

Distributed Deep Learning (part 1)

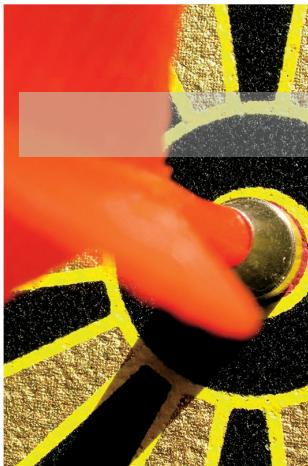
Joseph E. Gonzalez
Co-director of the RISE Lab
jegonzal@cs.berkeley.edu

What is the Problem Being Solved?

- Training models is time consuming
 - Convergence can be slow
 - Training is computationally intensive
- Not all models fit in single machine or GPU memory
- Less of a problem: big data
 - Problem for data preparation / management
 - Not a problem for training ... why?

On Dataset Size and Learning

- Data is a resource! (e.g., like processors and memory)
 - Is having lots of processors a problem?
- You don't have to use all the data!
 - Though using more data can often help
- More data often* dominates models and algorithms



EXPERT OPINION
Contact Editor: Brian Brannon, bbrannon@computer.org

The Unreasonable Effectiveness of Data

Alon Halevy, Peter Norvig, and Fernando Pereira, Google

*More data also enables more sophisticated.

What are the Metrics of Success?

- **Marketing Team:** Maximize number of GPUs/CPUs used
 - A bad metric ... why?
- **Machine Learning:** Minimize passes through the training data
 - Easy to measure, but not informative ... why?
- **Systems:** minimize time to complete a pass through the training data
 - Easy to measure, but not informative ... why?

Ideal Metric of Success

$$\frac{\text{“Learning”}}{\text{Second}} = \left(\frac{\text{“Learning”}}{\text{Record}} \times \frac{\text{Record}}{\text{Second}} \right)$$

How do we measure
Learning?

Convergence
Machine Learning
Property

Throughput
System
Property

Training and Validation



*If you are making modeling decisions based on this then it should be called validation error.

Metrics of Success

- Minimize training time to “best model”
 - Best model measured in terms of test error
- Other Concerns?
 - **Complexity:** Does the approach introduce additional training complexity (e.g., hyper-parameters)
 - **Stability:** How consistently does the system train the model?

Two papers

NIPS 2012 (Same Year as AlexNet)

Large Scale Distributed Deep Networks

Jeffrey Dean, Greg S. Corrado, Rajat Monga, Kai Chen,
Matthieu Devin, Quoc V. Le, Mark Z. Mao, Marc'Aurelio Ranzato,
Andrew Senior, Paul Tucker, Ké Yang, Andrew Y. Ng
`{jeff, gcorrado}@google.com`
Google Inc., Mountain View, CA

Abstract

Recent work in unsupervised feature learning and deep learning has shown that being able to train large models can dramatically improve performance. In this paper, we consider the problem of training a deep network with billions of parameters using tens of thousands of CPU cores. We have developed a software framework called *DistBelief* that can utilize computing clusters with thousands of machines to train large models. Within this framework, we have developed two algorithms for large-scale distributed training: (i) Downpour SGD, an asynchronous stochastic gradient descent procedure supporting a large number of model replicas, and (ii) Sandblaster, a framework that supports a variety of distributed batch optimization procedures, including a distributed implementation of L-BFGS. Downpour SGD and Sandblaster L-BFGS both increase the scale and speed of deep network training. We have successfully used our system to train a deep network 30x larger than previously reported in the literature, and achieves state-of-the-art performance on ImageNet, a visual object recognition task with 16 million images and 21k categories. We show that these same techniques dramatically accelerate the training of a more modestly-sized deep network for a commercial speech recognition service. Although we focus on and report performance of these methods as applied to training large neural networks, the underlying algorithms are applicable to any gradient-based machine learning algorithm.

1 Introduction

Deep learning and unsupervised feature learning have shown great promise in many practical applications. State-of-the-art performance has been reported in several domains, ranging from speech recognition [1, 2], visual object recognition [3, 4], to text processing [5, 6].

It has also been observed that increasing the scale of deep learning, with respect to the number of training examples, the number of model parameters, or both, can drastically improve ultimate

2018 (Unpublished on Arxiv)

Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour

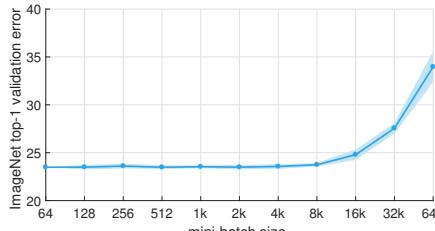
Priya Goyal Lukasz Wesolowski Piotr Dollár Ross Girshick Pieter Noordhuis
Aapo Kyrola Andrew Tulloch Yangqing Jia Kaiming He

Facebook

Abstract

Deep learning thrives with large neural networks and large datasets. However, larger networks and larger datasets result in longer training times that impede research and development progress. Distributed synchronous SGD offers a potential solution to this problem by dividing SGD minibatches over a pool of parallel workers. Yet to make this scheme efficient, the per-worker workload must be large, which implies nontrivial growth in the SGD minibatch size. In this paper, we empirically show that on the ImageNet dataset large minibatches cause optimization difficulties, but when these are addressed the trained networks exhibit good generalization. Specifically, we show no loss of accuracy when training with large minibatch sizes up to 8192 images. To achieve this result, we adopt a hyper-parameter-free linear scaling rule for adjusting learning rates as a function of minibatch size and develop a new warmup scheme that overcomes optimization challenges early in training. With these simple techniques, our Caffe2-based system trains ResNet-50 with a minibatch size of 8192 on 256 GPUs in one hour, while matching small minibatch accuracy. Using commodity hardware, our implementation achieves ~90% scaling efficiency when moving from 8 to 256 GPUs. Our findings enable training visual recognition models on internet-scale data with high efficiency.

Xiv:1706.02677v2 [cs.CV] 30 Apr 2018



minibatch size	ImageNet top-1 validation error
64	~22%
128	~22%
256	~22%
512	~22%
1k	~22%
2k	~22%
4k	~22%
8k	~22%
16k	~23%
32k	~25%
64k	~35%

Figure 1. ImageNet top-1 validation error vs. minibatch size. Error range of plus/minus two standard deviations is shown. We present a simple and general technique for scaling distributed synchronous SGD to minibatches of up to 8k images while maintaining the top-1 error of small minibatch training. For all minibatch sizes we set the learning rate as a linear function of the minibatch size and apply a simple warmup phase for the first few epochs of training. All other hyper-parameters are kept fixed. Using this simple approach, accuracy of our models is invariant to minibatch size (up to an 8k minibatch size). Our techniques enable a linear reduction in training time with ~90% efficiency as we scale to large minibatch sizes, allowing us to train an accurate 8k minibatch ResNet-50 model in 1 hour on 256 GPUs.

tation [8, 10, 28]. Moreover, this pattern generalizes: larger datasets and neural network architectures consistently yield

Large Scale Distributed Deep Networks

Described the system for the 2012 ICML Paper

Building High-level Features Using Large Scale Unsupervised Learning

Quoc V. Le
Marc'Aurelio Ranzato
Rajat Monga
Matthieu Devin
Kai Chen
Greg S. Corrado
Jeff Dean
Andrew Y. Ng

Abstract

We consider the problem of learning high-level, class-specific features from unlabeled data, without any labeled data. For example, it is possible to learn a face detector from a large collection of unlabeled images using unsupervised learning. To answer this, we train a deep neural network with a sparse autoencoder architecture. It consists of two layers of connected sparse autoencoders, with local contrast normalization [1]. The first layer has 1 million connections, the dataset has 10 million

but current experimental evidence suggests the possibility that some neurons in the temporal cortex are



NIPS 2012 (Same Year as AlexNet)

Large Scale Distributed Deep Networks

Jeffrey Dean, Greg S. Corrado, Rajat Monga, Kai Chen,
Matthieu Devin, Quoc V. Le, Mark Z. Mao, Marc'Aurelio Ranzato,
Andrew Senior, Paul Tucker, Ke Yang, Andrew Y. Ng
{jeff, gcorrado}@google.com
Google Inc., Mountain View, CA

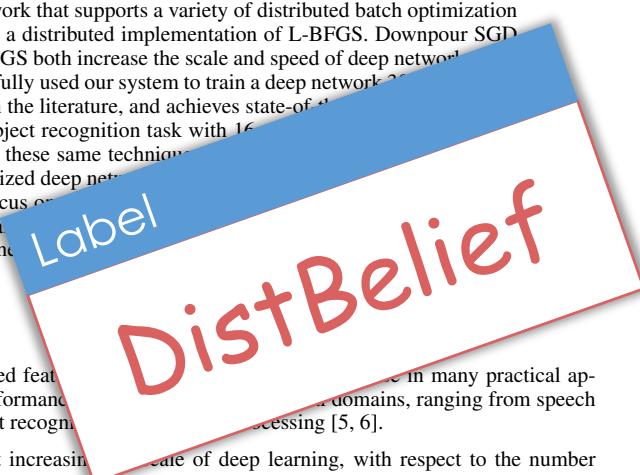
Abstract

Recent work in unsupervised feature learning and deep learning has shown that being able to train large models can dramatically improve performance. In this paper, we consider the problem of training a deep network with billions of parameters using tens of thousands of CPU cores. We have developed a software framework called *DistBelief* that can utilize computing clusters with thousands of machines to train large models. Within this framework, we have developed two algorithms for large-scale distributed training: (i) Downpour SGD, an asynchronous stochastic gradient descent procedure supporting a large number of model replicas, and (ii) Sandblaster, a framework that supports a variety of distributed batch optimization procedures, including a distributed implementation of L-BFGS. Downpour SGD and Sandblaster L-BFGS both increase the scale and speed of deep network training. We have successfully used our system to train a deep network 20 times larger than previously reported in the literature, and achieves state-of-the-art performance on ImageNet, a visual object recognition task with 16 million images and 1000 categories. We show that these same techniques can be applied to training a deep network of a more modestly-sized deep neural network with 100 million parameters. Although we focus on training large neural networks, our framework also applies to training large neural networks using gradient-based machine learning algorithms.

