



# Towards understanding today's and tomorrow's scheduling challenges in HPC systems

**Gonzalo P. Rodrigo** - [gonzalo@cs.umu.se](mailto:gonzalo@cs.umu.se)

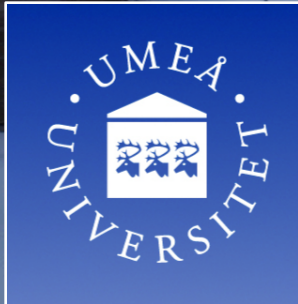
Erik Elmroth - [elmroth@cs.umu.se](mailto:elmroth@cs.umu.se)

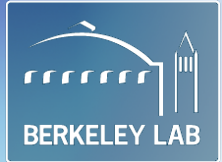
Lavanya Ramakrishnan - [iramakrishnan@lbl.gov](mailto:iramakrishnan@lbl.gov)

P-O Östberg - [p-o@cs.umu.se](mailto:p-o@cs.umu.se)

---

Distributed Systems Group – Umeå University, Sweden  
Data Science & Technology – Lawrence Berkeley National Lab





## Outline

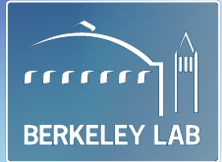
---

- Batch schedulers: Some basics
- Challenges:
  - “Exascale initiative” and “Data Explosion”
- Are schedulers ready?
- Takeaways



### **Disclaimer:**

This talk is about single site HPC scheduling!



# Is not Scheduling a “solved problem”?

[1] “And you have to realize that there are not very many things that have aged as well as the scheduler. Which is just another proof that scheduling is easy.”

**Censored**, 2001

[1] End of Dennard scaling = scheduler with an incredibly complex implementation:

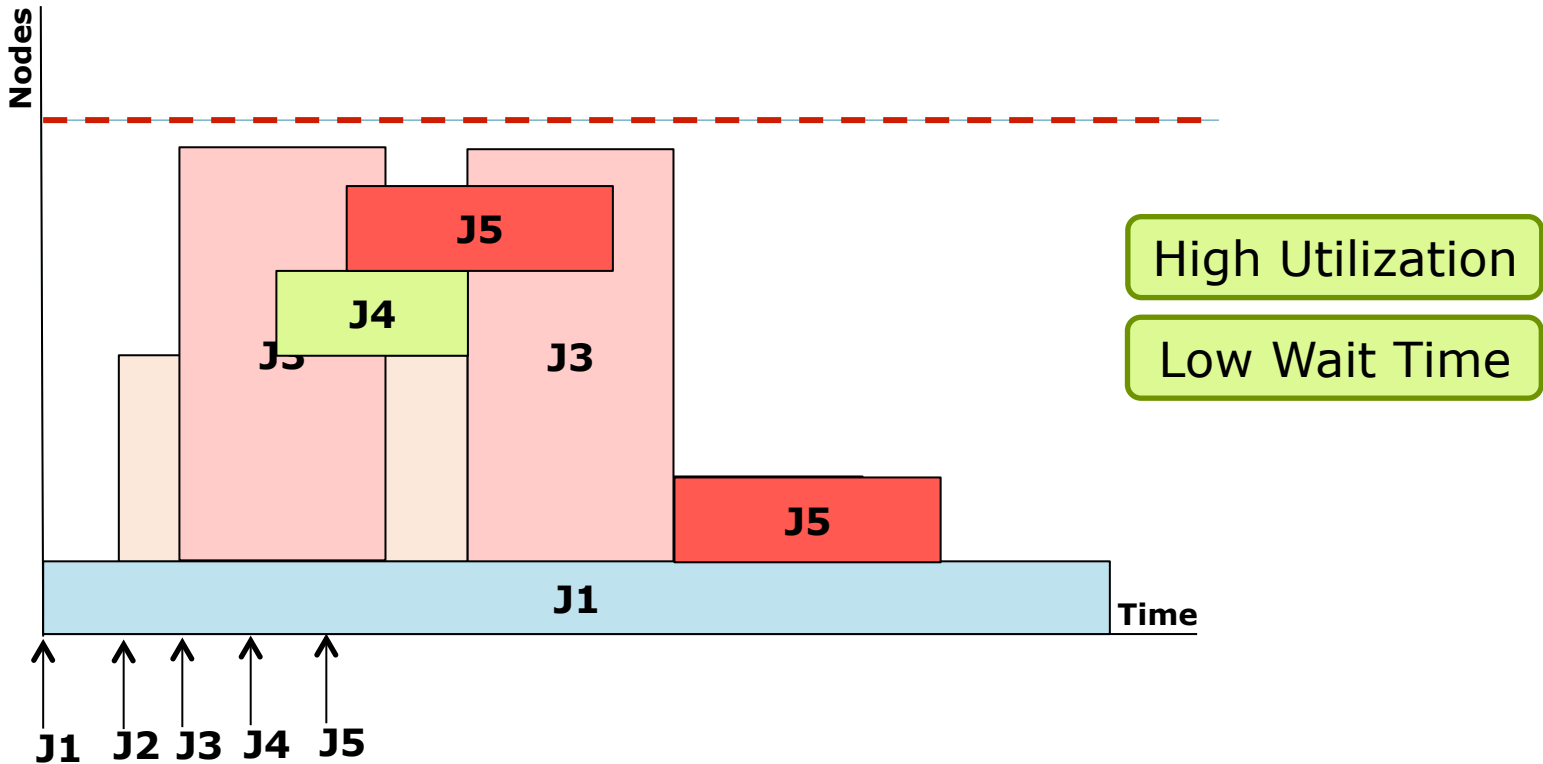
- Non-uniform memory access latencies (NUMA)
- High costs of cache coherency and synchronization
- Diverging CPU and memory latencies.

Non optimal schedulers

# Batch Schedulers: FCFS and Back-filling

**FCFS:** Jobs execute in arrival order

**Back-filling:** Job can start if it does not delay previous jobs.



# Batch Schedulers: Fairness and prioritization

Fairness

Don't starve jobs or users

Priority

Run more important jobs first

Placement?

Actually not so important(?)



# Episode 1: Upcoming challenges

---

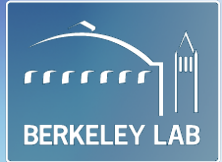
**Exascale**

**Data  
Explosion**



# Exascale: Achieve One Exaflop in 2020

---



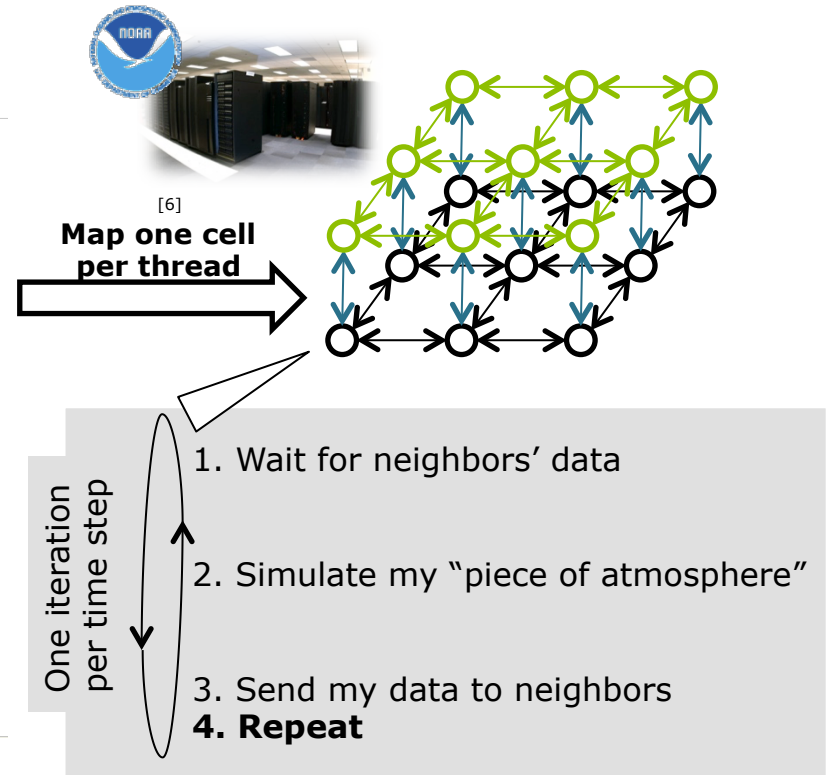
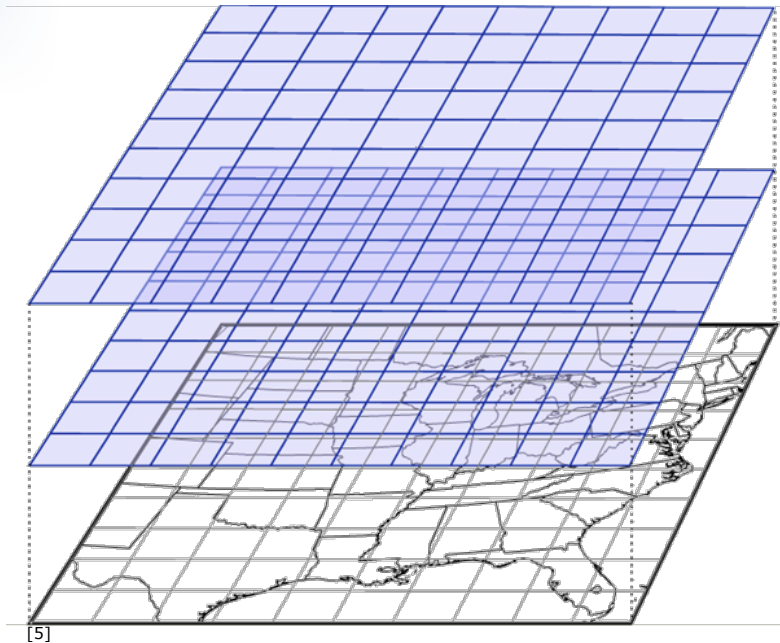
## Why?

Science is fueled by computation

Certain problems require better resolution



# Understanding Large Parallel Tightly-Coupled Jobs



[5] [https://www.e-education.psu.edu/worldofweather/sites/www.e-education.psu.edu.worldofweather/files/image/Section2/Three\\_Dimensional\\_grid%20\(Medium\).PNG](https://www.e-education.psu.edu/worldofweather/sites/www.e-education.psu.edu.worldofweather/files/image/Section2/Three_Dimensional_grid%20(Medium).PNG)

[6] NOAA Stratus and Cirrus NOAA supercomputers 2009, [http://www.noaanews.noaa.gov/stories2009/20090908\\_computer.html](http://www.noaanews.noaa.gov/stories2009/20090908_computer.html)

## It's all about power and cost

### Tianhe-2

33.86 PFLOPS

US\$390M

24 MW <sup>[7]</sup>



1 Exaflop

US\$12 870M

792 MW



~Operative Income  
Ericsson 2014 <sup>[8]</sup>

~Average Swedish  
Nuclear reactor <sup>[9]</sup>

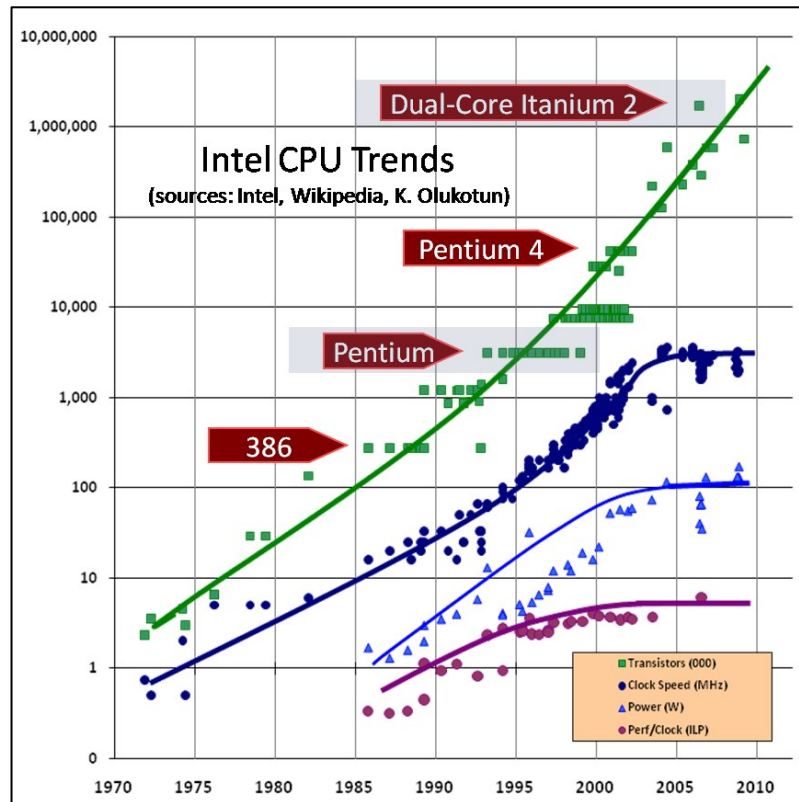
[7] <https://en.wikipedia.org/wiki/Tianhe-2>

[8] Fourth quarter and full-year report 2014 - Ericsson

[9] <http://world-nuclear.org/information-library/country-profiles/countries-o-s/sweden.aspx>

## It's all about power and cost

### Break down of Dennard scaling

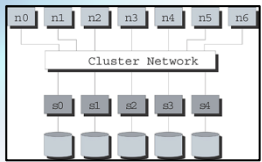


[10]

## Extreme parallelization

[10] <http://www.extremetech.com/computing/116561-the-death-of-cpu-scaling-from-one-core-to-many-and-why-were-still-stuck>

## Raw Exaflops are possible but...



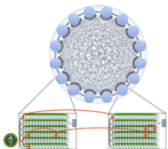
I/O

Only scalable in parallel!  
Not so good optimizations!



