# Scheduling for future HPC systems
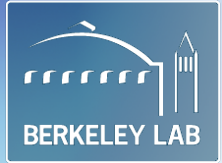
**Gonzalo P. Rodrigo** - **gonzalo@cs.umu.se**
Erik Elmroth – elmroth@cs.umu.se
Lavanya Ramakrishnan – lramakrishnan@lbl.gov
P-O Östberg – p-o@cs.umu.se

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

# What is this talk about?

We analyzed the lifetime workloads of three last Berkeley Lab's HPC systems: **Job and application heterogeneity is giving the scheduler a hard time.**

Surveyed trends toward Exascale: **How systems heterogeneity and extreme parallelization will challenge schedulers**.

We talked with HPC users using workflows:
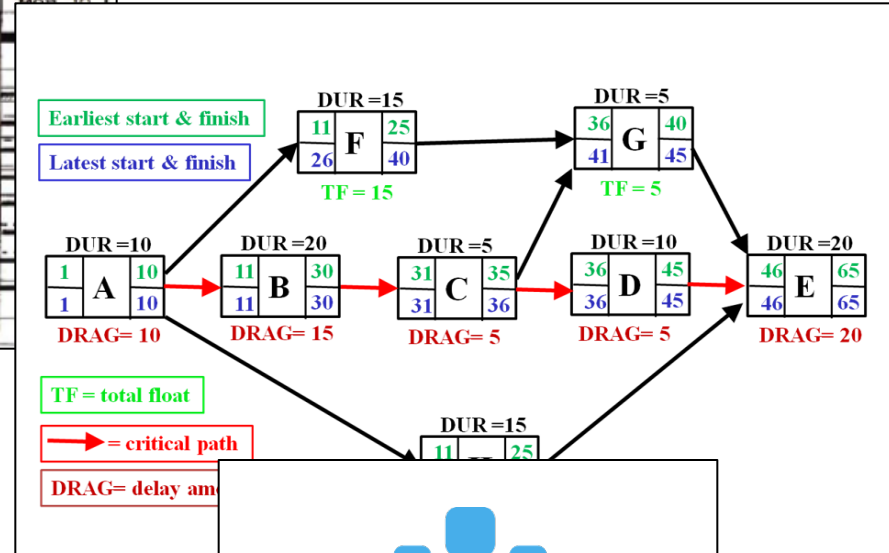**They wait forever or waste resources.
We propose a solution.**

We designed a cloud inspired **scheduling model to cope with systems and workload heterogeneity.**

**Present the toolset that made this scheduling research possible.**

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# Batch scheduling basics


[2]


[1]


[3]

Earliest start & finish
Latest start & finish
TF = total float
= critical path
DRAG= delay amount

TF = 15
TF = 5

DUR =15
11 F 25
26 40

DUR =5
36 G 40
41 45

DUR =10
1 A 10
1 10
DRAG= 10

DUR =20
11 B 30
11 30
DRAG= 15

DUR =5
31 C 35
31 36
DRAG= 5

DUR =10
36 D 45
36 45
DRAG= 5

DUR =20
46 E 65
46 65
DRAG= 20

DUR =15
11 25

slurm
workload manager

Man RECORD CHART FOR ____ DEPT.

| NAME | NO | Mon. 3 | Tues 4 | Wed. 5 | Thurs.6 | Frid. 7 | Sat 8 | Mon. 10 |
|------|-----|--------|--------|--------|---------|---------|-------|---------|
| PALEN | | | | 64% | | 5%A | 28% | |
| Griffen | 501 | | T | I | T | I | T | I |
| Palen | 503 | 6R | G | G | G | G | | |
| Millspaugh | 507 | | | | | | | |
| Owens | 514 | | | T | A | A | A | |
| Rogee | 517 | | | | R | | | |
| Williams | 519 | T | I | | T | | | |
| Martell | 527 | | | | | I | I | |
| Stewart | 535 | 6 | GR | G | G | G | G | |

Time Scale
Paper Strips
Clips to Hold the Strips
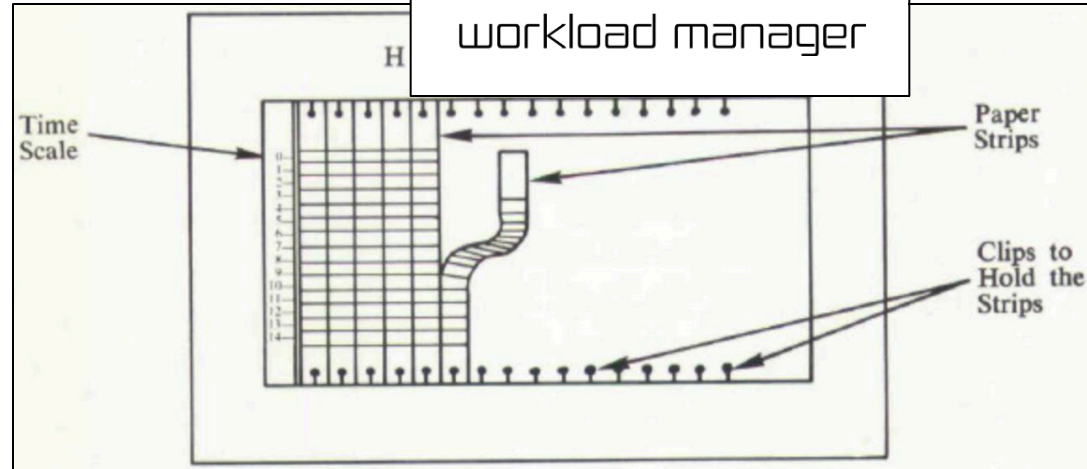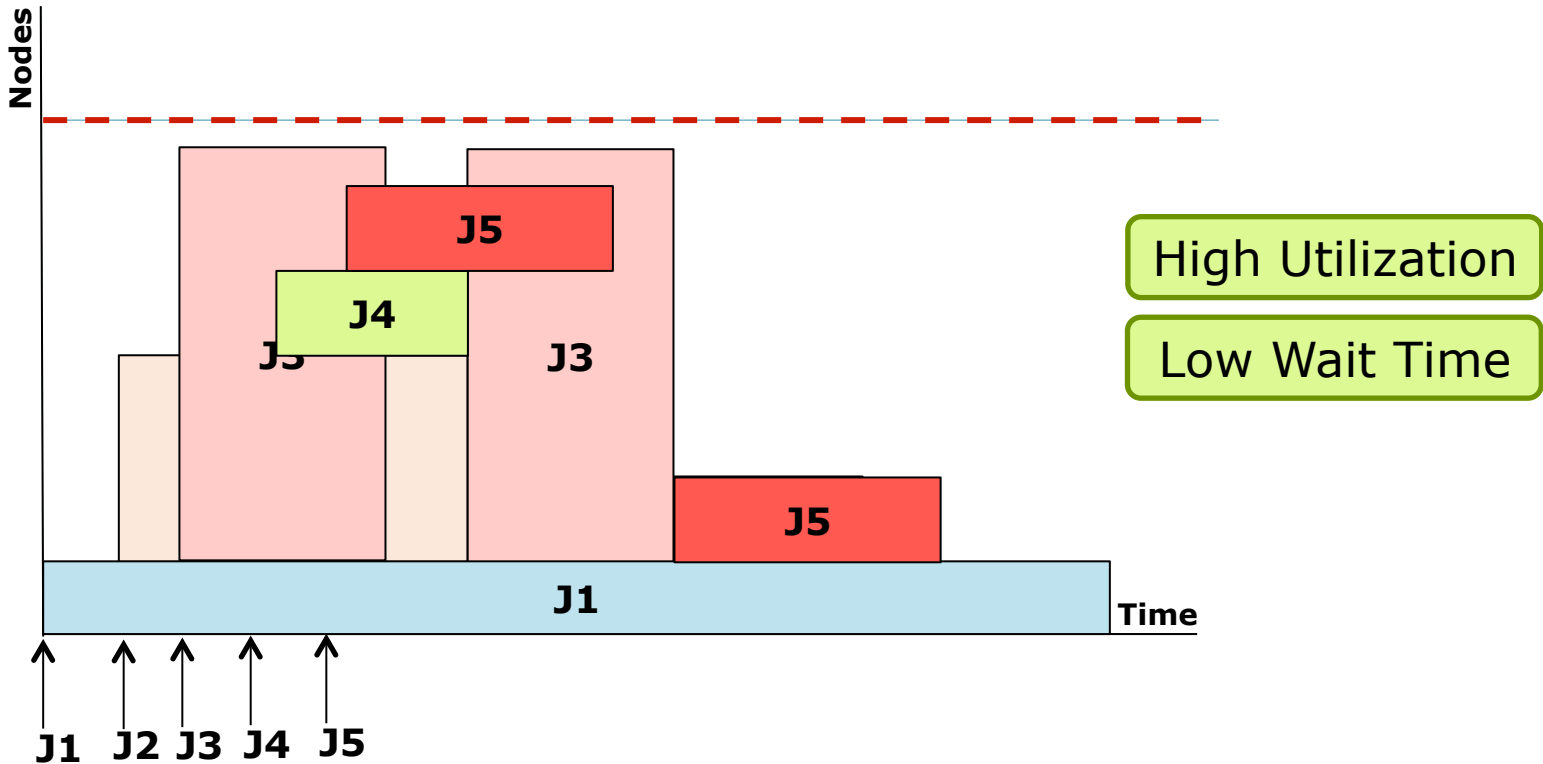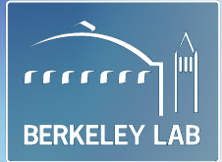H

[1] CPM example
[2] Original Gantt chart by Henry Gantt
[3] Harmonogram by Karol Adamiecki

