# Distributed Asynchronous Online Learning for Natural Language Processing

Kevin Gimpel     Dipanjan Das     Noah A. Smith

# Introduction

- Two recent lines of research in speeding up large learning problems:
  - ☐ Parallel/distributed computing
  - ☐ Online (and mini-batch) learning algorithms: stochastic gradient descent, perceptron, MIRA, stepwise EM

- How can we bring together the benefits of parallel computing and online learning?

# Introduction

- We use **asynchronous** algorithms
  (Nedic, Bertsekas, and Borkar, 2001;
  Langford, Smola, and Zinkevich, 2009)

- We apply them to structured prediction tasks:
  - Supervised learning
  - Unsupervised learning with both convex and non-convex objectives

- Asynchronous learning speeds convergence and works best with small mini-batches

# Problem Setting

- **Iterative learning**
  - Moderate to large numbers of training examples
  - Expensive inference procedures for each example
  - For concreteness, we start with gradient-based optimization

- **Single machine with multiple processors**
  - Exploit shared memory for parameters, lexicons, feature caches, etc.
  - Maintain one master copy of model parameters

# Single-Processor Batch Learning
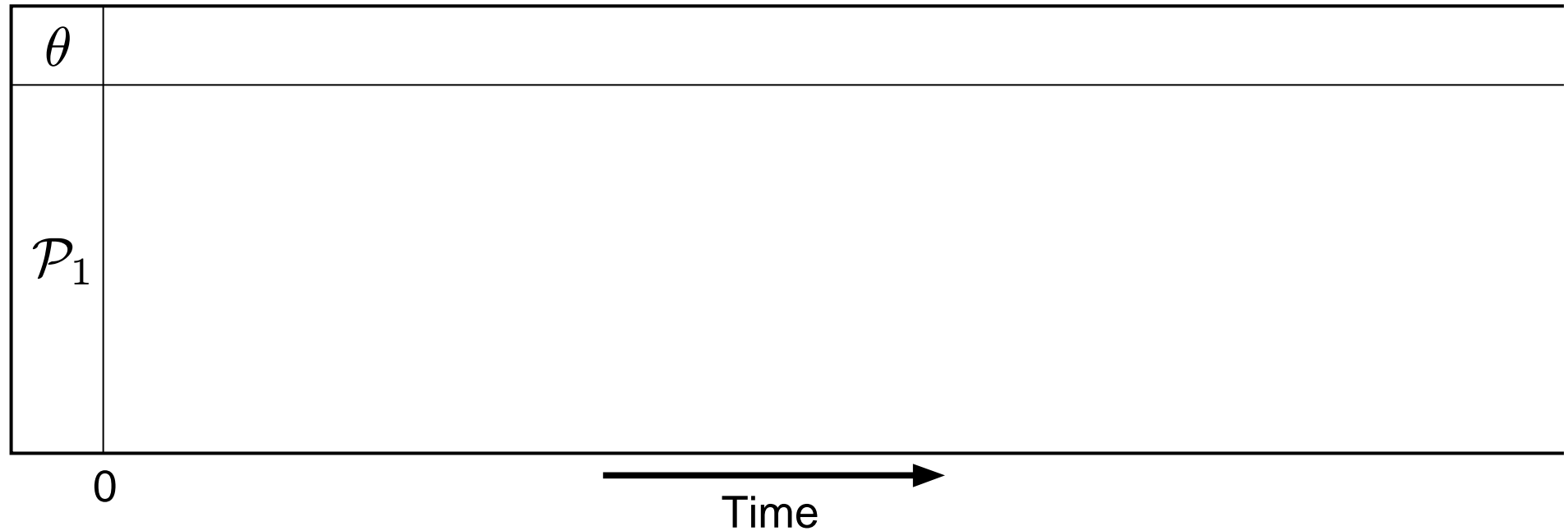
| | |
|---|---|
| Parameters: | $\theta_t$ |
| Processors: | $\mathcal{P}_i$ |
| Dataset: | $\mathcal{D}$ |

# Single-Processor Batch Learning

# Single-Processor Batch Learning

# Single-Processor Batch Learning

| $\theta$ | $\theta_0$ | | | $\theta_1$ |
|---|---|---|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}=\mathrm{calc}(\mathcal{D},\theta_0)$ | | | $\theta_1=\mathrm{up}(\theta_0,\mathbf{g})$ |

0

Time

$\boxed{\mathbf{g}=\mathrm{calc}(\mathcal{D},\theta)}$ :

   Calculate gradient $\mathbf{g}$ on data $\mathcal{D}$ using parameters $\theta$

$\boxed{\theta_1=\mathrm{up}(\theta_0,\mathbf{g})}$ :

   Update $\theta_0$ using gradient $\mathbf{g}$ to obtain $\theta_1$
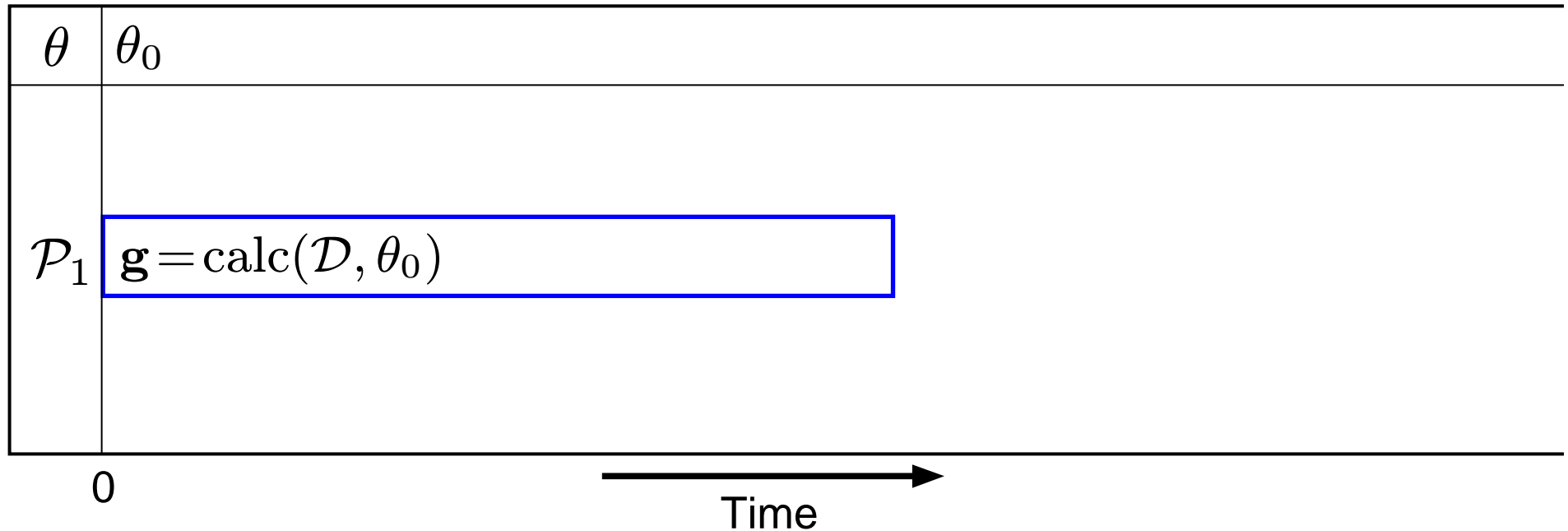
Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Dataset: $\mathcal{D}$

# Single-Processor Batch Learning



| $\theta$ | $\theta_0$ | | $\theta_1$ |
|---|---|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}=\mathrm{calc}(\mathcal{D},\theta_0)$ | $\theta_1=\mathrm{up}(\theta_0,\mathbf{g})$ | $\mathbf{g}=\mathrm{calc}(\mathcal{D},\theta_1)$ |

0      Time →

$\boxed{\mathbf{g}=\mathrm{calc}(\mathcal{D},\theta)}$ :

Calculate gradient $\mathbf{g}$ on data $\mathcal{D}$ using parameters $\theta$

$\boxed{\theta_1=\mathrm{up}(\theta_0,\mathbf{g})}$ :

