



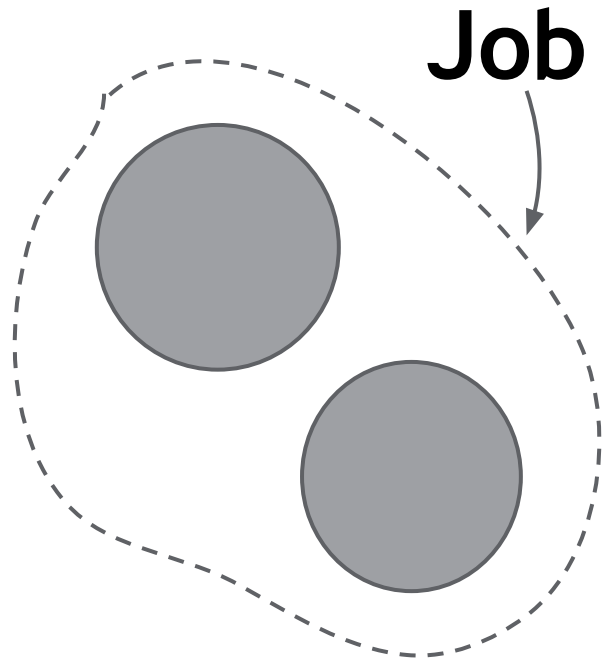
Omega: flexible, scalable schedulers for large compute clusters

Malte Schwarzkopf (University of Cambridge Computer Lab)

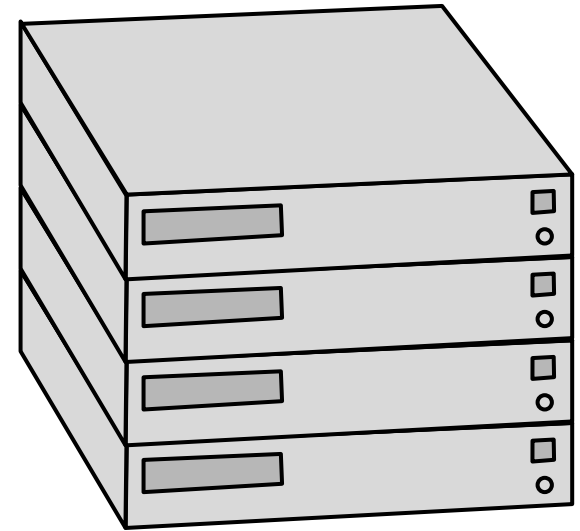
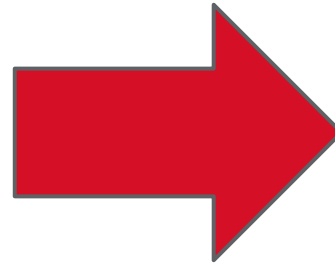
Andy Konwinski (UC Berkeley)

Michael Abd-El-Malek (Google)

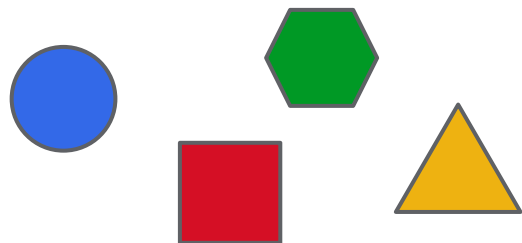
John Wilkes (Google)



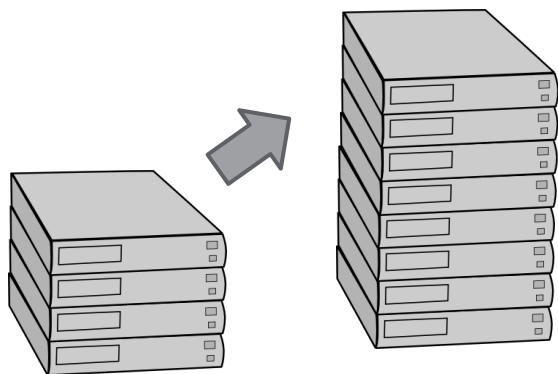
Tasks



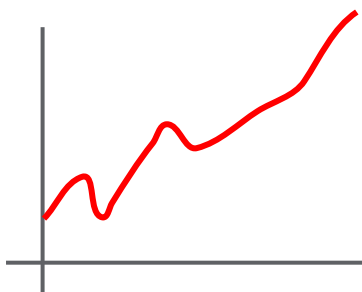
Machines



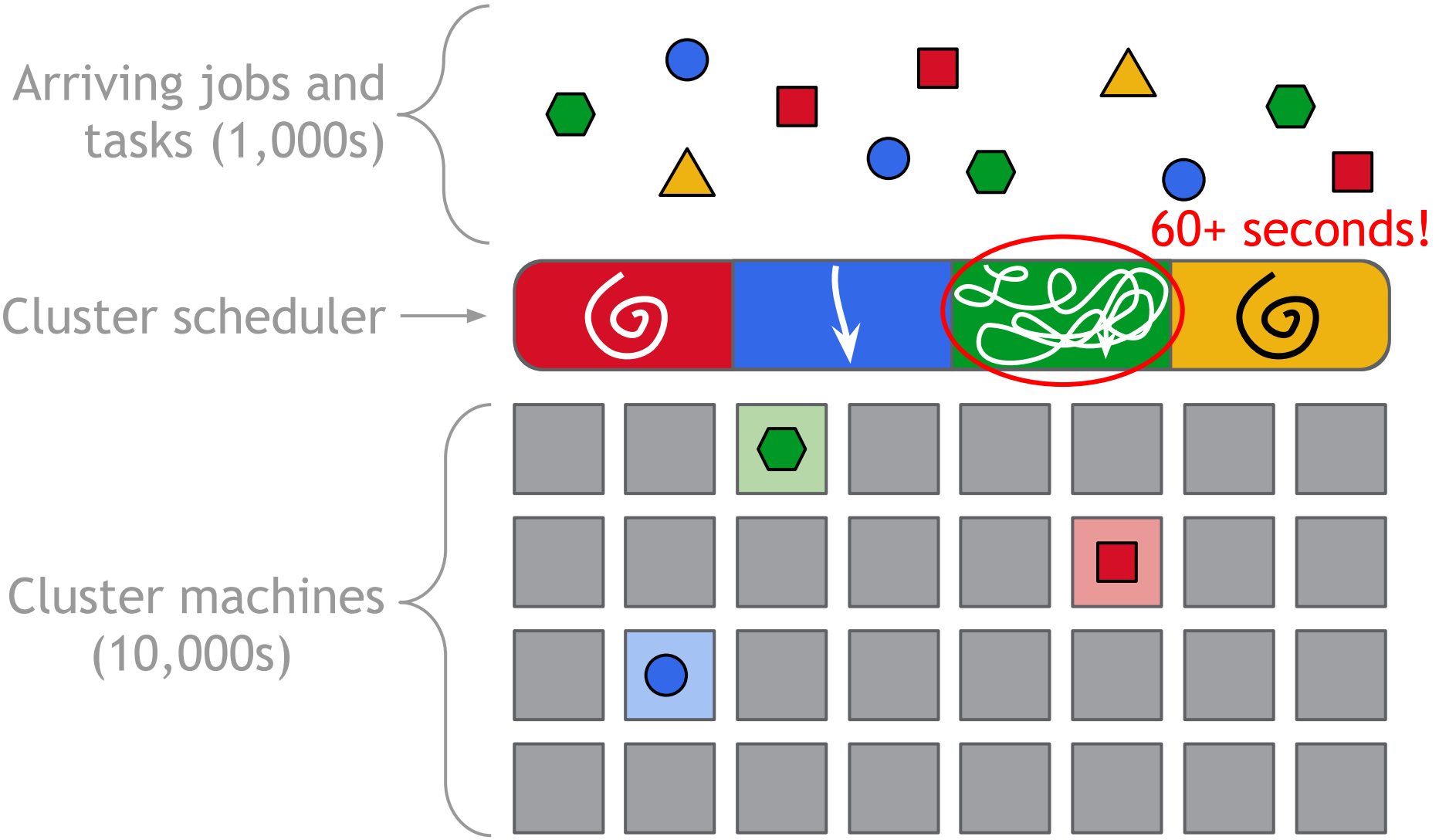
Diverse workloads



Increasing cluster sizes



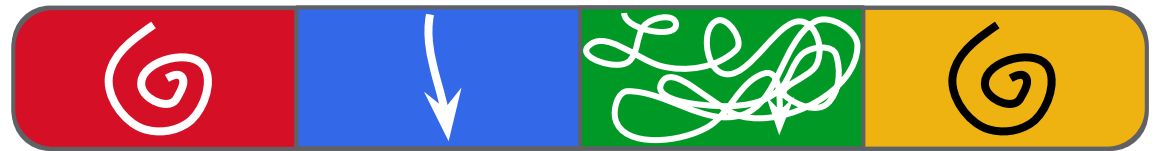
Growing job arrival rates



Arriving jobs and tasks (1,000s)

Increasing complexity!

Cluster scheduler

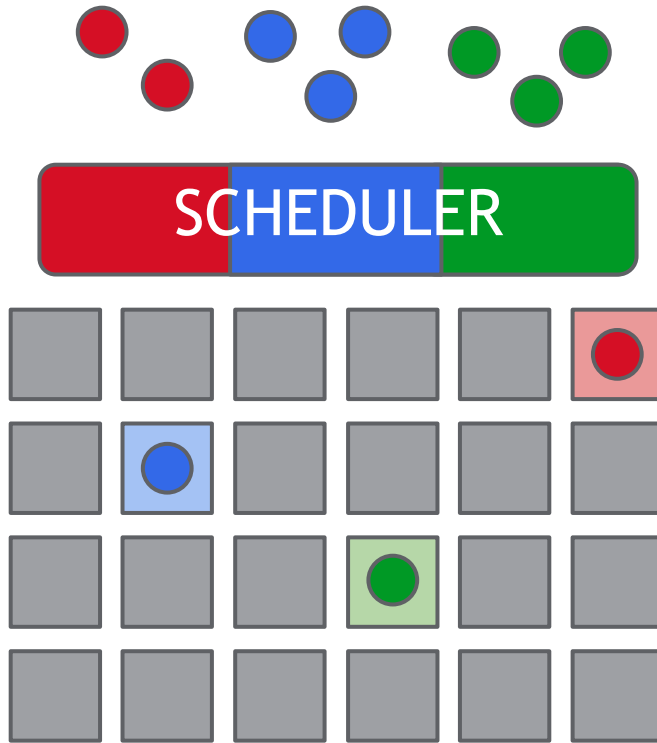


Hence:

Break up into independent schedulers.

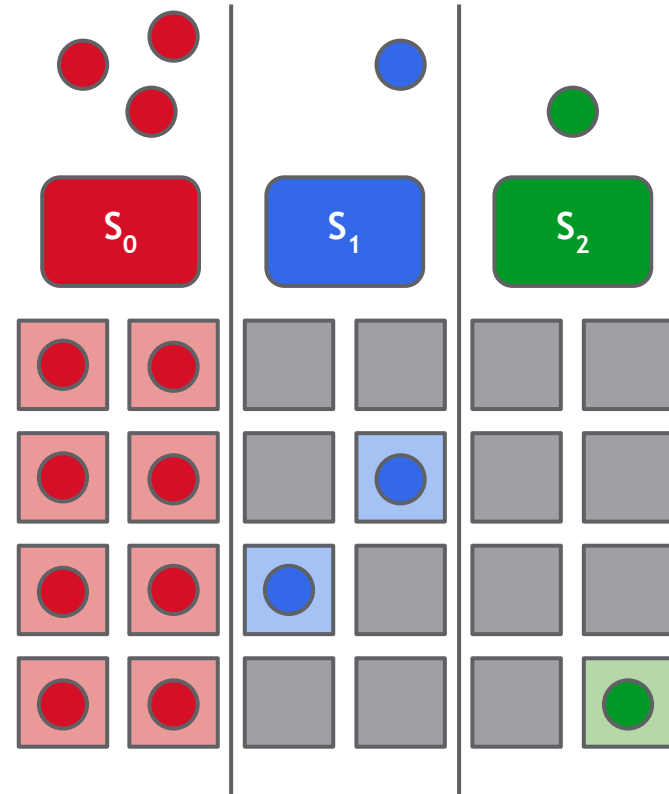
But: How do we arbitrate resources between schedulers?

monolithic scheduler



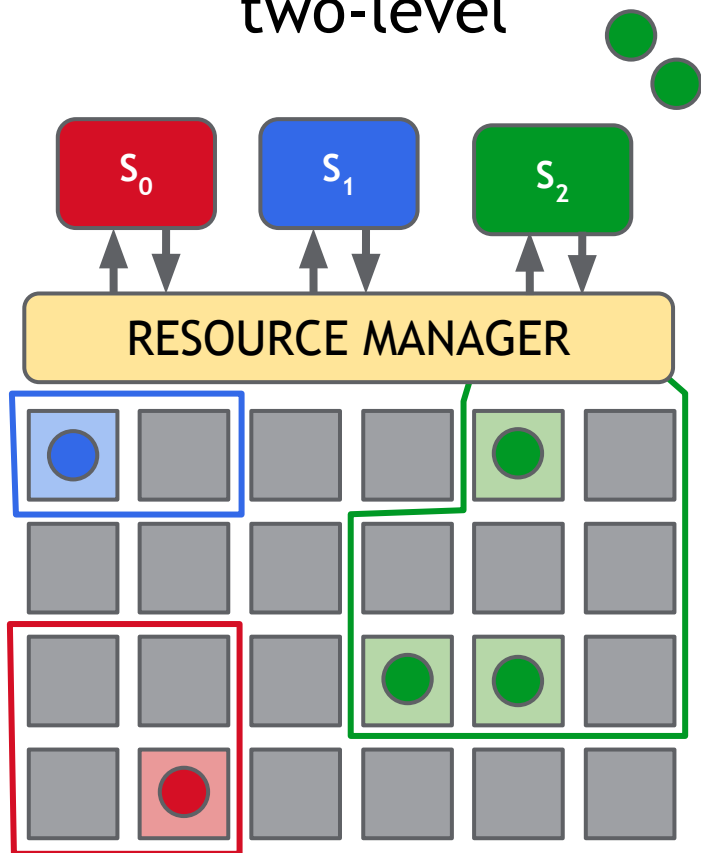
- hard to diversify
- code growth
- scalability bottleneck

