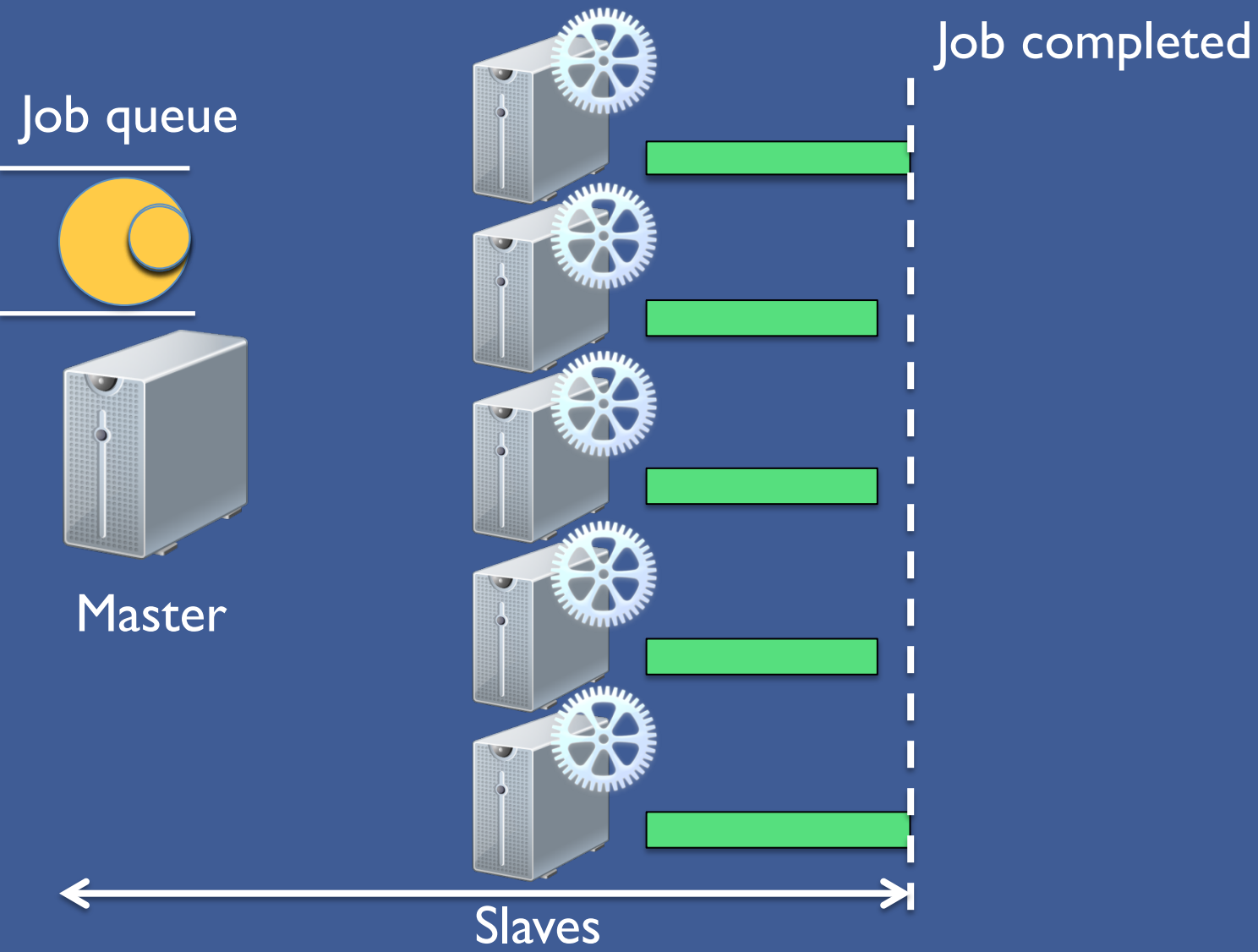# Wrangler: Predictable and Faster Jobs in Distributed Processing Systems using Machine Learning

