

Introduction to ML Systems

2022 OxML Summer School – ML Fundamentals

Dingwen Tao
Washington State University

Today's Agenda

- Key Trends in Hardware for ML 14:00 – 14:15
- Data Parallel Training & Its Challenges 14:15 – 15:15
- Break 15:15 – 15:30
- Pipeline Parallelism 15:30 – 15:45
- Model Parallelism 15:45 – 16:15
- Spatial Parallelism 16:15 – 16:20
- Summary & Close 16:20 – 16:25

Hardware for ML

Key Drivers for Neural Network Success

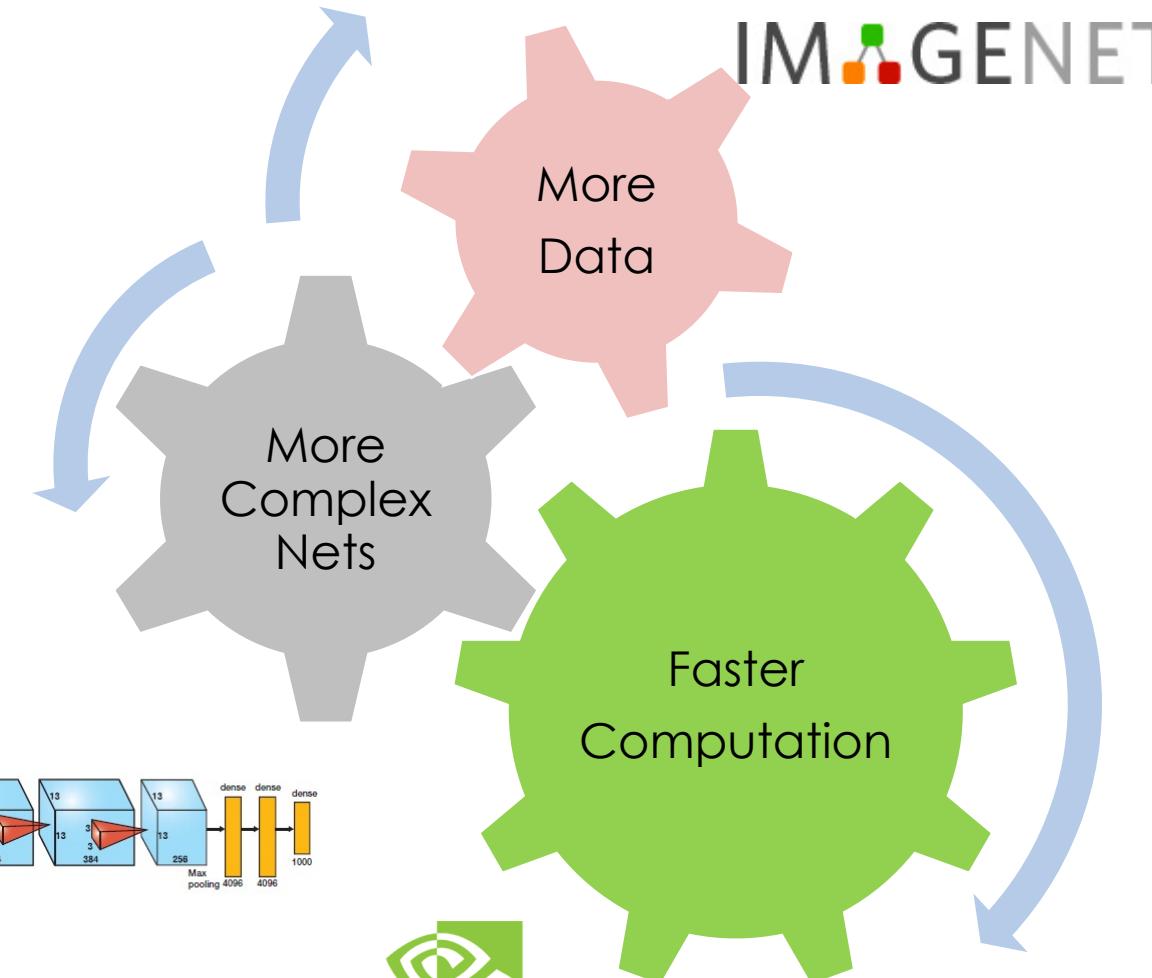
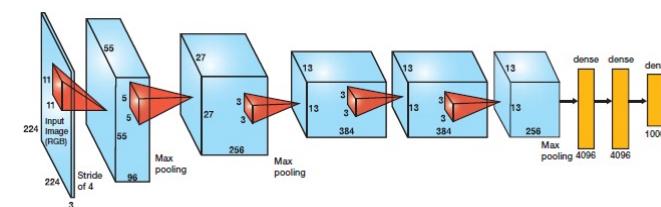
DARPA Neural Network Study Final Report (606 pages):

“After participating in this Study, my personal view is that **neural networks will provide the next major advance in computing technology.**”