1 Introduction

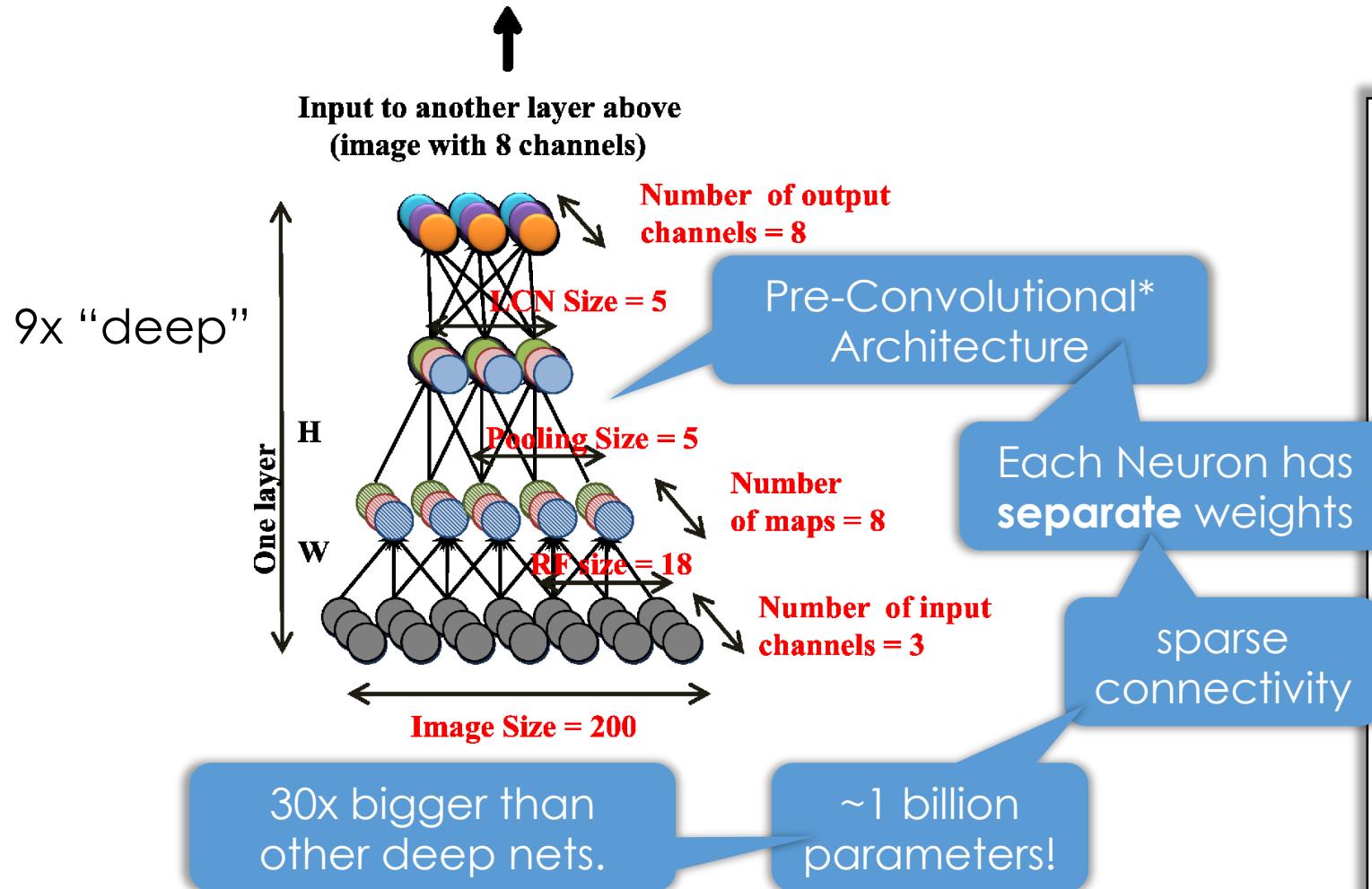
Deep learning and unsupervised feature learning have become increasingly popular in many practical applications. State-of-the-art performance has been achieved in a variety of domains, ranging from speech recognition [1, 2], visual object recognition [3, 4] to natural language processing [5, 6].

It has also been observed that increasing the scale of deep learning, with respect to the number of training examples, the number of model parameters, or both, can drastically improve ultimate classification accuracy [3, 4, 7]. These results have led to a surge of interest in scaling up the training and inference algorithms used for these models [8] and in improving applicable optimization procedures [7, 9]. The use of GPUs [1, 2, 3, 8] is a significant advance in recent years that makes the training of such large models tractable. A large limitation of the GPU approach is



Building High-Level Features Using Large Scale Unsupervised Learning

ICML 2012



*This pre-dates AlexNet but is two decades after LeNet.

Building High-level Features
Using Large Scale Unsupervised Learning

Quoc V. Le
Marc'Aurelio Ranzato
Rajat Monga
Mathieu Devin
Kai Chen
Greg S. Corrado
Jeff Dean
Andrew Y. Ng

QUOCLE@CS.STANFORD.EDU
RANZATO@GOOGLE.COM
RAJATMONGA@GOOGLE.COM
MDEVIN@GOOGLE.COM
KAICHEN@GOOGLE.COM
GCORRADO@GOOGLE.COM
JEFF@GOOGLE.COM
ANG@CS.STANFORD.EDU

Abstract

We consider the problem of building high-level, class-specific feature detectors from *unlabeled* images. For instance, we would like to understand if it is possible to build a face detector from only unlabeled images. This approach is inspired by the neuroscientific conjecture that there exist highly class-specific neurons in the human brain, generally and informally known as “grandmother neurons.” The extent of class-specificity of neurons in the brain is an area of active investigation, but current experimental evidence suggests the possibility that some neurons in the temporal cortex are highly selective for object categories such as faces or hands (Desimone et al., 1984), and perhaps even specific people (Quiroga et al., 2005).

1. Introduction

Contemporary computer vision methodology typically emphasizes the role of *labeled* data to obtain these class-specific feature detectors. For example, to build a face detector, one needs a large collection of images labeled as containing faces, often with a bounding box around the face. The need for large labeled sets poses a significant challenge for problems where labeled data are rare. Although approaches that make use of inexpensive unlabeled data are often preferred, they have not been shown to work well for building high-level features.

This work investigates the feasibility of building high-level features from only *unlabeled* data. A positive result in this setting will indicate a significant

More Context (ML circa 2012)

- Focus of existing distributed ML research
 - Convex Optimization (e.g., SVMs, Lasso)
 - Matrix Factorization
 - Graphical Models
- Key systems at the time
 - Map Reduce → not great for iterative computation (why?)
 - Spark really wasn't visible to ML community
 - GraphLab → A truly wonderful system*
 - we worked with Quoc/Andrew to get their model running on GraphLab but wasn't as performant.

*Developed by the speaker.

Key Problems Addressed in DistBelief Paper

Main Problem

- **Speedup training for large models**

Sub Problems

- How to partition models and data
- Variance in worker performance → Stragglers
- Failures in workers → Fault-Tolerance

Crash Course on Stochastic Gradient Descent



The Gradient Descent Algorithm

$\theta^{(0)} \leftarrow$ initial model parameters (random, warm start)

For τ from 1 to convergence:

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_t \left(\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

Learning Rate

Average Gradient of
Over Training Dataset

How do we distribute this computation?

$\theta^{(0)} \leftarrow$ initial model parameters (random, warm start)For τ from 1 to convergence:

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_t \left(\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

Learning Rate

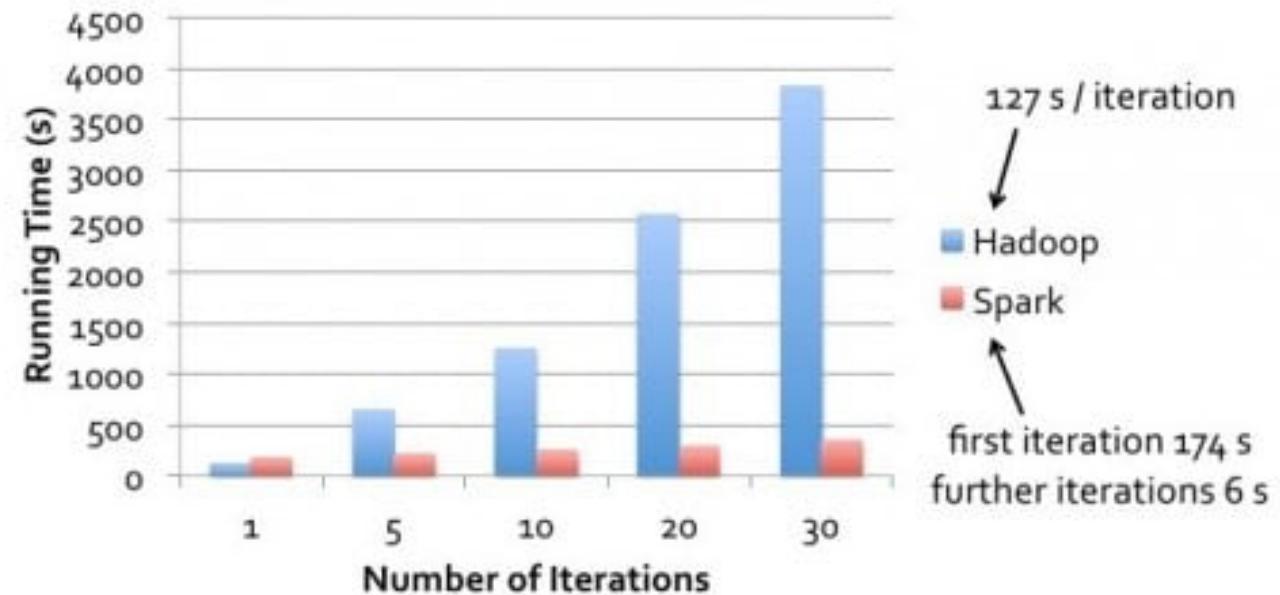
Average Gradient of
Over Training Dataset

How do we distribute this computation?

Data parallelism: divide data across machines, compute local gradient sums and then aggregate across machines. repeat.**Issues?** Repeatedly scanning the data... what if we cache it?



Logistic Regression Performance



also substantially improved the programming API over Hadoop.

The Gradient Descent Algorithm

$\theta^{(0)} \leftarrow$ initial model parameters (random, warm start)

For t from 1 to convergence:

Can we use statistics to improve this algorithm?

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_t \left(\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

Learning Rate

Average Gradient of Over Training Dataset

The empirical gradient is an approximation of what I really want:

$$\left(\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \right) \underset{\theta=\theta^{(t)}}{\approx} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\nabla_{\theta} \mathbf{L}(y, f(x; \theta))]$$

Law of large numbers → more data provides a better approximation (variance in the estimator decreases linearly)

Do I really need to use all the data?

The empirical gradient is an approximation of what I really want:

$$\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \underset{\theta=\theta^{(t)}}{\approx} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\nabla_{\theta} \mathbf{L}(y, f(x; \theta))]$$

Law of large numbers → more data provides a better approximation (variance in the estimator decreases linearly)

$$\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \approx \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta))$$

Random subset of
the data

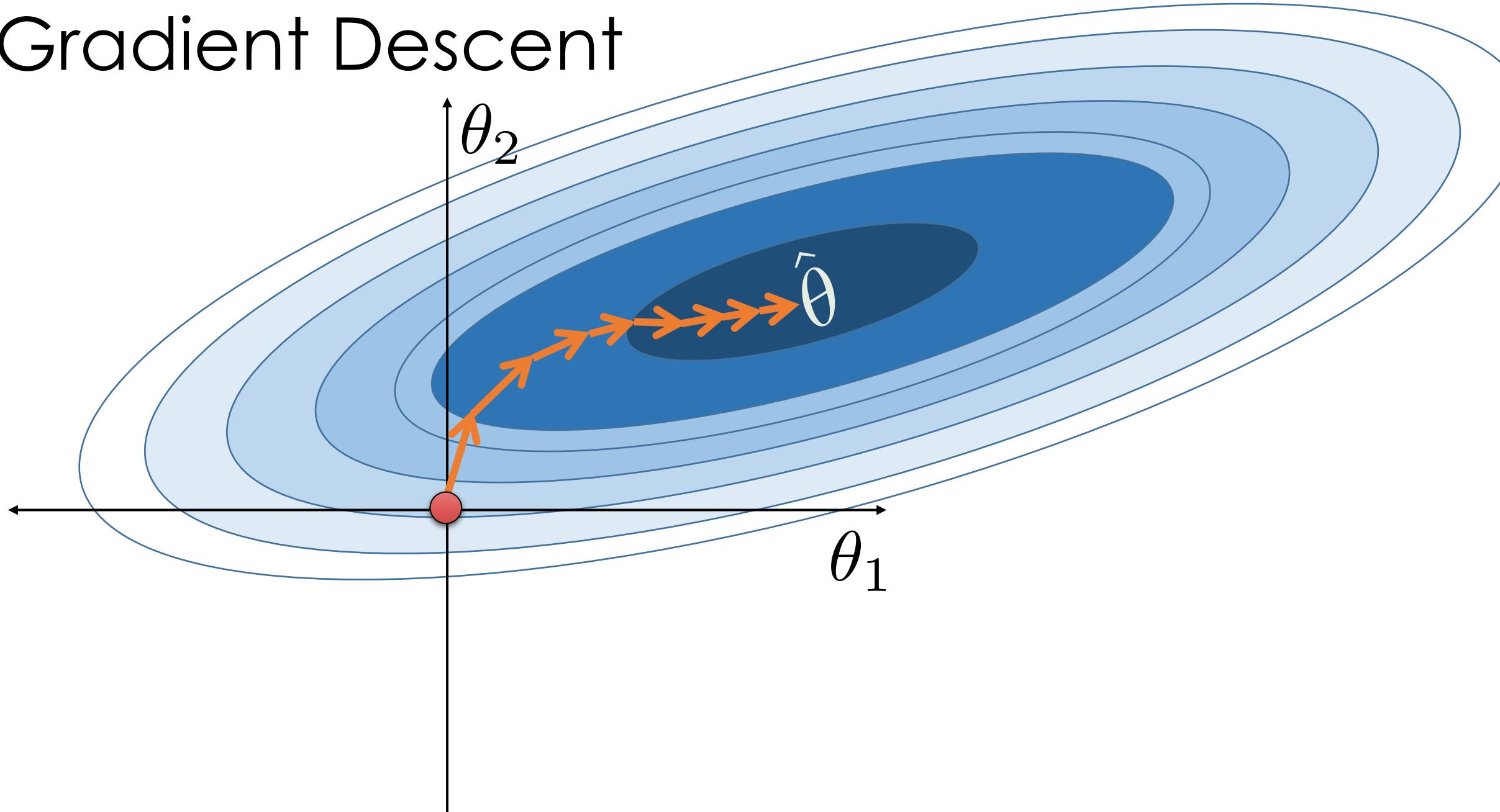
Small B : fast but less accurate
Large B : slower but more accurate