Neeraja J. Yadwadkar (neerajay@eecs.berkeley.edu),
Bharath Hariharan,
Ganesh Ananthanarayan,
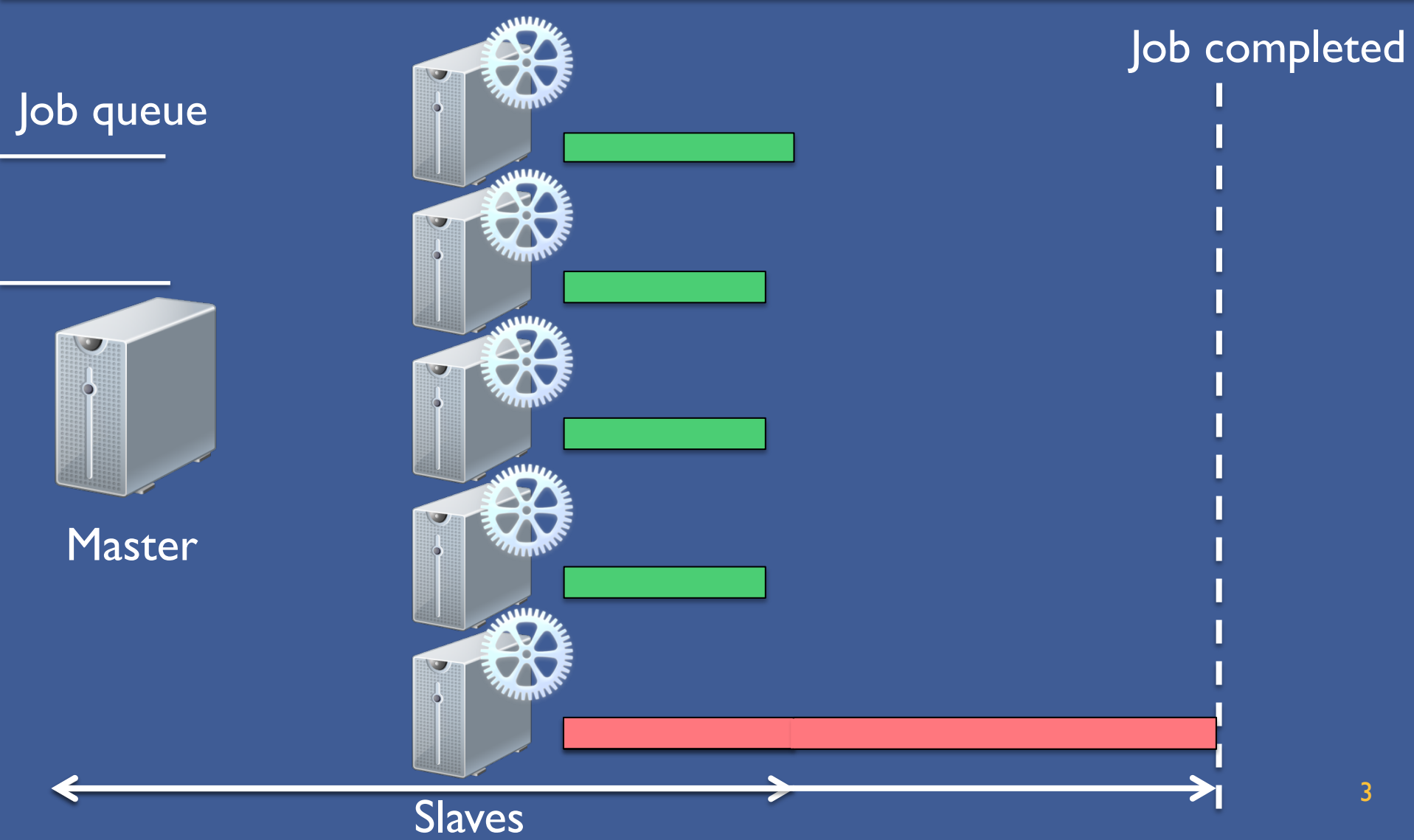Joseph Gonzalez, and Randy Katz

