

ECE826: Part 2

So far we talked about

- Generalization: why would expect our models to work
- SGD/GD: under what conditions do they work
- Why can we optimize neural networks?

Not a lot of algorithmic design principles for
large-scale learning

Advances and Challenges in Distributed Machine Learning

We'll see a lot of principles for scaling up, and designing the “plumbing” of deep learning systems

Overview

- Multicore vs. Distributed
- Algorithms of choice
- Open challenges with Performance Gains/Analysis
- Communication bottlenecks
- Straggler Nodes
- Robustness

Stochastic Gradient Descent

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{z}_i)$$

loss for data point i

- Idea ('50s, '60s [Robbins, Monro], [Widrow, Hoff]):
Sample a data point + locally optimize.

SGD: An *Über*-algorithm

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \gamma \cdot \nabla \ell(\mathbf{w}_k; \mathbf{z}_{i_k})$$

Stochastic Gradient Descent

Different names and flavors

ML / Optimization / Statistics / EE

Perceptron

Incremental Gradient

Back Propagation (NNs)

Oja's iteration (PCA)

LMS Filter

...

Has been around for a while, for good reasons:

- Robust to noise

- Simple to implement

- Near-optimal learning performance *

- Small computational foot-print

Stochastic Gradient Descent

GPT-3 would take 288 years to train on a single Tesla
V100 GPU

[source: <https://arxiv.org/pdf/2104.04473.pdf>]

Goal:

Speed up Machine Learning

Idea:
Train at scale



System Setup

Parallel vs. Distributed

- Parallel (CPU/multicore GPU)
 - Single machine, many cores (usually up to 10-100s)
 - Shared memory (all cores have access to RAM)
 - Comm. to RAM is cheap

- Distributed
 - Many machines (usually up to 100-1000s) connected via network
 - Shared-nothing architecture (each node has its own resources)
 - Communication costs non-negligible

Scaling up vs. Scaling out

Scaling up vs. out

What we'd ideally like

- Cores ∞
- RAM ∞
- Comm. Cost 0
- Cost to build 0

Feasible solutions:

Scaling up:

- Getting the largest machine possible, with maxed out RAM

Scaling out:

- Getting a bunch of machines, and linking them together

Scaling up vs. out

Scaling up

Pros:

RAM comm. cheap (no network)
Less impl. overheads
Less power/smaller footprint

Cons:

100s cores/machine = expensive
Smaller fault tolerance
Limited upgradability

Scaling out

Pros:

Much cheaper (especially on Ec2)
Can replace faulty parts
Better fault tolerance (if it matters)

Cons:

Network bound
Major implementation overheads
Large power footprint

Should I buy or rent?

Price Comparisons for 4 GPUs racks

- Sim

Flops/\$ (1-yr analysis with 50% occupancy)

| | Hyperplane-A100 (25% annual depreciation) | Hyperplane-A100 (100% annual depreciation) | Scalar-A100 (25% annual depreciation) | Scalar-A100 (100% annual depreciation) | p4d (qll upfront) | p4d (partial upfront) | p4d (no upfront) | p4d (on demand) |
|----------------------|---|--|---------------------------------------|--|-------------------|-----------------------|------------------|-----------------|
| Total Cost 1-yr | \$55,534 | \$160,534 | \$47,279 | \$138,779 | \$164,955 | \$168,321 | \$176,737 | \$143,544 |
| Total Petaflops 1-yr | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 |
| Petaflops/\$ | 44.3 | 15.3 | 52.0 | 17.7 | 14.9 | 14.6 | 13.9 | 17.1 |

- Alt



Google Cloud Platform

Price Comparisons for 4 GPUs racks

- Sim

Flops/\$ (1-yr analysis with 50% occupancy)

| | Hyperplane-A100 (25% annual depreciation) | Hyperplane-A100 (100% annual depreciation) | Scalar-A100 (25% annual depreciation) | Scalar-A100 (100% annual depreciation) | p4d (qll upfront) | p4d (partial upfront) | p4d (no upfront) | p4d (on demand) |
|----------------------|---|--|---------------------------------------|--|-------------------|-----------------------|------------------|-----------------|
| Total Cost 1-yr | \$55,534 | \$160,534 | \$47,279 | \$138,779 | \$164,955 | \$168,321 | \$176,737 | \$143,544 |
| Total Petaflops 1-yr | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 | 2,459,808 |
| Petaflops/\$ | 44.3 | 15.3 | 52.0 | 17.7 | 14.9 | 14.6 | 13.9 | 17.1 |

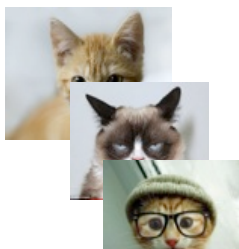
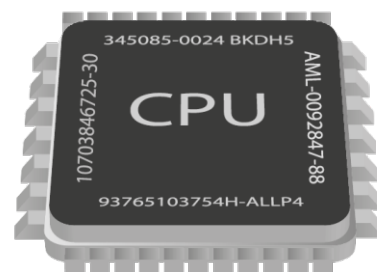
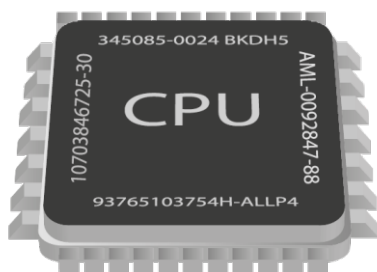
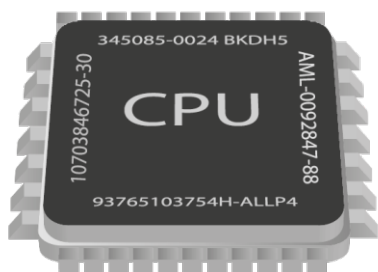
- Alt

Doesn't sound fun



How to distributed the
compute effort?

Distribute the effort!



Several issues

- Many models weaker than one
- Delays and Slow Nodes
- Communication Costs

Theory

- Barriers to entry / High Cost
- Implementation Overhead
- Nontrivial choice of ML-framework

Practice

How to Parallelize

$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k}(w_k; x_i)$$

- Parallelize computation of one update?

Issue:

computing $\nabla \ell_i(w; x_i)$ cheap (even for deep nets) [$O(d)$]

- Parallel Updates?

Issue: SGD is inherently serial....

Q: Can we parallelize inherently serial algorithms?

Simple idea: mini-batch SGD

- Compute multiple gradients in parallel

$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^1}(w_k; x_i)$$

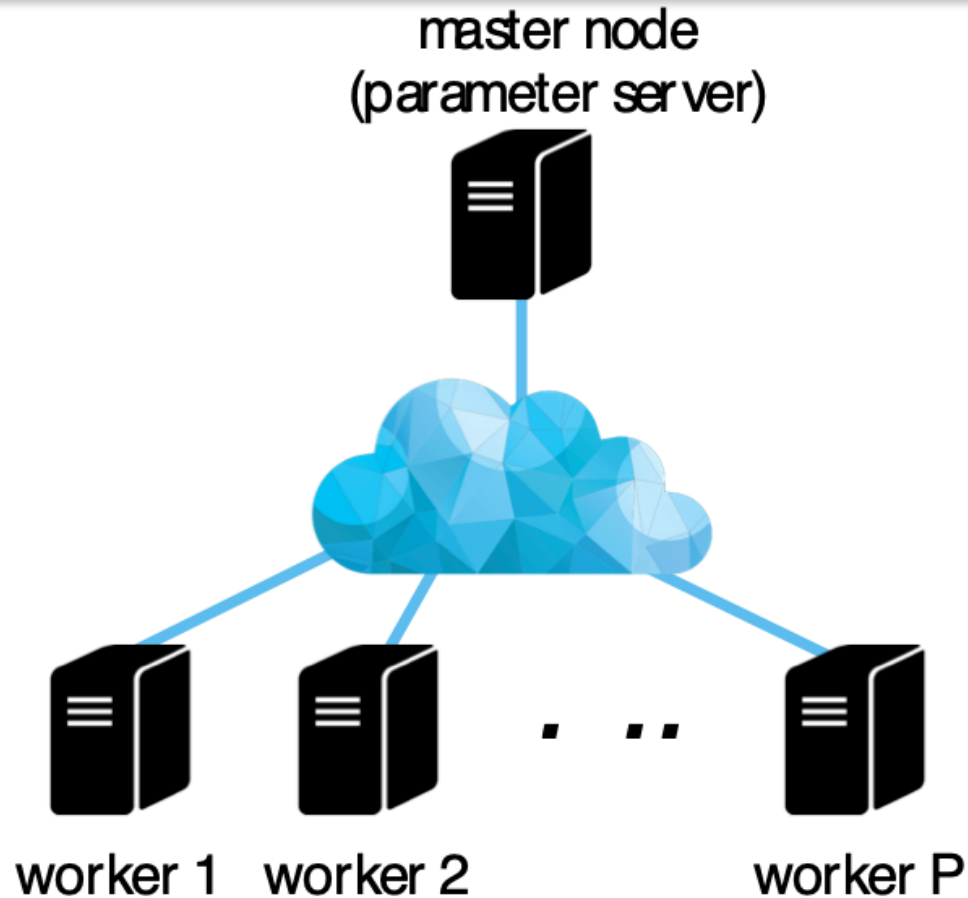
$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^2}(w_k; x_i)$$

$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^3}(w_k; x_i)$$

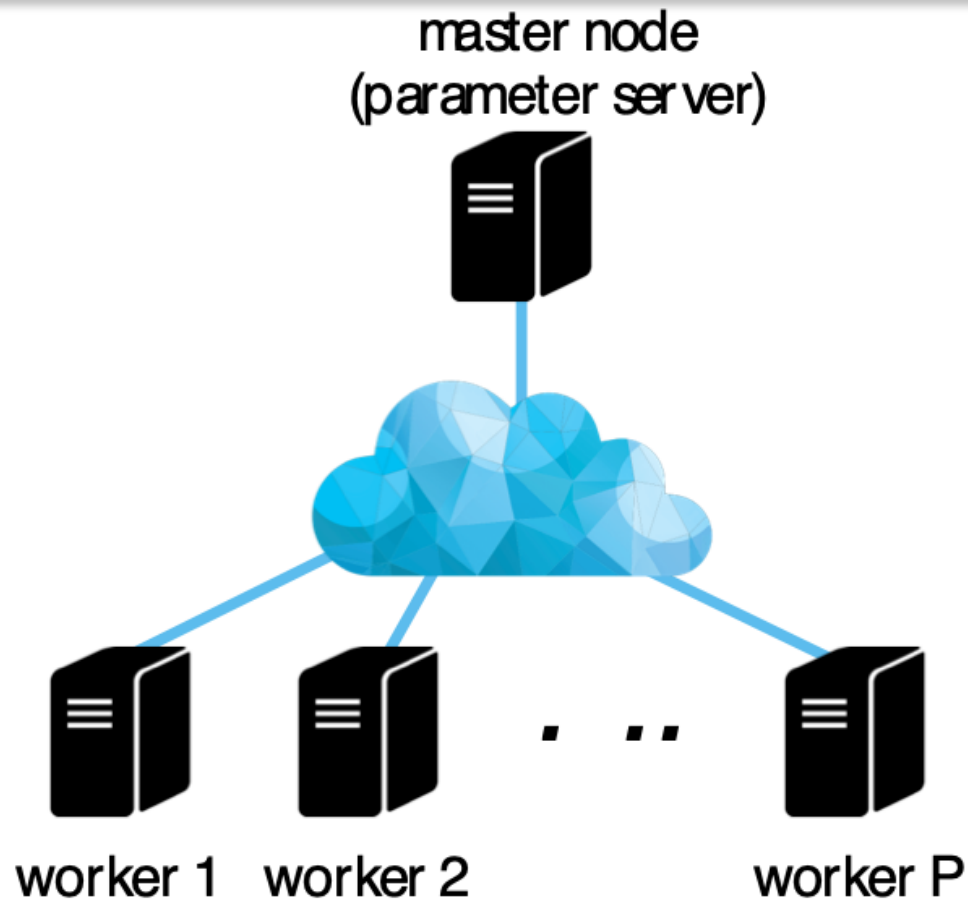
$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^4}(w_k; x_i)$$

Q: Does it perform the same as SGD?