$$\theta^{(0)} \leftarrow \text{initial vector (random, zeros ...)}$$

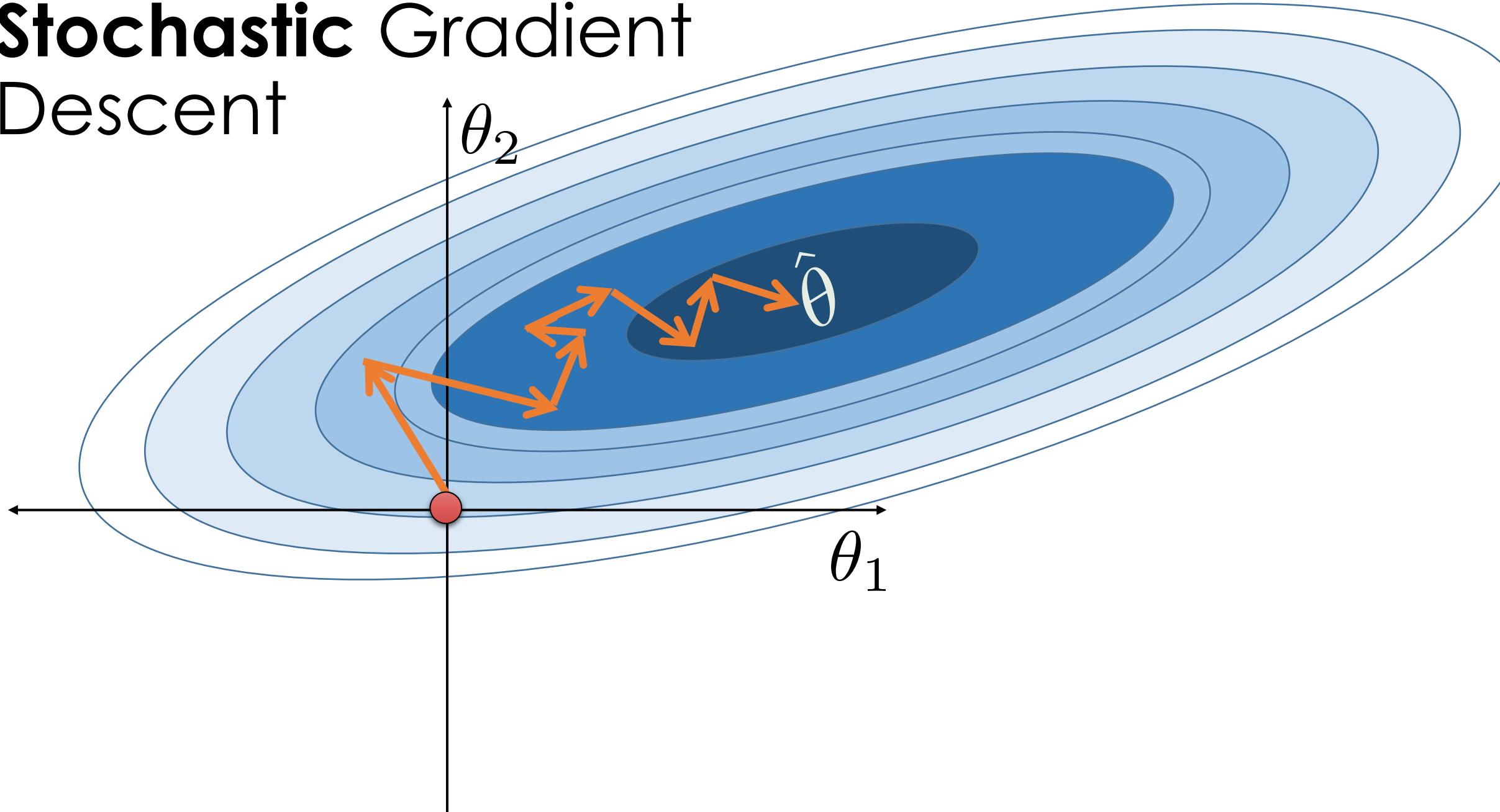
For t from 1 to convergence:

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_t \left(\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

Gradient Descent



Stochastic Gradient Descent



How do you distribute SGD?

Stochastic Gradient Descent

$\theta^{(0)} \leftarrow$ initial vector (random, zeros ...)

For t from 0 to convergence:

$\mathcal{B} \sim$ Random subset of indices

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_t \left(\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

Model Parallelism

Speed up Gradient.
Depends on Model

**Data
Parallelism**

Speed up Sum.
Depends on size of B

(End digression ... for now)

Key Innovations in Large Scale Distributed Deep Networks

NIPS 2012 (Same Year as AlexNet)

Large Scale Distributed Deep Networks

Jeffrey Dean, Greg S. Corrado, Rajat Monga, Kai Chen,
Matthieu Devin, Quoc V. Le, Mark Z. Mao, Marc'Aurelio Ranzato,
Andrew Senior, Paul Tucker, Ke Yang, Andrew Y. Ng
`{jeff, gcorrado}@google.com`
Google Inc., Mountain View, CA

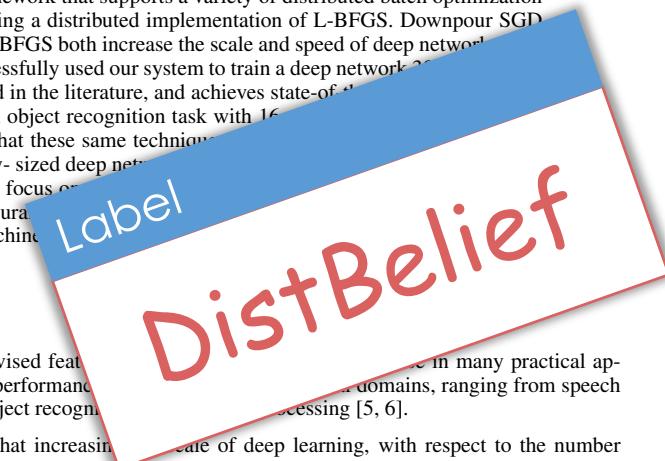
Abstract

Recent work in unsupervised feature learning and deep learning has shown that being able to train large models can dramatically improve performance. In this paper, we consider the problem of training a deep network with billions of parameters using tens of thousands of CPU cores. We have developed a software framework called *DistBelief* that can utilize computing clusters with thousands of machines to train large models. Within this framework, we have developed two algorithms for large-scale distributed training: (i) Downpour SGD, an asynchronous stochastic gradient descent procedure supporting a large number of model replicas, and (ii) Sandblaster, a framework that supports a variety of distributed batch optimization procedures, including a distributed implementation of L-BFGS. Downpour SGD and Sandblaster L-BFGS both increase the scale and speed of deep networking. We have successfully used our system to train a deep network 20 times larger than previously reported in the literature, and achieves state-of-the-art performance on ImageNet, a visual object recognition task with 16 million training images and 1000 categories. We show that these same techniques can be applied to training a deep network of a more modestly-sized deep neural network with only 100 million parameters. Although we focus on distributed training of deep networks, our framework is also applicable to training large neural networks using gradient-based machine learning.

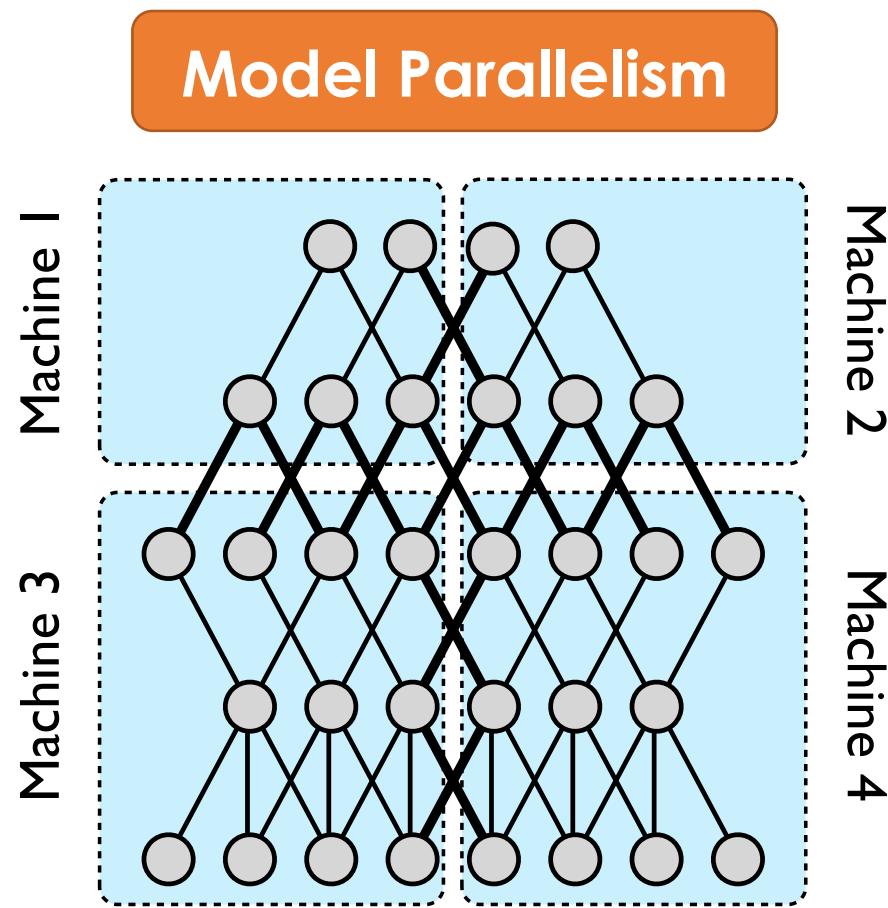
1 Introduction

Deep learning and unsupervised feature learning have had a significant impact in many practical applications. State-of-the-art performance has been achieved in a variety of domains, ranging from speech recognition [1, 2], visual object recognition [3, 4], and natural language processing [5, 6].