Update $\theta_0$ using gradient $\mathbf{g}$ to obtain $\theta_1$

Parameters: $\theta_t$
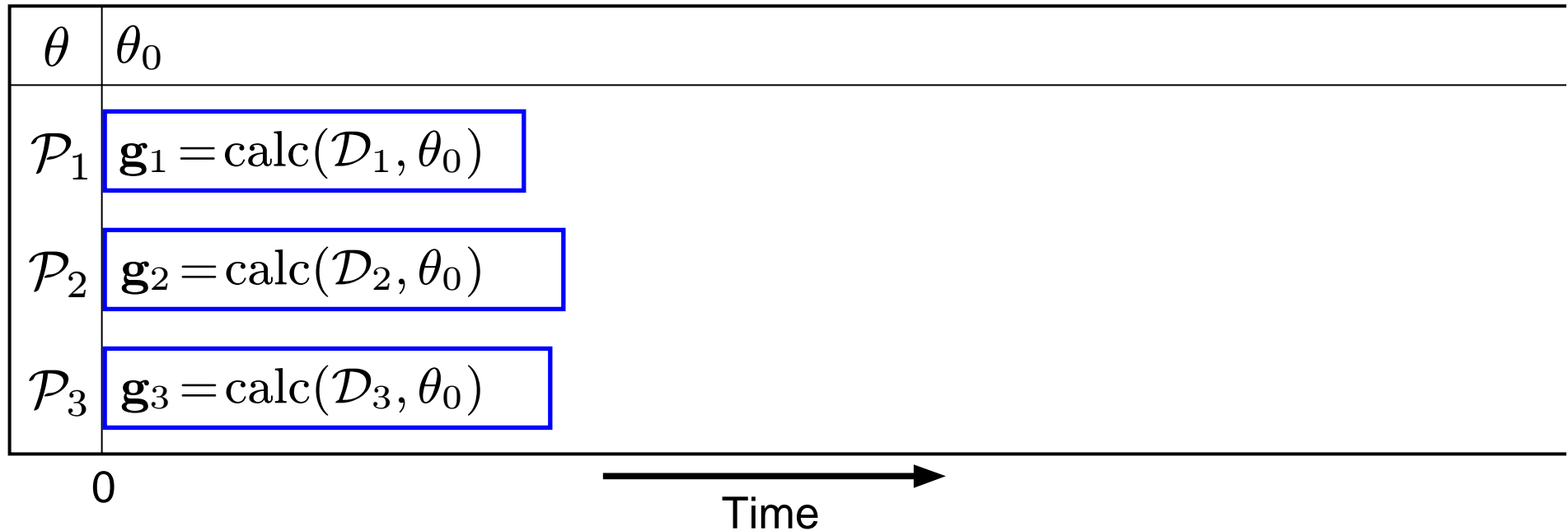
Processors: $\mathcal{P}_i$

Dataset: $\mathcal{D}$

# Parallel Batch Learning

| $\theta$ | $\theta_0$ |
|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = \mathrm{calc}(\mathcal{D}_1, \theta_0)$ |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = \mathrm{calc}(\mathcal{D}_2, \theta_0)$ |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = \mathrm{calc}(\mathcal{D}_3, \theta_0)$ |

0 →

Time

- Divide data into parts, compute gradient on parts in parallel
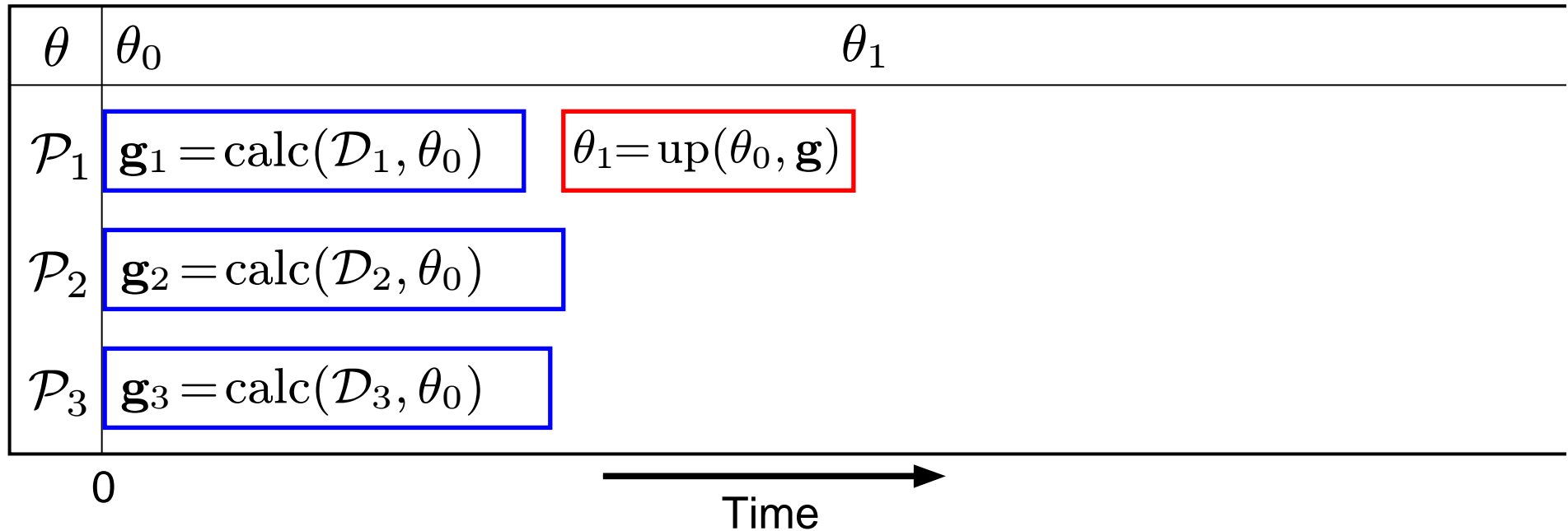
Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Dataset: $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \mathcal{D}_3$

Gradient: $\mathbf{g} = \mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3$

ARK  *lti*  **Carnegie Mellon**

# Parallel Batch Learning

| $\theta$ | $\theta_0$ | $\theta_1$ |
|---|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = \mathrm{calc}(\mathcal{D}_1, \theta_0)$ | $\theta_1 = \mathrm{up}(\theta_0, \mathbf{g})$ |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = \mathrm{calc}(\mathcal{D}_2, \theta_0)$ | |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = \mathrm{calc}(\mathcal{D}_3, \theta_0)$ | |

0

Time

- Divide data into parts, compute gradient on parts in parallel
- One processor updates parameters

Parameters: $\theta_t$
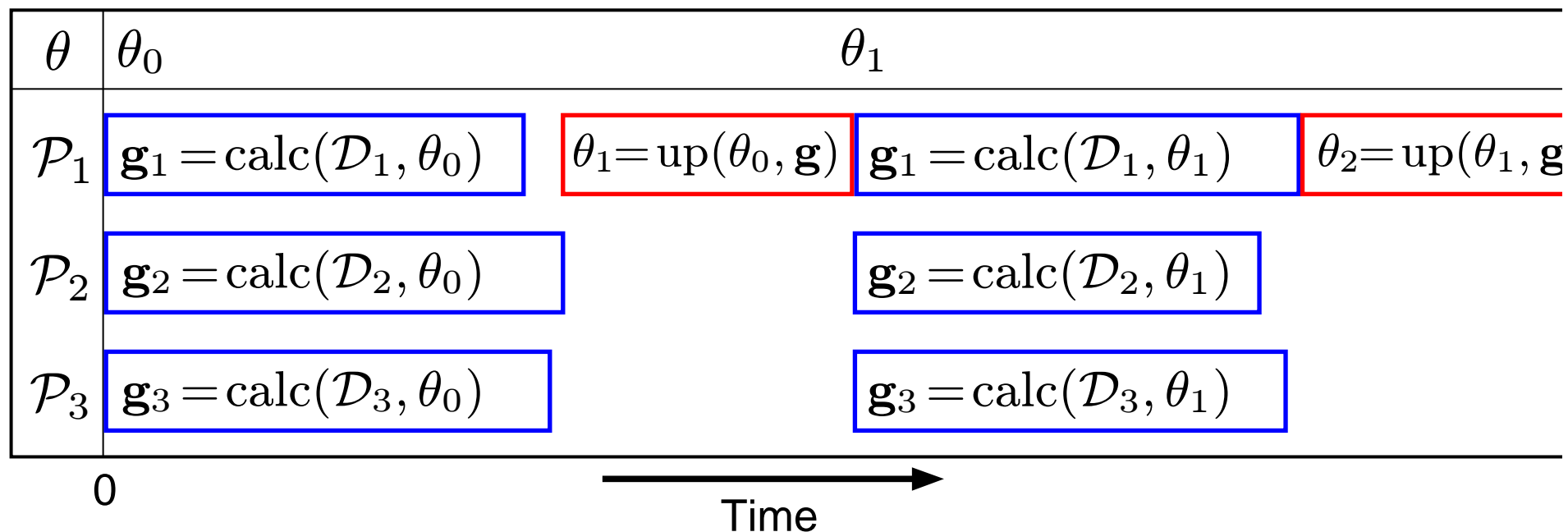
Processors: $\mathcal{P}_i$

