

Distributed Asynchronous Online Learning for Natural Language Processing

Kevin Gimpel Dipanjan Das Noah A. Smith



Carnegie Mellon

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



Carnegie Mellon

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



Carnegie Mellon

Single-Processor Batch Learning

Parameters: θ_t

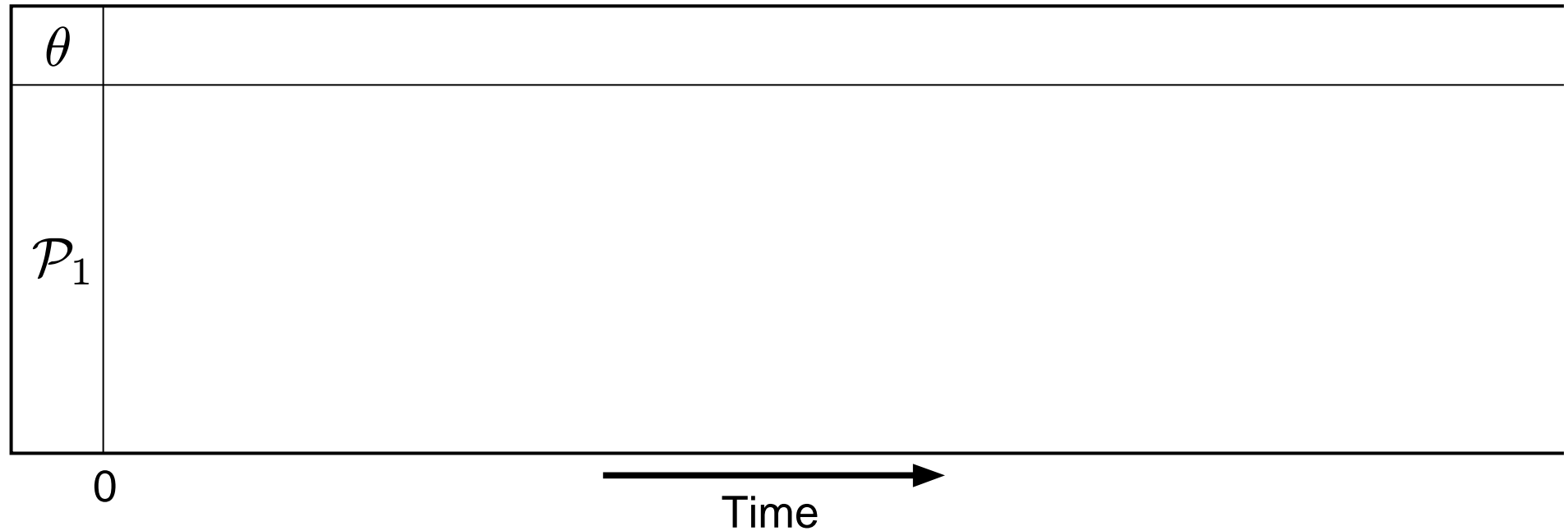
Processors: \mathcal{P}_i

Dataset: \mathcal{D}



Carnegie Mellon

Single-Processor Batch Learning

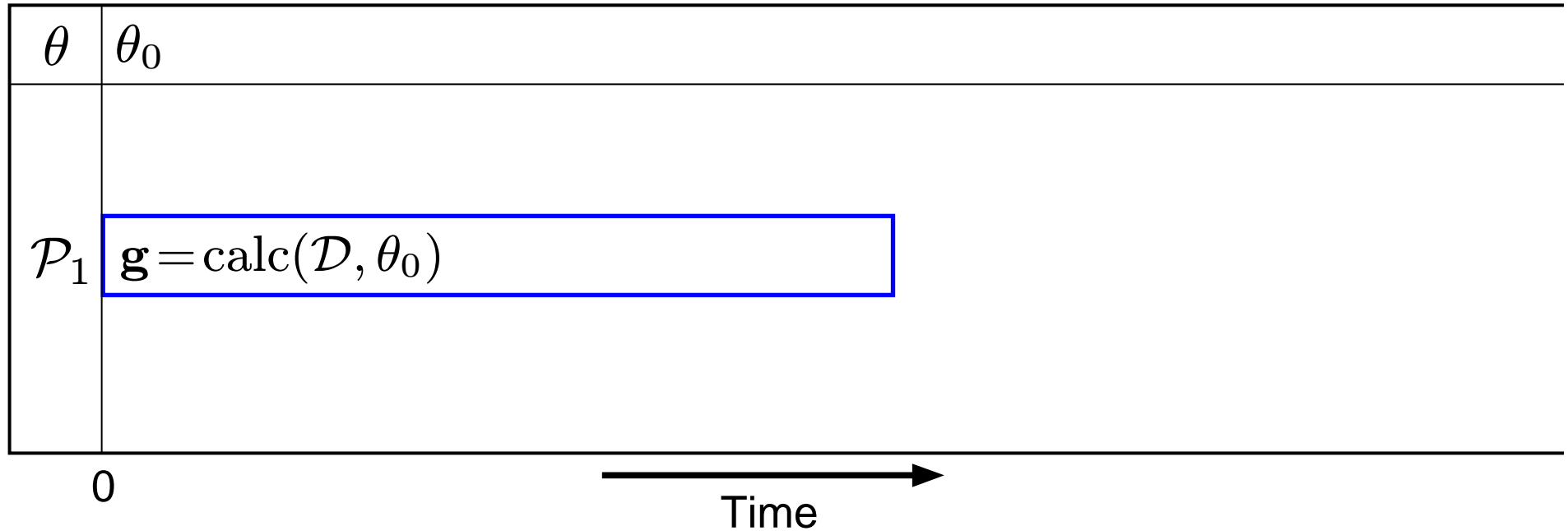


Parameters:	θ_t
Processors:	\mathcal{P}_i
Dataset:	\mathcal{D}



Carnegie Mellon

Single-Processor Batch Learning



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

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

Parameters: θ_t

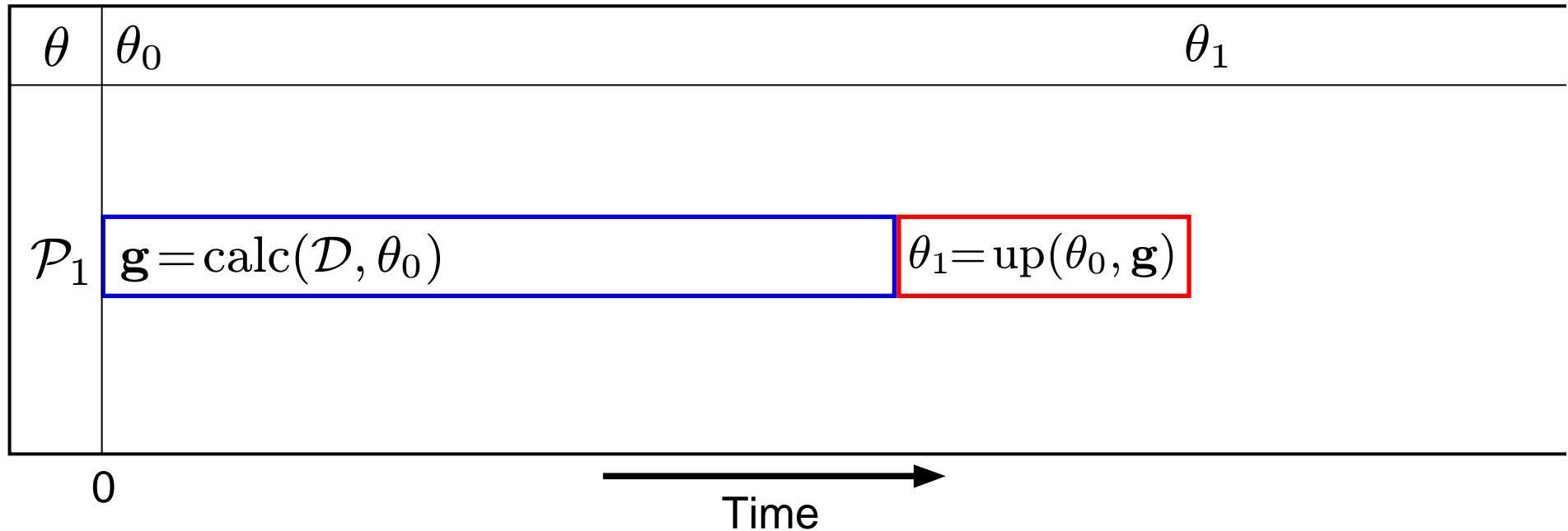
Processors: \mathcal{P}_i

Dataset: \mathcal{D}



Carnegie Mellon

Single-Processor Batch Learning



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

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

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

Update θ_0 using gradient \mathbf{g} to obtain θ_1

Parameters: θ_t

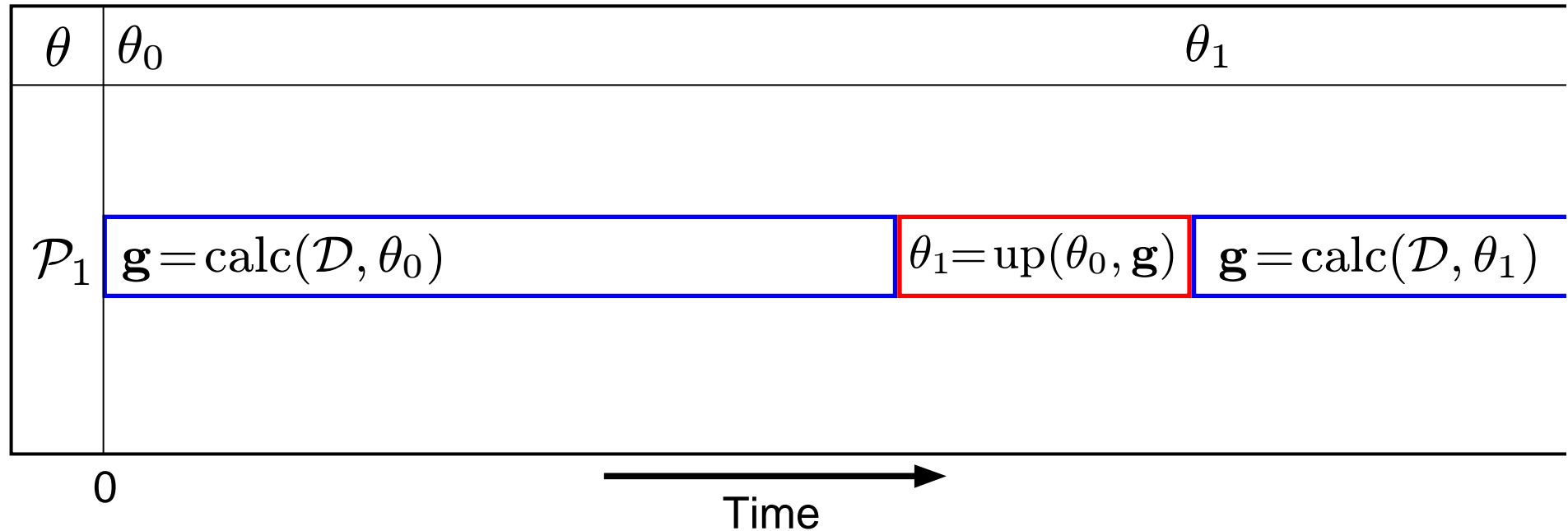
Processors: \mathcal{P}_i

Dataset: \mathcal{D}



Carnegie Mellon

Single-Processor Batch Learning



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

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

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

Update θ_0 using gradient \mathbf{g} to obtain θ_1

Parameters: θ_t

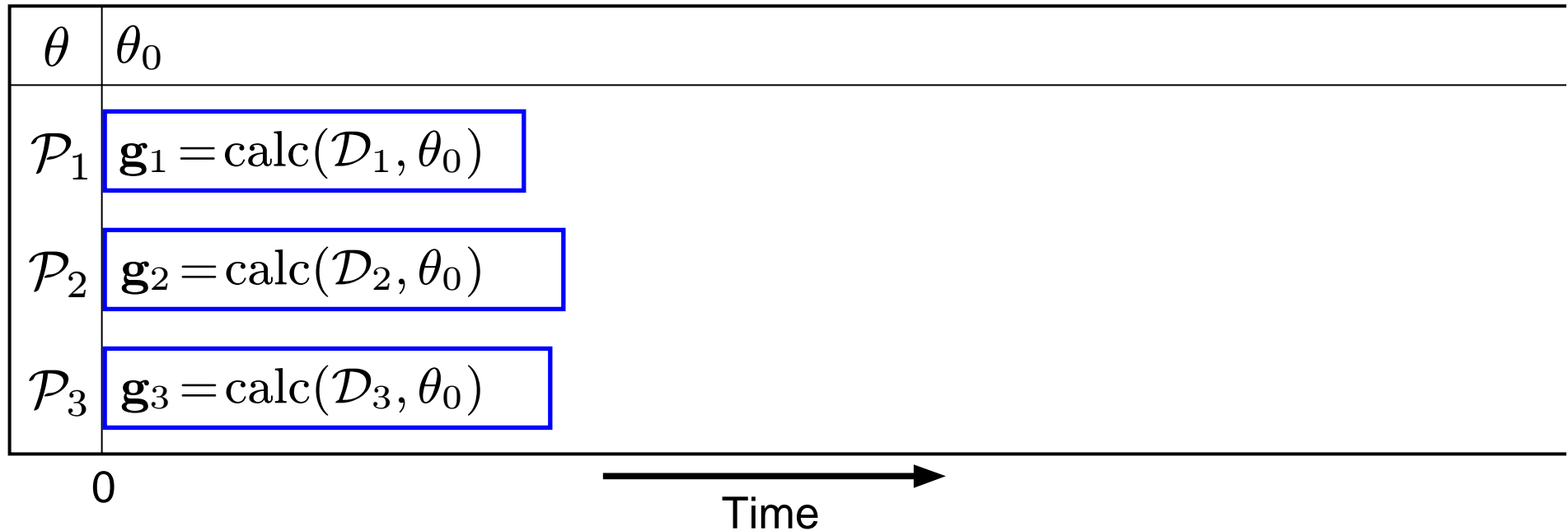
Processors: \mathcal{P}_i

Dataset: \mathcal{D}