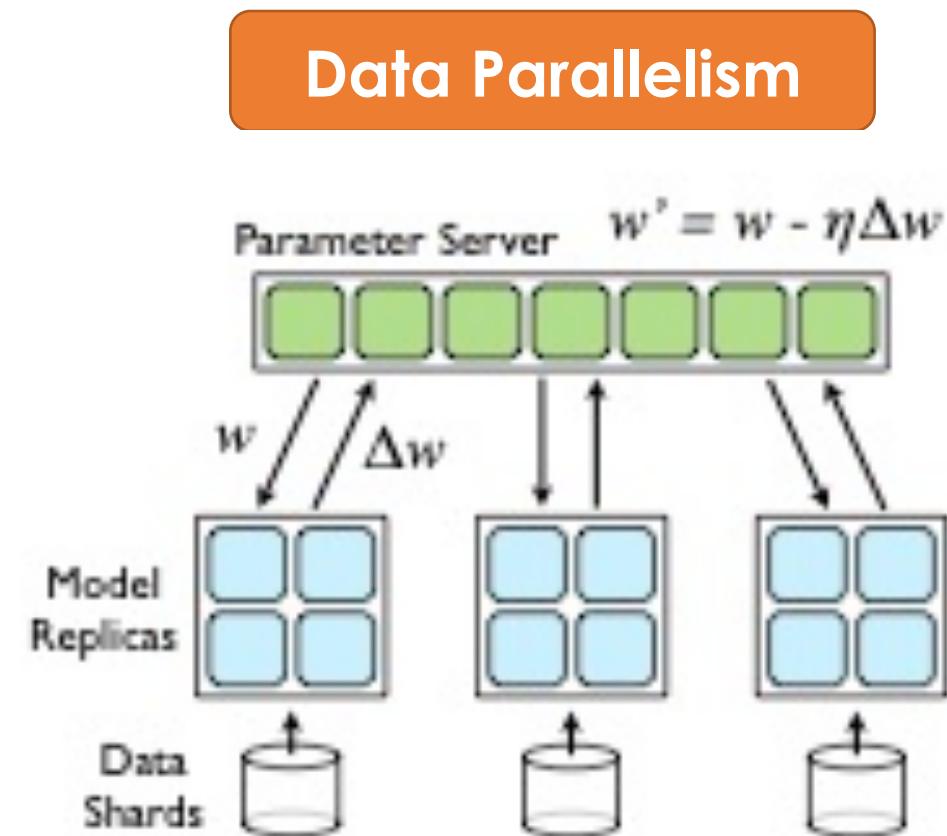
It has also been observed that increasing the scale of deep learning, with respect to the number of training examples, the number of model parameters, or both, can drastically improve ultimate classification accuracy [3, 4, 7]. These results have led to a surge of interest in scaling up the training and inference algorithms used for these models [8] and in improving applicable optimization procedures [7, 9]. The use of GPUs [1, 2, 3, 8] is a significant advance in recent years that makes the training of such large models feasible. Advances in the utilization of the GPU, in conjunction with



Combine Model and Data Parallelism



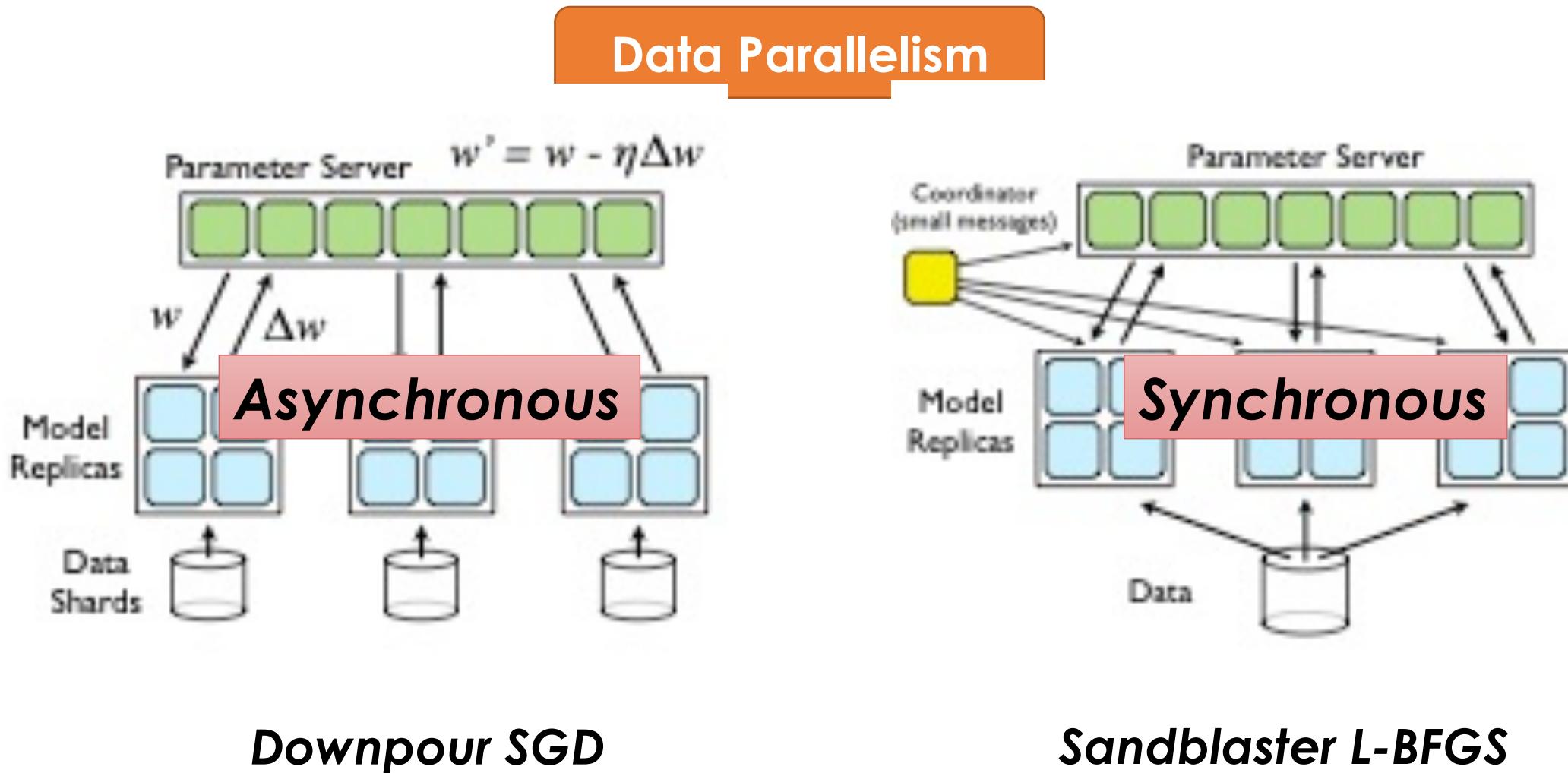
This appears in earlier work on graph systems ...



