

# ScSF: A **S**cheduling **S**imulation **F**ramework

JSSPP'17

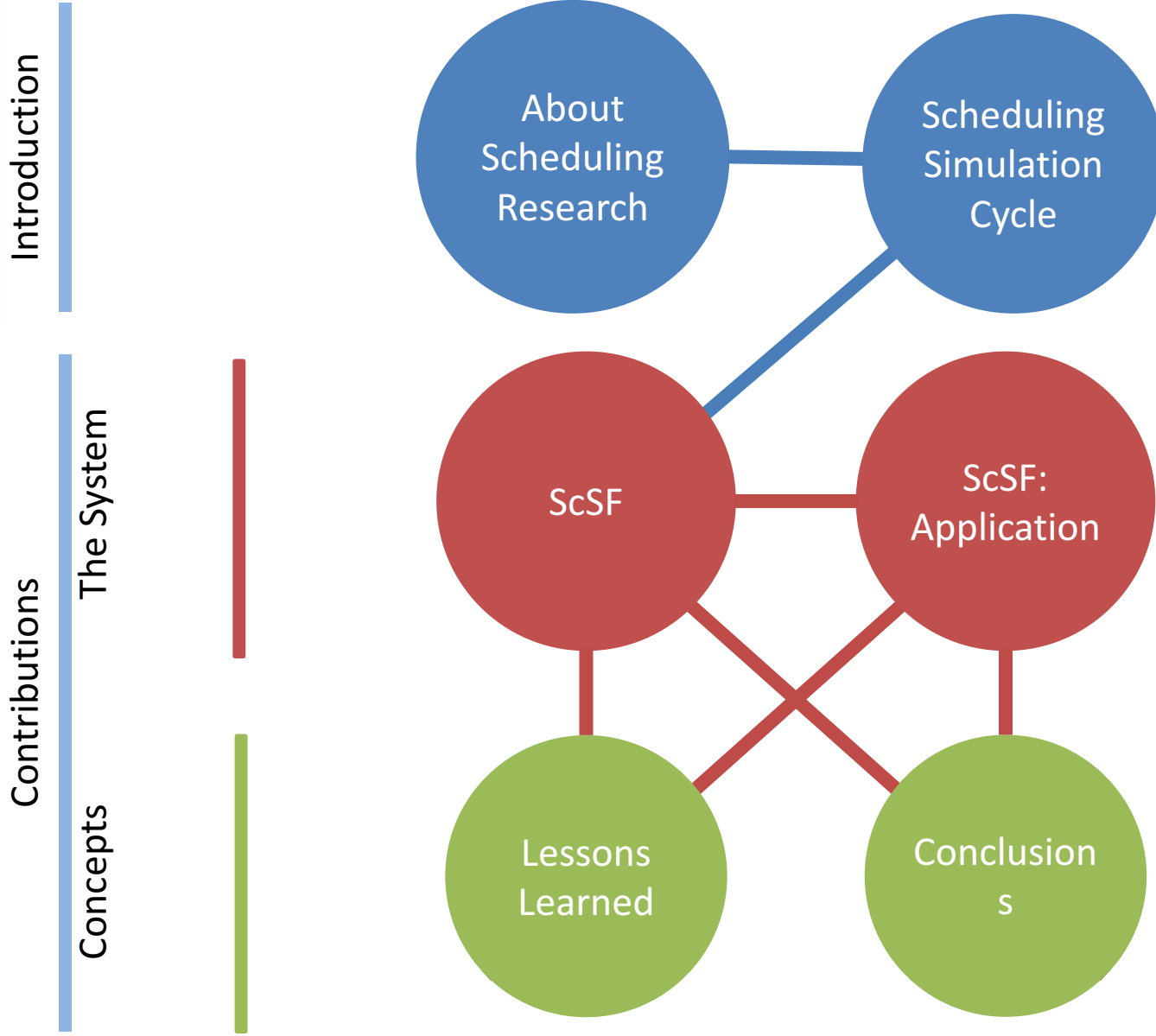
**Gonzalo P. Rodrigo Álvarez**

`gprodrigoalvarez@lbl.gov`

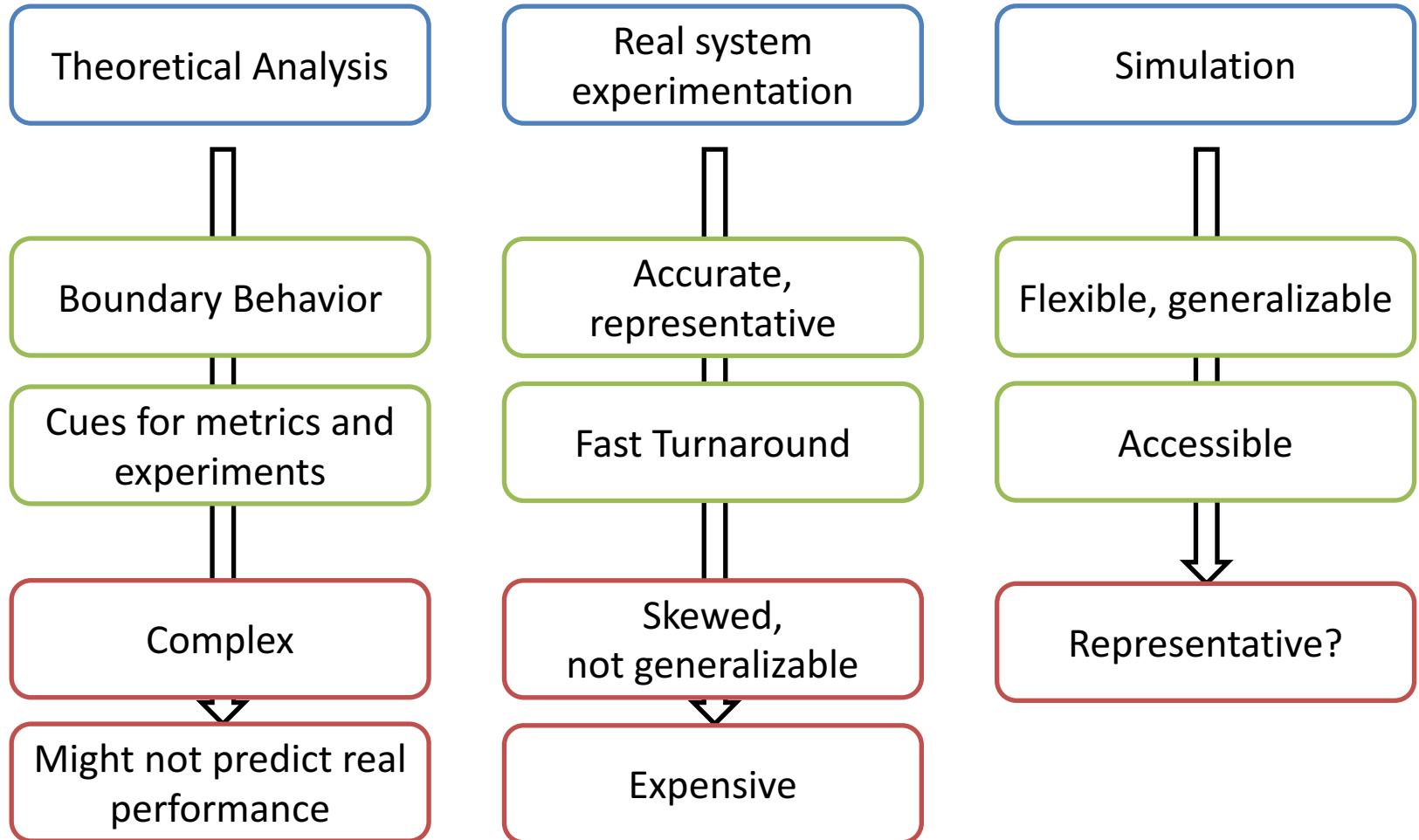
ScSCF is open and **available at:**  
<http://frieda.lbl.gov/download>