Dataset: $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \mathcal{D}_3$

Gradient: $\mathbf{g} = \mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3$

ARK  lti  **Carnegie Mellon**

# Parallel Batch Learning

| $\theta$ | $\theta_0$ | | $\theta_1$ | | |
|---|---|---|---|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = \mathrm{calc}(\mathcal{D}_1, \theta_0)$ | $\theta_1 = \mathrm{up}(\theta_0, \mathbf{g})$ | $\mathbf{g}_1 = \mathrm{calc}(\mathcal{D}_1, \theta_1)$ | $\theta_2 = \mathrm{up}(\theta_1, \mathbf{g}$ | |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = \mathrm{calc}(\mathcal{D}_2, \theta_0)$ | | $\mathbf{g}_2 = \mathrm{calc}(\mathcal{D}_2, \theta_1)$ | | |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = \mathrm{calc}(\mathcal{D}_3, \theta_0)$ | | $\mathbf{g}_3 = \mathrm{calc}(\mathcal{D}_3, \theta_1)$ | | |

0

Time

- Divide data into parts, compute gradient on parts in parallel
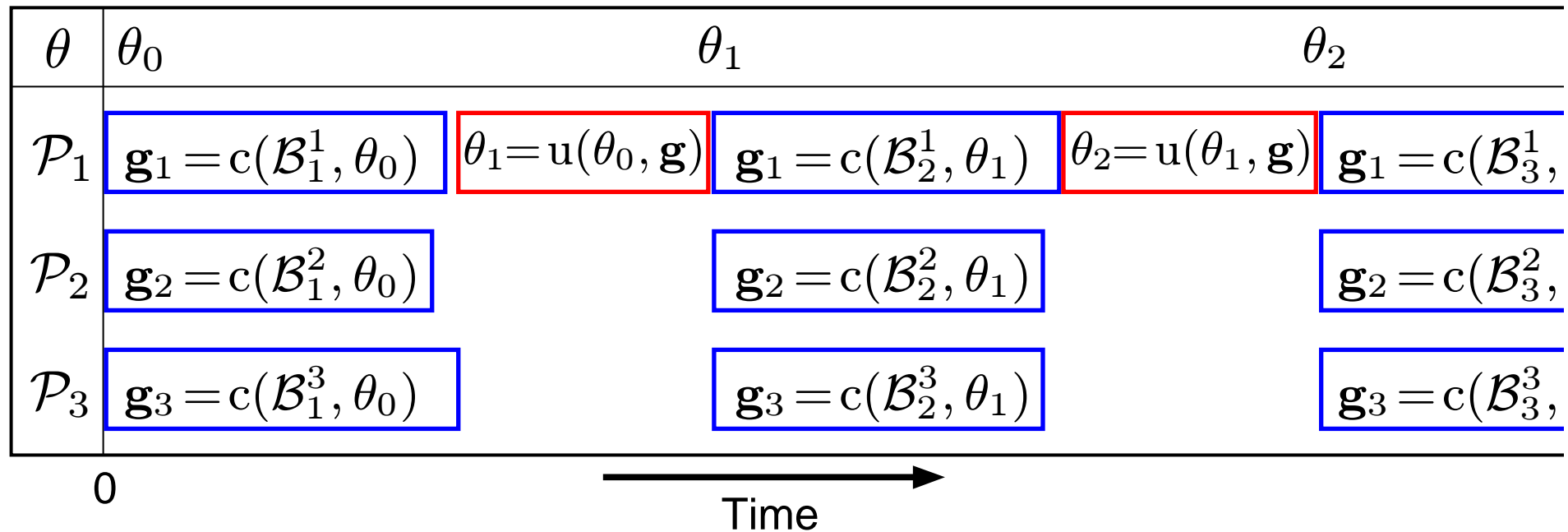- One processor updates parameters

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Dataset: $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \mathcal{D}_3$

Gradient: $\mathbf{g} = \mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3$

ARK lti **Carnegie Mellon**

# Parallel Synchronous Mini-Batch Learning

Finkel, Kleeman, and Manning (2008)

| $\theta$ | $\theta_0$ | $\theta_1$ | | | $\theta_2$ |
|---|---|---|---|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_1^1, \theta_0)$ | $\theta_1 = \mathrm{u}(\theta_0, \mathbf{g})$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_2^1, \theta_1)$ | $\theta_2 = \mathrm{u}(\theta_1, \mathbf{g})$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_3^1,$ |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_1^2, \theta_0)$ | | $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_2^2, \theta_1)$ | | $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_3^2,$ |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_1^3, \theta_0)$ | | $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_2^3, \theta_1)$ | | $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_3^3,$ |

0 →  Time

- Same architecture, just more frequent updates

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_t = \mathcal{B}_t^1 \cup \mathcal{B}_t^2 \cup \mathcal{B}_t^3$
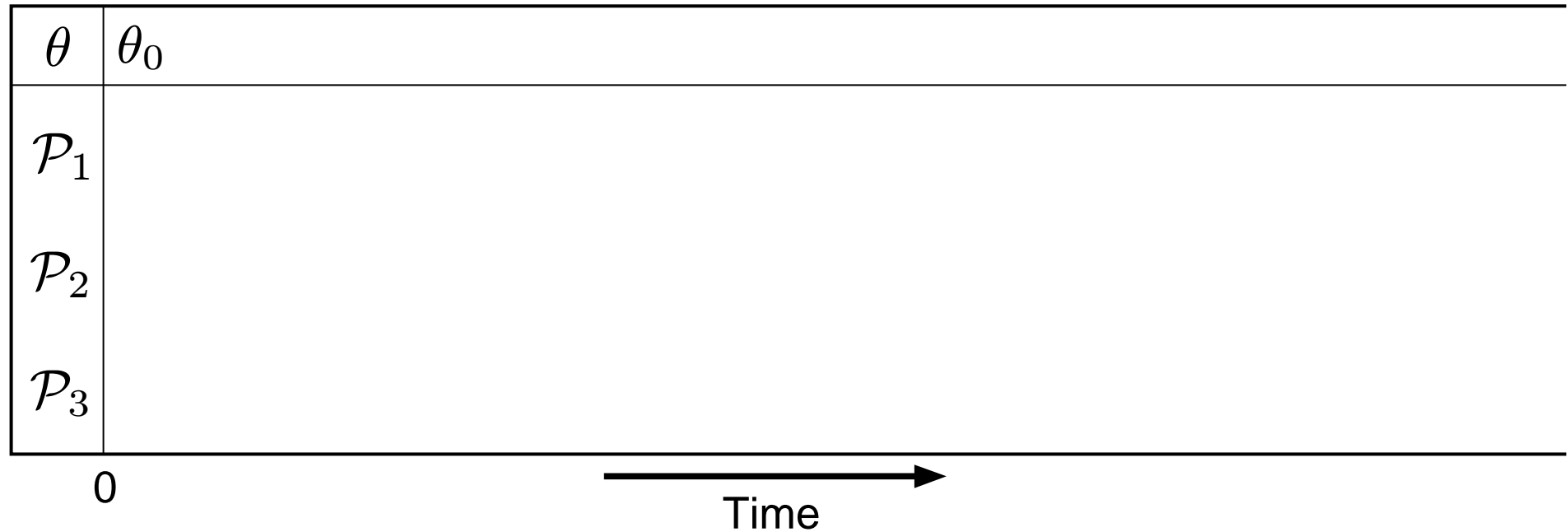
Gradient: $\mathbf{g} = \mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3$

ARK lti **Carnegie Mellon**

# Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)

| $\theta$ | $\theta_0$ |
|---|---|
| $\mathcal{P}_1$ | |
| $\mathcal{P}_2$ | |
| $\mathcal{P}_3$ | |