static partitioning



- poor utilization
- inflexible

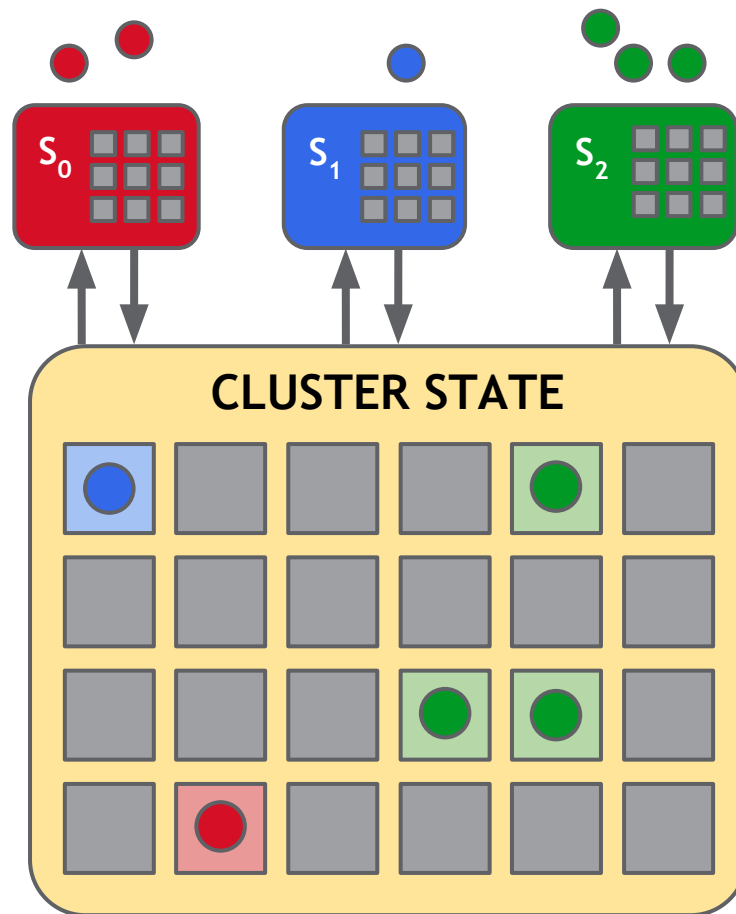
two-level

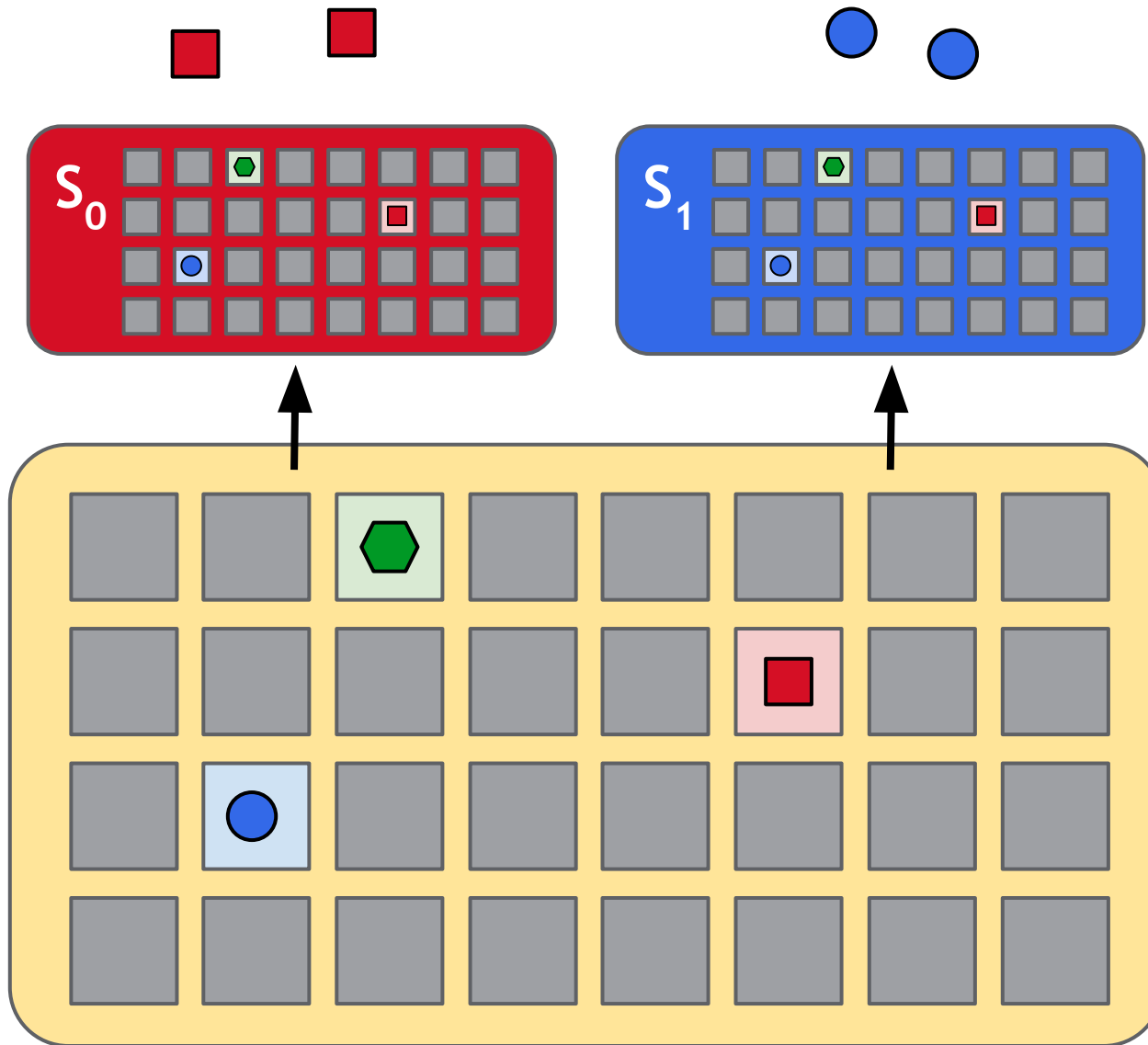


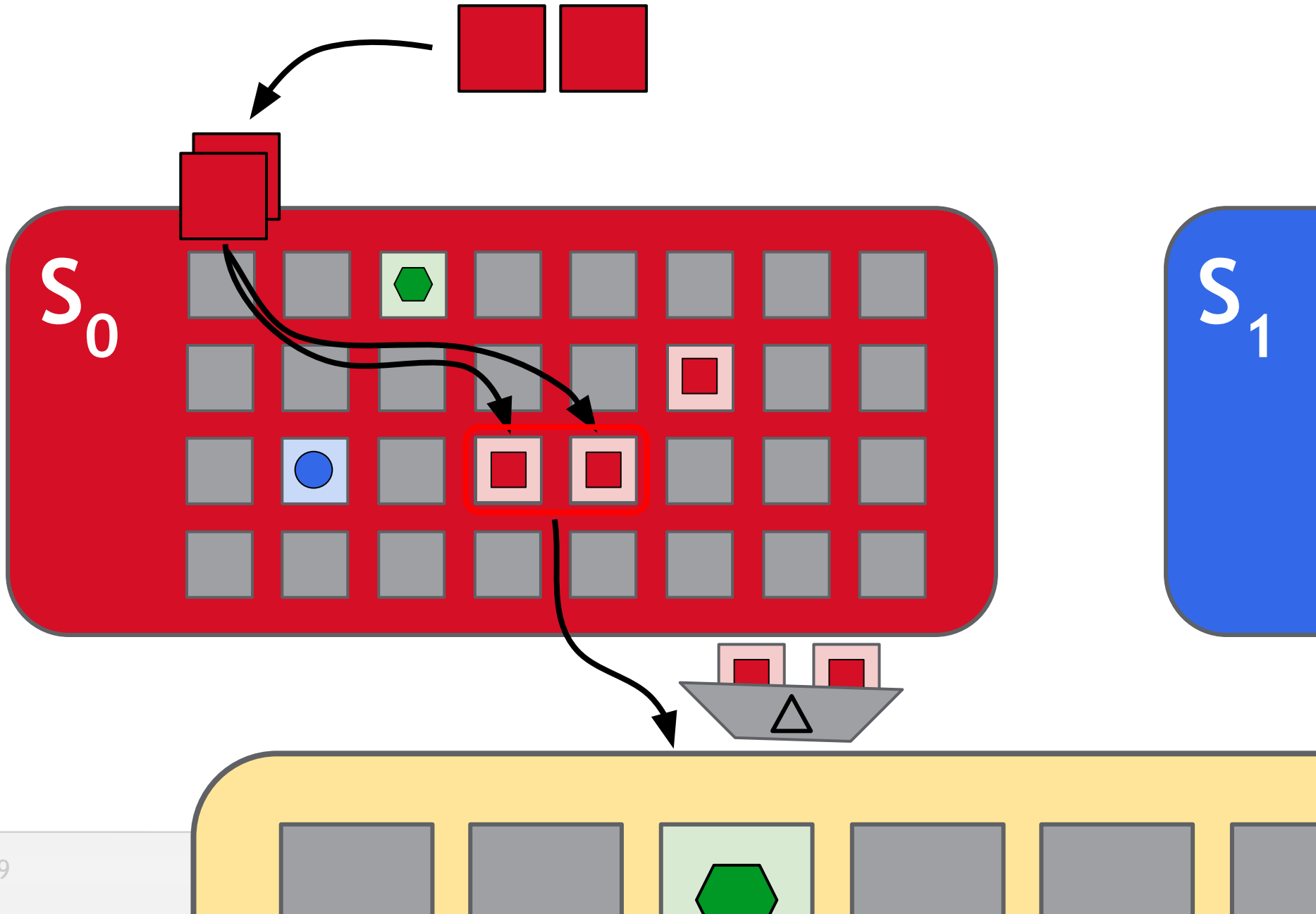
- hoarding
- information hiding

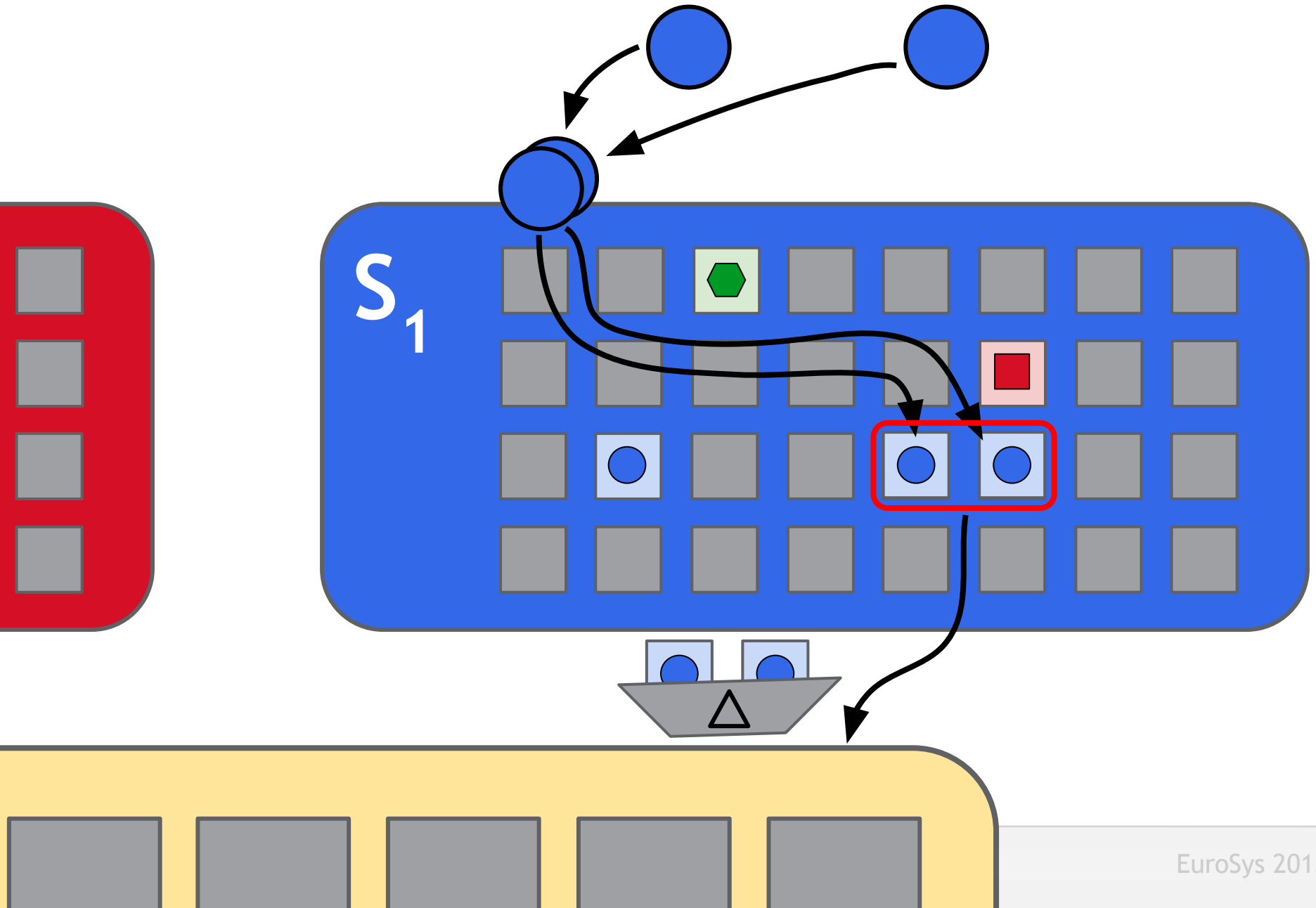
e.g. UCB Mesos [NSDI 2011]

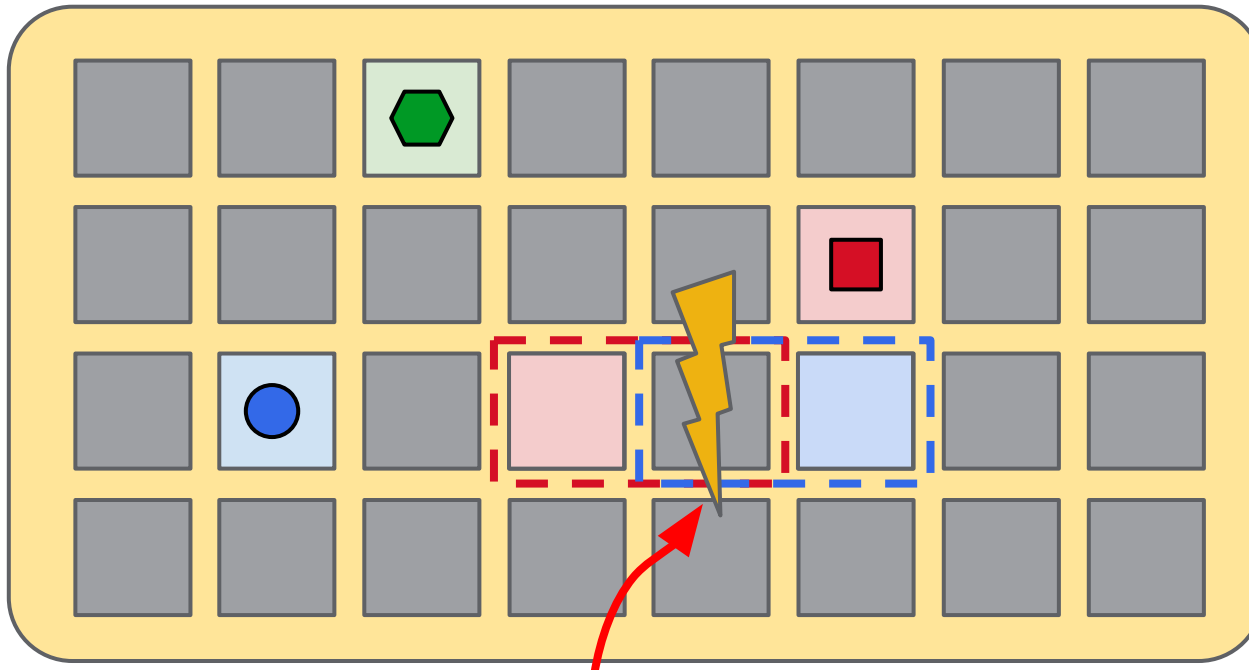
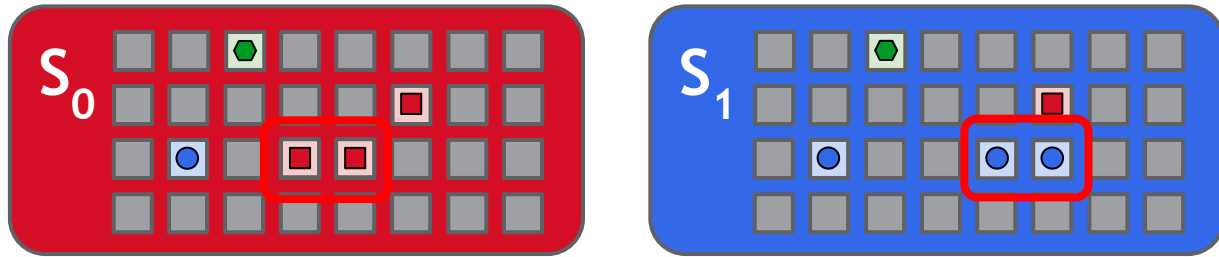
shared-state



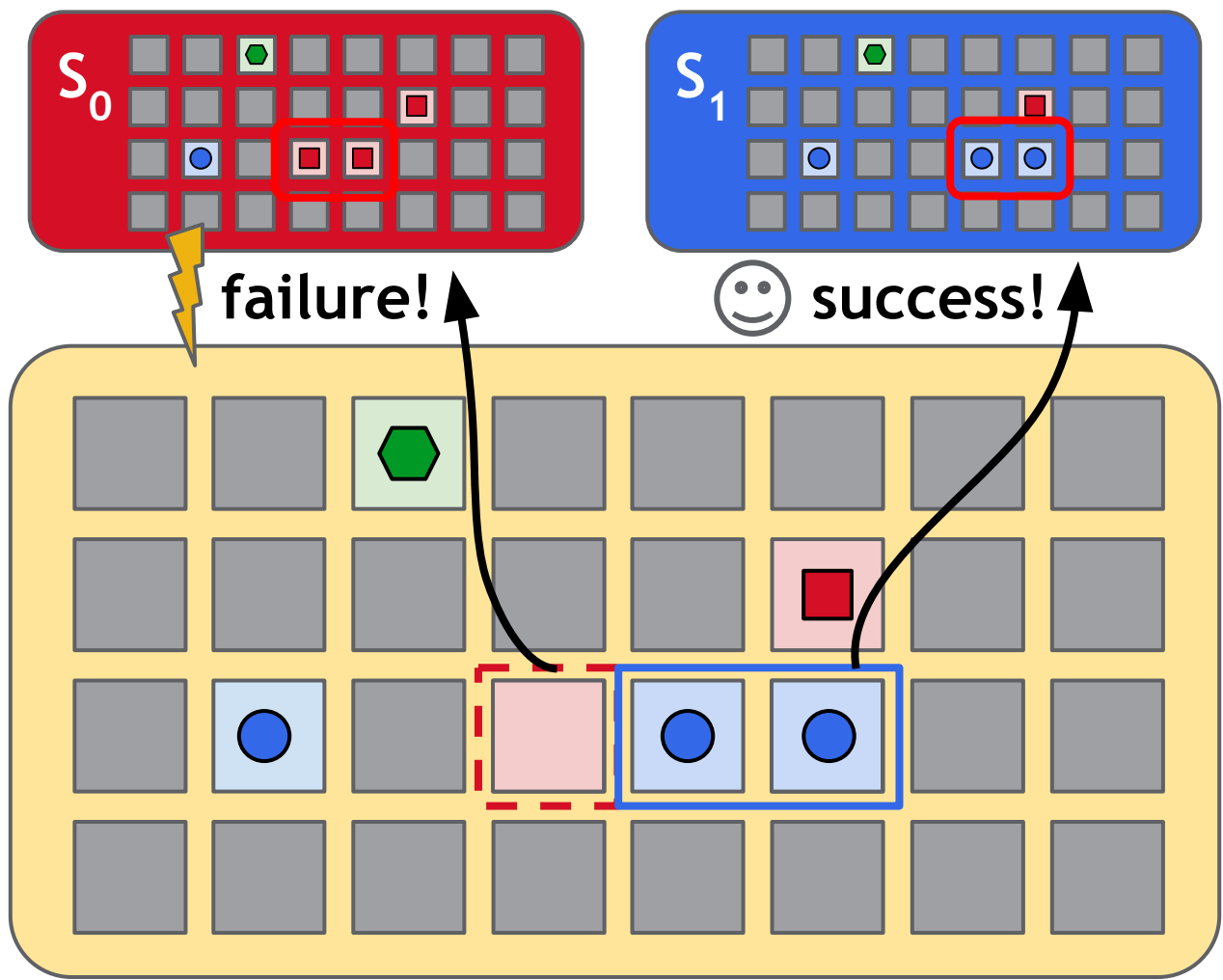









Conflict!