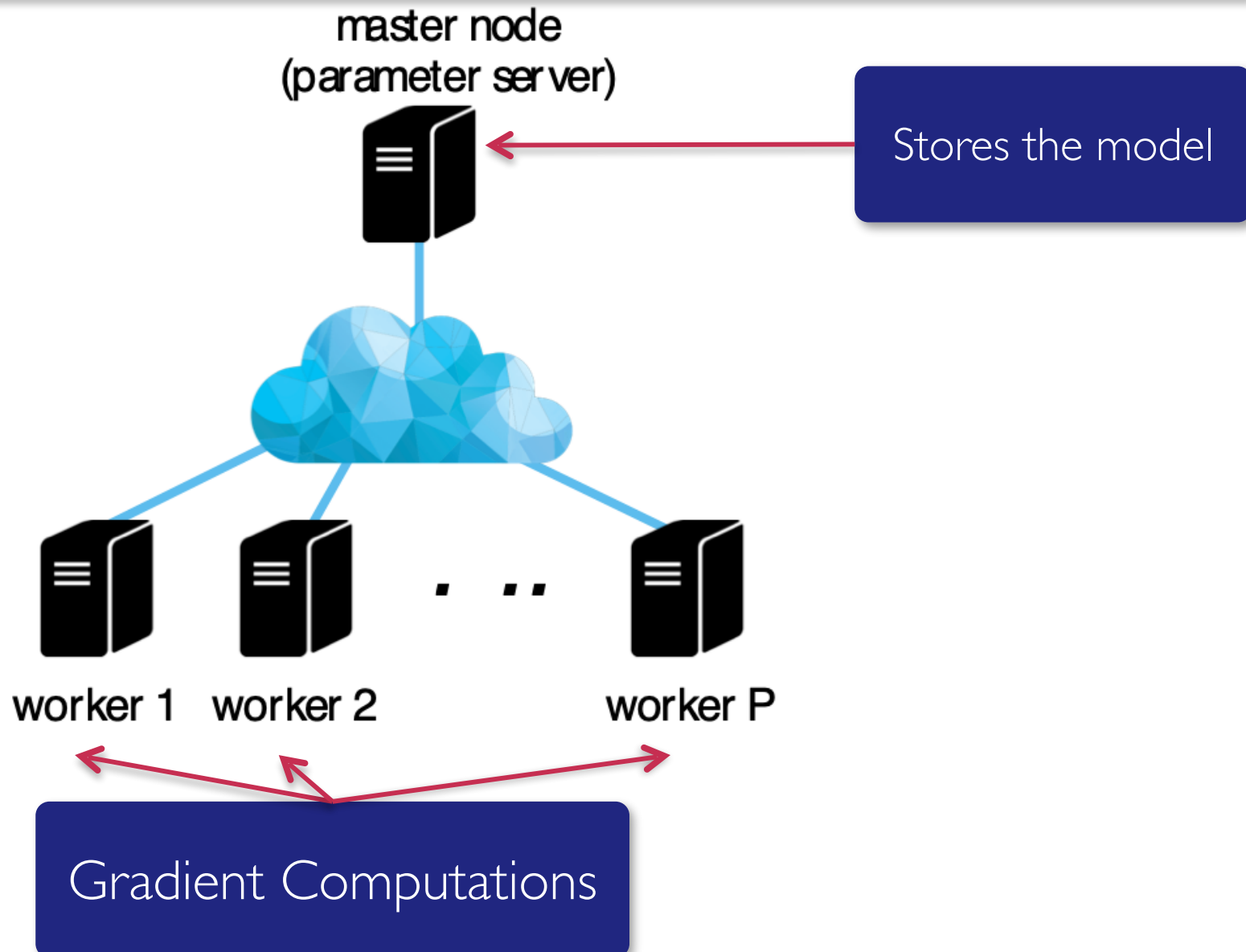
The Master-worker Setting



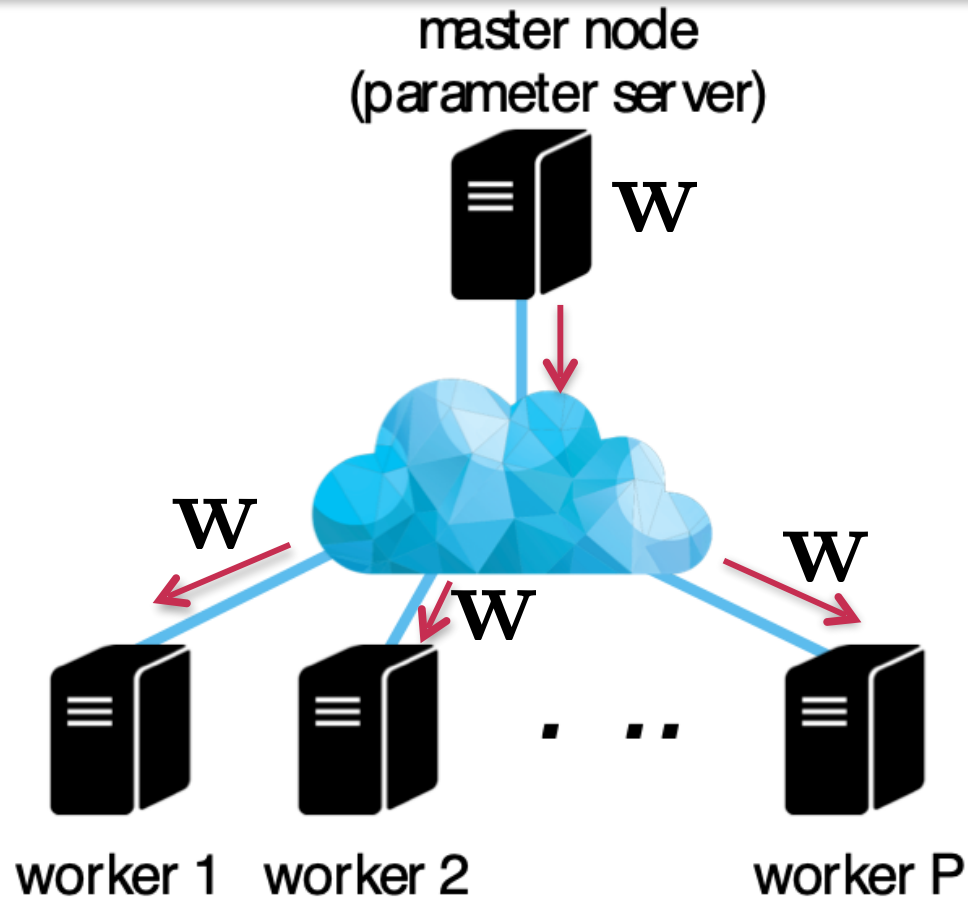
Algorithm of choice: minibatch SGD



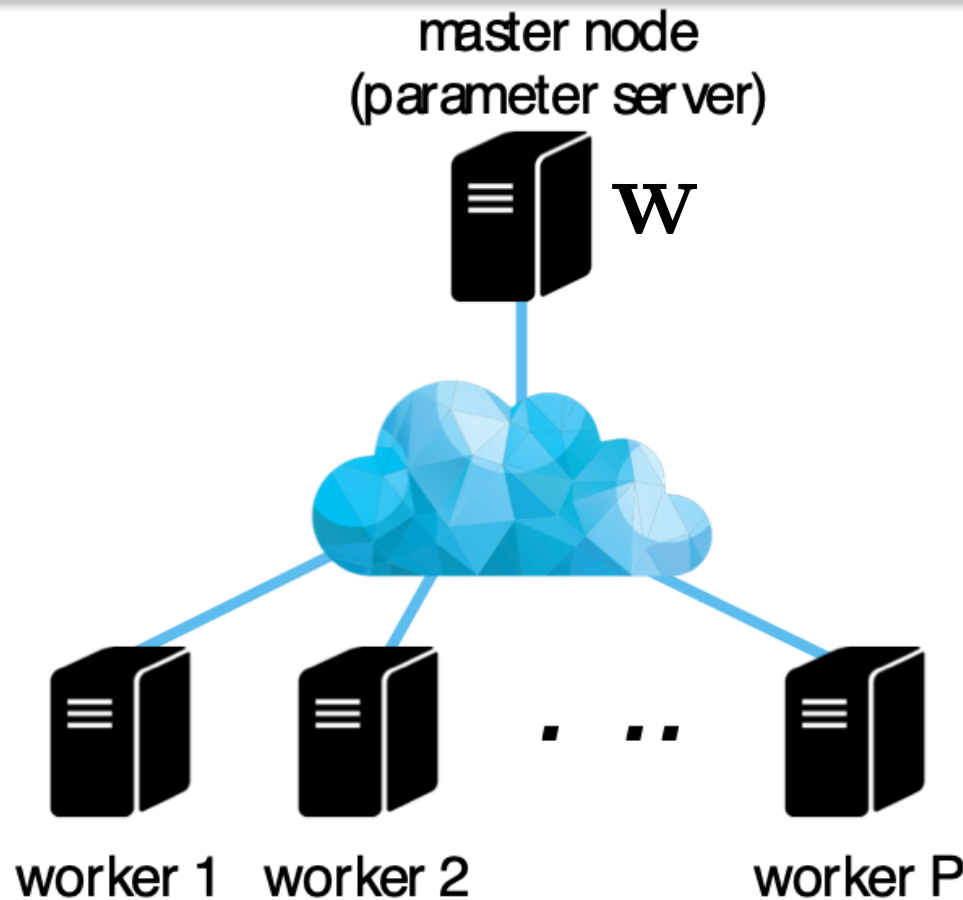
Algorithm of choice: minibatch SGD



Algorithm of choice: minibatch SGD



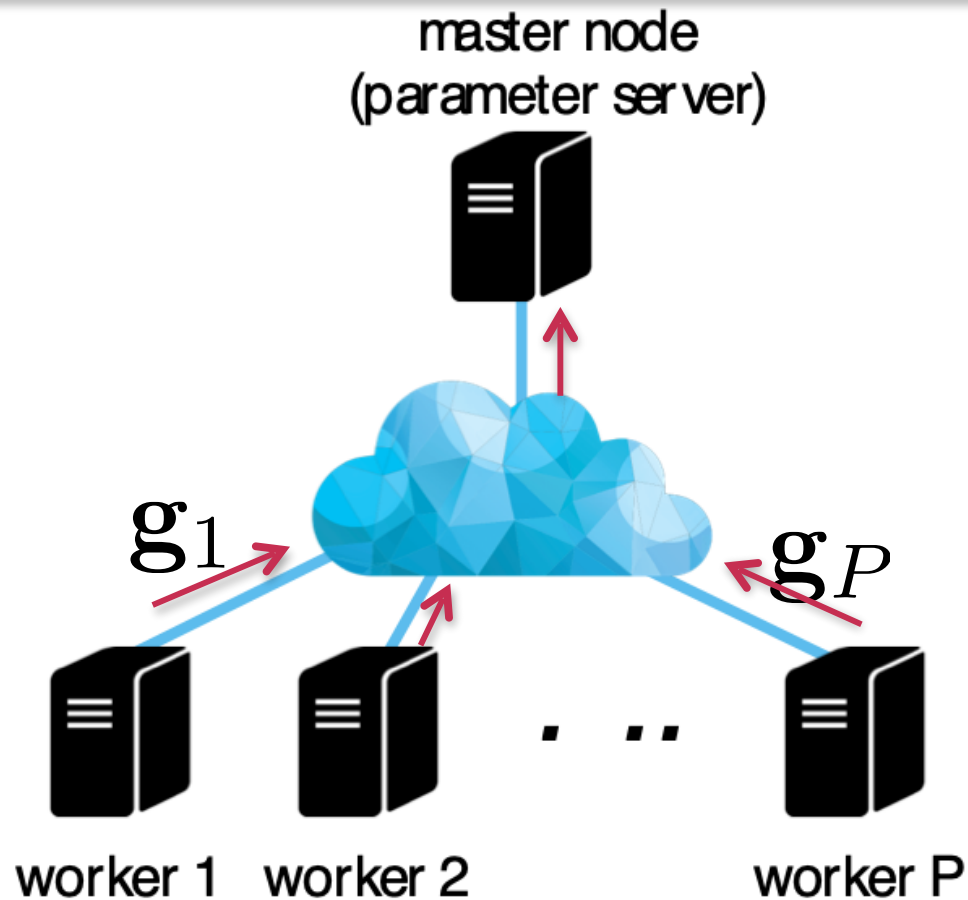
Algorithm of choice: minibatch SGD



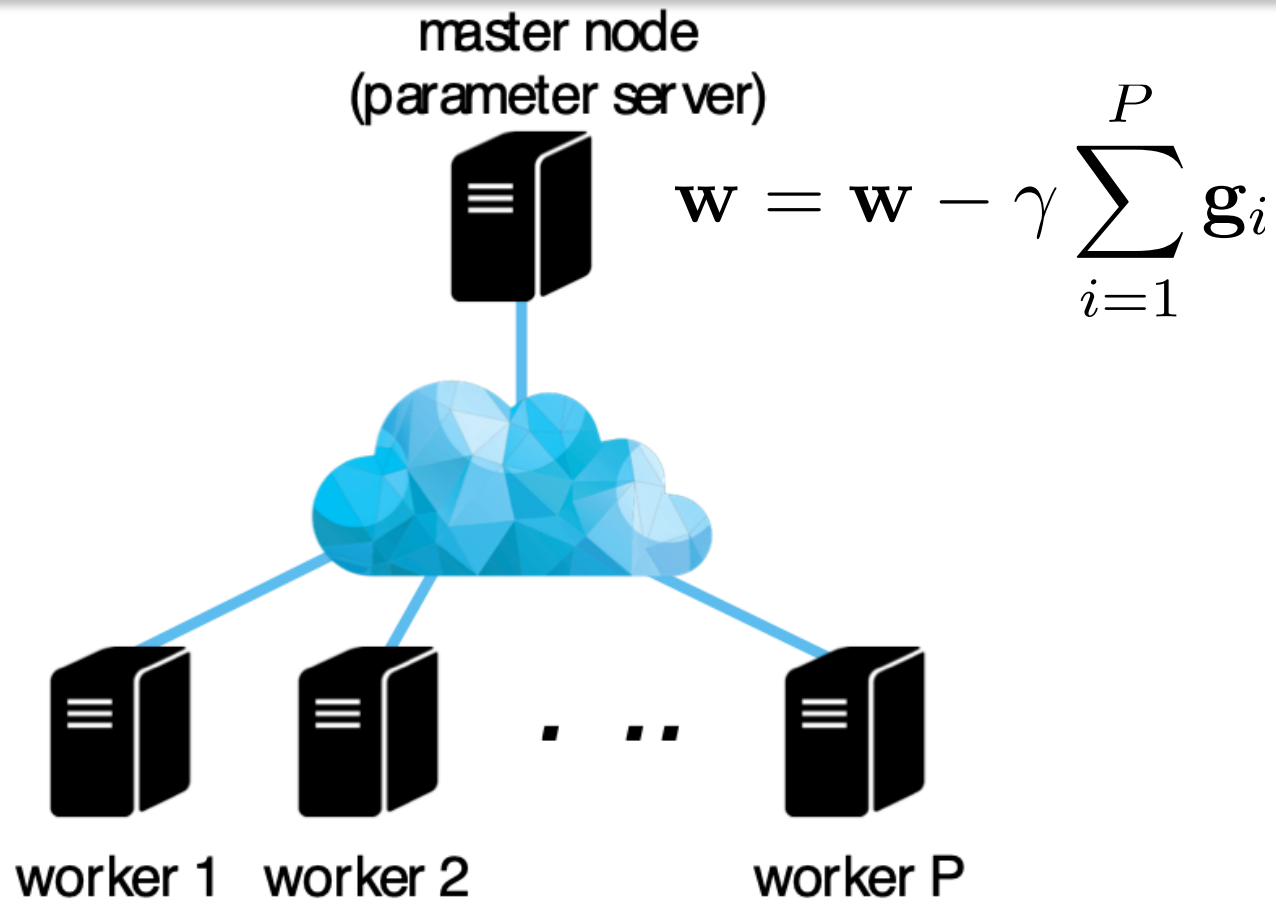
$$\mathbf{g}_1 = \sum_{i \in \mathcal{S}_1} \nabla \ell((\mathbf{x}_i, y_i); \mathbf{w})$$

$$\mathbf{g}_P = \sum_{i \in \mathcal{S}_P} \nabla \ell((\mathbf{x}_i, y_i); \mathbf{w})$$