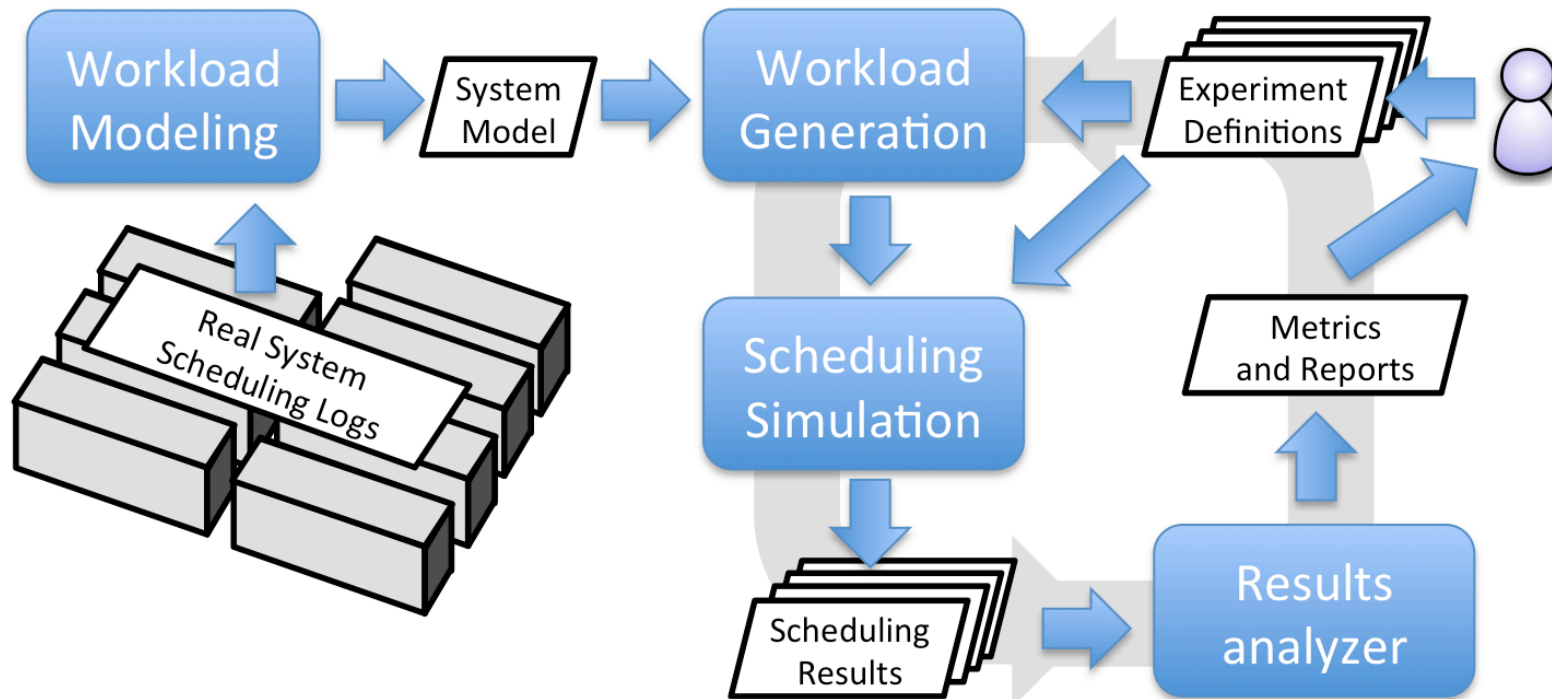
June 2<sup>nd</sup>, Orlando, Florida



# HPC Scheduling Research approaches

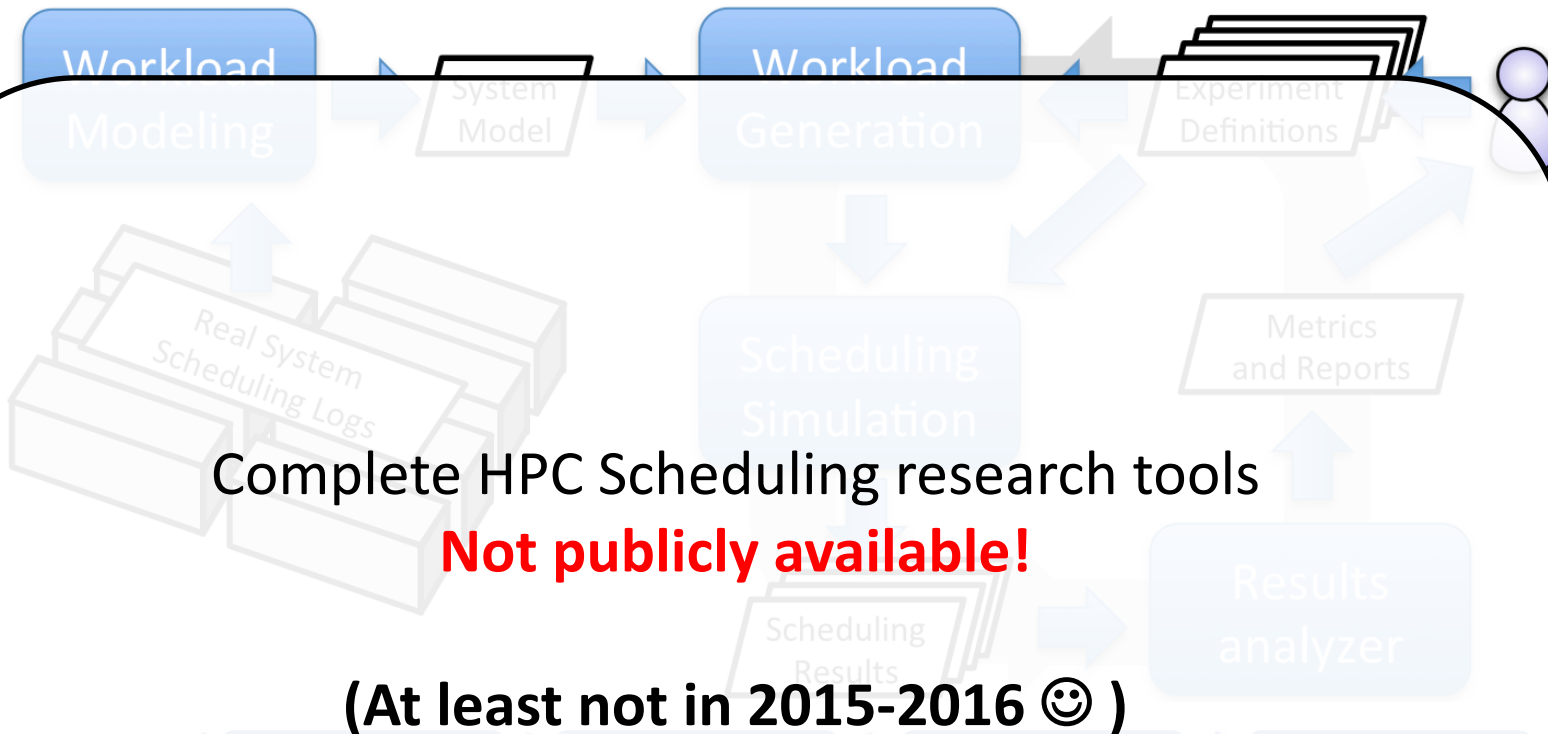


# HPC Scheduling Simulation: Research cycle



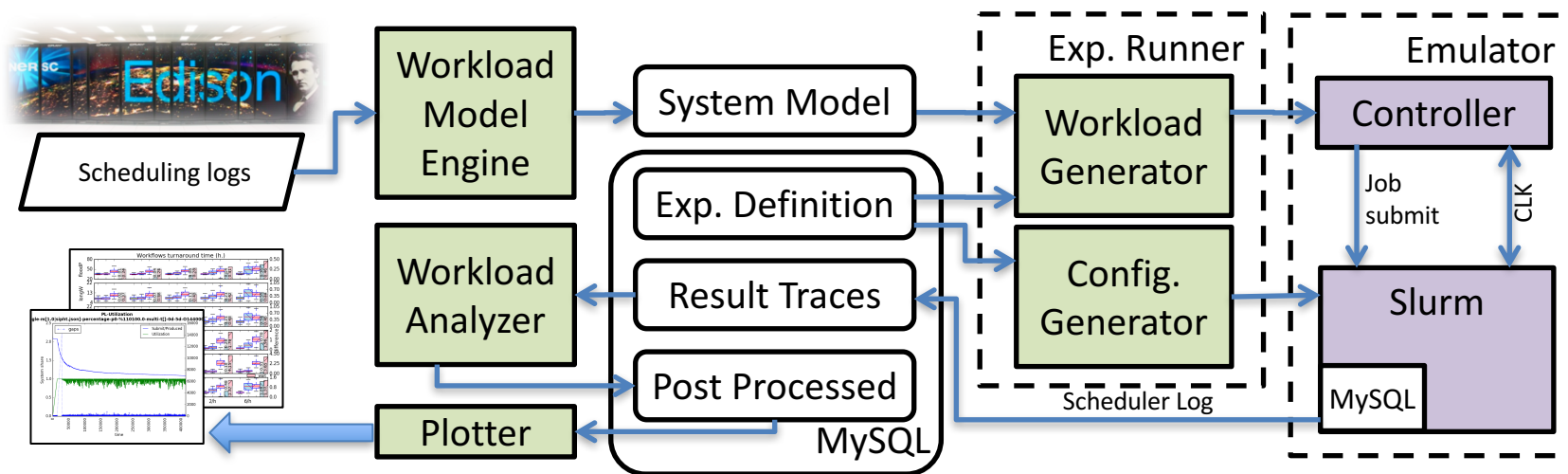
	Modeling	Generation	Simulation	Analysis
<b>Alea/BatSim</b>	✗	✗	✓ Not prod. scheduler	✗
<b>Slurm Simulator</b>	✗	✗	✓ Slow	✗
<b>Parallel Archive</b>	✗	✓ Old Small	✗	✗

# HPC Scheduling Simulation: Research cycle



	Modeling	Generation	Simulation	Analysis
Alea/BatSim	✗	✗	✓ Not prod. scheduler	✗
Slurm Simulator	✗	✗	✓ Slow	✗
Parallel Archive	✗	✓ Old Small	✗	✗