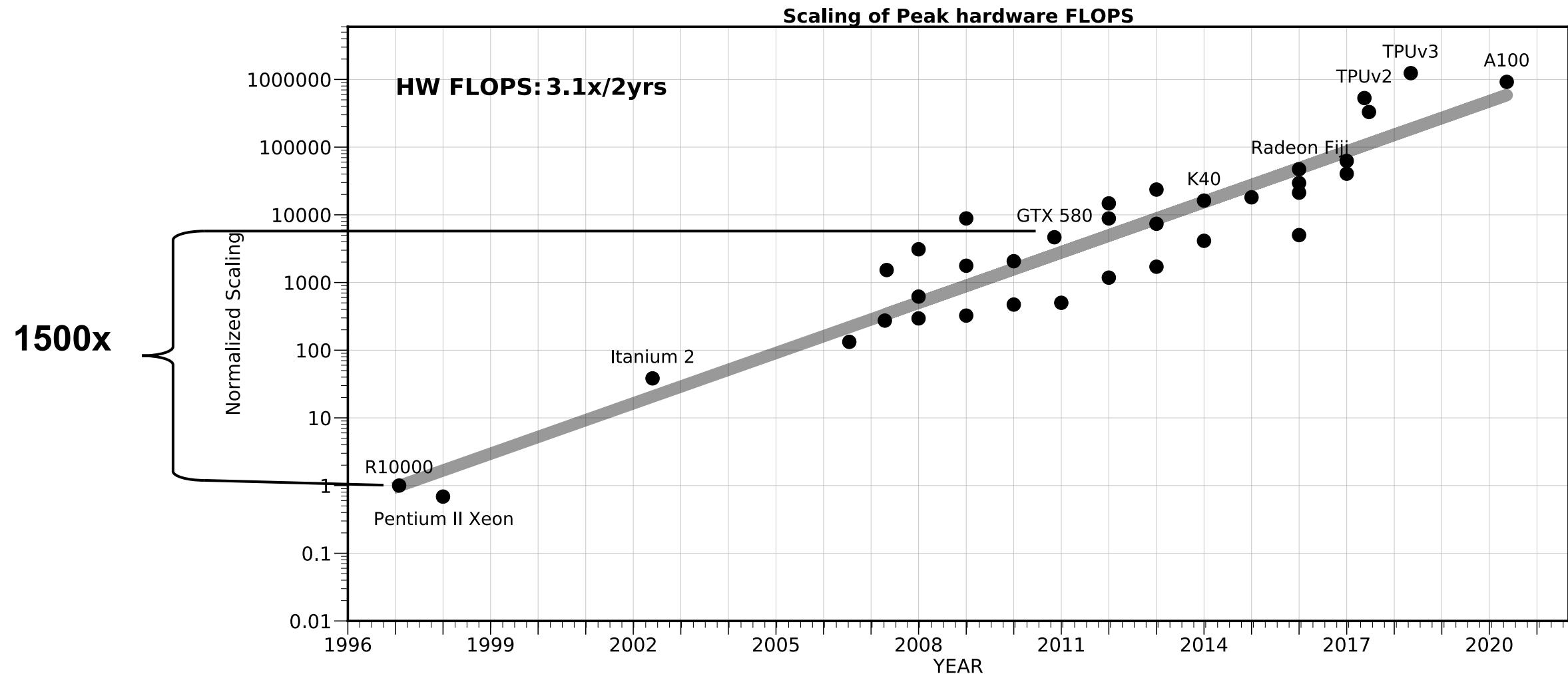
Dr. Jasper Lupo

DARPA, Washington, DC

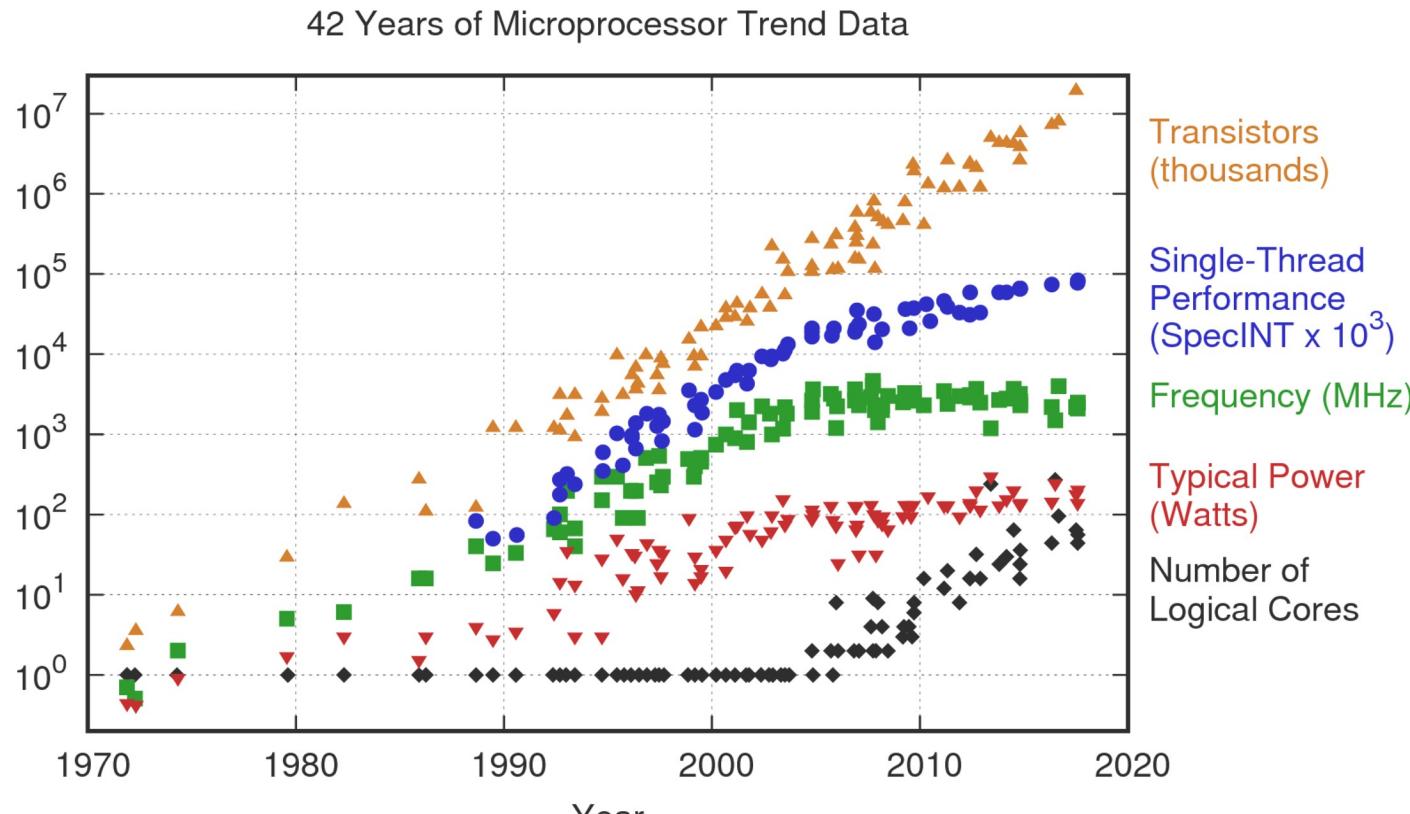
June, 1988



AlexNet vs Lenet5: 1000x More Compute



General Purpose Hardware Trend

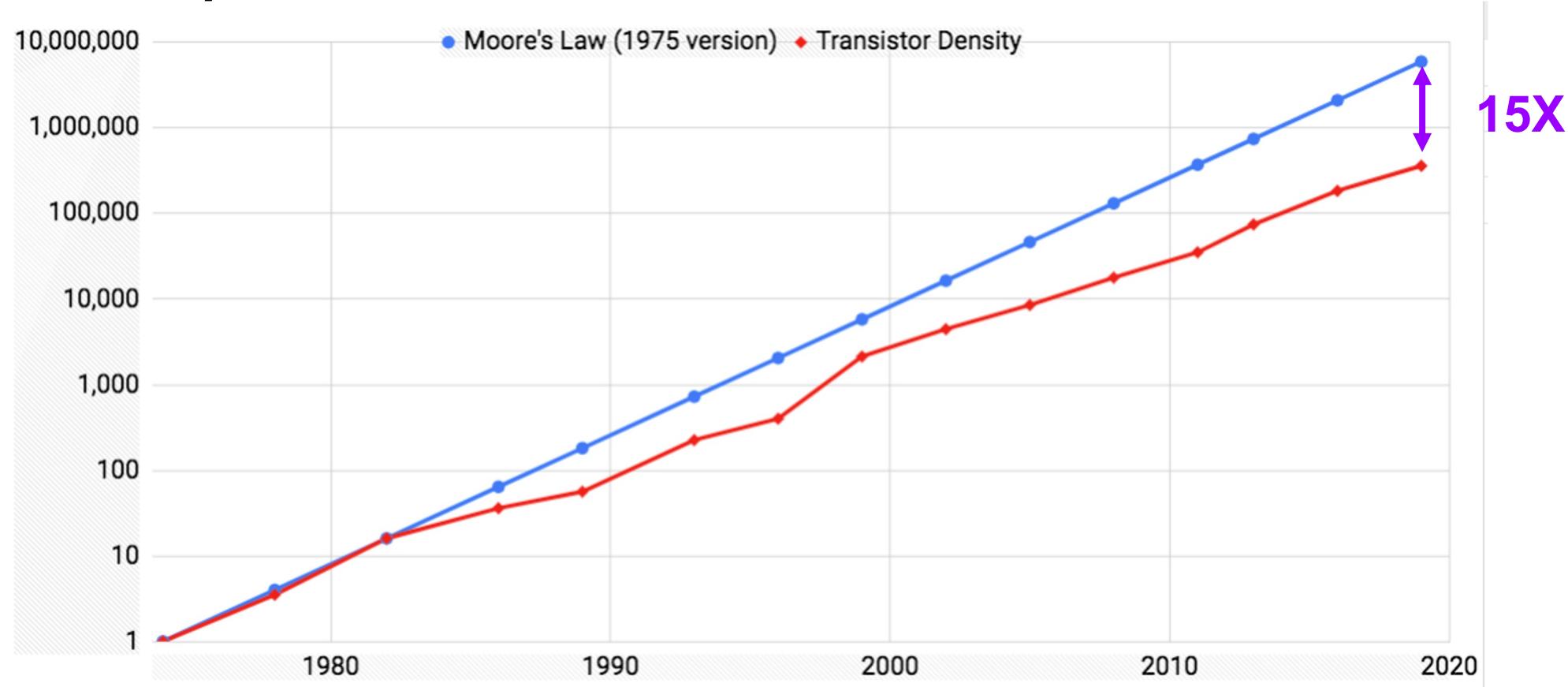


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

Key Observations

- # Transistors still increasing
- Single Core Performance Plateauing
- End of Dennard Scaling
- Distributed Computing

Common Fallacy: Moore's Law is Dead (it's not)



Moore, Gordon E. "No exponential is forever: but 'Forever' can be delayed!"
Solid-State Circuits Conference, 2003.