Carnegie Mellon

Parallel Batch Learning



- Divide data into parts, compute gradient on parts in parallel

Parameters: θ_t

Processors: \mathcal{P}_i

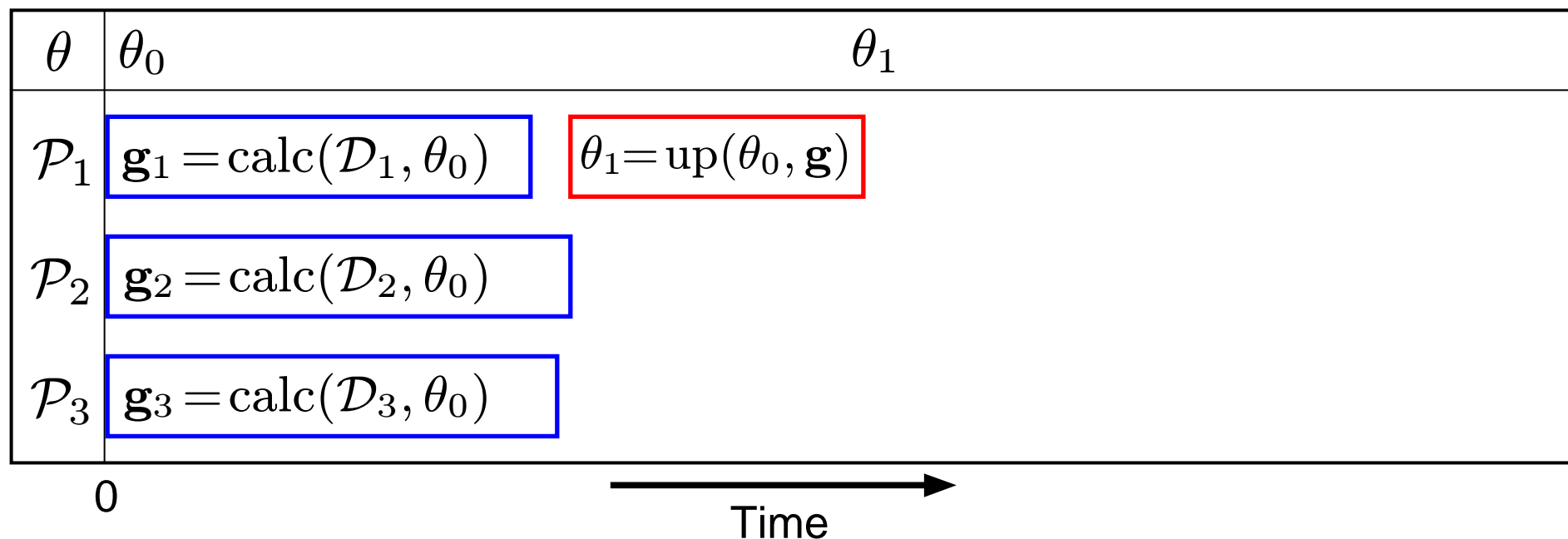
Dataset: $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \mathcal{D}_3$

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



Carnegie Mellon

Parallel Batch Learning



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

Parameters: θ_t

Processors: \mathcal{P}_i

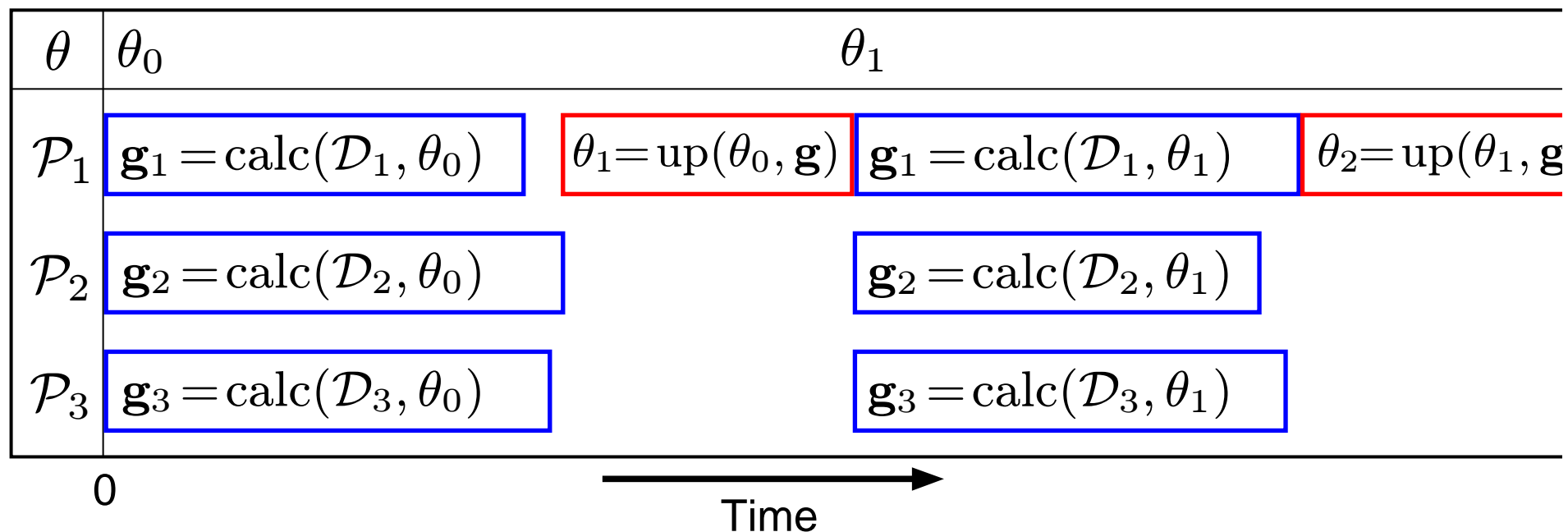
Dataset: $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \mathcal{D}_3$

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



Carnegie Mellon

Parallel Batch Learning

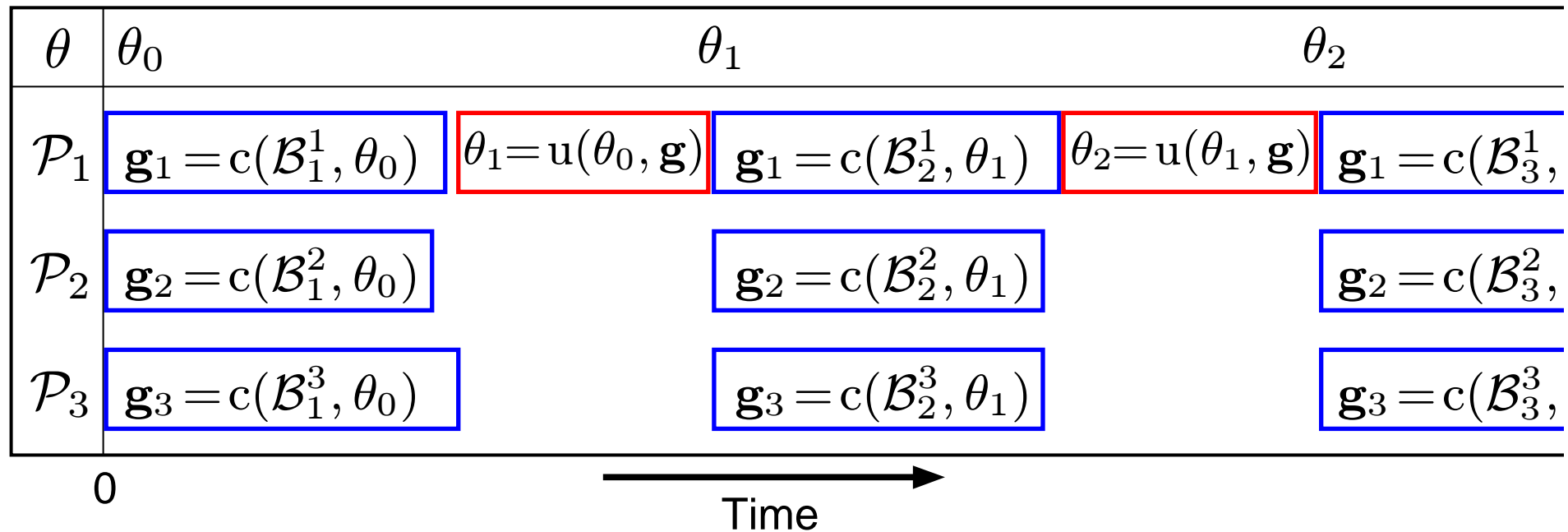


- Divide data into parts, compute gradient on parts in parallel
- One processor updates parameters

Parameters: θ_t
 Processors: \mathcal{P}_i
 Dataset: $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \mathcal{D}_3$
 Gradient: $\mathbf{g} = \mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3$

Parallel Synchronous Mini-Batch Learning

Finkel, Kleeman, and Manning (2008)



- Same architecture, just more frequent updates

Parameters: θ_t

Processors: \mathcal{P}_i

Mini-batches: $\mathcal{B}_t = \mathcal{B}_t^1 \cup \mathcal{B}_t^2 \cup \mathcal{B}_t^3$

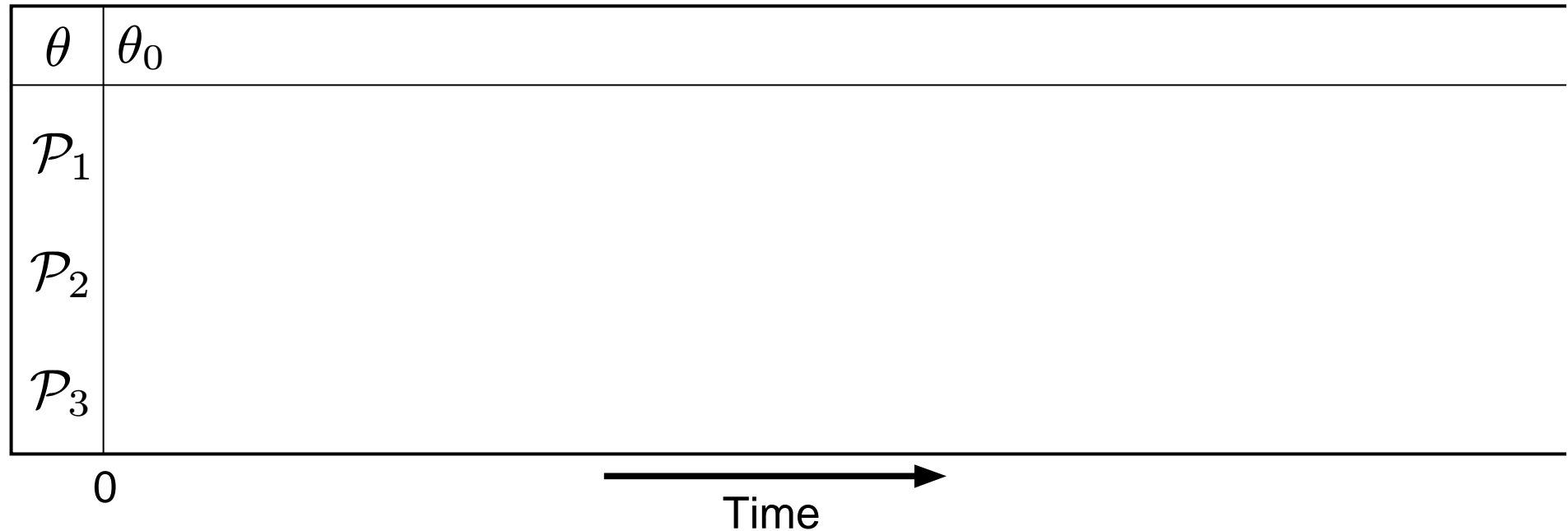
Gradient: $\mathbf{g} = \mathbf{g}_1 + \mathbf{g}_2 + \mathbf{g}_3$



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



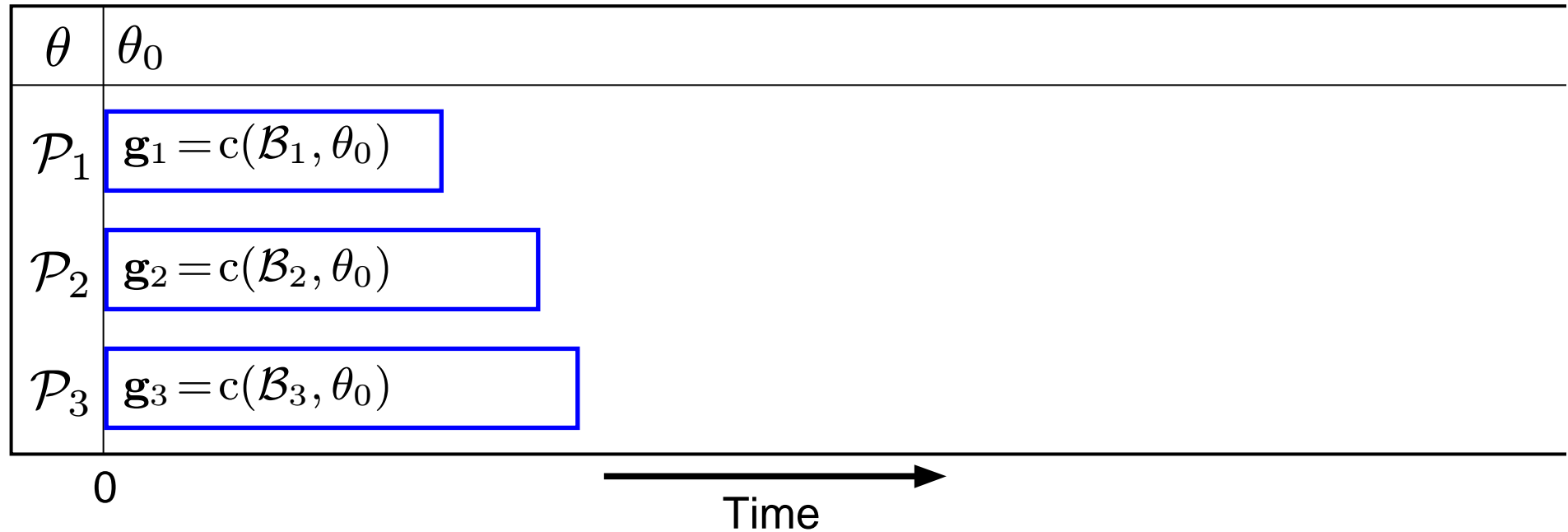
Parameters: θ_t
Processors: \mathcal{P}_i
Mini-batches: \mathcal{B}_j
Gradient: \mathbf{g}_k