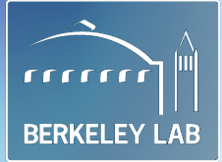
RAM

Power hungry!

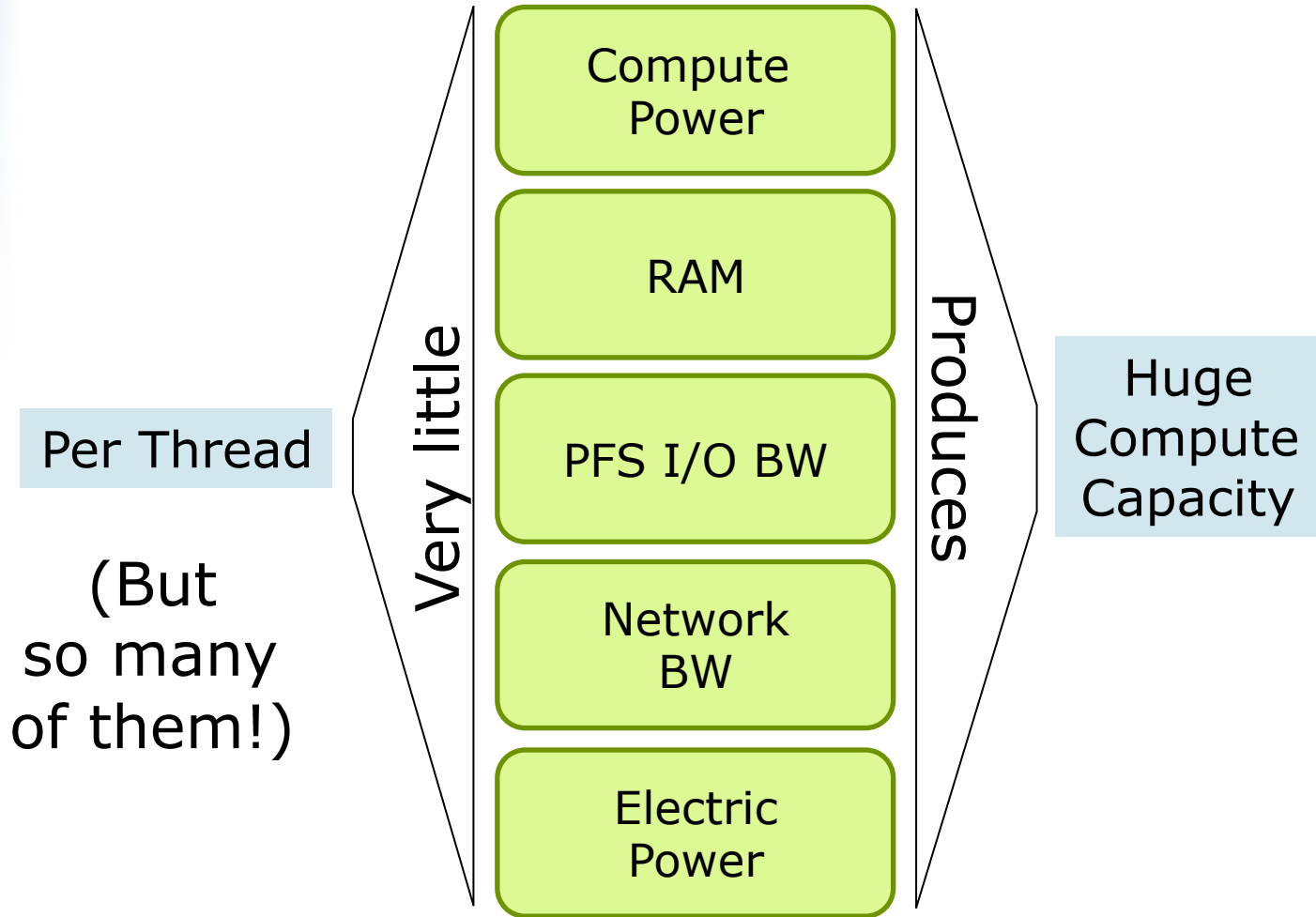


Interconnect

More parallelism => More complexity  
**Less uniform latency**

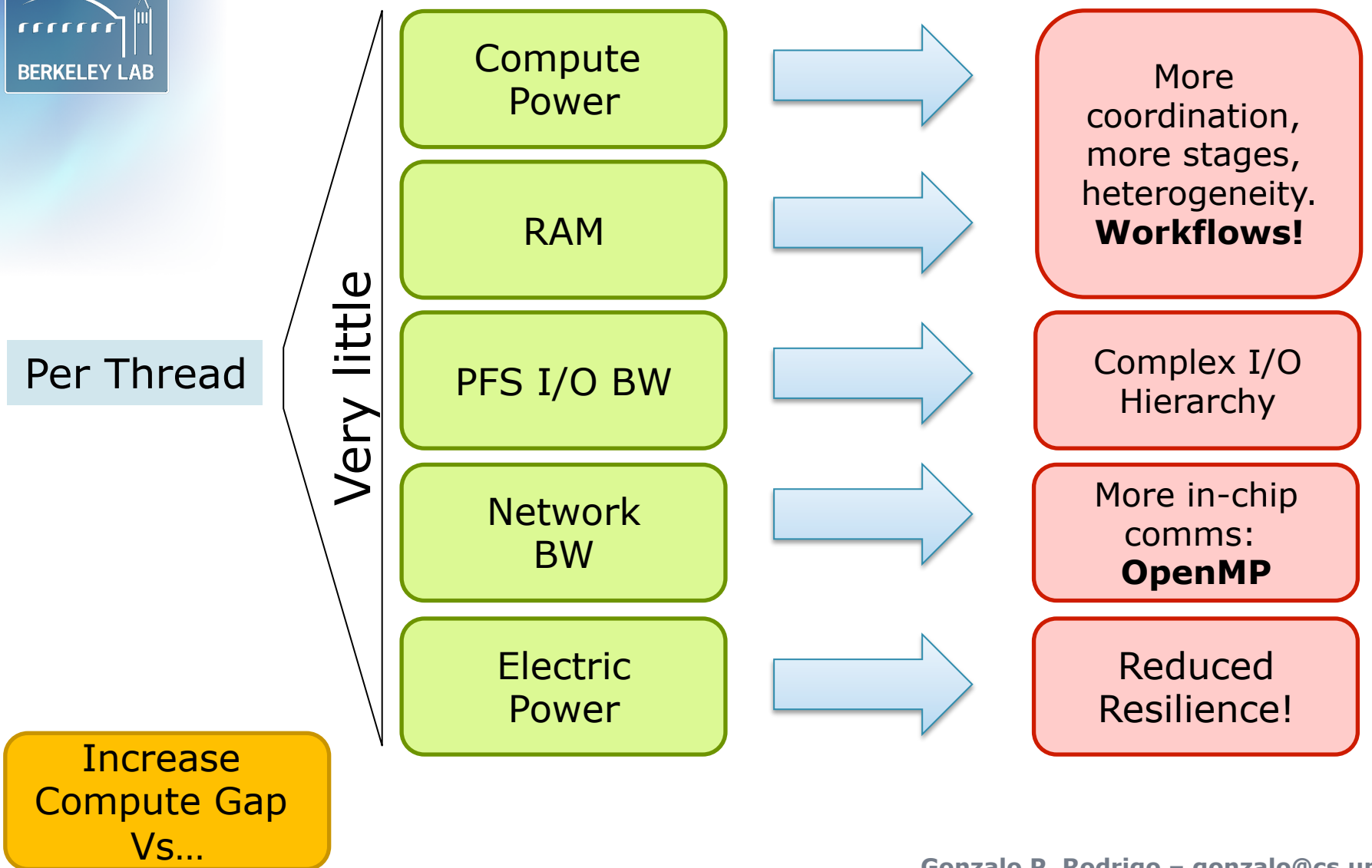


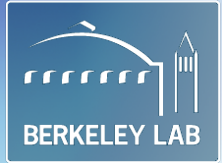
# Exascale strategy: Paradox





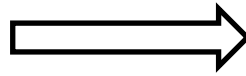
# The Exascale paradox





# Data Explosion Challenge: 4<sup>th</sup> paradigm of Science<sup>[11]</sup>

**Science**



More data than ever

**More  
compute  
power**



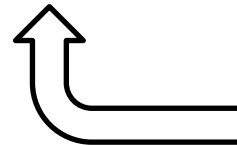
**More  
simulations**



**More  
Data**

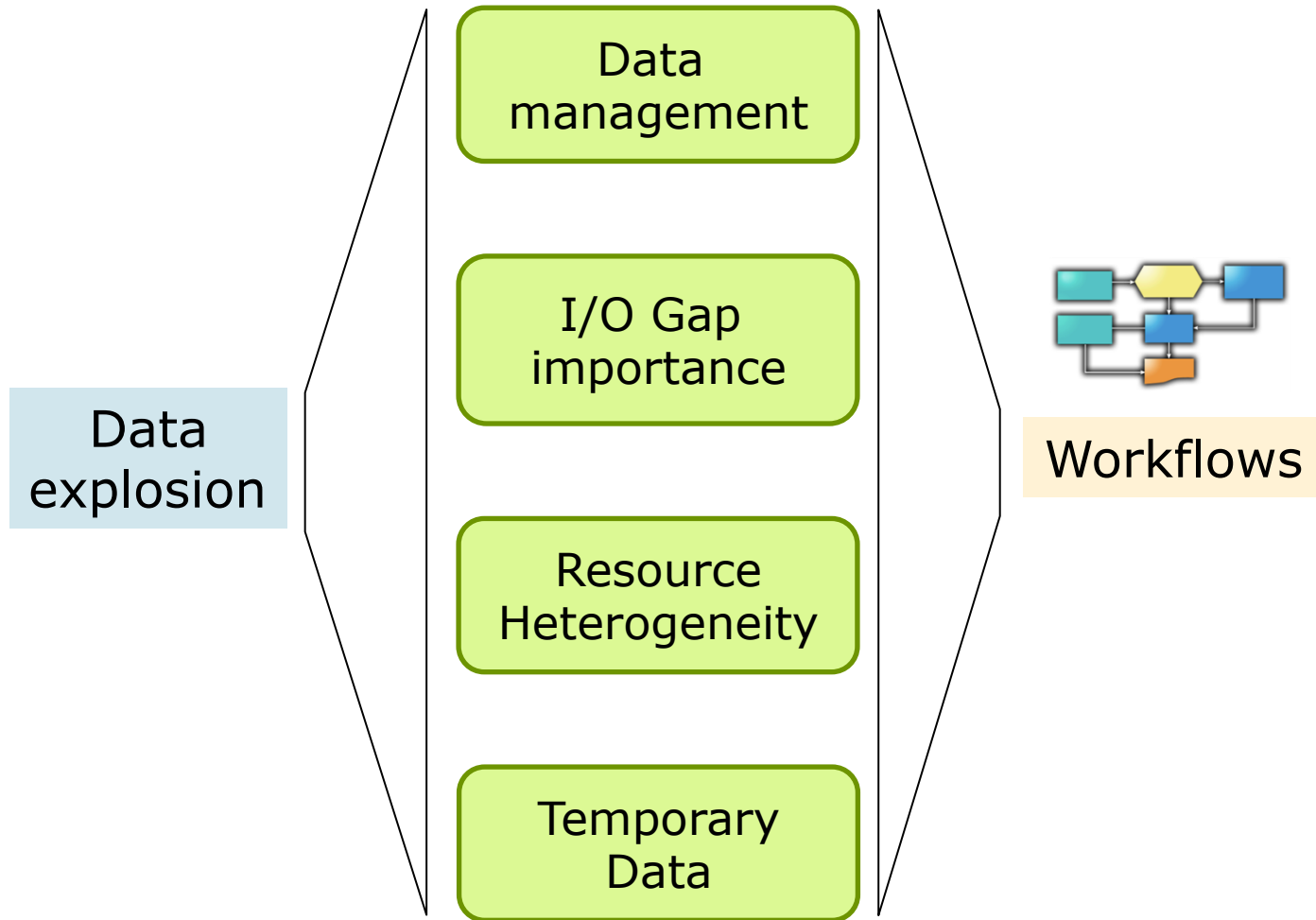


**Data  
Analysis**

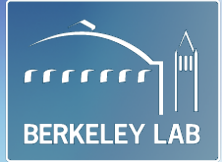


[11] Tansley, Stewart, and Kristin Michele Tolle, eds. The fourth paradigm: data-intensive scientific discovery. Vol. 1. Redmond, WA: Microsoft research, 2009.