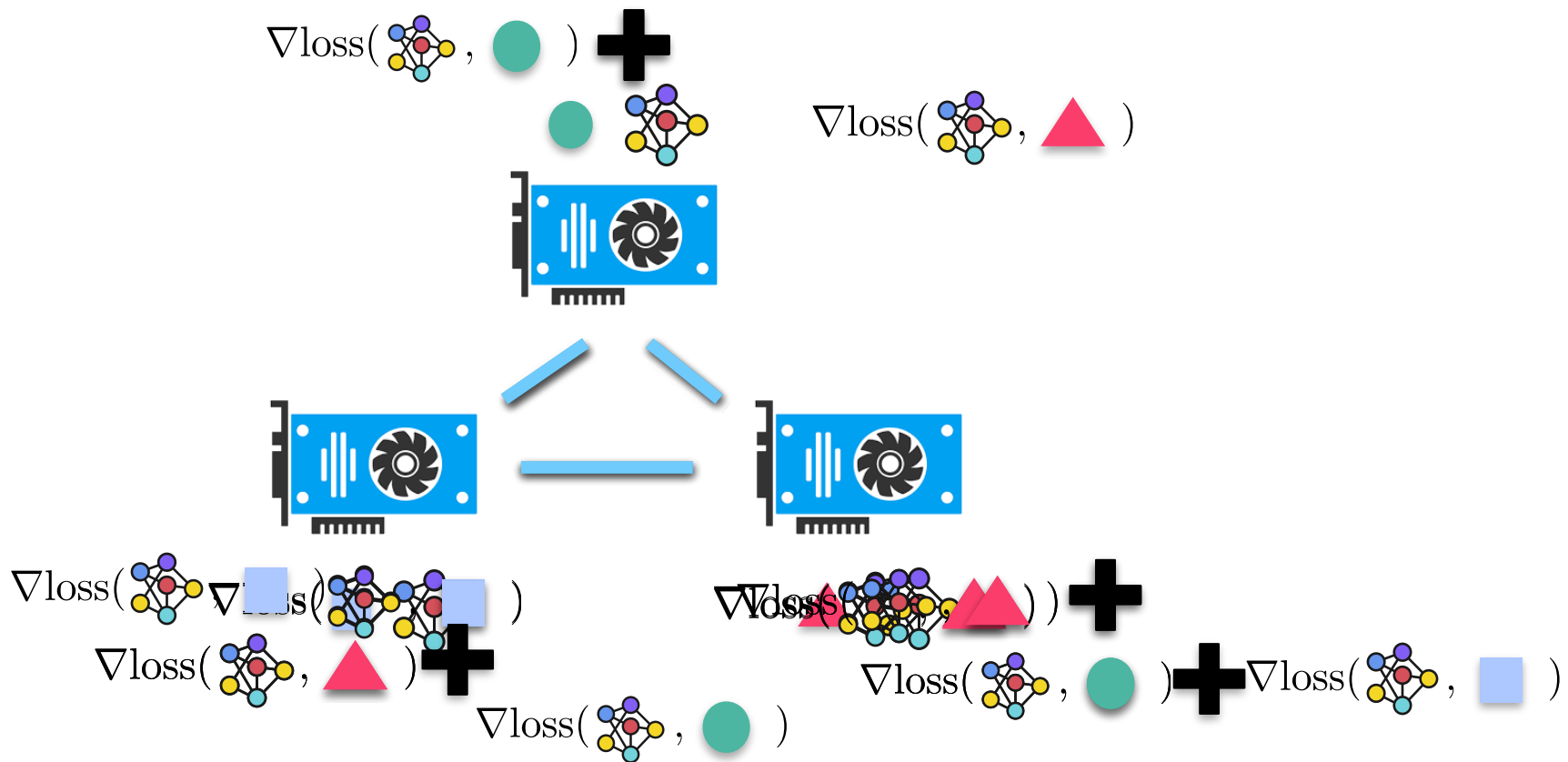
Algorithm of choice: minibatch SGD



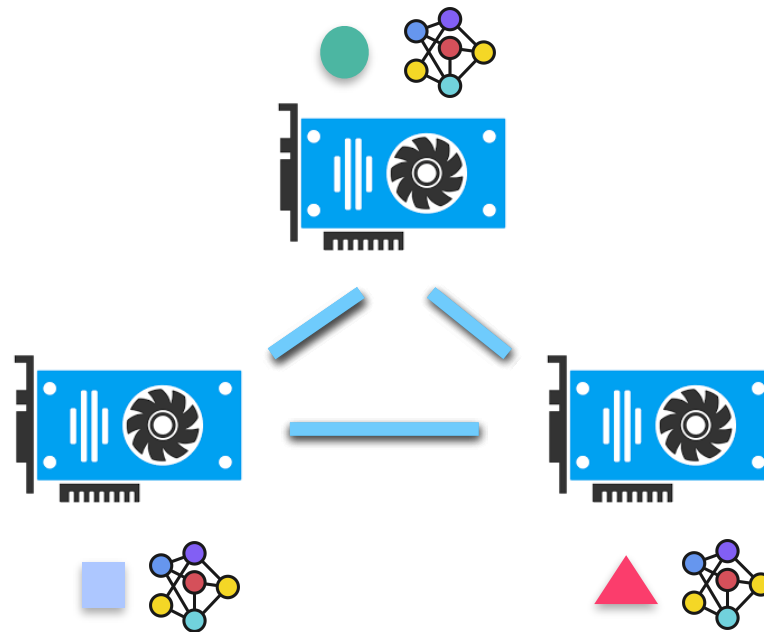
Algorithm of choice: minibatch SGD



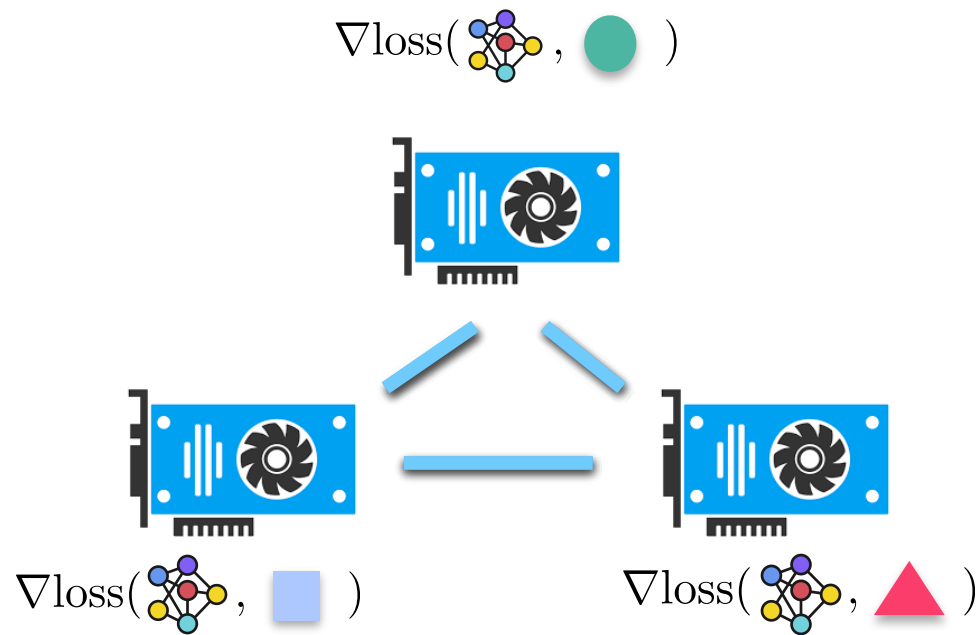
All-reduce



All-reduce

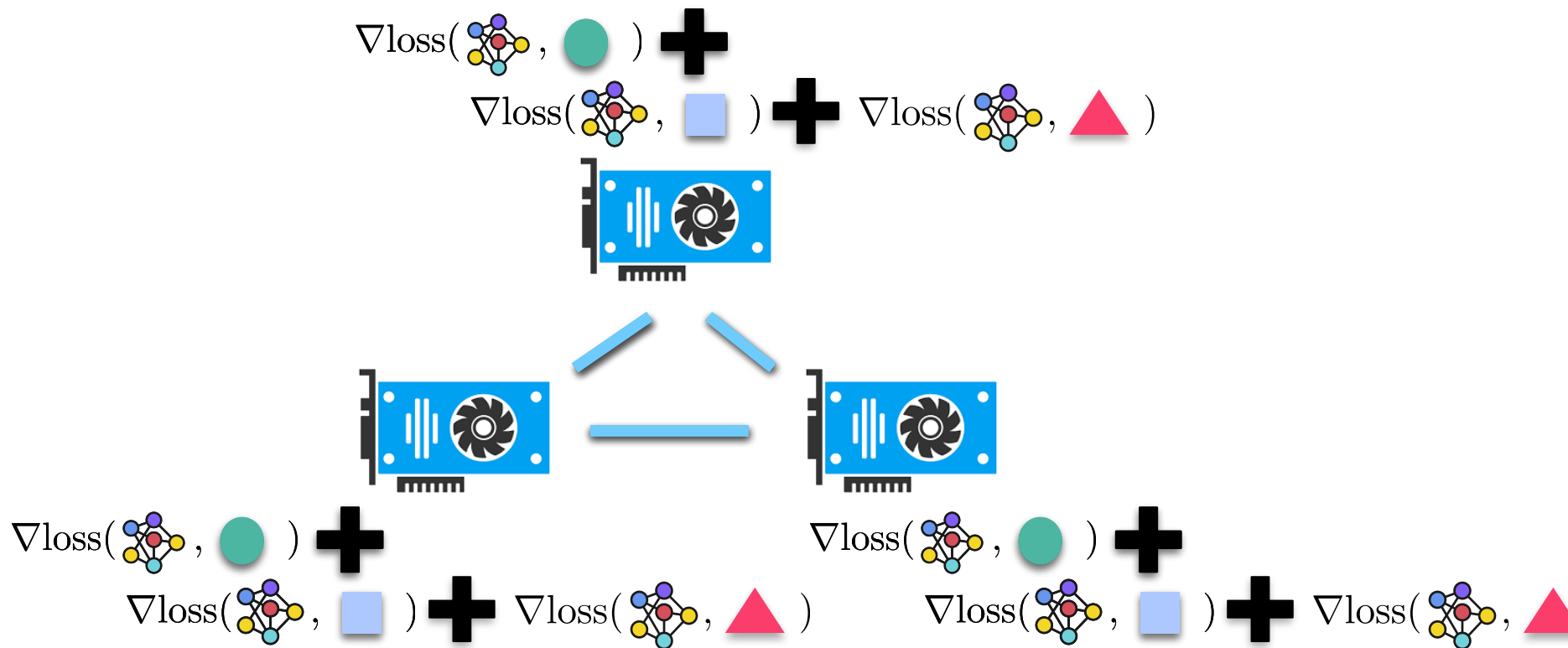


All-reduce



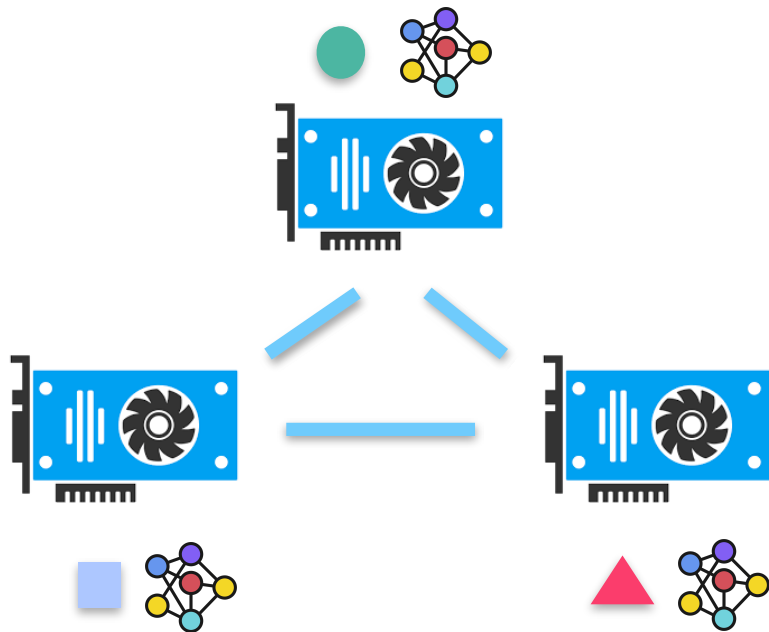
All-reduce

after $K-1$ push rounds



All-reduce, the server-free case

All-reduce SGD



Repeat until convergence

[Cotter, Shamir, Srebro, Sridharan, NIPS 11]

[Dekel, Gilad-Bachrach, Shamir, Xiao, JMLR 2012]

[Friedlander and Schmidt, SIAM JSC 2012]

[Takác, Bijral, Richtárik, Srebro, ICML 2013]

[Li, Zhang, Chen, Smola, KDD 2014]

[Jain, Kakade, Kidambi, Netrapalli, Sidford, arxiv'16]

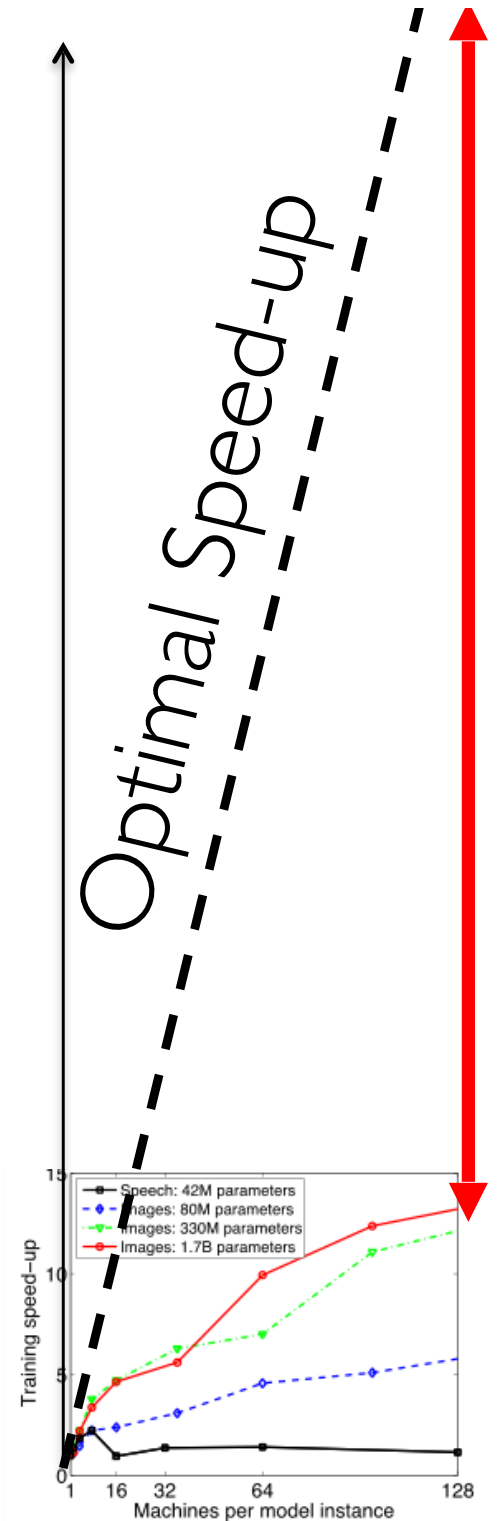
[De, Yadav, Jacobs, Goldstein, arxiv'16]

The Ideal Speedup Should be Proportional to
#Compute Nodes

Many Questions....

