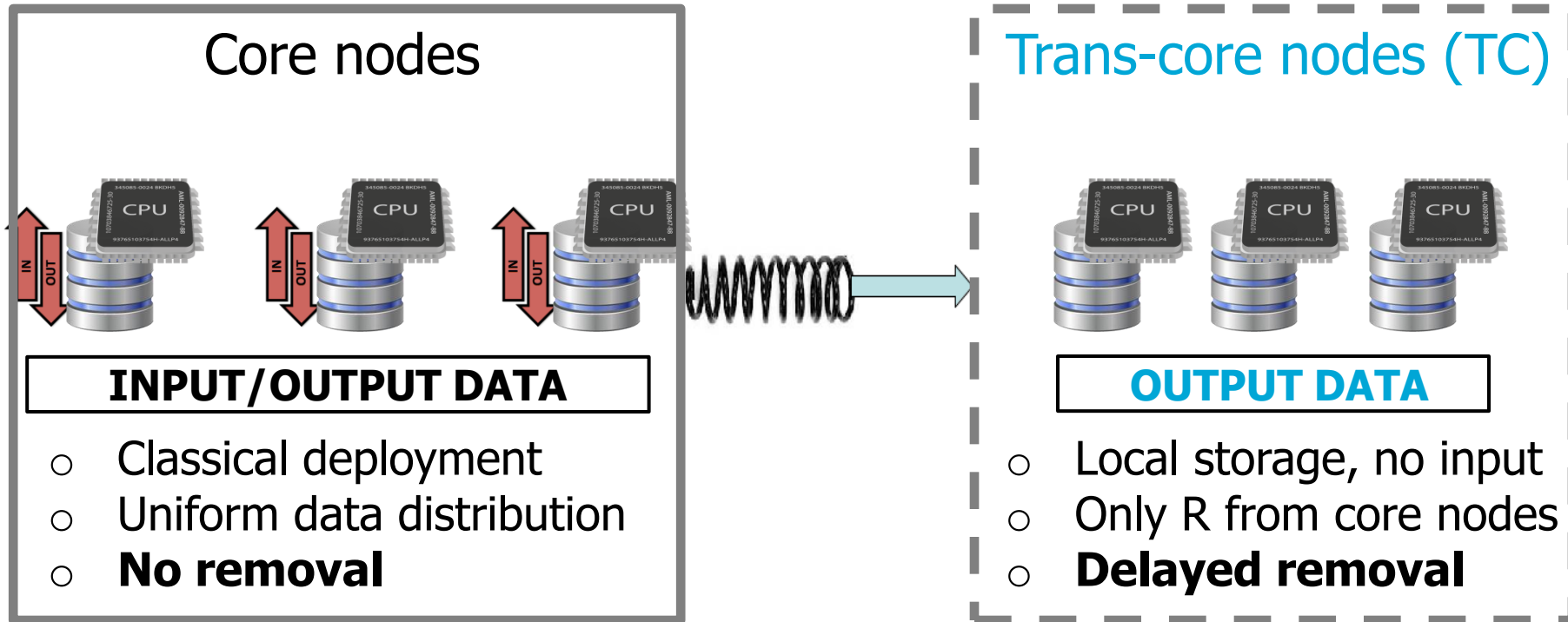


No data locality



Performance?

Relaxed data locality



Better performance?

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

FAWKES in a nutshell

1. Size of MapReduce cluster

- Changes dynamically
- Balanced by weight
- Weight by demand, usage, actual service