http://www.istc-cc.cmu.edu/

**Intel Science & Technology
Center for Cloud Computing**

# Parallel Data Analytics

Job completed

Job queue

Master

Slaves

# Stragglers

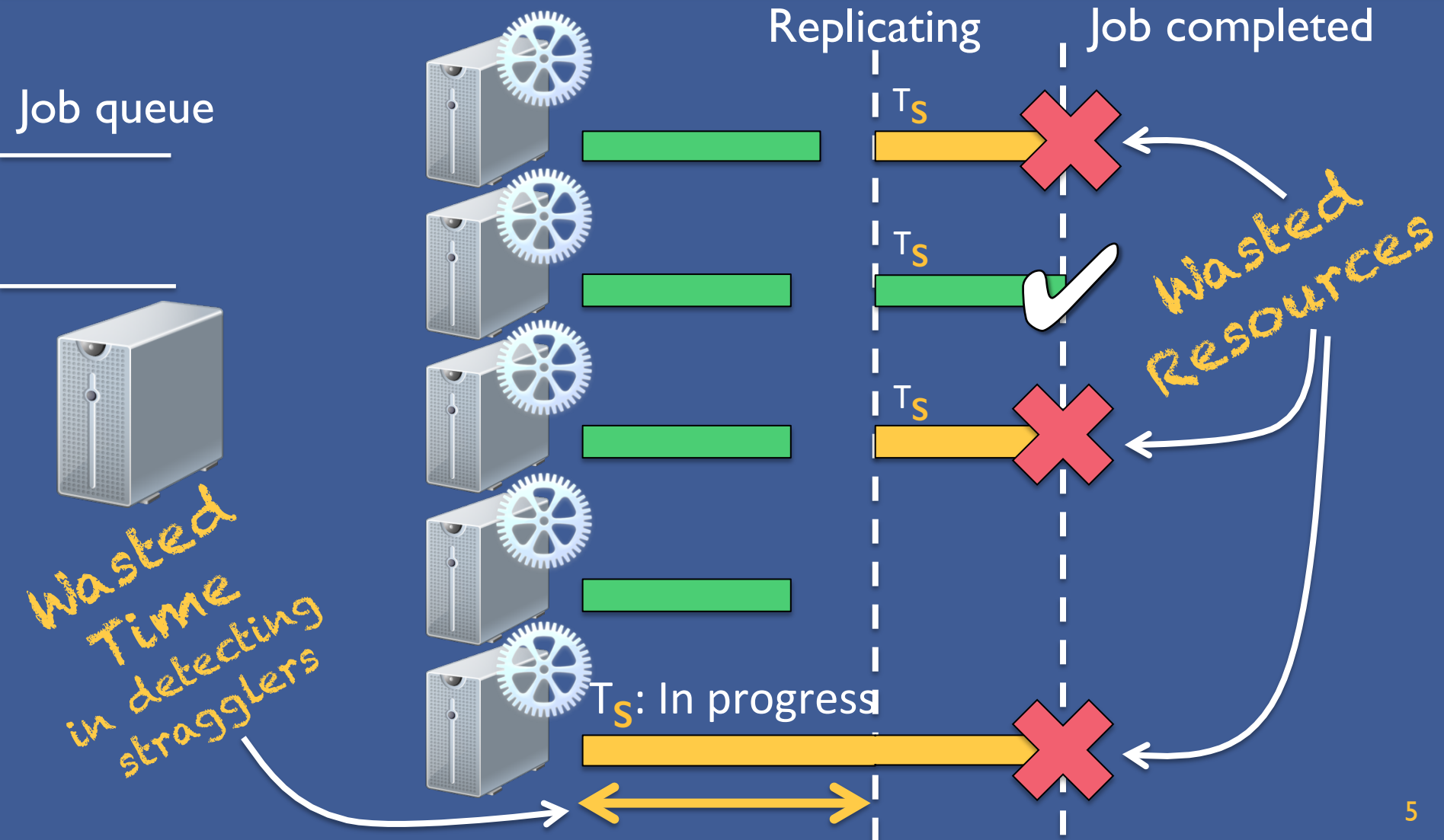Job queue