**FCFS**: Jobs execute in arrival order

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

Mechanisms to reorder
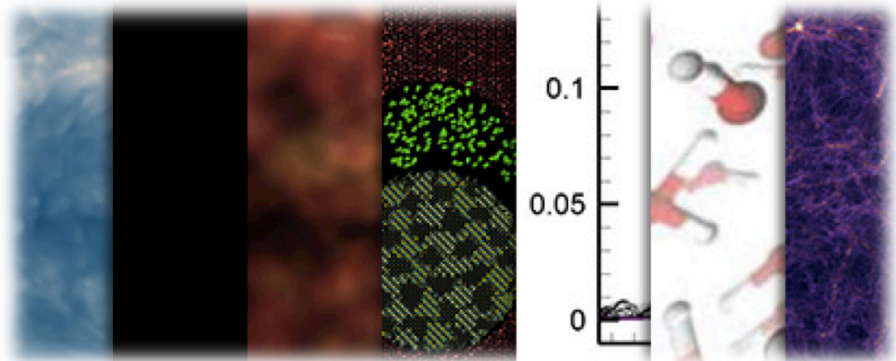the waiting queue

Fairness

Don't starve jobs or users

Priority

Run more important jobs first

**Gonzalo P. Rodrigo – gonzalo@cs.umu.se**

# Workloads and Systems

[1]

[2]

[1] Aurora Supercomputer: http://aurora.alcf.anl.gov
[2] Visualization elements from climate science, design accelerator design, biological research, transportation improvement, chemistry, and cosmology: http://aurora.alcf.anl.gov

# Present and future workload: Diversity and evolution



Job diversity

Data intensive Applications

Application diversity

Job heterogeneity in queues

Different performance Metrics

Different performance Model
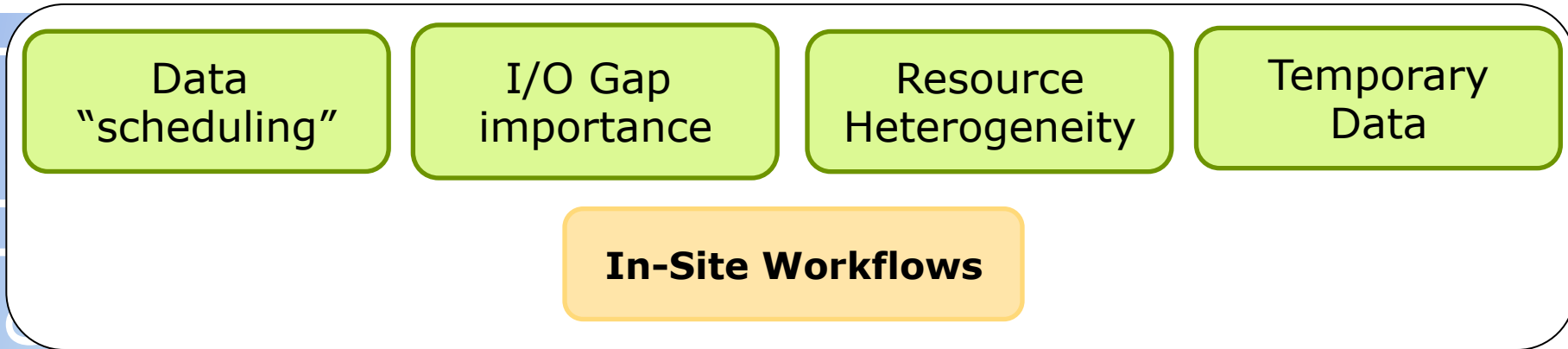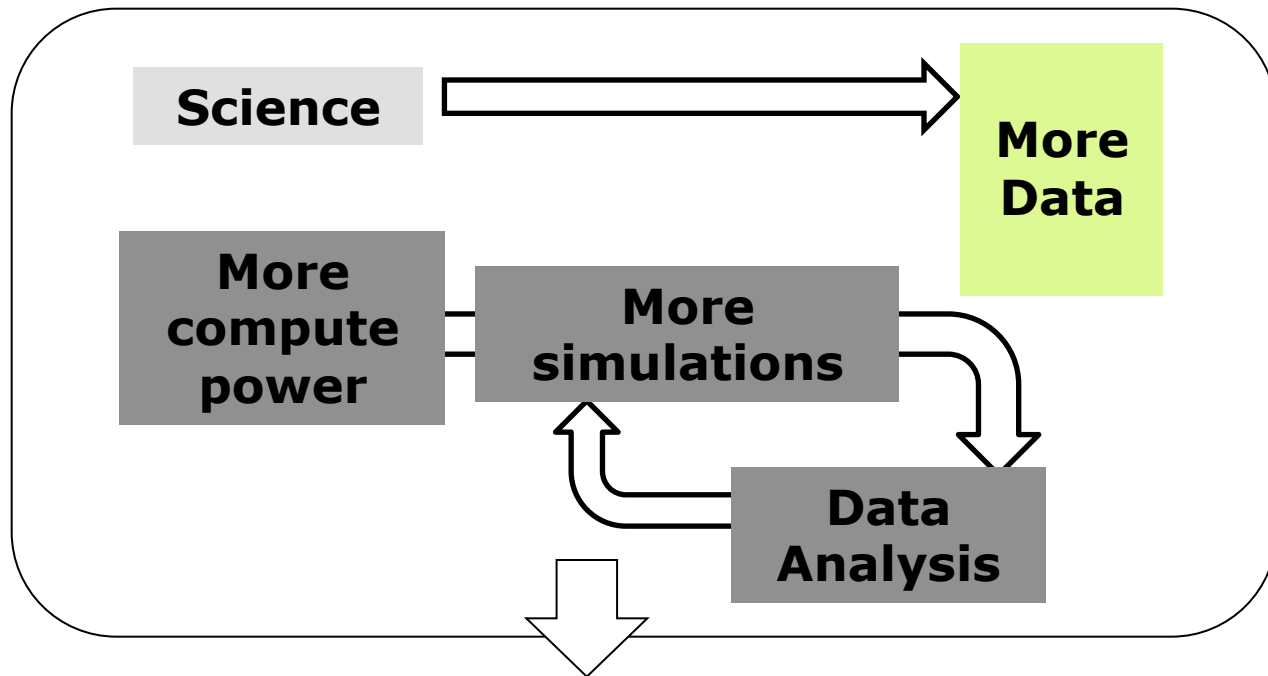
**Hard to predict wait times**

**Scheduler Comparing Apples & Pears**

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, 2016.
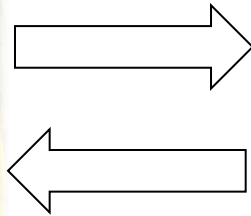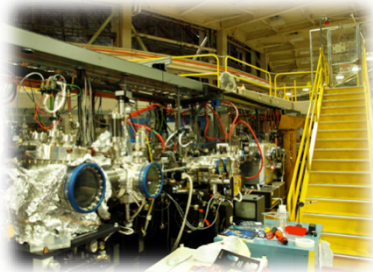
G. Rodrigo, P-O. Östberg, E. Elmroth, K. Antypas, R. Gerber, and L. Ramakrishnan. (2015, June). HPC System Lifetime Story: Workload Characterization and Evolutionary Analyses on NERSC Systems. In Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing (pp. 57-60)

Science → More Data

More compute power → More simulations → Data Analysis

Data "scheduling" | I/O Gap importance | Resource Heterogeneity | Temporary Data

**In-Site Workflows**

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

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

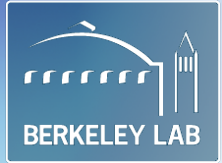**Real Time Applications**

**Bigger Systems**

**Grid workflows = In-Site Workflows**

# **Why Exascale?**
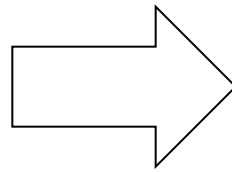
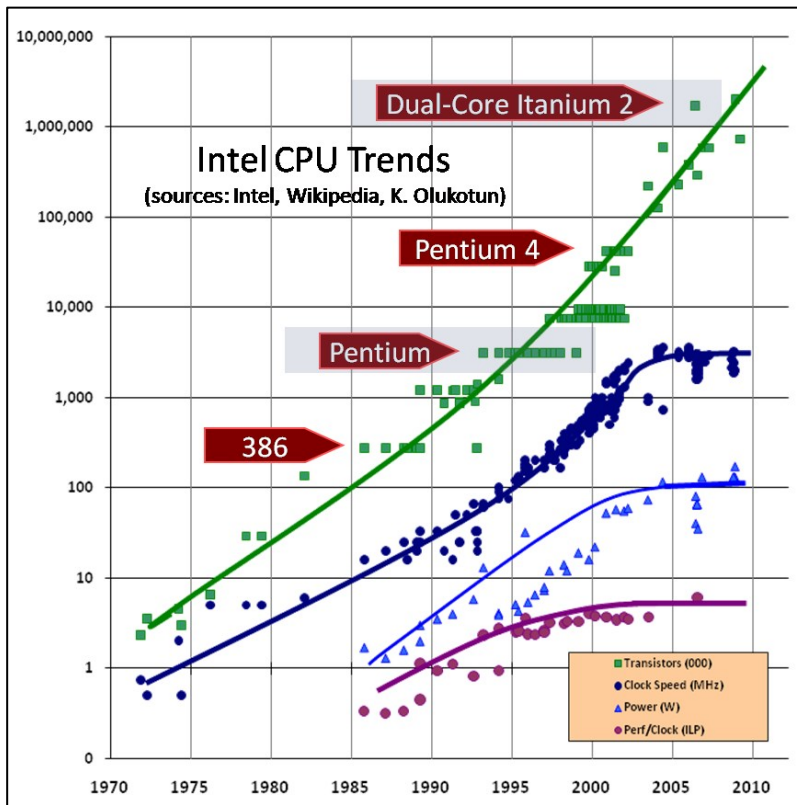Science is fueled by computation:
More power, more science.