# ScSF: Scheduling Simulation Framework



Modeling

Real  
System Logs

Generation

Synthetic Jobs  
System Model

Experiment  
Conditions

Simulation

Real Scheduler:  
Slurm

20x Time

Analysis

Metrics  
calculation

Aggregate  
repetitions

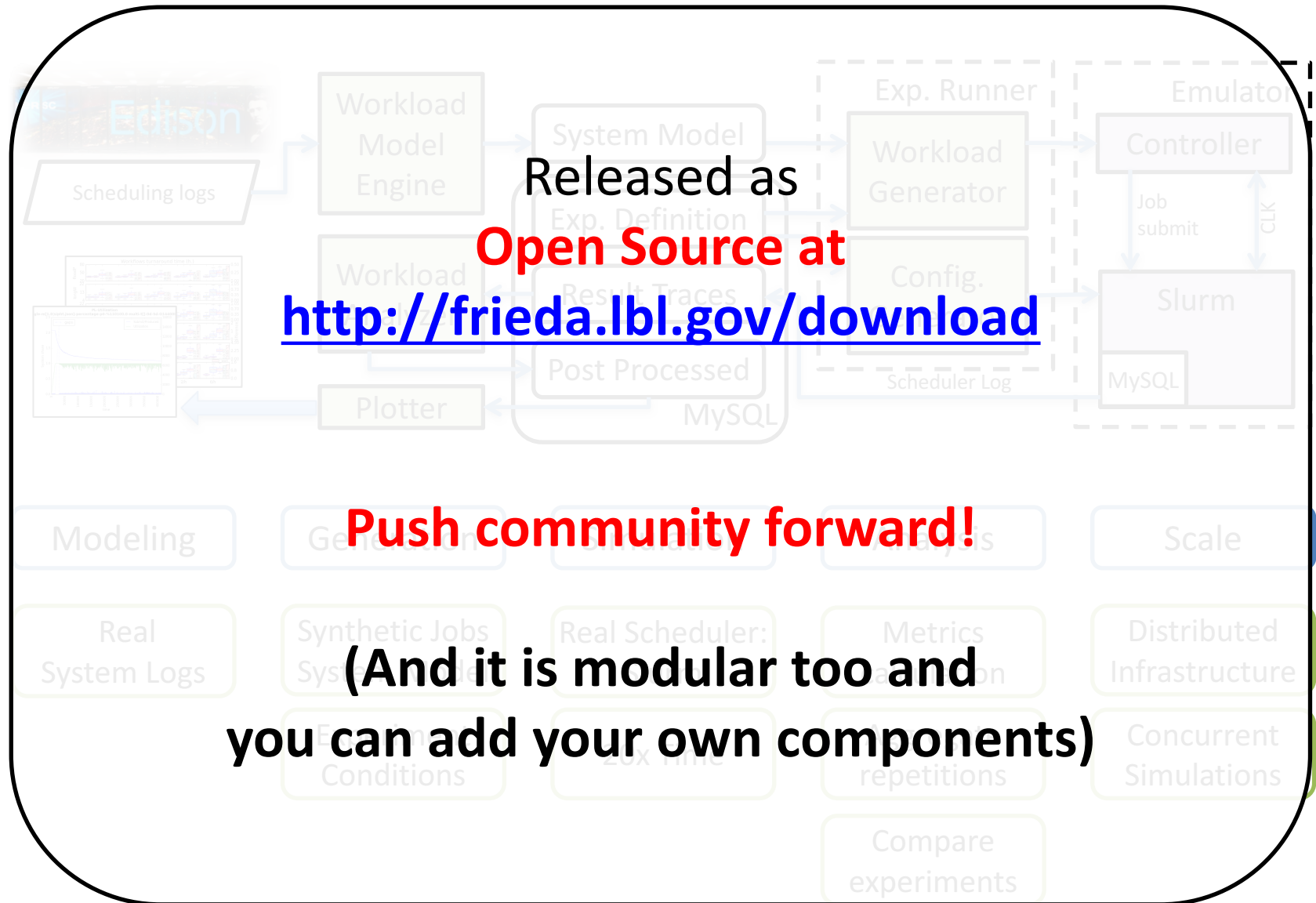
Compare  
experiments

Scale

Distributed  
Infrastructure

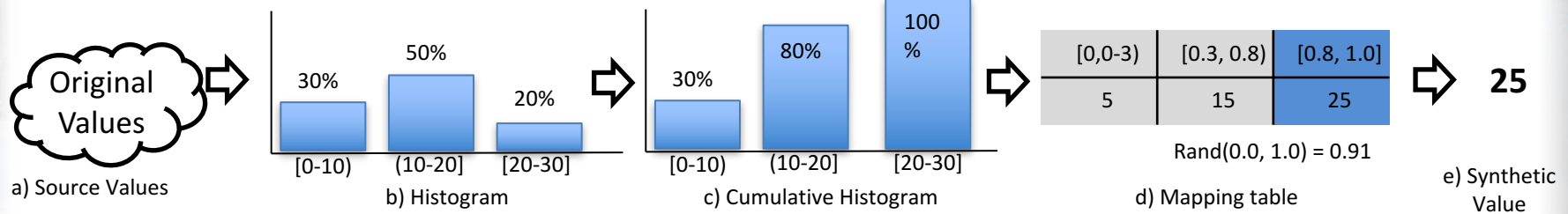
Concurrent  
Simulations

# ScSF: Scheduling Simulation Framework



# ScSF: Workload Modeling & Generation

Empirical Distribution





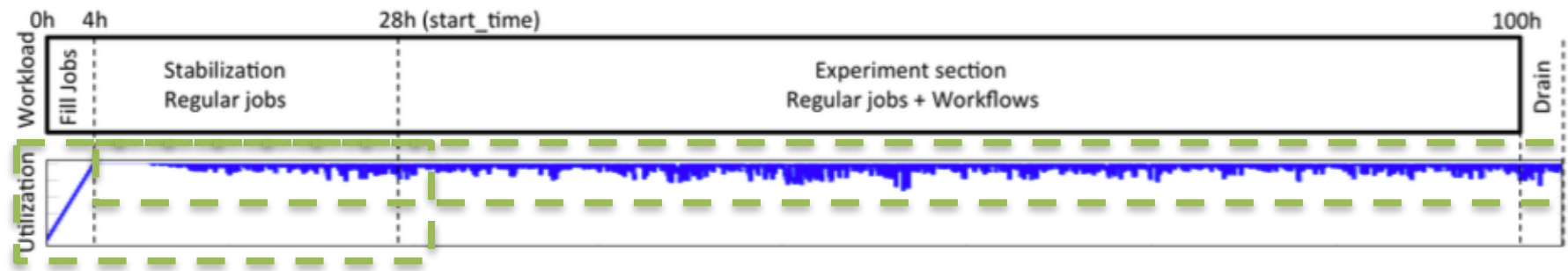