- 1) intro & motivation
- 2) workload characterization** 
- 3) comparison of approaches
- 4) trace-based simulation
- 5) flexibility case study

Batch

Service

Cluster A

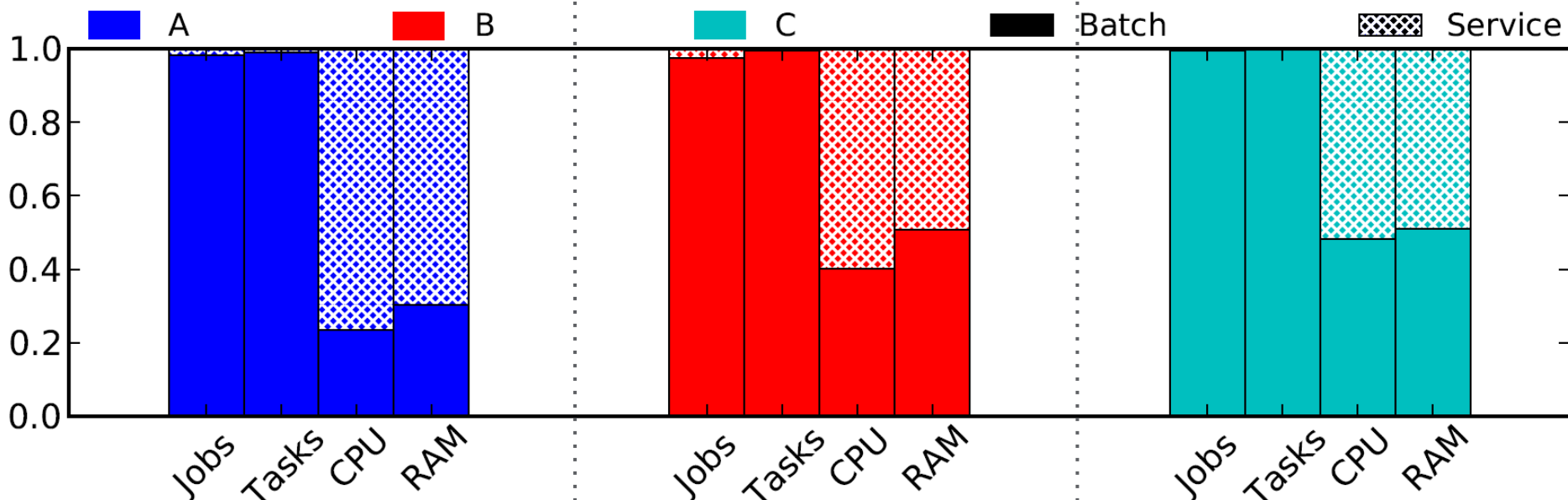
Medium size
Medium utilization

Cluster B

Large size
Medium utilization

Cluster C

Medium (12k mach.)
High utilization
Public trace

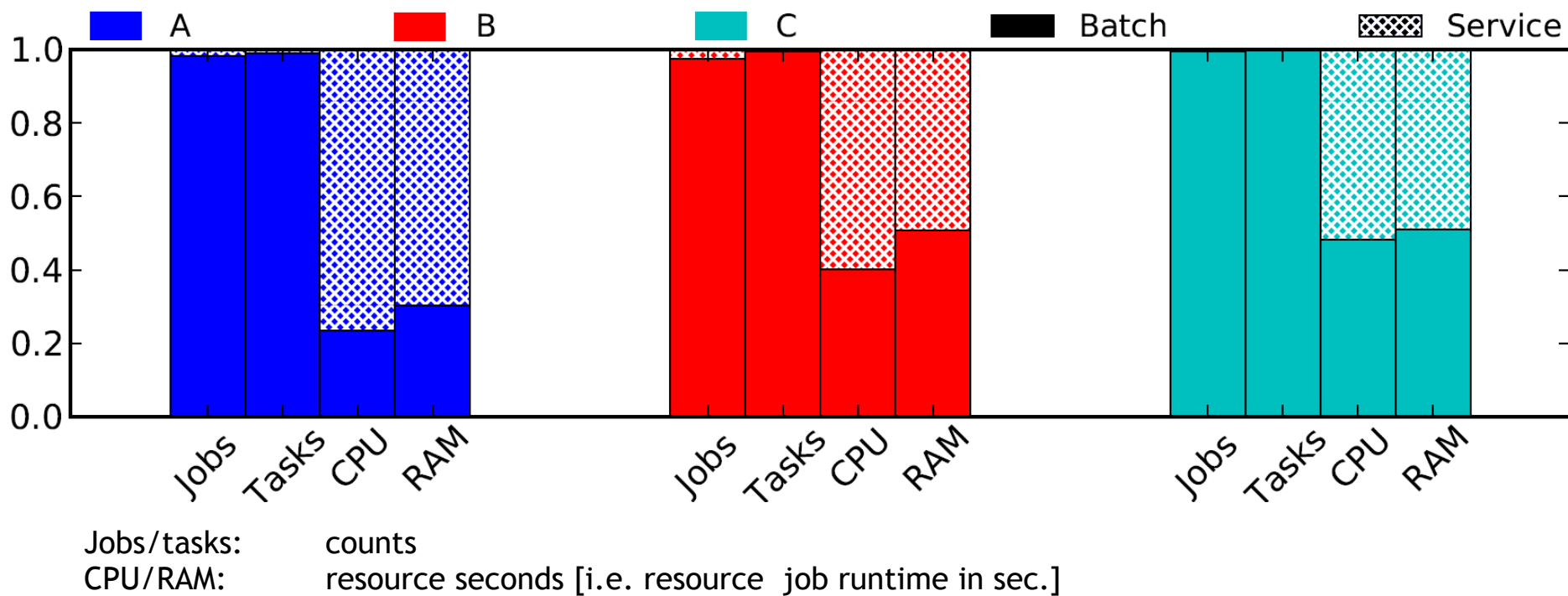


Jobs/tasks:
CPU/RAM:

counts
resource seconds [i.e. resource job runtime in sec.]

TAKEAWAY

Most jobs are batch, but most resources are consumed by service jobs.



Batch jobs

Service jobs

80th %ile runtime


12-20 min.

29 days

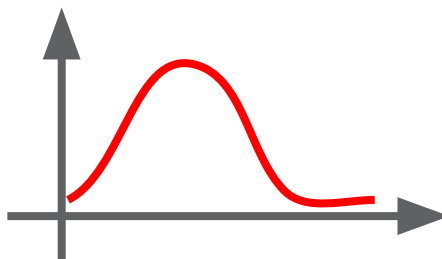
80th %ile inter-arrival time

4-7 sec.

2-15 min.

- 1) intro & motivation
- 2) workload characterization
- 3) comparison of approaches** 
- 4) trace-based simulation
- 5) flexibility case study

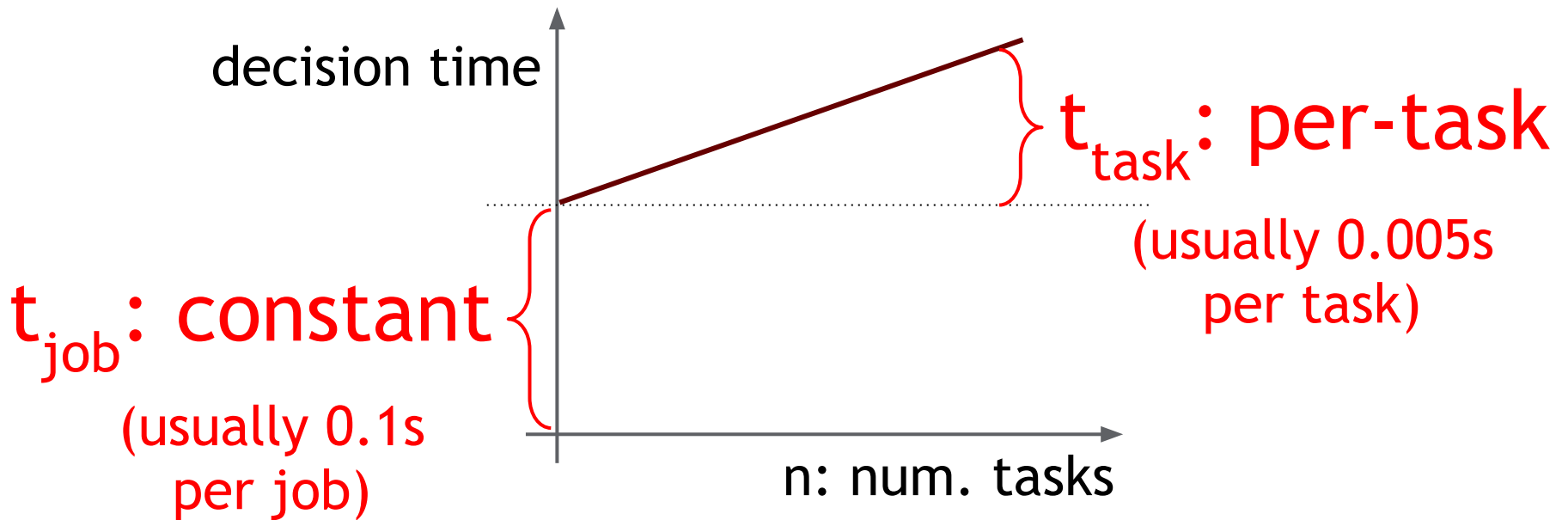
simulation using empirical workload parameters distributions



Code [soon to be] available:

<http://code.google.com/p/cluster-scheduler-simulator>

Scheduler decision time

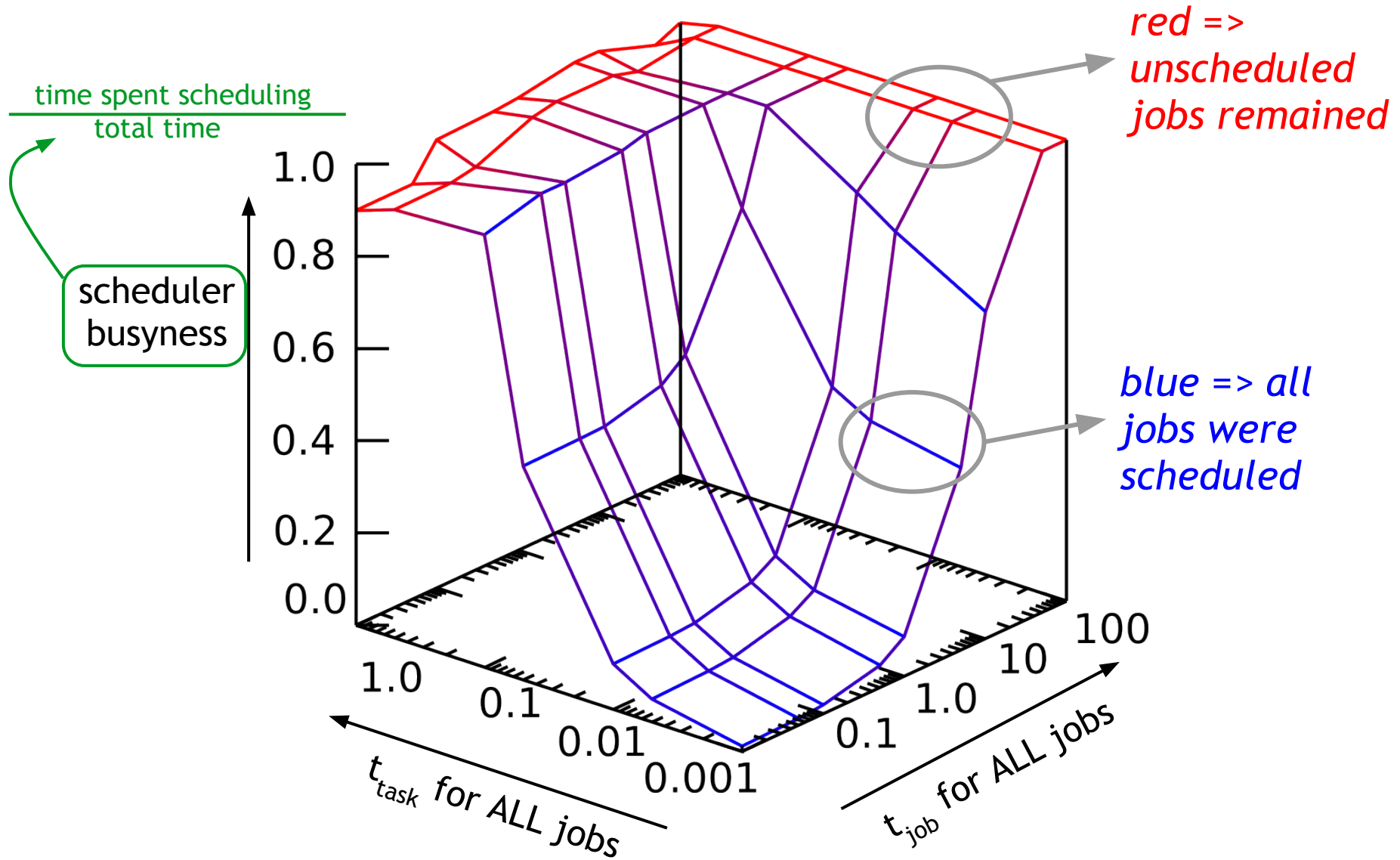


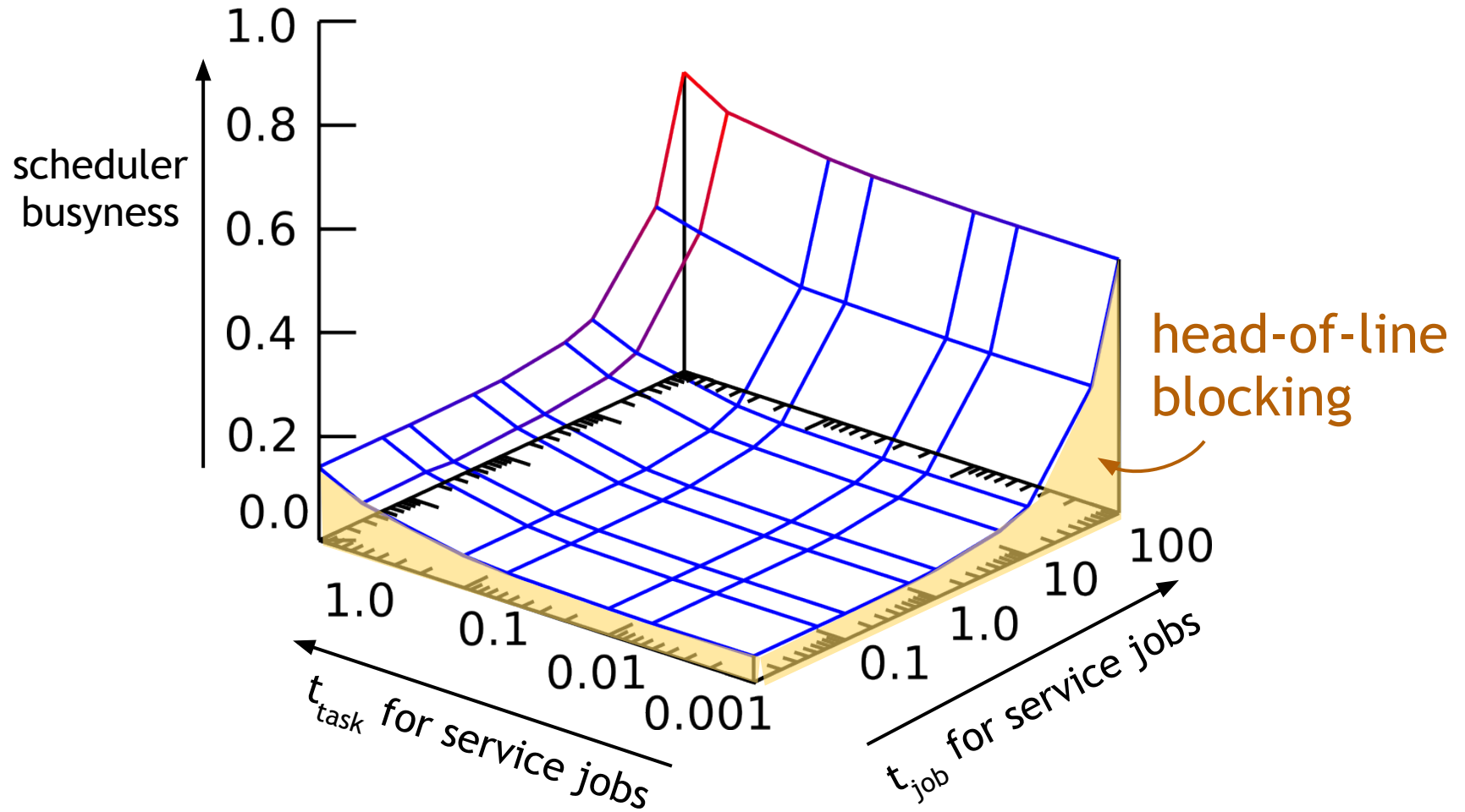
Experiment 1:

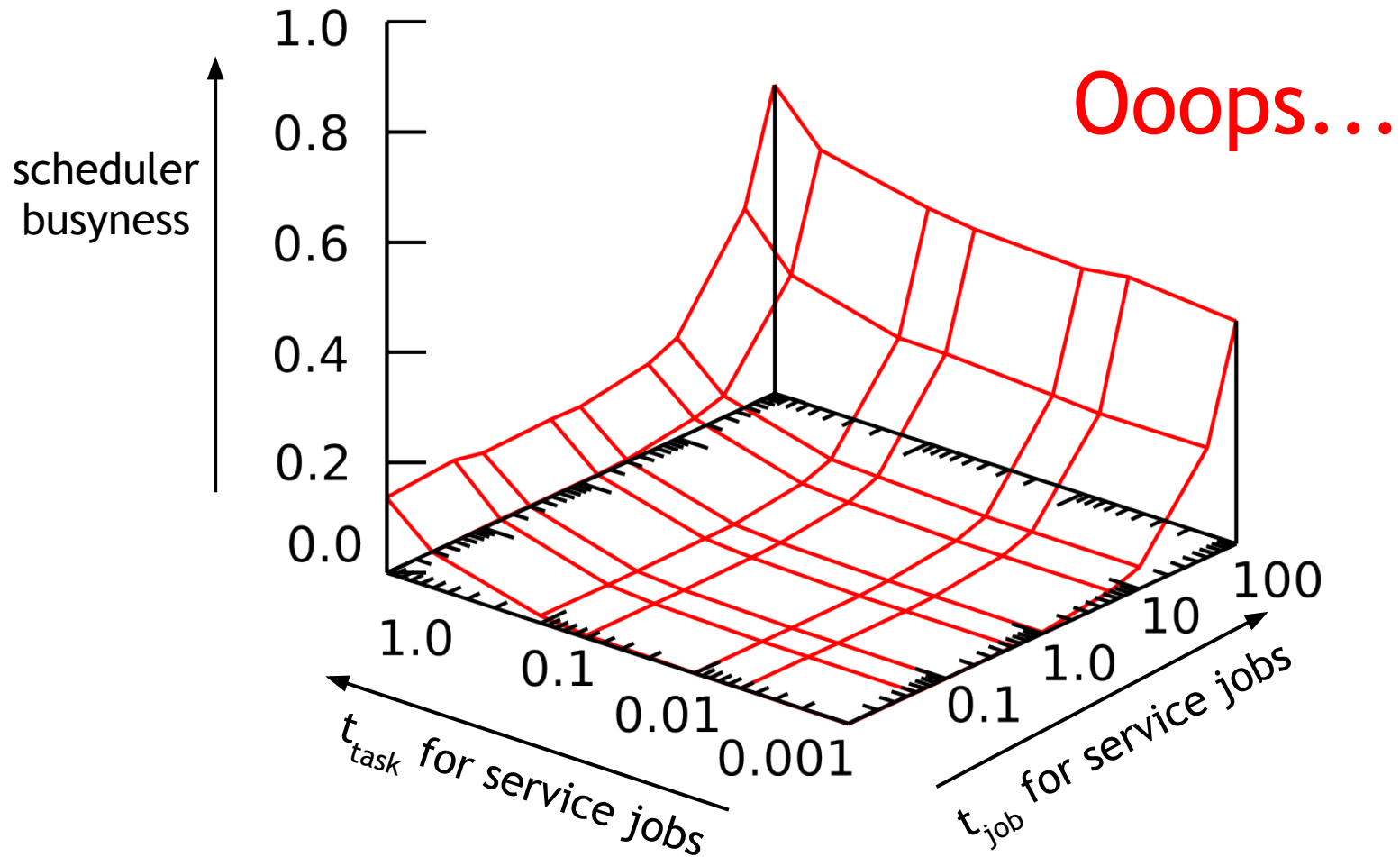
How do does the shared-state design compare with other architectures?

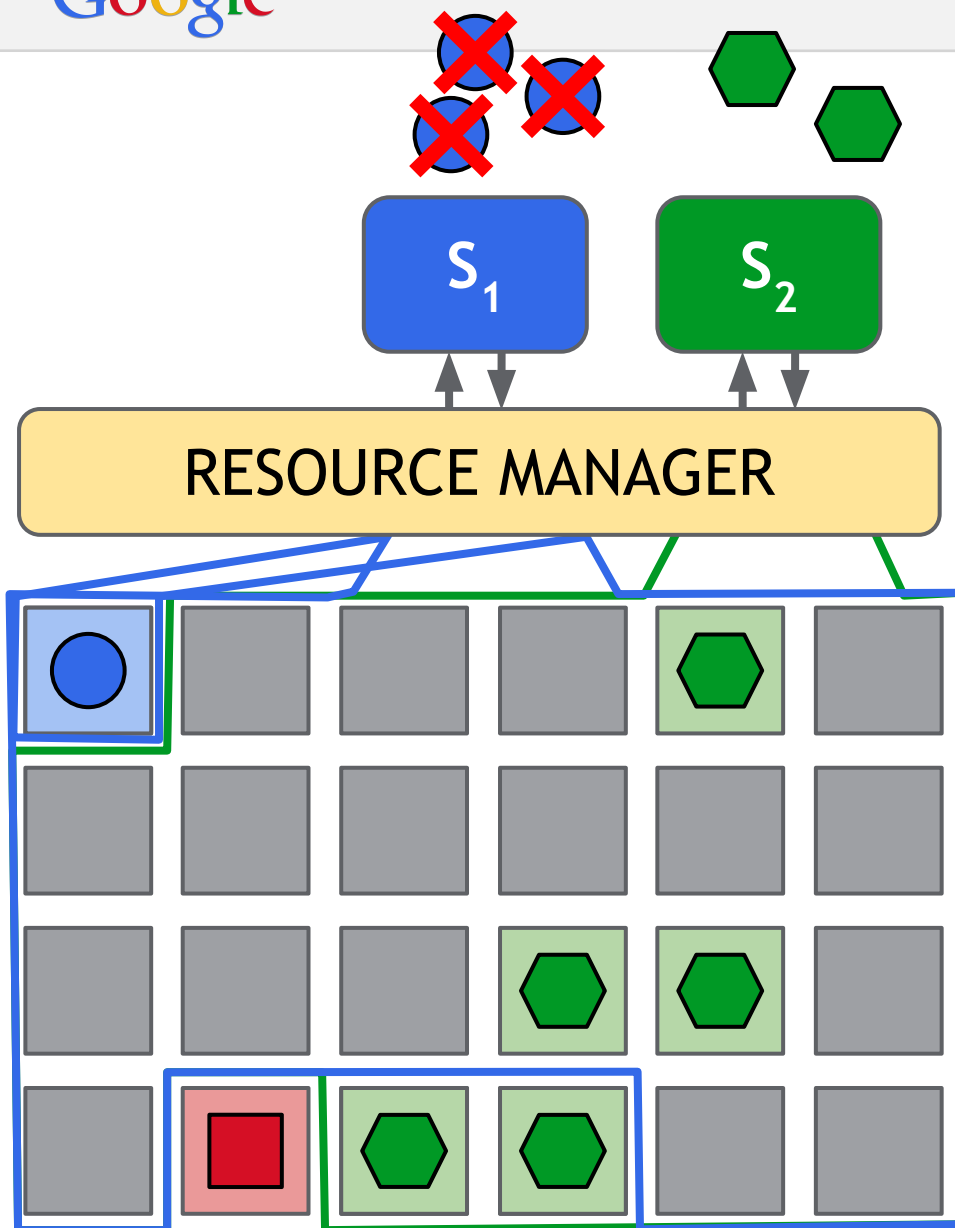
Experiment details:

- all clusters, 7 simulated days
- 2 schedulers
- varying **Service** scheduler

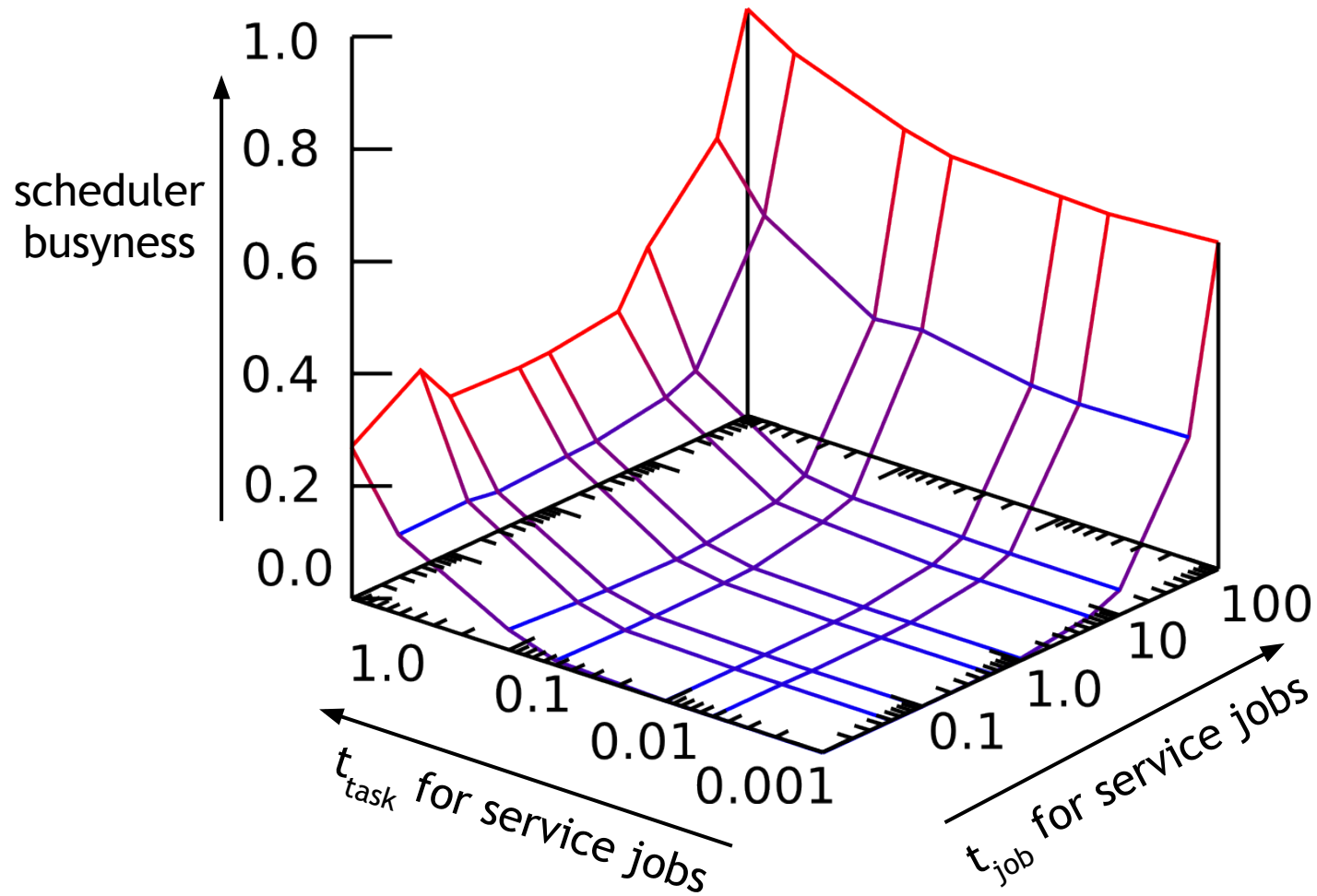


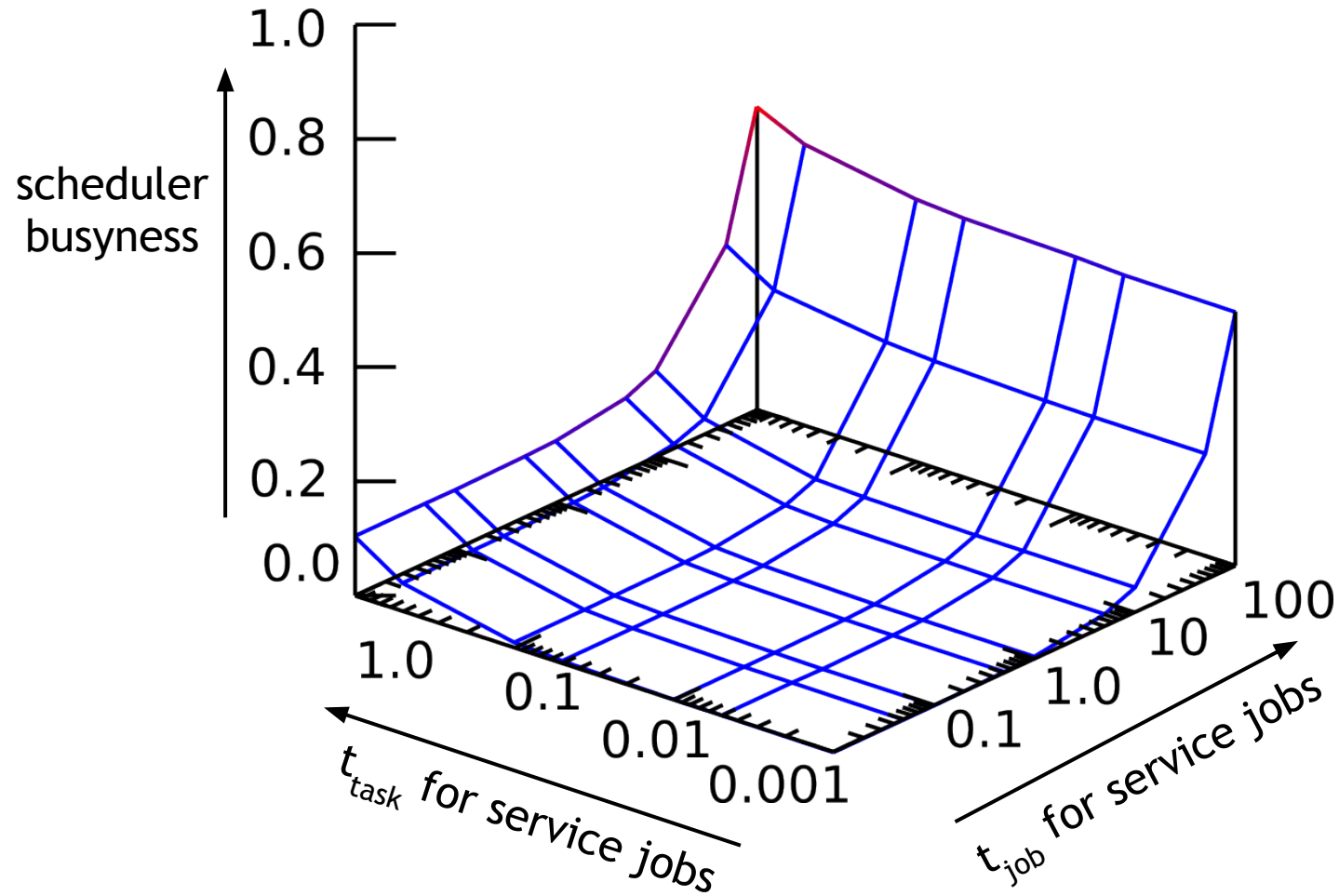






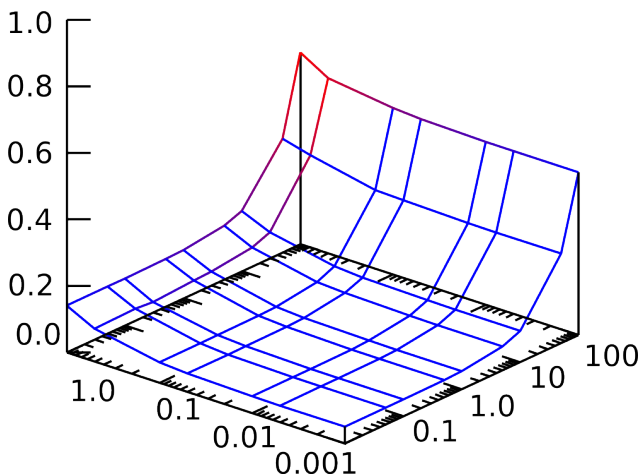
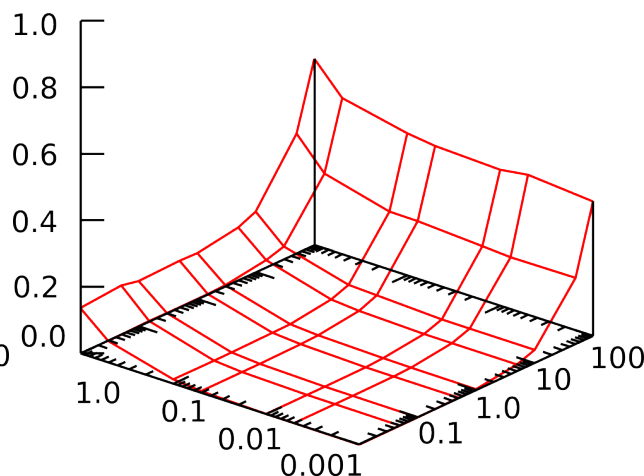
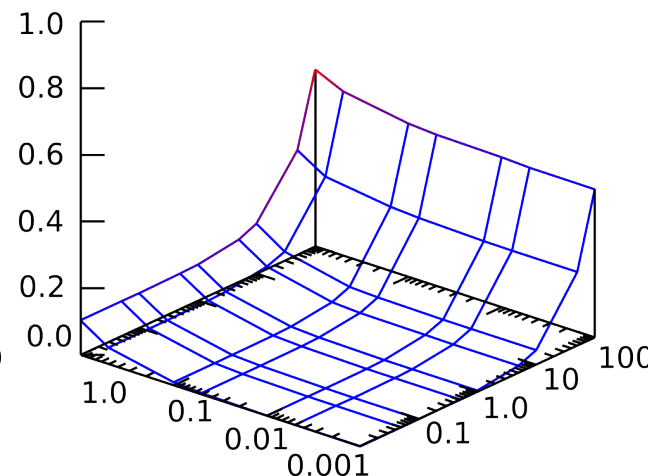
1. **Green** receives offer of all available resources.
 2. **Blue's** task finishes.
 3. **Blue** receives tiny offer.
 4. **Blue** cannot use it.
- [repeat many times]*
5. **Green** finishes scheduling.
 6. **Blue** receives large offer.
- By now, it has given up.**





TAKEAWAY

The Omega shared-state model performs as well as a (complex) monolithic multi-path scheduler.