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



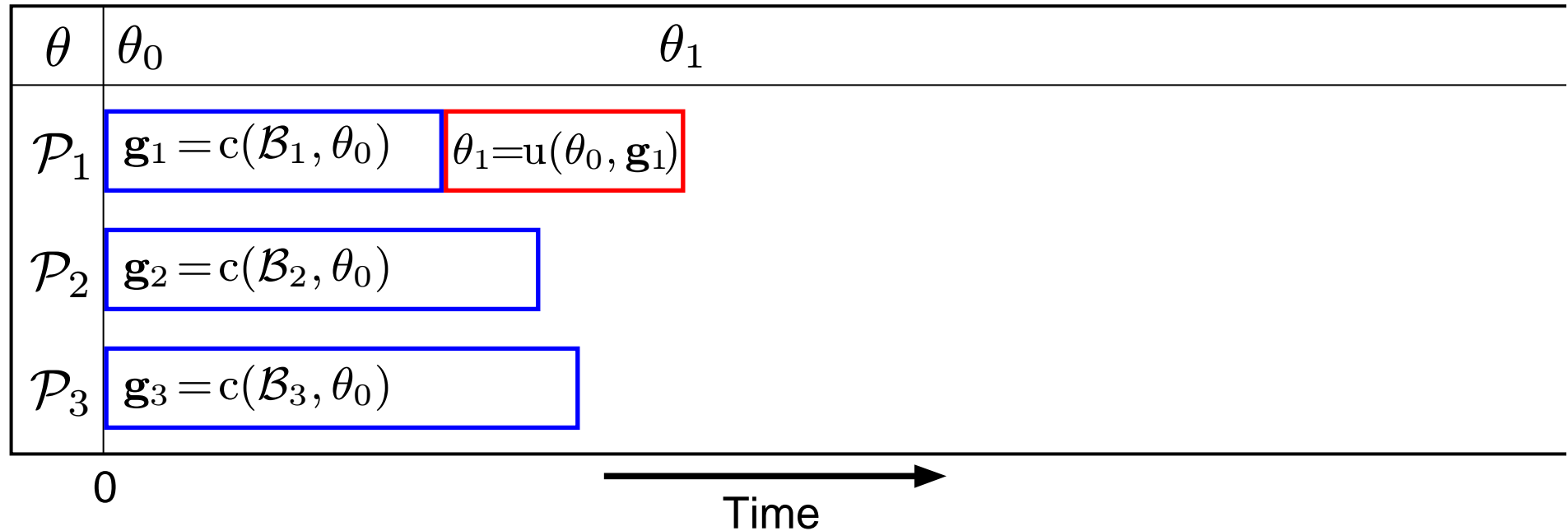
Parameters:	θ_t
Processors:	\mathcal{P}_i
Mini-batches:	\mathcal{B}_j
Gradient:	\mathbf{g}_k



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



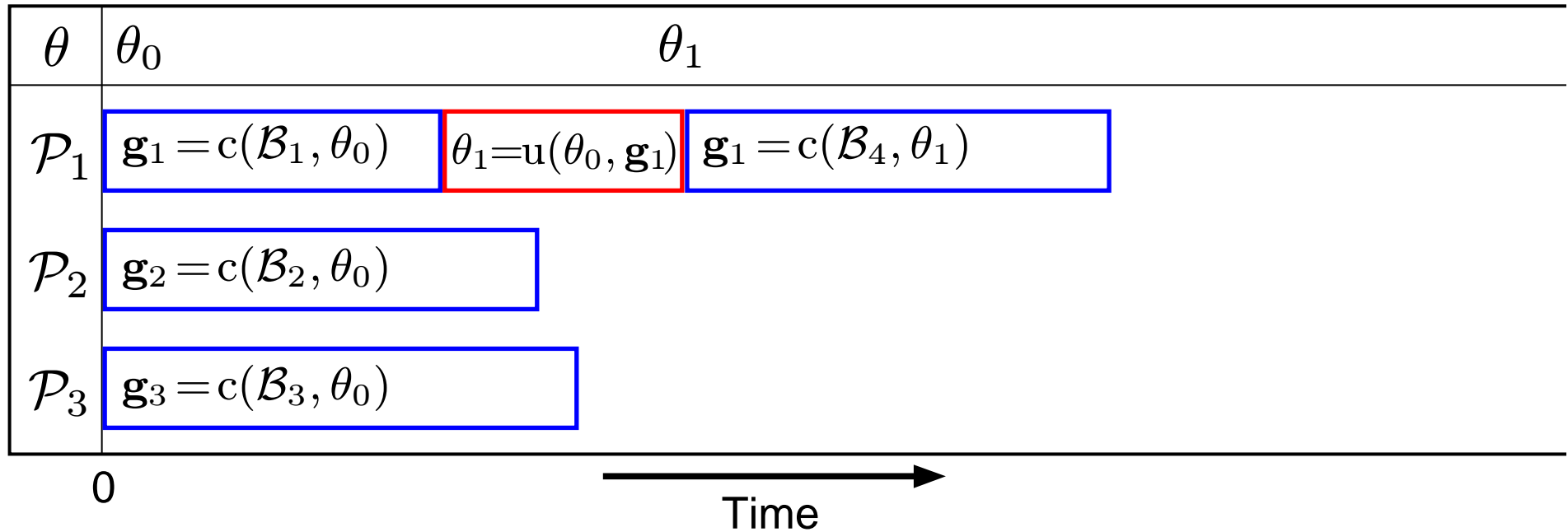
Parameters: θ_t
Processors: \mathcal{P}_i
Mini-batches: \mathcal{B}_j
Gradient: \mathbf{g}_k



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



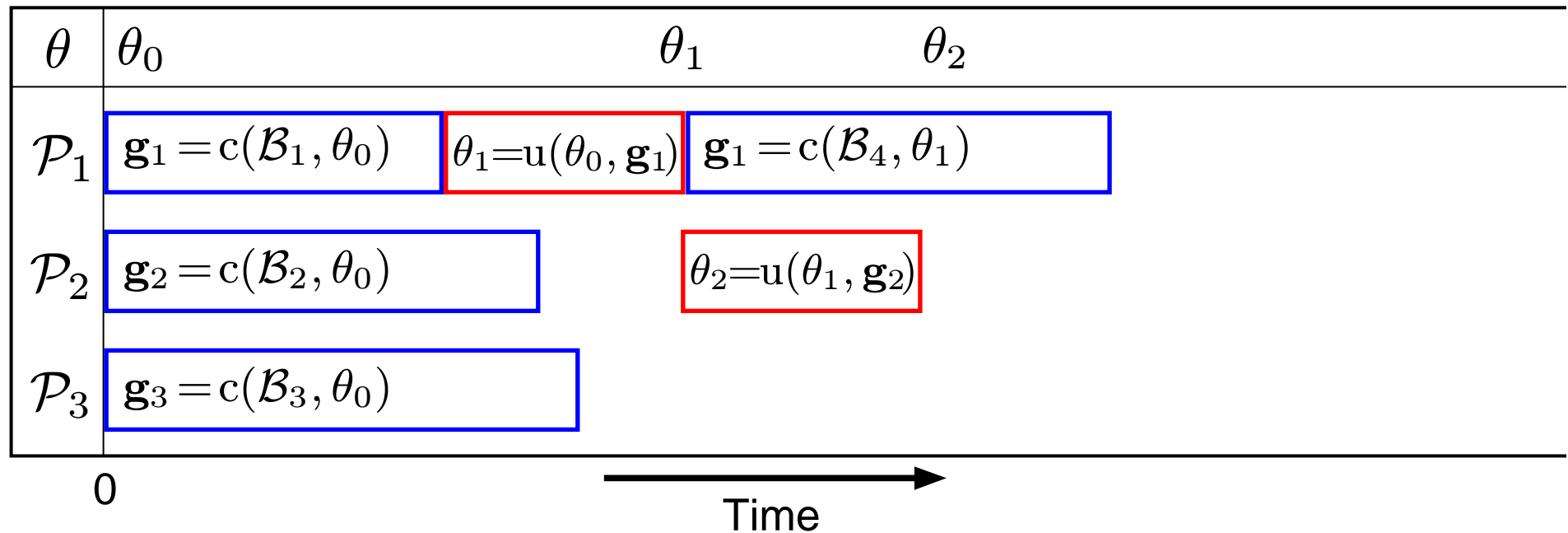
Parameters:	θ_t
Processors:	\mathcal{P}_i
Mini-batches:	\mathcal{B}_j
Gradient:	\mathbf{g}_k



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



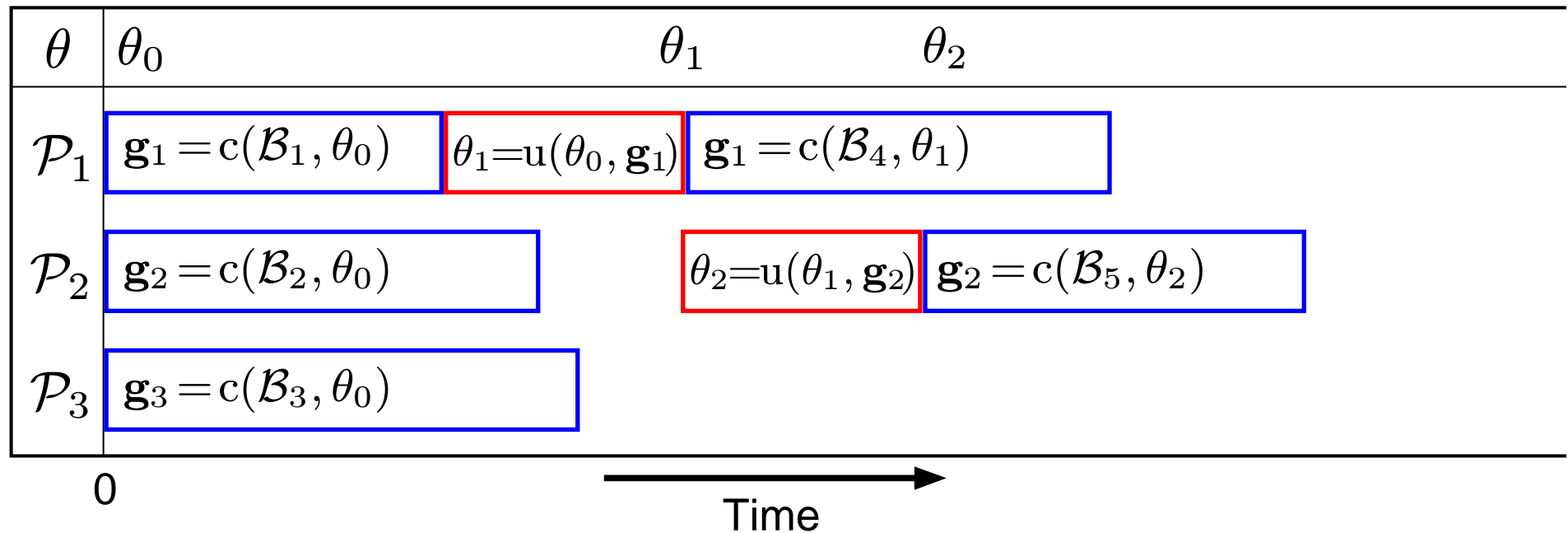
Parameters:	θ_t
Processors:	\mathcal{P}_i
Mini-batches:	\mathcal{B}_j
Gradient:	\mathbf{g}_k



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



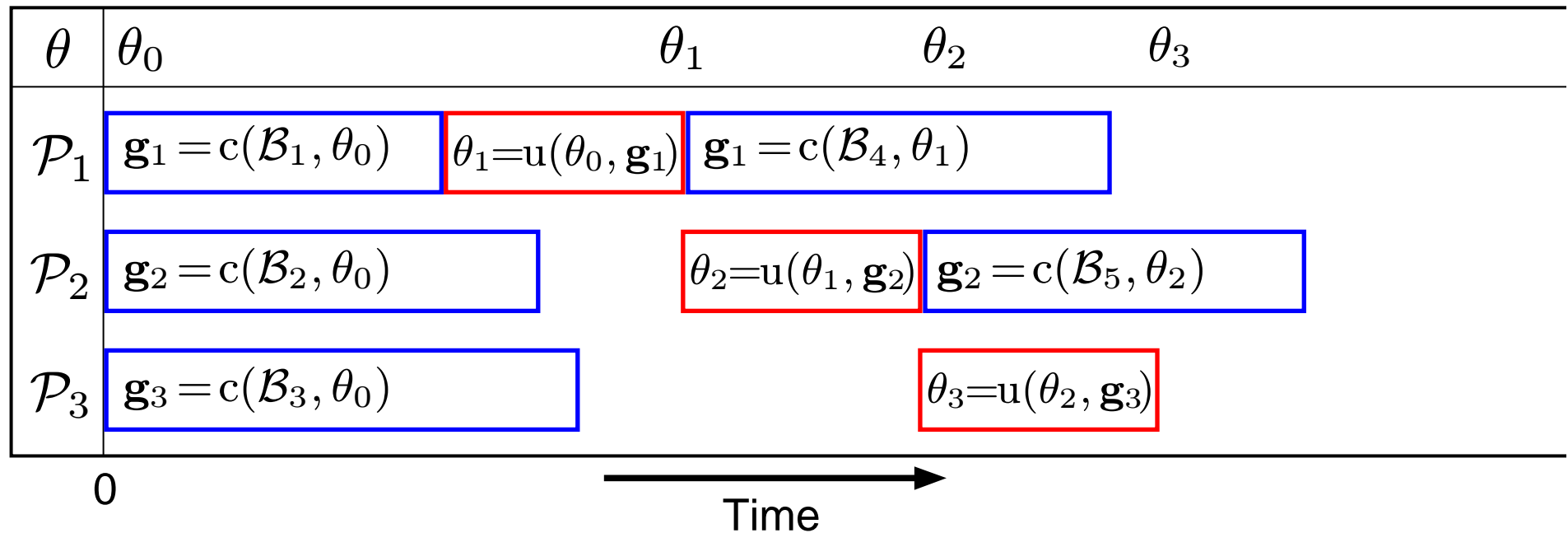
Parameters: θ_t
Processors: \mathcal{P}_i
Mini-batches: \mathcal{B}_j
Gradient: \mathbf{g}_k