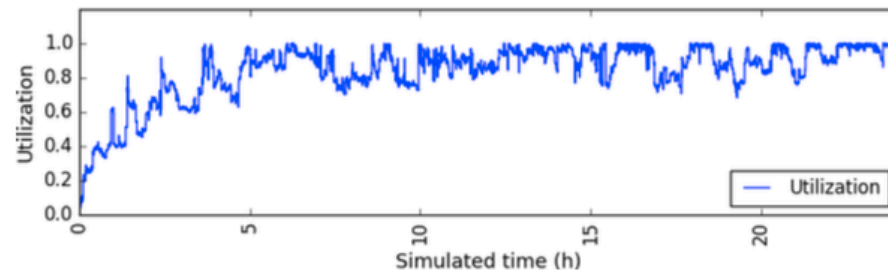
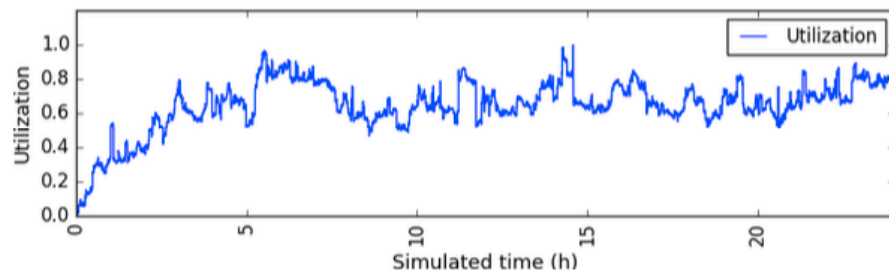
# ScSF: Workload Modeling & Generation

Empirical Distribution

Open Loop

Job Pressure Control

Cold start stabilization



# ScSF: Workload Modeling & Generation

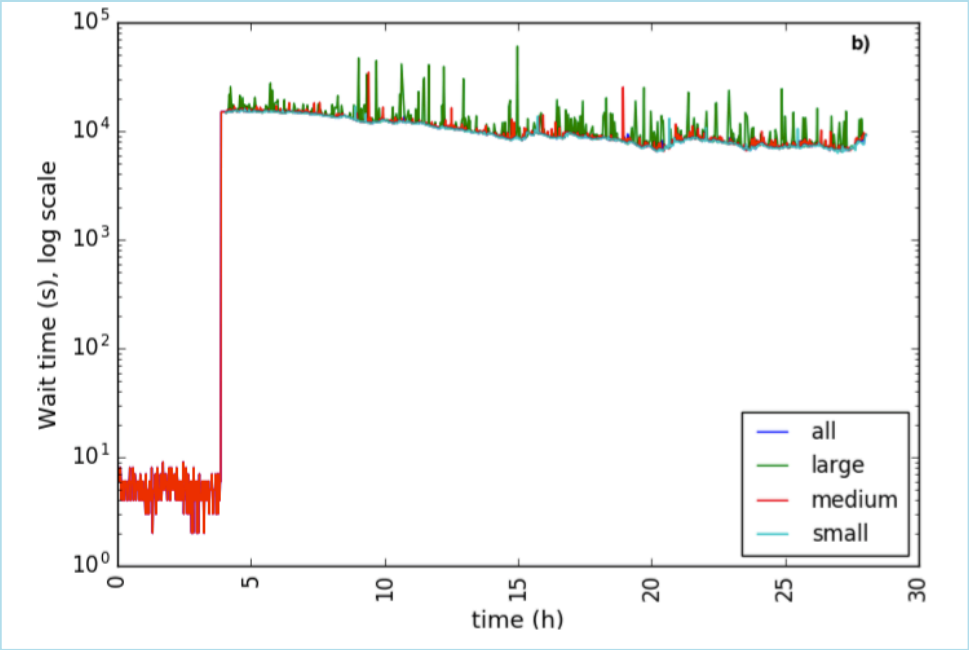
Empirical Distribution

Open Loop

Job Pressure Control

Cold start stabilization

Baseline wait time



# ScSF: Workload Modeling & Generation

Empirical Distribution

Open Loop

Job Pressure Control

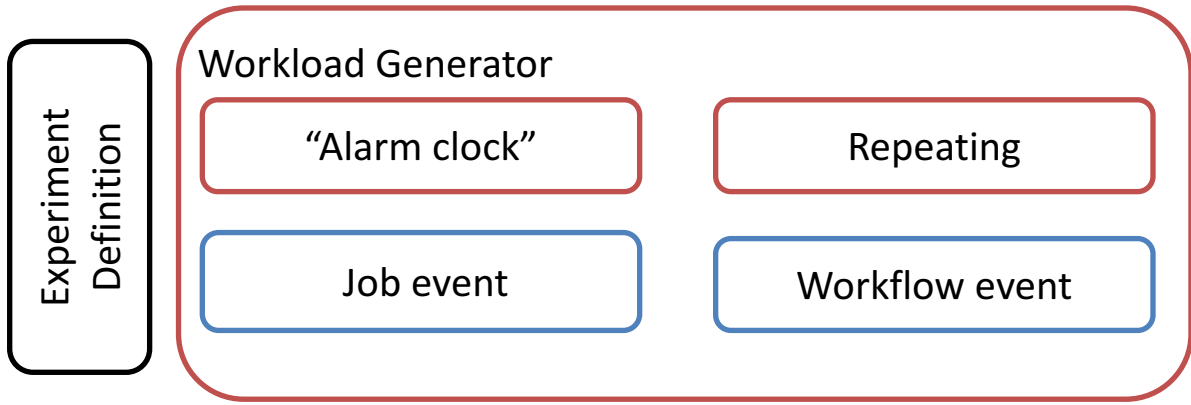
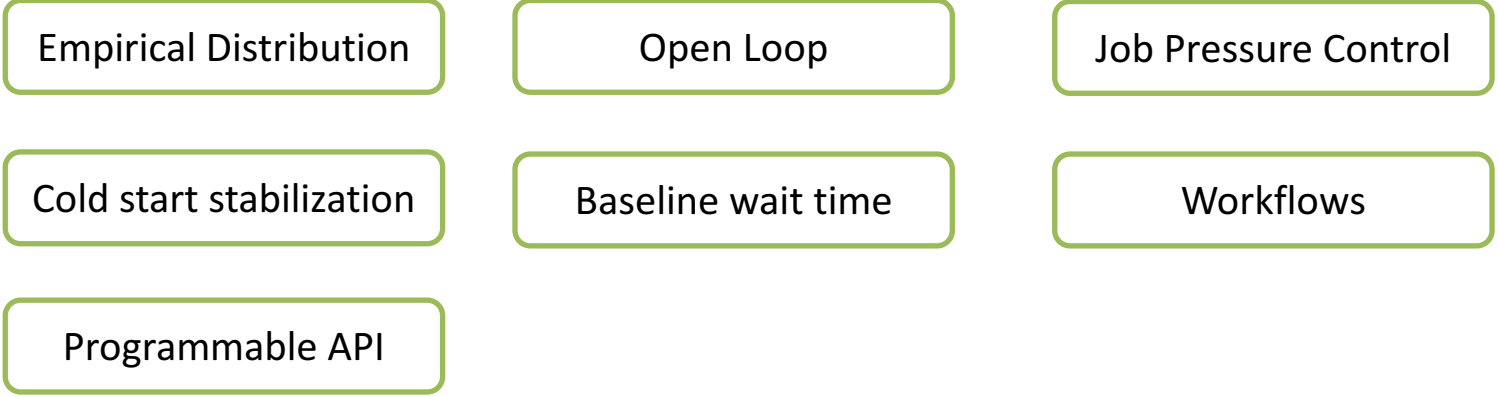
Cold start stabilization