# Data Explosion consequences







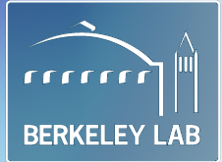
## Episode 2: challenges vs. Schedulers

Are schedulers ready for current workloads?

Are other scheduling models possible?

Can we schedule workflows better?

Performance?



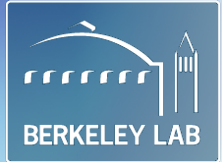
# Are schedulers ready for current workloads?

**Understanding how workloads have evolved in the past**

**Detailed analysis of current workloads**

**Observations on the performance**

# Workloads we studied



## Supercomputers

### Hopper

Deployed January 2010

Cray XE6

Gemini Network

6,384 Nodes, 24 cores/node  
**154,216 cores**

1.28 Pflops/s

Torque + Moab



### Edison

Deployed January 2014

Cray XC30

Aries Network

5,576 Nodes, 24 cores/node  
**133,824 cores**

2.57 Pflops/s

Torque + Moab



## Cluster

### Carver

Deployed 2010

IBM iDataPlex

Infiniband (fat-tree)

1,120 Nodes, 8/12/32 cores/node,  
**9,984 cores**

106.5 Tflops

Torque + Moab



# First step: System's lifetime workload evolution

**Hypothesis:** Job geometry has changed during the system's lifespan.

**Method:** Workload analysis

## Job variables

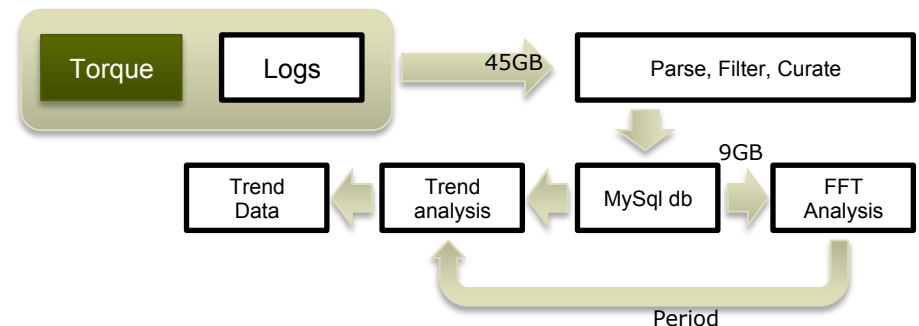
- Wall clock, number of cores (allocated), compute time, wait time, and wall clock time estimation.

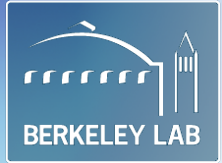
## Dataset

- 2010 – 2014: Torque logs
- 4.5M (Hopper) and 9.3M (Carver) Jobs
- Raw data 45 GB. Filtered data 9.3GB

## Analysis

- Period slicing
- Period analysis
- Comparison





# First step: System's lifetime workload evolution

Two machines with very different starting workloads, become more similar towards the end.

Most jobs are not very long and very parallel

Systems get "more loaded" in time

Users' estimations are really inaccurate.

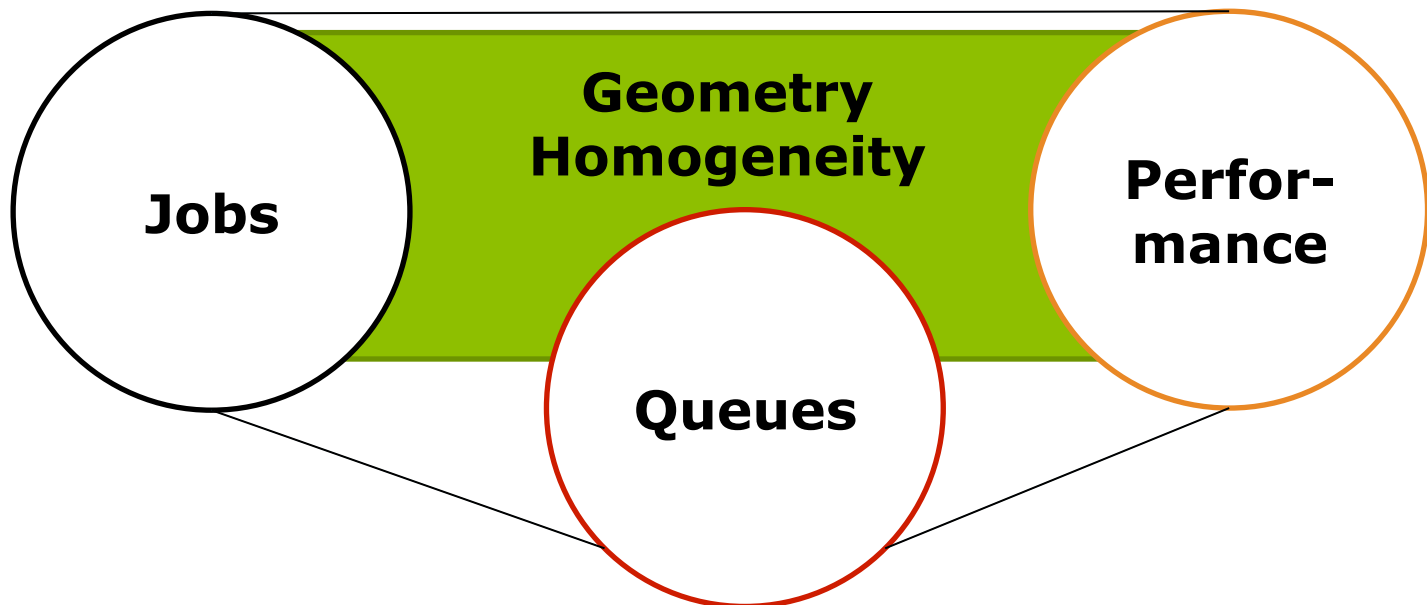
	2010		2014	
(medians)	Hopper	Carver	Hopper	Carver
<b>Wall Clock</b>	< 1 min	20 min	12 min	6 min
<b>Number of Cores</b>	100 cores	5 cores	30 cores	1 core
<b>Core Hours</b>	4 c.h.	0.9 c.h.	11 c.h.	0.09 c.h.
<b>Wait time</b>	100 s	10 min	20 min	20 min
<b>Wall clock accuracy</b>	0.2	0.25	0.21	< 0.1

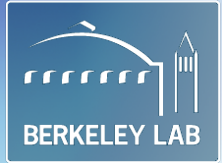
## Second step: Job Heterogeneity

**Hypothesis:** Job heterogeneity affects the scheduler performance.

**Method:** Detailed workload analysis of a year Dataset

- 2014 Torque logs
- Hopper, Edison, and Carver Jobs
- Define a method heterogeneity analysis





# Job geometry + Job priority + Job Wait time

Job Geometry

Bigger = Longer Wait

Job Priority

Higher = Shorter Wait

Queue busy

Higher = Longer Wait

Queue Homog.

Low = Predictable?

Do wait time expectation  
hold in  
Heterogeneous Queues?

Machine learning technique to detect clusters (k-means)

## Wall clock time + allocated cores

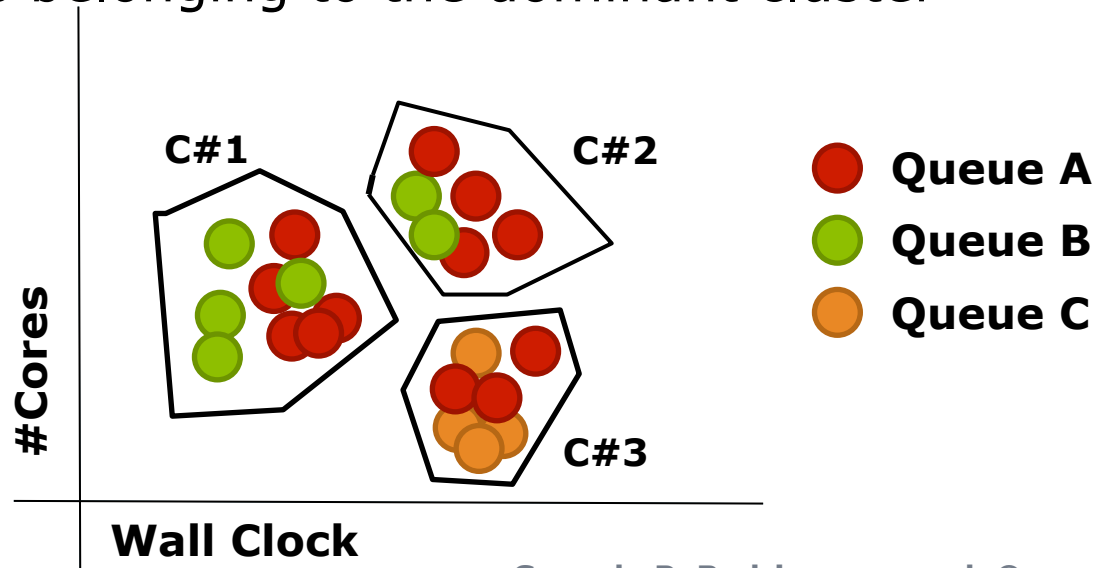
### Dominant Cluster

Cluster to which most queue jobs belong

### Queue homogeneity index

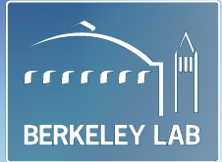
% of jobs belonging to the dominant cluster