Downpour SGD

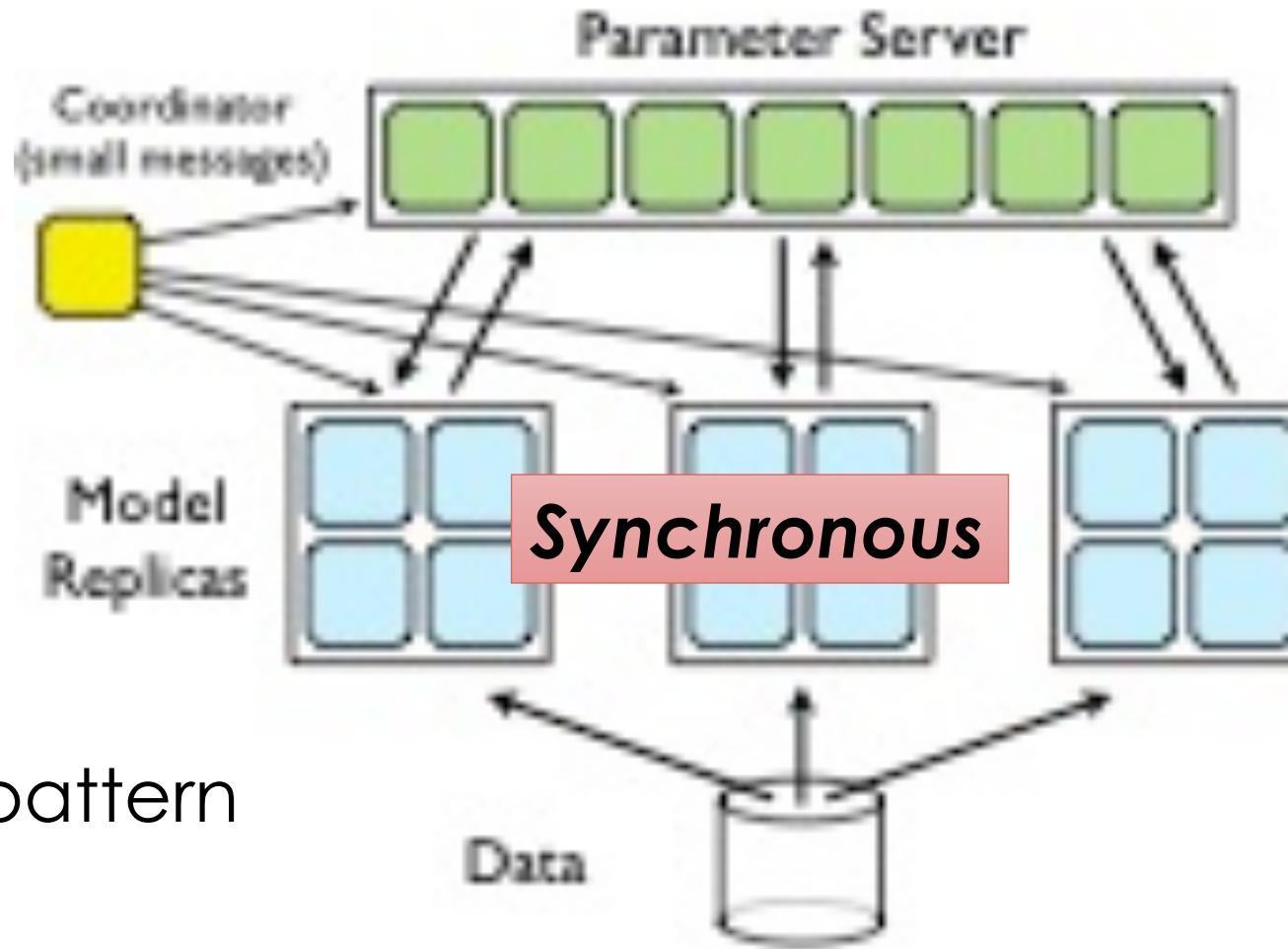
Combine Model and Data Parallelism

Machine 2 Machine 4



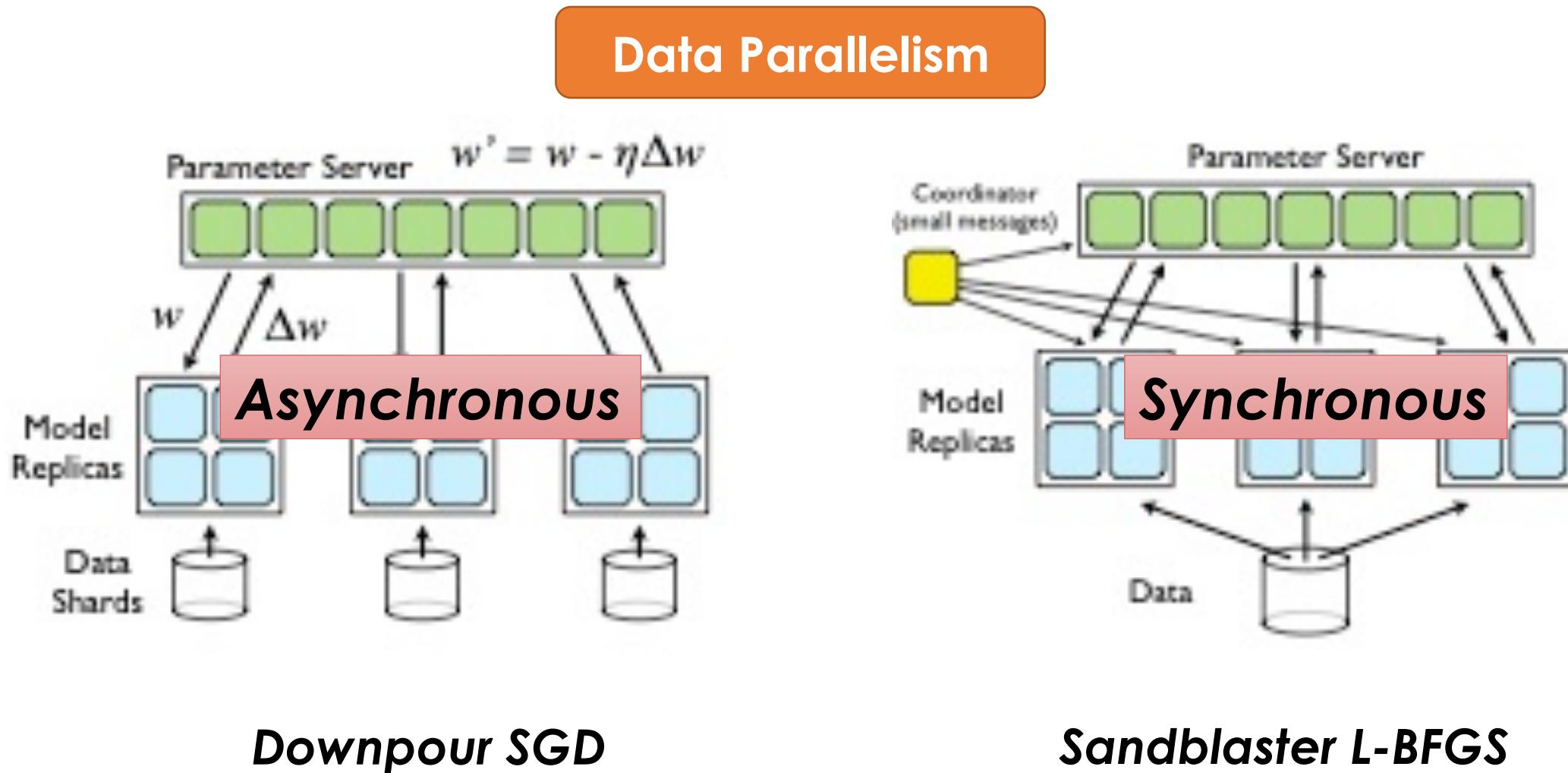
Sandblaster L-BFGS

- L-BFGS
 - Commonly used for convex opt. problems
 - Requires repeated scans of all data
 - Robust, minimal tuning
- Naturally fits map-reduce pattern
- **Innovations:**
 - accumulate gradients and store outputs in a sharded key value store (parameter server)
 - Tiny tasks + backup tasks to mitigate stragglers



Combine Model and Data Parallelism

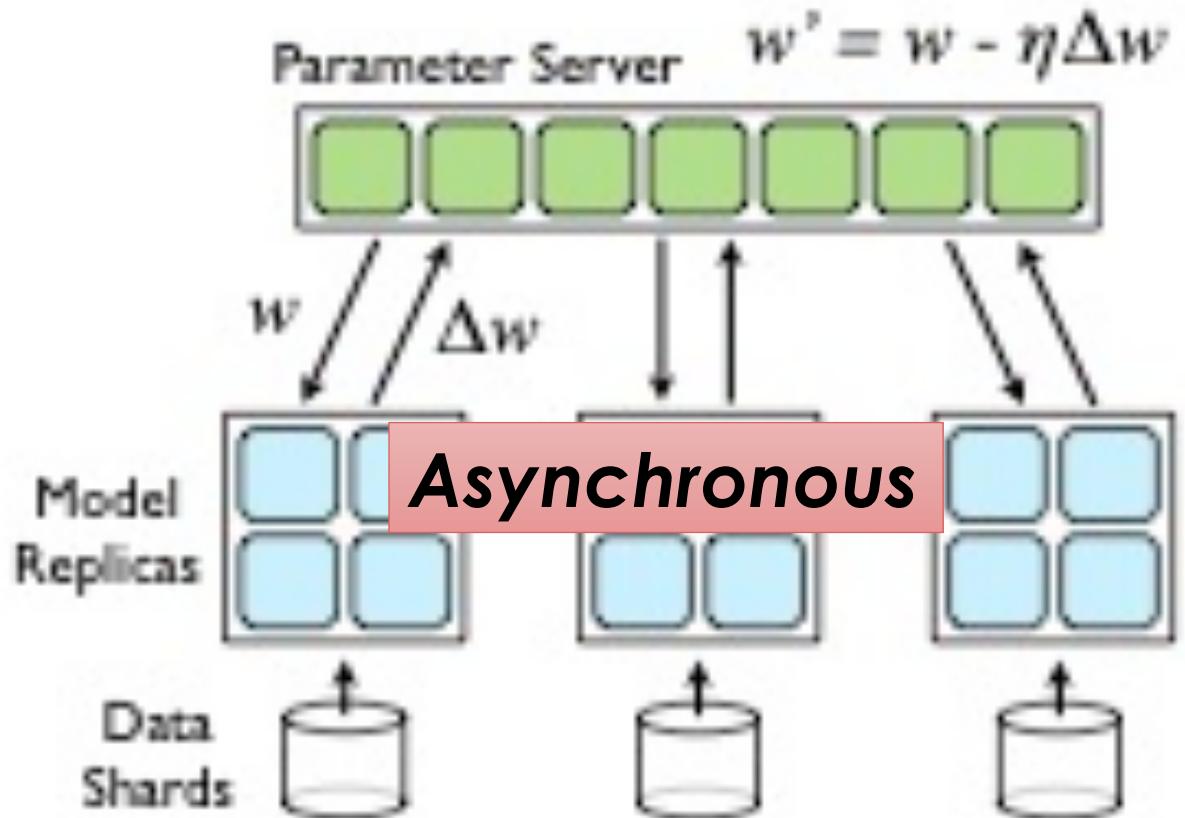
Machine 2 Machine 4



Downpour SGD

Claimed Innovations

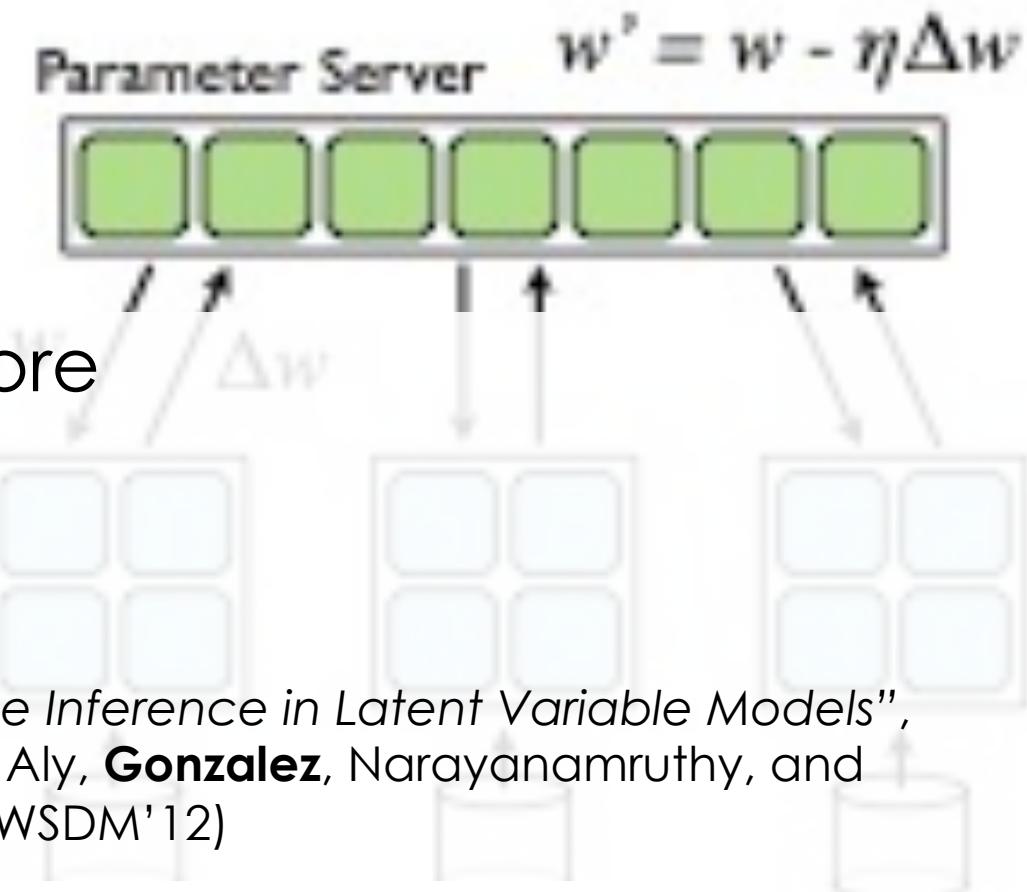
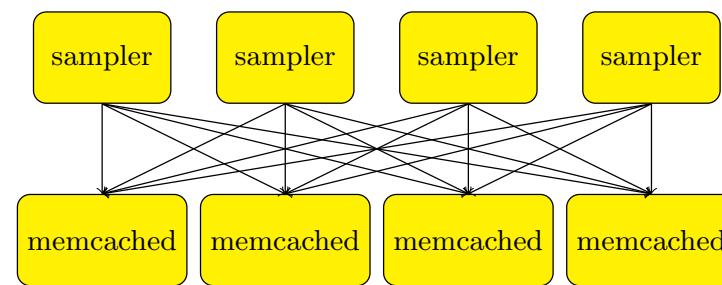
- Parameter Server
- Combine model and data parallelism in an async. execution.
- Adagrad stabilization
- Warmstarting



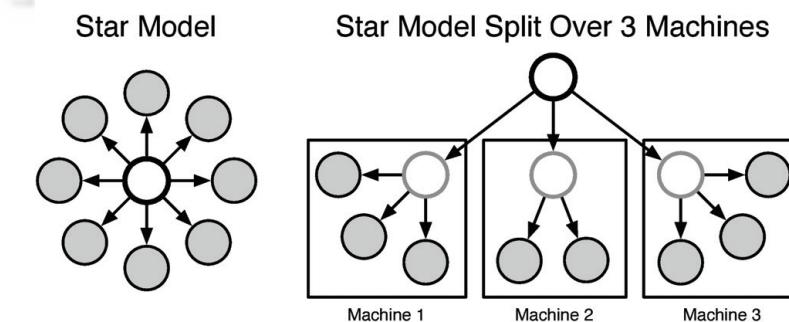
Parameter Servers

- Essentially a **sharded** key-value store
 - support for put, get, **add**
- Idea appears in earlier papers:

"An Architecture for Parallel Topic Models", Smola and Narayananamruthy.
(VLDB'10)