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



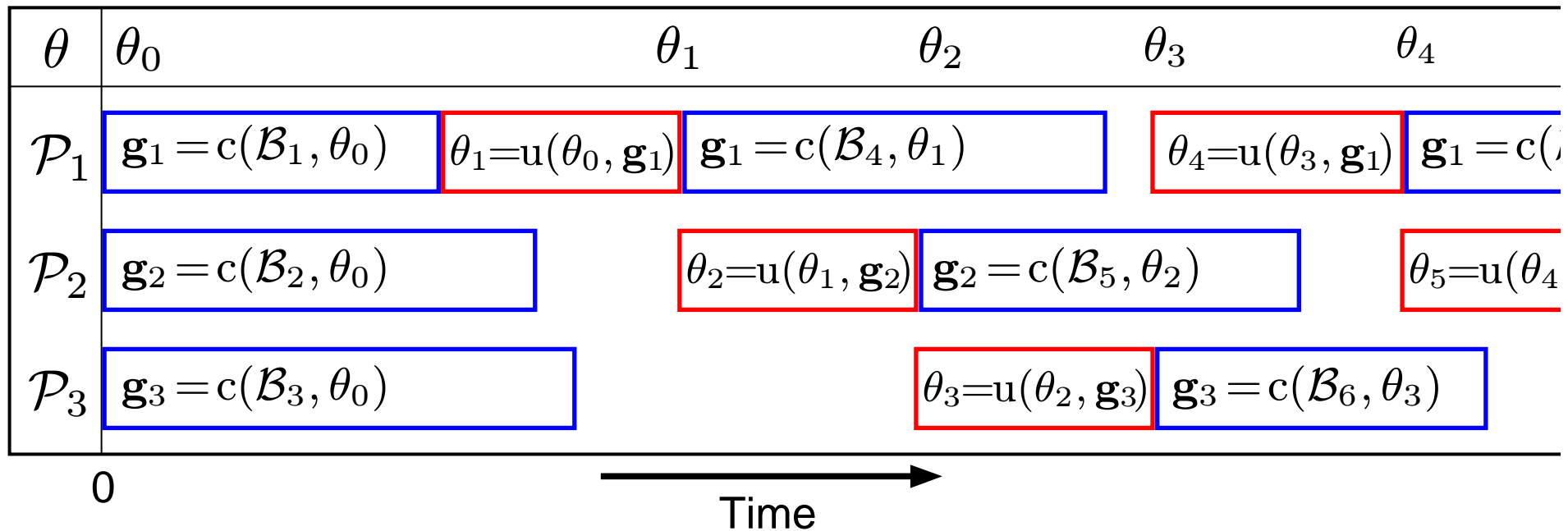
Parameters: θ_t
 Processors: \mathcal{P}_i
 Mini-batches: \mathcal{B}_j
 Gradient: \mathbf{g}_k



Carnegie Mellon

Parallel Asynchronous Mini-Batch Learning

Nedic, Bertsekas, and Borkar (2001)



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

Parameters:	θ_t
Processors:	\mathcal{P}_i
Mini-batches:	\mathcal{B}_j
Gradient:	\mathbf{g}_k



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](#))



Carnegie Mellon

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



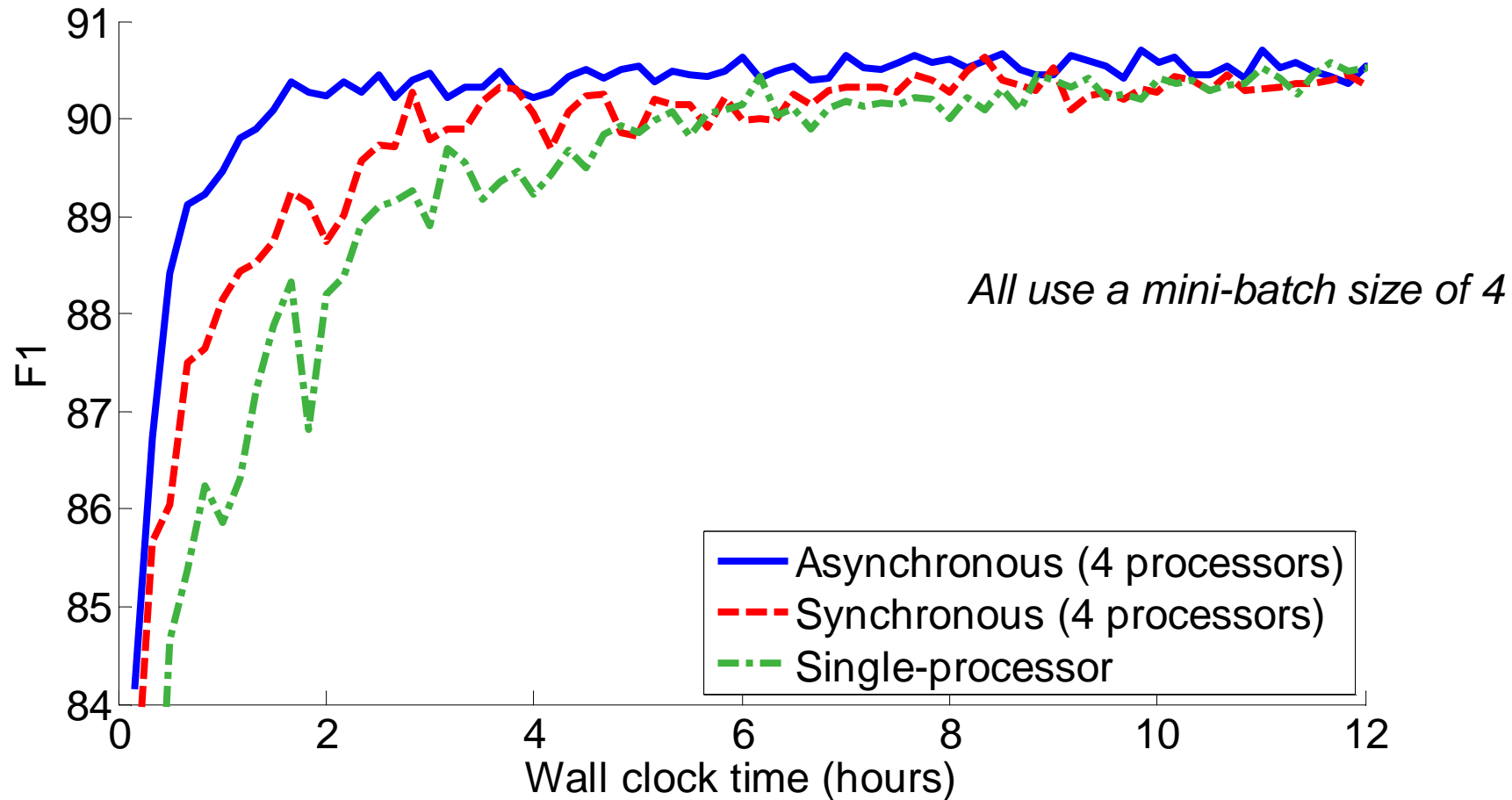
Carnegie Mellon

Experiments

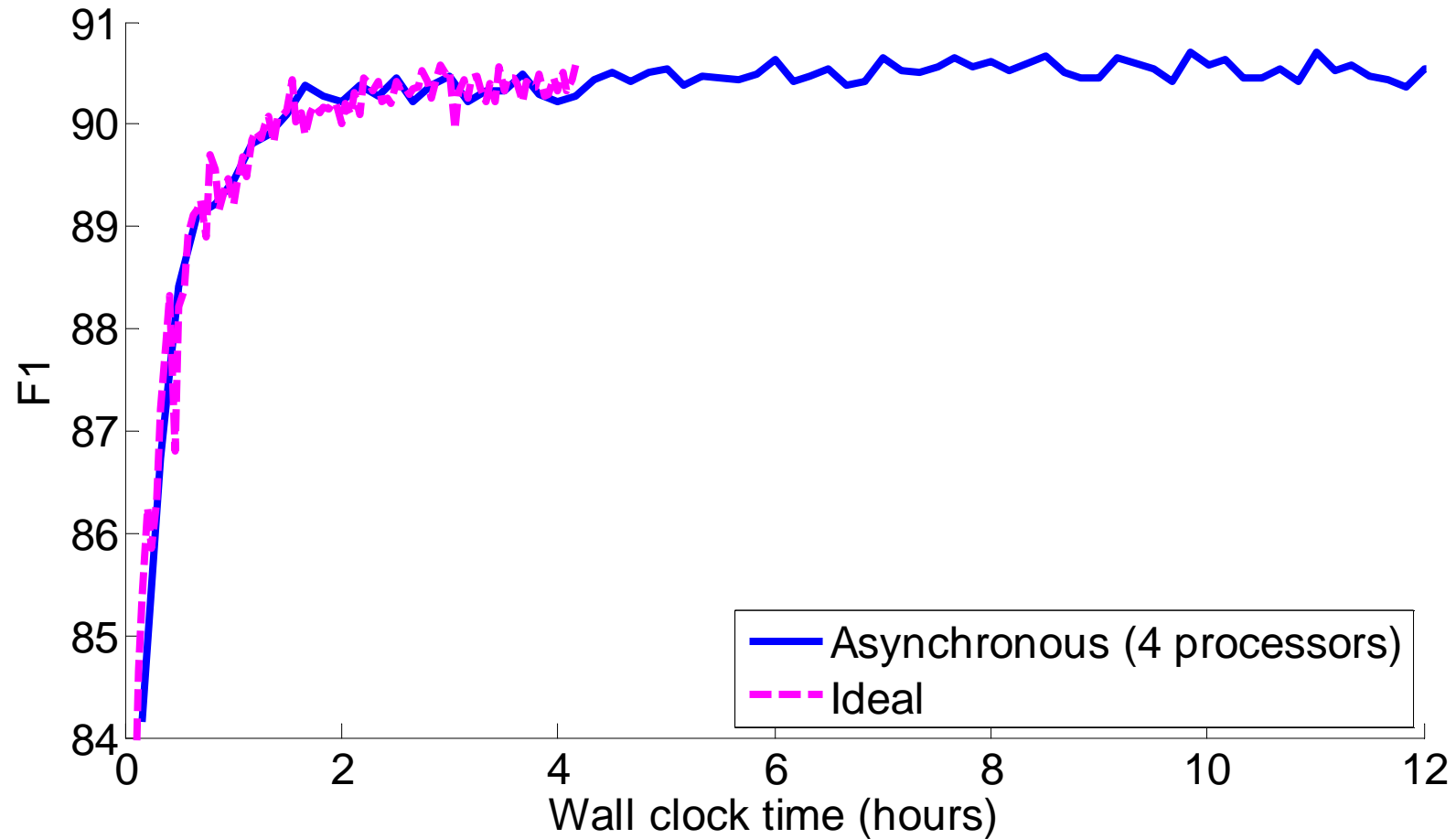
Task	Model	Method	Convex?	$ \mathcal{D} $	$ \theta $	m
Named-Entity Recognition	CRF	Stochastic Gradient Descent	Y	15k	1.3M	4

- 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



Comparison with Ideal Speed-up



Carnegie Mellon

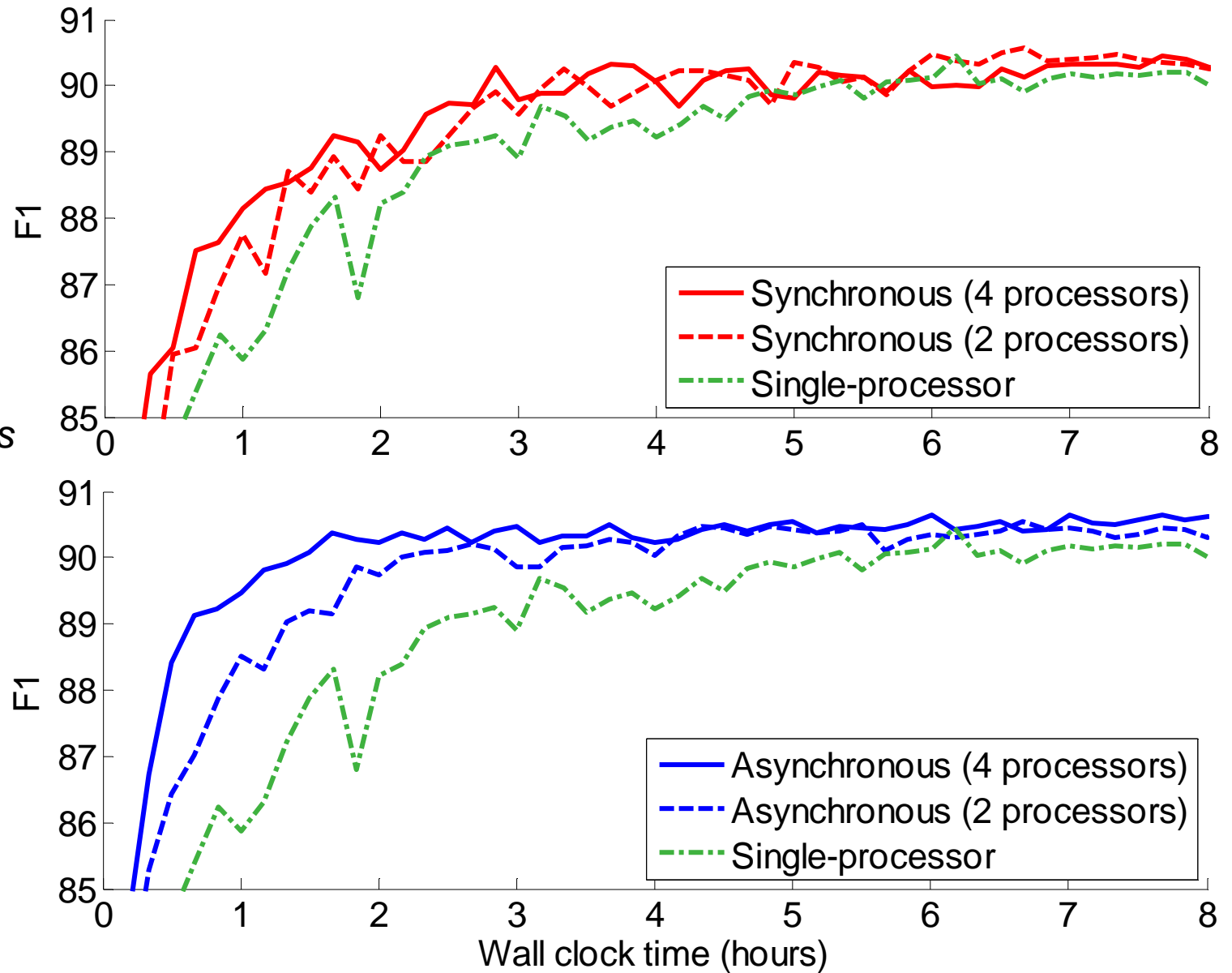
Why Does Asynchronous Converge Faster?

