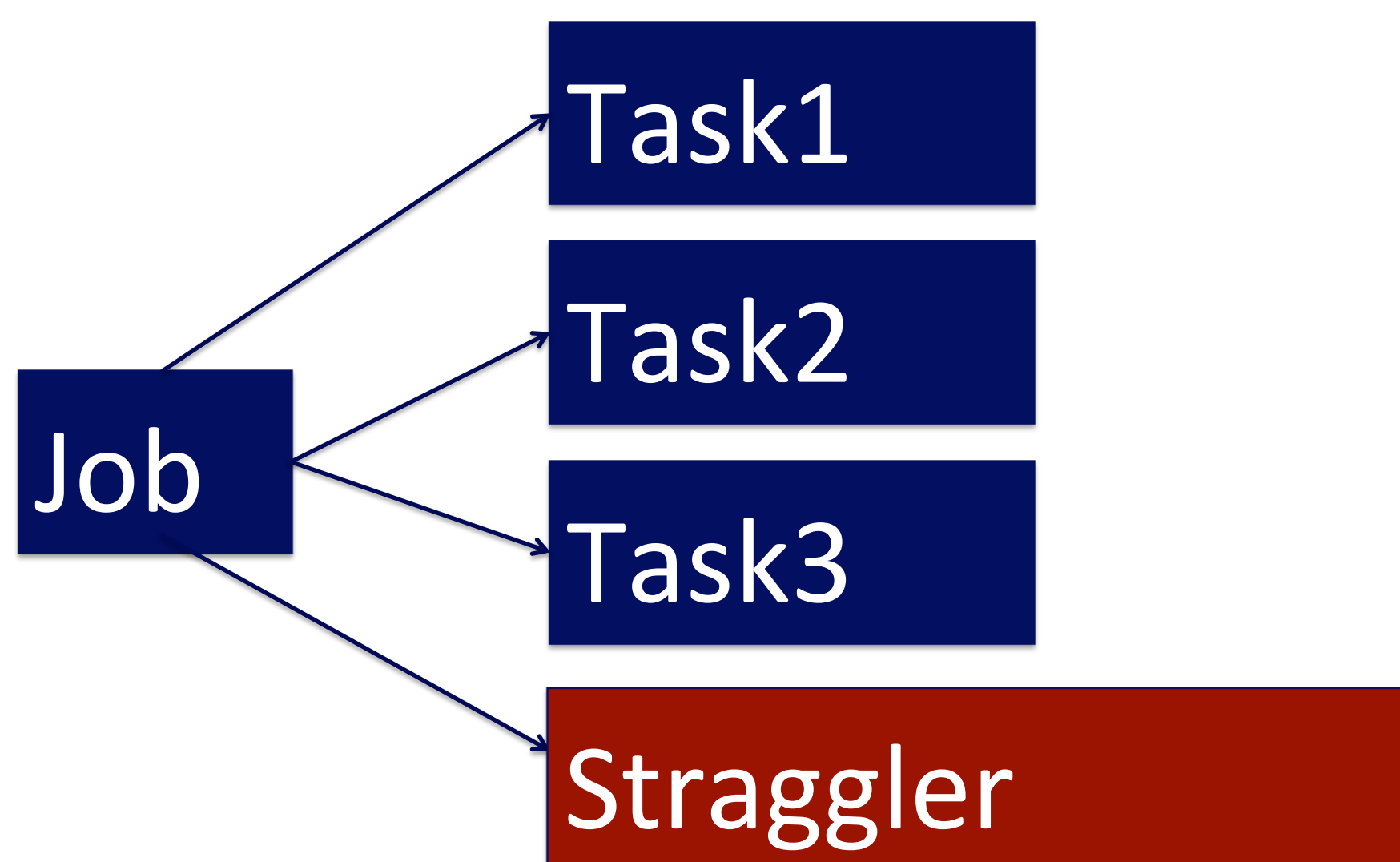# Wrangler: Predictable and Faster Jobs in Distributed