"Scalable Inference in Latent Variable Models", Ahmed, Aly, **Gonzalez**, Narayananamruthy, and Smola. (WSDM'12)

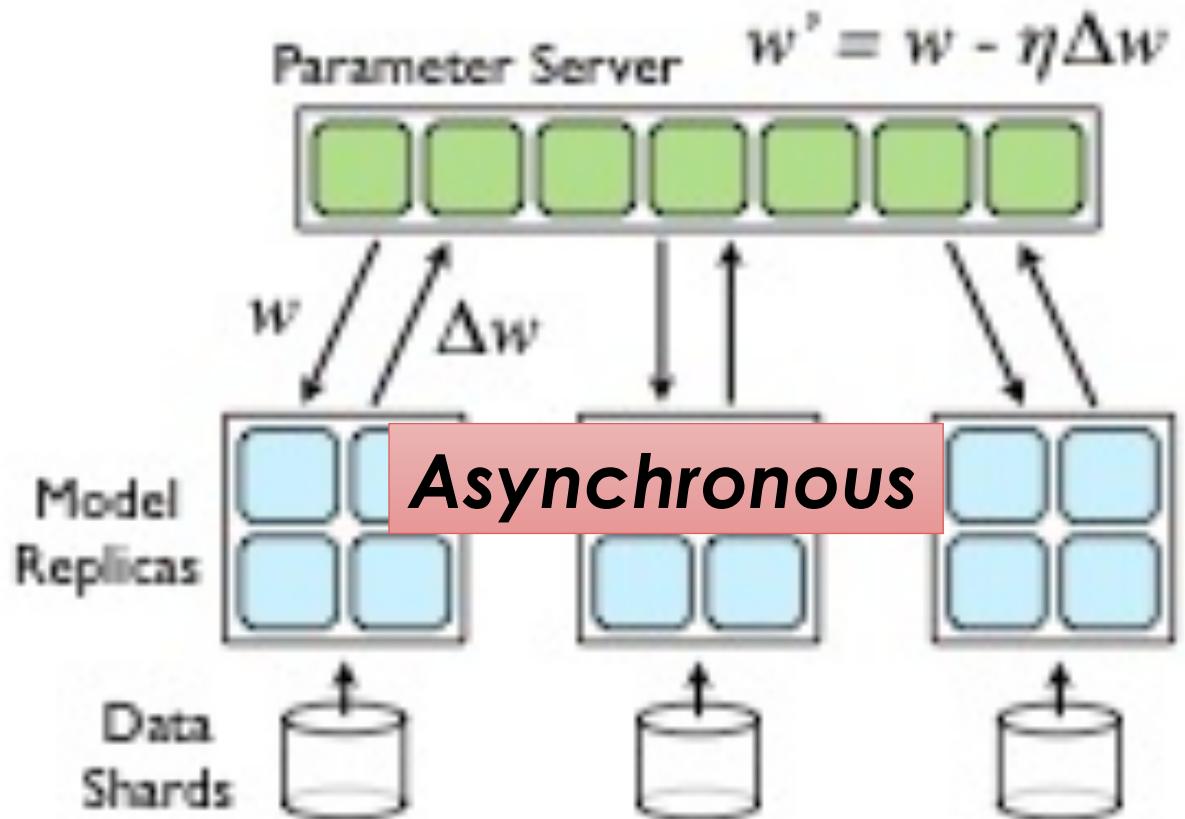


DistBelief was probably the first paper to call a sharded key-value store a Parameter Server.

Downpour SGD

Claimed Innovations

- Parameter Server
- Combine model and data parallelism in an **async. execution.**
- Adagrad stabilization
- Warmstarting



Adagrad Stabilization

- Address *large variability in magnitudes of gradients*
 - Rescale the gradients by an estimate of the diagonal variance
 - More recently superseded by Adam
- Stability is needed here to support asynchronous gradient updates

Warmstarting

- Recall

Starting closer to a solution can help!

$$\theta^{(0)} \leftarrow \text{initial vector (random, zeros ...)}$$

For t from 0 to convergence:

$\mathcal{B} \sim$ Random subset of indices

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_t \left(\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

Initial Value

Solution path for stochastic gradient on a 2d logistic regression problem.

Steep Gradient

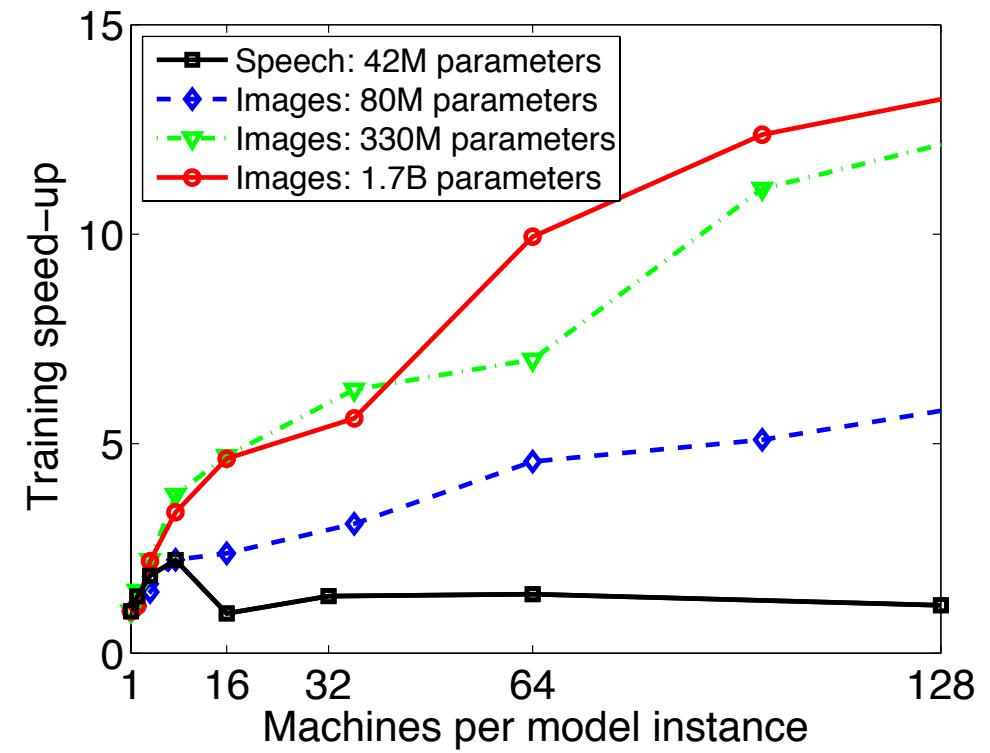
Closer to the solution gradient is smoother.



Key Results

Model Parallelism

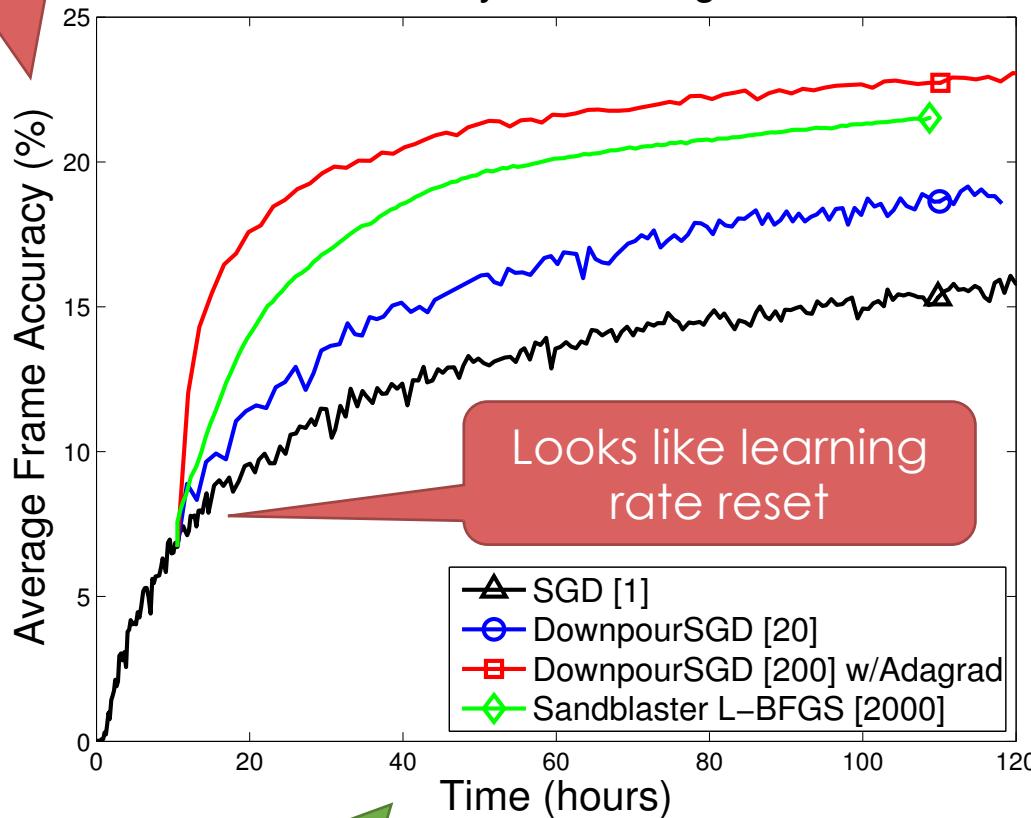
- Measured speedup to compute a single mini-batch
 - Is this a good metric?
- Results are not that strong...



Key Results: Training and Test Error

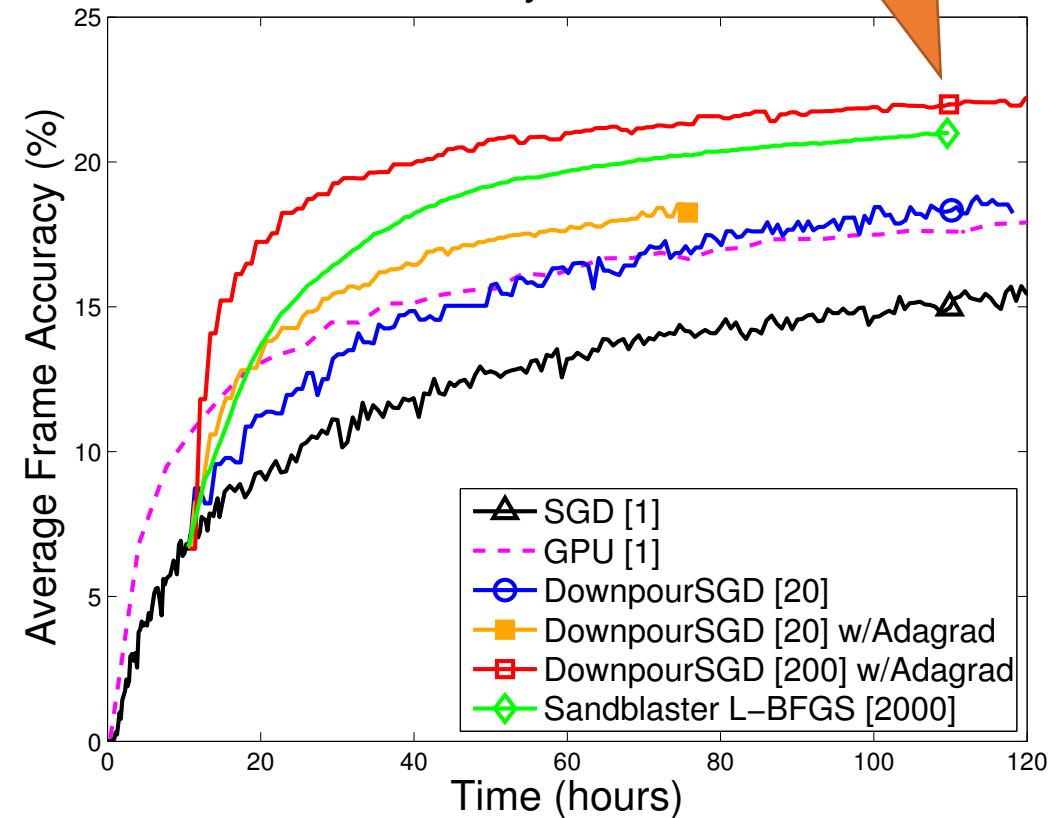
Weird 20K
Error Metric

Accuracy on Training Set



Which would
you use?