Q	Dom. C	Hom. Idx
A	1	41%
B	1	71%
C	3	100%

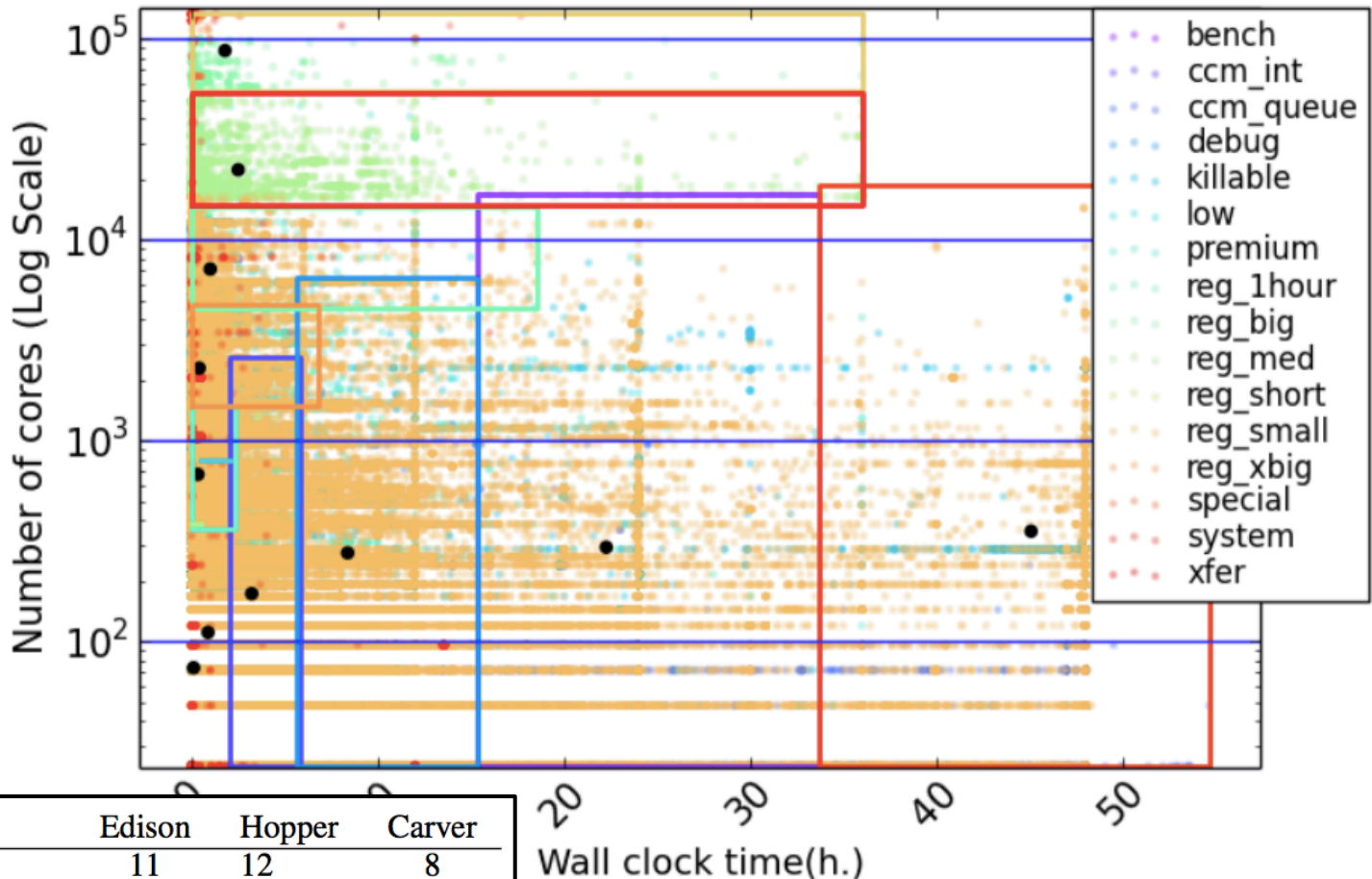




# Queue homogeneity: Cluster mapping



**Edison-Number of cores vs. Wall clock time  
Submission Queues**



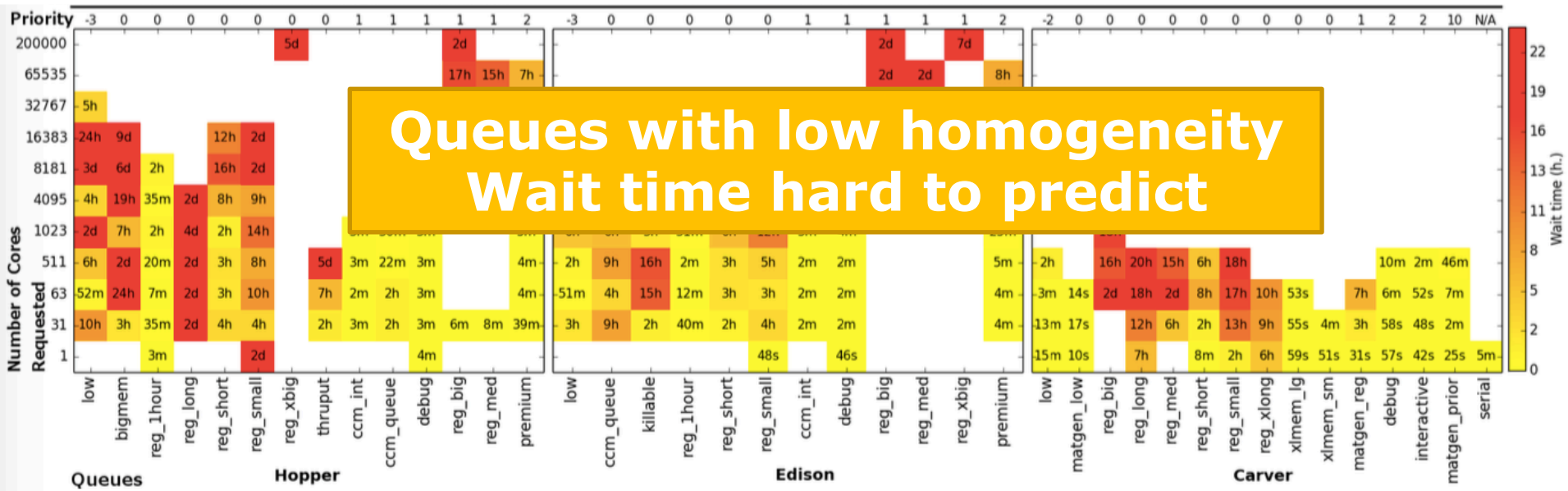
	Edison	Hopper	Carver
Clusters	11	12	8
Job homogeneity idx.	0.51	0.57	0.82
Time homogeneity idx.	0.64	0.49	0.51

20 30 40 50  
Wall clock time(h.)  
6 7 8 9 10



# Performance + Queues + Homogeneity

Wait time median in queues by requested cores



# Conclusions

(1) job geometries were fairly diverse, including a significant number of sma

**Job diversity is high**

The low per queue homogeneity indexes, show that (2) single priority policies are affecting jobs with a fairly diverse geome

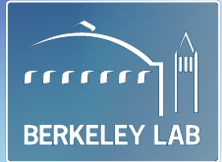
**Deal with it, or your system wait time might be hard to predict**

The w low homogeneity indexes present poor correlation between job's wait time and geometry.

Finally clock time accuracy (fundamental for the performance of backfilling functions) is very low.

**Maybe queues should be re-ordered**

**Let's do something about run time prediction**



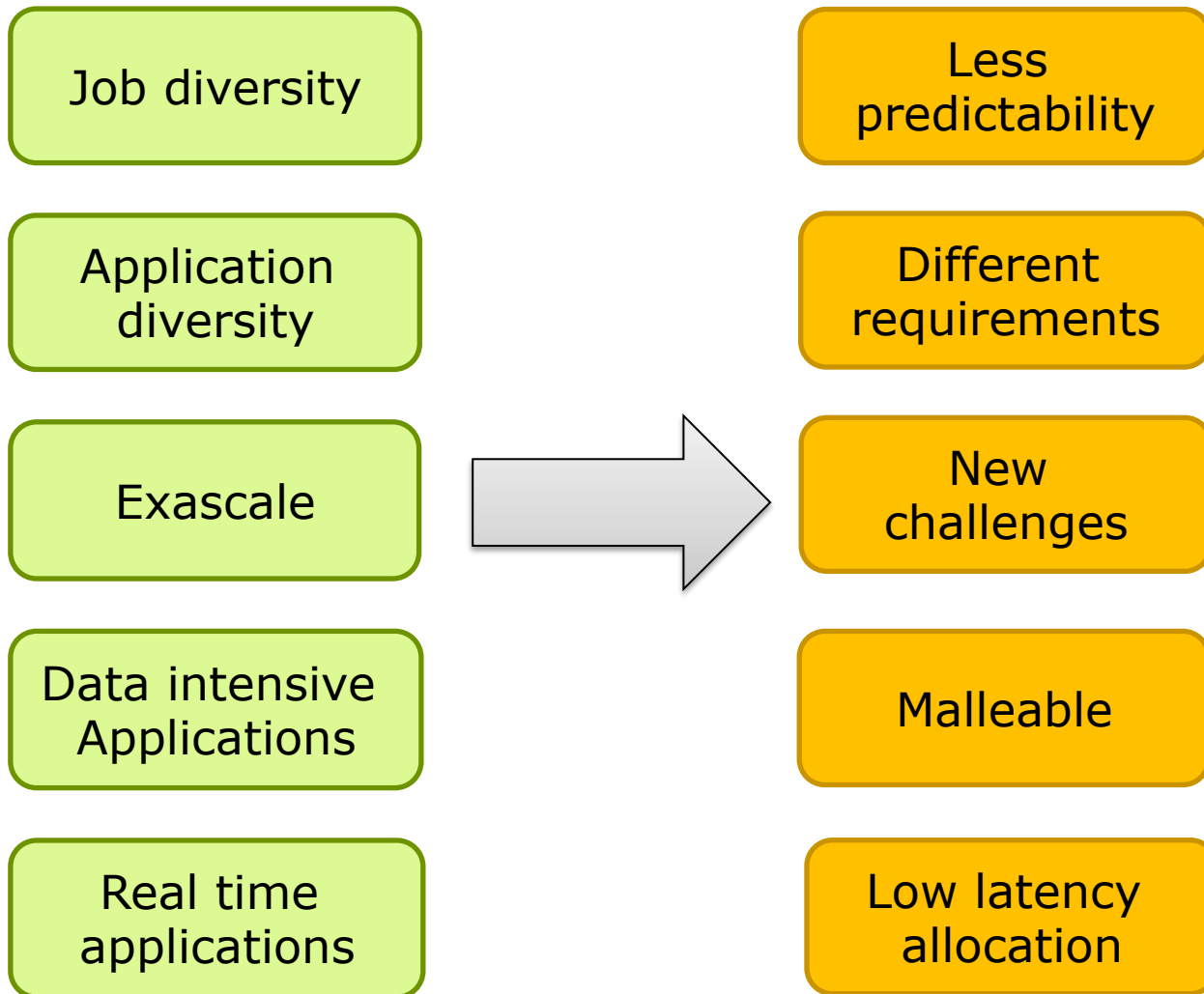
So..

---

Are schedulers ready for the current  
(and future) workload?

**Other challenges?**

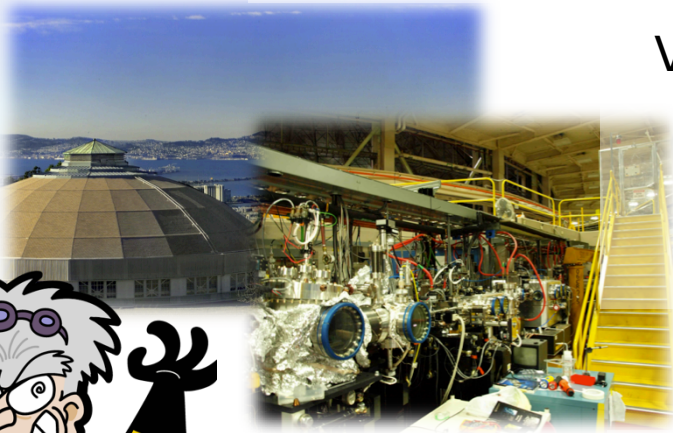
# Game changers vs. Schedulers



# Game changer: Live experiments data processing (stream)



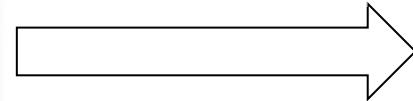
## Advance Light Source



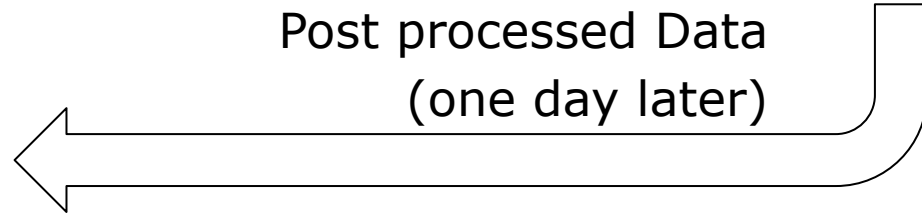
## Carver (IBM iDataPlex)



Video recording  
(data)  
3D Scanner of  
materials



Post processed Data  
(one day later)



- Live experiment
- Produces data (large amounts)
- Required to be processed on a super computer
- Processed results one day later

- Experiment would benefit of live feedback!
- Reservations are hard to align to reality!



# Looking for inspiration... in the clouds.

---

Cloud infrastructures have faced similar challenges...



**Hypothesis:** Cloud scheduling techniques can be applied to tackled new HPC challenges.

**Method:** Compared study on techniques and application circumstances (Survey)

# Similarities



## Applications



### Heterogeneous Workload

Batch Jobs

Data is Key

Wait Time is important

### Heterogeneous Workload

Many non tightly coupled

Non-classical HPC

Response time

### SSDs on Nodes

### Distributed Filesystems

### Heterogeneous resources

### Burst Buffer

Accelerator HW

BB nodes

Compute nodes

## Infrastructure



**Application aware scheduling:** Aware of characteristics, performance models, different rules for different types of job.