0

Time
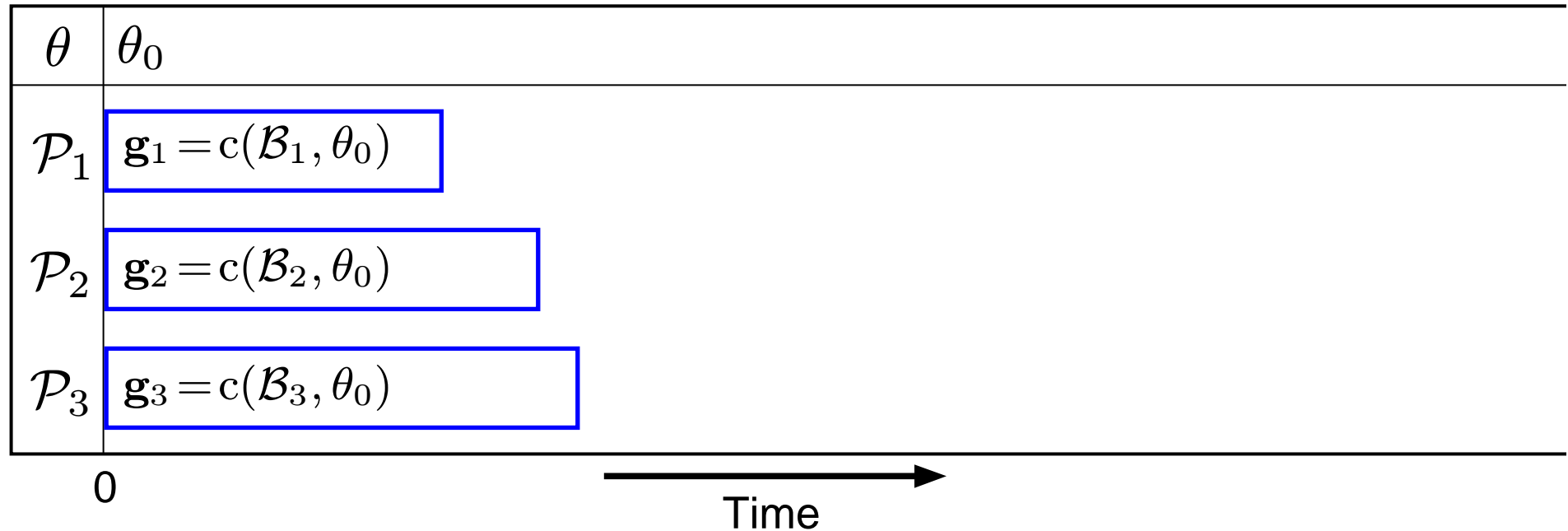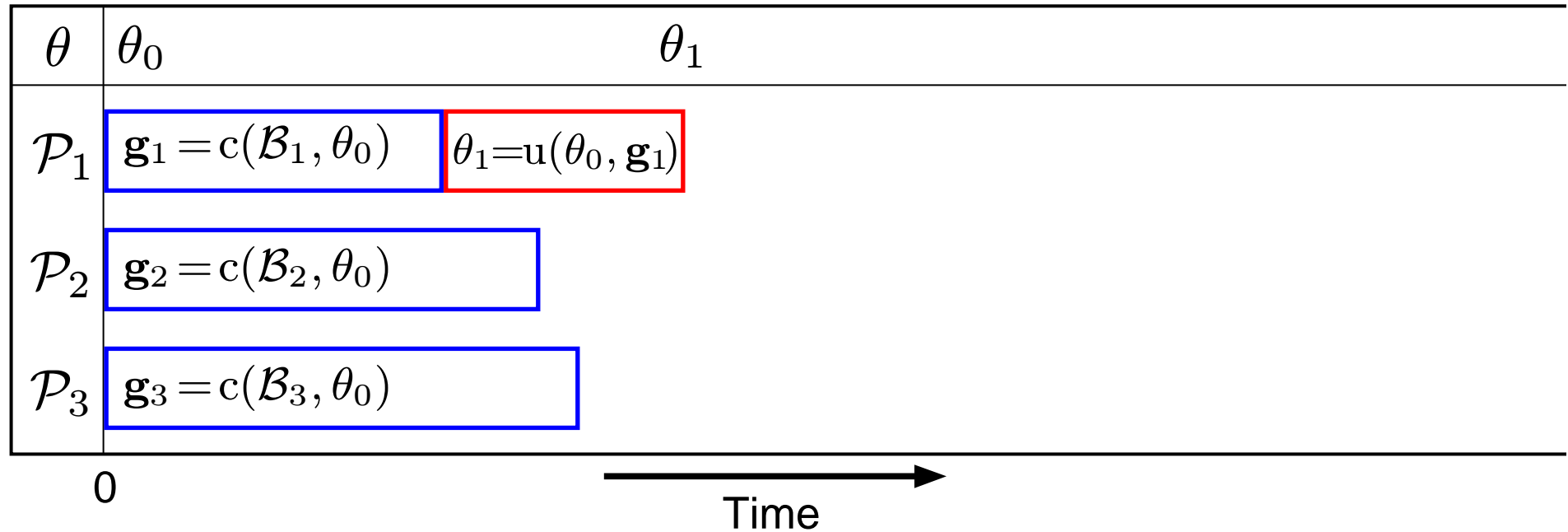
Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_j$

Gradient: $\mathbf{g}_k$

# Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)

| $\theta$ | $\theta_0$ |
|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = c(\mathcal{B}_1, \theta_0)$ |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = c(\mathcal{B}_2, \theta_0)$ |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = c(\mathcal{B}_3, \theta_0)$ |

0

Time

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_j$

Gradient: $\mathbf{g}_k$

ARK lti Carnegie Mellon

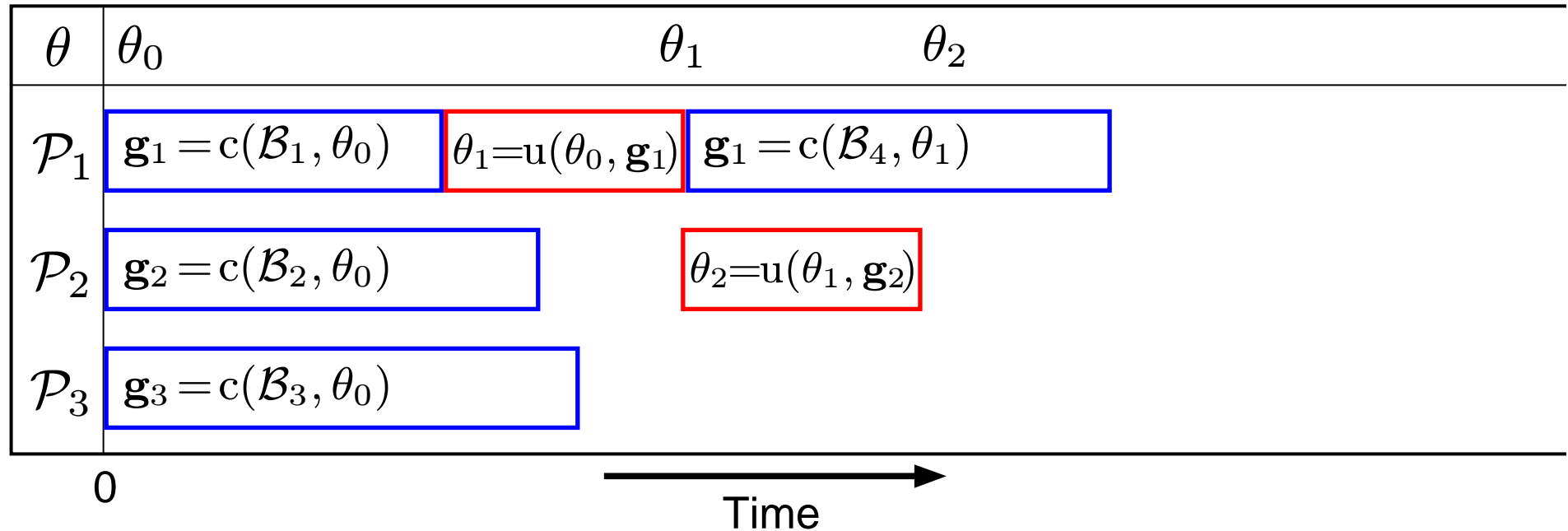# Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)