- How fast does distributed-SGD converge?
- How can we measure speed?
- How can we update the model faster?
- How can we reduce communication?
- What happens with delayed nodes?
- Does fault tolerance matter?

Compute and Communication Bottlenecks



Reality \sim 100x worse than optimal



“Large Scale Distributed Deep Networks” [Dean et al., NIPS 2012]

How to analyze parallel algorithms?

- Main measure of performance

$$\text{speedup} = \frac{\text{Time of serial } \mathcal{A} \text{ to accuracy } \epsilon}{\text{Time of parallel } \mathcal{A} \text{ to accuracy } \epsilon}$$

Example: Gradient Descent

Serial

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_k)$$

Parallel

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma \cdot \frac{1}{n} \left(\sum_{i_1=1}^P \nabla f_{i_1}(\mathbf{x}_k) + \dots + \sum_{i_P=n-P+1}^n \nabla f_{i_P}(\mathbf{x}_k) \right)$$

- Convergence is invariant of allocation
- Both algorithms reach to same accuracy after T iterations
- Speedup is independent of convergence rate

How to analyze parallel algorithms?

- Main measure of performance

$$\text{speedup} = \frac{\text{Time of serial } \mathcal{A} \text{ to accuracy } \epsilon}{\text{Time of parallel } \mathcal{A} \text{ to accuracy } \epsilon}$$

Example: Gradient Descent

Serial

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma \cdot \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_k)$$

Embarrassingly Parallel $O(\#\text{cores})$ speedup

Parallel

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma \cdot \frac{1}{n} \left(\sum_{i_1=1}^P \nabla f_{i_1}(\mathbf{x}_k) + \dots + \sum_{i_P=n-P+1}^n \nabla f_{i_P}(\mathbf{x}_k) \right)$$

- Convergence is invariant of allocation
- Both algorithms converge after T iterations
- Not true for mini-batch SGD
- Speedup is independent of convergence rate

speedups for minibatch SGD

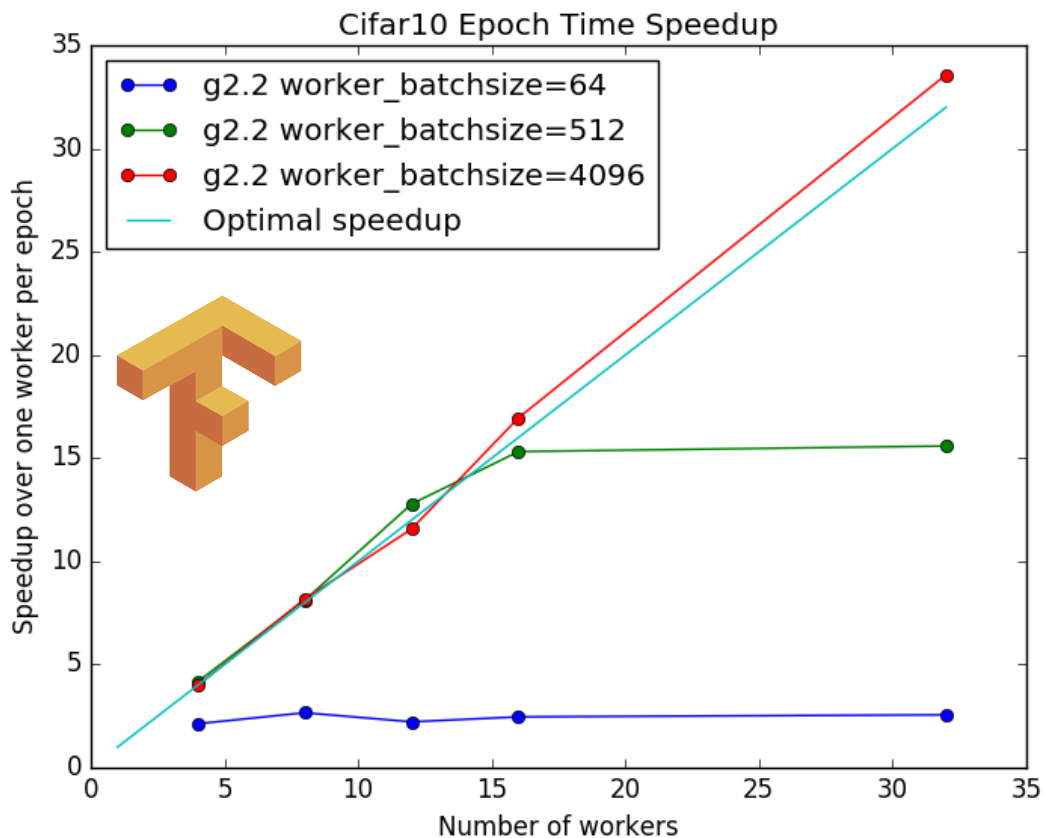
Two factors control run-time

$$\text{Time to accuracy } \varepsilon = [\text{time per data pass}] \times [\#\text{passes to accuracy } \varepsilon]$$

speedups for minibatch SGD

“[...] on more than 8 machines [...] network overhead starts to dominate [...]”

- TL;DR: Communication is the Bottleneck
- Why?

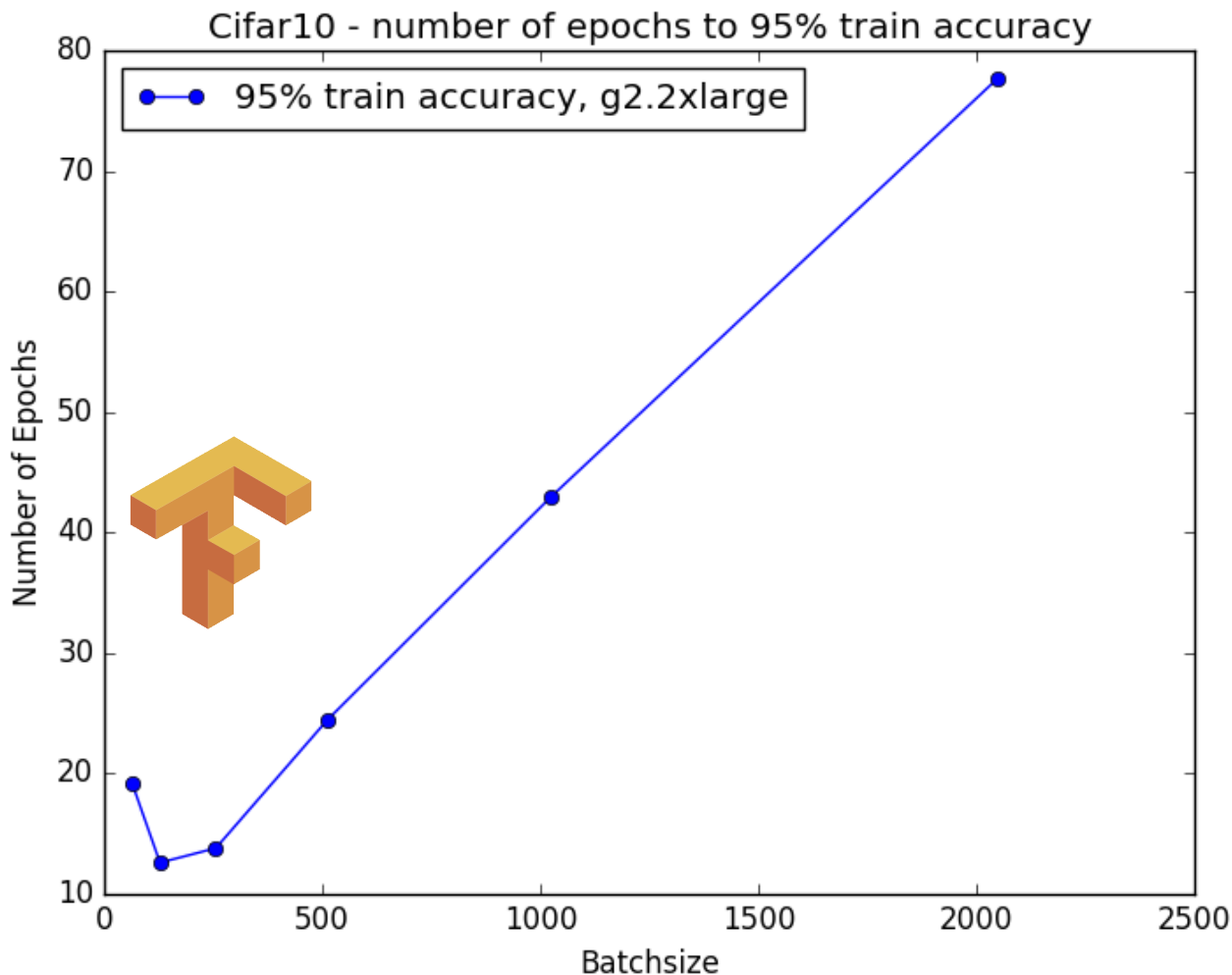


Time per pass:
time for $\text{dataset_size}/\text{batch_size}$
distributed iterations

⇒ Bigger Batch
Less Communication
(smaller time per epoch)

What's wrong with Large Batches?

- If small batch is bad, then maximize it



⇒

Large Batch
worse train error
(more #passes to accuracy ϵ)

⇒

Large Batch
some generalization issues

[Keskar, Mudigere, Nocedal,
Smelyanskiy, Tang, ICLR 2017]

How to Analyze mini-batch?

- Measure of performance

$$\text{worst case speedup} = \frac{\text{bound on \#iter of SGD to } \epsilon}{\text{bound on \#iter of Parallel SGD to } \epsilon}$$

Main Question:

How does minibatch SGD compare against serial SGD?

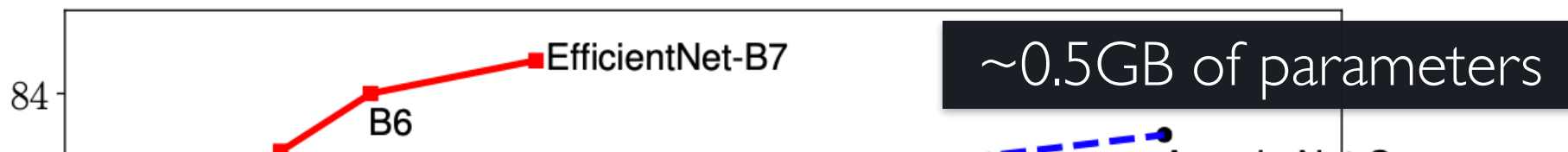
Main questions:

Convergence after T gradient computations

Answer is Complicated: Depends on Problem
Generally if batch $B_0 > B(\text{data}, \text{loss})$
Minibatch SGD offers no speedup.

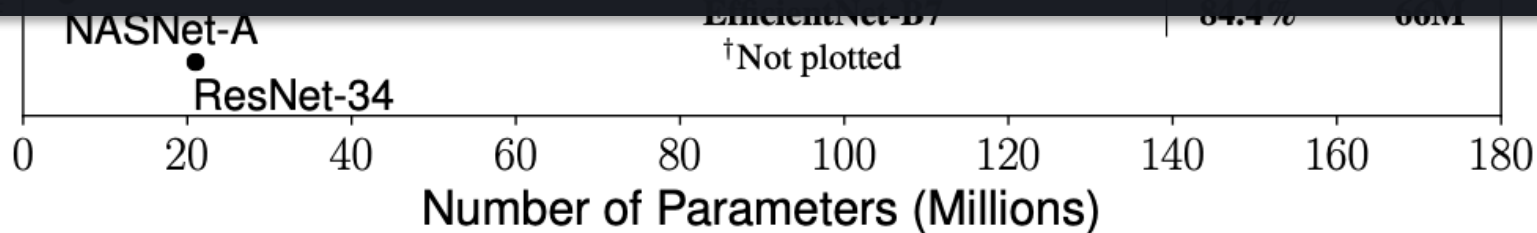
Modern Architectures are huge
make everything slower

Model sizes (vision)



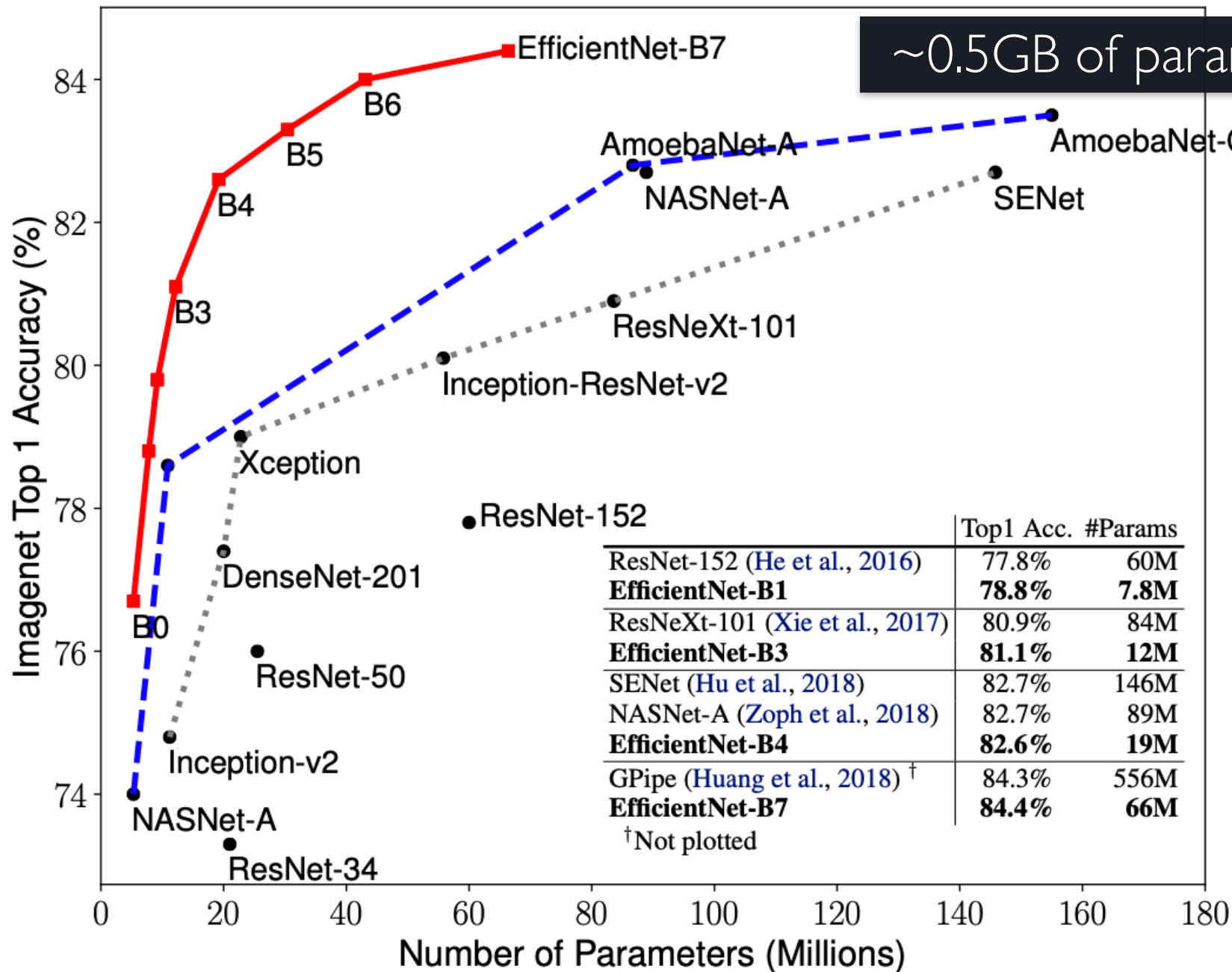
Q: is 0.5GB that large to cause a bottleneck?