**Dynamically malleable management:** runtime re-scaling of jobs, performance based allocation.

**Flexible backfilling:** for better utilization

**Low latency allocation:** To allow allocation of jobs a short time after submission (stream job)





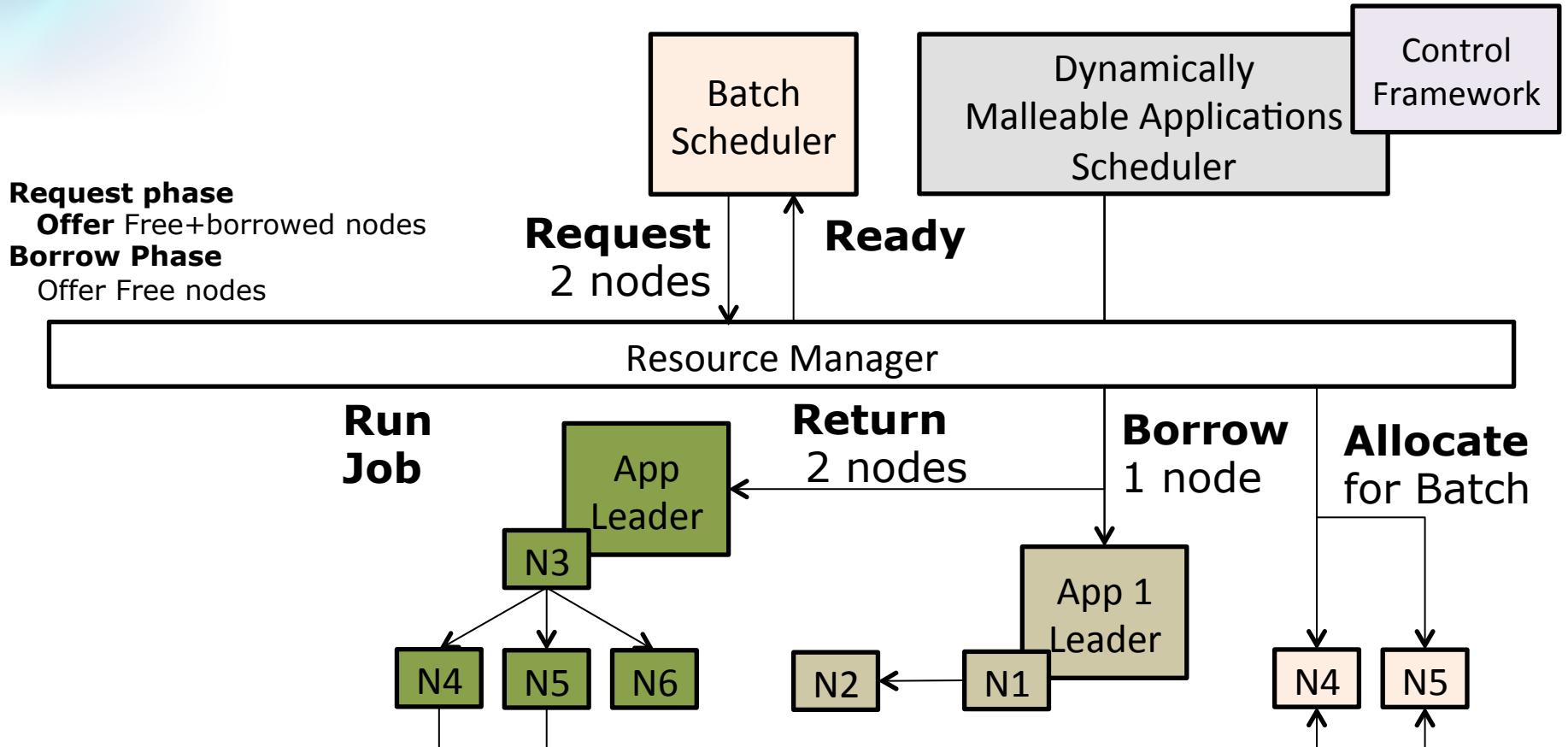
# Scheduler model

Cloud borrowed solution: **Two level scheduling**

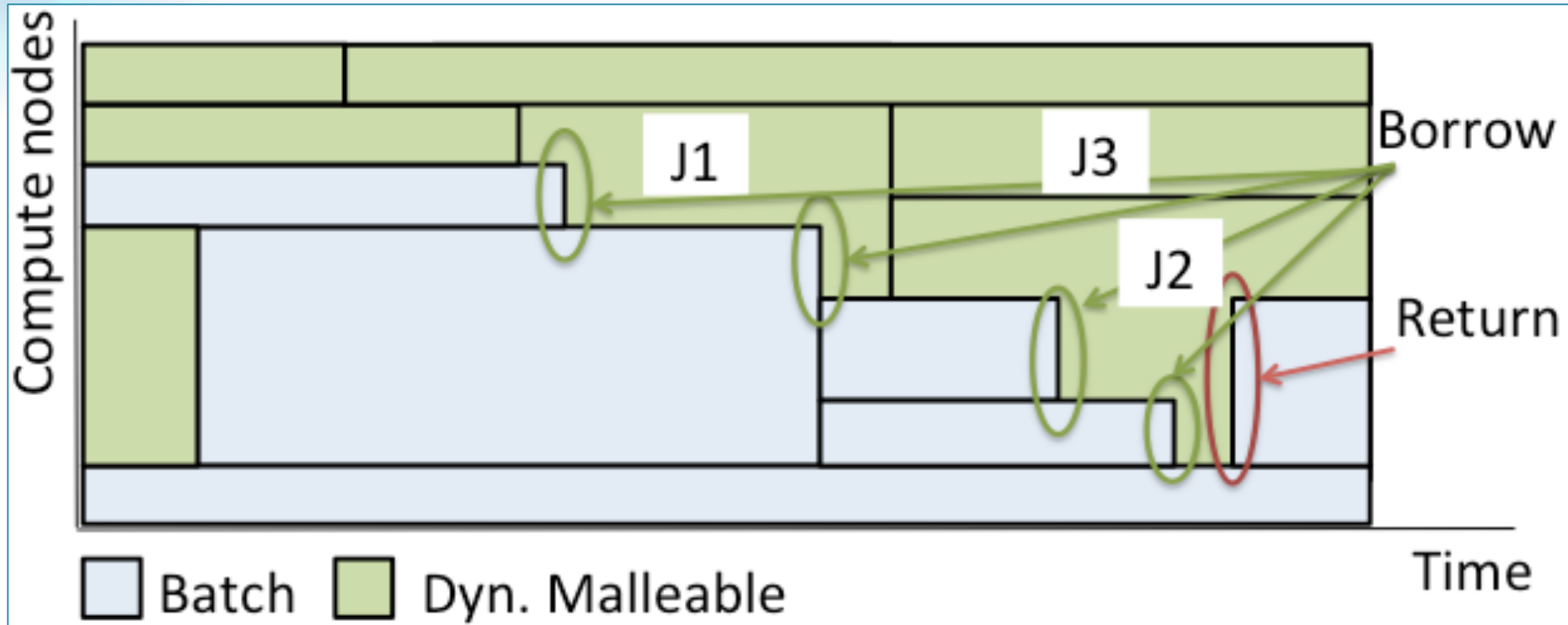
One scheduler per application + smart RM

Malleable Applications: Dynamic allocation

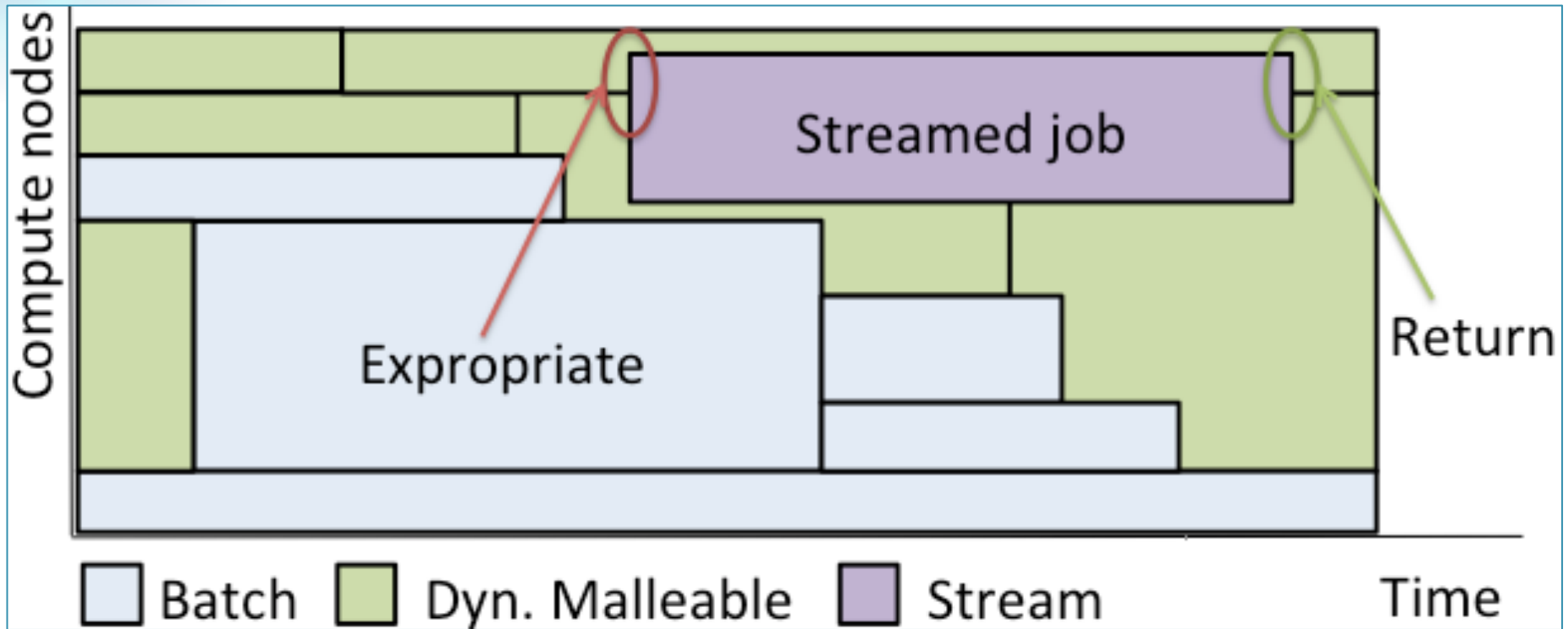
Low latency allocation



# Flexible backfilling



# Resource Expropriation: Low latency allocation



## A2L2: Conclusions

Application heterogeneity and the mix of both cloud and HPC

Application Aware

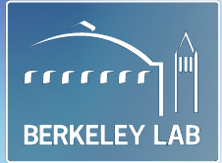
Flexibility (i.e. applications can be maybe enough malleable to make be useful)

Experimentation in progress

Application Management

Better utilization

Stream job allocation

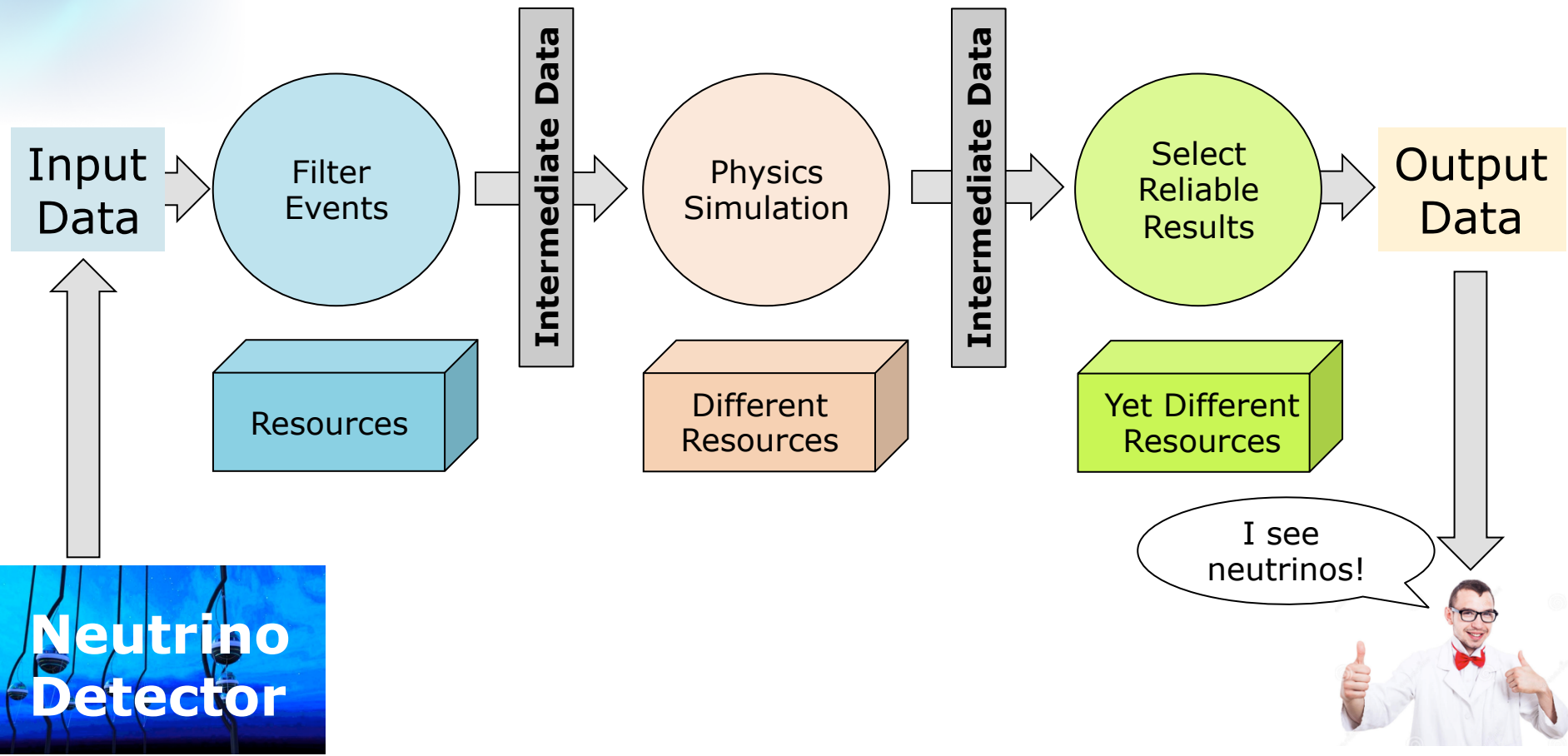


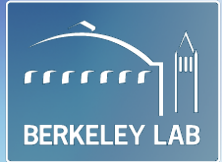
But..

---

How do scheduler deal with  
**Workflows?**

# But before... What is a workflow?





## But before... What is the problem?

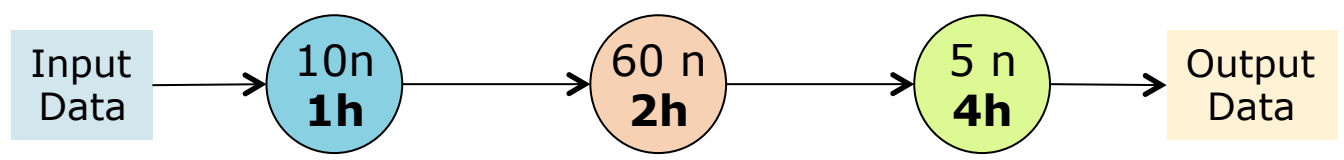
---

Classical schedulers are not optimized to manage workflows within the cluster.