Job completed

Master
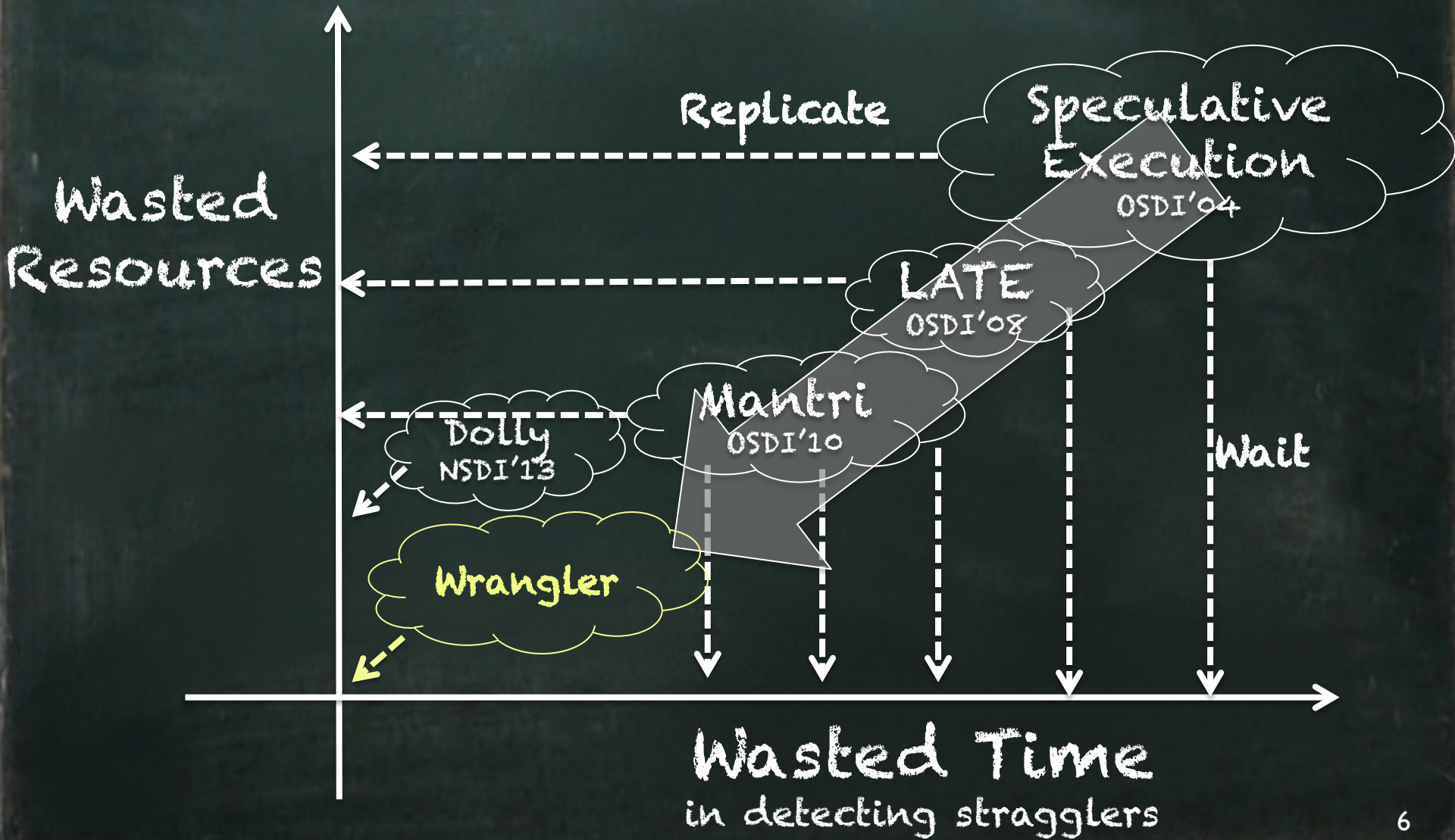
Slaves

# Impact of **Stragglers**

**Impact of Stragglers:** We measure the potential in speeding up jobs in the trace using the following