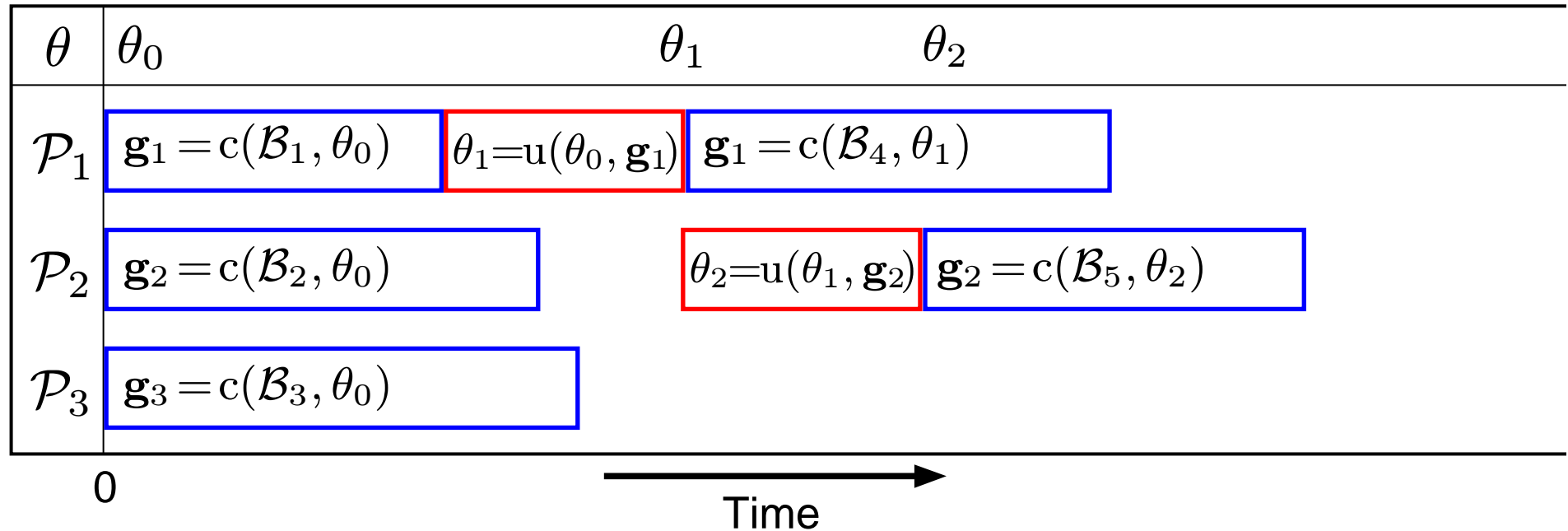| $\theta$ | $\theta_0$ | $\theta_1$ |
|---|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_1, \theta_0)$ | $\theta_1 = \mathrm{u}(\theta_0, \mathbf{g}_1)$ |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_2, \theta_0)$ | |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_3, \theta_0)$ | |

0

Time

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_j$

Gradient: $\mathbf{g}_k$

ARK  lti  **Carnegie Mellon**

# Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)

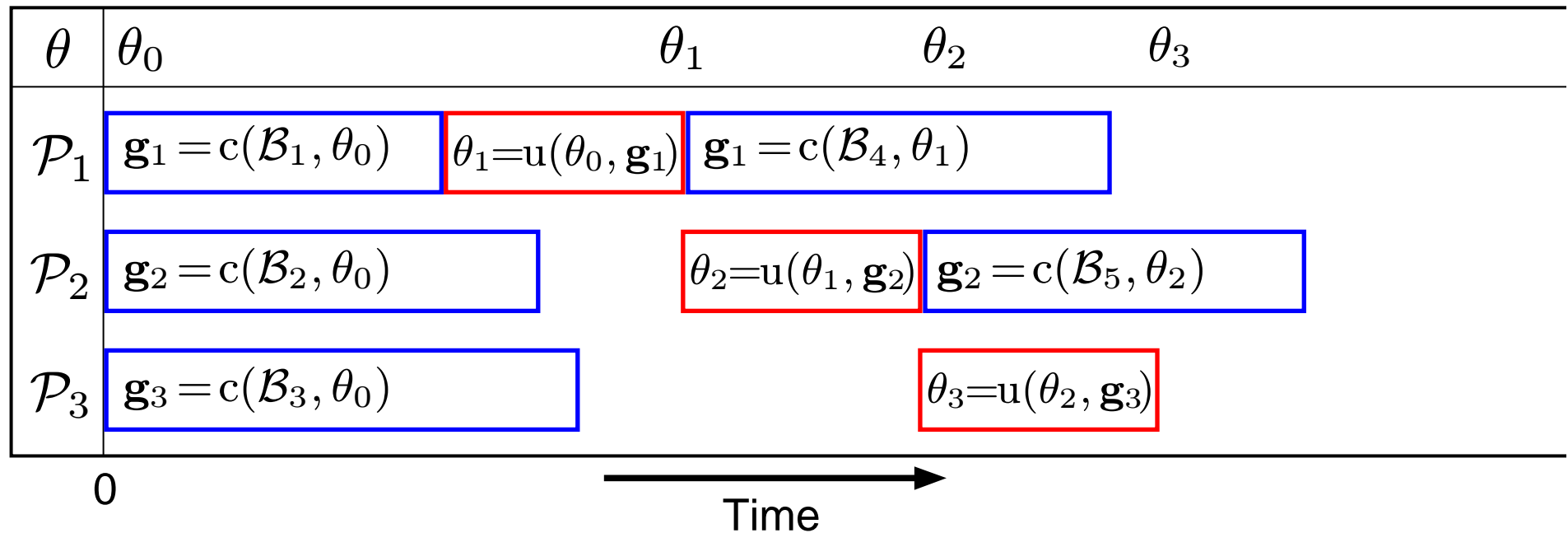| $\theta$ | $\theta_0$ | $\theta_1$ | |
|---|---|---|---|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_1, \theta_0)$ | $\theta_1 = \mathrm{u}(\theta_0, \mathbf{g}_1)$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_4, \theta_1)$ |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_2, \theta_0)$ | | |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_3, \theta_0)$ | | |

0

Time

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_j$

Gradient: $\mathbf{g}_k$

ARK lti **Carnegie Mellon**

# Parallel Asynchronous Mini-Batch Learning

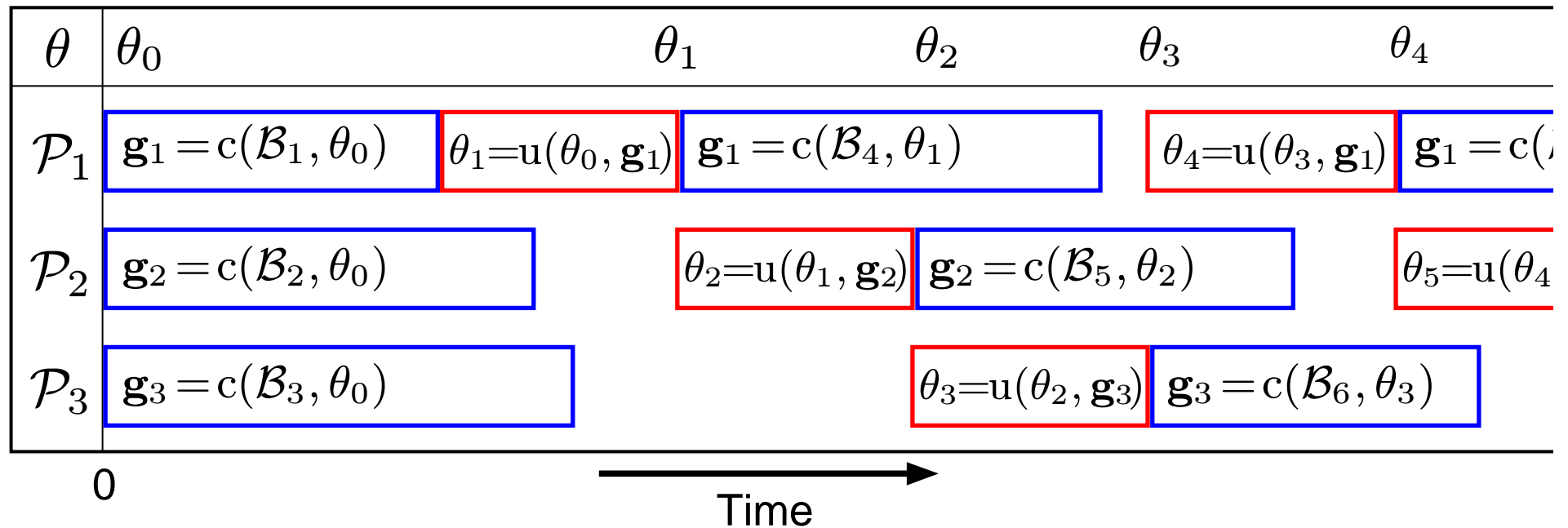| $\theta$ | $\theta_0$ | | $\theta_1$ | $\theta_2$ |
|----------|------------|---|------------|------------|
| $\mathcal{P}_1$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_1, \theta_0)$ | $\theta_1 = \mathrm{u}(\theta_0, \mathbf{g}_1)$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_4, \theta_1)$ | |
| $\mathcal{P}_2$ | $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_2, \theta_0)$ | | $\theta_2 = \mathrm{u}(\theta_1, \mathbf{g}_2)$ | |
| $\mathcal{P}_3$ | $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_3, \theta_0)$ | | | |