Processing Systems using Machine Learning
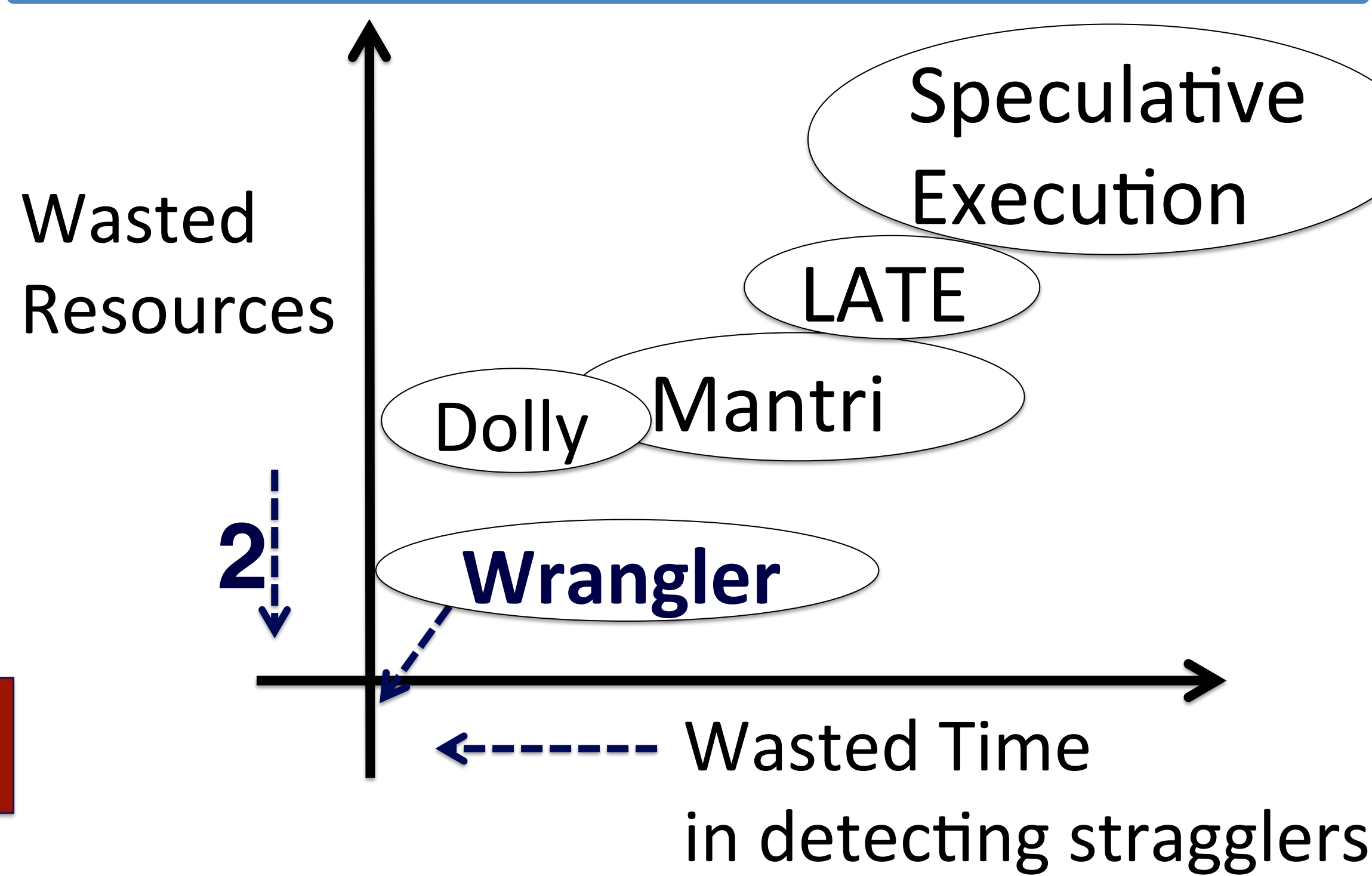