A: When training we compute gradients on B data points = $B * \text{model size}$ RAM required



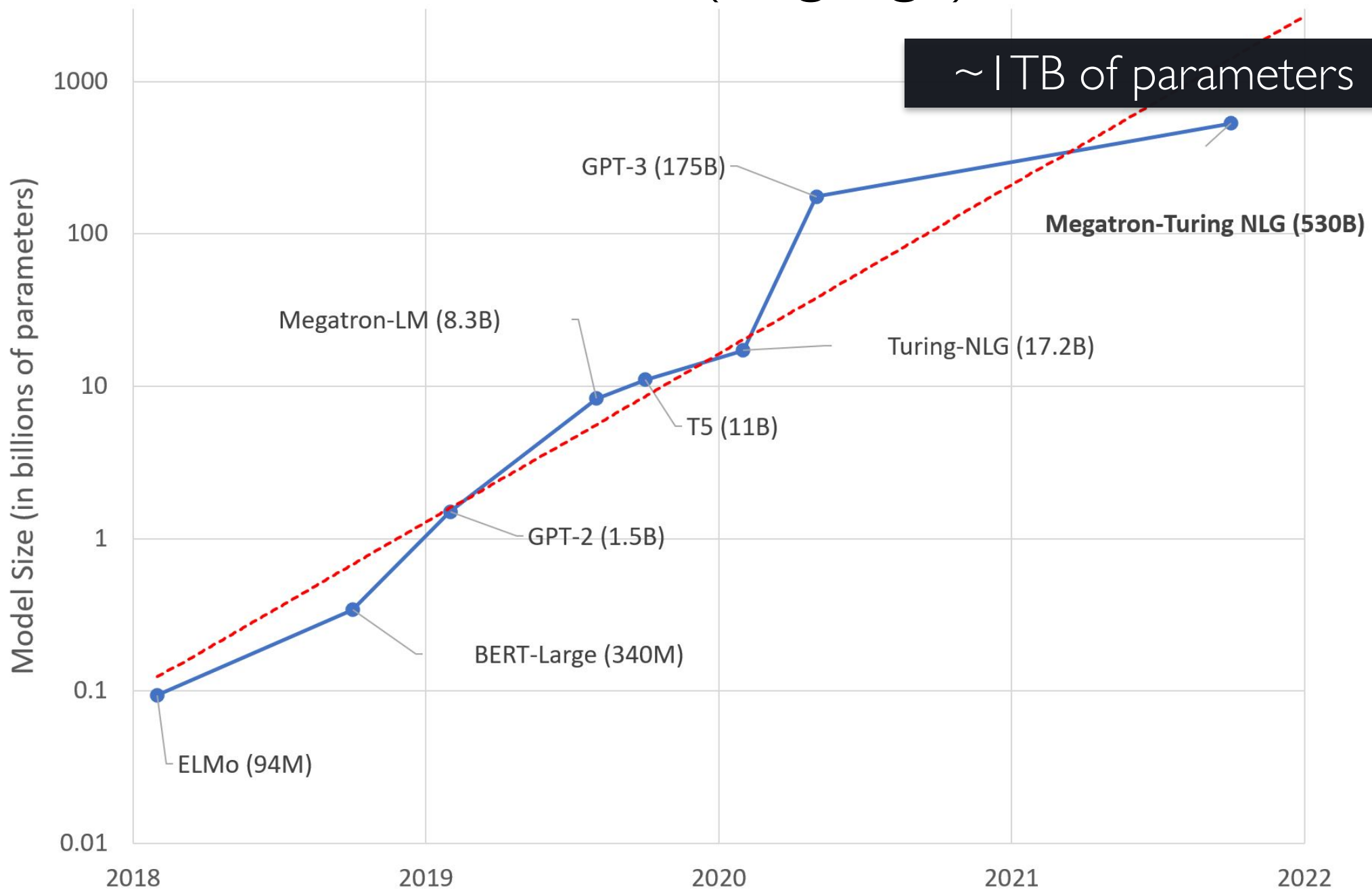
Model sizes (vision)

~0.5GB of parameters



[Source: <http://proceedings.mlr.press/v97/tan19a/tan19a.pdf>]

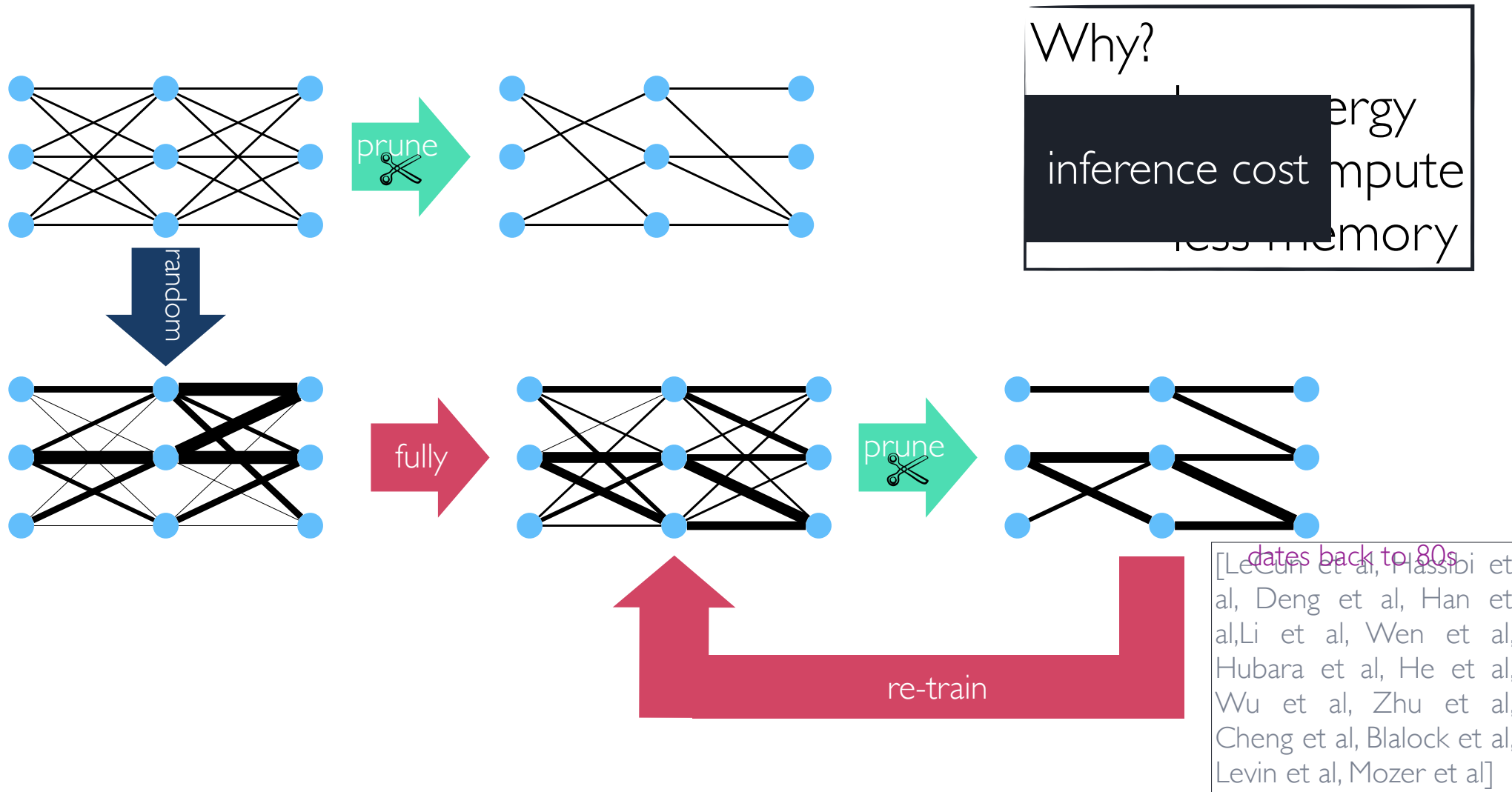
Model sizes (language)



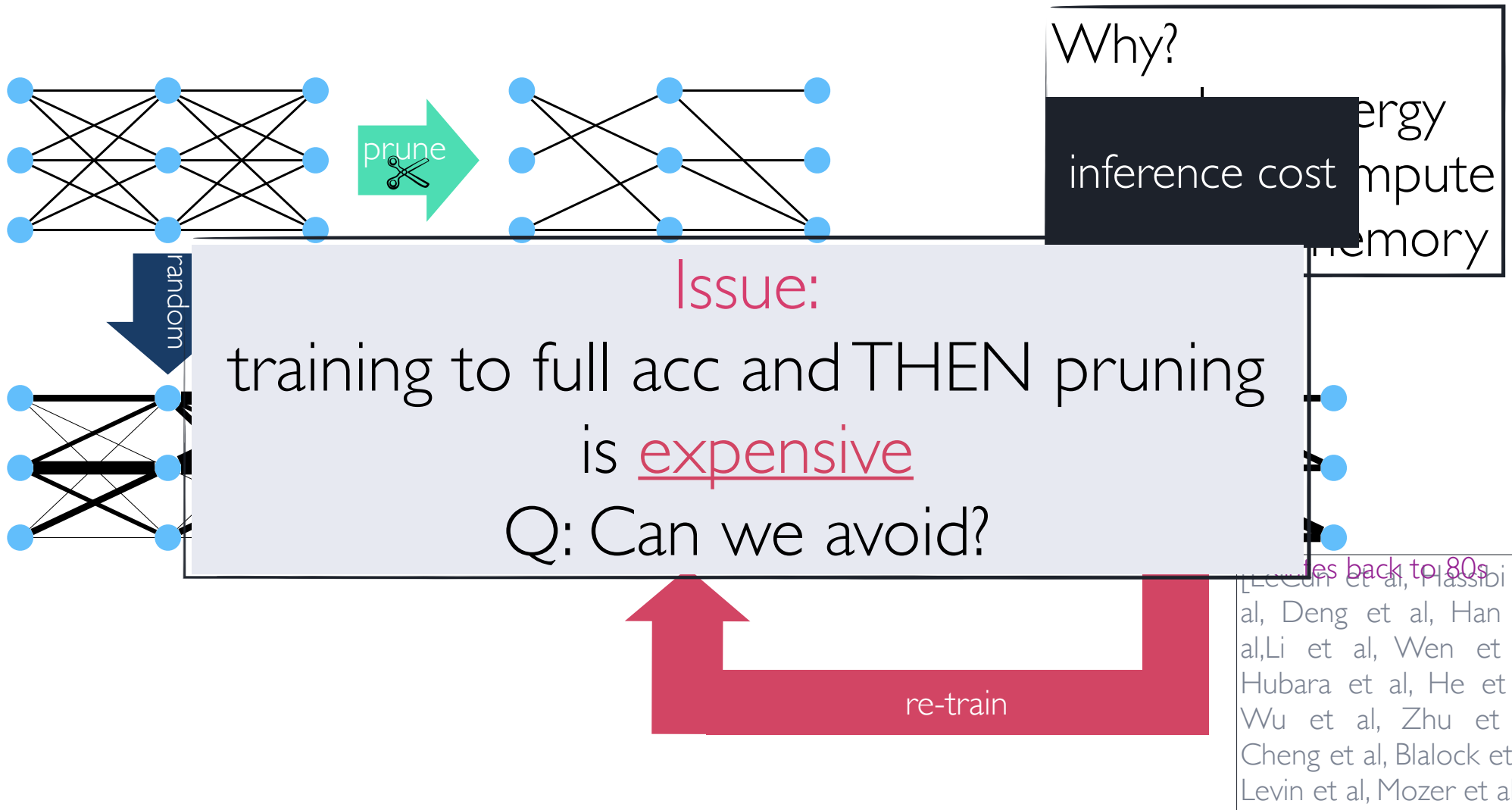
[Source: <https://tinyurl.com/Megatron-Turing-NLG>]