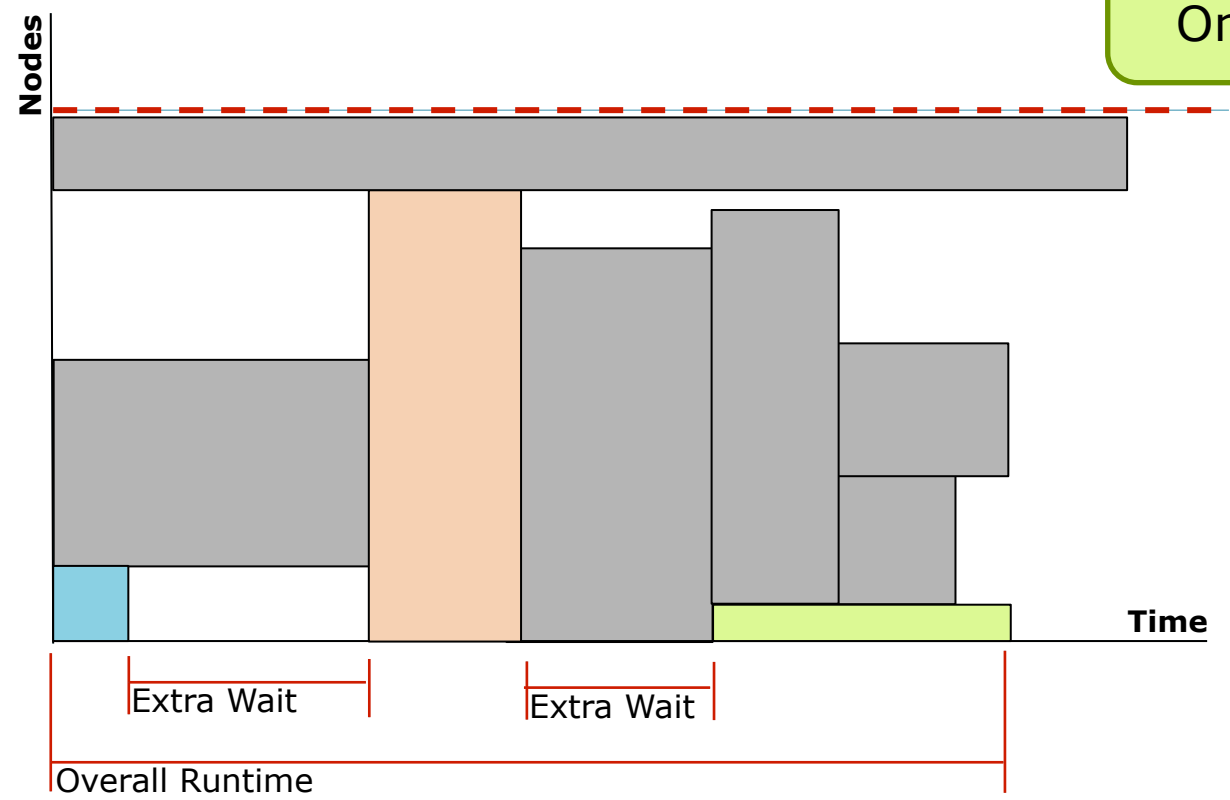
Is that so bad?



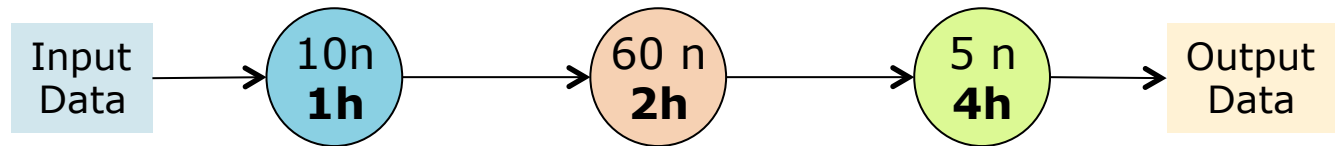
# Submitting a workflow: Wait! (approach)



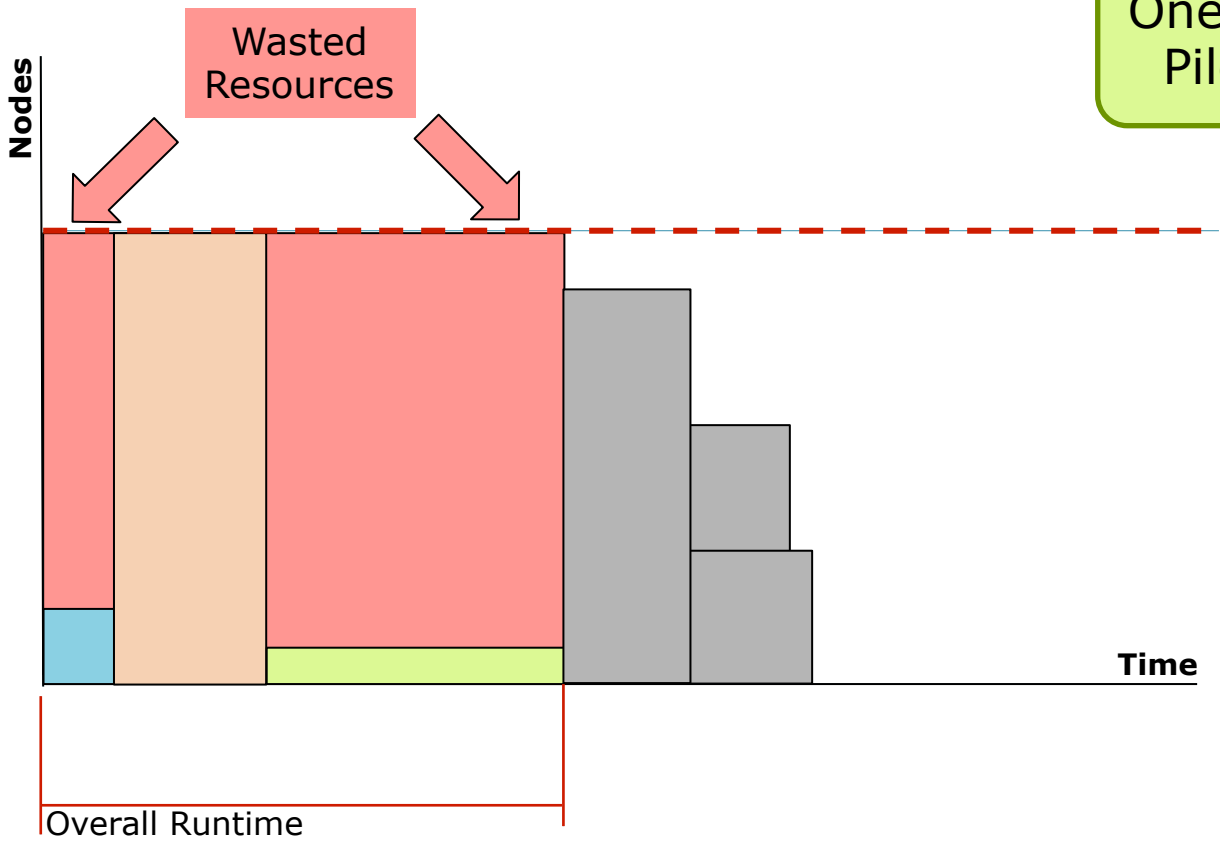
One stage  
One Job

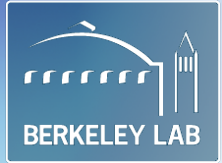


# Submitting a workflow: Waste! (approach)



One single Pilot Job





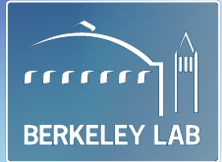
# The reality of current workflow scheduling

---

Users either waste resources...  
...or wait long time.

Something in between could be done!

(...To be published next fall)



## A final note on scale and performance

### **Exascale = More parallel jobs**

How many flops will need the scheduler alone?  
A mini cluster to manage the cluster?

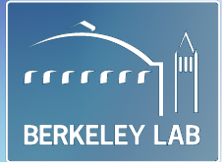
Distributed Scheduler

Multiple schedulers

Partitions

Smart RMs

New programming models?



## Takeaways!

Workloads have **changed**

Observations of **job heterogeneity** possibly affecting schedulers performance

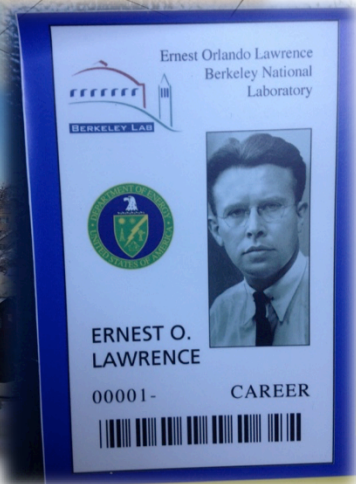
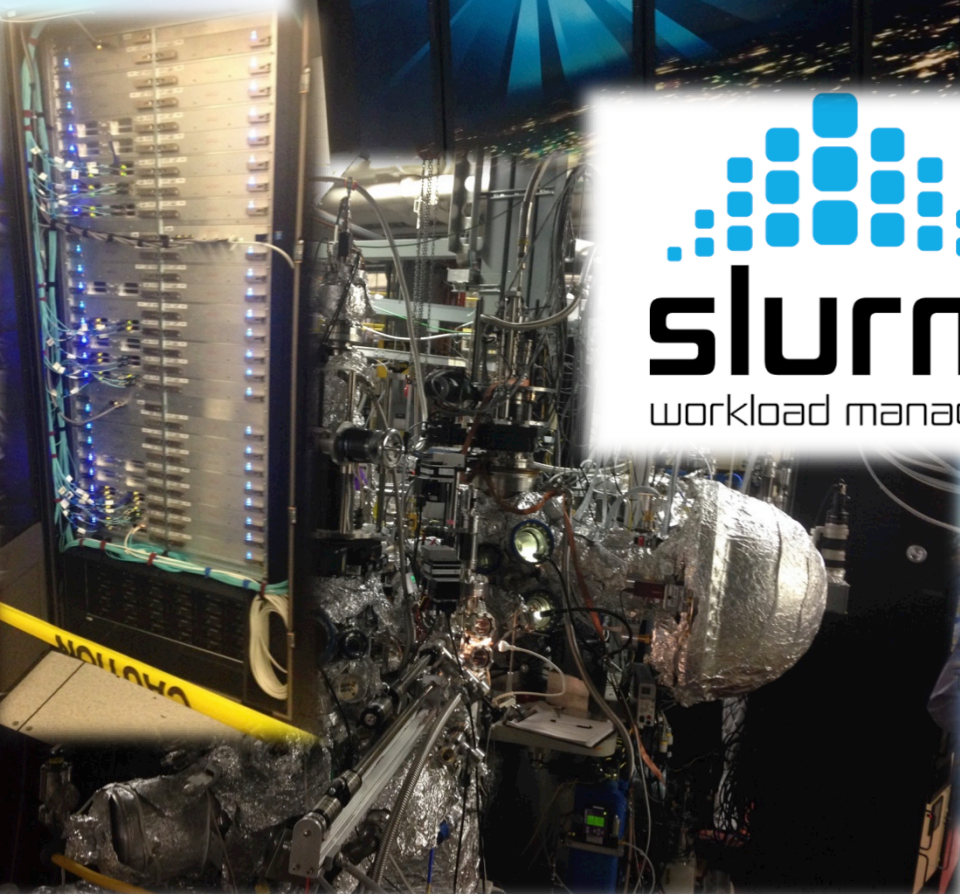
**Alternate models of scheduling** should be explored to address new challenges

**Workflows are more important than ever:**  
Scheduler should address them accordingly.

Big systems: more scheduling load...  
**performance!**

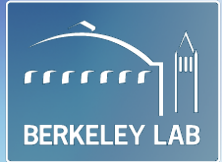


Thanks for time... questions?





To know more....



Contact:

[gonzalo@cs.umu.se](mailto:gonzalo@cs.umu.se) - [gprodrigoalvarez@lbl.gov](mailto:gprodrigoalvarez@lbl.gov)

G. Rodrigo, P-O. Östberg, E. Elmroth, K. Antypas, R. Gerber, and L. Ramakrishnan. Towards Understanding Job Heterogeneity in HPC: A NERSC Case Study. CCGrid 2016 - The 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, Accepted, 2016.