Neeraja J. Yadwadkar, Bharath Hariharan, Ganesh Ananthanarayan, Joseph Gonzalez, and Randy Katz
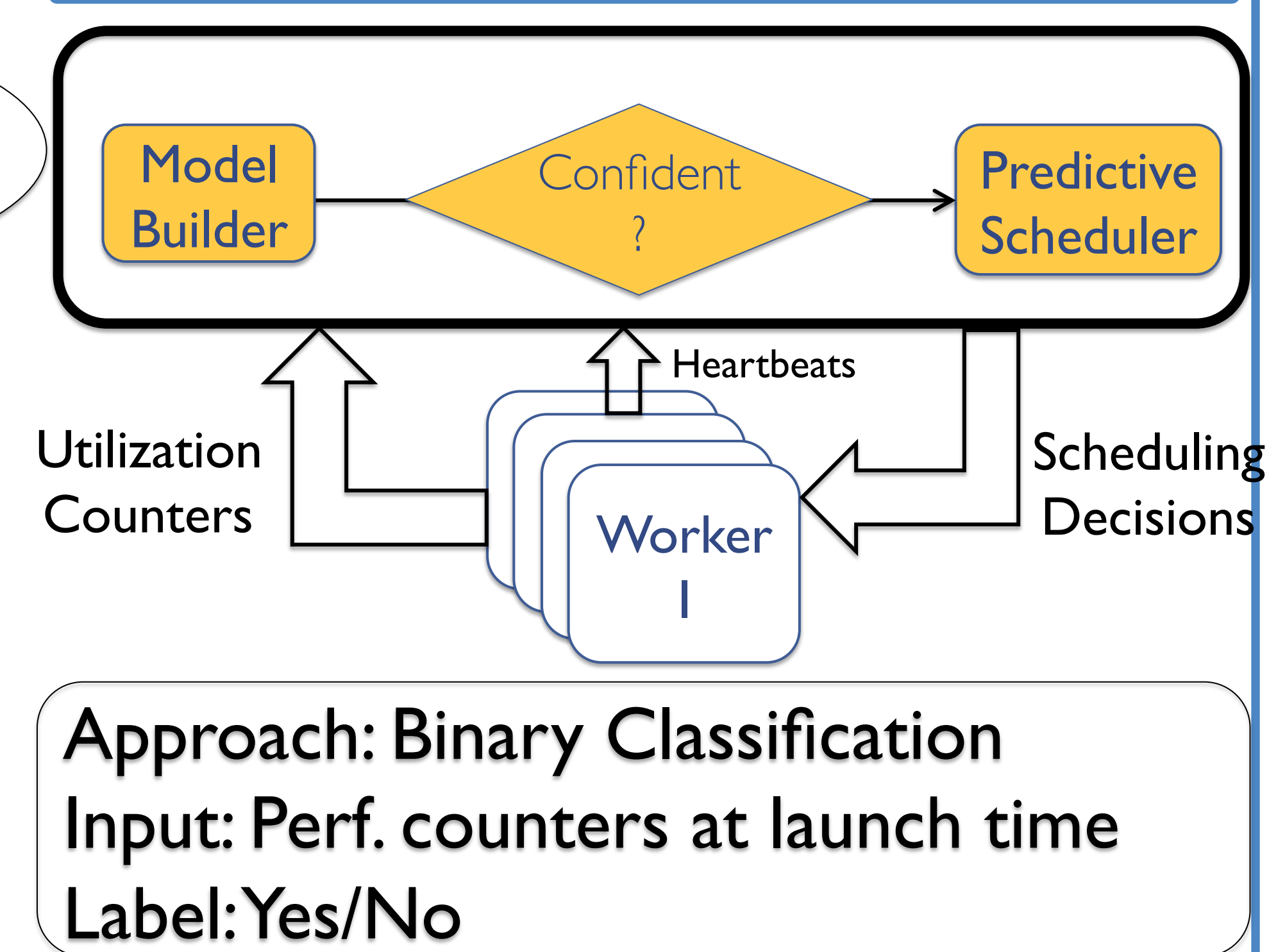
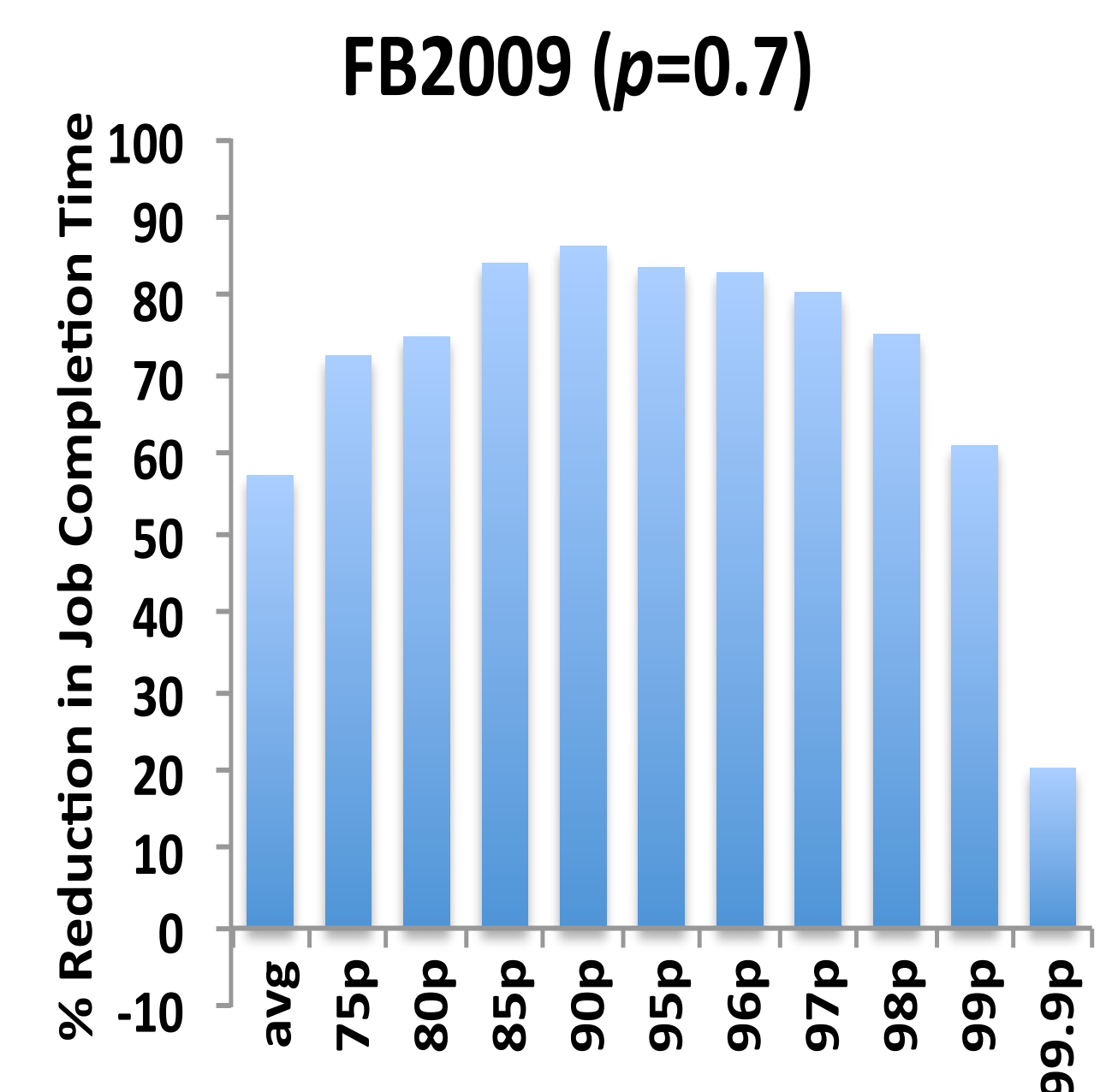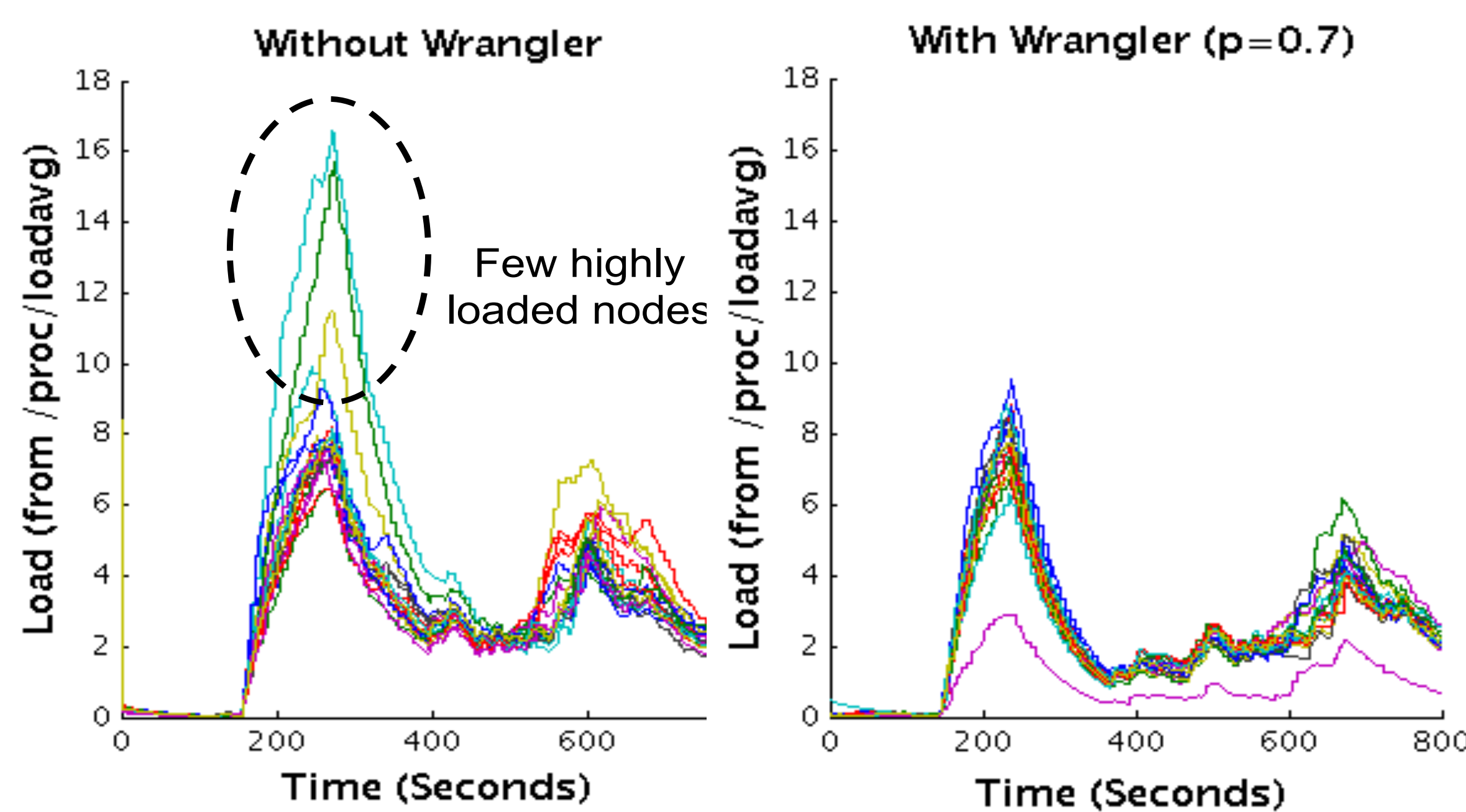