Grid based simulations (e.g. climate)
require more resolution:
More parallelism.

Systems

# It's all about power and cost



Intel CPU Trends
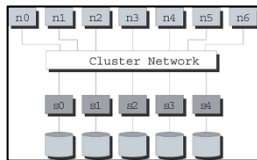(sources: Intel, Wikipedia, K. Olukotun)

Dual-Core Itanium 2

Pentium 4

Pentium

386

Transistors (000)
Clock Speed (MHz)
Power (W)
Perf/Clock (ILP)

[1]

**Higher degree parallelization**

## Break down of Dennard scaling

Systems

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# Raw Exaflops are possible by increasing the number of CPUs but...

Systems

| | | |
|---|---|---|
|  | **I/O** | Only scalable in parallel! |
|  | **RAM** | Power hungry! |
|  | **Interconnect** | More parallelism => More complexity **Less uniform latency** |

# The Exascale paradox

Heterogeneous HW

Per Thread

Very little

| | |
|---|---|
| Compute Power | → Operation Specific HW |
| RAM | → Coordination, more stages. **Workflows!** |
| PFS I/O BW | → Complex I/O Hierarchy: **Burst Buffers** |
| Network BW | → More in-chip Comms. **Placement** |
| Electric Power | → Reduced Resilience! |

Systems

# A pre-Exascale system: HPC2N's Kebnekaise

**Kebnekaise**

*Freshly deployed*

**432 classical compute nodes (12 096 cores)**

**20 large memory nodes (3 terabytes/node)**

**32 2xGPU Nodes (319 488 gpu cores)**

**4 4xGPU Nodes (79 872 gpu cores)**

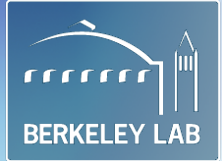**36 KNL Nodes (9 792 threads)**

**416 352 Cores**

**437 232 Threads**

**128 TBytes RAM**

**4 types of CPUS**

**Non Uniform Memory BW**

## Heterogeneity: memory, compute, interconnect, and programming models.

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

[1] **Sunway TaihuLight**
93.014 PFLOPS
US$273M
15 MW (No cooling)

**X 11** →

1 Exaflop
US$3003M
165 MW

**256+4 cores/CPU**

**Scratchpad memory**

**In Chip network**

**X3 Gflops/Watt**

**Bad HPCG Benchmark**

**Hard to program**

**Slow memory**

**Modest interconnect**

Systems

[1] Dongarra, Jack. "Report on the Sunway TaihuLight system." www.netlib.org. Retrieved June 20 (2016).

**Gonzalo P. Rodrigo – gonzalo@cs.umu.se**

System's Heterogeneity

Job diversity

**Real Time Applications**

Coordination, more stages,.

Application diversity

More in-chip Comms.

**Workflow management & scheduling**

In-Site Workflows

Data explosion

Reduced Resilience!

**Job and Application diversity**

Complex I/O

**Data scheduling (workflows?)**

Grid workflows = In-Site Workflows

Burst Buffers

Comparing Apples & Pears

**New time requirements**

Hard to predict wait times

# No Wait, No Waste: Workflow aware scheduling

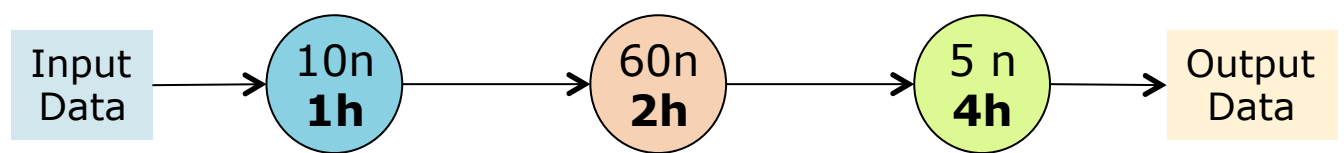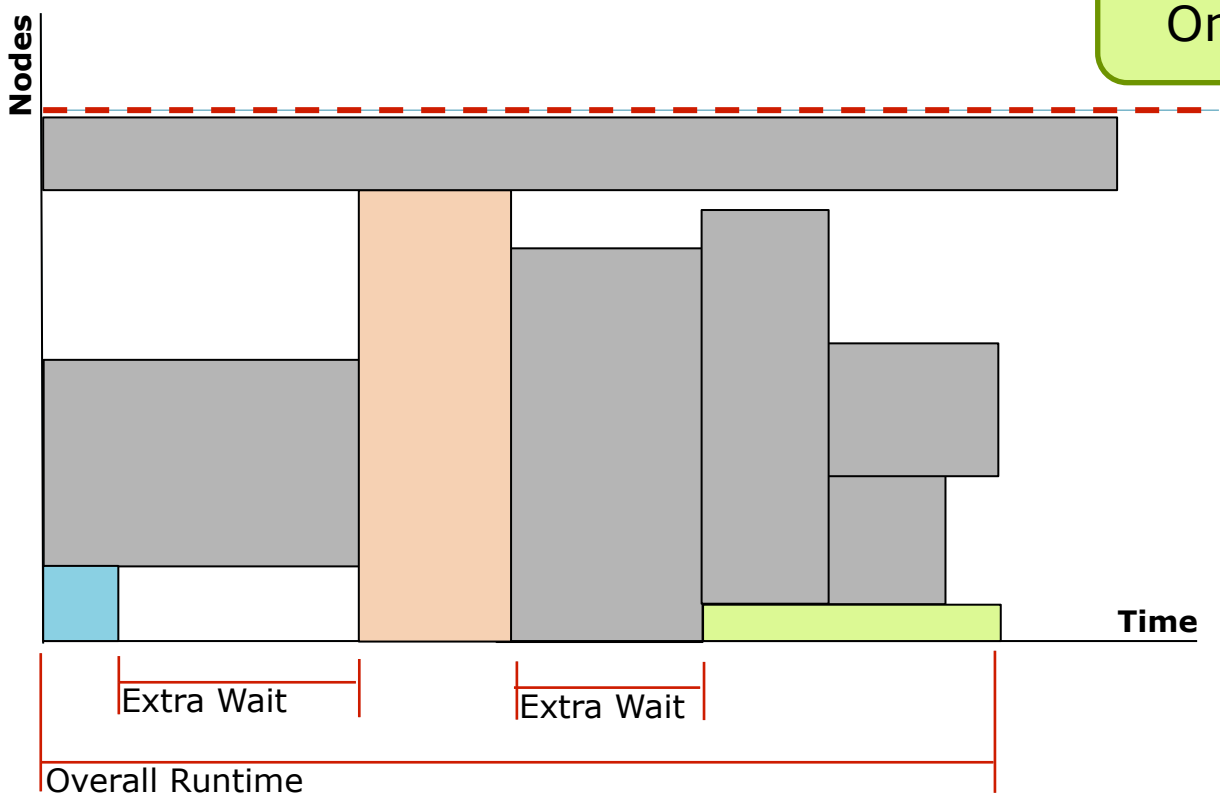# But before… What is a workflow?
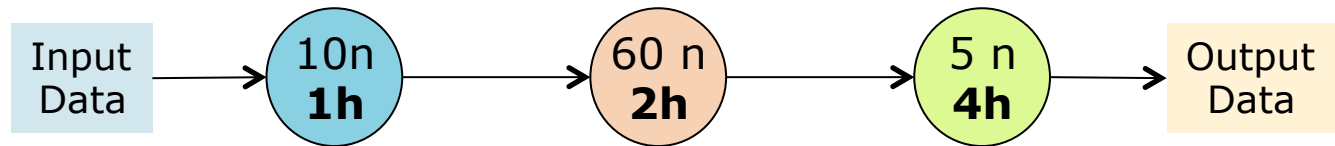
# But..

How does a scheduler deal with
**Workflows**?

# Submitting a workflow: Wait! (approach)

Input Data → 10n **1h** → 60n **2h** → 5 n **4h** → Output Data

One stage One Job

Nodes

Time

Extra Wait

Extra Wait

Overall Runtime

# Submitting a workflow: Waste! (approach)

Input Data → **10n 1h** → **60 n 2h** → **5 n 4h** → Output Data

One single Pilot Job

Nodes

Wasted Resources

Time

Overall Runtime

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# Workflow aware scheduling: Backfilling