0

Time

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_j$

Gradient: $\mathbf{g}_k$

ARK  lti  **Carnegie Mellon**

# Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



| $\theta$ | $\theta_0$ | | $\theta_1$ | $\theta_2$ |
|---|---|---|---|---|

$\mathcal{P}_1$: $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_1, \theta_0)$ | $\theta_1 = \mathrm{u}(\theta_0, \mathbf{g}_1)$ | $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_4, \theta_1)$

$\mathcal{P}_2$: $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_2, \theta_0)$ | $\theta_2 = \mathrm{u}(\theta_1, \mathbf{g}_2)$ | $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_5, \theta_2)$

$\mathcal{P}_3$: $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_3, \theta_0)$

0

Time

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_j$

Gradient: $\mathbf{g}_k$

# Parallel Asynchronous Mini-Batch Learning

# Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)

| $\theta$ | $\theta_0$ | | $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ |
|---|---|---|---|---|---|---|

$\mathcal{P}_1$  $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_1, \theta_0)$  $\theta_1 = \mathrm{u}(\theta_0, \mathbf{g}_1)$  $\mathbf{g}_1 = \mathrm{c}(\mathcal{B}_4, \theta_1)$  $\theta_4 = \mathrm{u}(\theta_3, \mathbf{g}_1)$  $\mathbf{g}_1 = \mathrm{c}($

$\mathcal{P}_2$  $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_2, \theta_0)$  $\theta_2 = \mathrm{u}(\theta_1, \mathbf{g}_2)$  $\mathbf{g}_2 = \mathrm{c}(\mathcal{B}_5, \theta_2)$  $\theta_5 = \mathrm{u}(\theta_4$

$\mathcal{P}_3$  $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_3, \theta_0)$  $\theta_3 = \mathrm{u}(\theta_2, \mathbf{g}_3)$  $\mathbf{g}_3 = \mathrm{c}(\mathcal{B}_6, \theta_3)$

0

Time

- Gradients computed using stale parameters
- Increased processor utilization
- Only idle time caused by lock for updating parameters

Parameters: $\theta_t$

Processors: $\mathcal{P}_i$

Mini-batches: $\mathcal{B}_j$

Gradient: $\mathbf{g}_k$

ARK  lti  Carnegie Mellon

# Theoretical Results

- How does the use of stale parameters affect convergence?


- Convergence results exist for convex optimization using stochastic gradient descent
  - Convergence guaranteed when max delay is bounded (Nedic, Bertsekas, and Borkar, 2001)
  - Convergence rates linear in max delay (Langford, Smola, and Zinkevich, 2009)

# Experiments

| Task | Model | Method | Convex? | $\|\mathcal{D}\|$ | $\|\theta\|$ | $m$ |
|---|---|---|---|---|---|---|
| Named-Entity Recognition | CRF | Stochastic Gradient Descent | Y | 15k | 1.3M | 4 |
| Word Alignment | IBM Model 1 | Stepwise EM | Y | 300k | 14.2M | 10k |
| Unsupervised Part-of-Speech Tagging | HMM | Stepwise EM | N | 42k | 2M | 4 |

- To compare algorithms, we use wall clock time (with a dedicated 4-processor machine)
- $m$ = mini-batch size

# Experiments

| Task | Model | Method | Convex? | $|\mathcal{D}|$ | $|\theta|$ | $m$ |
|------|-------|--------|---------|-----------------|------------|-----|
| Named-Entity Recognition | CRF | Stochastic Gradient Descent | Y | 15k | 1.3M | 4 |

- CoNLL 2003 English data

- Label each token with entity type (person, location, organization, or miscellaneous) or non-entity

- We show convergence in F1 on development data

# Asynchronous Updating Speeds Convergence



*All use a mini-batch size of 4*

Legend:
- Asynchronous (4 processors)
- Synchronous (4 processors)
- Single-processor

Comparison with Ideal Speed-up

# Why Does Asynchronous Converge Faster?

- Processors are kept in near-constant use

- Synchronous SGD leads to idle processors → need for load-balancing

*Clearer improvement for asynchronous algorithms when increasing number of processors*

# Artificial Delays



After completing a mini-batch, 25% chance of delaying

Delay (in seconds) sampled from
$$\max(\mathcal{N}(\mu, (\mu/5)^2), 0)$$

Legend:
- Asynchronous, no delay
- Asynchronous, $\mu = 5$
- Single-processor, no delay
- Asynchronous, $\mu = 10$
- Asynchronous, $\mu = 20$

Y-axis: F1

X-axis: Wall clock time (hours)

Avg. time per mini-batch = 0.62 s

# Experiments

| Task | Model | Method | Convex? | $|\mathcal{D}|$ | $|\theta|$ | $m$ |
|---|---|---|---|---|---|---|
| Word Alignment | IBM Model 1 | Stepwise EM | Y | 300k | 14.2M | 10k |

- **Given parallel sentences, draw links between words:**

konnten   sie   es   übersetzen   ?

could   you   translate   it   ?

- We show convergence in log-likelihood (convergence in AER is similar)

# Stepwise EM
(Sato and Ishii, 2000; Cappe and Moulines, 2009)

- Similar to stochastic gradient descent in the space of sufficient statistics, with a particular scaling of the update

- More efficient than incremental EM

  (Neal and Hinton, 1998)