Accuracy on Test Set



Why are they in the NY Times

- Trained a 1.7 billion parameter model (30x larger than state-of-the-art) (was it necessary?)
- Using 16,000 cores (efficiently?)
- Achieves 15.8 accuracy on ImageNet 20K (70% improvement over state of the art).
 - Non-standard benchmark
- Qualitatively interesting results



Figure 6. Visualization of the cat face neuron (left) and human body neuron (right).

Long-term Impact

- The **parameter server** appears in many later machine learning systems
- Downpour (**Asynchronous**) SGD has been largely **replaced by synchronous systems** for supervised training
 - Asynchrony is still popular in RL research
 - Why?
- Model parallelism is still used for large language models
 - Predated this work
- The neural network architectures studied here have been largely replaced by convolutional networks

Second Paper

- Generated a lot of press
 - Recently (Aug) surpassed by Fast.ai: “Now anyone can train ImageNet in 18 minutes for \$40.” blog post
- Popularized linear learning rate scaling

2018 (Unpublished on Arxiv)

**Accurate, Large Minibatch SGD:
Training ImageNet in 1 Hour**

Priya Goyal Piotr Dollár Ross Girshick Pieter Noordhuis
Lukasz Wesolowski Aapo Kyrola Andrew Tulloch Yangqing Jia Kaiming He

Facebook

Abstract

Deep learning thrives with large neural networks and large datasets. However, larger networks and larger datasets result in longer training times that impede research and development progress. Distributed synchronous SGD offers a potential solution to this problem by dividing SGD minibatches over a pool of parallel workers. Yet to make this scheme efficient, the per-worker workload must be large, which implies nontrivial growth in the SGD minibatch size. In this paper, we empirically show that on the ImageNet dataset large minibatches cause optimization difficulties, but when these are addressed the trained networks exhibit good generalization. Specifically, we show no loss of accuracy when training with large minibatch sizes up to 8192 images. To achieve this result, we adopt a hyper-parameter-free linear scaling rule for adjusting learning rates as a function of minibatch size and develop a new warmup scheme that overcomes optimization challenges early in training. With these simple techniques, our Caffe2-based system trains ResNet-50 with a minibatch size of 8192 on 256 GPUs in one hour, while matching small minibatch accuracy. Using commodity hardware, our implementation achieves ~90% scaling efficiency when moving from 8 to 256 GPUs. Our findings enable training visual recognition models on internet-scale data with high efficiency.

arXiv:1706.02677v2 [cs.CV] 30 Apr 2018

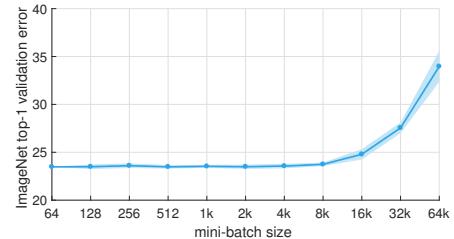


Figure 1. **ImageNet top-1 validation error vs. minibatch size.** Error range of plus/minus two standard deviations is shown. We present a simple and general technique for scaling distributed synchronous SGD to minibatches of up to 8k images while maintaining the top-1 error of small minibatch training. For all minibatch sizes we set the learning rate as a linear function of the minibatch size and apply a simple warmup phase for the first few epochs of training. All other hyper-parameters are kept fixed. Using this simple approach, accuracy of our models is invariant to minibatch size (up to an 8k minibatch size). Our techniques enable a linear reduction in training time with ~90% efficiency as we scale to large minibatch sizes, allowing us to train an accurate 8k minibatch ResNet-50 model in 1 hour on 256 GPUs.

1. Introduction

Scale matters. We are in an unprecedented era in AI research history in which the increasing data and model scale is rapidly improving accuracy in computer vision [22, 41, 34, 35, 36, 16], speech [17, 40], and natural language processing [7, 38]. Take the profound impact in computer vision as an example: visual representations learned by deep convolutional neural networks [23, 22] show excellent performance on previously challenging tasks like ImageNet classification [33] and can be transferred to difficult perception problems such as object detection and segmen-

tation [8, 10, 28]. Moreover, this pattern generalizes: larger datasets and neural network architectures consistently yield improved accuracy across all tasks that benefit from pre-training [22, 41, 34, 35, 36, 16]. But as model and data scale grow, so does training time; discovering the potential and limits of large-scale deep learning requires developing novel techniques to keep training time manageable.

The goal of this report is to demonstrate the feasibility of, and to communicate a practical guide to, large-scale training with distributed *synchronous* stochastic gradient descent (SGD). As an example, we scale ResNet-50 [16] training, originally performed with a minibatch size of 256 images (using 8 Tesla P100 GPUs, training time is 29 hours), to larger minibatches (see Figure 1). In particular, we show that with a large minibatch size of 8192, we can train ResNet-50 in 1 hour using 256 GPUs while maintaining

Contrasting to the first paper

- **Synchronous SGD**
 - Much of the recent work has focused on synchronous setting
 - Easier to reason about
- Focus exclusively on data parallelism: **batch-size scaling**
- Focuses on the **generalization gap problem**

How do you distribute SGD?

Stochastic Gradient Descent

$\theta^{(0)} \leftarrow$ initial vector (random, zeros ...)

For t from 0 to convergence:

$\mathcal{B} \sim$ Random subset of indices

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta_t \left(\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

Data
Parallelism

Slow? (~150ms)
Depending on size of B

Batch Size Scaling

- Increase the batch size by adding machines

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \hat{\eta} \left(\frac{1}{k} \sum_{j=1}^k \frac{1}{|\mathcal{B}_j|} \sum_{i \in \mathcal{B}_j} \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

- Each server processes a fixed batch size (e.g., n=32)
- As more servers are added (k) the effective overall batch size increases linearly
- Why do these additional servers help?

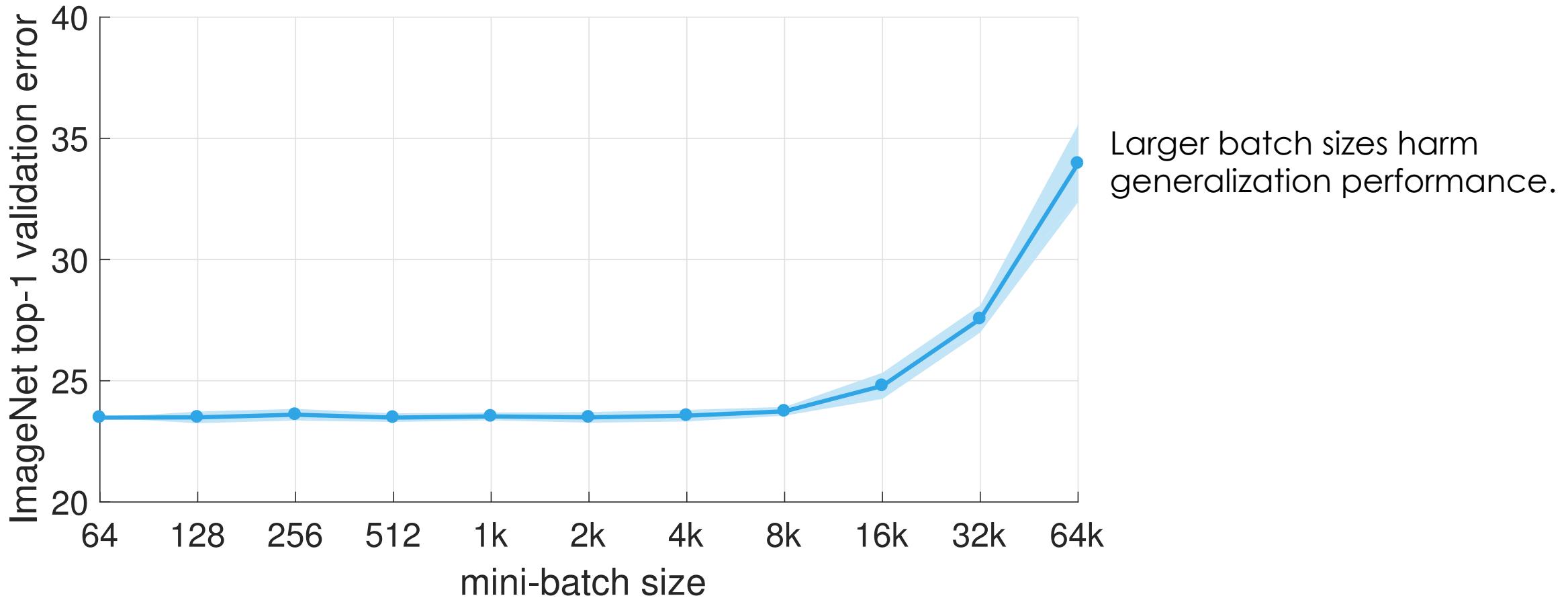
Bigger isn't Always Better

- Motivation for larger batch sizes
 - More opportunities for parallelism → but is it useful?
 - Recall ($1/n$ variance reduction):

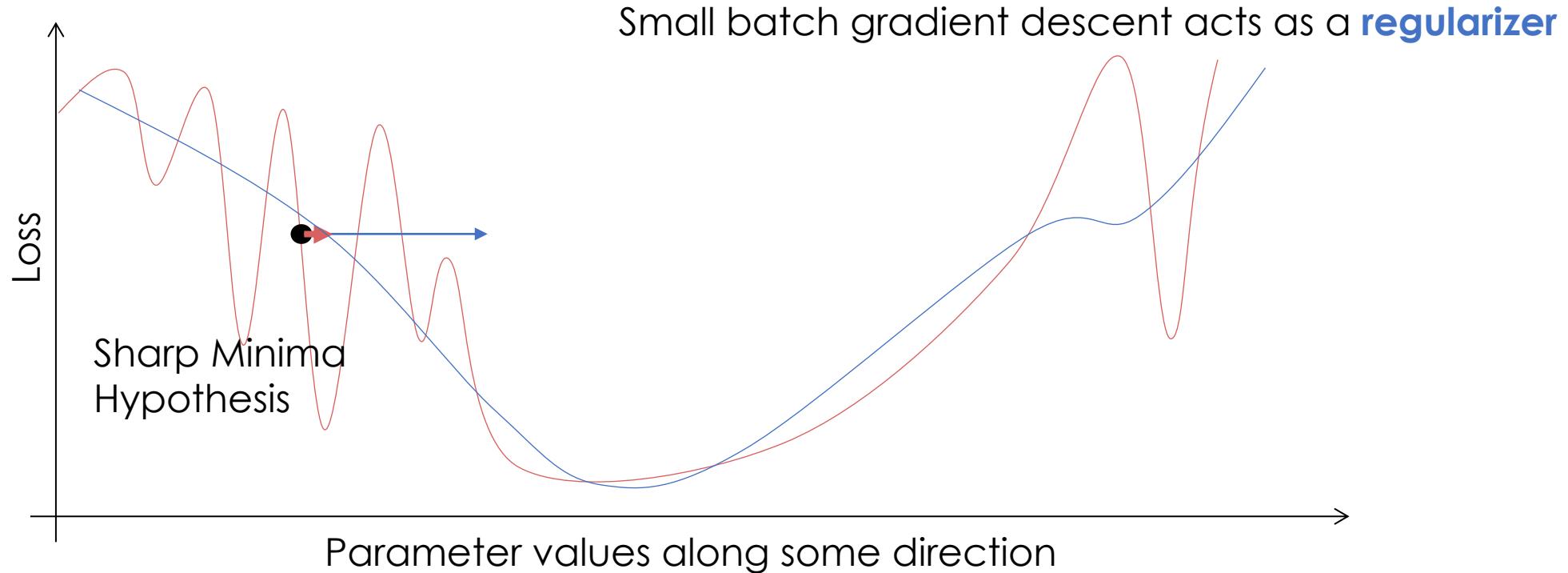
$$\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \approx \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta))$$

- Is a variance reduction helpful?
 - Only if it lets you take bigger steps (move faster)
 - Doesn't affect the final answer...