Input Data → 10n **1h** → 60 n **2h** → 5 n **4h** → Output Data
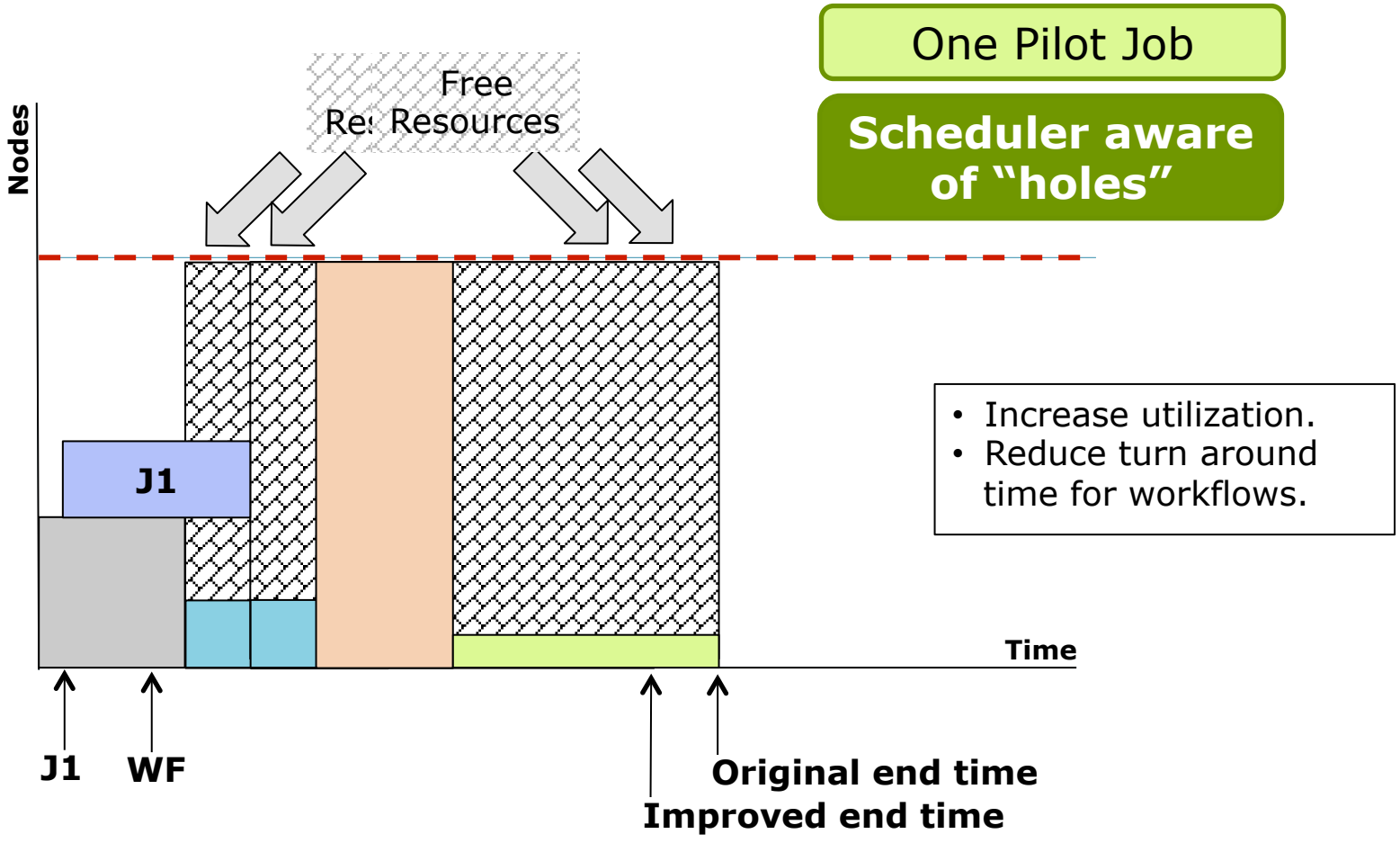
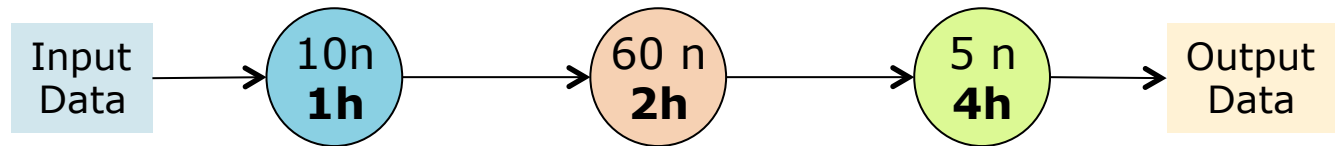One Pilot Job

Scheduler aware of "holes"

- Increase utilization.
- Reduce wait time of other jobs.

# Workflow aware scheduling: FCFS



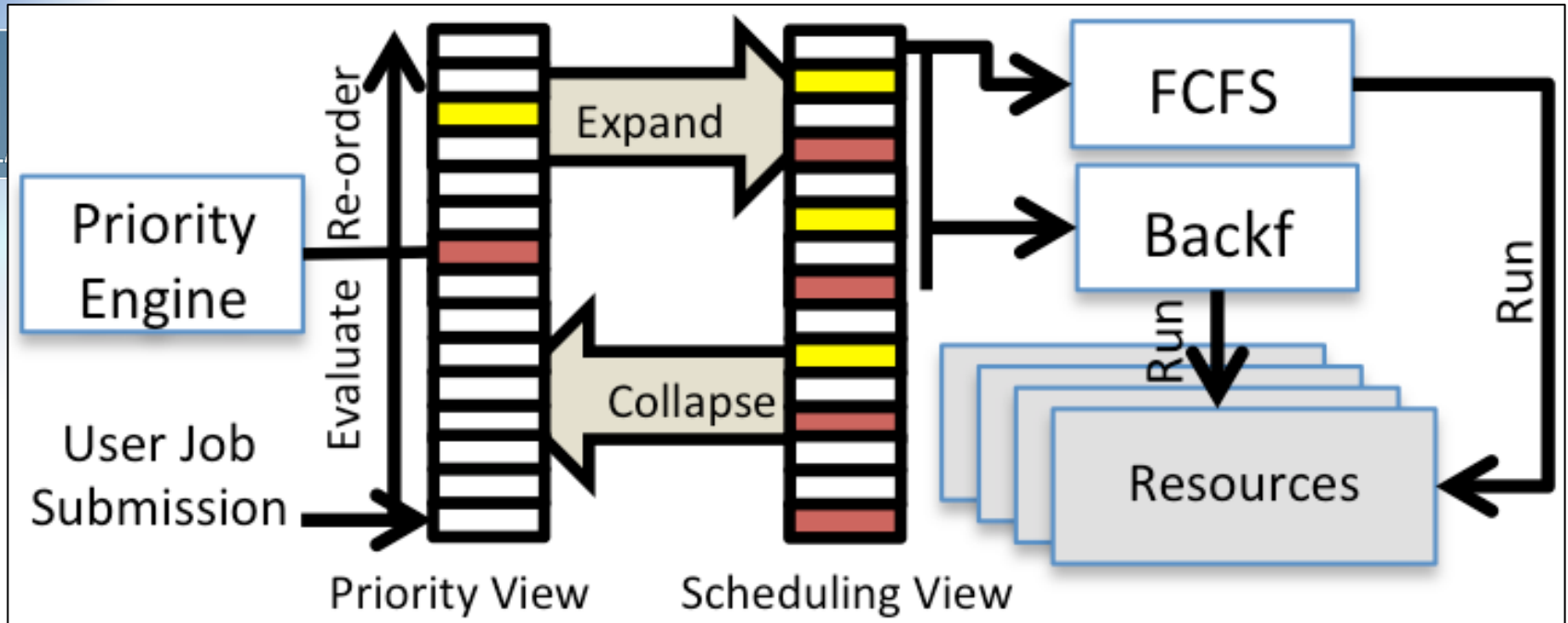Input Data → (10n **1h**) → (60 n **2h**) → (5 n **4h**) → Output Data

**One Pilot Job**

**Scheduler aware of "holes"**

Free Resources

Nodes

J1

- Increase utilization.
- Reduce turn around time for workflows.

Time

J1    WF

**Original end time**
**Improved end time**

**Scheduler Independent**

**Avoid System Gaming**

Wrap Scheduling Algorithms

Priority Workflow by Bounding job

# Experimentation: Modeling NERSC's Edison

**Edison's workload model** → **Workload Generator** ← **Workflow's description**

**Slurm Simulator** — Workflow Aware Slurm

**Workload Analyzer**

Cray XC30

Aries Network

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

2.57 Pflops/s

SLURM

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

First job vs. Workflow: turnaround time (s)

[1]

**Monolithic**     **Two-level**     **Shared state**

scheduling logic

subset

cluster state information

cluster machines

**no concurrency**     **pessimistic concurrency (offers)**     **optimistic concurrency (transactions)**

full state

[3]

# A2L2: Long term proposal


[2]

High priority

Real-time    $P_9$ → $P_7$ → $P_3$ → $P_1$    preempted

System    P → P → P → P    P    preempted

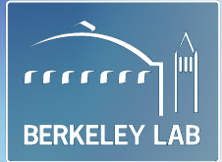Interactive    P → P → P → P    preempted

Batch    P → P → P → P    preempted

Low priority

[1] Scheduling taxonomy:  Schwarzkopf, Malte, et al. "Omega: flexible, scalable schedulers for large compute clusters." Proceedings of the 8th ACM European Conference on Computer Systems. ACM, 2013.

[2] Multilevel Queues: https://www.cs.rutgers.edu/~pxk/416/notes/07-scheduling.html

[3] 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)

Cloud infrastructures have faced similar challenges…



**Hypothesis:** Cloud scheduling techniques can be applied to tackle new HPC challenges.
**Method:** Compared study on techniques and application circumstances (Survey)

# Similarities

**Cloud**

**Edison HPC**

| Heterogeneous Workload | Heterogeneous Workload |
|---|---|

| Batch Jobs | Data is Key | Many non tightly coupled | Non-classical HPC |
|---|---|---|---|

| Wait Time is important | Response time |
|---|---|

| SSDs on Nodes | Burst Buffer |
|---|---|
| Distributed Filesystems | |

| Heterogeneous resources | Accelerator HW | BB nodes | Compute nodes |
|---|---|---|---|