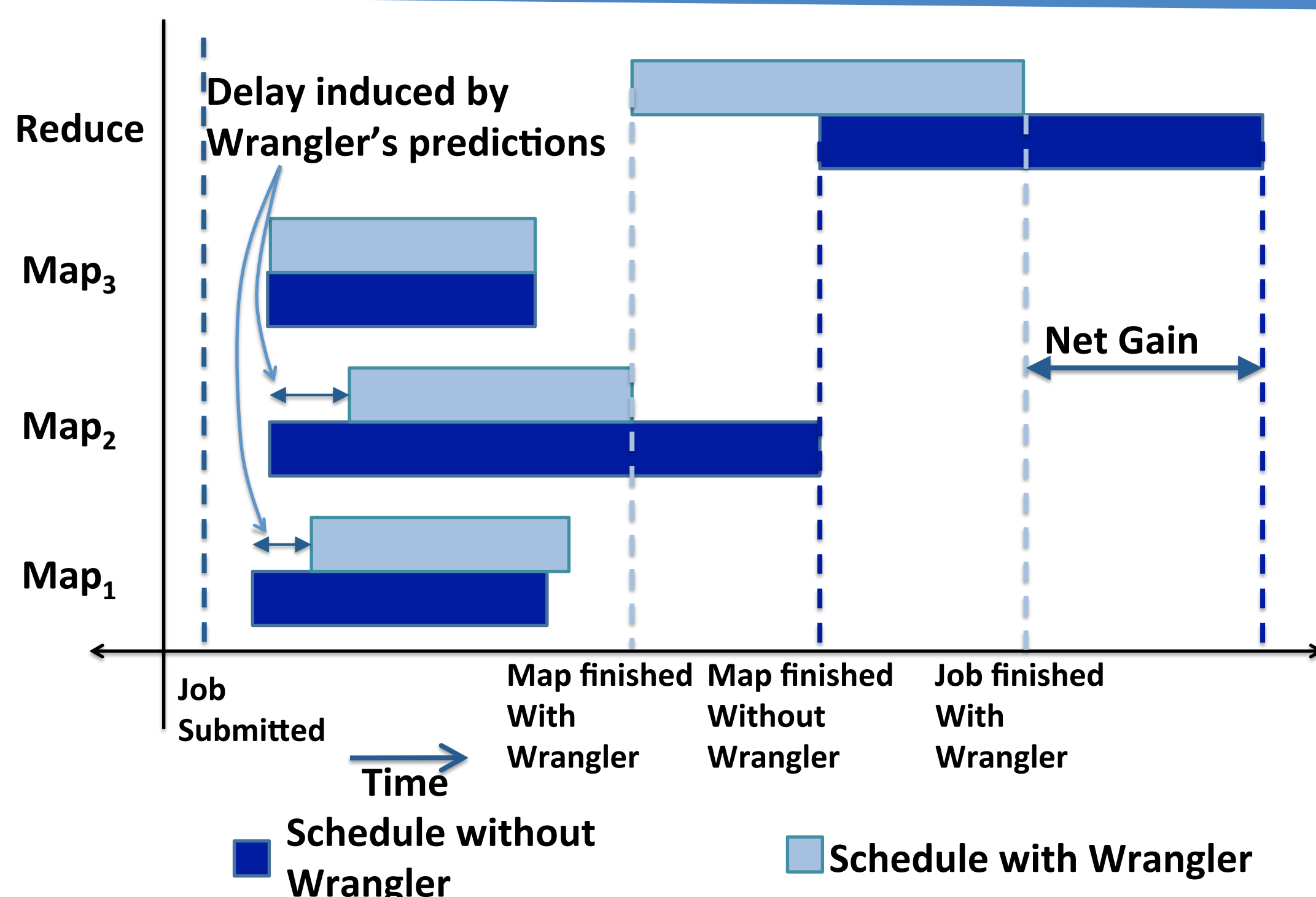## Parallel Data analytics and stragglers