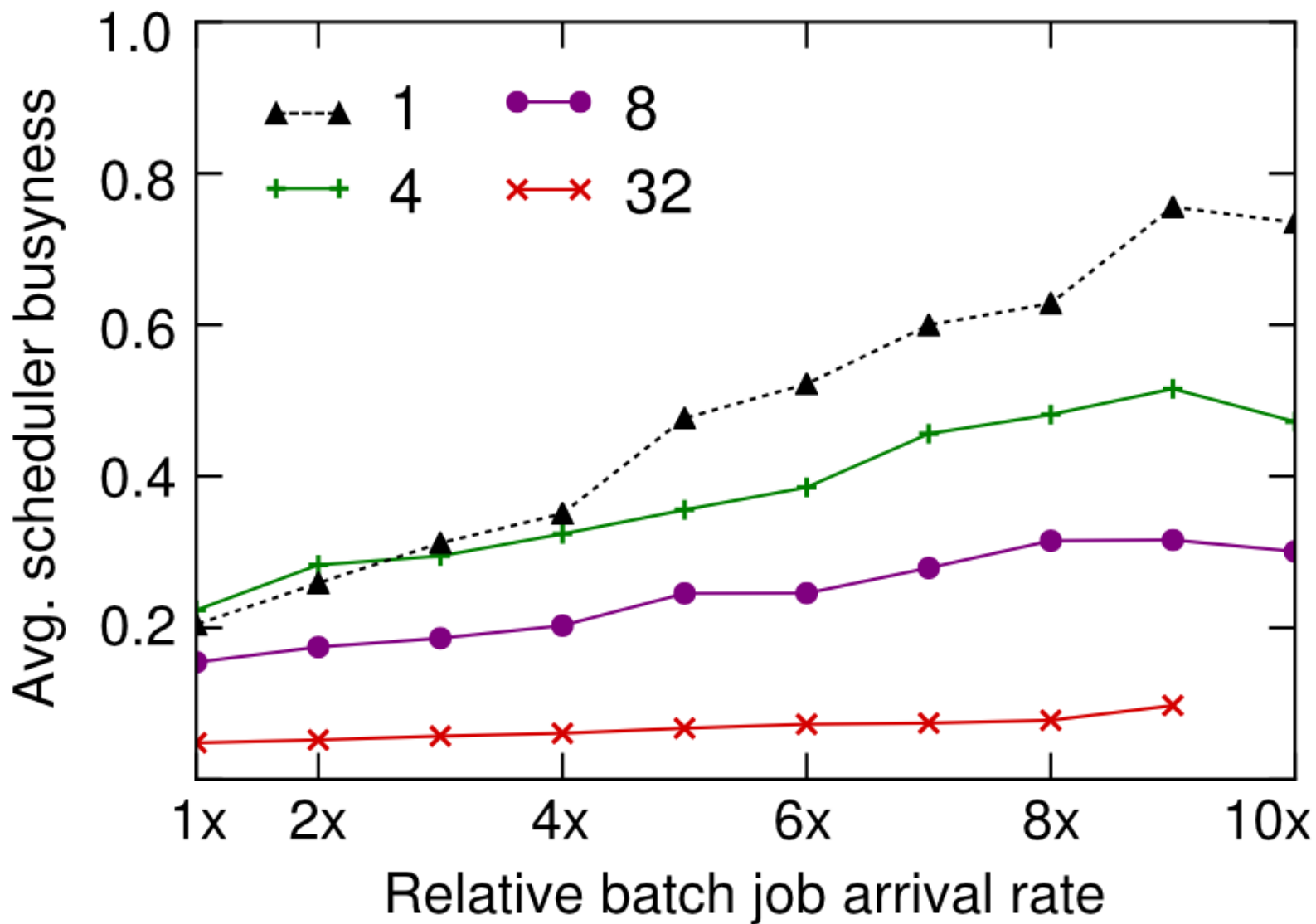
**Monolithic****Mesos****Omega**


Experiment 2:

Does the shared-state design scale to many schedulers?

Experiment details:

- cluster B, 7 simulated days
- 2 schedulers
- varying job arrival rate and number of schedulers



- 1) intro & motivation
- 2) workload characterization
- 3) comparison of approaches
- 4) trace-based simulation** 
- 5) flexibility case study

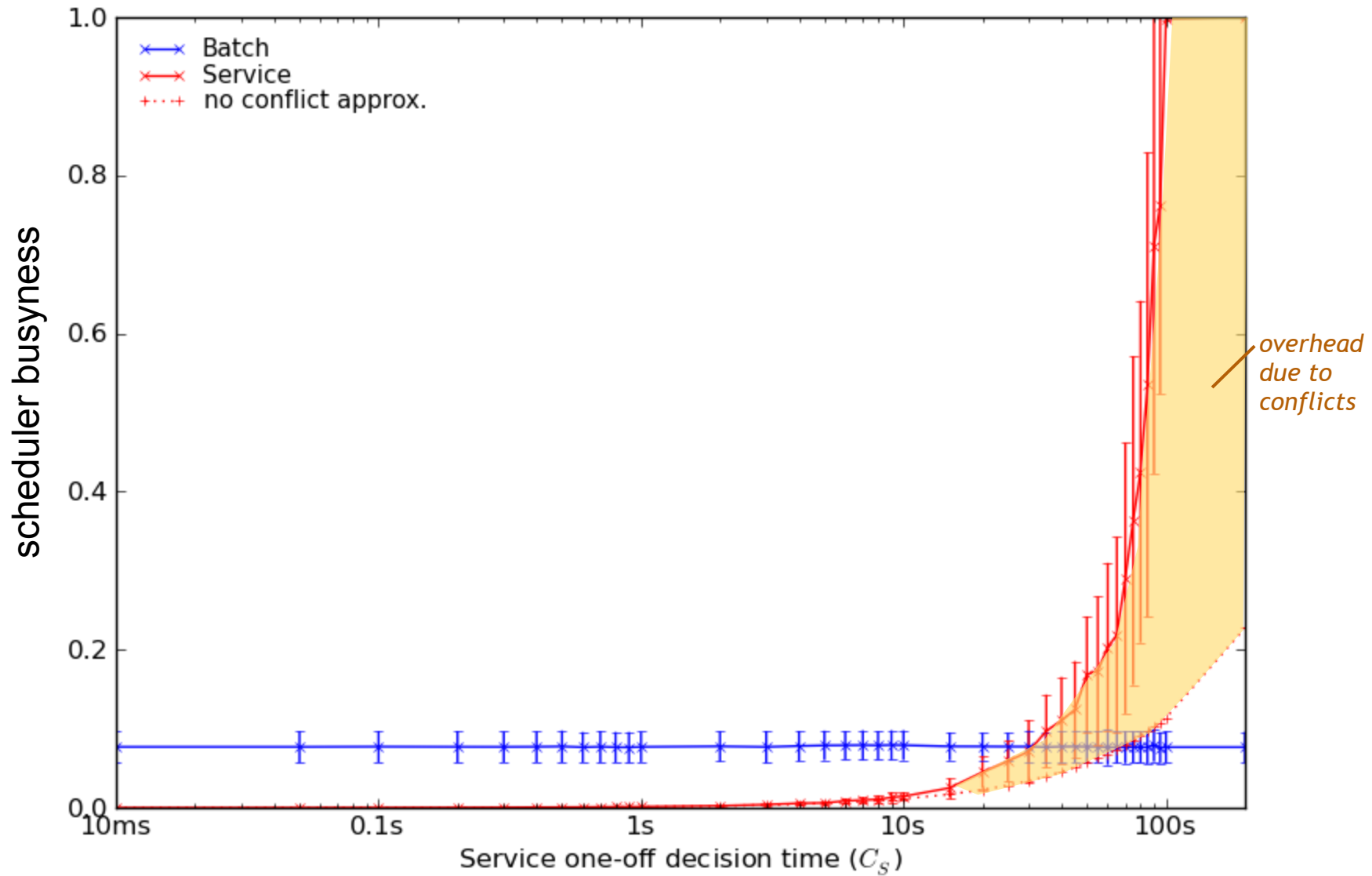
	<i>lightweight simulator</i>	<i>high-fidelity simulator</i>
<i>machines</i>	homogeneous	real-world
<i>job parameters</i>	empirical distribution	workload trace
<i>constraints</i>	not supported	supported
<i>scheduling algorithm</i>	random first fit	Google algorithm
<i>runtime</i>	fast (24h \approx 5min)	slow (24h \approx 2h)

Experiment 3:

How much scheduler interference do we see with real Google workloads?

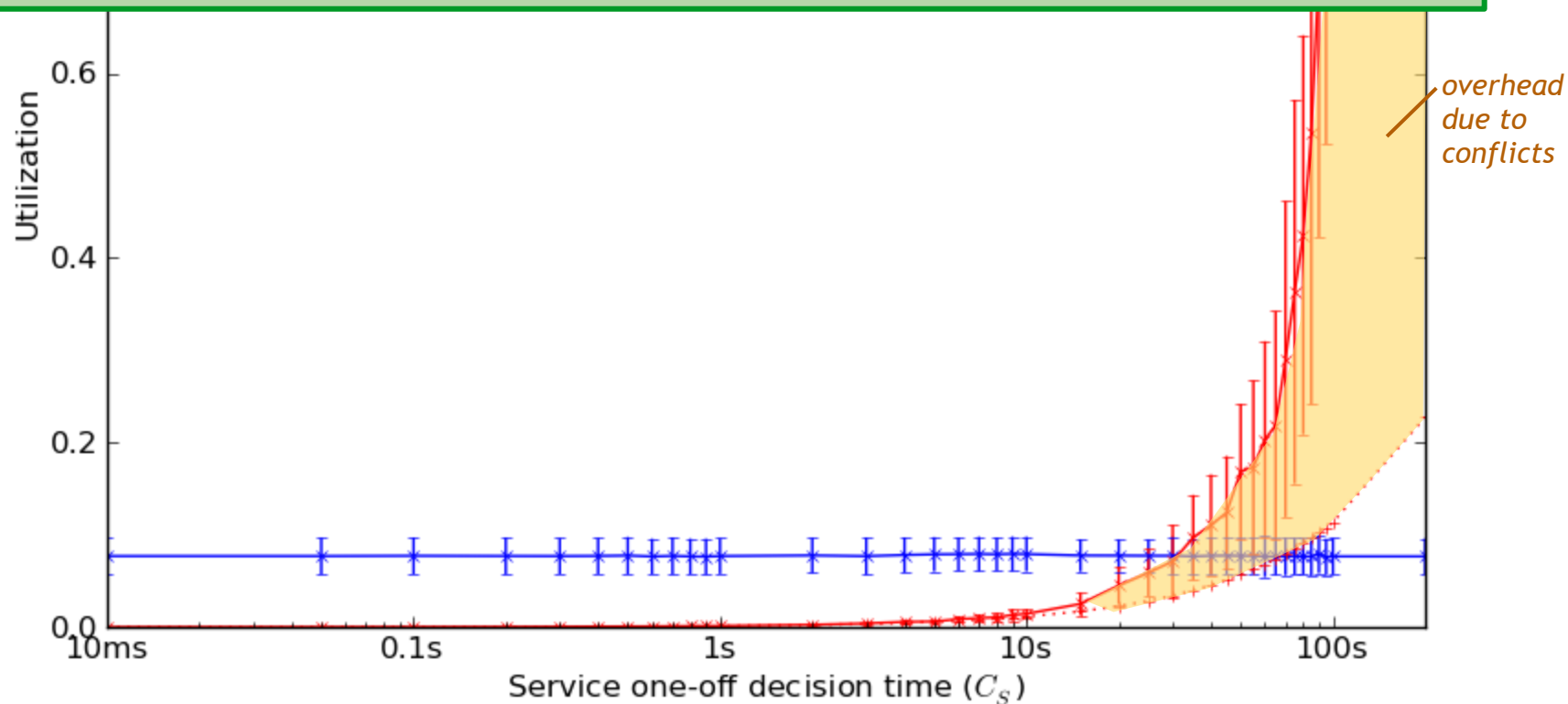
Experiment details:

- cluster C, 29 days
- 2 schedulers,
non-uniform decision time
- varying **Service** scheduler

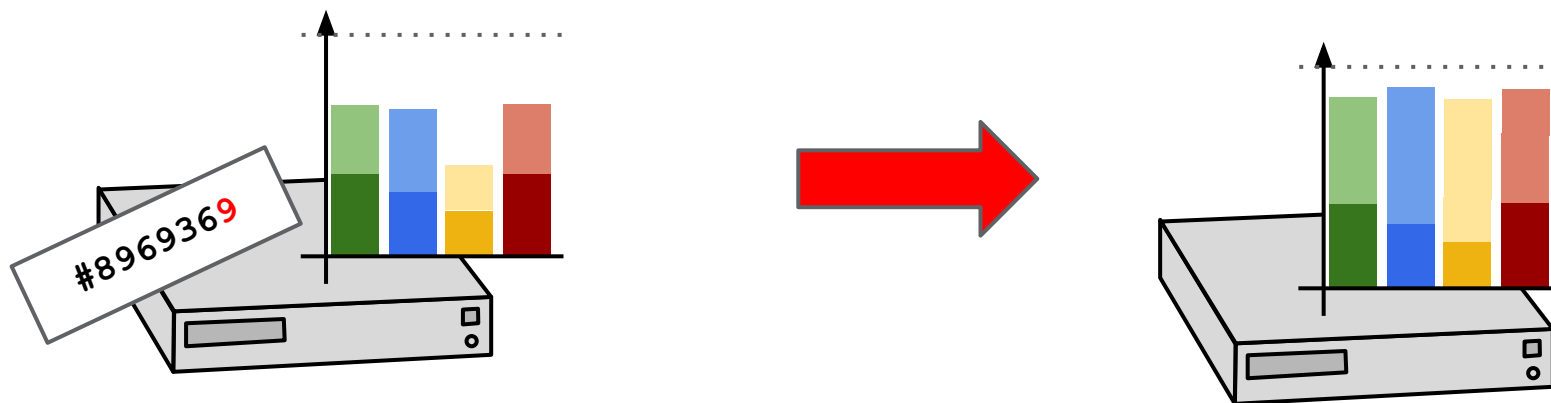


TAKEAWAY

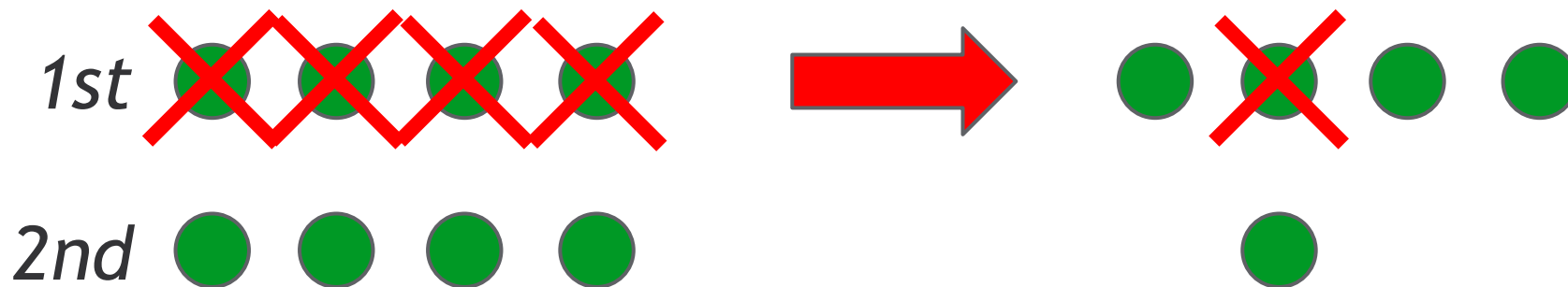
Interference is higher for real-world settings.