**Infrastructure**

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

**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)

Position Paper
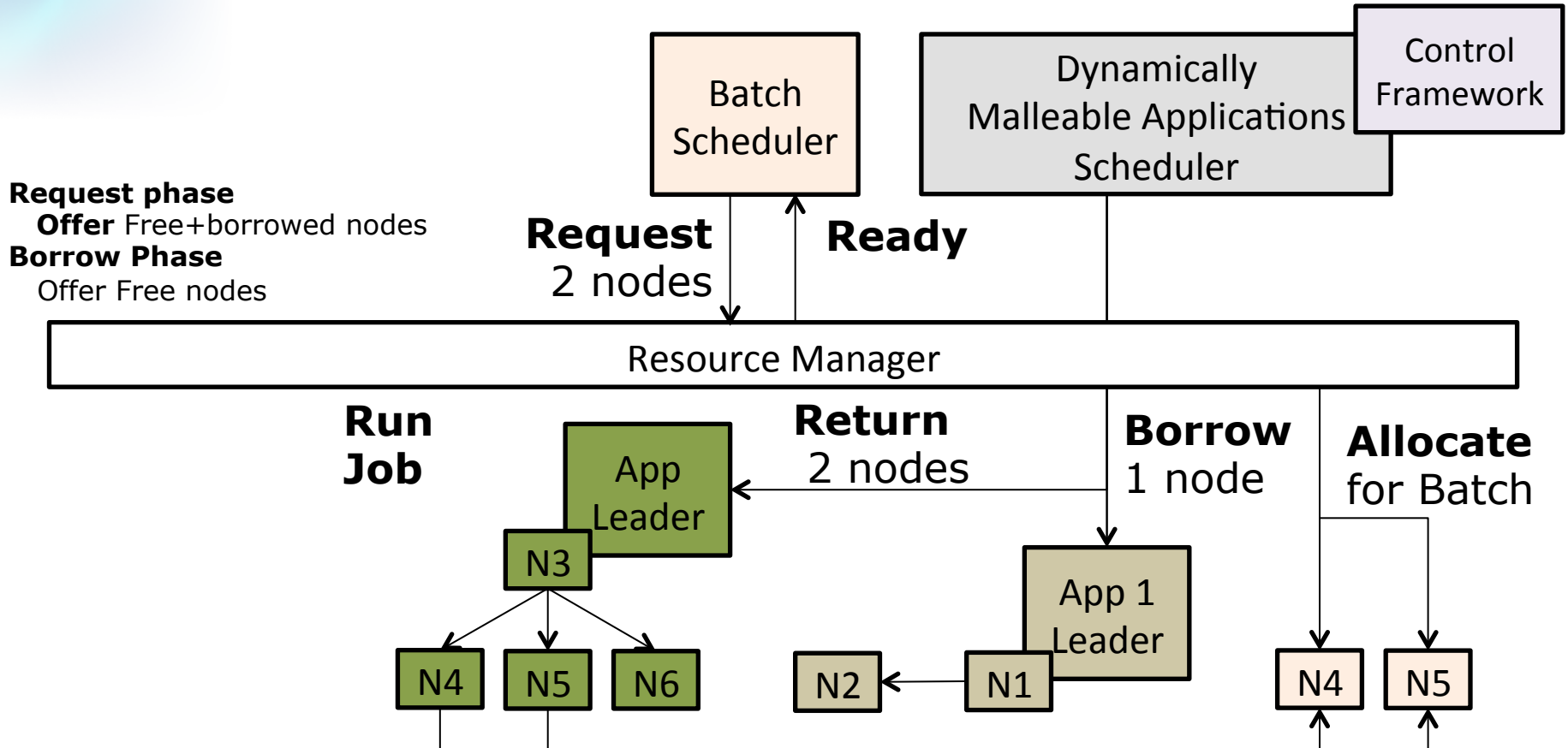
Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# Scheduler model

Cloud borrowed solution: **Two level scheduling**
One scheduler per application + smart RM
Malleable Applications: Dynamic allocation
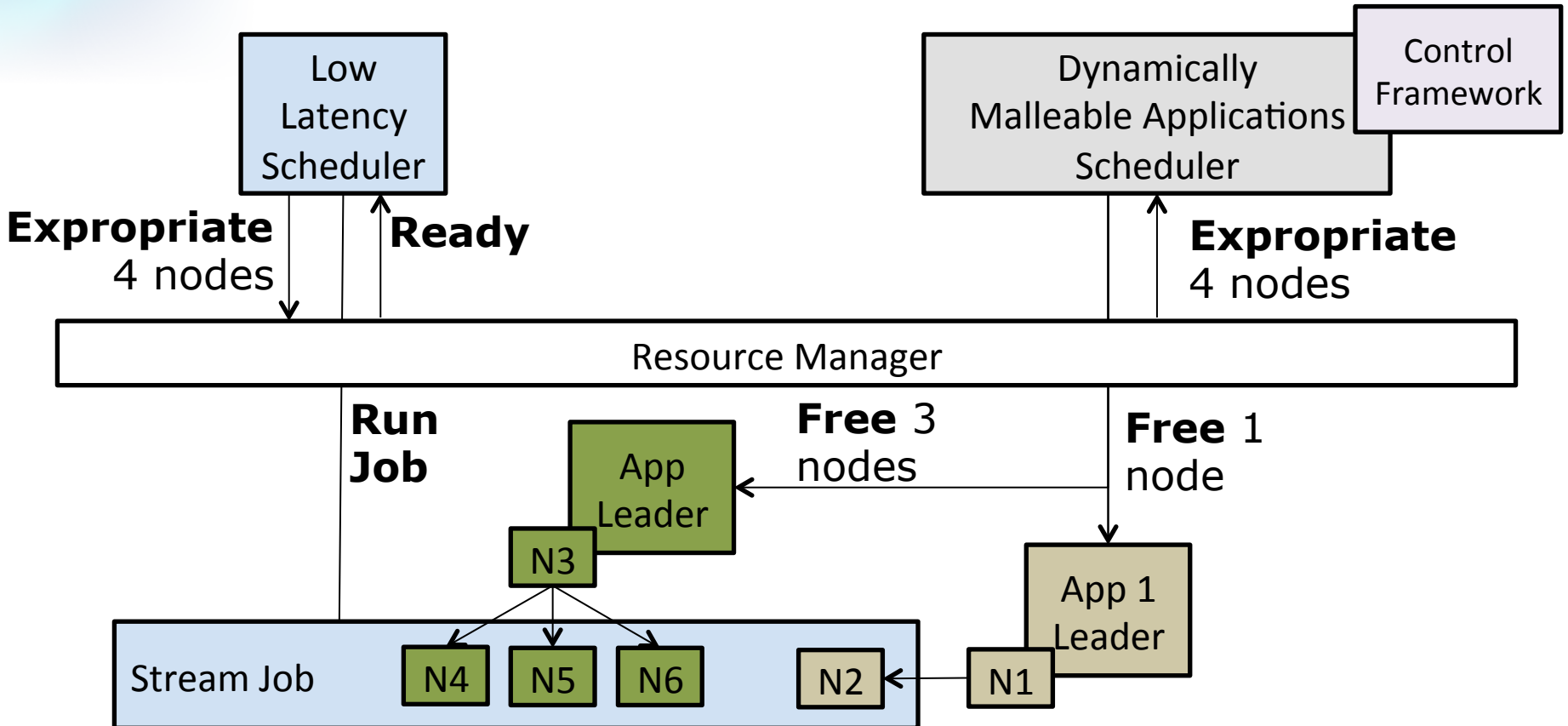Low latency allocation

**Request phase**
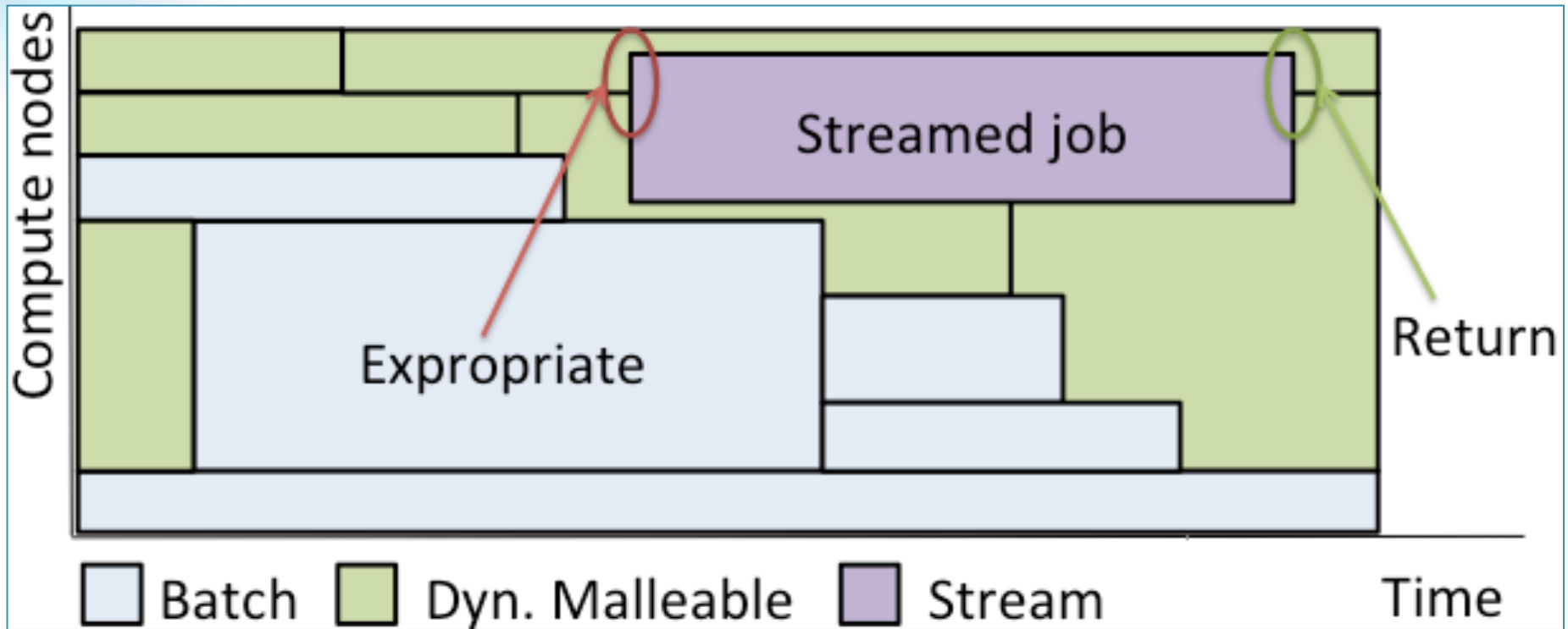   **Offer** Free+borrowed nodes
**Borrow Phase**
   Offer Free nodes