## Design Space



- Wasted Resources
- Speculative Execution
- LATE
- Dolly
- Mantri
- **Wrangler**
- 2
- Wasted Time in detecting stragglers

## Wrangler: Architecture



Model Builder — Confident? — Predictive Scheduler

Utilization Counters — Heartbeats — Scheduling Decisions — Worker

Approach: Binary Classification
Input: Perf. counters at launch time
Label: Yes/No

## Intuition



## Load-Balancing



Without Wrangler — Few highly loaded nodes

With Wrangler (p=0.7)

## Faster Job Completion



FB2009 (*p*=0.7)

## Model Builder



Node 1   Node 2   Node 3
FB2009
FB2010
CC_e

**Prediction Accuracy: 70-80%**

**Scalability!**
Train too many models separately
Why? Heterogeneity across nodes and tasks
Prohibitively long data capture time

## Proposal

- Underlying modeling task remains the same
- Learning from other similar tasks
  - Reduce training data capture time
  - Improve accuracy by generalizing better

Regularized MTL [KDD'04]:

$$\mathbf{w}_t = \mathbf{w}_0 + \mathbf{v}_t$$
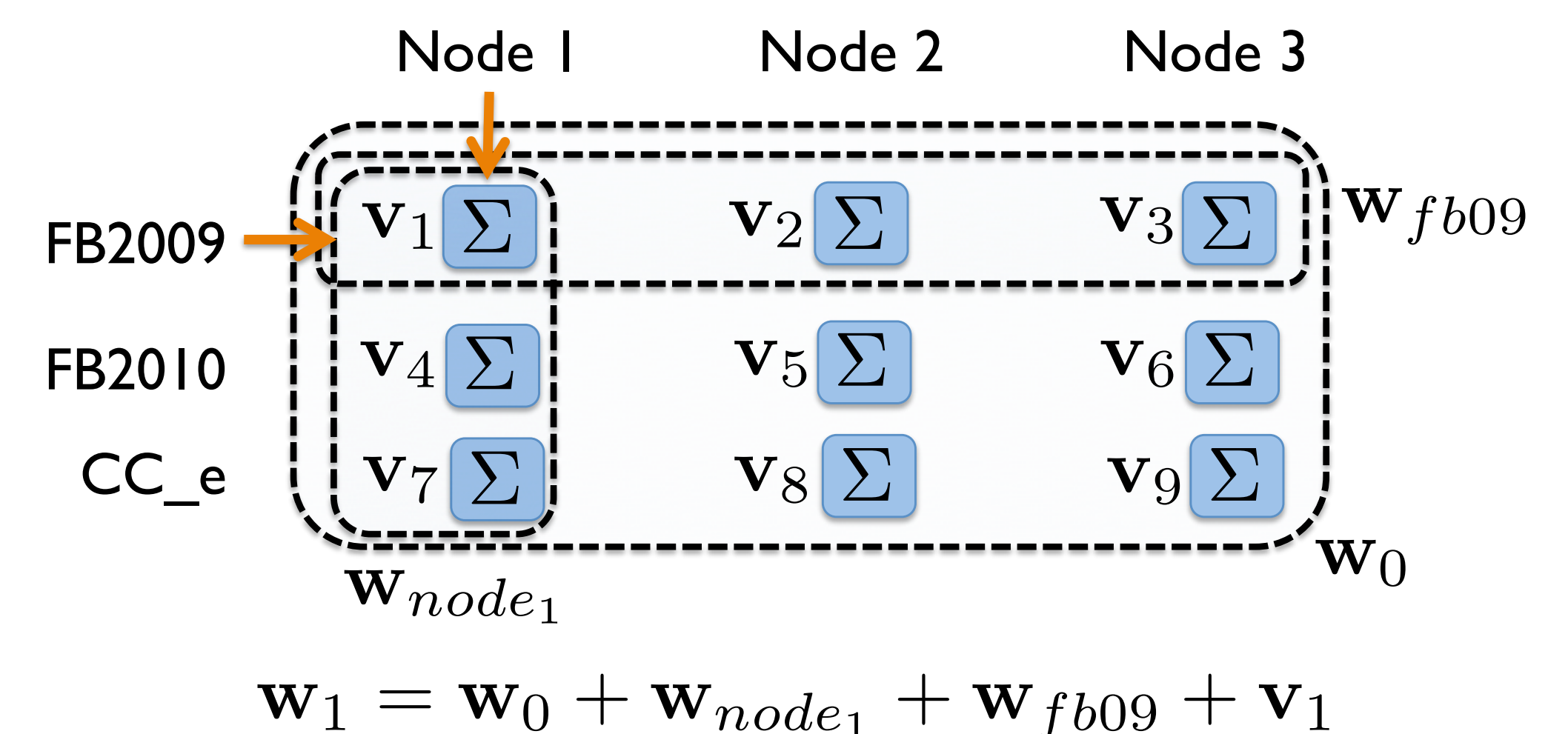
Common across all the learning tasks