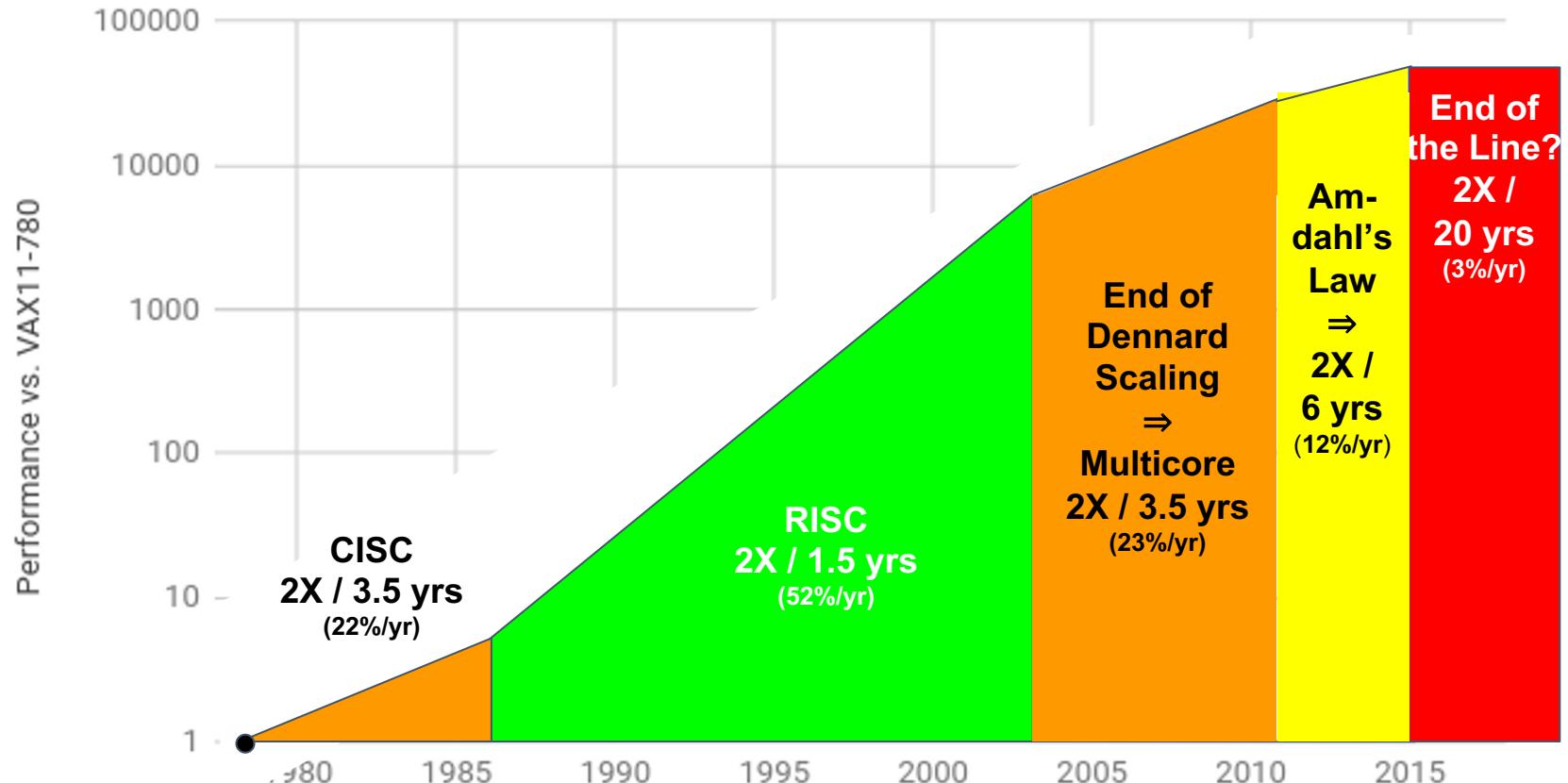
It is becoming increasingly difficult to push the boundary

Building a 3nm fab costs around \$20B. This is still economical given the \$600B ARR for the semi-conductor industry, but it is questionable how much farther we can push the limit.



But It has Slowed Down

40 years of Processor Performance



Based on SPECintCPU. Source: John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e. 2018

Domain Specific Accelerators

- John Hennessy and David Patterson,
“A New Golden Age for Computer Architecture,”
Communications of the ACM, February 2019



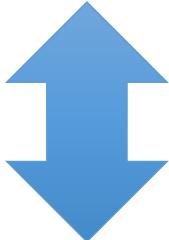
Domain Specific Accelerators



SnapDragon 835 (\rightarrow 845 \rightarrow 855)

$\sim 3.23 - 4$ MOPS/mW (835)

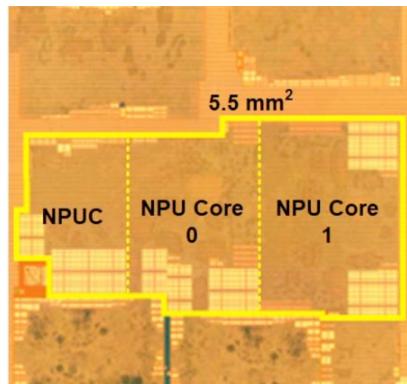
11– 16.6 GFLOPS SGEMM (835)



2.5K – 30K X increase in MOPs/mW

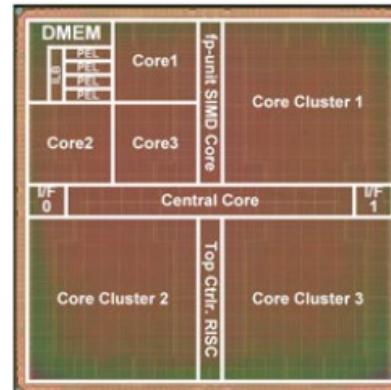
11,500 MOPS/mW

7.1 An 11.5TOPS/W 1024-MAC Butterfly Structure Dual-Core Sparsity-Aware Neural Processing Unit in 8nm Flagship Mobile SoC



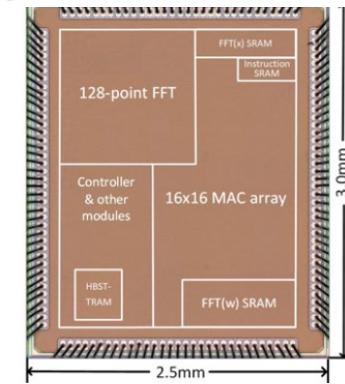
25,300 MOPS/mW

7.7 LNPU: A 25.3TFLOPS/W Sparse Deep-Neural-Network Learning Processor with Fine-Grained Mixed Precision of FP8-FP16



140,300 MOPS/mW

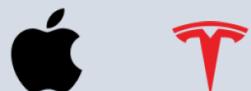
7.5 A 65nm 0.39-to-140.3TOPS/W 1-to-12b Unified Neural-Network Processor Using Block-Circulant-Enabled Transpose-Domain Acceleration with 8.1 \times Higher TOPS/mm² and 6T HBST-TRAM-Based 2D Data-Reuse Architecture