Batch Scheduler

Dynamically Malleable Applications Scheduler

Control Framework

**Request** 2 nodes

**Ready**

Resource Manager

**Run Job**

**Return** 2 nodes

**Borrow** 1 node

**Allocate** for Batch

App Leader

N3

N4  N5  N6

App 1 Leader

N2  N1

N4  N5

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# Flexible backfilling



Batch    Dyn. Malleable

Temporary "expropriation" of resources assigned assigned to dynamically malleable applications

## Expropriate and return actions

Application heterogeneity is a trait of both cloud and HPC applications

Application Aware

Two level scheduling

Flexible nature of malleable applications can be useful (and there maybe enough malleable workload to make be useful)
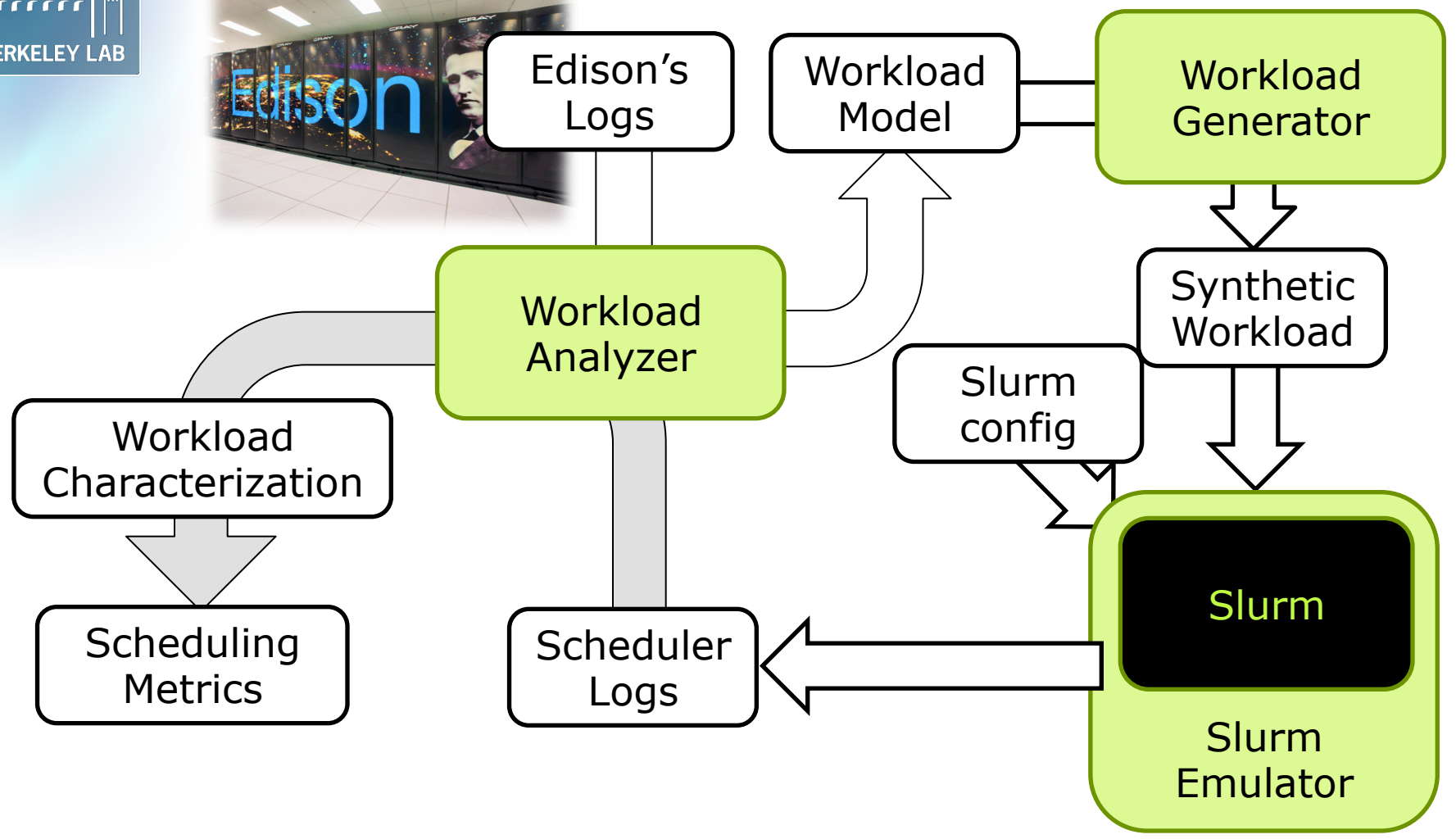
Application Management

Better utilization

Stream job allocation
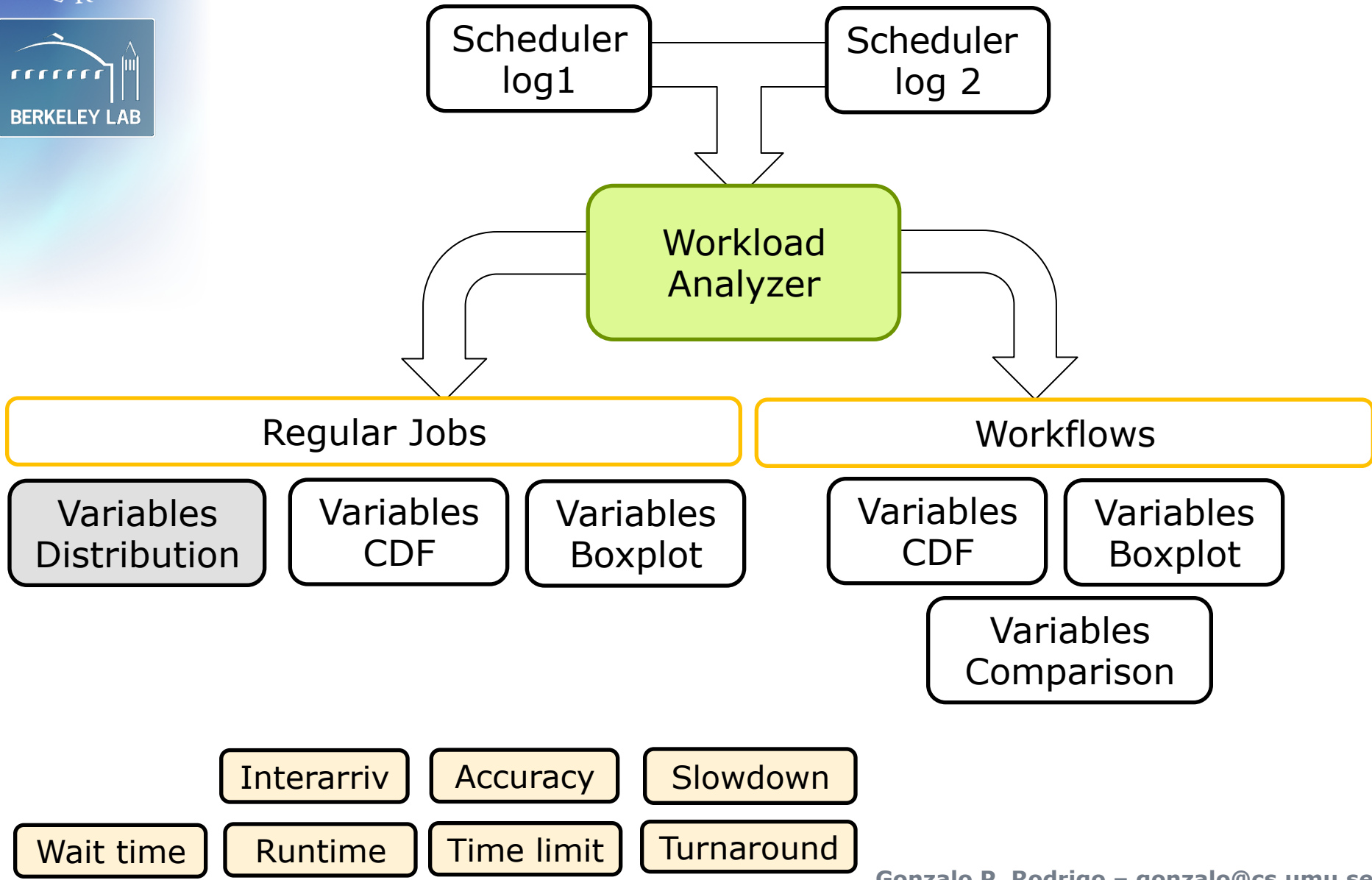
# Scheduling: Challenges Research & Operational