Baseline wait time

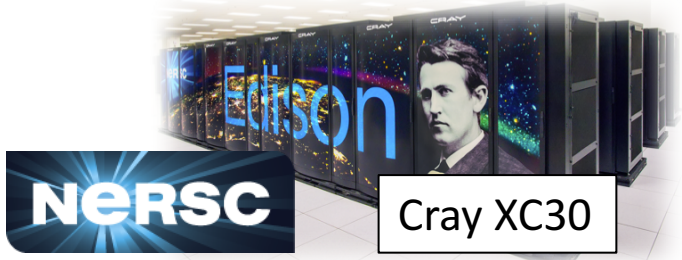
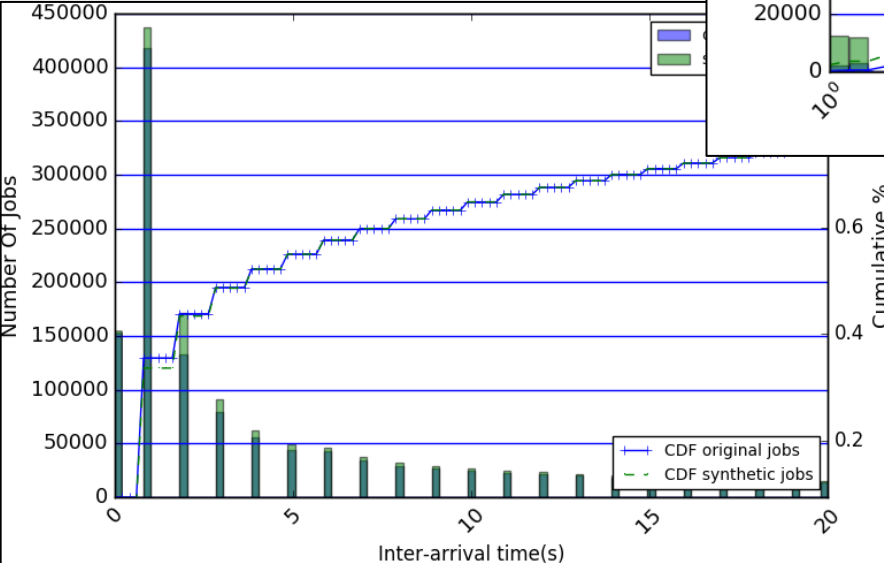
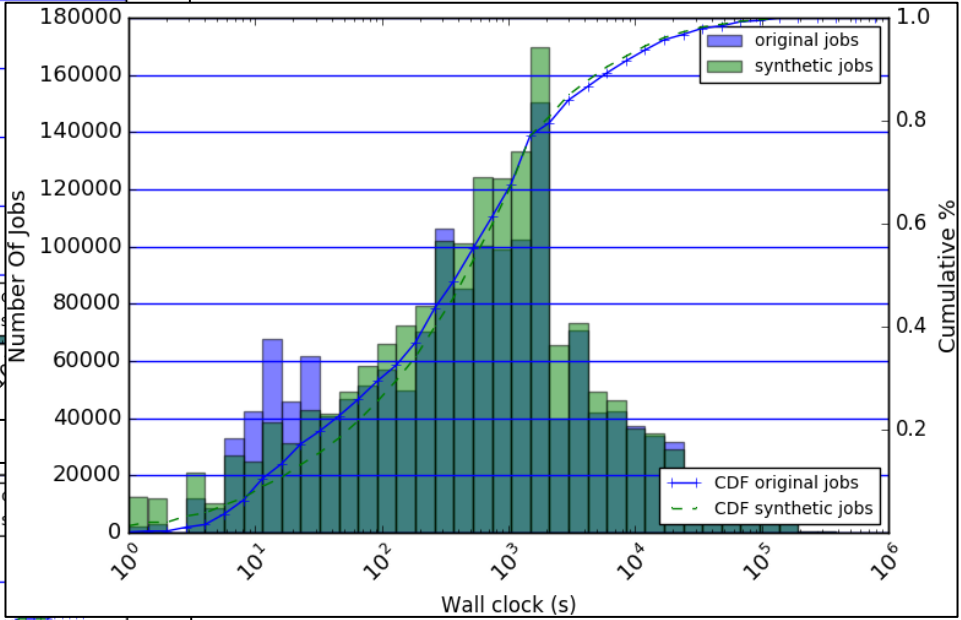
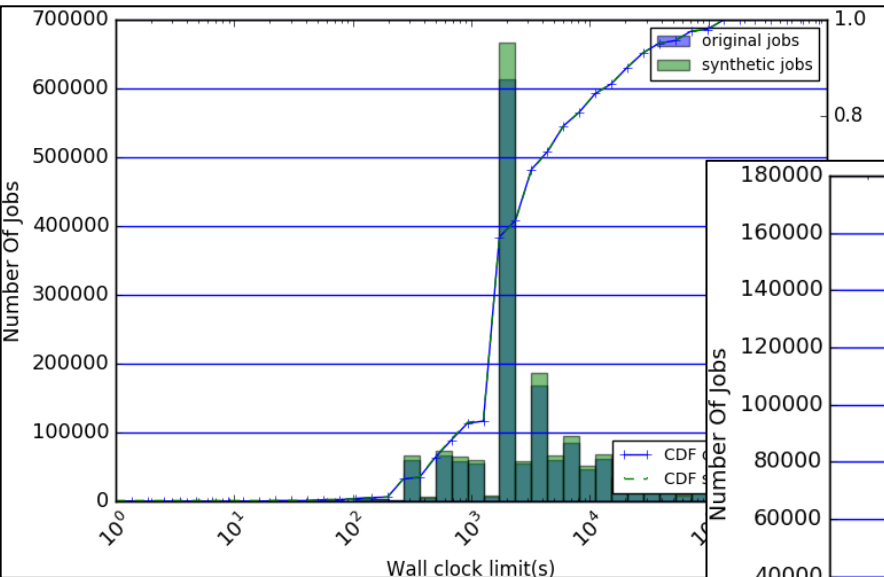
Workflows

```
1 {"tasks": [  
2   {"id": "SWide", "cmd": "./W.py", "cores": 480, "rtime": 360.0},  
3   {"id": "SLong", "cmd": "./L.py", "cores": 48, "rtime": 1440.0,  
4     "deps": ["SWide"]}]}]
```