AI Chip Landscape

basicmi.github.io/AI-chip

Tech Giants/Systems



TOSHIBA



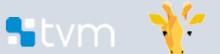
IC Vender/Fabless



IP/Design Service



Compiler



Startup in China



Startup Worldwide



more on <https://basicmi.github.io/AI-Chip/>

Benchmarks



中国人工智能产业链发展联盟
Artificial Intelligence Industry Alliance

Designing an accelerator

1) Accelerators are ONLY the First 80% of the Problem

The remaining 20%: SW development + Full system design

2) HW design shouldn't be about what can be built, rather what can be programmed

<https://eecs.wsu.edu/~dtao/download/Distributed-DL-PyTorch-Zhang.pdf>

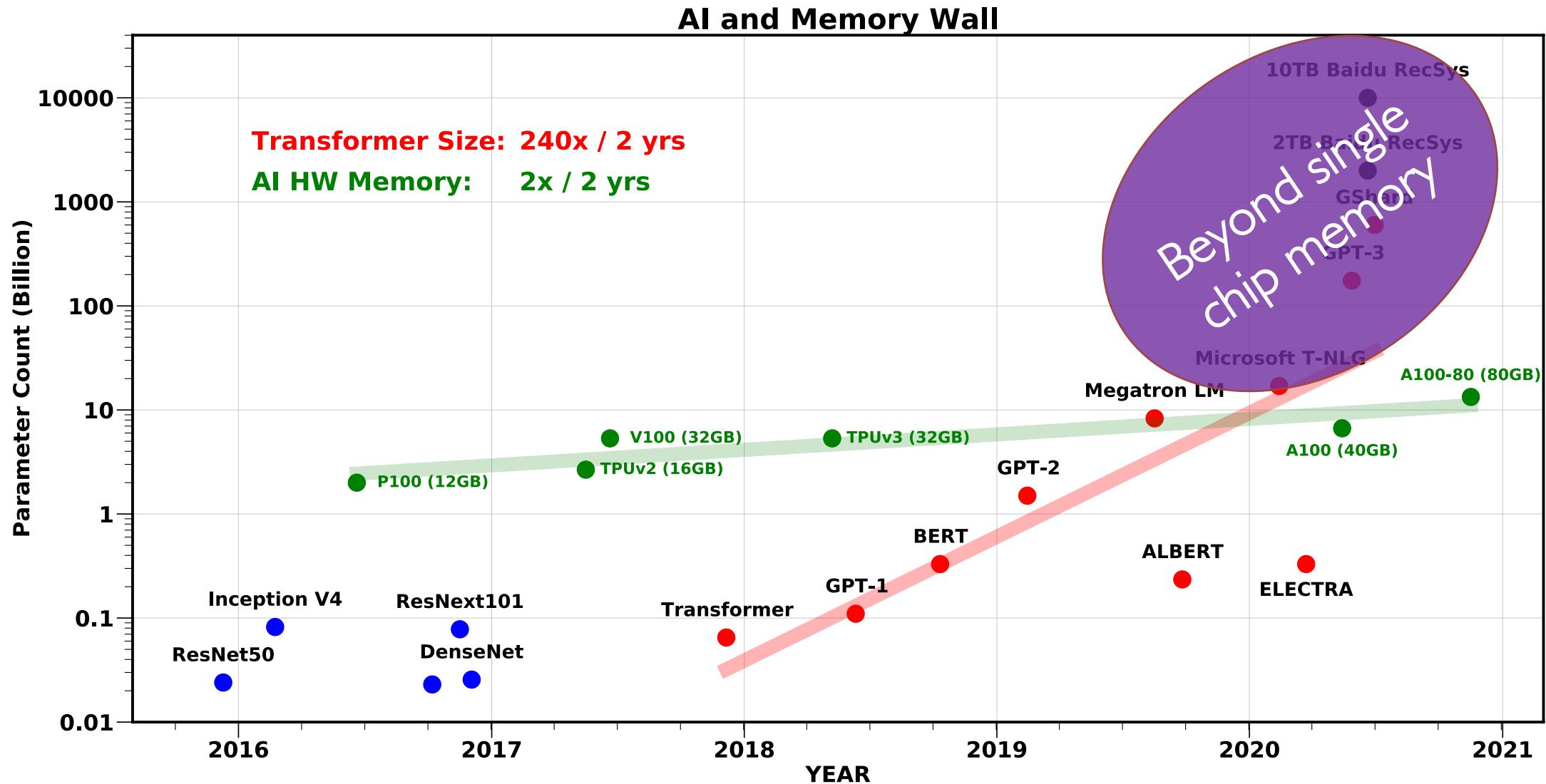
3) Deploy at scale? Distributed Deep Learning