- Processors are kept in near-constant use
- Synchronous SGD leads to idle processors → need for load-balancing



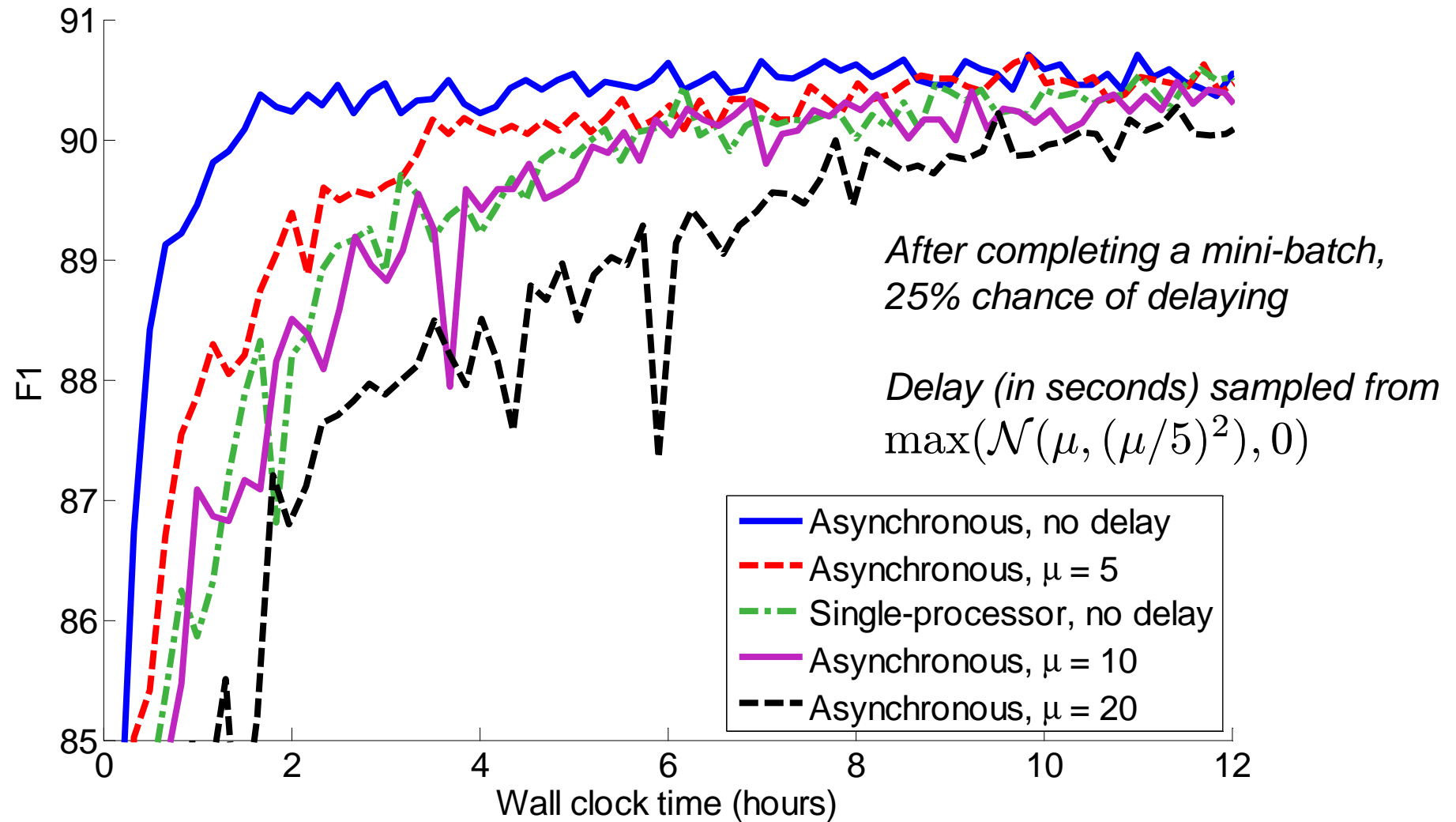
Carnegie Mellon

Clearer improvement for asynchronous algorithms when increasing number of processors



Carnegie Mellon

Artificial Delays



Avg. time per mini-batch = 0.62 s

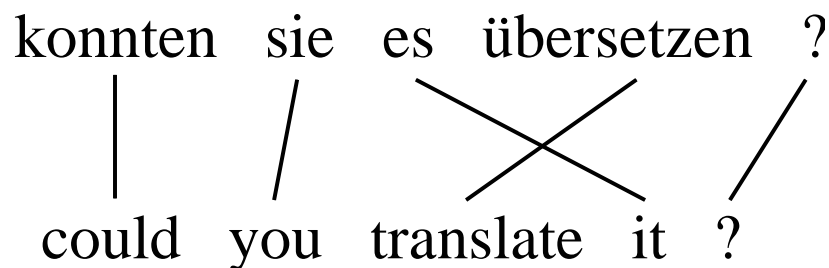


Carnegie Mellon

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:



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



Carnegie Mellon

Stepwise EM

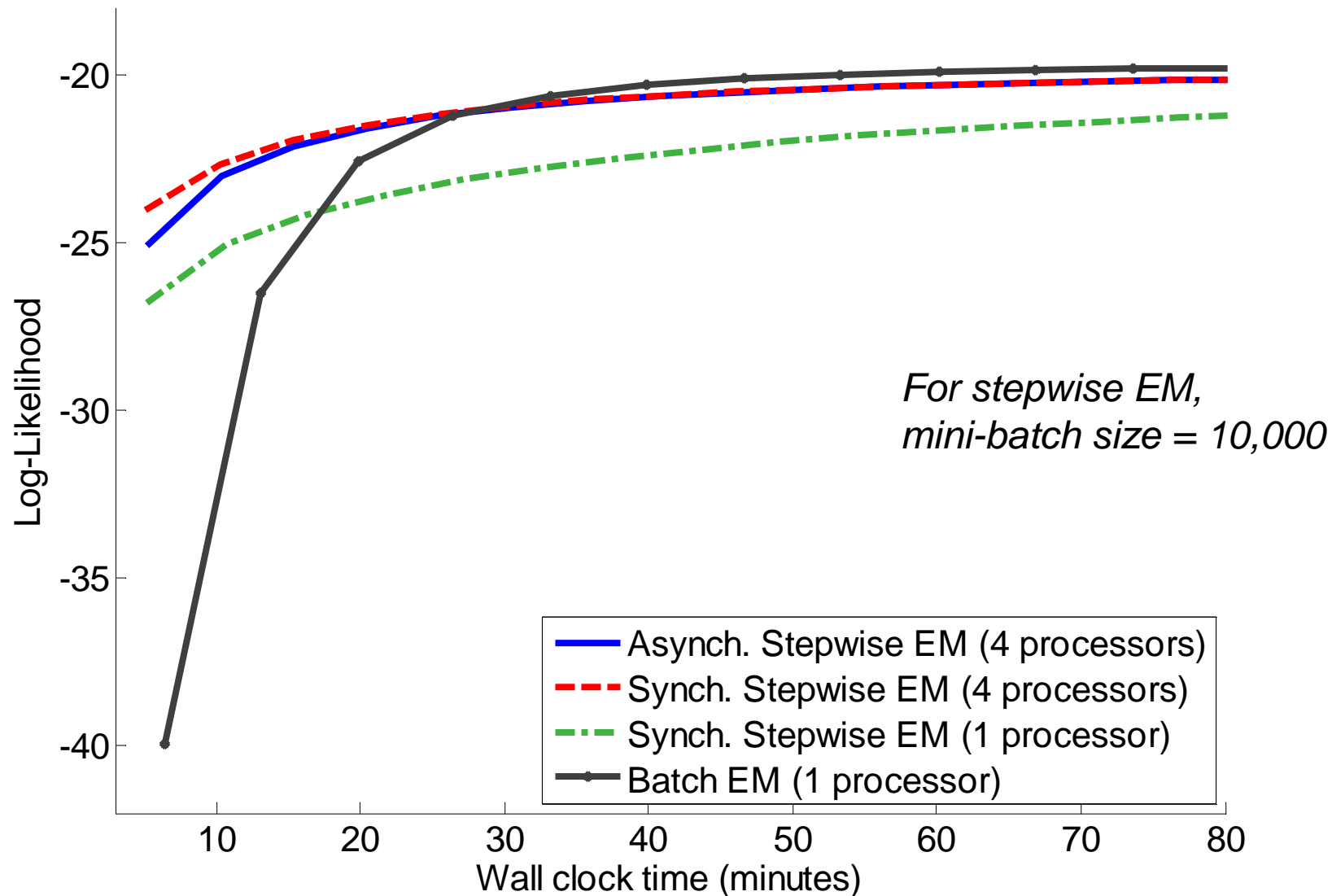
(Sato and Ishii, 2000; Cappe and Moulines, 2009)

- 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)