1. Fine-grained conflict detection



2. Incremental commits

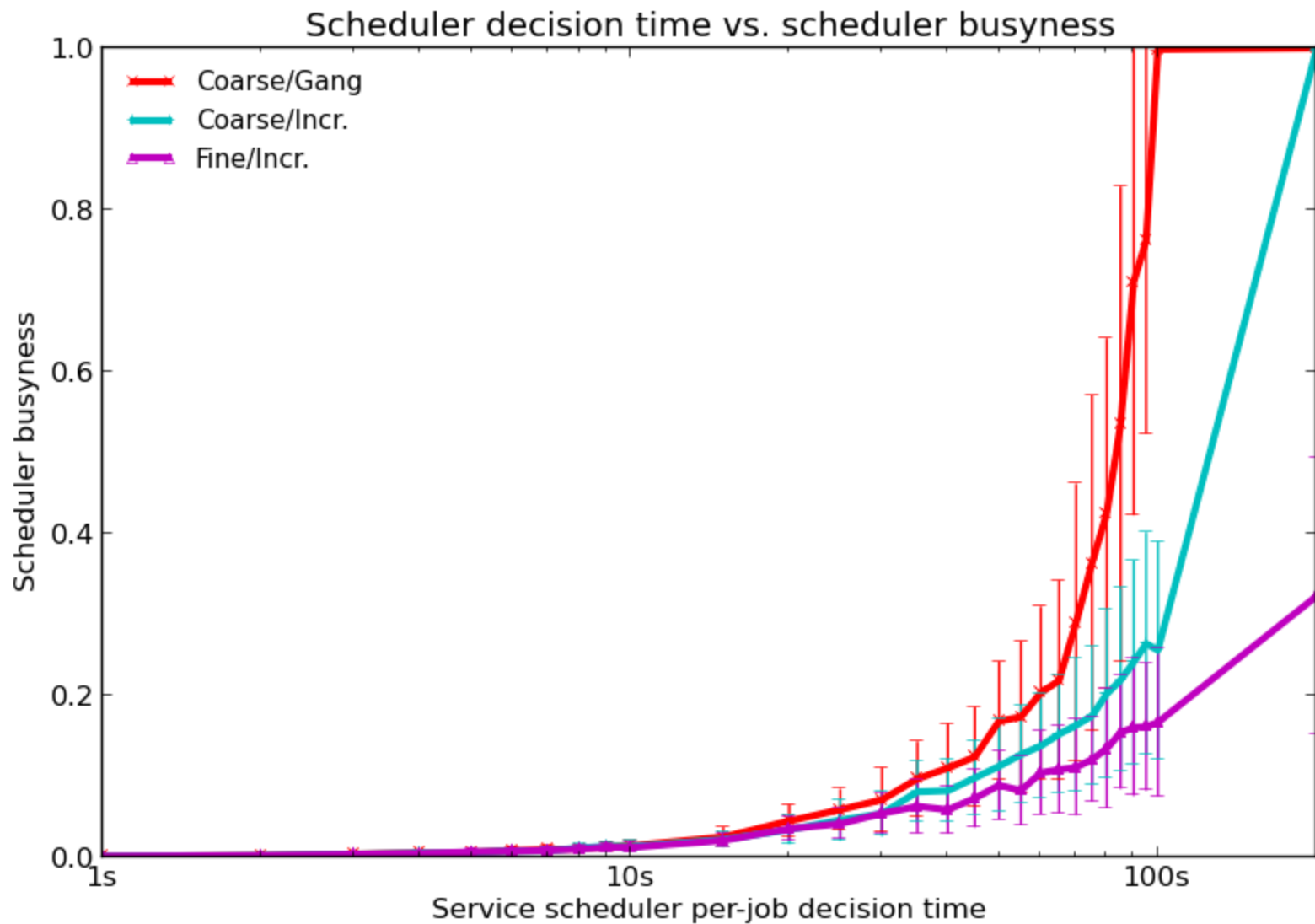


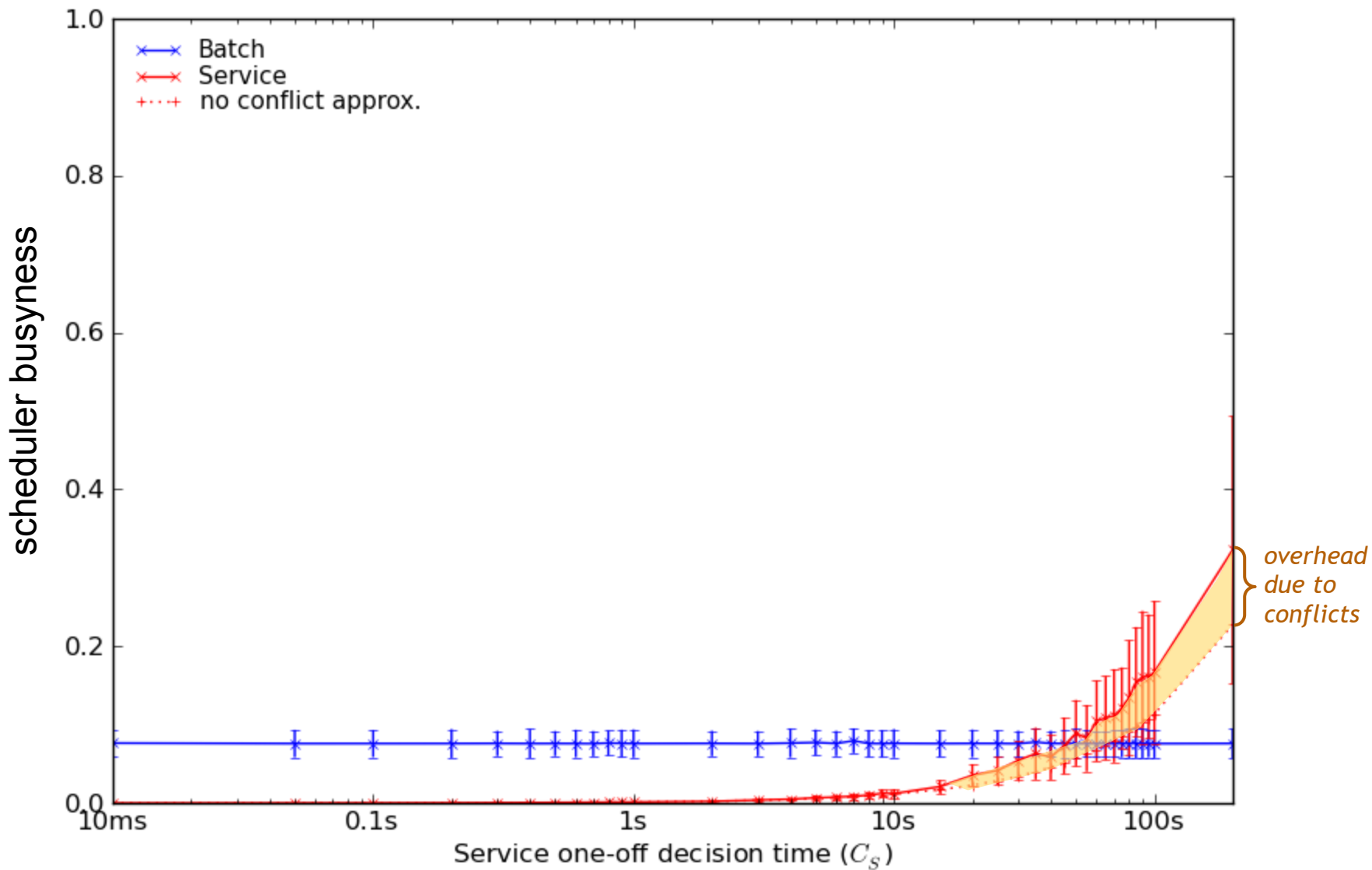
Experiment 4:

How do the optimizations affect performance?

Experiment details:

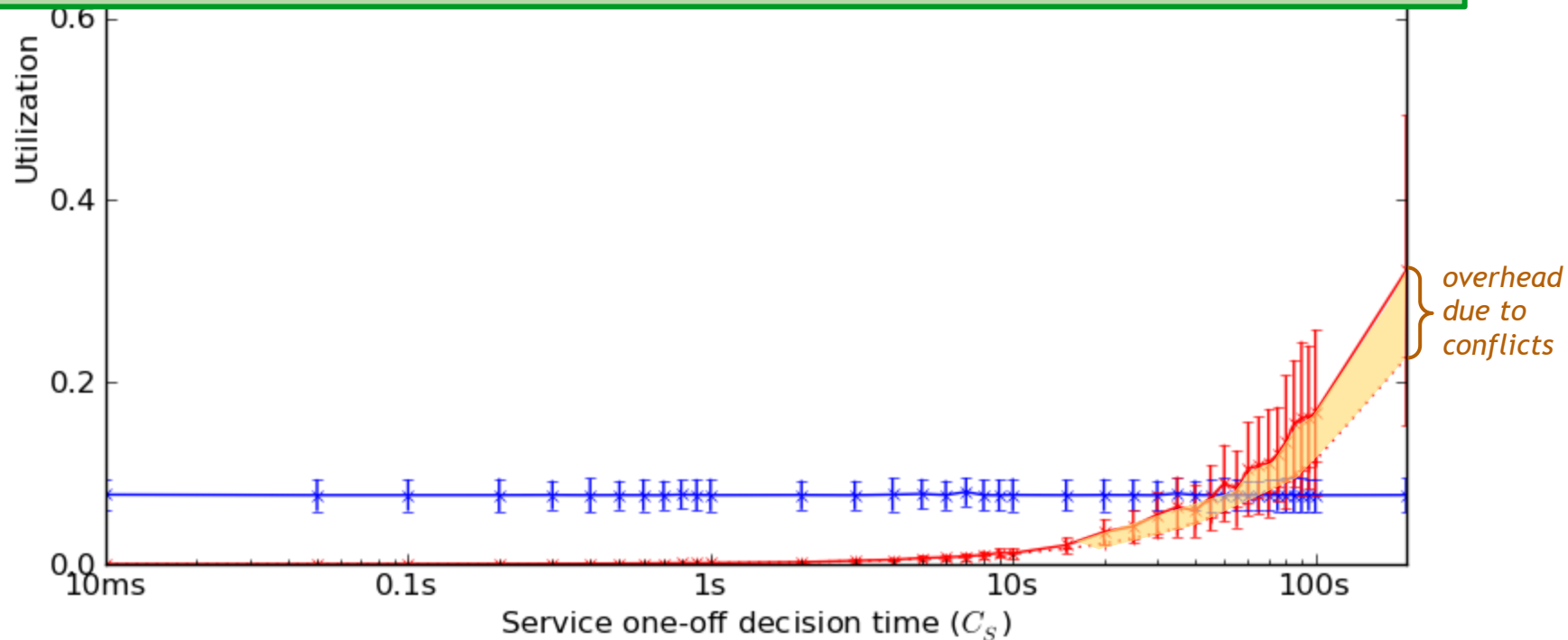
- cluster C, 29 days
- 2 schedulers,
non-uniform decision time
- varying **Service** scheduler





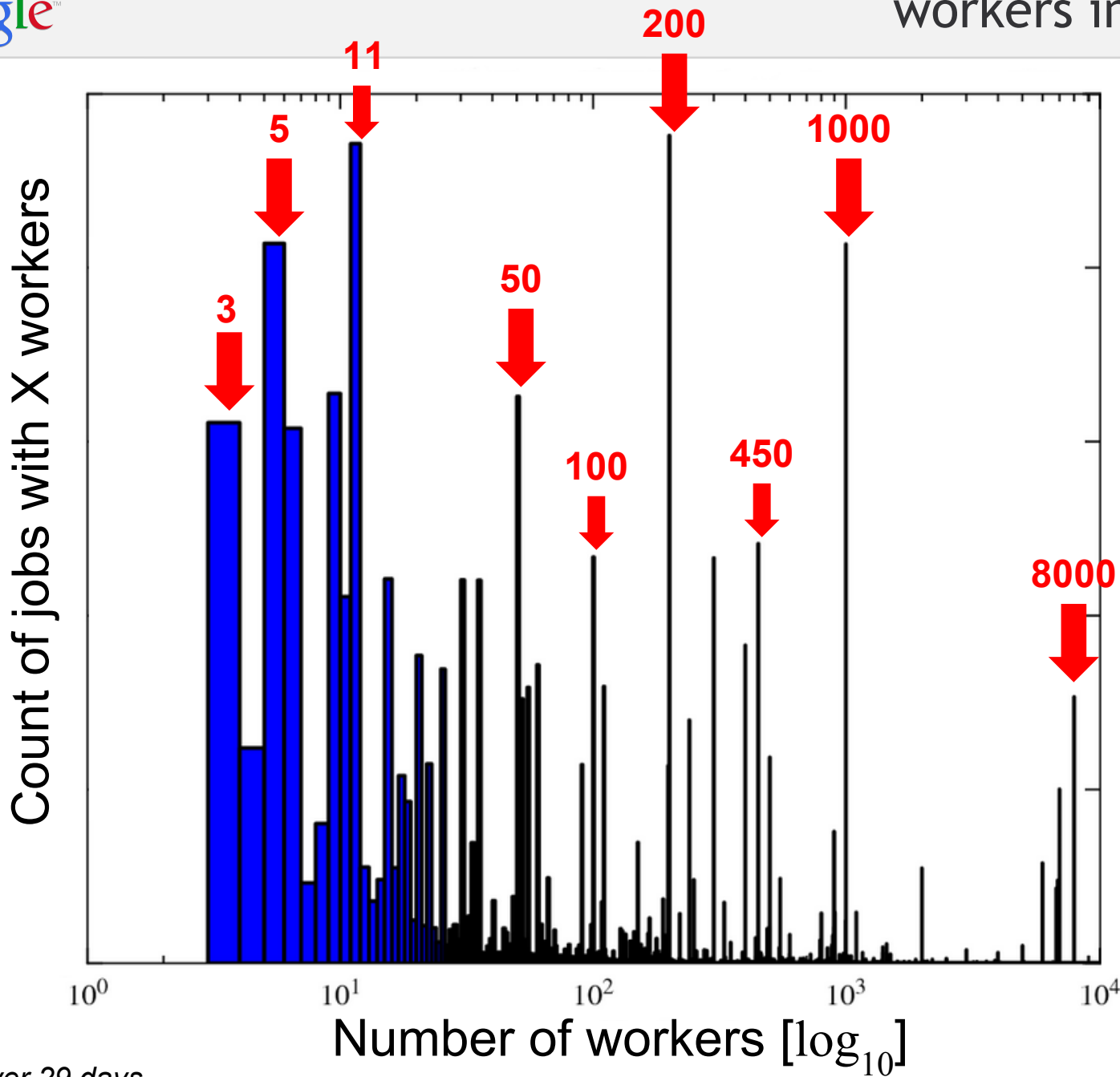
TAKEAWAY

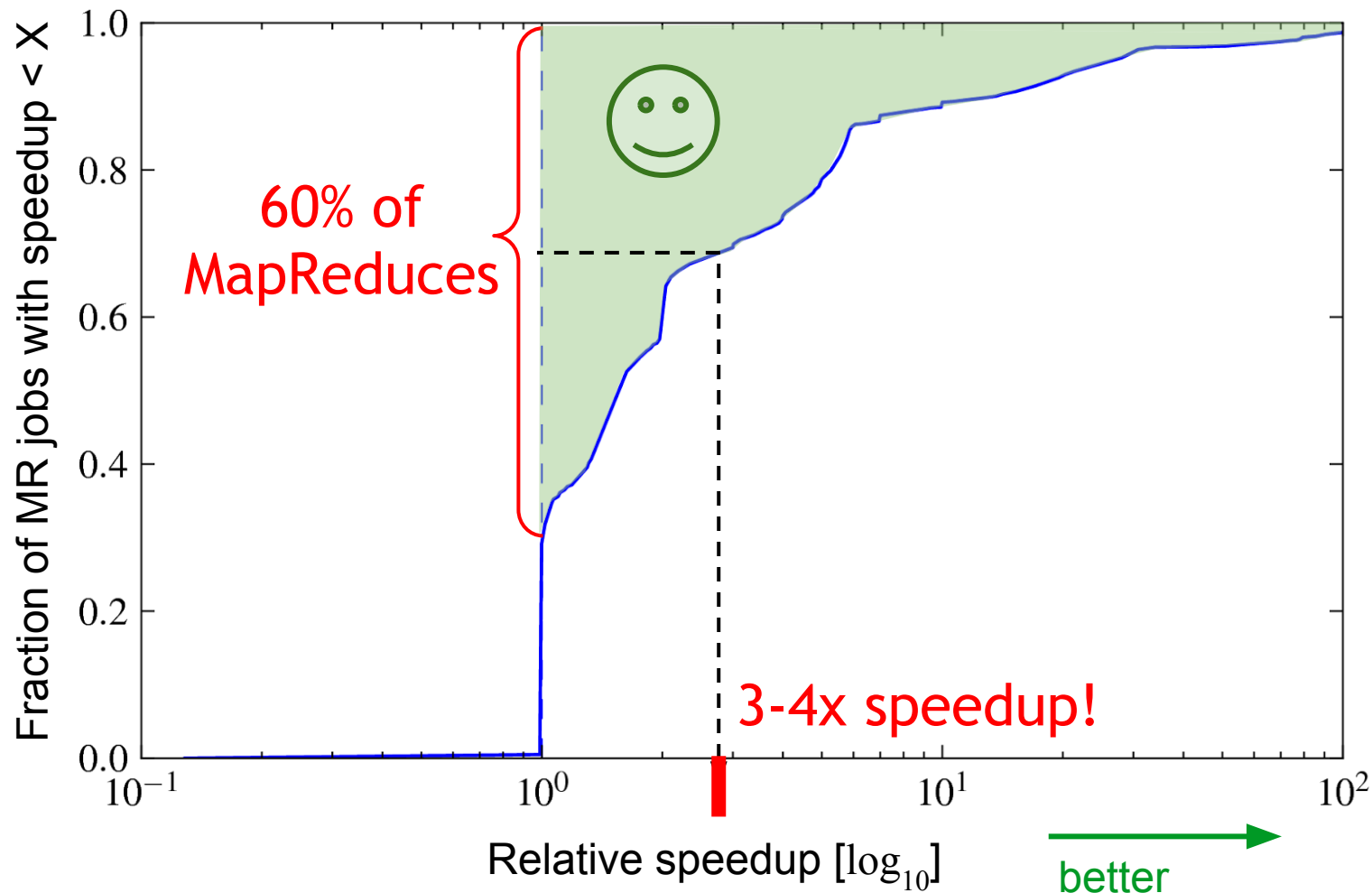
We can make simple improvements that significantly improve scalability.