# ScSF: Workload Modeling & Generation



# ScSF: Workload Modeling & Generation

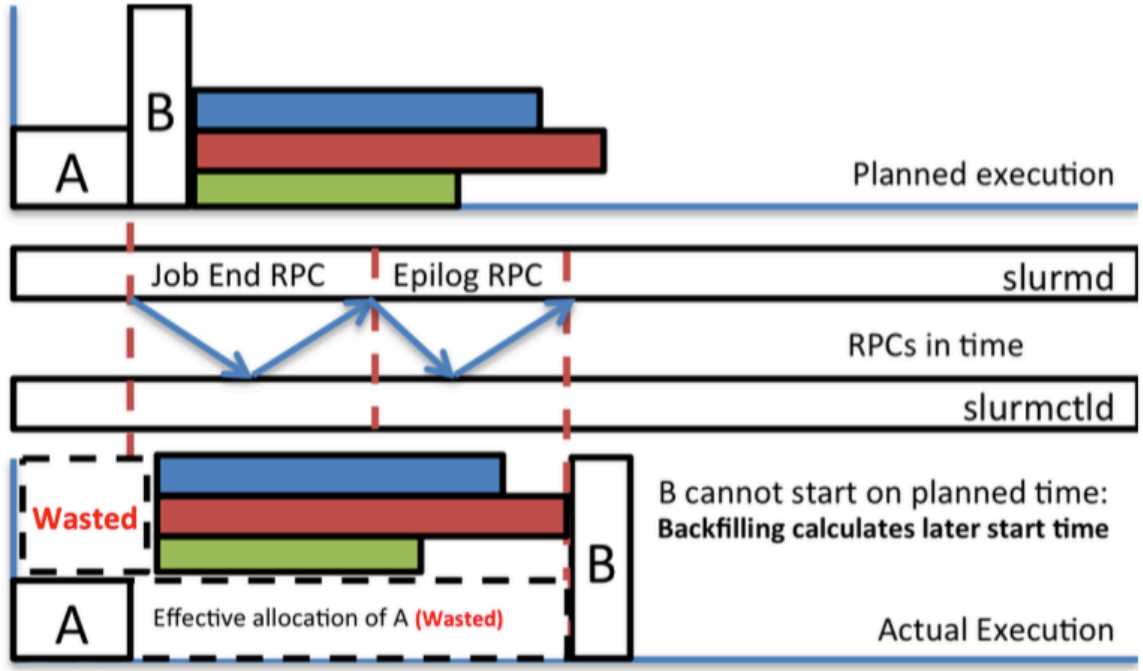


Peak: 2.57 Petaflops/s.

# ScSF: Slurm Simulator

- Wraps real Slurm Scheduler
- Emulates system and job execution
- Emulates job submission (replay)

Original Implementation: Slow (1 to 1), no determinism



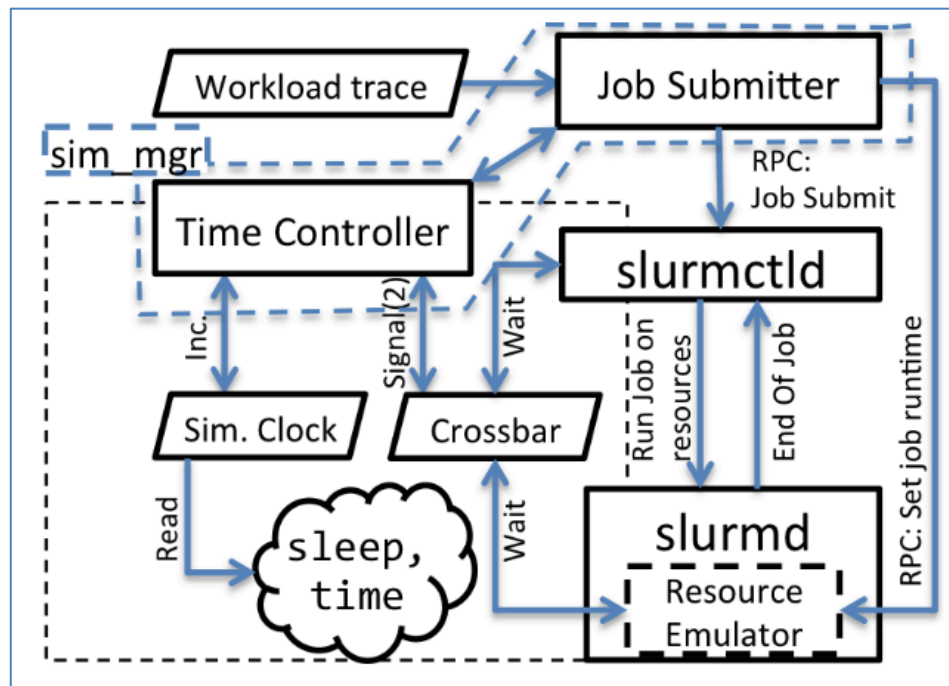
# ScSF: Slurm Simulator

Wraps real Slurm Scheduler

Emulates system and job execution

Emulates job submission (replay)

Original Implementation: Slow (1 to 1), no determinism



**Slurm simulator improved by synchronizing scheduling threads**

Faster (20x speed-up)

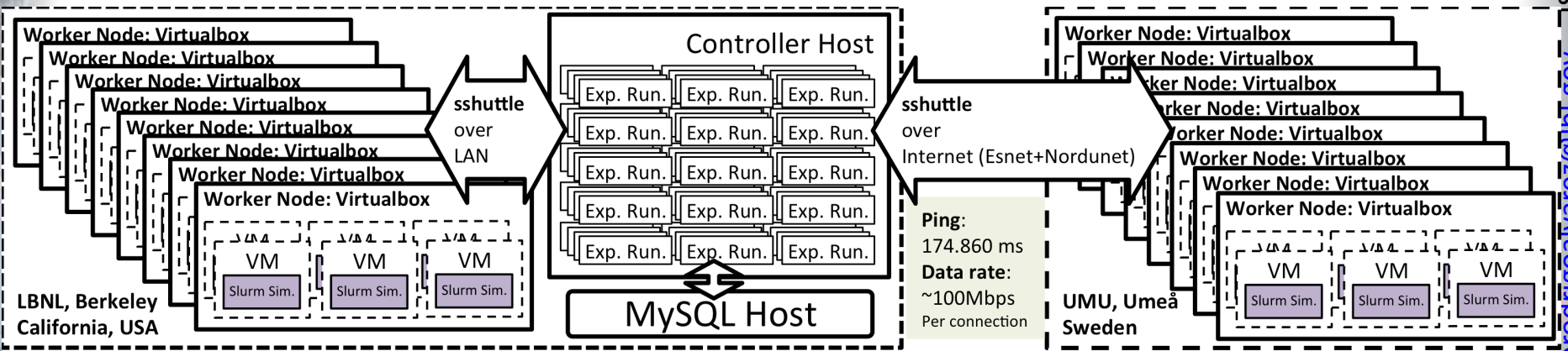
Time consistent

Achieve good utilization with out-of-the-box scheduler

# ScSF: Running experiments in scale

Many degrees of freedom, many experiments:  
ScSF, years to complete

**Orchestrator**      Concurrent experiments      Distributed resources



Simulation worker VM

Controller

MySQL Database



# ScSF: Running experiments in scale

Many degrees of freedom, many experiments:  
1000s of experiments, years to complete