Generalization Gap Problem



Rough “Intuition”



Key problem: Addressing the generalization gap for large batch sizes.

Solution: Linear Scaling Rule

- Scale the learning rate linearly with the batch size

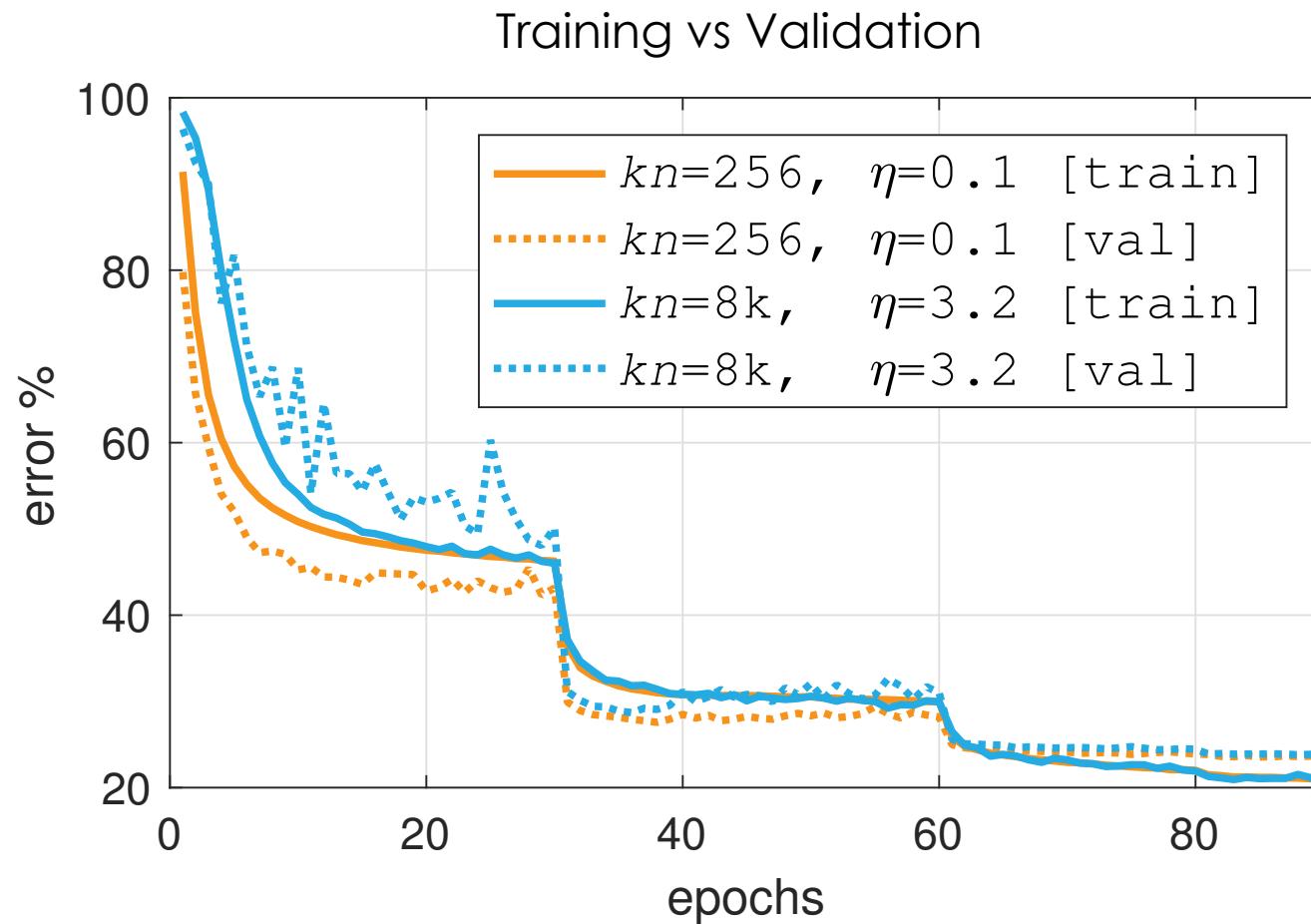
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \hat{\eta} \left(\frac{1}{k} \sum_{j=1}^k \frac{1}{|\mathcal{B}_j|} \sum_{i \in \mathcal{B}_j} \nabla_{\theta} \mathbf{L}(y_i, f(x_i; \theta)) \Big|_{\theta=\theta^{(t)}} \right)$$

- Addresses generalization performance by **taking larger steps** (also improves training convergence)
- **Sub-problem:** Large learning rates can be destabilizing in the beginning. Why?
 - **Gradual warmup solution:** increase learning rate scaling from constant to linear in first few epochs
 - Doesn't help for very large k...

Other Details

- **Independent Batch Norm:** Batch norm calculation applies only to local batch size (n).
- **All-Reduce:** Recursive halving and doubling algorithm
 - Used instead of popular ring reduction (fewer rounds)
- **Gloo** a library for efficient collective communications
- **Big Basin GPU Servers:** Designed for deep learning workloads
 - Analysis of communication requirements → latency bound
- **No discussion on straggler or fault-tolerance**
 - Why?!

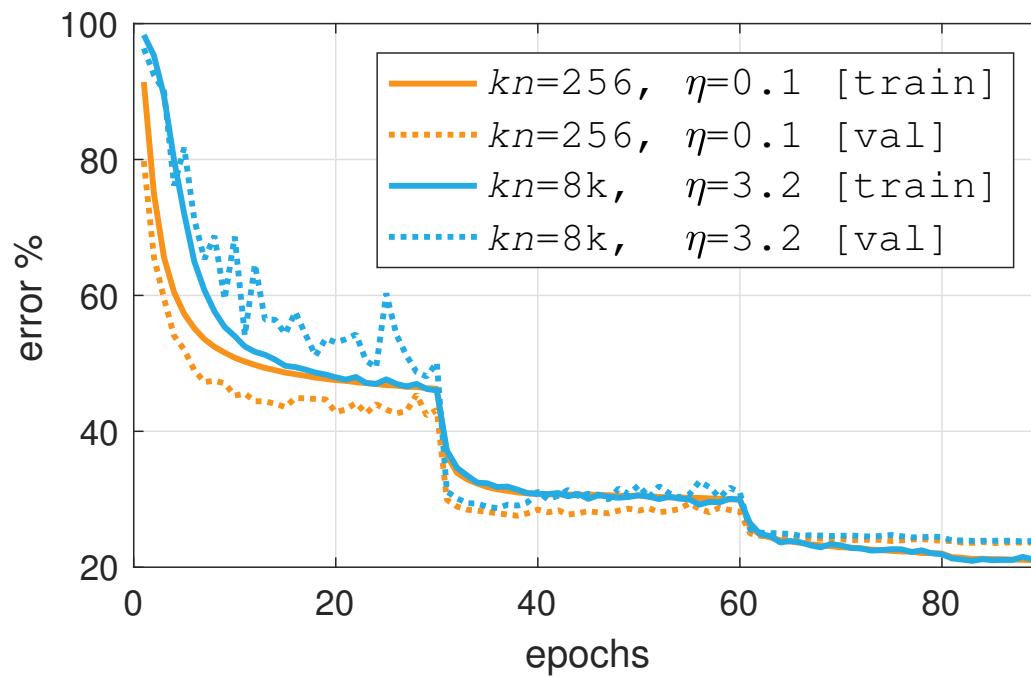
Key Results



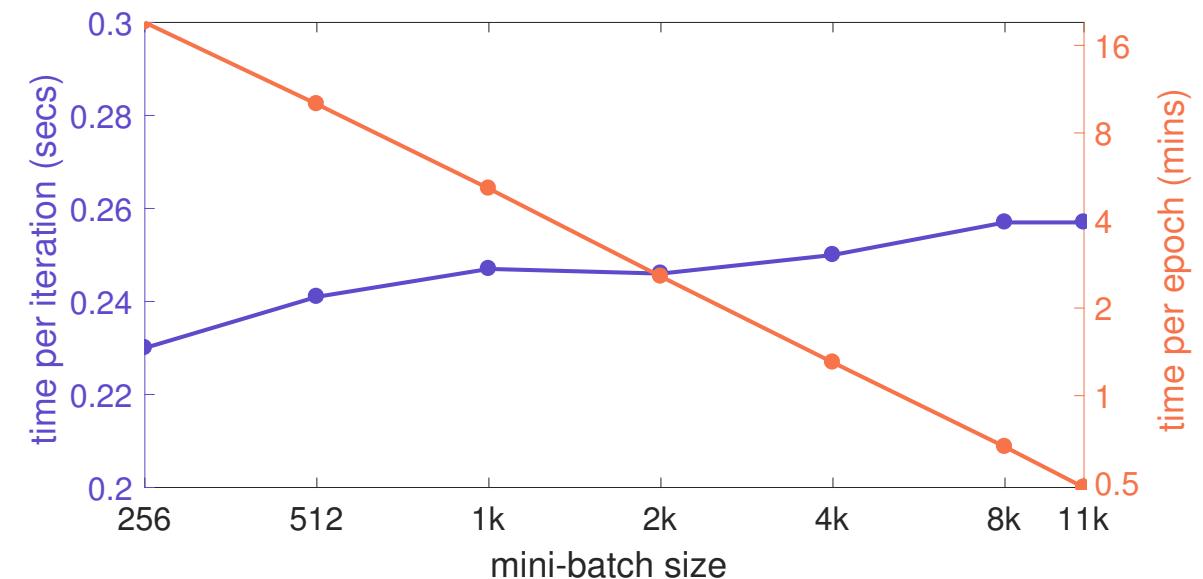
All curves closely match using the linear scaling rule.

Note learning rate schedule drops.

Key Results



“Learning”
Epoch
Machine Learning



Epoch
Second
System

Key Results

- Train ResNet-50 to state-of-the-art on 256 GPUs in 1 hour
 - 90% scaling efficiency
- Fairly careful study of the linear scaling rule
 - Observed limits to linear scaling do not depend on dataset size
 - Cannot scale parallelism with dataset size

Long-term Impact

- Still early (this paper is not yet published)
- Ideas that will appear in other papers
 - Linear rate scaling
 - Independent batch norm
- Paper points to limits of Synchronous SGD → need new methods to accelerate training
 - Eg.: Fast.ai → curriculum learning though scale variation

Questions for Discussion

- Distributed model training is not very common. Why?
- Should we return to asynchrony?
 - What is needed to address issues with asynchronous training?
- How will changes in hardware affect distributed training
 - E.g., faster GPUs → larger batches, faster networks → smaller batches ...
- How will the emergence of “dynamic models” and large “mixture of expert models” affect distributed training?