Case study

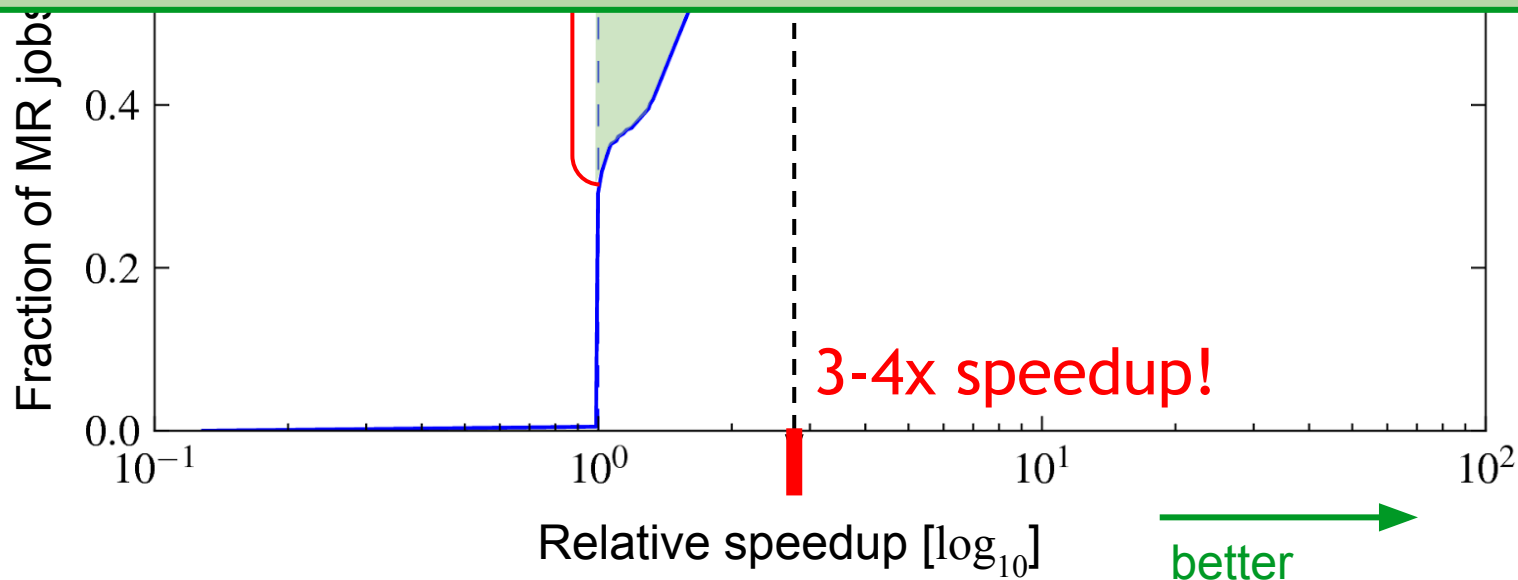
MapReduce scheduler with
opportunistic extra resources





TAKEAWAY

The Omega approach gives us the flexibility to easily support custom policies.



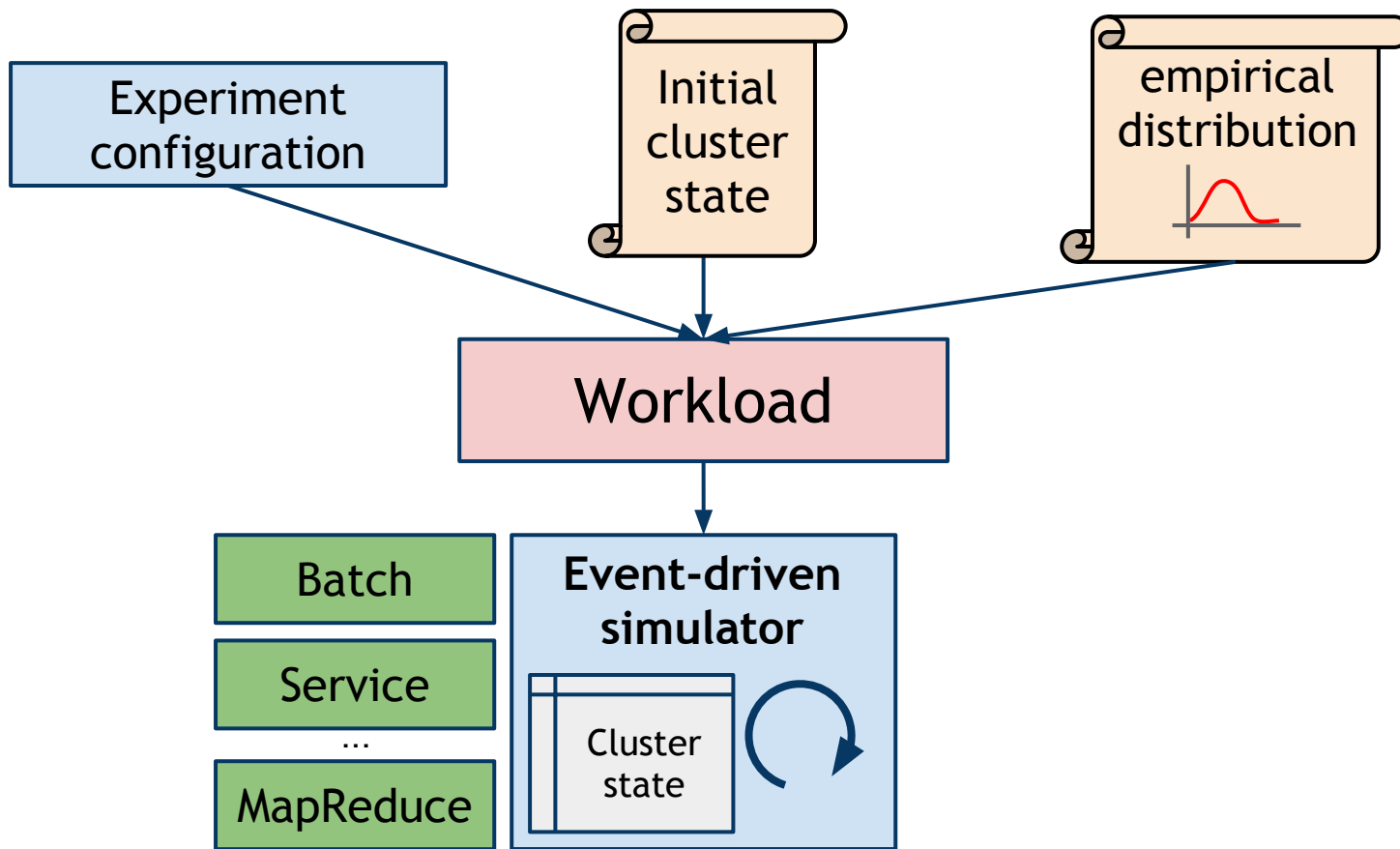
TAKEAWAYS

Flexibility and scale require parallelism,
parallel scheduling works *if you do it right*,
and
using shared state is the way to do it right!

BACKUP SLIDES

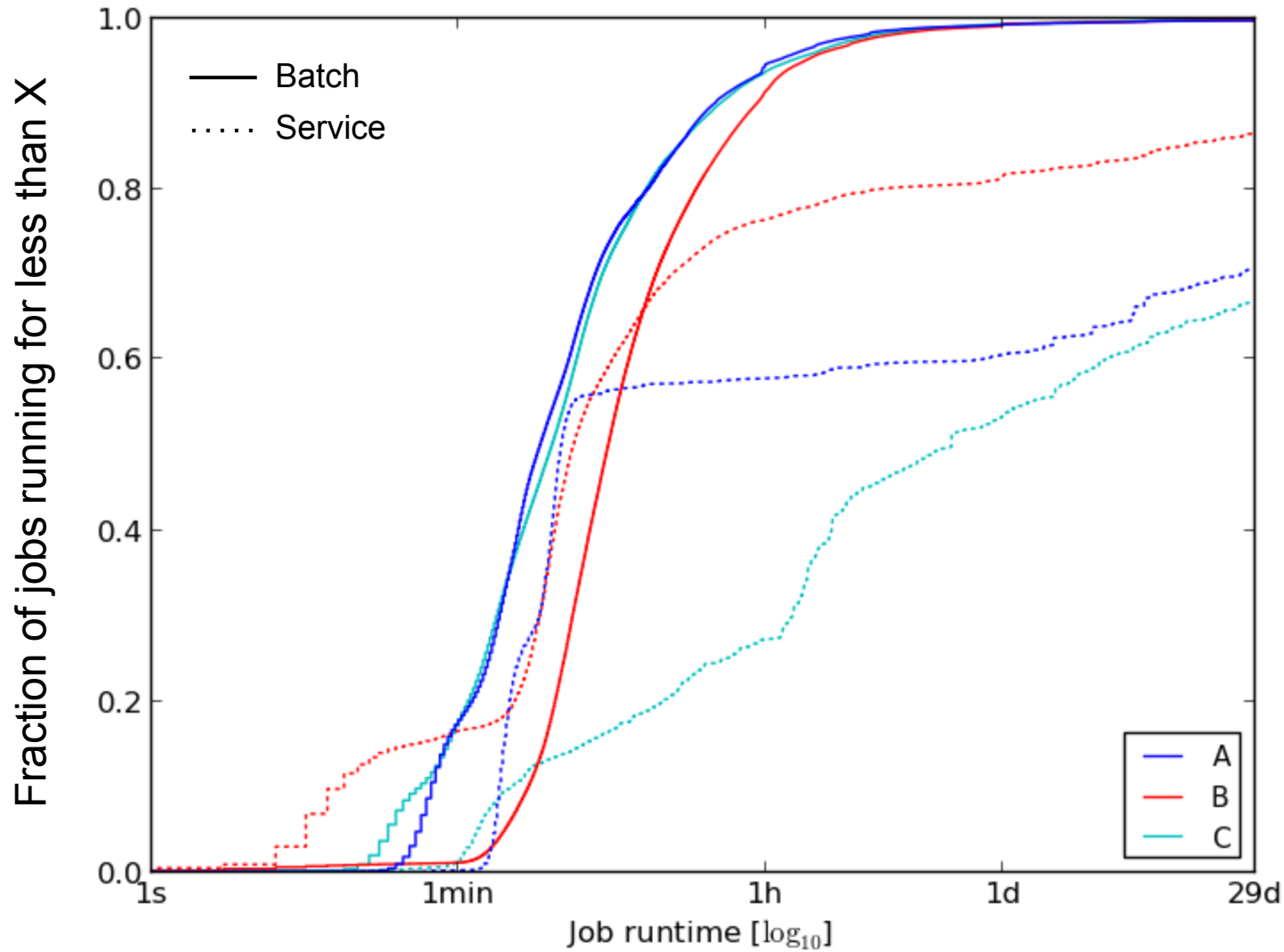
Why might scheduling take 60 seconds?

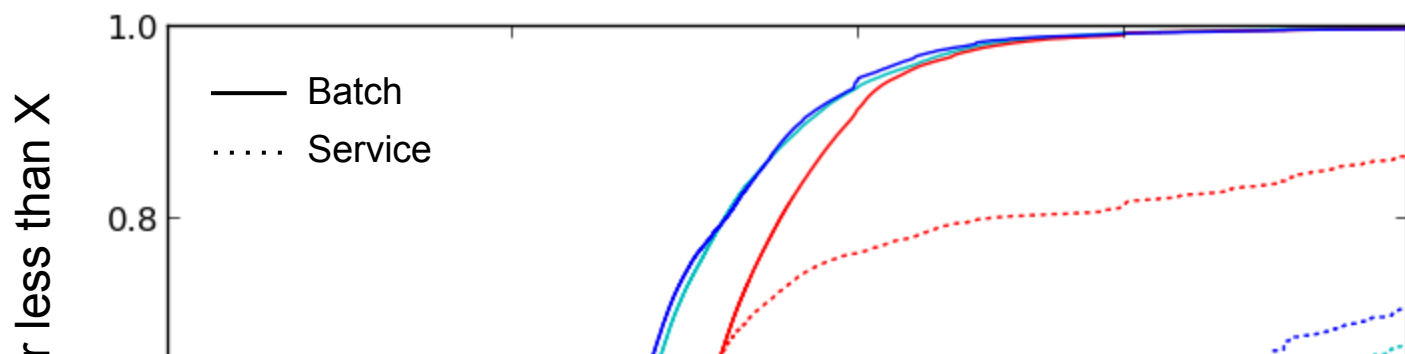
- Large jobs (1,000s of tasks)
- Optimization algorithms (constraint solving, bin packing)
- Very picky jobs in a full cluster (preemption consequences)
- Monte Carlo simulations (fault tolerance)



Code [soon to be] available:

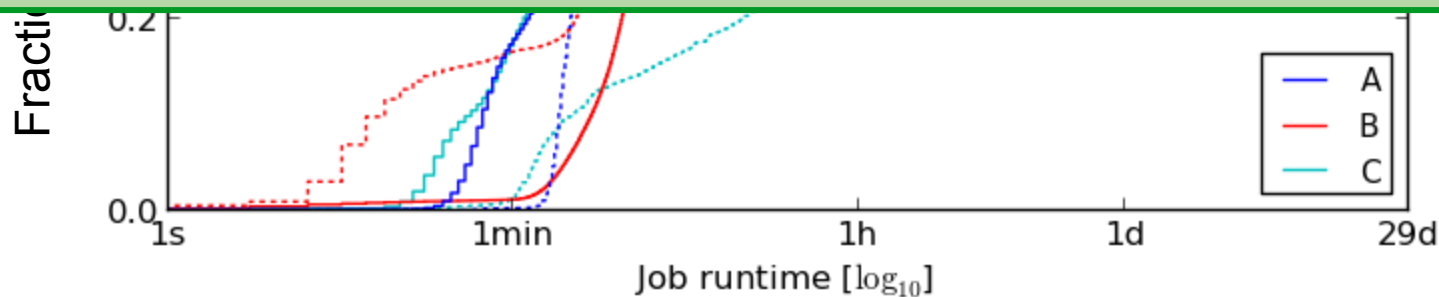
<http://code.google.com/p/cluster-scheduler-simulator>