Distributed Deep Learning

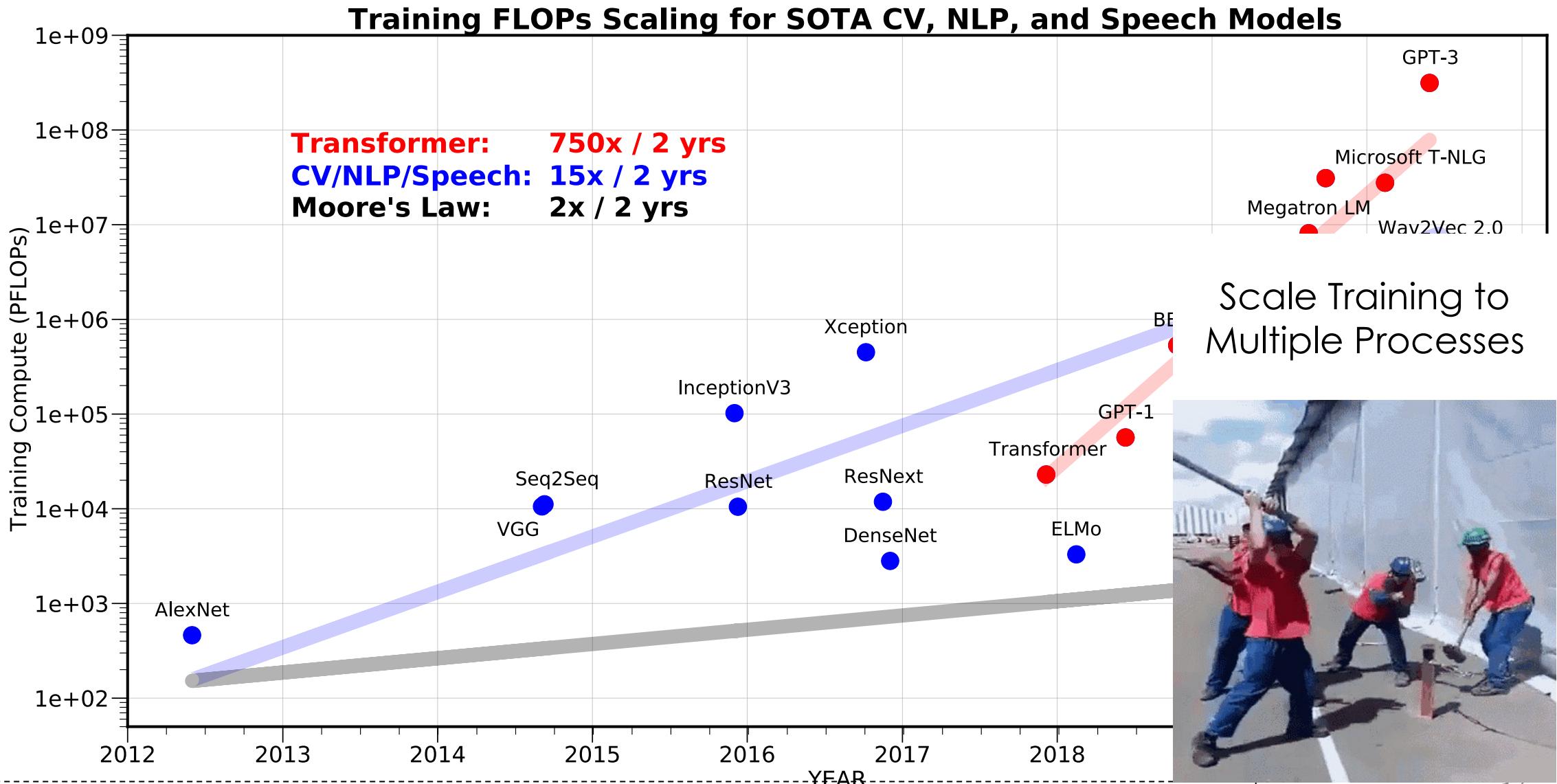
Distributed Training: What is it? & Why?

- **Distributed Training*** ~ Training across multiple devices
 - Different local and remote memory speeds / network
- Why do we need distributed training?
 - **Additional memory** (memory bandwidth) for larger model
 - “Need” to store weights + activations
 - Faster training by leveraging **parallel computation**
 - Reduce or eliminate **data movement**
 - Privacy → Federated Learning
 - Limited bandwidth to edge devices