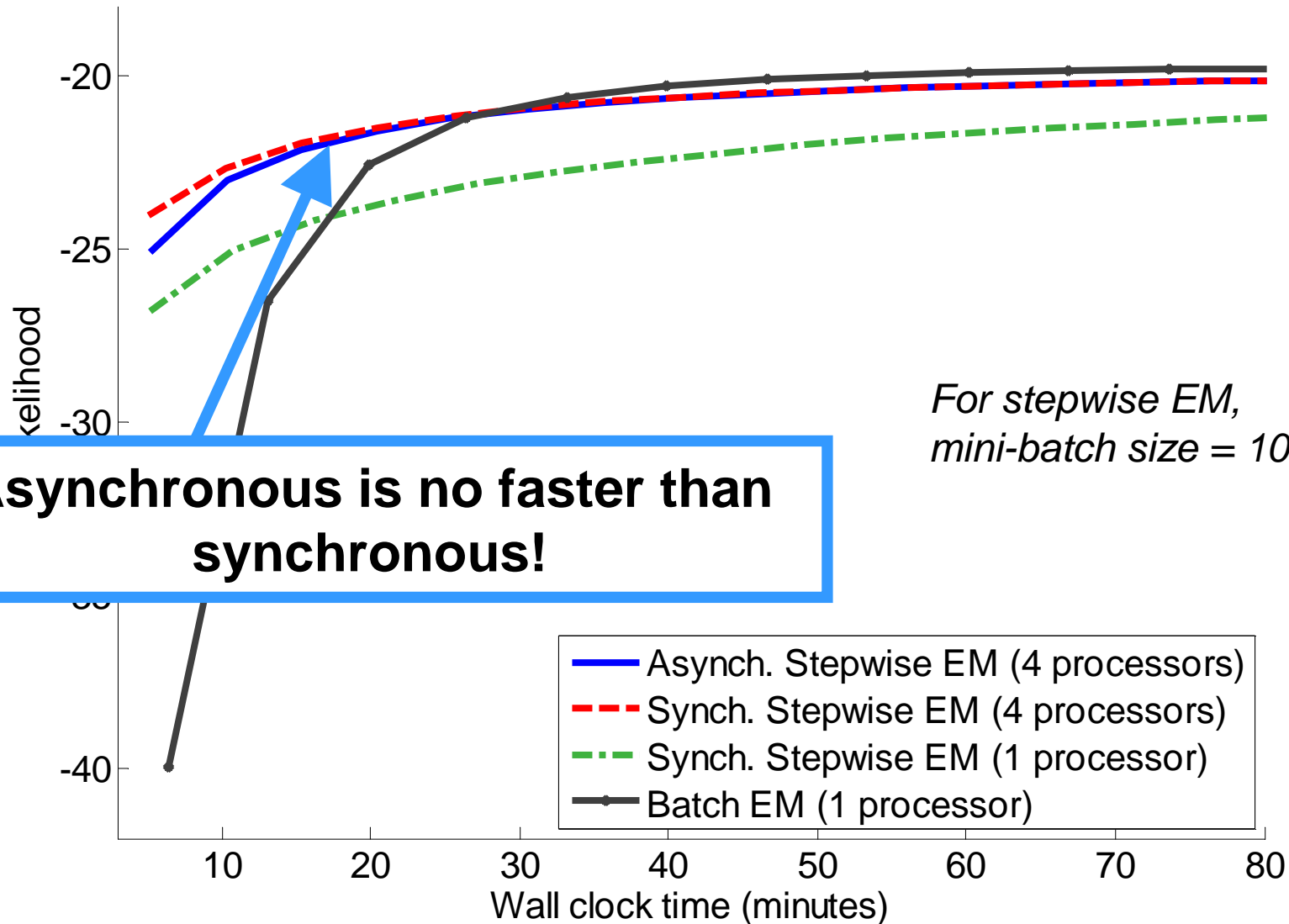
Carnegie Mellon

Word Alignment Results



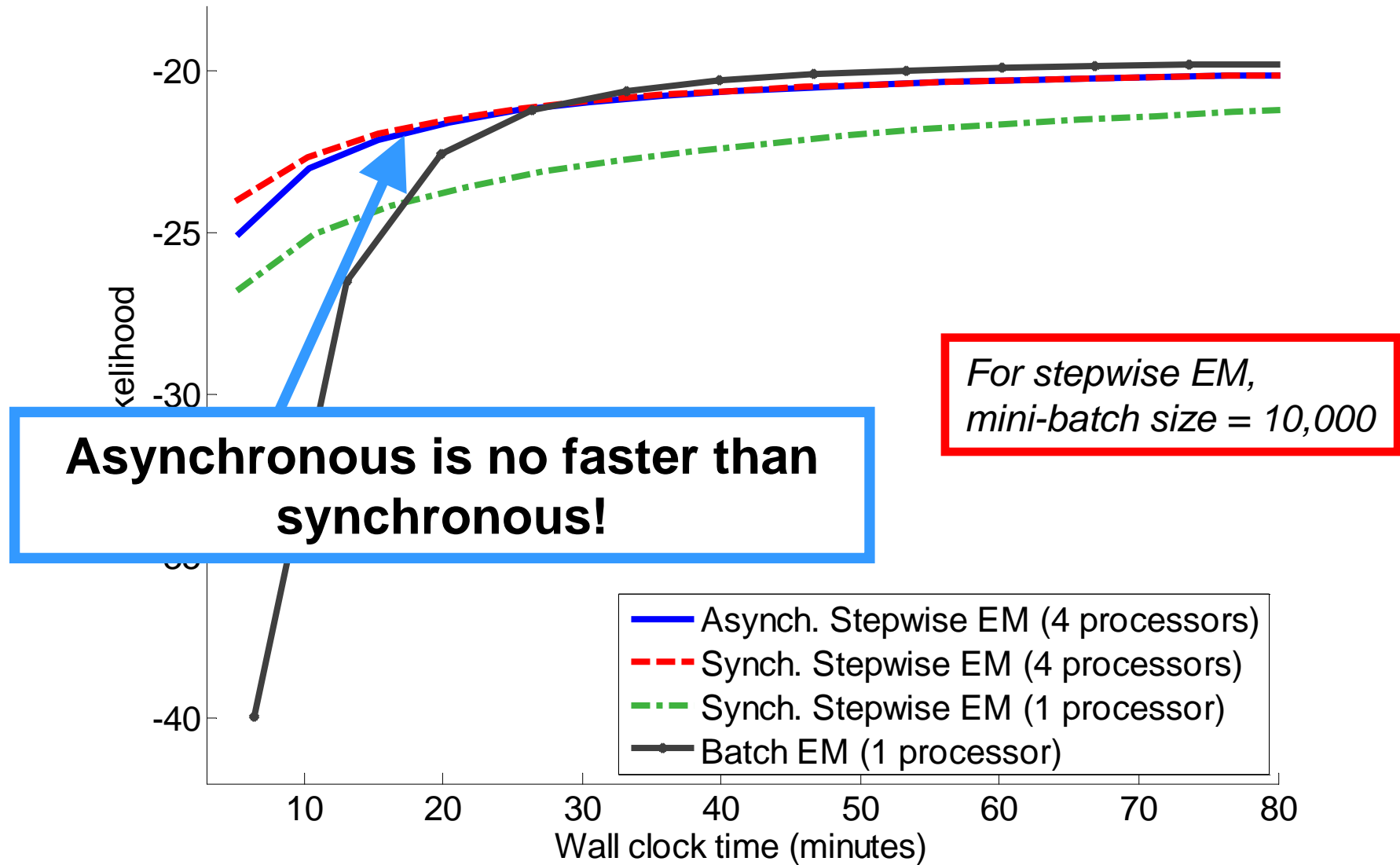
Carnegie Mellon

Word Alignment Results



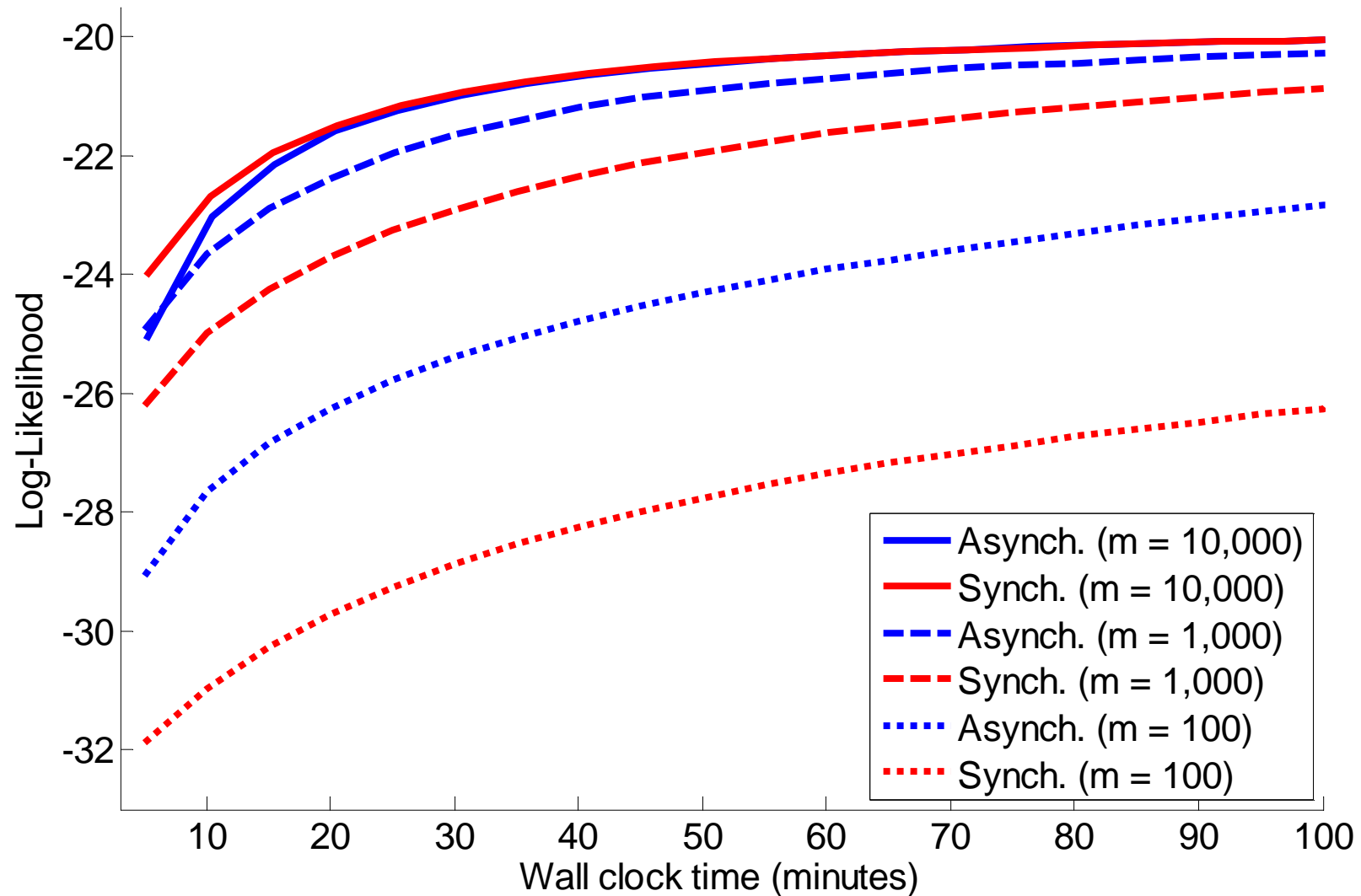
Carnegie Mellon

Word Alignment Results



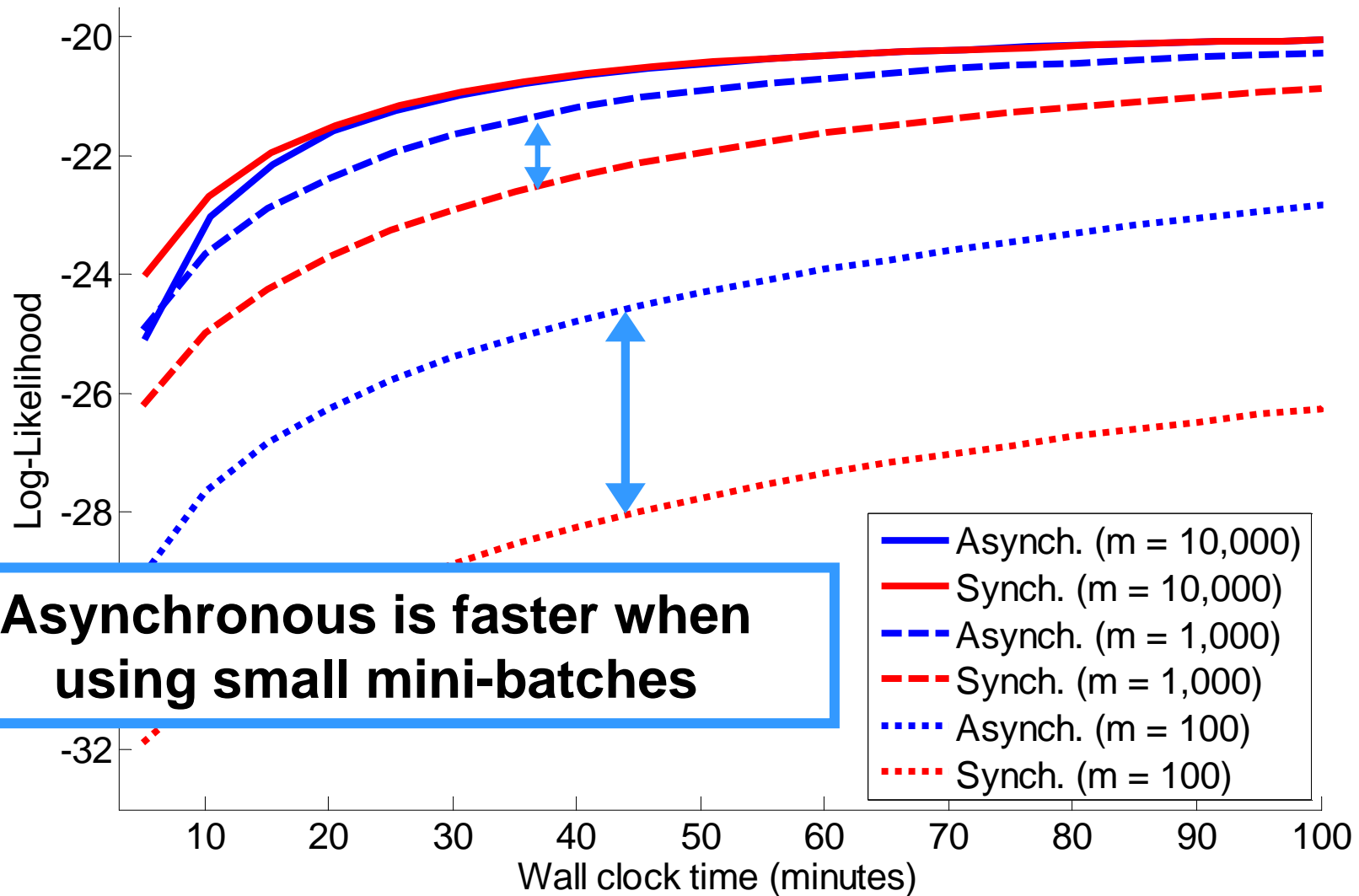
Carnegie Mellon

Comparing Mini-Batch Sizes



Carnegie Mellon

Comparing Mini-Batch Sizes