170 Worker VMs  
17 Hosts  
Two continents

Simulate  
30 years  
of Edison's Life

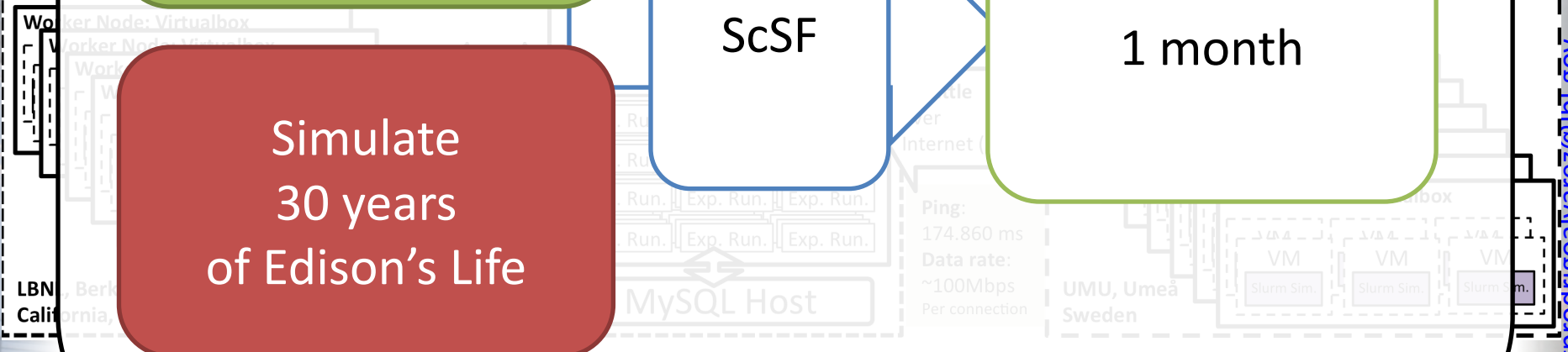
ScSF

1 month

Simulation worker VM

Controller

MySQL Database



# ScSF: Analysis Capabilities

Job Variables

Histogram+ CDF

Workflow variables

Average, Median, Percentiles

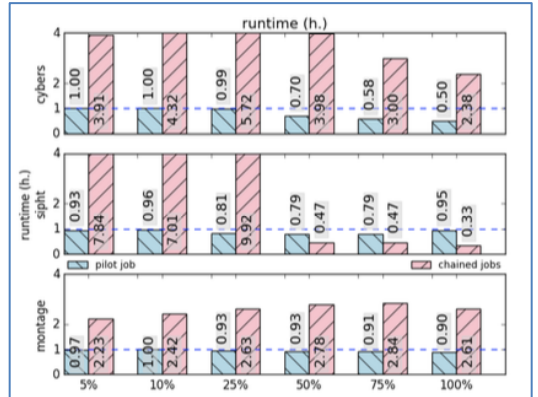
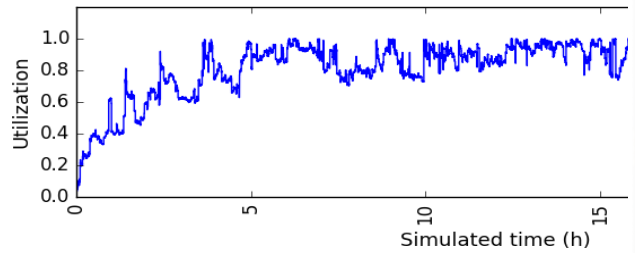
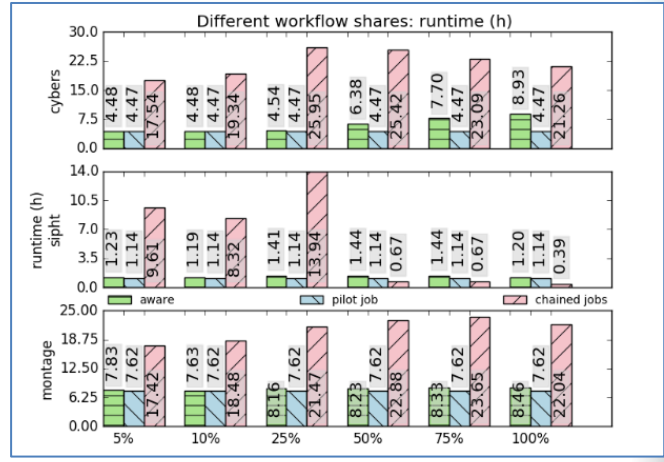
Utilization

Aggregated

Per Minute

Aggregating results from repetitions

Compare experiments with similar configuration



# ScSF: Research Use case, WoAS\*

**WoAS:** An scheduling technique to minimize workflows turnaround time without over allocating resources

## Questions to answer with ScSF

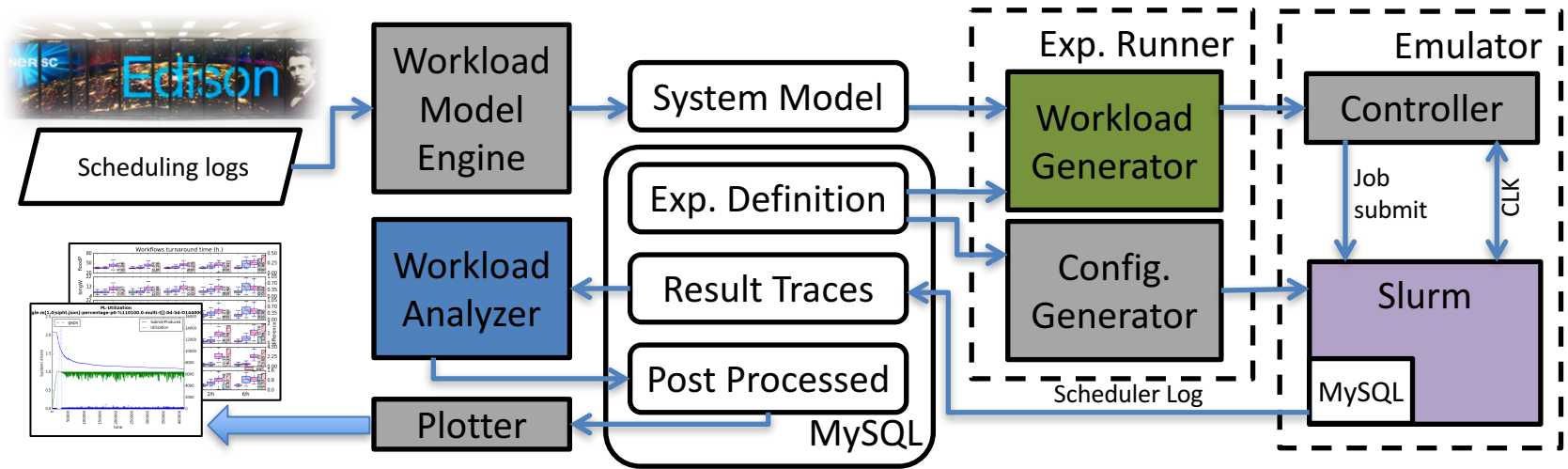
Can the WoAS scheduling technique minimize workflows turnaround time without over allocating resources?

How much shorter is turnaround time with WoAS? (Compared to existing techniques)

Does WoAS impact negatively on non-workflow jobs? (Slowdown)

\* Results to be published in **HPDC'17**: "Enabling Workflow-Aware Scheduling on HPC systems"