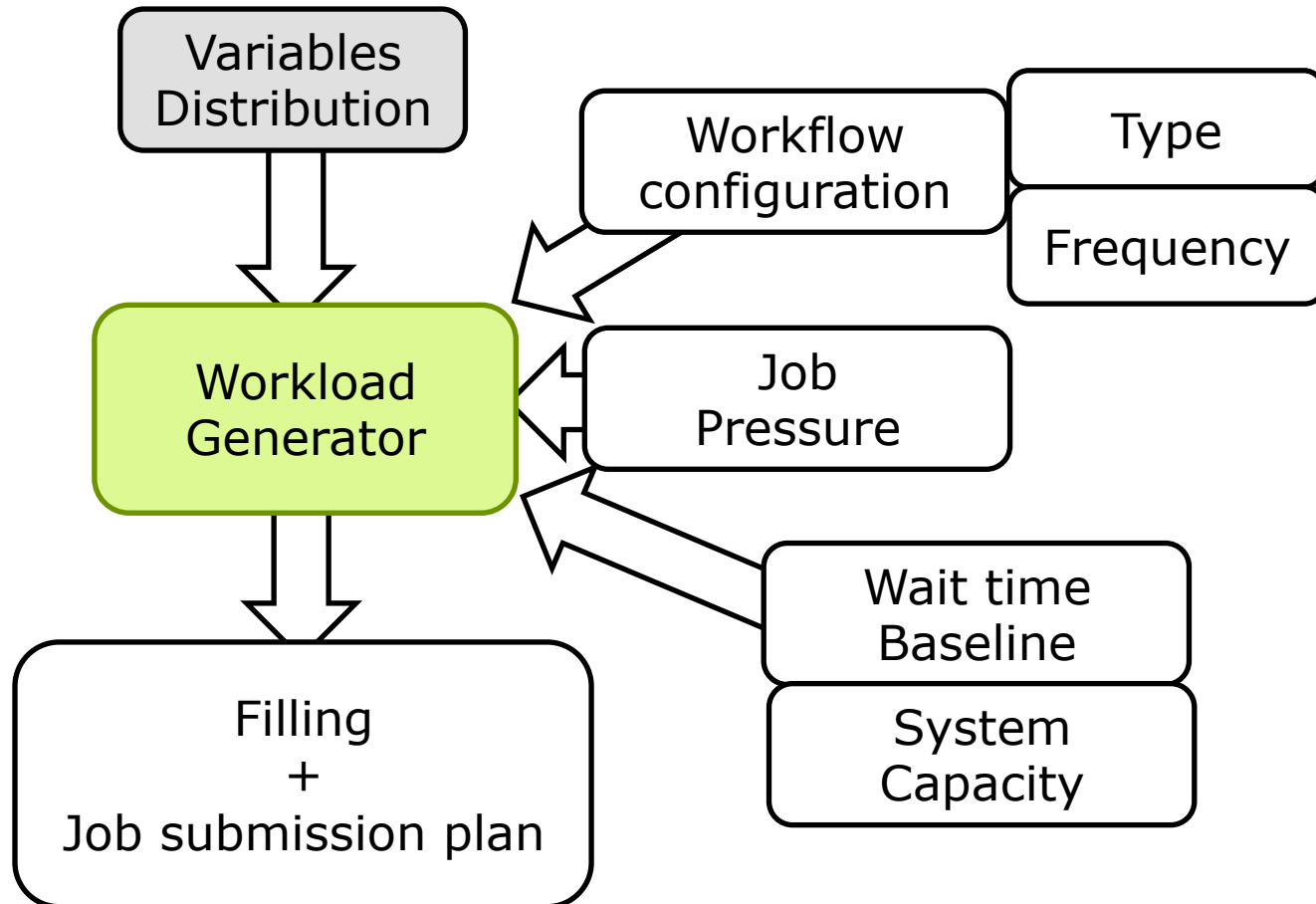# Our work required too much engineering…

# Workload analyzer



Scheduler log1 — Scheduler log 2

Workload Analyzer

Regular Jobs | Workflows

Variables Distribution | Variables CDF | Variables Boxplot

Variables CDF | Variables Boxplot

Variables Comparison

Interarriv | Accuracy | Slowdown

Wait time | Runtime | Time limit | Turnaround

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# Workload generator

# Slurm Simulator

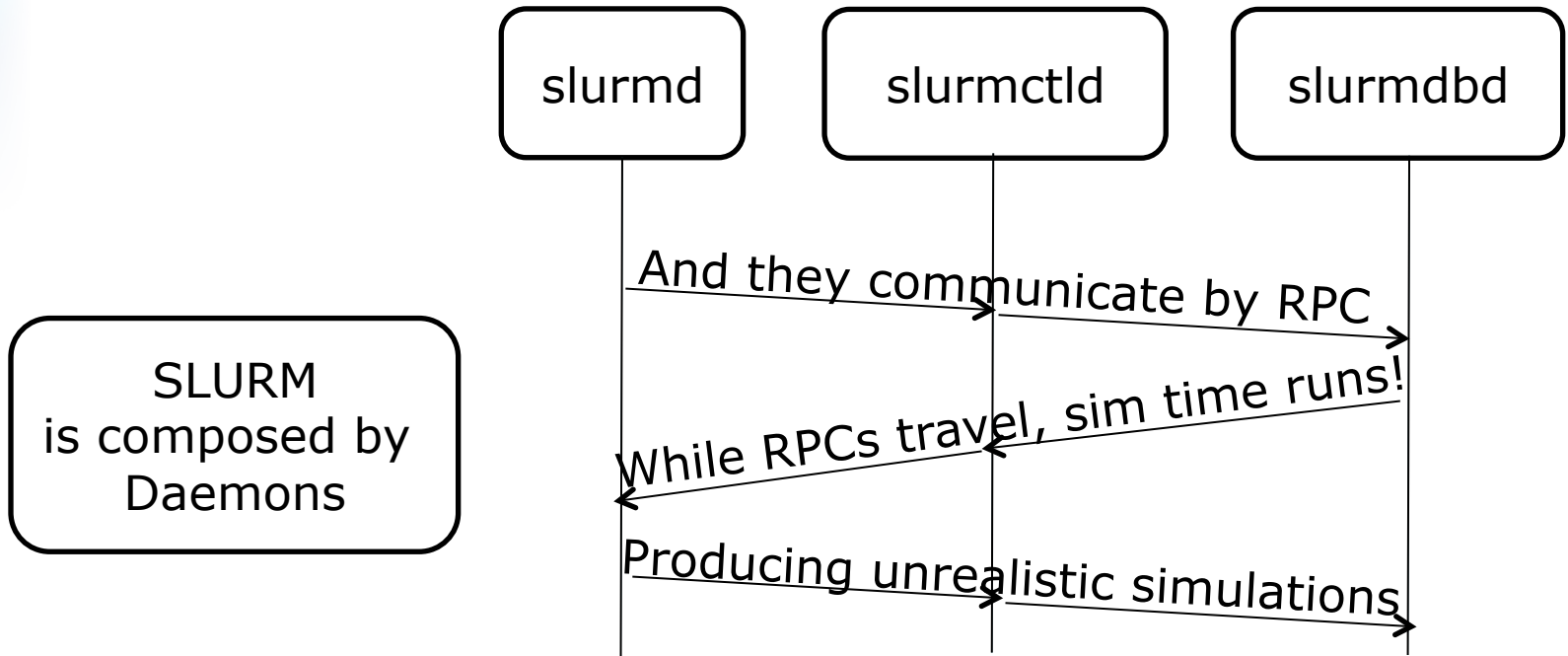Slurm

**Slurm Emulator**

Based on BSC original work

Based on CSCS work

Still, it required quite some work!

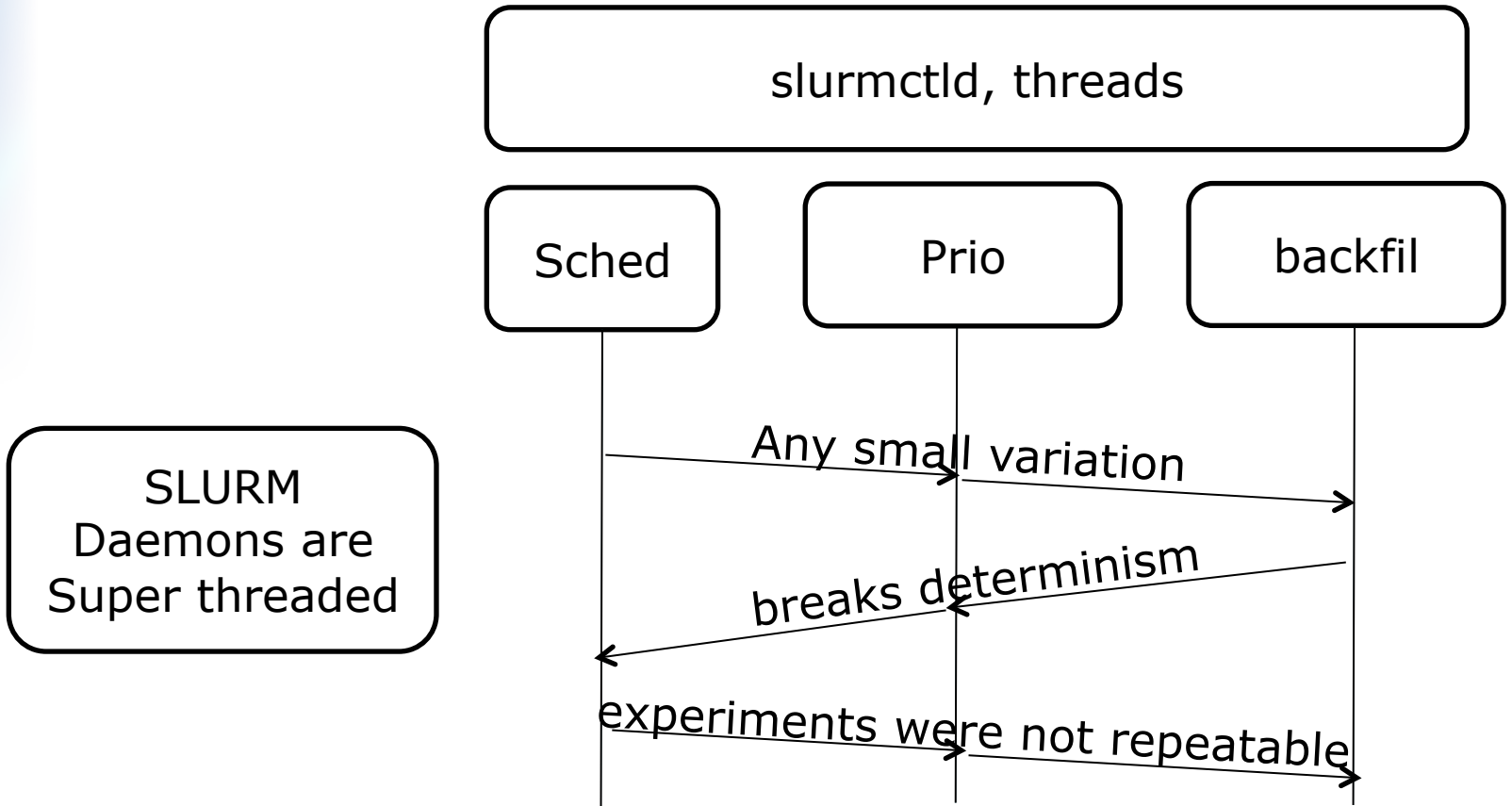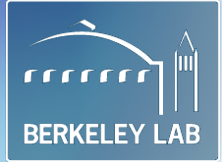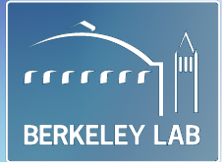## Functions