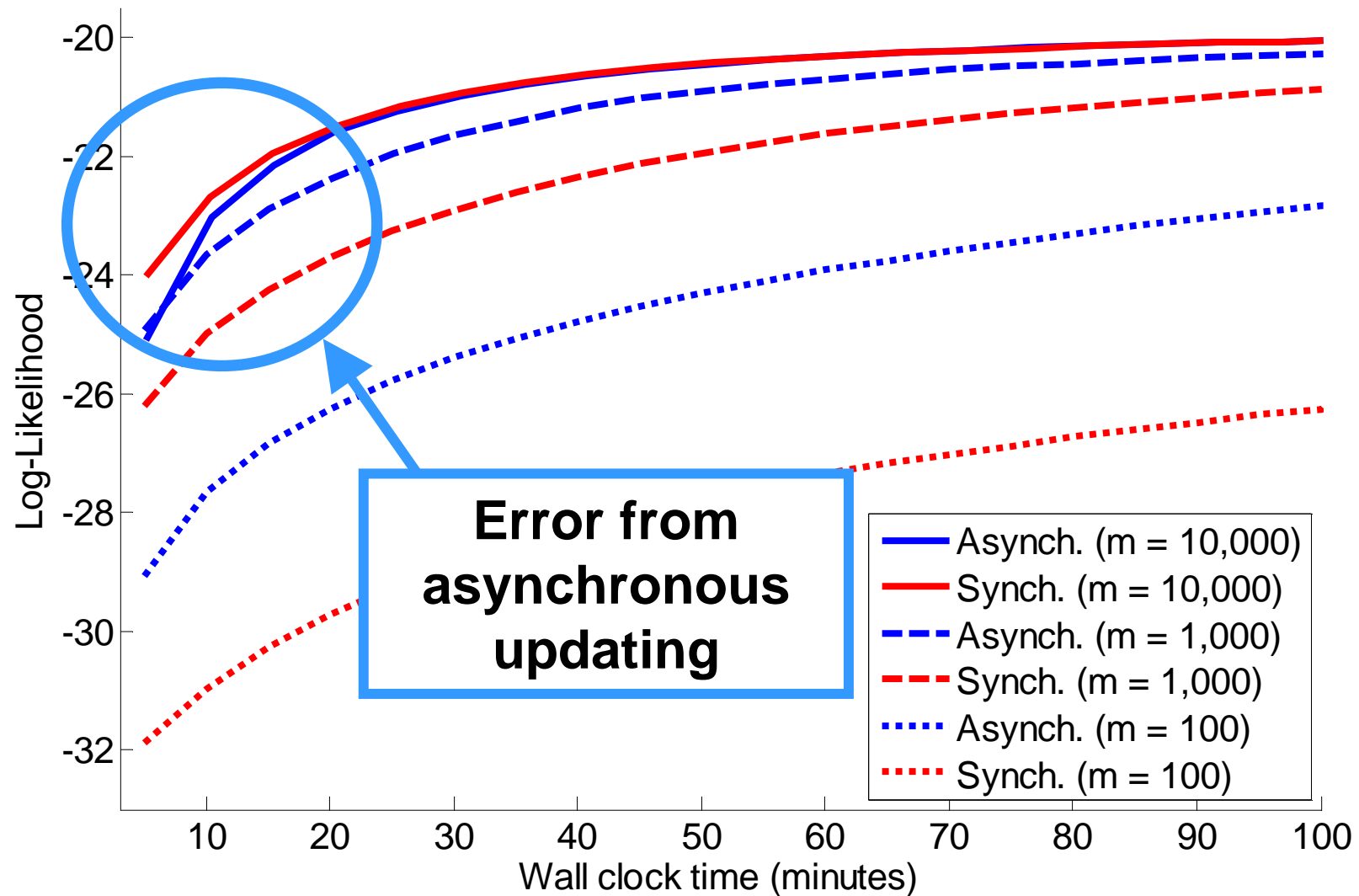


Asynchronous is faster when using small mini-batches



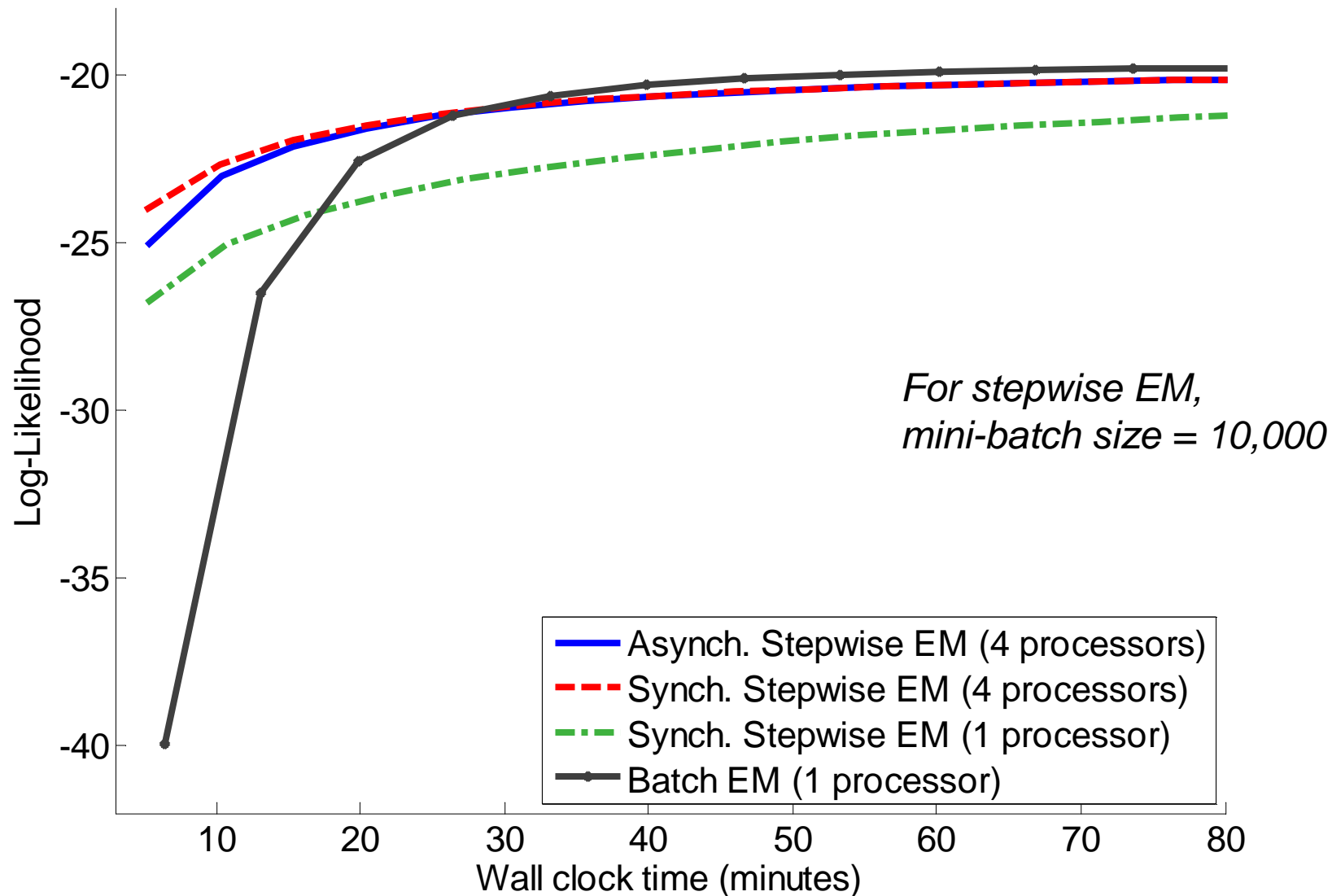
Carnegie Mellon

Comparing Mini-Batch Sizes



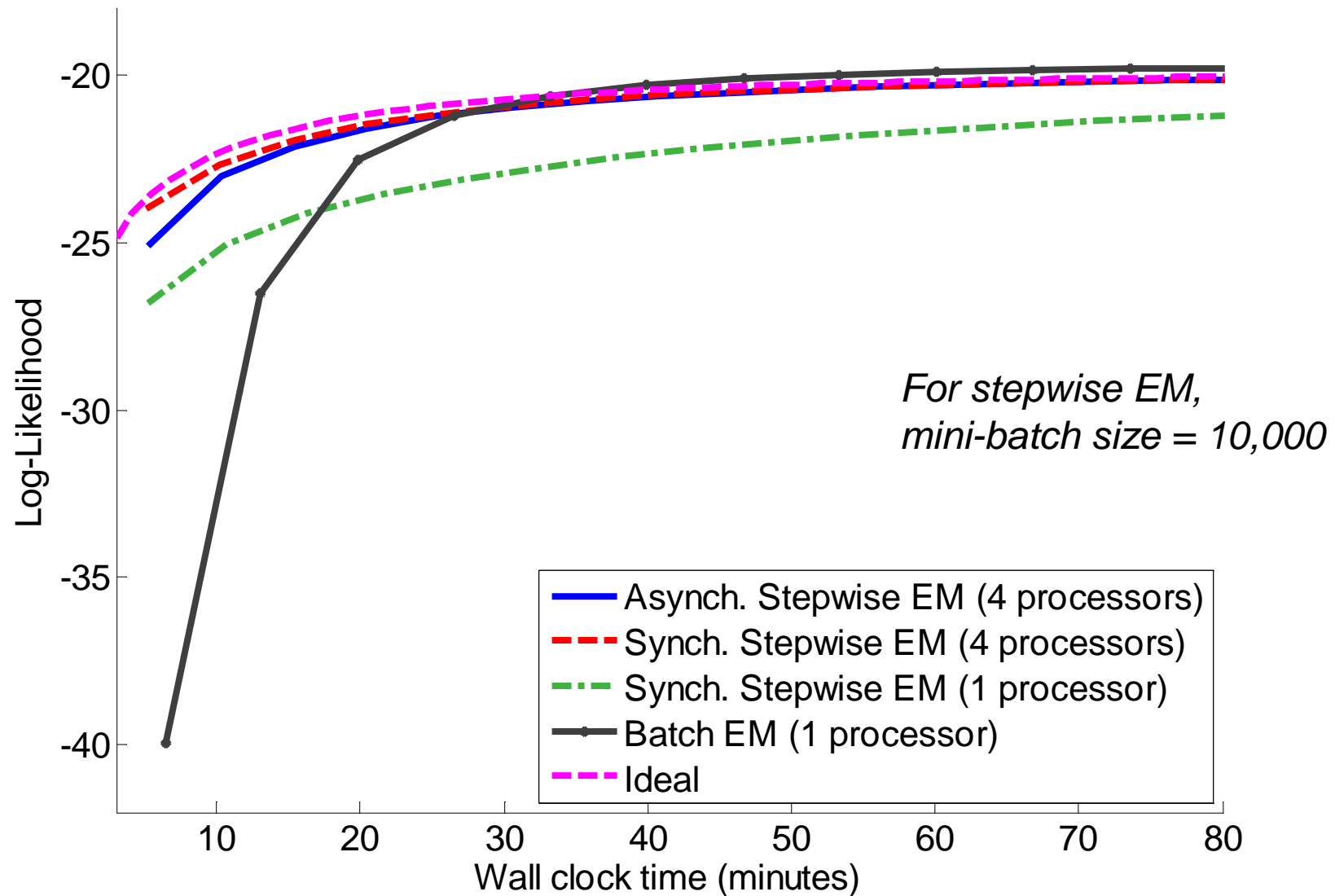
Carnegie Mellon

Word Alignment Results



Carnegie Mellon

Comparison with Ideal Speed-up



Carnegie Mellon

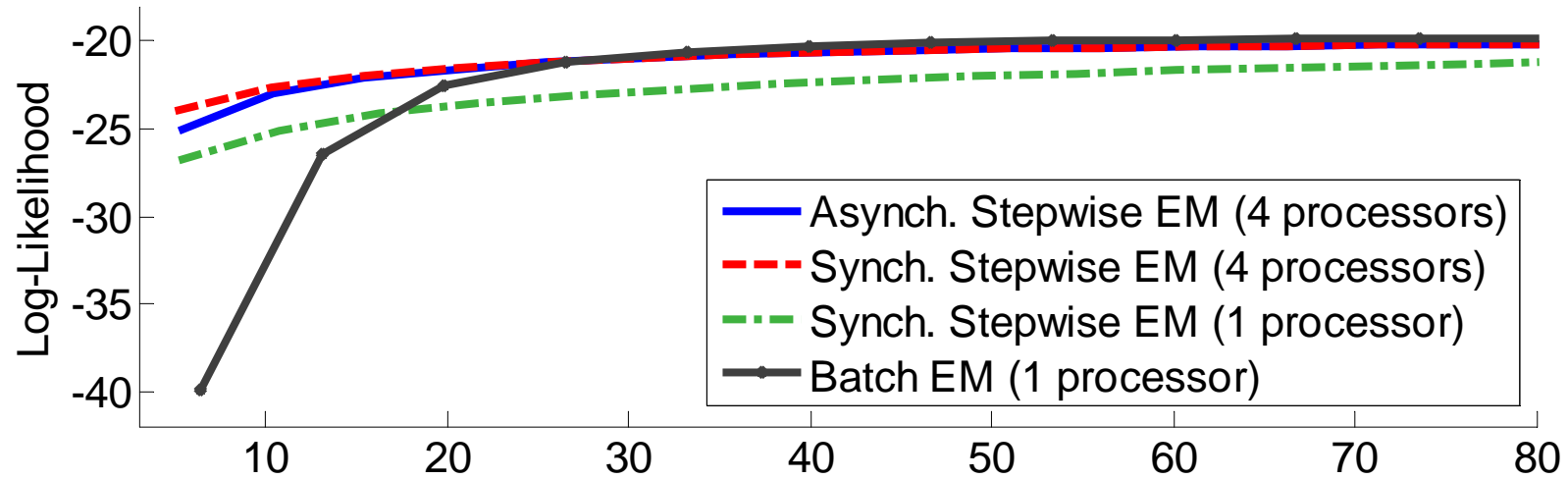
MapReduce?