compression by sparsity

Network Pruning, 1980-2018

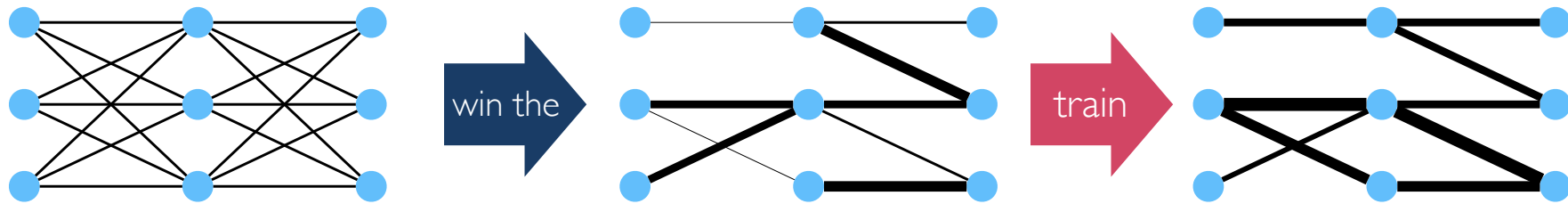


Network Pruning, 1980-2018

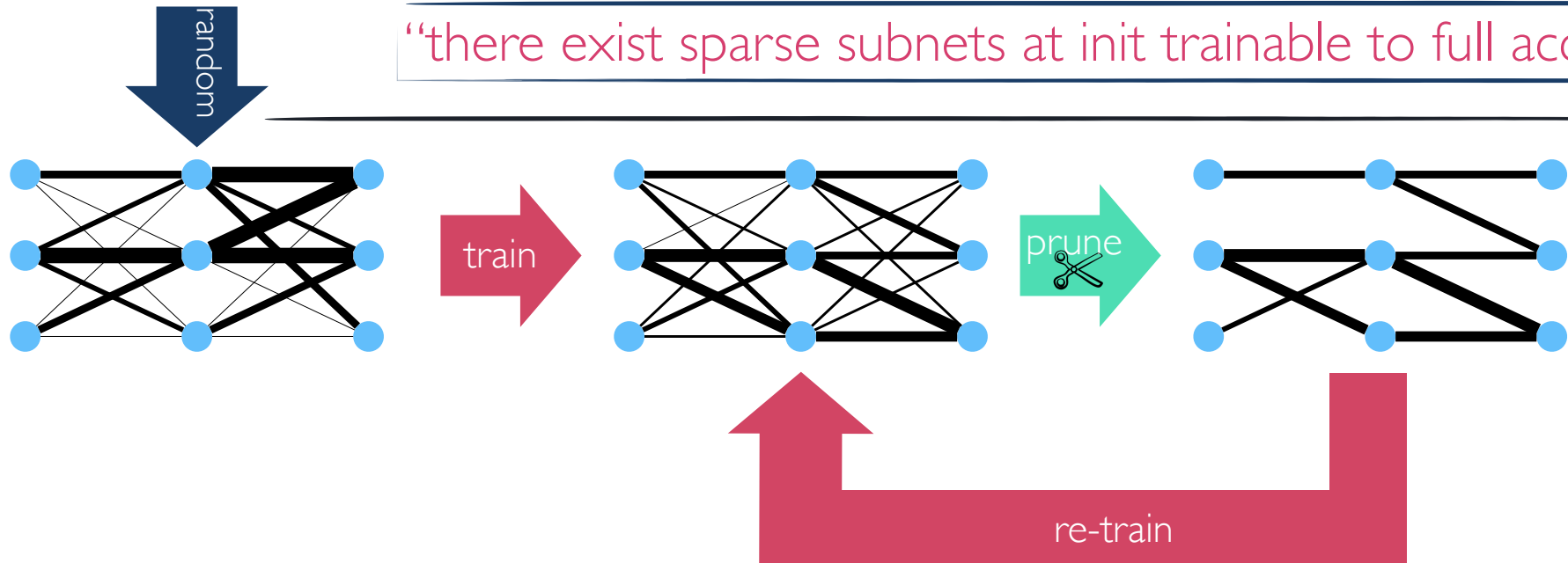


Lottery Ticket Hypothesis (LTH)

Frankle, Carbin, ICLR 2019 



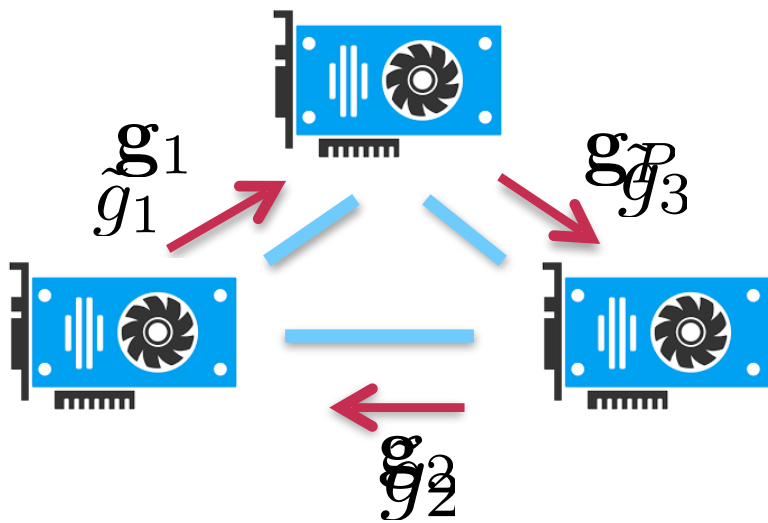
“there exist sparse subnets at init trainable to full accuracy”*



Gradient Compression

Communication / worker:

$$O(\text{size of gradient}) * 32 \text{ bits}$$



Gradient Quantization

Quantize to precision:

$$O(\text{size of gradient}) * 2/4/8 \text{ bits}$$

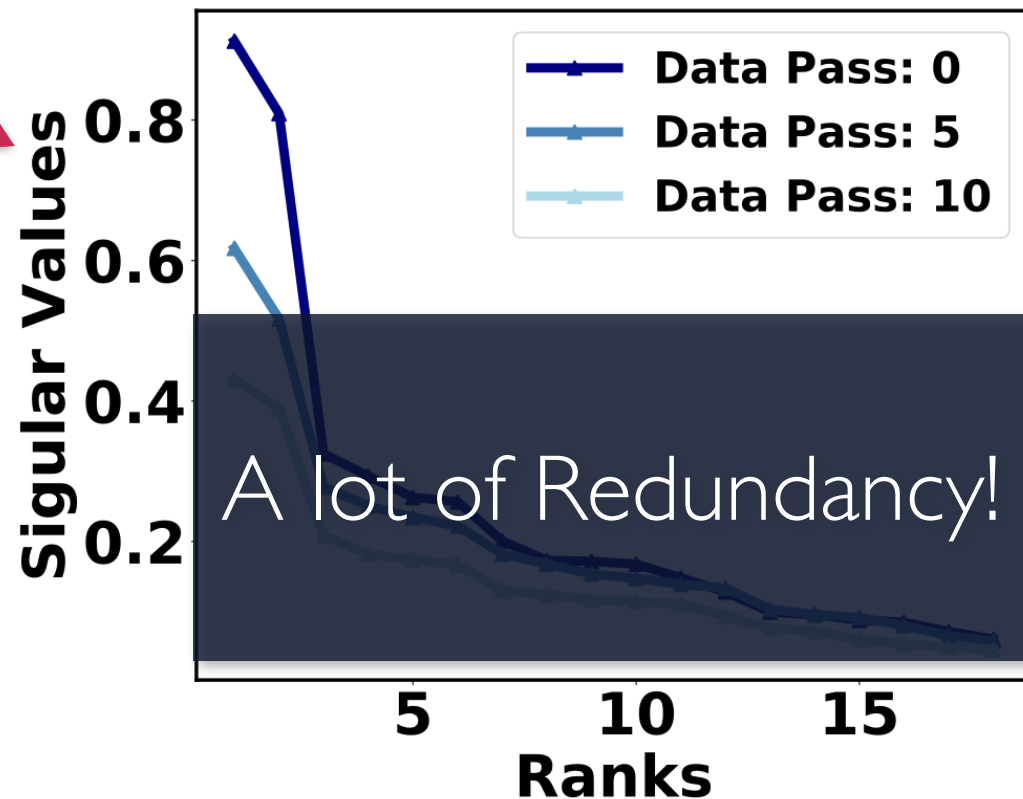
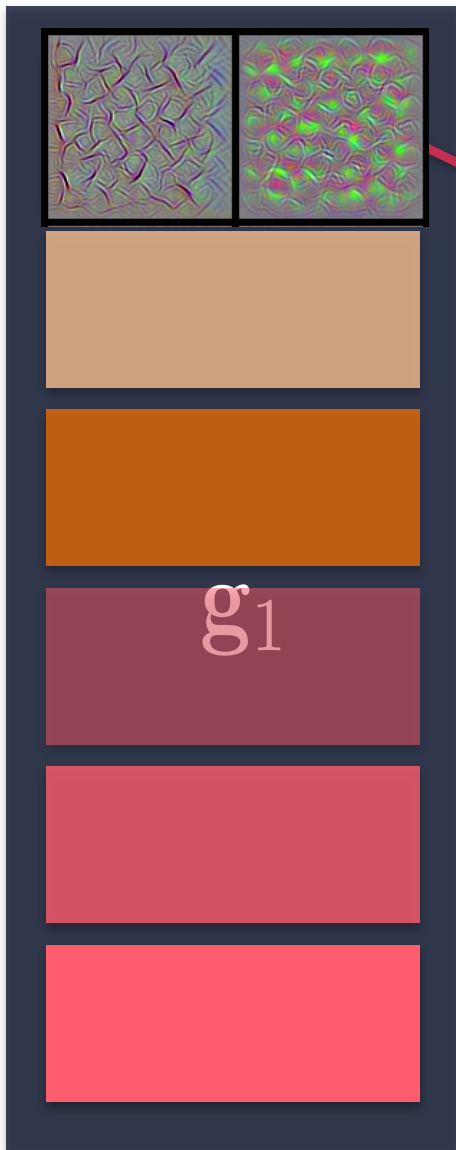
Gradient Sparsification

Sparsified Gradients (k-sparse):

$$O(k)$$

Gradients nearly-sparse in SVD

- Observation: Fast decaying spectral profile (SVD)



- *What if we compress in the spectral domain?*

LoRA: Low-Rank Adaptation of Large Language Models

Edward Hu* Yelong Shen* Phillip Wallis
Zeyuan Allen-Zhu Yuanzhi Li Shean Wang Weizhu Chen
Microsoft Corporation
{edwardhu, yeshe, phwallis, zeyuana, yuanzhil
swang, wzchen}@microsoft.com

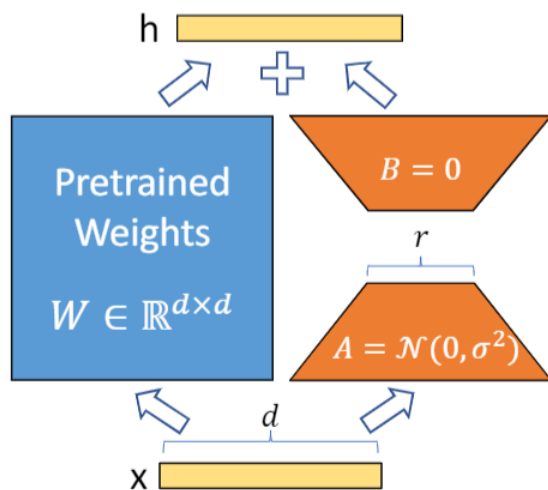
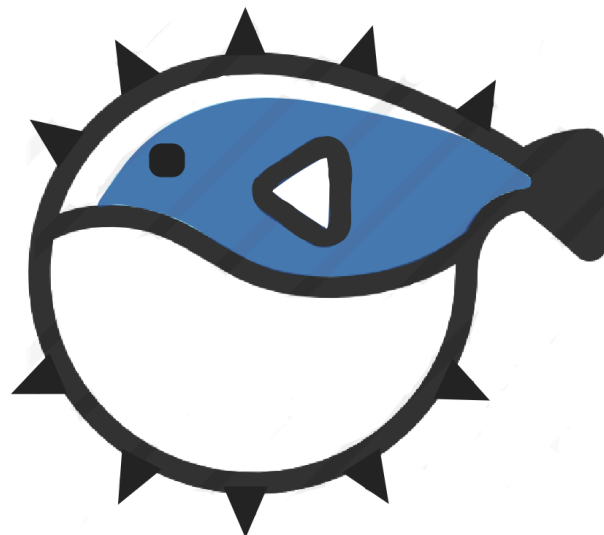


Figure 1: Our reparametrization. We only train A and B .

PUFFERFISH: COMMUNICATION-EFFICIENT MODELS AT NO EXTRA COST

Hongyi Wang,¹ Saurabh Agarwal,¹ Dimitris Papailiopoulos²



compression by quantization

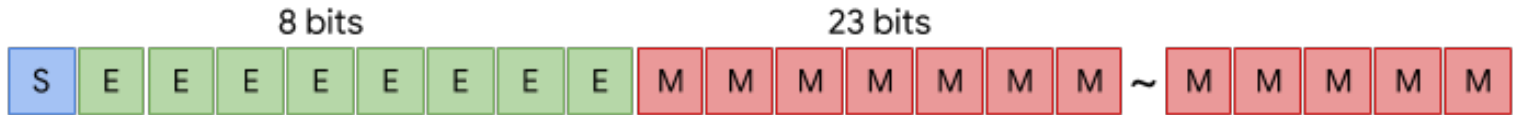
bfloat16

range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



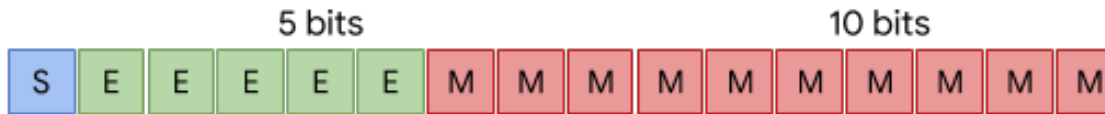
float32

range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



float16

range: $\sim 5.9e^{-8}$ to $6.5e^4$

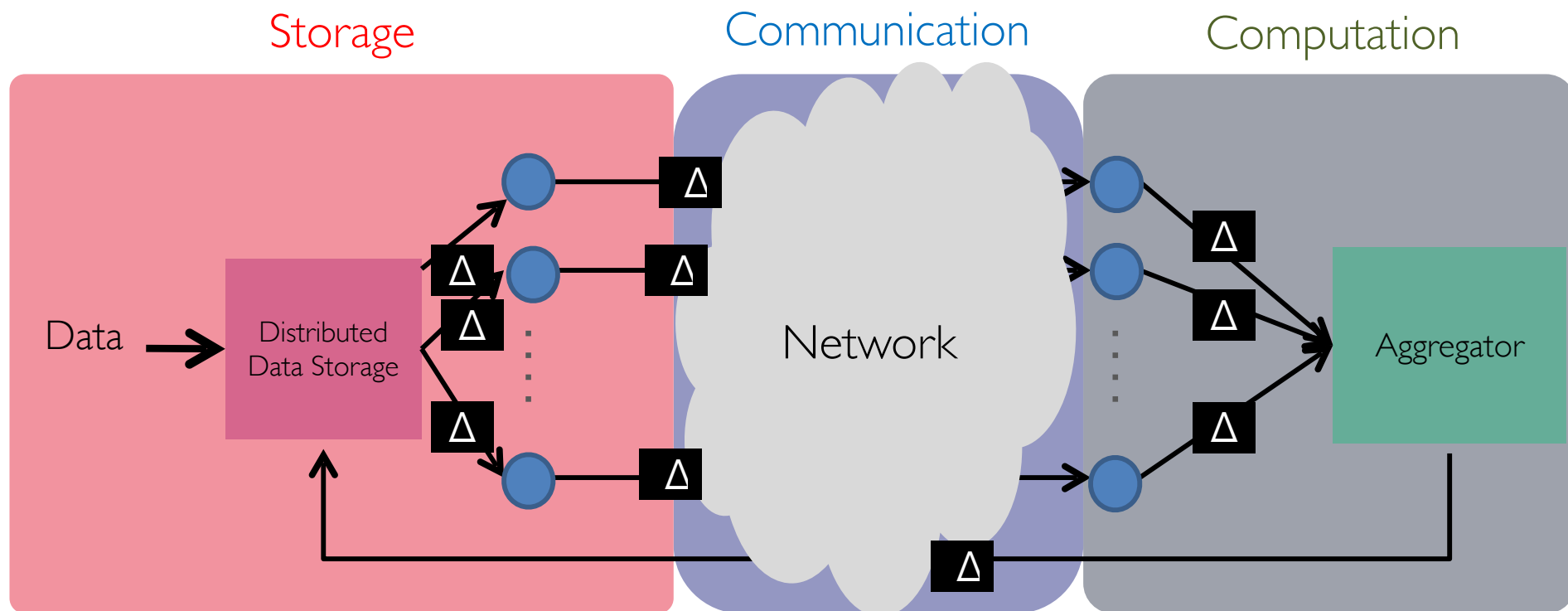


Makes model smaller!