crude analysis: replace the progress rate of every task of a phase that is slower than the median task with the median task's rate. If this were to happen, the average completion time of jobs improves by 47%, 29% and 36% in the Facebook, Bing and Yahoo! traces, respectively; small jobs ($\leq$ 10 tasks) improve by 49%, 38% and 41%.
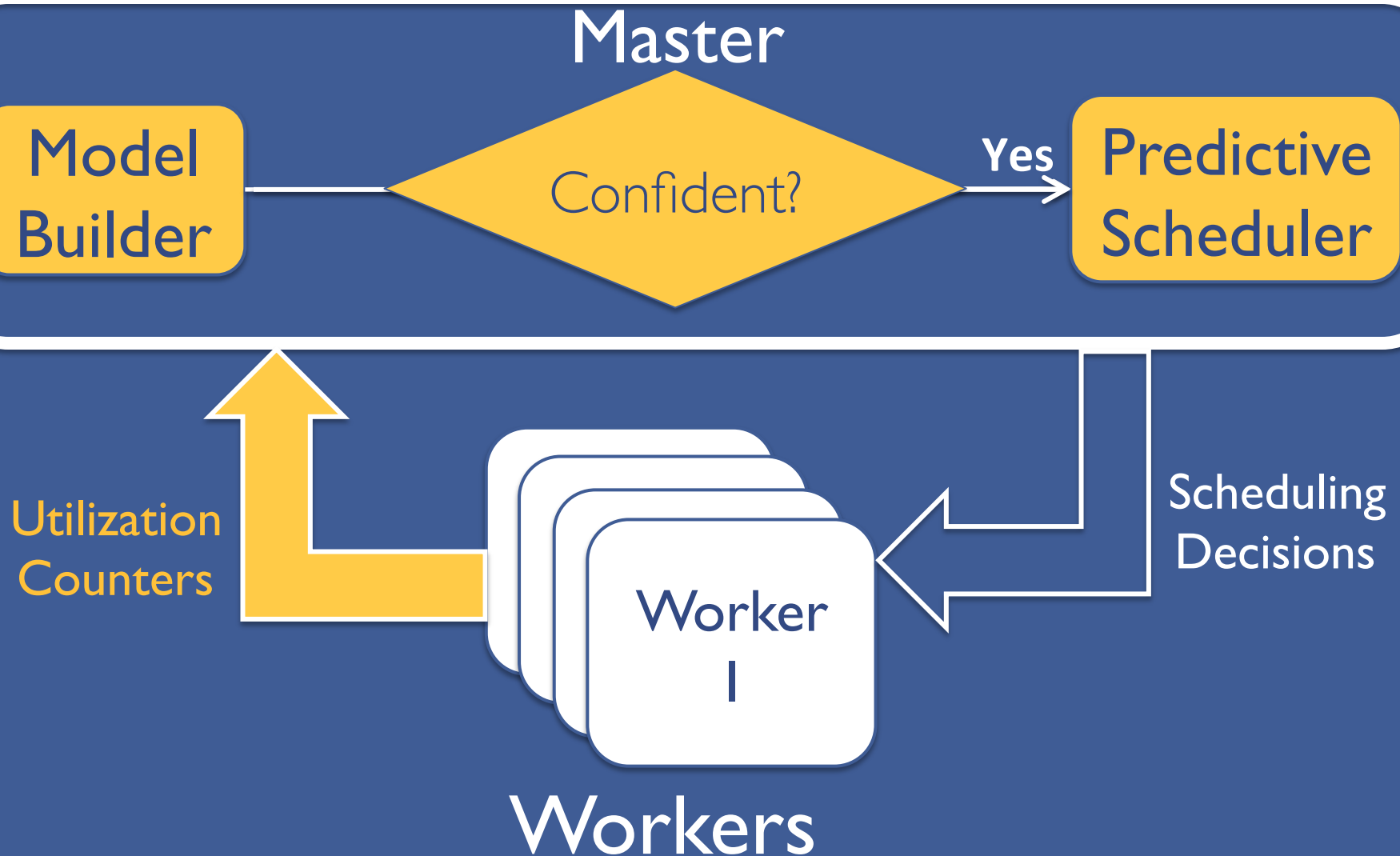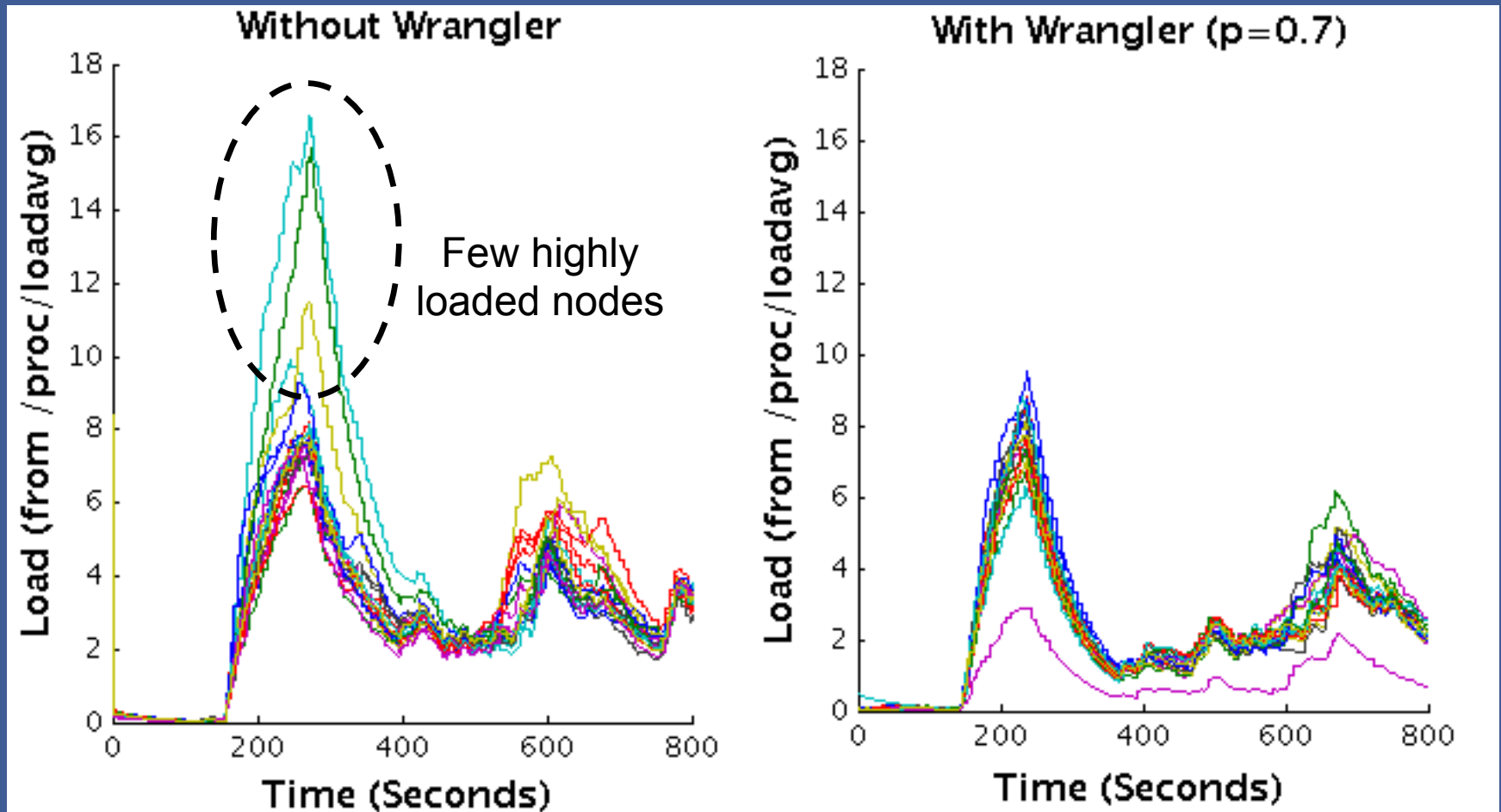
Dolly, NSDI'13