Rodrigo Álvarez, G. P., Östberg, P. O., Elmroth, E., Antypas, K., Gerber, R., & Ramakrishnan, L. (2015, June). HPC System Lifetime Story: Workload Characterization and Evolutionary Analyses on NERSC Systems. In Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing (pp. 57-60). ACM. Citation

Rodrigo Álvarez, G. P., Östberg, P. O., Elmroth, E., & Ramakrishnan, L. (2015, June). A2L2: An Application Aware Flexible HPC Scheduling Model for Low-Latency Allocation. In Proceedings of the 8th International Workshop on Virtualization Technologies in Distributed Computing (pp. 11-19). ACM. Citation

Rodrigo, G. P., Östberg, P-O. & Elmroth, E. (2014). Priority Operators for Fairshare Scheduling. 18th Workshops on Job Scheduling Strategies for Parallel Processing (JSSPP 2014) hosted at the IPDPS-2014 conference. Full Text

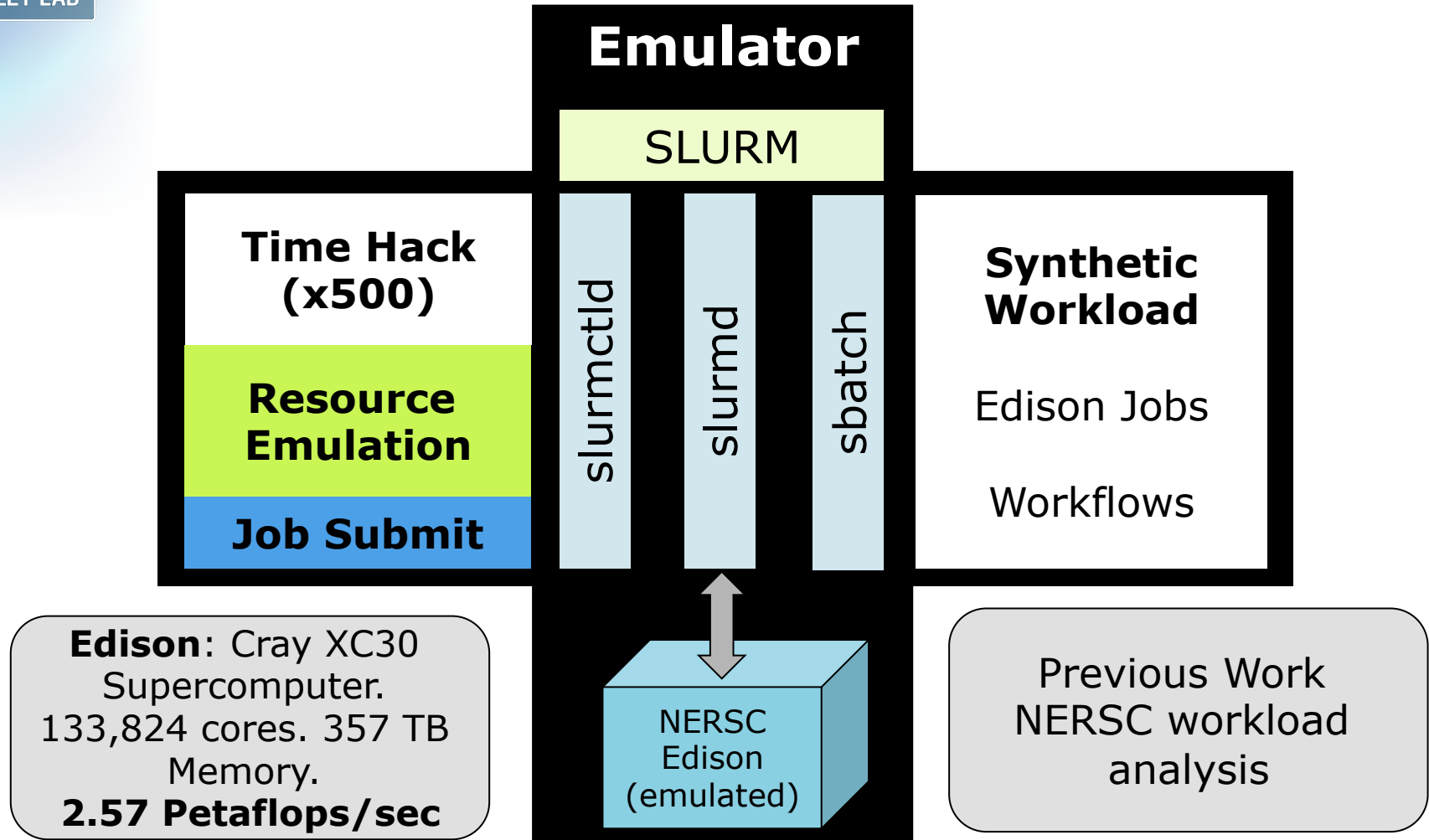
Rodrigo, G. P. Establishing the equivalence between operators: theorem to establish a sufficient condition for two operators to produce the same ordering in a Fairshare prioritization system. January 2014. Full Text

Rodrigo, G. P. Proof of compliance for the relative operator on the proportional distribution of unused share in an ordering fairshare system. January 2014. Full Text



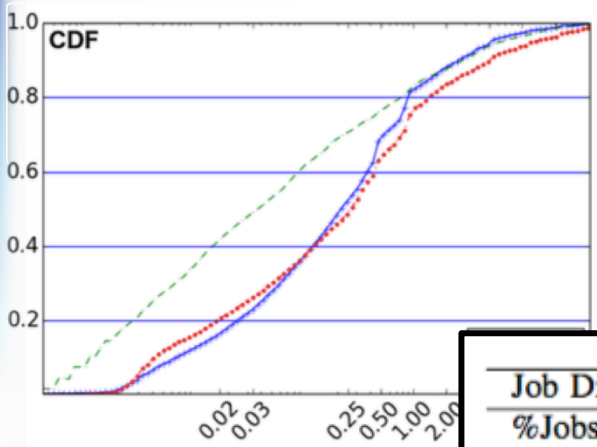
# Slurm as a scheduling research platform

Emulator based work on Barcelona Supercomputing Center (BSC) and Swiss National Supercomputing Center (CSCS)... **but our own timing routines**

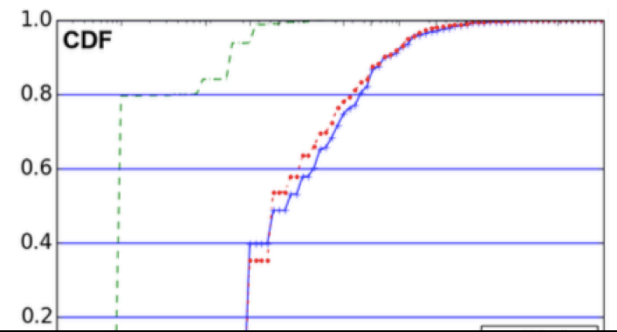




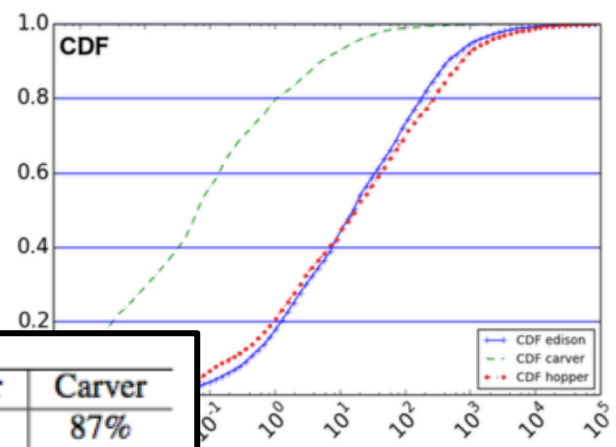
# Second step: Current Jobs



(a) Wall clock Time (h.)

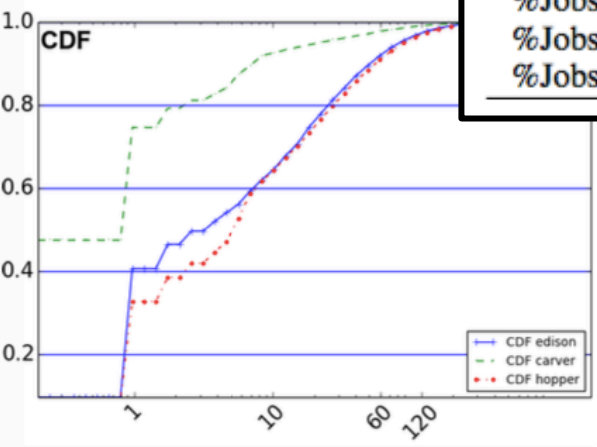


(b) Wall clock time accuracy

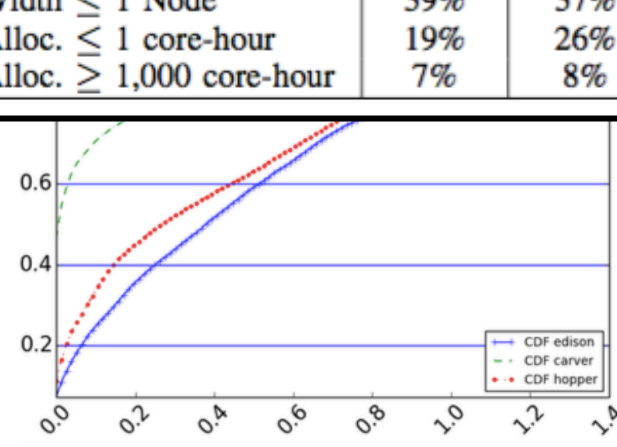


(c) Core-hours

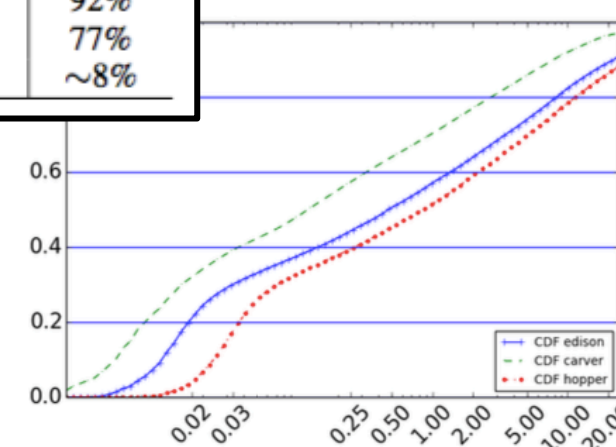
Job Distribution	Edison	Hopper	Carver
%Jobs Wall Clock < 2 h.	88%	86%	87%
%Jobs Width < 240 Cores	69%	75%	99%
%Jobs Width ≤ 1 Node	39%	37%	92%
%Jobs Alloc. ≤ 1 core-hour	19%	26%	77%
%Jobs Alloc. ≥ 1,000 core-hour	7%	8%	~8%



(a) Inter-arrival time (s.)



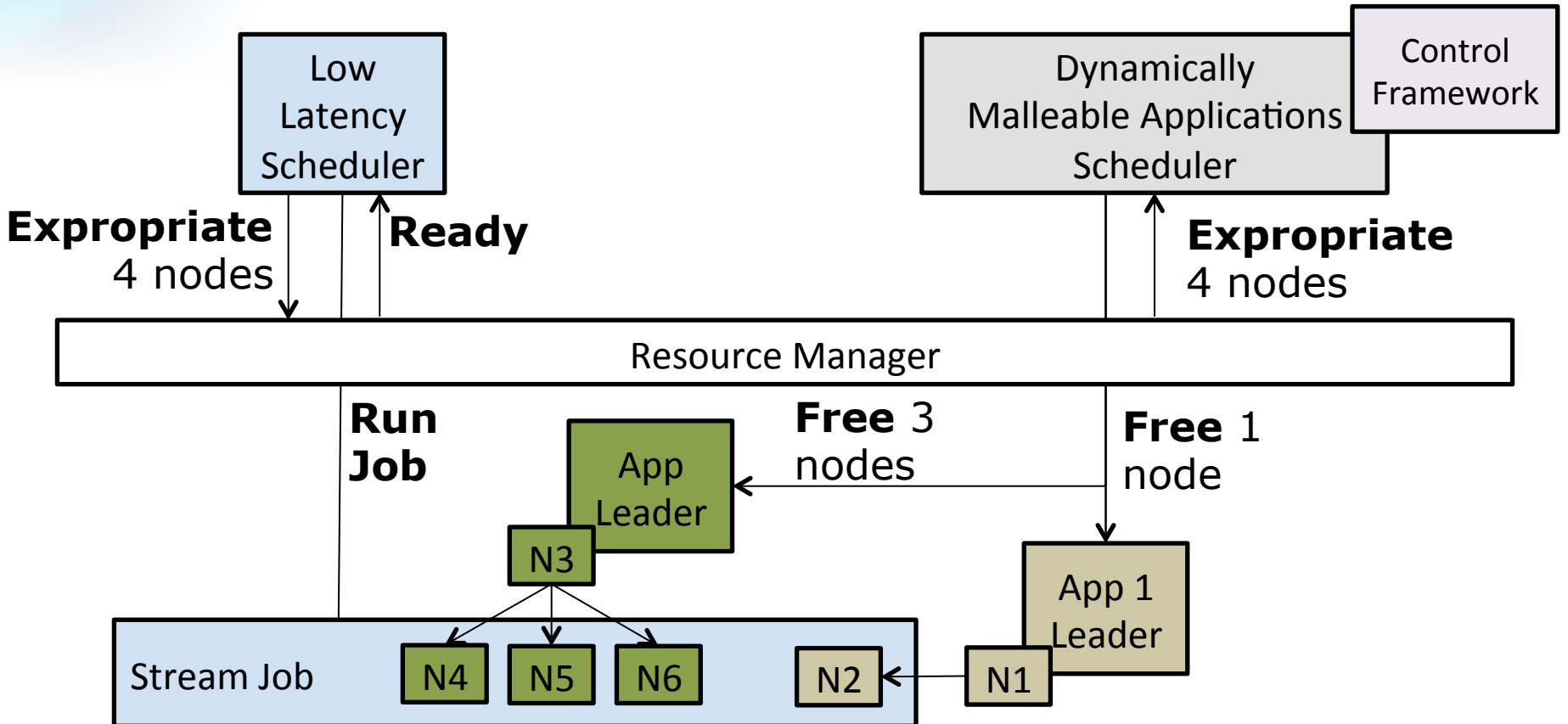
(c) Wait time (h.)