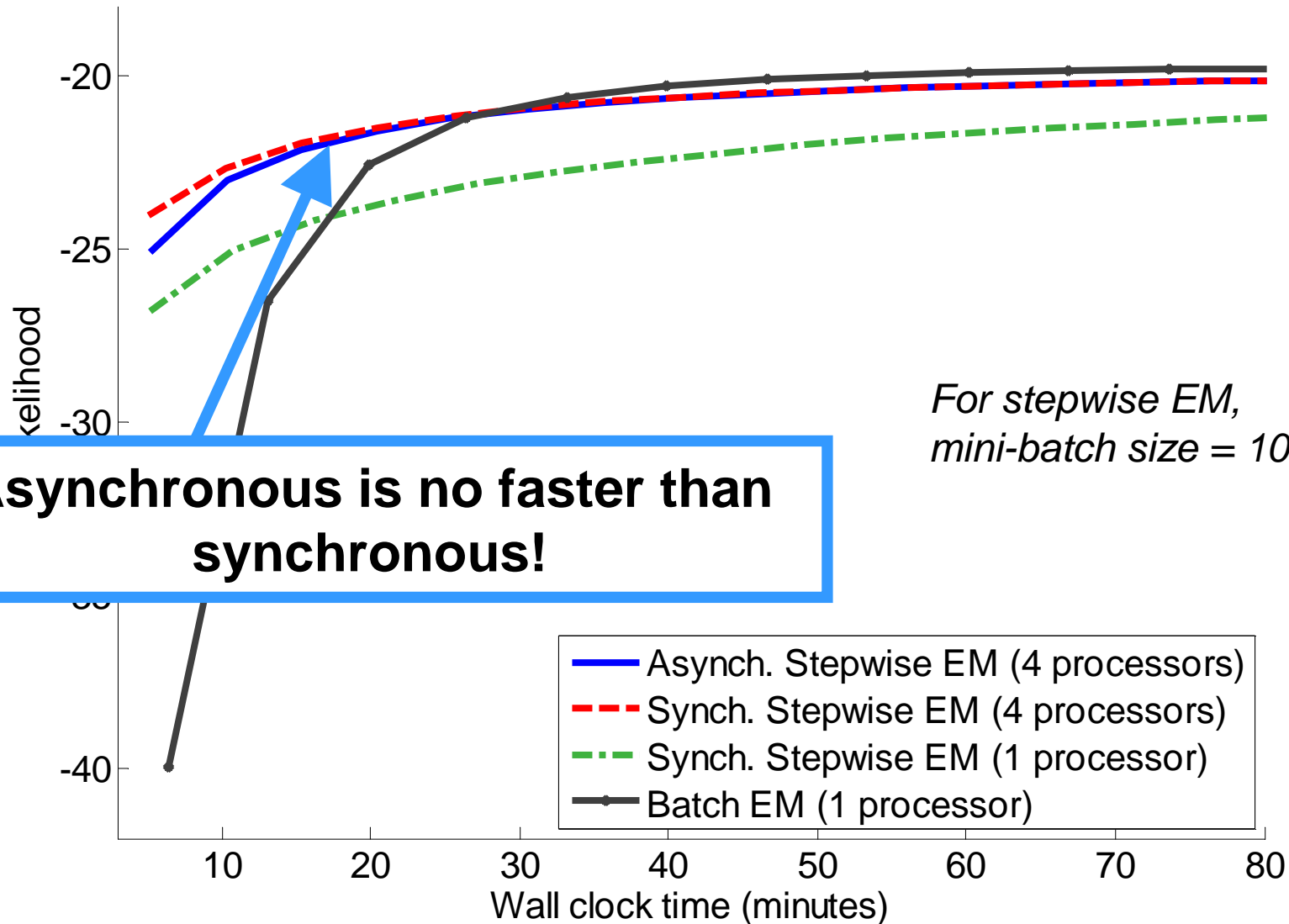- Found to converge much faster than batch EM (Liang and Klein, 2009)
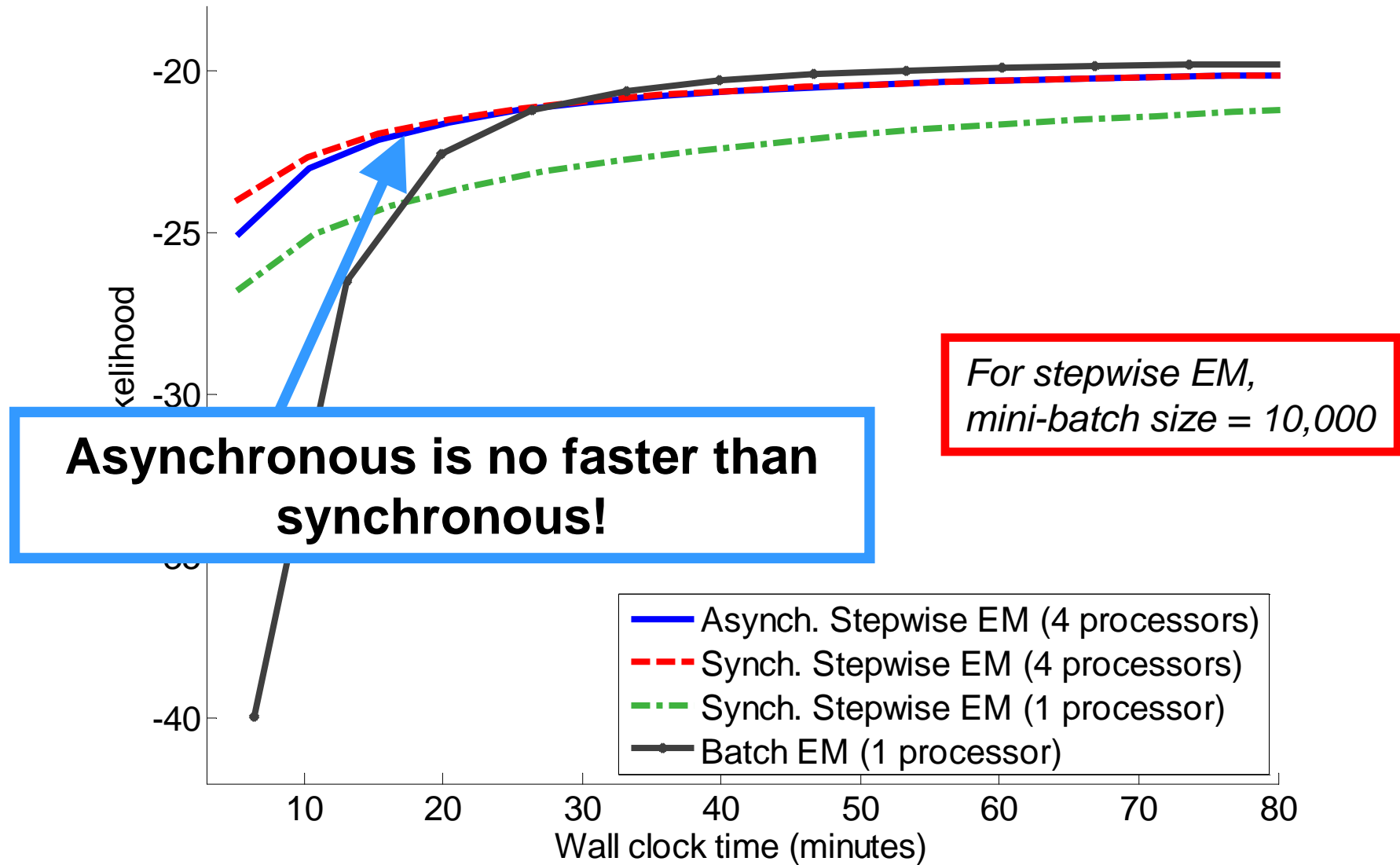
# Word Alignment Results

*For stepwise EM,*
*mini-batch size = 10,000*

Log-Likelihood

Wall clock time (minutes)

- Asynch. Stepwise EM (4 processors)
- Synch. Stepwise EM (4 processors)
- Synch. Stepwise EM (1 processor)
- Batch EM (1 processor)

ARK  *lti*  **Carnegie Mellon**

# Word Alignment Results



*For stepwise EM, mini-batch size = 10,000*

**Asynchronous is no faster than synchronous!**

Likelihood

Wall clock time (minutes)

Asynch. Stepwise EM (4 processors)
Synch. Stepwise EM (4 processors)
Synch. Stepwise EM (1 processor)
Batch EM (1 processor)

# Word Alignment Results



For stepwise EM, mini-batch size = 10,000

**Asynchronous is no faster than synchronous!**

Legend:
- Asynch. Stepwise EM (4 processors)
- Synch. Stepwise EM (4 processors)
- Synch. Stepwise EM (1 processor)
- Batch EM (1 processor)

y-axis: kelihood

x-axis: Wall clock time (minutes)

# Comparing Mini-Batch Sizes

Log-Likelihood vs. Wall clock time (minutes)

Legend:
- Asynch. (m = 10,000)
- Synch. (m = 10,000)
- Asynch. (m = 1,000)
- Synch. (m = 1,000)
- Asynch. (m = 100)
- Synch. (m = 100)

# Comparing Mini-Batch Sizes



Log-Likelihood

Wall clock time (minutes)

**Asynchronous is faster when using small mini-batches**

- Asynch. (m = 10,000)
- Synch. (m = 10,000)
- Asynch. (m = 1,000)
- Synch. (m = 1,000)
- Asynch. (m = 100)
- Synch. (m = 100)

# Comparing Mini-Batch Sizes

Error from asynchronous updating

Log-Likelihood

Wall clock time (minutes)

Asynch. (m = 10,000)
Synch. (m = 10,000)
Asynch. (m = 1,000)
Synch. (m = 1,000)
Asynch. (m = 100)
Synch. (m = 100)

# Word Alignment Results



*For stepwise EM, mini-batch size = 10,000*

Log-Likelihood (y-axis)

Wall clock time (minutes) (x-axis)

Legend:
- Asynch. Stepwise EM (4 processors)
- Synch. Stepwise EM (4 processors)
- Synch. Stepwise EM (1 processor)
- Batch EM (1 processor)

# Comparison with Ideal Speed-up



For stepwise EM,
mini-batch size = 10,000

Log-Likelihood

Wall clock time (minutes)

— Asynch. Stepwise EM (4 processors)
-- Synch. Stepwise EM (4 processors)
-·- Synch. Stepwise EM (1 processor)
—•— Batch EM (1 processor)
-- Ideal

# MapReduce?

- **We also ran these algorithms on a large MapReduce cluster (M45 from Yahoo!)**
- **Batch EM**
  - Each iteration is one MapReduce job, using 24 mappers and 1 reducer
- **Asynchronous Stepwise EM**
  - 4 mini-batches processed simultaneously, each run as a MapReduce job
  - Each uses 6 mappers and 1 reducer

# MapReduce?

# MapReduce?

# Experiments

| Task | Model | Method | Convex? | $|\mathcal{D}|$ | $|\theta|$ | $m$ |
|---|---|---|---|---|---|---|
| Unsupervised Part-of-Speech Tagging | HMM | Stepwise EM | N | 42k | 2M | 4 |