Training Large Models

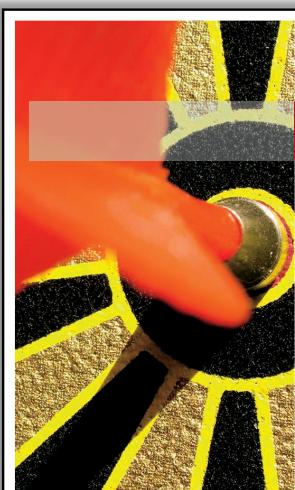


Faster Processing



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

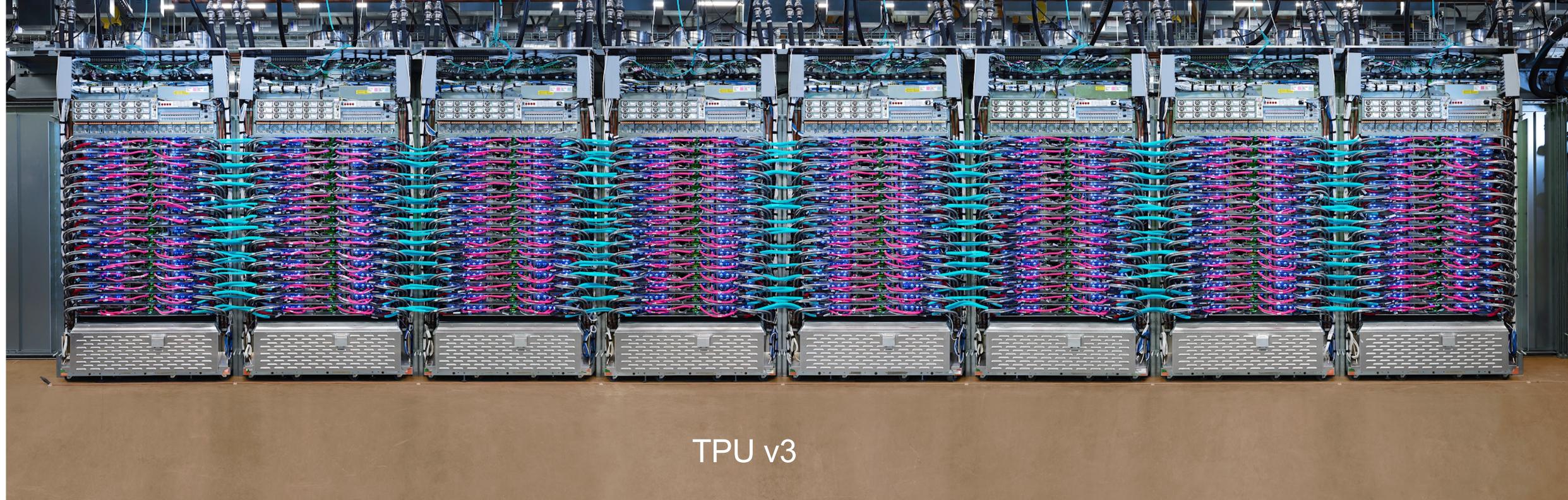


Example:
Scale is TPU's Primary Value Proposition



TPU Pod
64 2nd-gen TPUs
11.5 petaflops
4 terabytes of HBM memory

TPUv3



Ideal Metric of Success for Efficient Training

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

*Convergence
Machine Learning
Property*

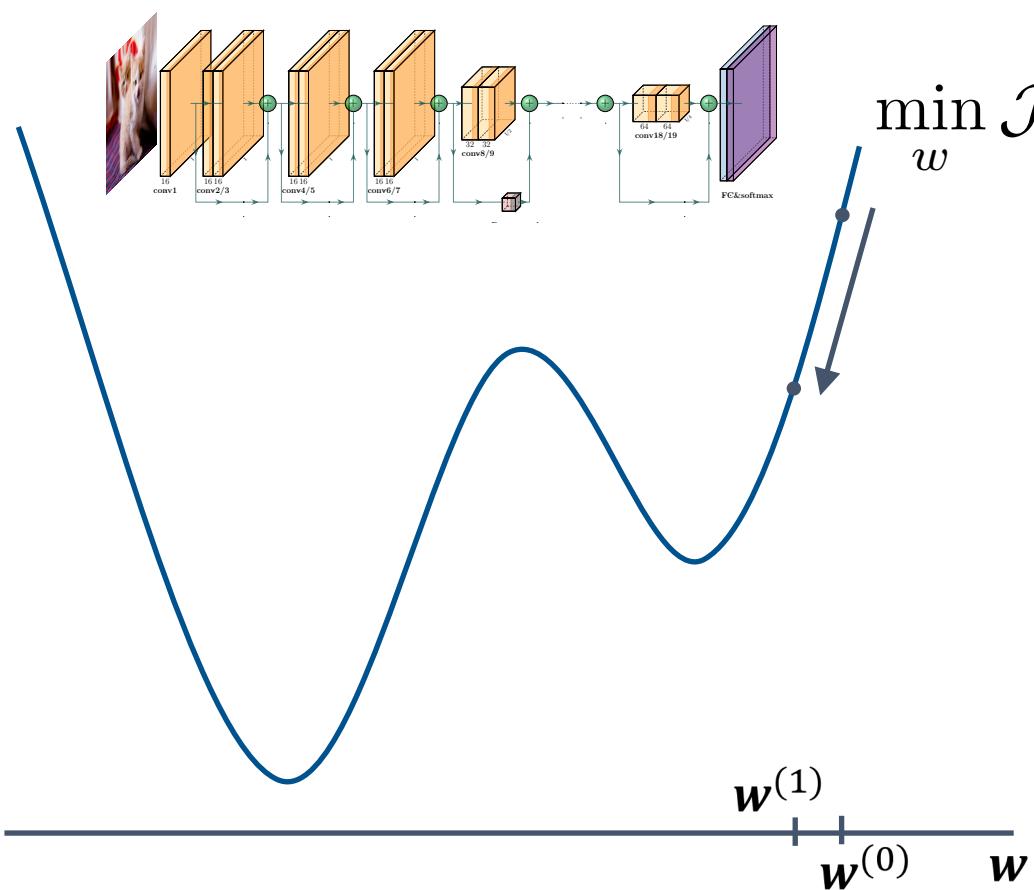
*Throughput
System
Property*

*Somewhat of a simplistic linear model. As we will later see there are many more moving parts to this

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?
 - **Cost:** Will obtaining a faster solution cost more money (power)?