- Emulate Hardware
- Time speed-up: By hacking time/sleep calls
- Job submission from submission plan

slurmctld, threads

Sched

Prio

backfil

SLURM
Daemons are
Super threaded

Any small variation

breaks determinism
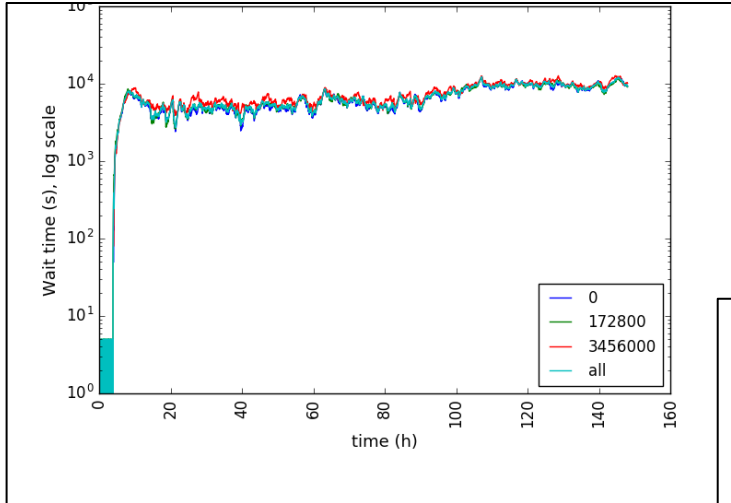
experiments were not repeatable

# Slurm Simulator: The work

- sim_mgr and scheduling loops RPC synchronized
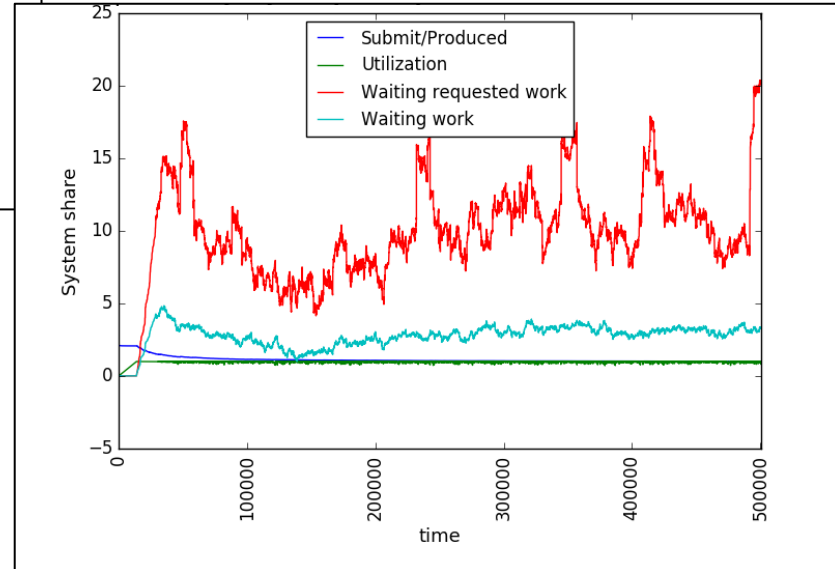- Every significant thread is synced with the simulation controller.

- Real vs simulation time: x10-x20 speedup
- Priority, Scheduling, and Accouting are synced.
- Scheduler can achieve high utilization.
- Results are "semi" deterministic.

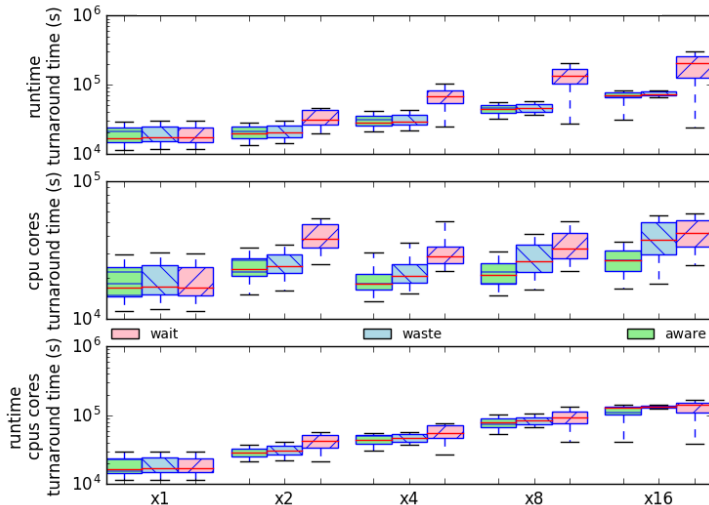# Workload analyzer: Evaluating scheduler

**Wait time evolution**

**Utilization, submitted, pending**
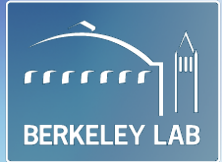
**Workflows behavior**

# Schedulling tools: Value

- To be fully described in upcoming paper
- Analysis, generation, running, and workflow aware scheduler **will be open sourced**

- Scheduling research: save "engineering hours"
- Admins: capacity to play with configurations of **their own systems and their own workloads.**
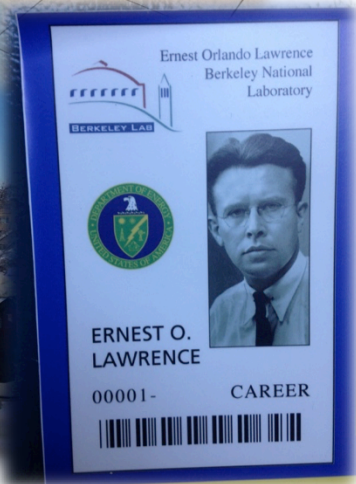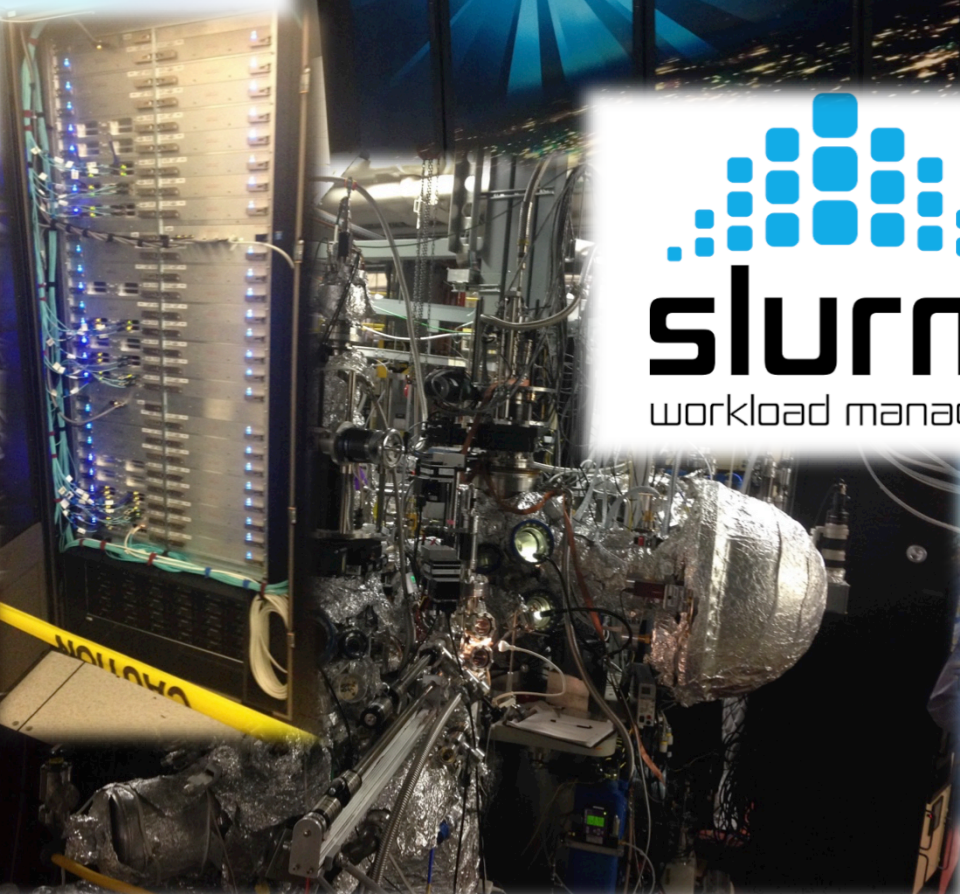
# Summary of take-aways

Systems and workloads require new scheduling: **We propose two-level cloud inspired model.**

In-site workflows are very important: **There is a better way to schedule them.**

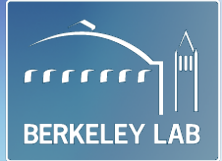Good **tools** are fundamental **for good research and operations**.

Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# Thanks for your time… questions?



Gonzalo P. Rodrigo – gonzalo@cs.umu.se

# To know more....

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

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

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.

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

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