# Resource Expropriation: Low latency allocation

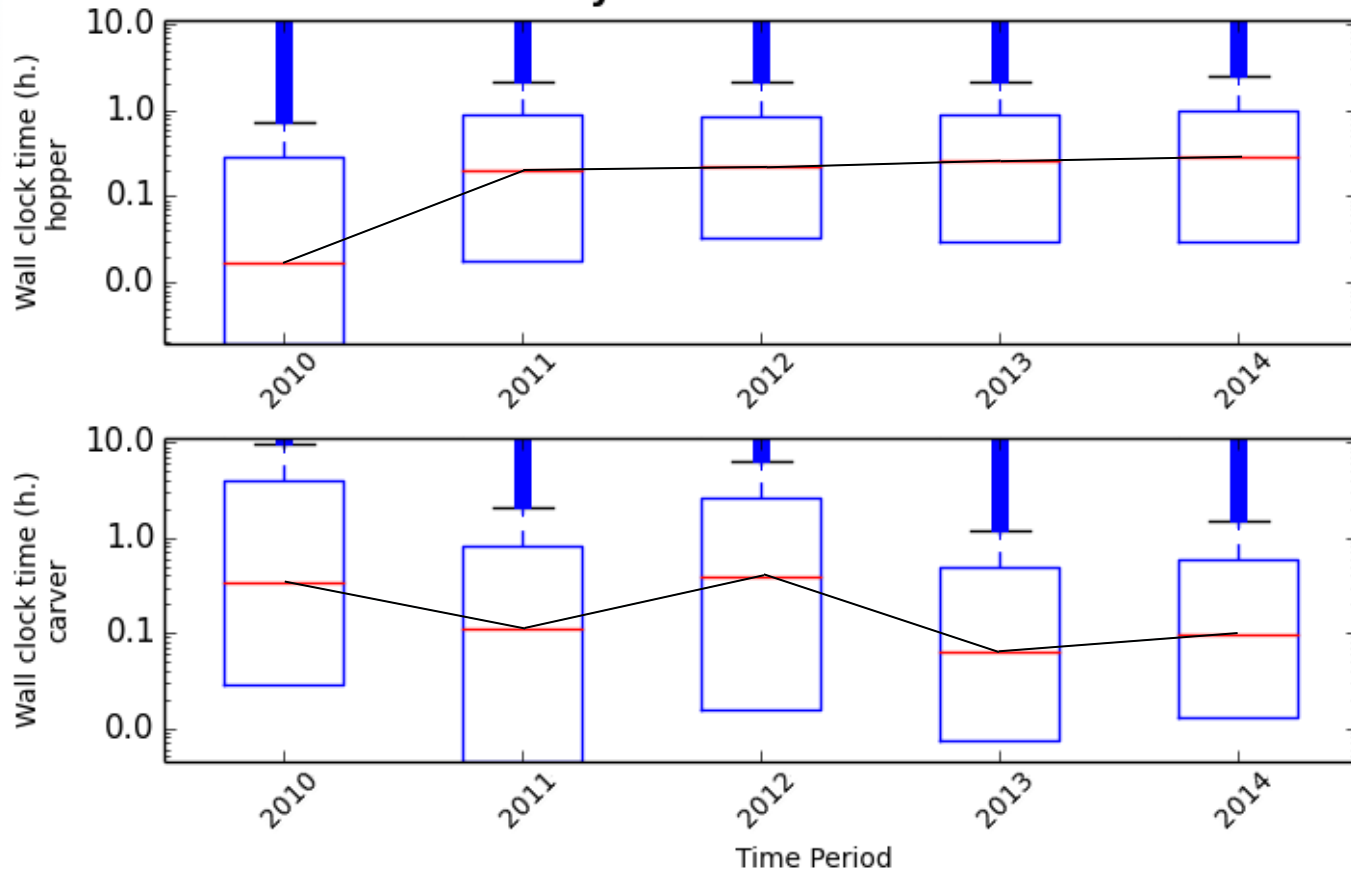
Temporary “expropriation” of resources assigned to dynamically malleable applications

## Expropriate and return actions



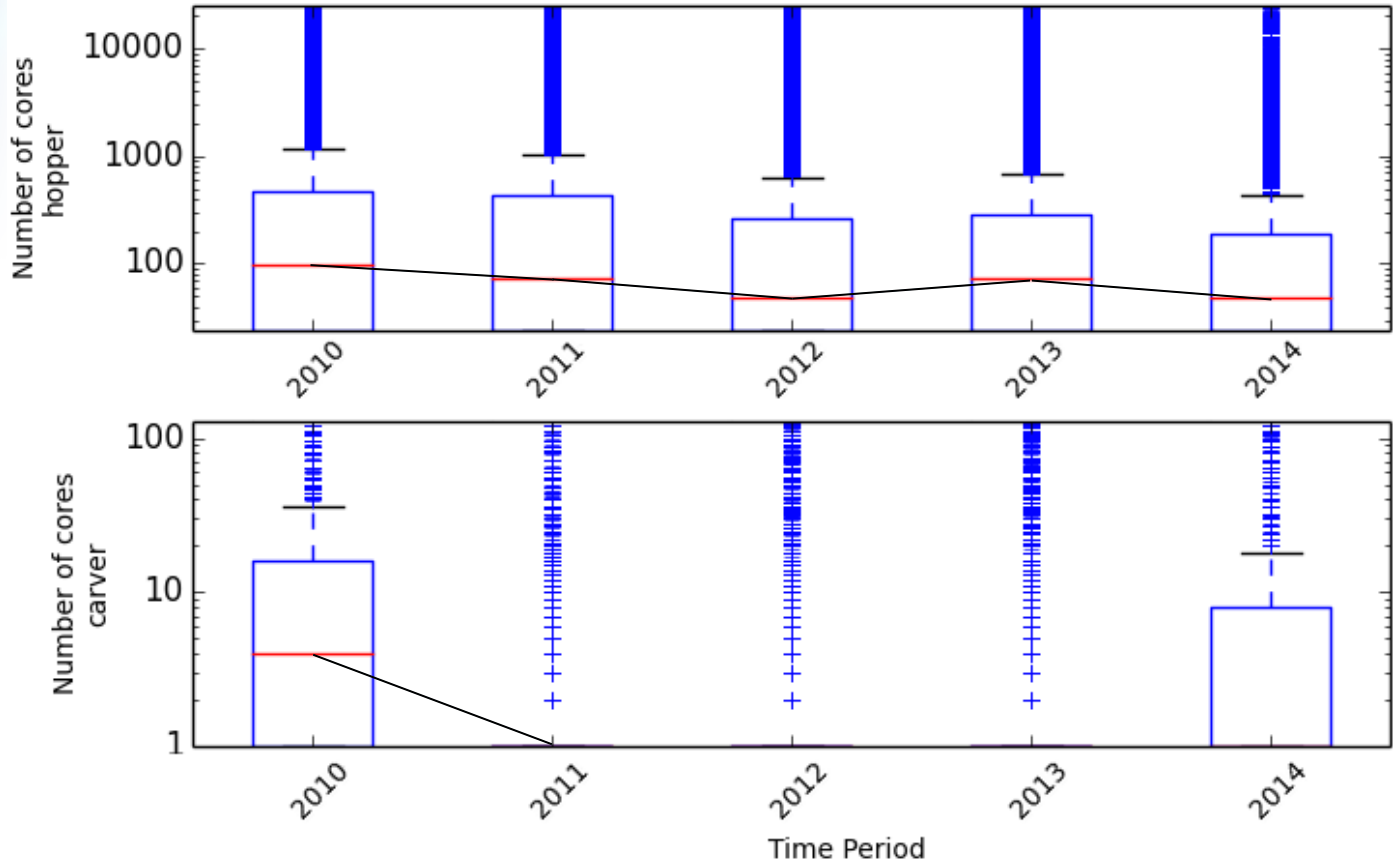
# Wall clock time

### Hopper Carver -Wall clock time per Job year evolution



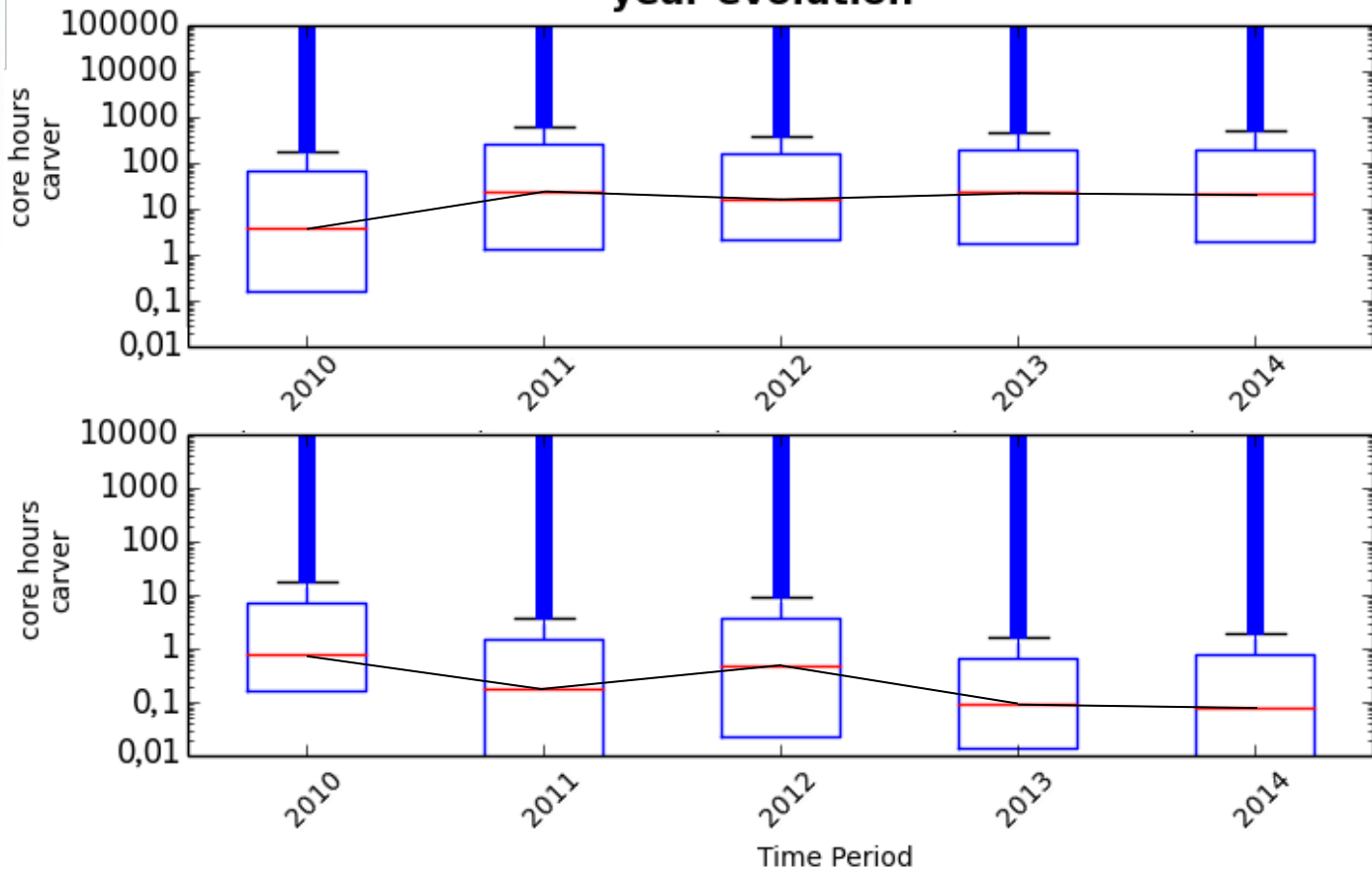
# Number of cores per job

### Hopper Carver - Number of cores per Job year evolution



# Core hours per job

### Hopper Carver - Core hours per Job year evolution



# Wait time per job