# Speculative Execution



Job queue

Replicating   Job completed

$T_S$

$T_S$

$T_S$

Wasted Resources

Wasted Time in detecting stragglers

$T_S$: In progress

# Existing Approaches



Wasted Resources

Replicate

Speculative Execution
OSDI'04

LATE
OSDI'08

Mantri
OSDI'10

Dolly
NSDI'13

Wrangler

Wait

Wasted Time
in detecting stragglers

# Our proposal: Wrangler [SoCC'14]

# Load-Balancing with Wrangler



Workload: FB2010

# Wrangler Improves Job Completions



Workload: CC_b

Percentage Reduction

- Prediction without confidence measure
- Prediction with confidence measure (p=0.8)

43.59
58.87
62.10
43.13
22.46

-4.22
-5.53
-5.39
-5.51

avg   95p   97p   99p   99.9p

But, we built a model for every node!

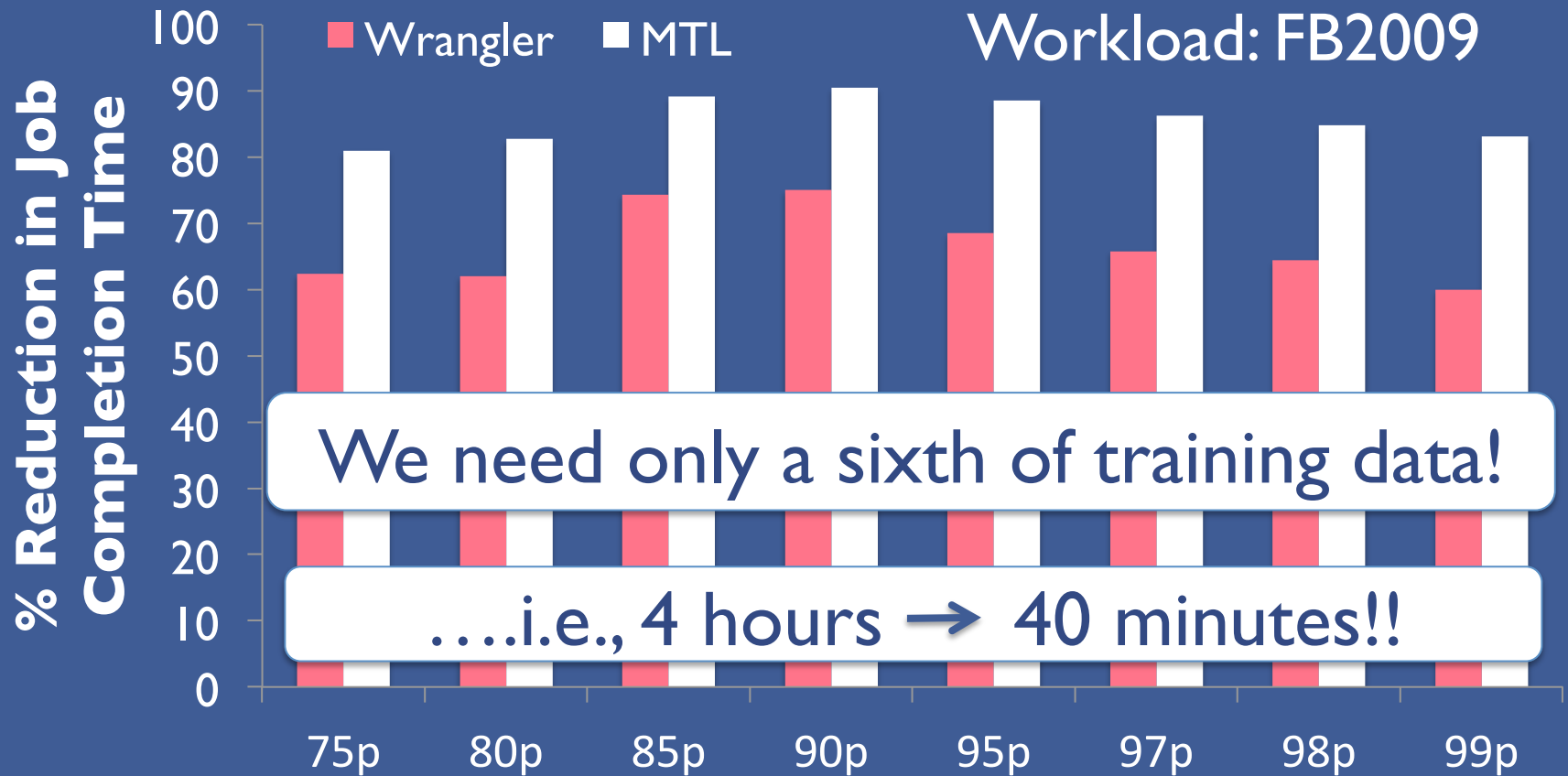Baseline: Speculative Execution

# However….

Real-world production clusters could contain over 1000 nodes

- Scalability!
  - Need to train too many models separately
  - Prohibitively long training data capture duration

## Idea

Share data across nodes and workloads: Multi Task Learning [SDM'15]

# Wrangler: Predictable and Faster Jobs in Distributed Processing Systems using Machine Learning

Neeraja J. Yadwadkar (neerajay@eecs.berkeley.edu),
Bharath Hariharan,
Ganesh Ananthanarayan,
Joseph Gonzalez, and Randy Katz

http://www.istc-cc.cmu.edu/