- 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



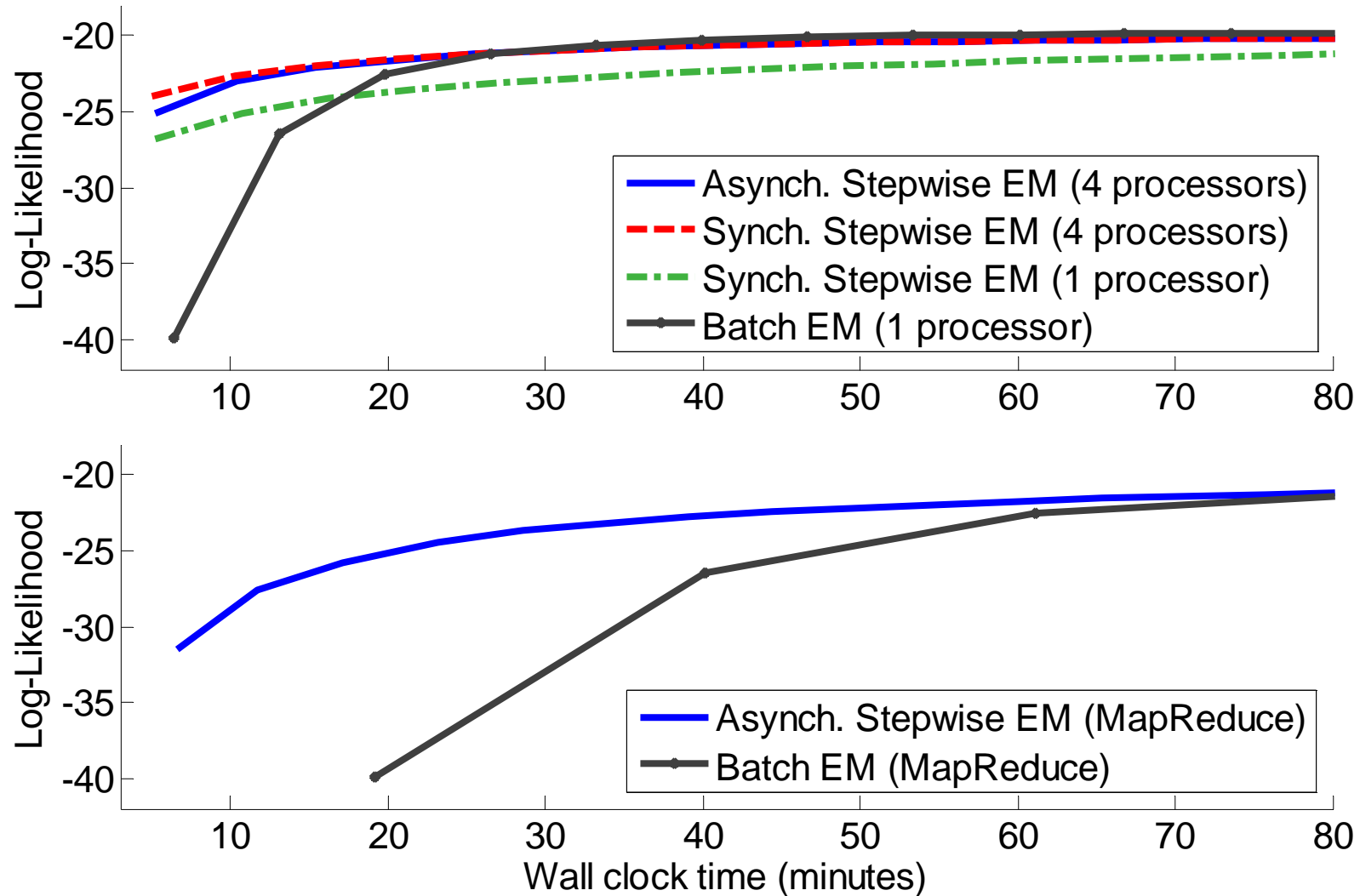
Carnegie Mellon

MapReduce?



Carnegie Mellon

MapReduce?



Carnegie Mellon

Experiments

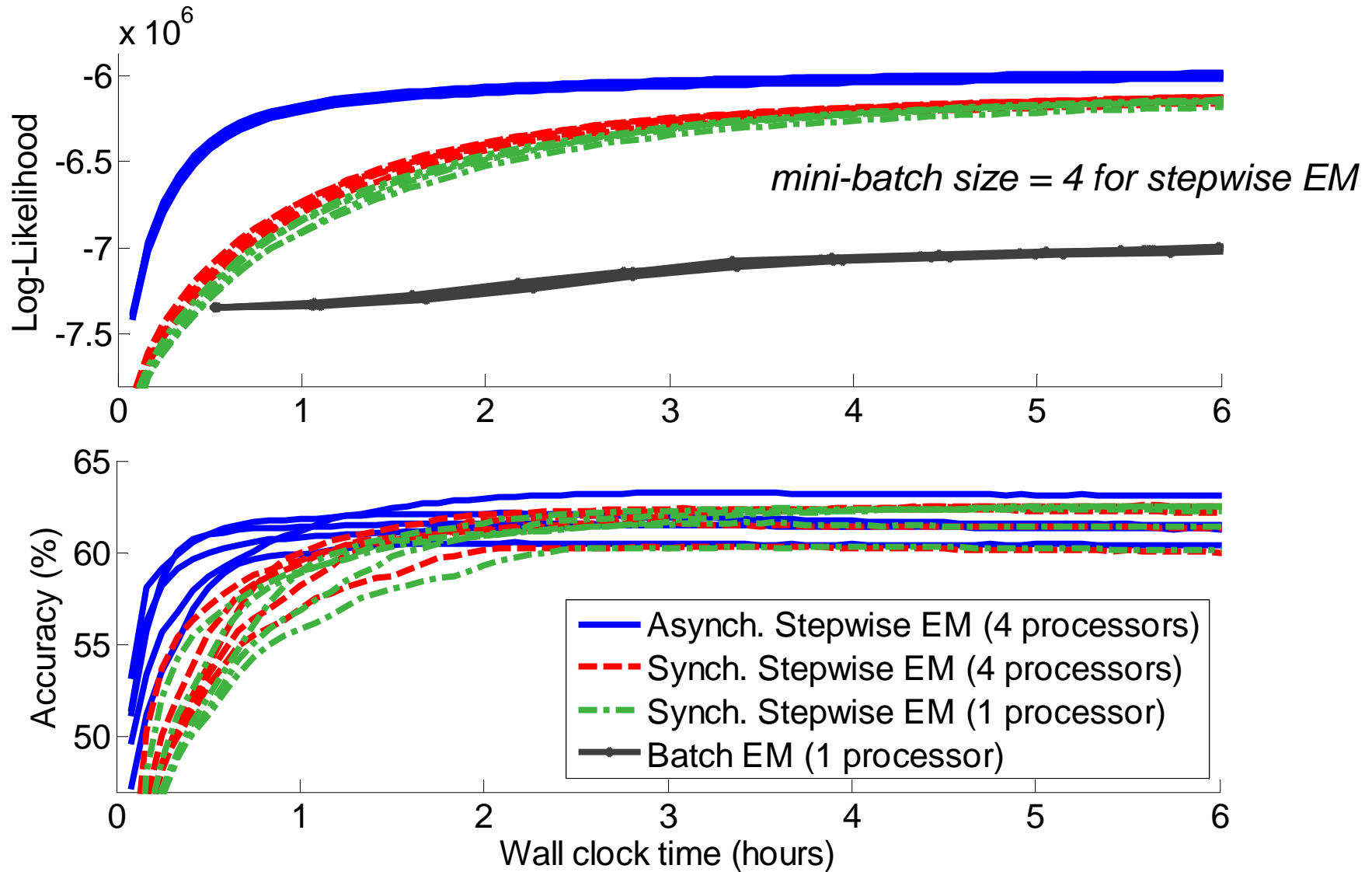
Task	Model	Method	Convex?	$ \mathcal{D} $	$ \theta $	m
Unsupervised Part-of-Speech Tagging	HMM	Stepwise EM	N	42k	2M	4

- Bigram HMM with 45 states
- We plot convergence in likelihood and many-to-1 accuracy



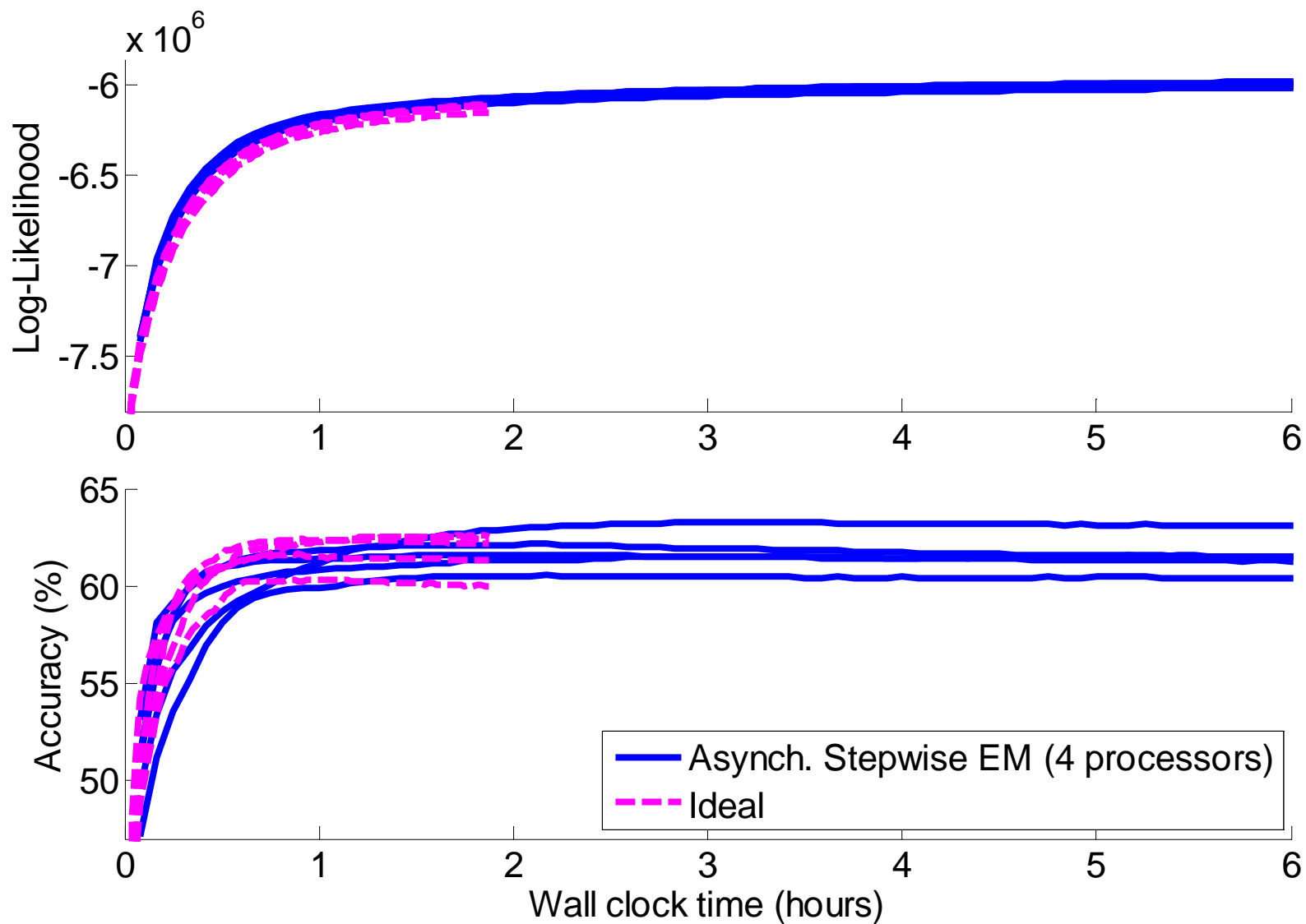
Carnegie Mellon

Part-of-Speech Tagging Results



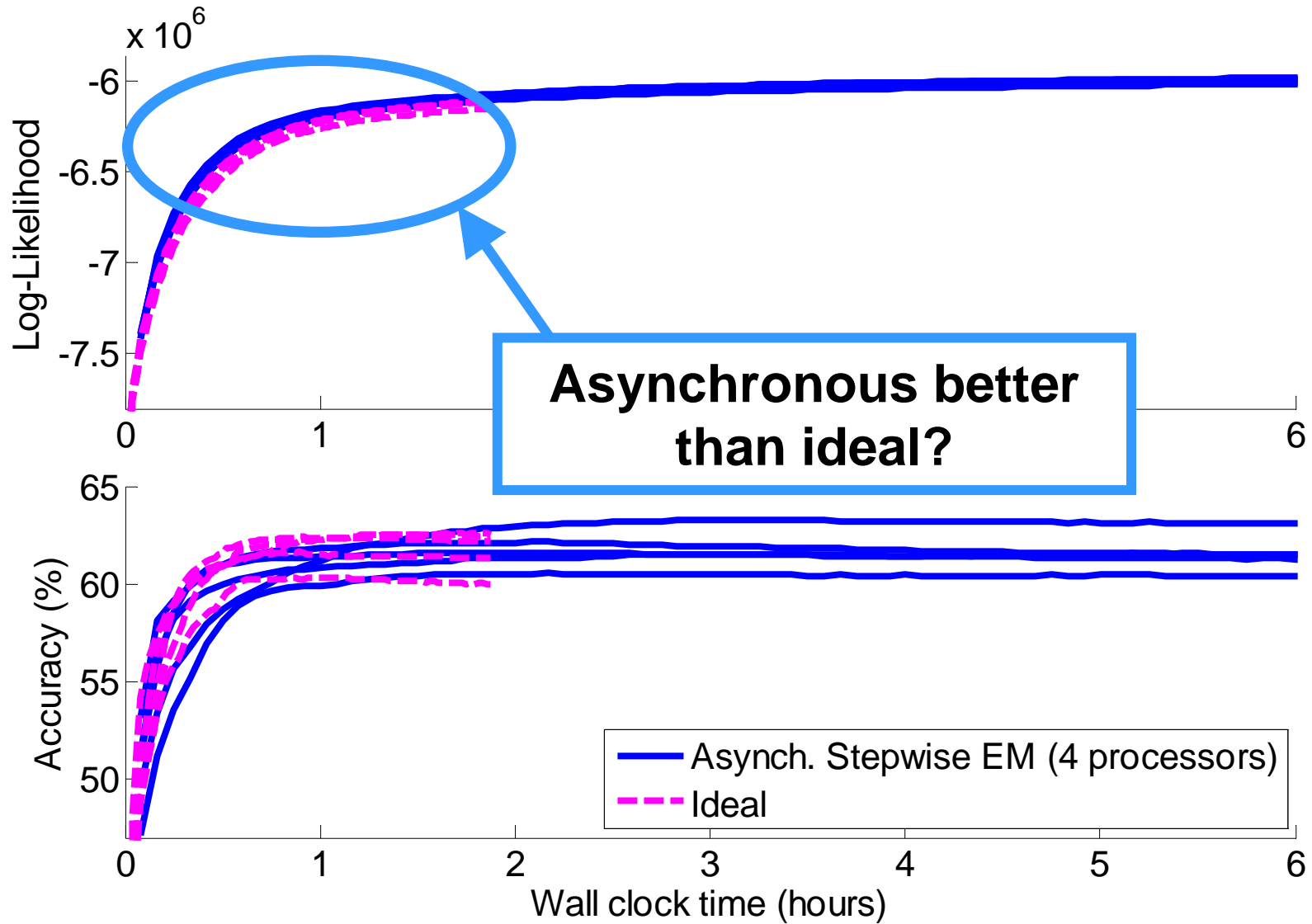
Carnegie Mellon

Comparison with Ideal



Carnegie Mellon

Comparison with Ideal



Carnegie Mellon

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 θ)



Carnegie Mellon

Thanks!