Can affect the training dynamics

Straggler Nodes

Large-scale Distributed Machine Learning Systems

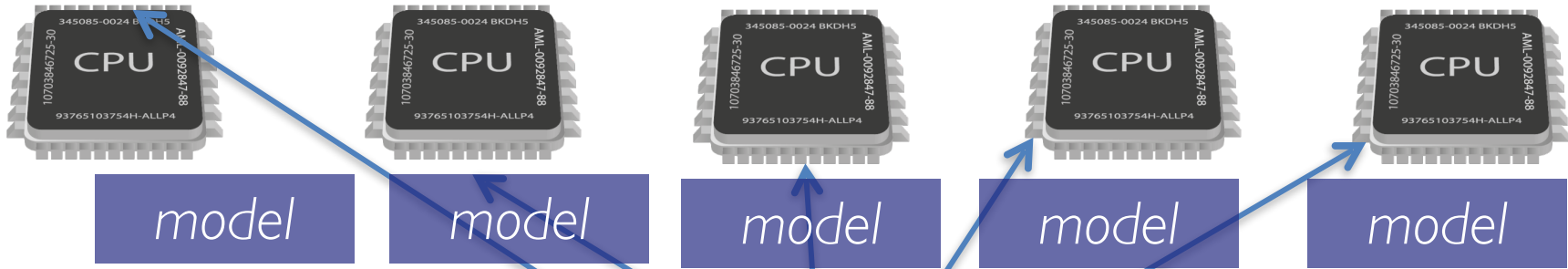
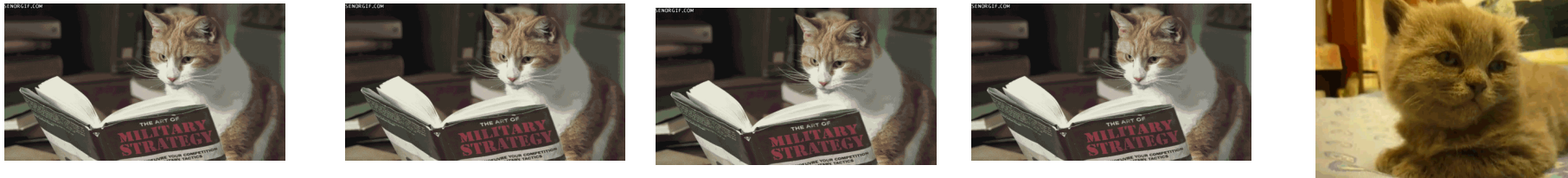


“The scale and complexity of modern Web services make it infeasible to eliminate all latency variability.”

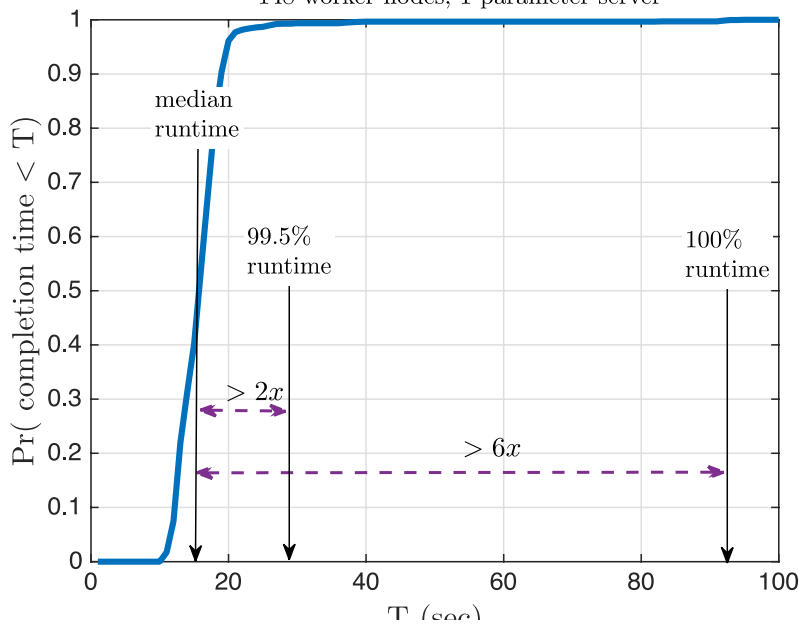
Jeff Dean, Google.



Bottleneck 2: Straggling Learners



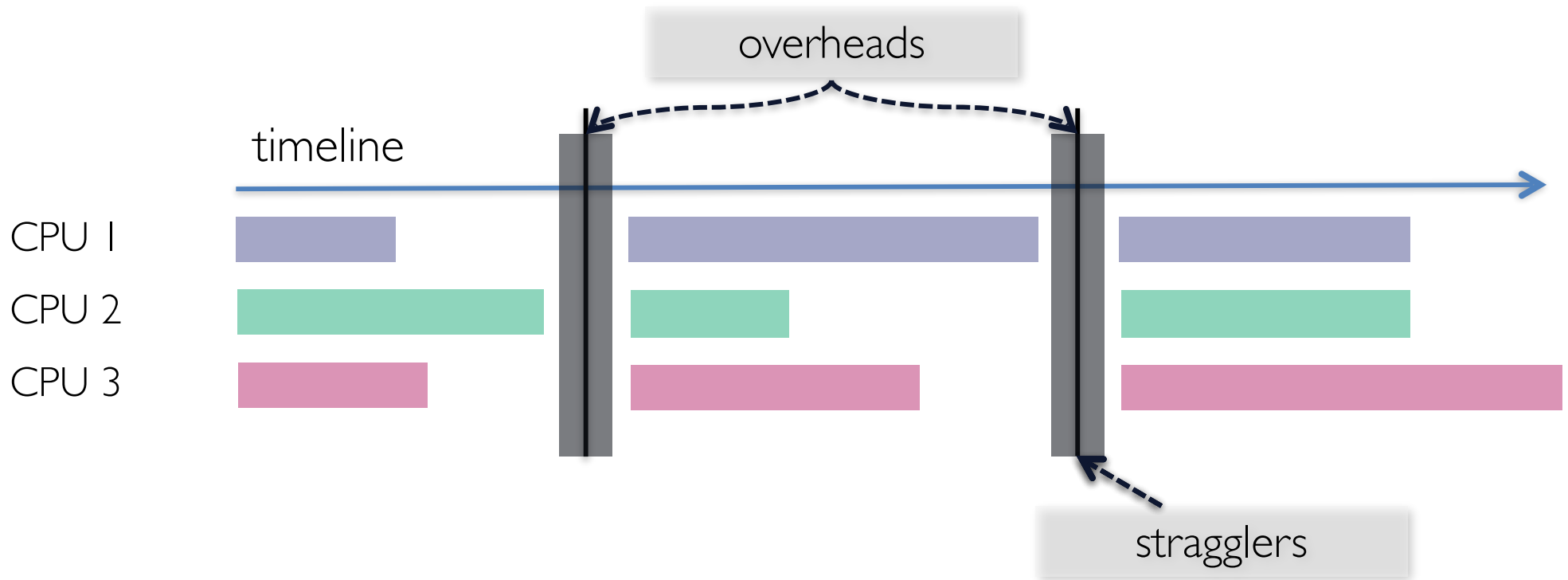
Iteration Completion Time per worker
Data set = CIFAR-10
t2.small EC2 instances
148 worker nodes, 1 parameter server



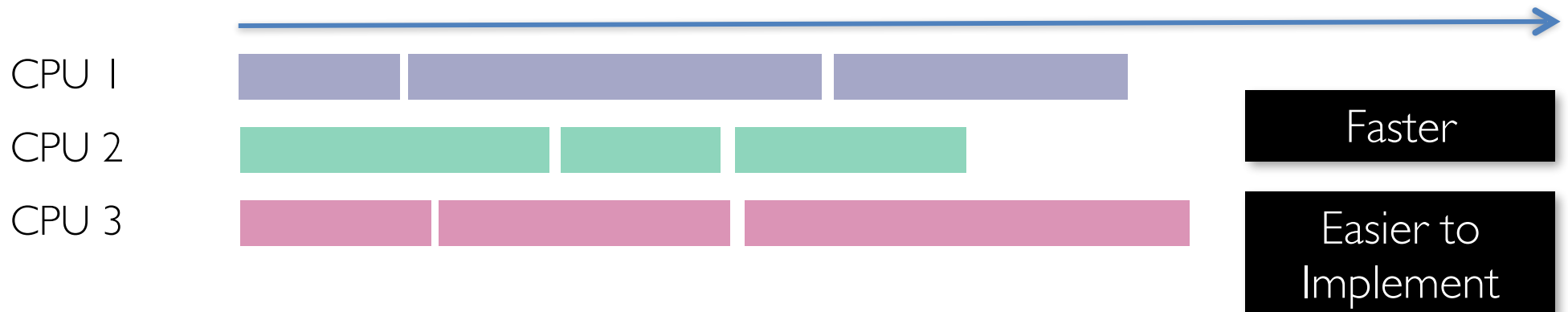
Can we “robustify” distributed ML against stragglers?

Measured on Amazon AWS

A case against Synchronization



Asynchronous World



HOGWILD! 2011

“Run parallel lock-free SGD without synchronization”



Niu



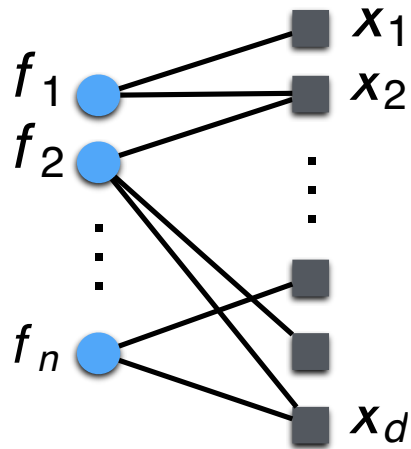
Recht



Ré



Wright



Each processor in parallel

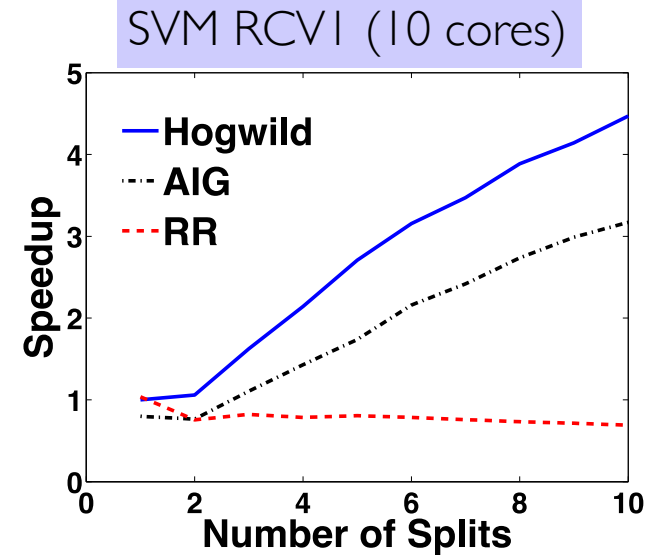
sample function f_i

$x = \text{read shared memory}$

$g = -\gamma \cdot \nabla f_i(x)$

for v in the support of f **do**

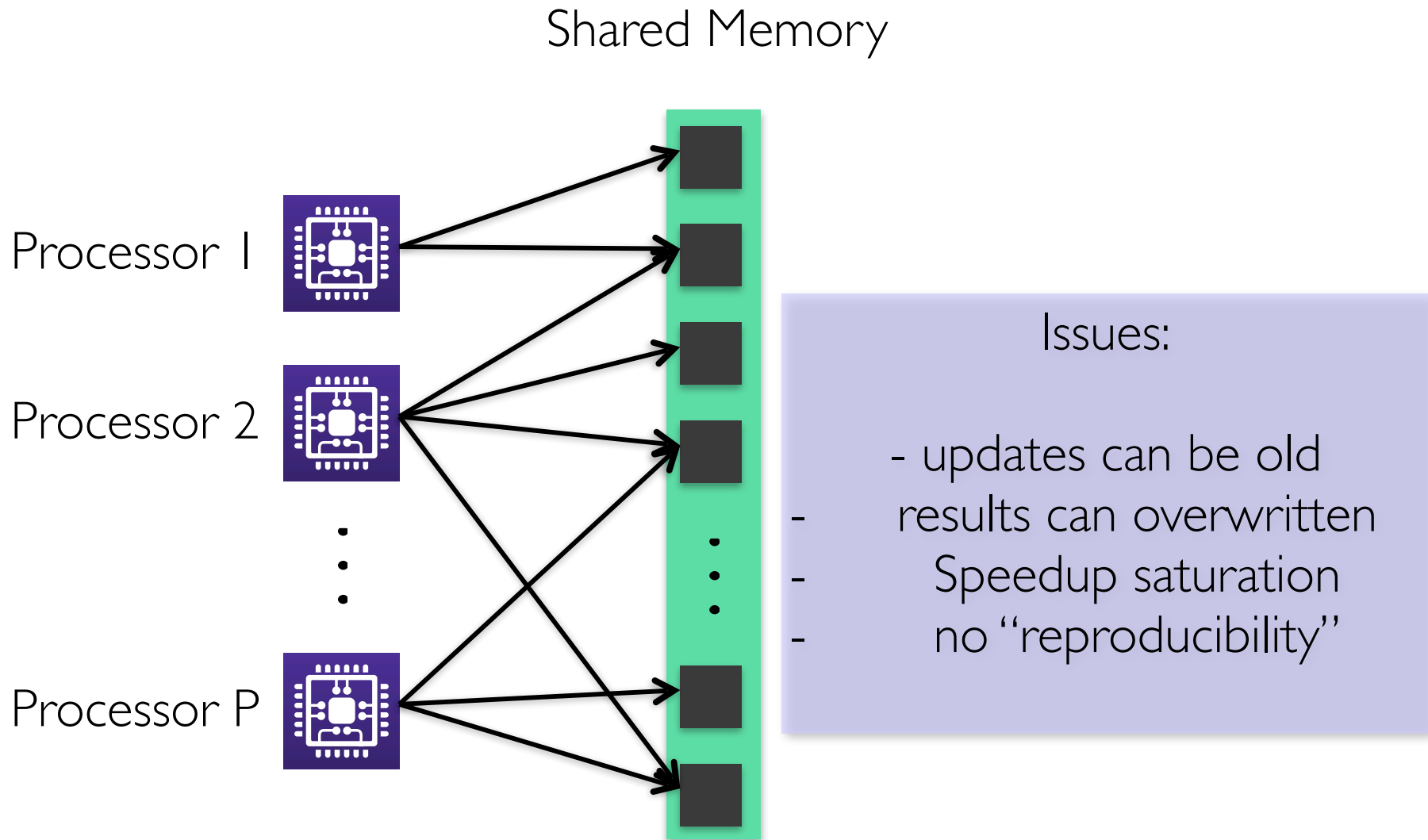
$x_v \leftarrow x_v + g_v$



Impact

Google Downpour SGD, Microsoft Project Adam use HOGWILD!
Renewed interest on async. optimization

Challenges in Hogwild!