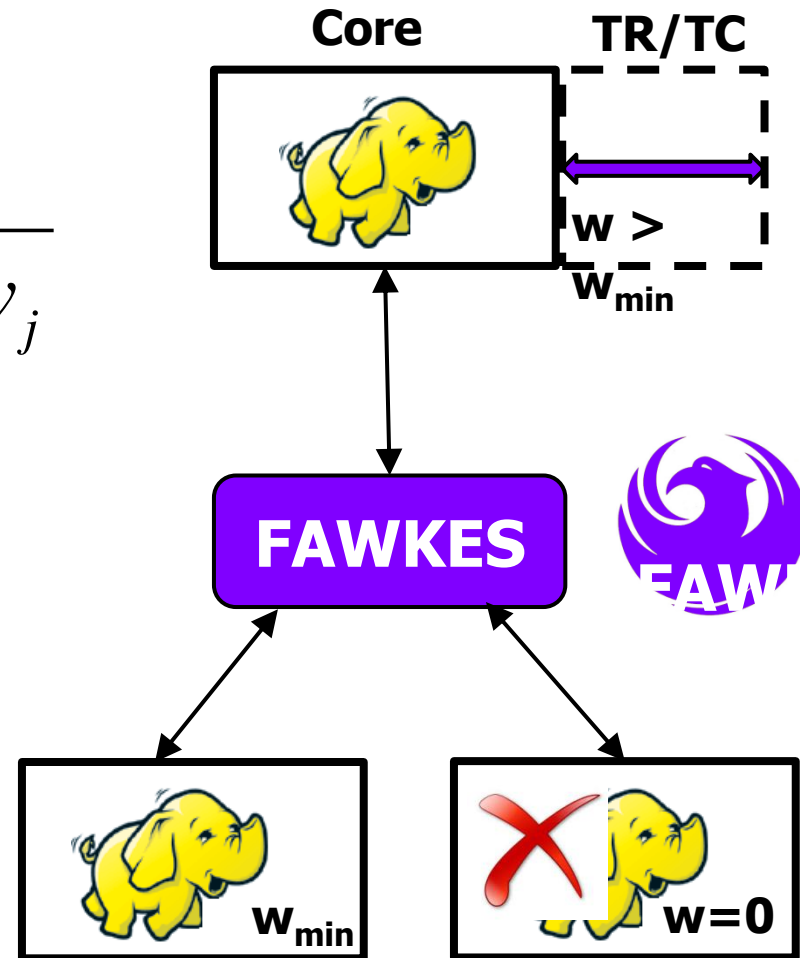
$$s_i = \frac{w_i}{\sum w_j}$$

2. Updates dynamic weights when

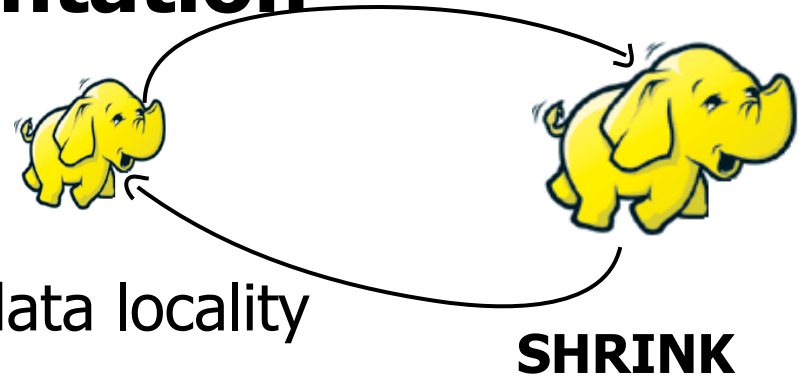
- New frameworks arrive
- Framework states change

3. Shrinks and grows frameworks to

- Allocate new frameworks (min. shares)
- Give fair shares to existing ones



Real-World Experimentation **GROW**



1. Dynamic MapReduce relaxes data locality
2. FAWKES policies can reduce imbalance between frameworks



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



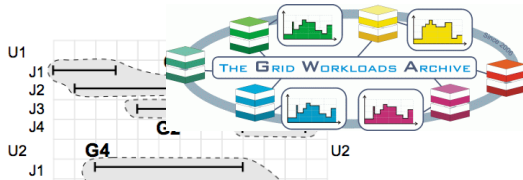
Agenda



Everyone is a Scientist



Can we afford it?



Workloads



Scheduling

Scheduling

Scheduling



Conclusion

~~Conclusion~~ Take-Home Message

- **“Everyone is a Scientist!”**
Our vision of a growing, leading Europe
- **Grand challenge**
 - Managing the datacentre
 - Helping demanding users
- **In this talk**
 - Business-critical workloads
 - Scheduling with portfolio scheduling
 - Scheduling by guessing size
 - Scheduling by load-balancing across elastic-MapReduce frameworks



Thank you for your attention! Questions? Suggestions? Observations?

More Info:



- <http://www.st.ewi.tudelft.nl/~iosup/>
- <http://www.pds.ewi.tudelft.nl/>
- <http://research.spec.org>

Alexandru Iosup

A.Iosup@tudelft.nl

(or google "iosup")

Parallel and Distributed Systems Group

Delft University of Technology

Do not hesitate
to contact me...