- **Bigram HMM with 45 states**

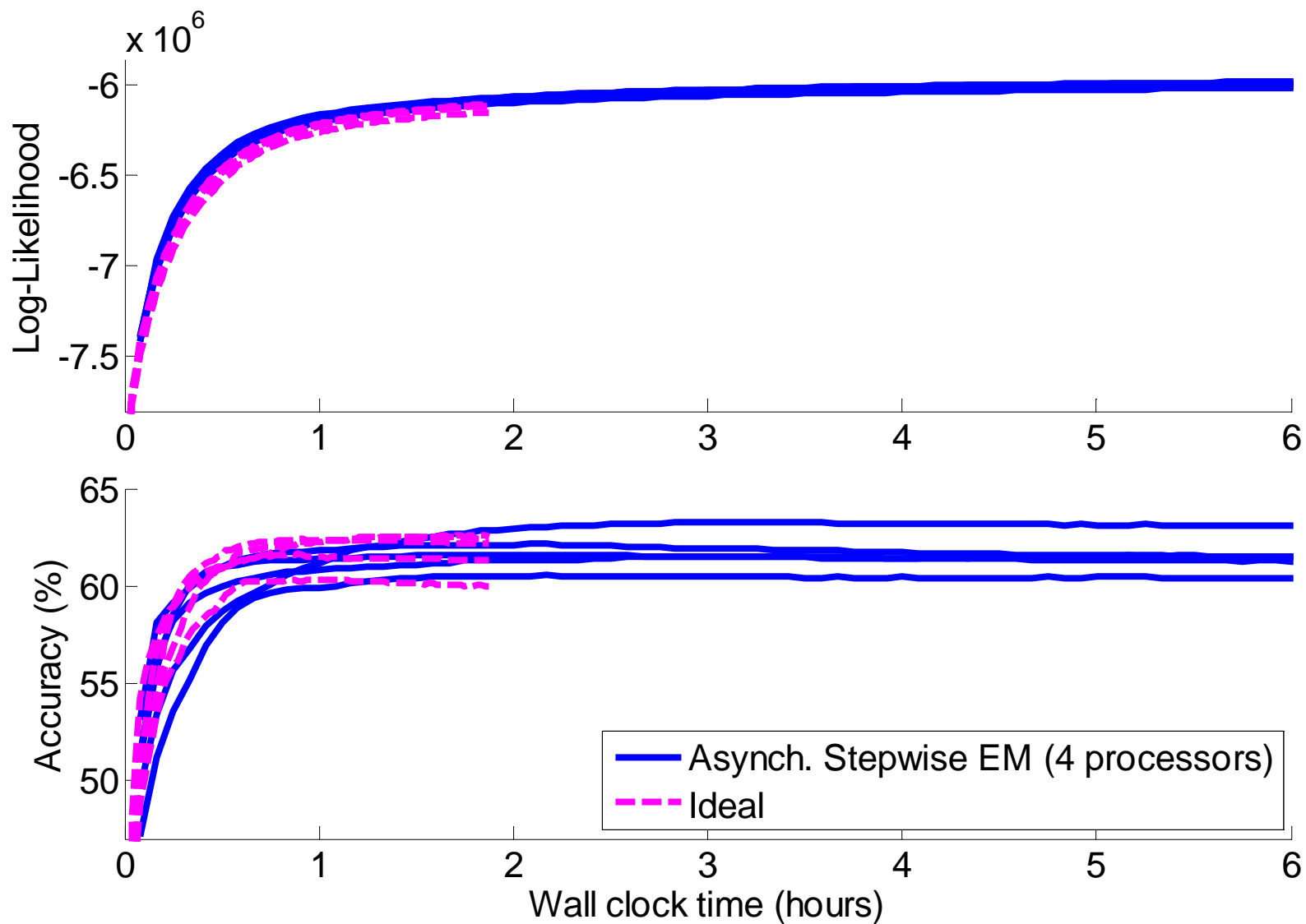- We plot convergence in likelihood and many-to-1 accuracy

**ARK** **lti** **Carnegie Mellon**

Part-of-Speech Tagging Results

*mini-batch size = 4 for stepwise EM*

Log-Likelihood (x 10^6) vs Wall clock time (hours); Accuracy (%) vs Wall clock time (hours)

Legend:
- Asynch. Stepwise EM (4 processors)
- Synch. Stepwise EM (4 processors)
- Synch. Stepwise EM (1 processor)
- Batch EM (1 processor)

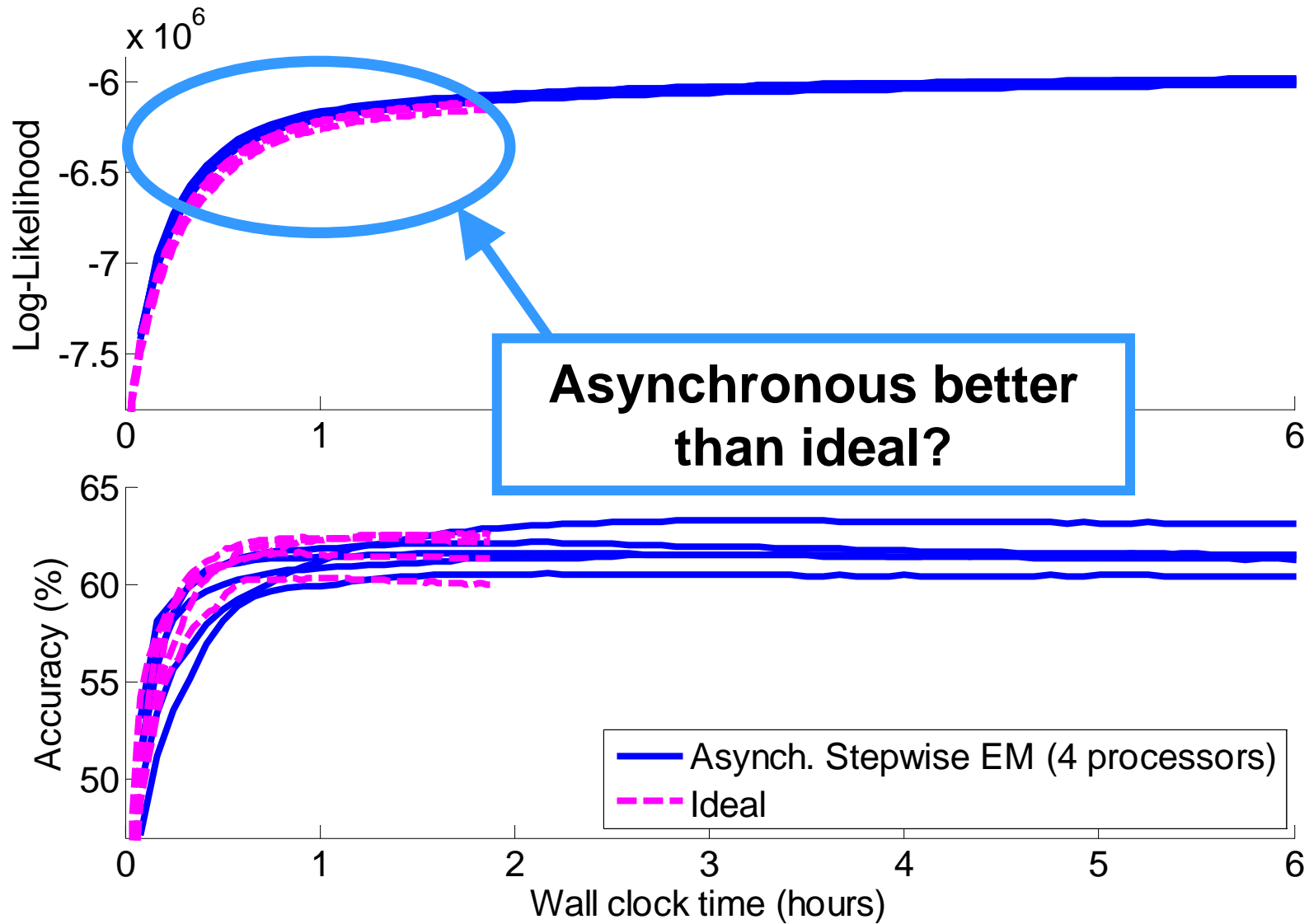ARK  *lti*  **Carnegie Mellon**

Comparison with Ideal

# Comparison with Ideal

# Conclusions and Future Work

- **Asynchronous algorithms speed convergence and do not introduce additional error**
- **Effective for unsupervised learning and non-convex objectives**
- **If your problem works well with small mini-batches, try this!**

- **Future work**
  - ☐ Theoretical results for non-convex case
  - ☐ Explore effects of increasing number of processors
  - ☐ New architectures (maintain multiple copies of $\theta$)

Thanks!