Specific for a learning tasks, $t$

Our Formulation:

$$\mathbf{w}_t = \mathbf{w}_0 + \mathbf{v}_t + \mathbf{w}_{g(t)}$$

Common across the tasks in a group, denoted by $g$

Share data across nodes and workloads:
Multi Task Learning



Node 1   Node 2   Node 3

FB2009  $\mathbf{v}_1\Sigma$   $\mathbf{v}_2\Sigma$   $\mathbf{v}_3\Sigma$  $\mathbf{w}_{fb09}$
FB2010  $\mathbf{v}_4\Sigma$   $\mathbf{v}_5\Sigma$   $\mathbf{v}_6\Sigma$
CC_e    $\mathbf{v}_7\Sigma$   $\mathbf{v}_8\Sigma$   $\mathbf{v}_9\Sigma$

$\mathbf{w}_{node_1}$   $\mathbf{w}_0$

$$\mathbf{w}_1 = \mathbf{w}_0 + \mathbf{w}_{node_1} + \mathbf{w}_{fb09} + \mathbf{v}_1$$

## Training Problem

$$\min_{\mathbf{w}_0, \mathbf{v}_t, b} \lambda_0 \|\mathbf{w}_0\|^2 + \frac{\lambda_1}{T} \sum_{t=1}^{T} \|\mathbf{v}_t\|^2 + \sum_{t=1}^{T} \sum_{i=1}^{m_t} \xi_{it}$$

$$\text{s.t} \quad y_{it}\left( (\mathbf{w}_0 + \mathbf{v}_t)^T \mathbf{x}_{it} + b \right) \geq 1 - \xi_{it} \quad \forall i, t$$

$$\xi_{it} \geq 0 \quad \forall i, t$$

## Evaluation



Wrangler — Our formulation — Workload: FB2009



Wrangler — MTL