# ScSF: Getting ready for WoAS \*



1: Generate workload with WoAS workflows

2: Implement WoAS in Slurm

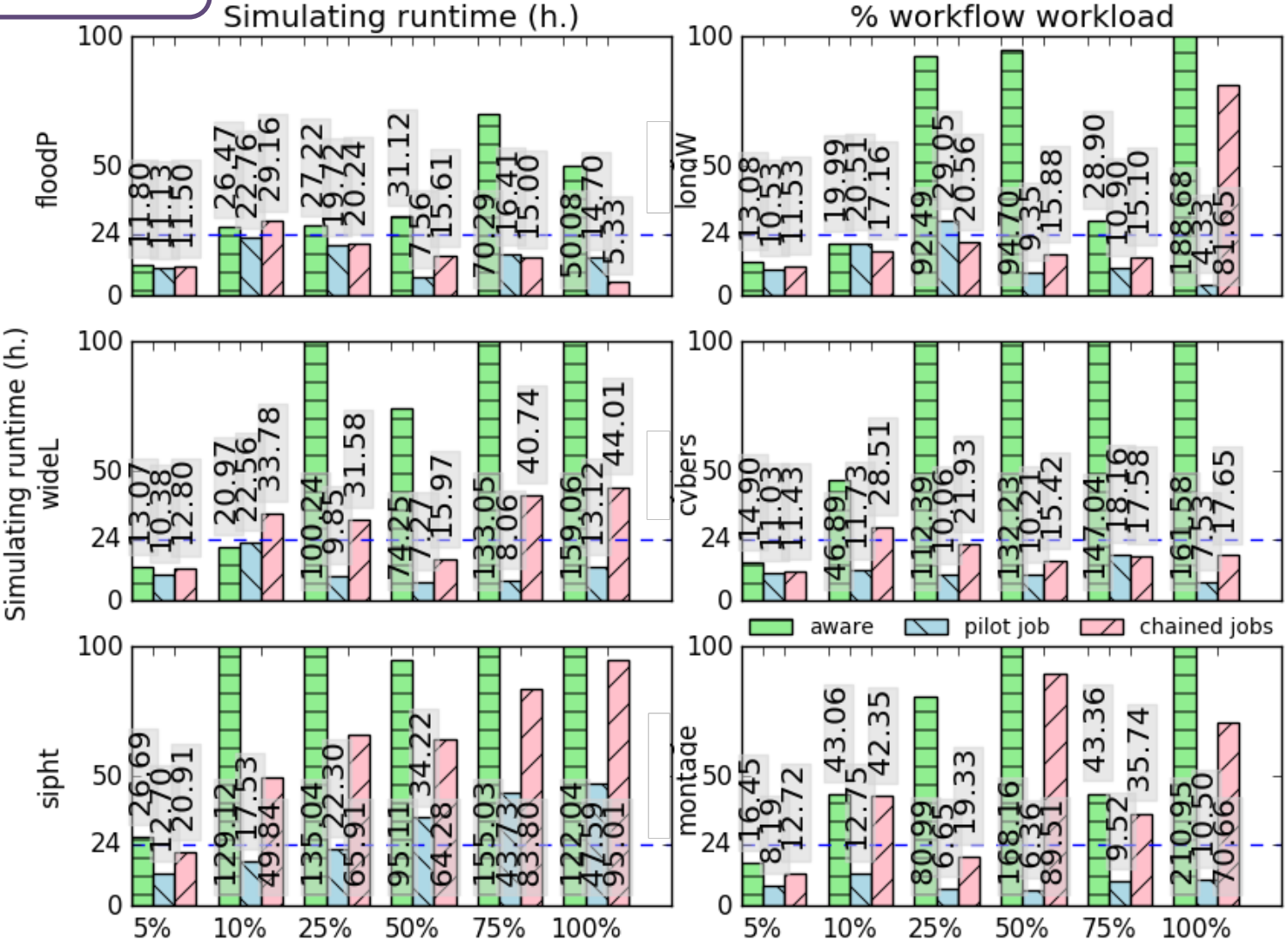
3: Identify WoAS workflows

271 scenarios,  
1728 experiments  
33 years of supercomputer time!  
38 Billion Core-hours (simulated)

\* Results to be published in **HPDC'17**: "Enabling Workflow-Aware Scheduling on HPC systems"

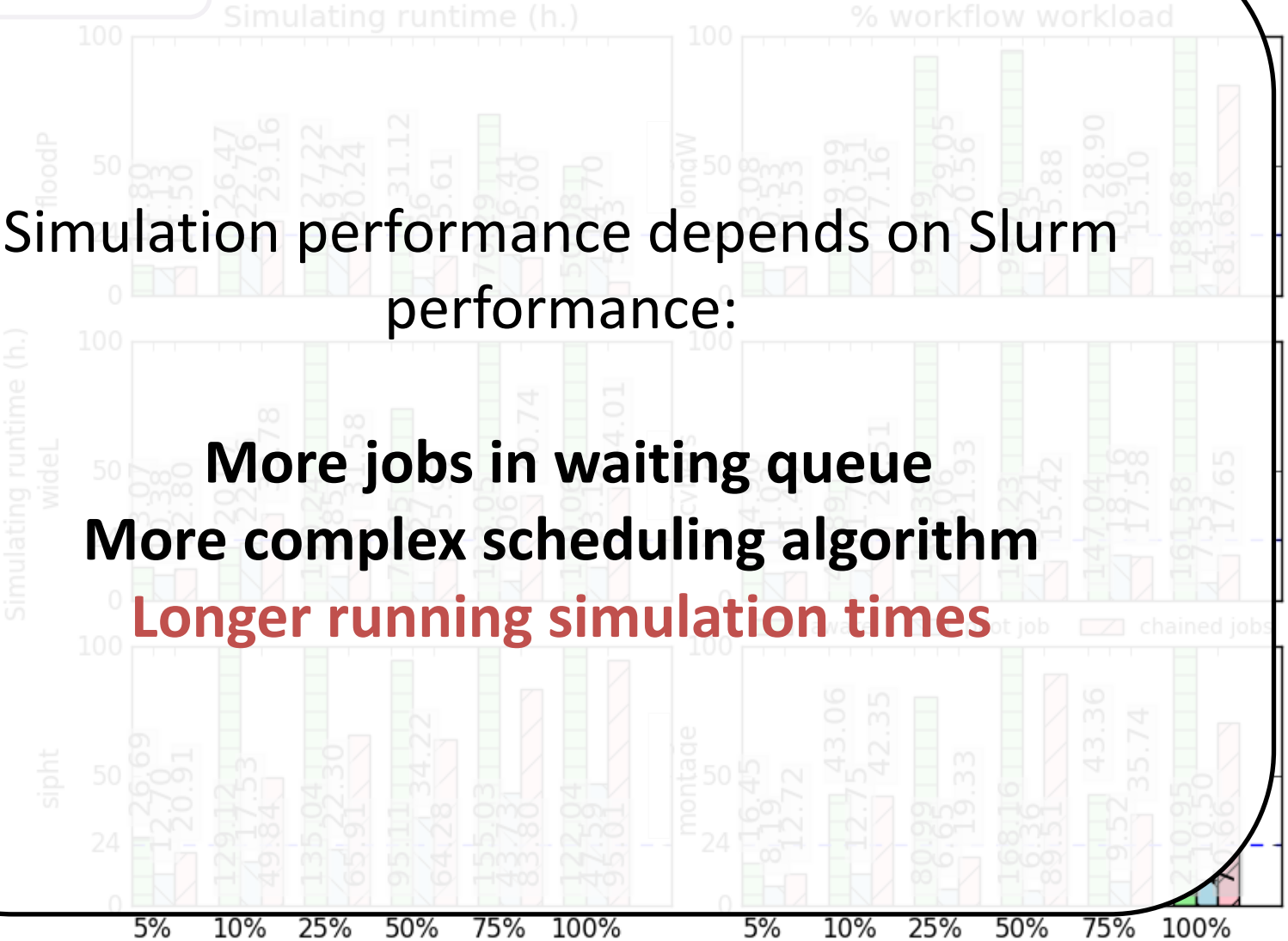
# ScSF: Simulation performance on WoAS

6 Days of simulated time.  
Different scenarios



# ScSF: Simulation performance on WoAS

6 Days of simulated time.  
Different scenarios



Simulation performance depends on Slurm performance:

**More jobs in waiting queue**  
**More complex scheduling algorithm**  
**Longer running simulation times**

# ScSF: Lessons learned

Slurm is a complex old-fashion-SWE package:  
expensive to modify

**Loss-less experiment restart is needed**  
Specially if experiment runtime are long (e.g. 5 days)

**Monitoring is important**  
To debug why something fails (and things will fail)

**Loaded systems network fail**  
So harden your comms

**The system is as weak as its weakest link**  
Single point of failure

# ScSF: Lessons learned

Slurm is a complex old-fashion-SWE package:  
expensive to modify

Loss-less experiment restart is needed  
Specially if experiment runtime are long (e.g. 5 days)

HPC scheduling requires a lot of simulation  
To debug **Think big from the beginning!** (anything will fail)

Loaded systems network fail  
So harden your comms

The system is as weak as its weakest link  
Single point of failure



# ScSF: Summary

HPC Scheduling research cycle:  
Model/generate workloads -> scheduling emulation -> analysis

Tools to run experiments in scale

Slurm simulator in its core: A production HPC simulator

Open Source! **Use it!** (<http://frieda.lbl.gov/download>)  
And **modular**: free to replace any of the parts

(Keep scale in mind!)

# THANKS

For any questions, please contact:  
[gprodrigoalvarez@lbl.gov](mailto:gprodrigoalvarez@lbl.gov)

To know more, read: Rodrigo Álvarez, G.P, Elmroth, E., Östberg, P.O., Ramakrishnan, L. ScSF: A Scheduling Simulation Framework. 21th Workshops on Job Scheduling Strategies for Parallel Processing (JSSPP 2017)

ScSF is open source and **available at:**  
<http://frieda.lbl.gov/download>

Supported by U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research (ASCR), the National Energy Research Scientific Computing Center, a DOE Office of Science User Facility, is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231. Financial support has been provided in part by the Swedish Government's strategic effort eSSSENCE and the Swedish Research Council (VR) under contract number C0590801 (Cloud Control).