Convergence analysis is usually a pain

A case against Asynchrony

Under review as a conference paper at ICLR 2017

REVISITING DISTRIBUTED SYNCHRONOUS SGD

Jianmin Chen[‡], Xinghao Pan[‡], Rajat Monga, Samy Bengio
Google Brain
Mountain View, CA, USA
{jmchen, xinghao, rajatmonga, bengio}@google.com

Rafal Jozefowicz
OpenAI
San Francisco, CA, USA
rafal@openai.com

ABSTRACT

Distributed training of deep learning models on large-scale training data is typically conducted with *asynchronous* stochastic optimization to maximize the rate of updates, at the cost of additional noise introduced from asynchrony. In contrast, the *synchronous* approach is often thought to be impractical due to idle time wasted on waiting for straggling workers. We revisit these conventional beliefs in this paper, and examine the weaknesses of both approaches. We demonstrate that a third approach, synchronous optimization with backup workers, can avoid asynchronous noise while mitigating for the worst stragglers. Our approach is empirically validated and shown to converge *faster* and to *better* test accuracies.

A case against Asynchrony

Serializability:

Hogwild Model \neq SGD Model

early conducted with *asynchronous* stochastic optimization to maximize the rate of updates, at the cost of additional noise introduced from asynchrony. In contrast, the *synchronous* approach is often thought to be impractical due to idle time wasted on waiting for straggling workers. We revisit these conventional beliefs in this paper, and examine the weaknesses of both approaches. We demonstrate that a third approach, synchronous optimization with backup workers, can avoid asynchronous noise while mitigating for the worst stragglers. Our approach is empirically validated and shown to converge *faster* and to *better* test accuracies.

Robustness During Inference

The Failures of Deep Learning

panda
58% confidence

+ .007 ×

=

gibbon
99% confidence

WHY?

The Failures of Deep Learning



WHY?

Robustness During Training

Robustness: a key Challenge

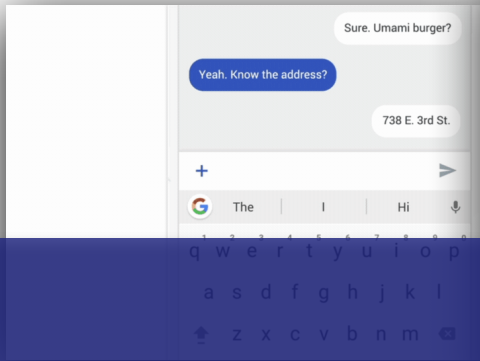


Fig. 5. Facebook global data center locations as of December 2017.

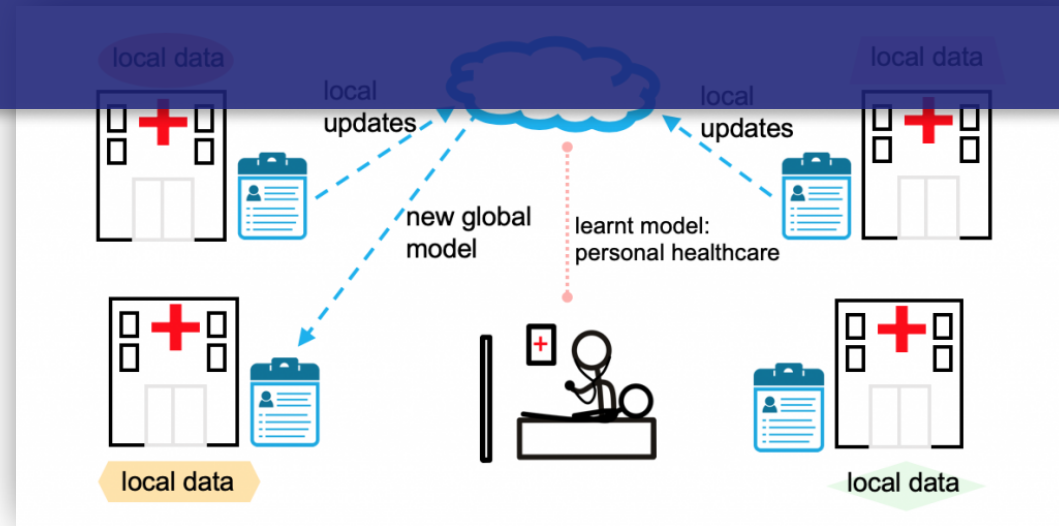
“For [...] training and inference [...] the importance of disaster-readiness cannot be underestimated.”

“Adversaries are constantly searching for new [...] ways to bypass our identifiers [...]”

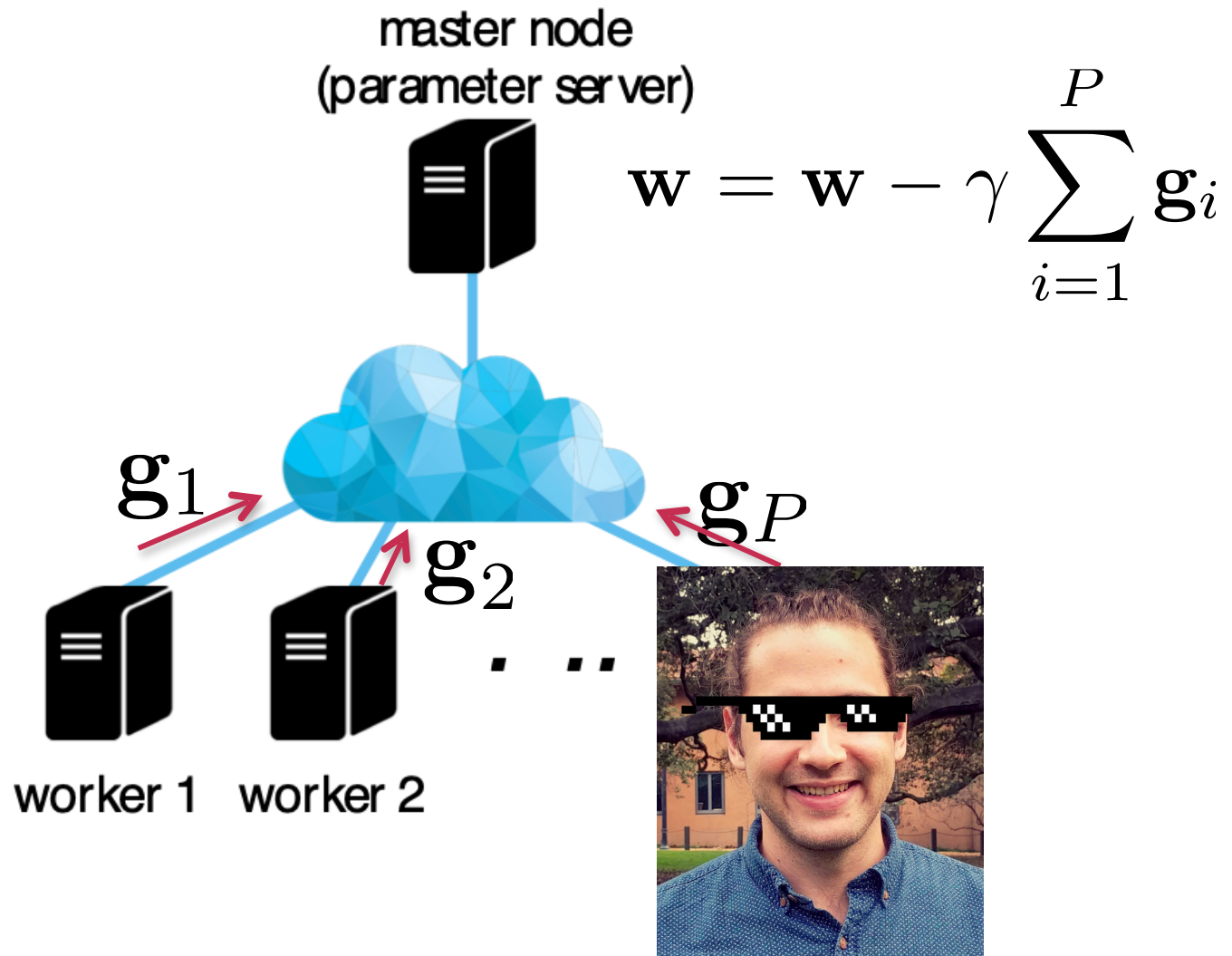
Federated Learning



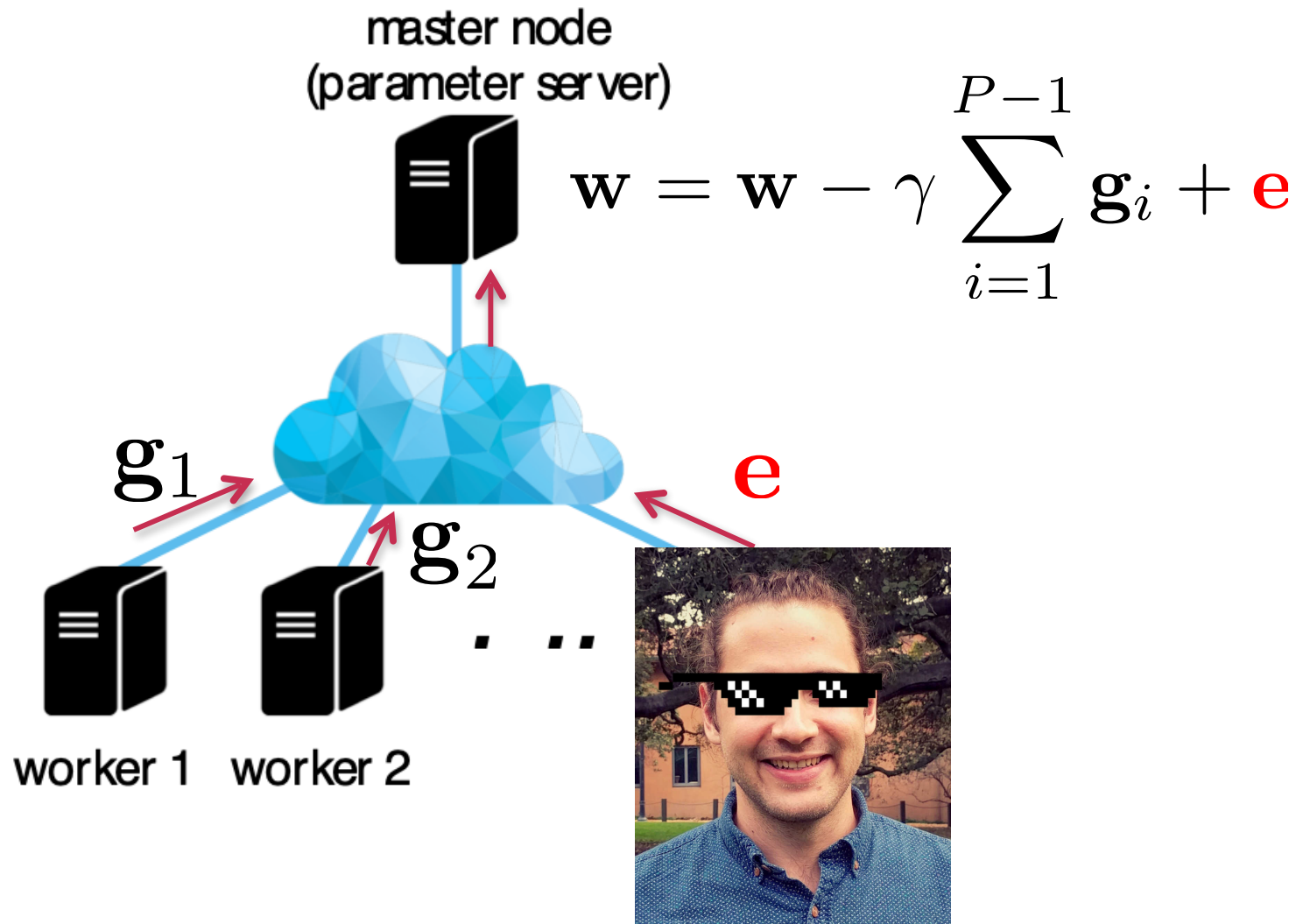
Data & models not inspected by central authority



is SGD Robust?



is SGD Robust?



minibatch SGD is not Robust

Can we build a robust version of SGD that is:

- + cheap
- + easy convergence
- + vanilla SGD model = robust SGD model

What's coming next

- Understanding mini-batch SGD performance
- Hogwild and theoretical challenges
- Model/Gradient Compression
- Low communication schemes
- Straggler Nodes and ways to overcome them
- Adversarial attacks, why they happen, how to defend

You want to be here