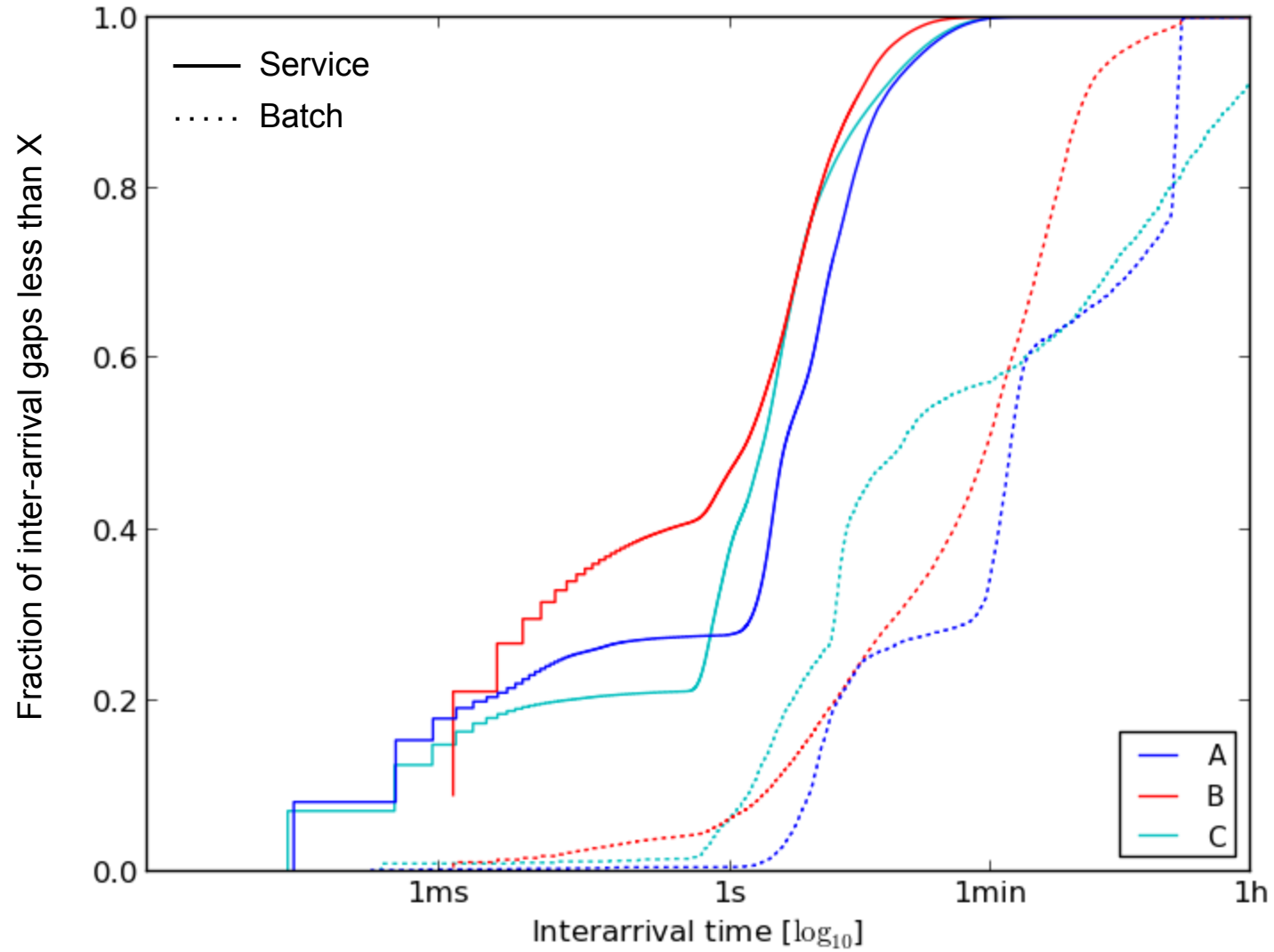


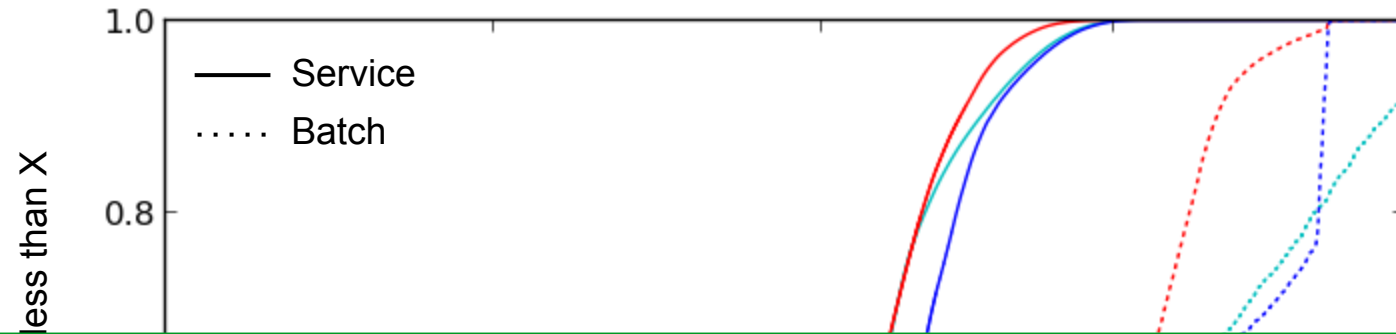


TAKEAWAY

*Service jobs, once scheduled, run for **much longer** than batch jobs do.*

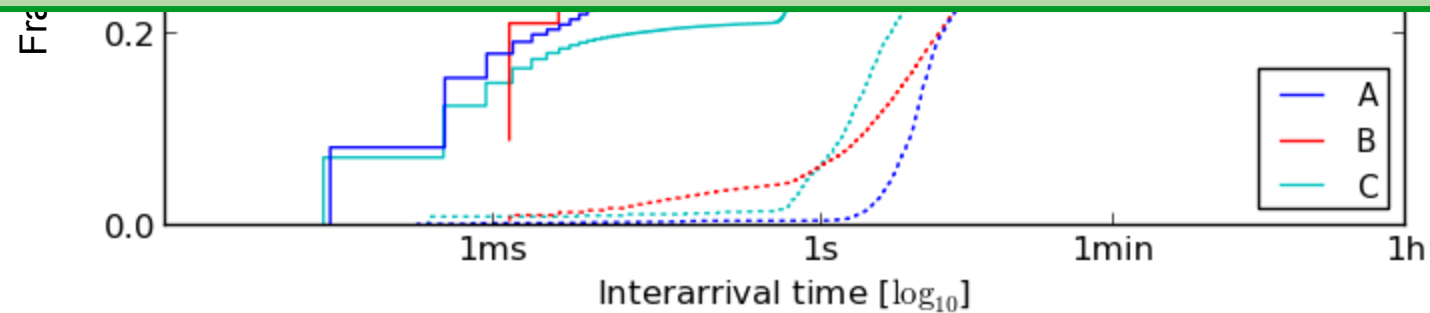






TAKEAWAY

Service jobs arrive much less frequently than batch jobs do.

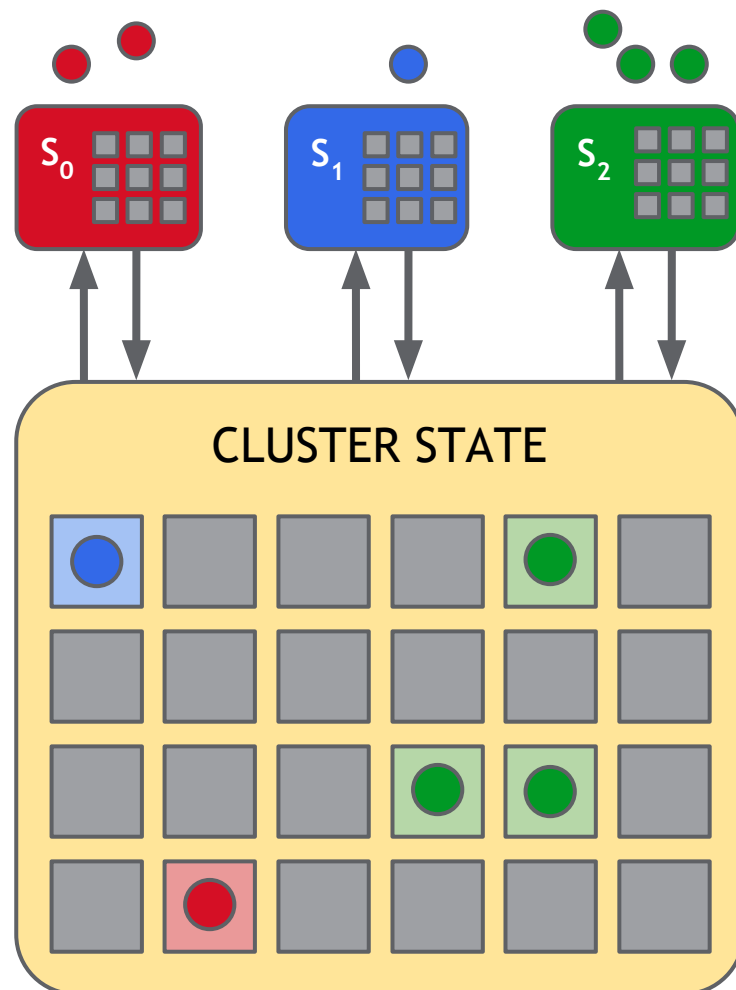


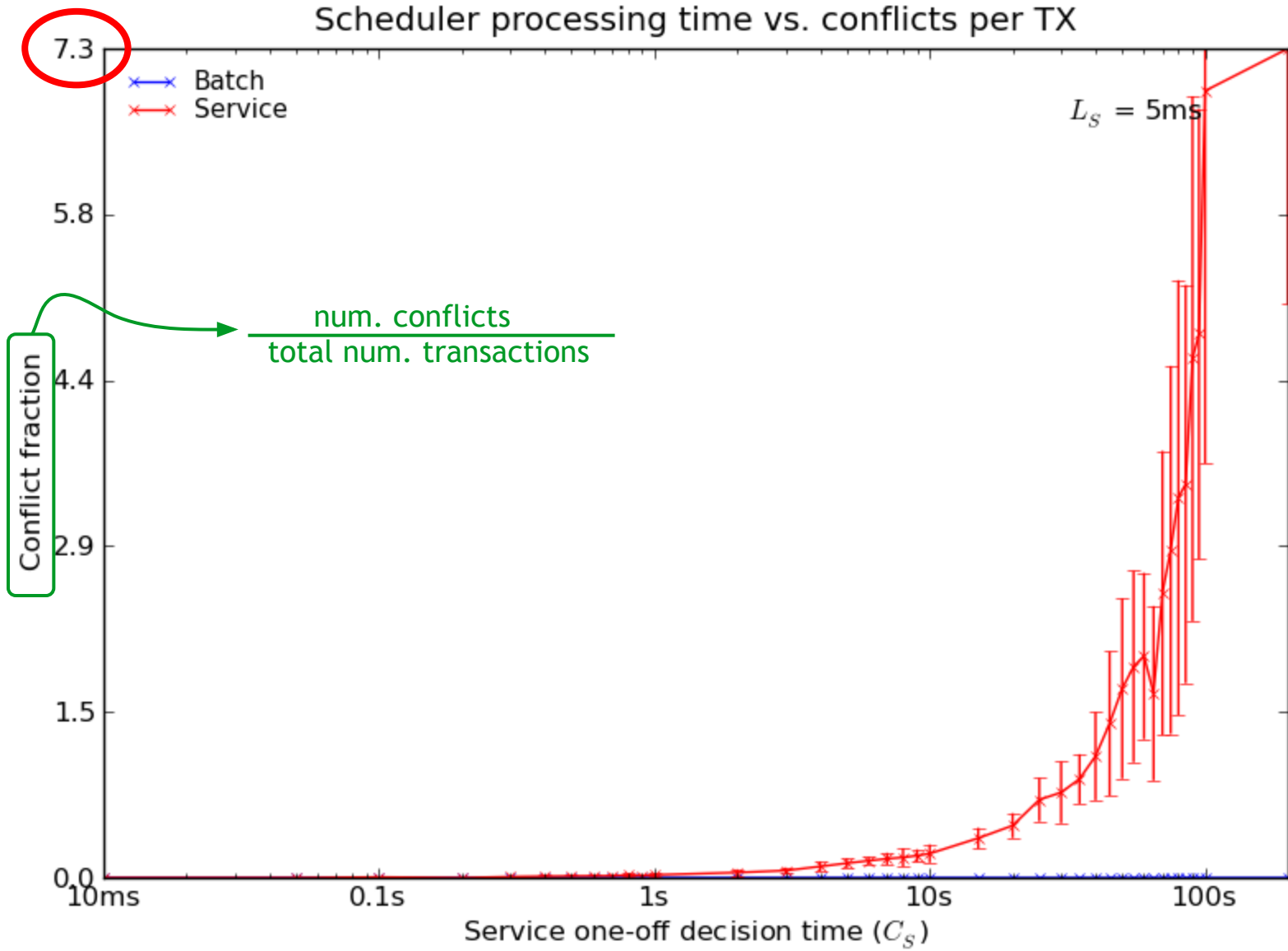
Shared state

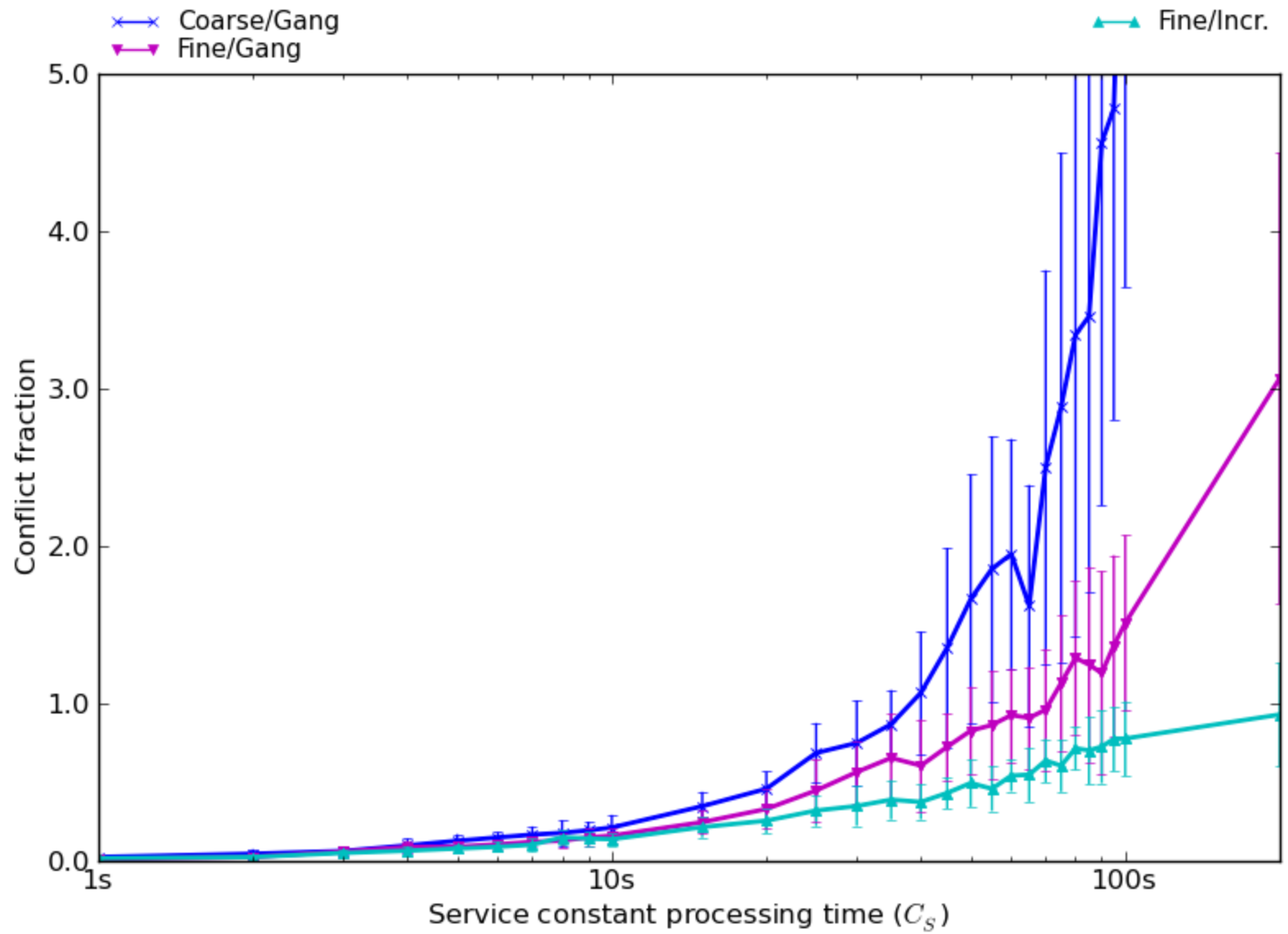
- Deltas against shared state
- Easy to develop & maintain
- Heterogeneous scheduling logic supported

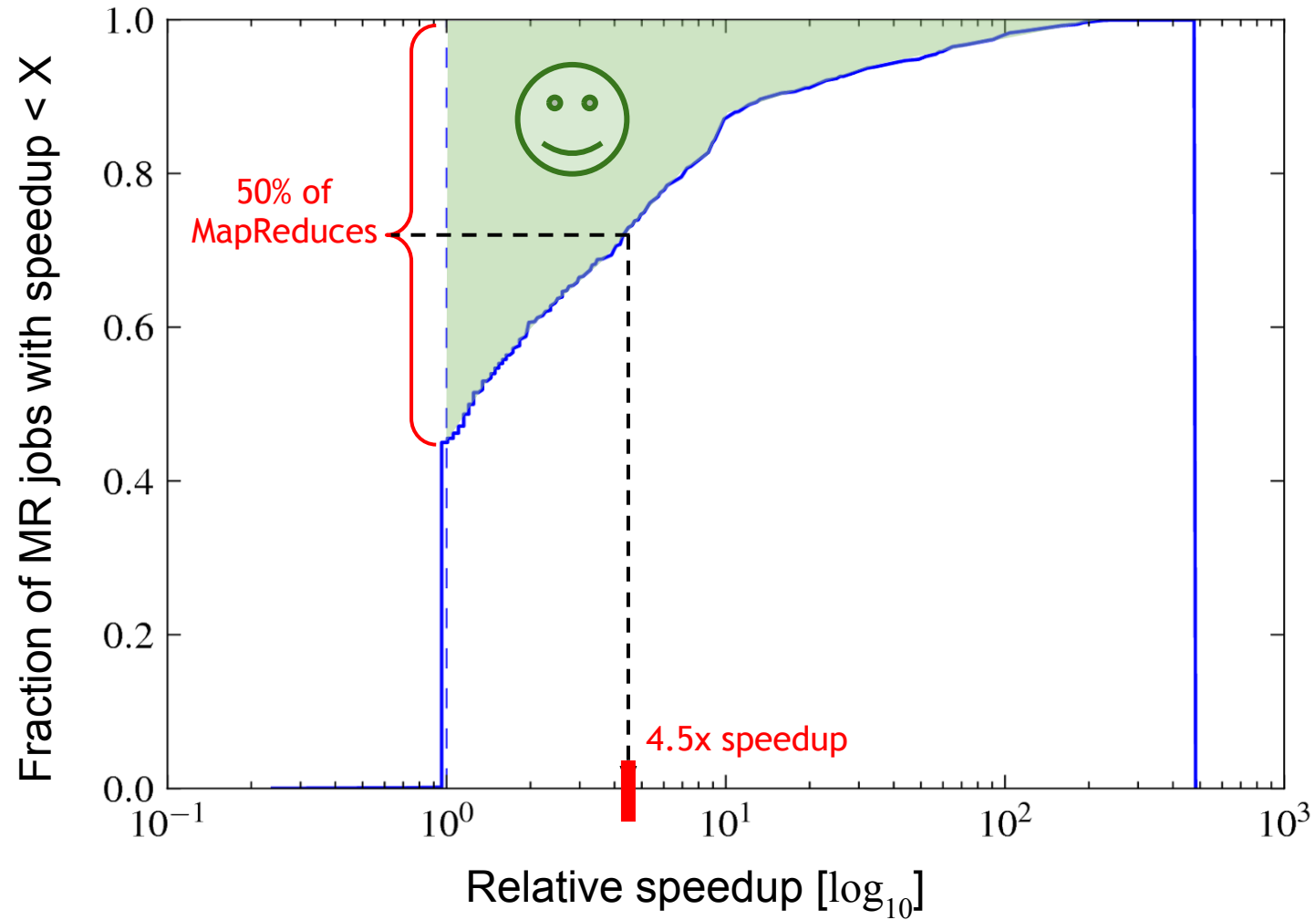
Optimistic concurrency

- No explicit coordination required
- Post-hoc interference resolution
- Scales well



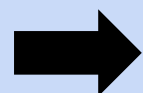






Possible problems...

- aggressive, systematically adverse workloads or schedulers
- small clusters with high overcommit



deal with using out-of-band or post-facto enforcement mechanisms