Gradient Descent



$$\min_w \mathcal{J}(w) = \frac{1}{N} \sum_{i=1}^N cost(w, x_i)$$

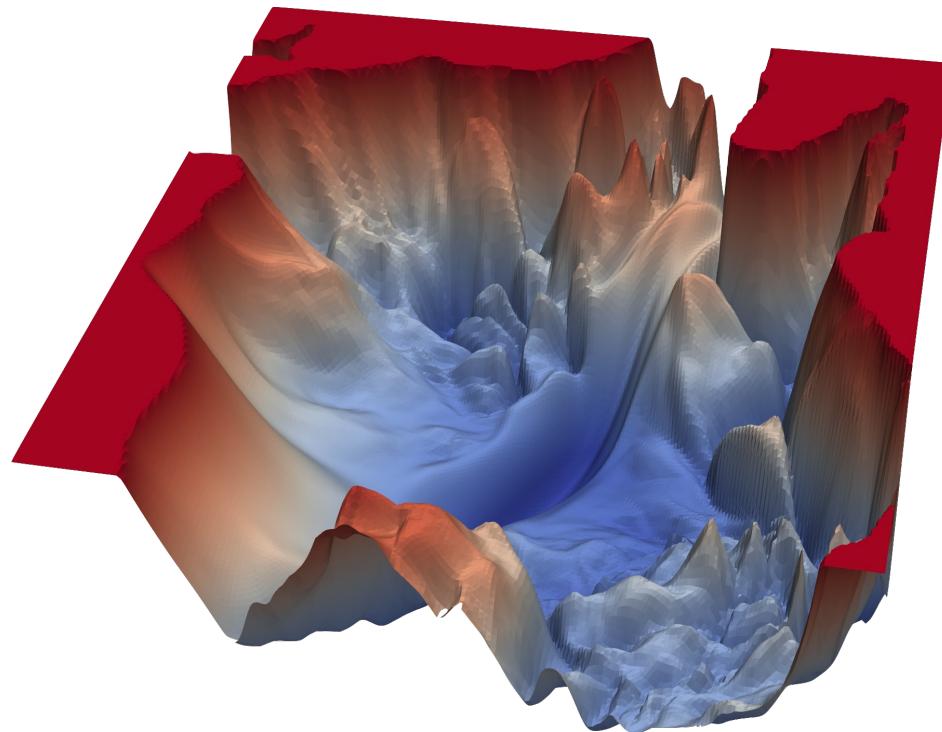
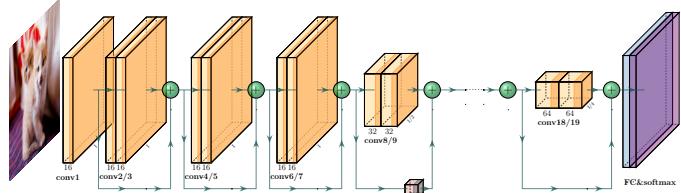
$$w^1 = w^0 - \alpha \frac{\partial \mathcal{J}(w^0)}{\partial w}$$

Learning rate Δw

Two key elements:

- The computed gradient: the direction
- The learning rate: how big a step do we take?

Stochastic Gradient Descent



$$\min_w \mathcal{J}(w) = \frac{1}{N} \sum_{i=1}^N \text{cost}(w, x_i)$$
$$w^1 = w^0 - \frac{\alpha}{B} \sum_{i=1}^B \frac{\partial \mathcal{J}(w^0)}{\partial w}$$

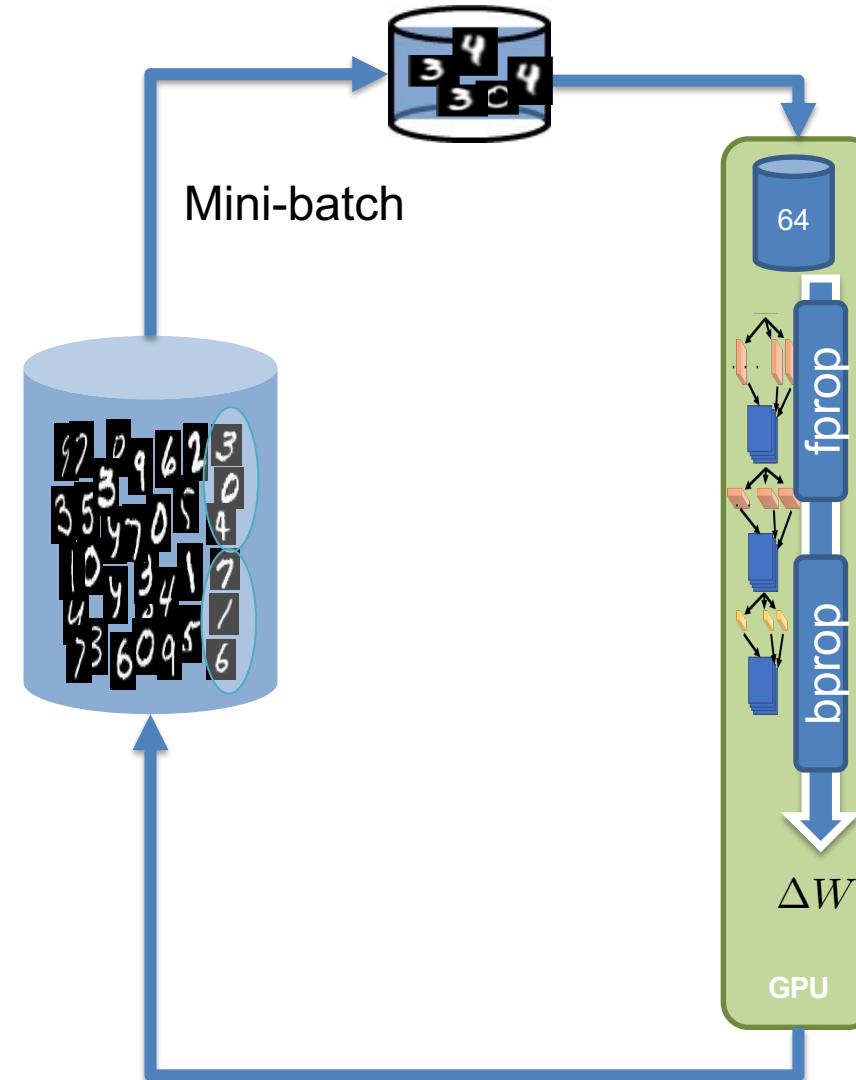
Learning rate Δw

Two key elements:

- The computed gradient: the direction
- The learning rate: how big a step do we take?

Synchronous Stochastic Gradient Descent

In every iteration of SGD we load a **random mini-batch of training data**, and compute the gradient.



$$\min_w \mathcal{J}(w) = \frac{1}{N} \sum_{i=1}^N cost(w, x_i)$$

$$w^1 = w^0 - \frac{\alpha}{B} \sum_{i=1}^B \frac{\partial \mathcal{J}(w^0)}{\partial w}$$

Δw

Parallelization Opportunities

Data Parallelism: Distribute the processing of data to multiple PEs.

$$w^1 = w^0 - \frac{\alpha}{B} \sum_{i=1}^B \frac{\partial \mathcal{J}(w^0)}{\partial w}$$

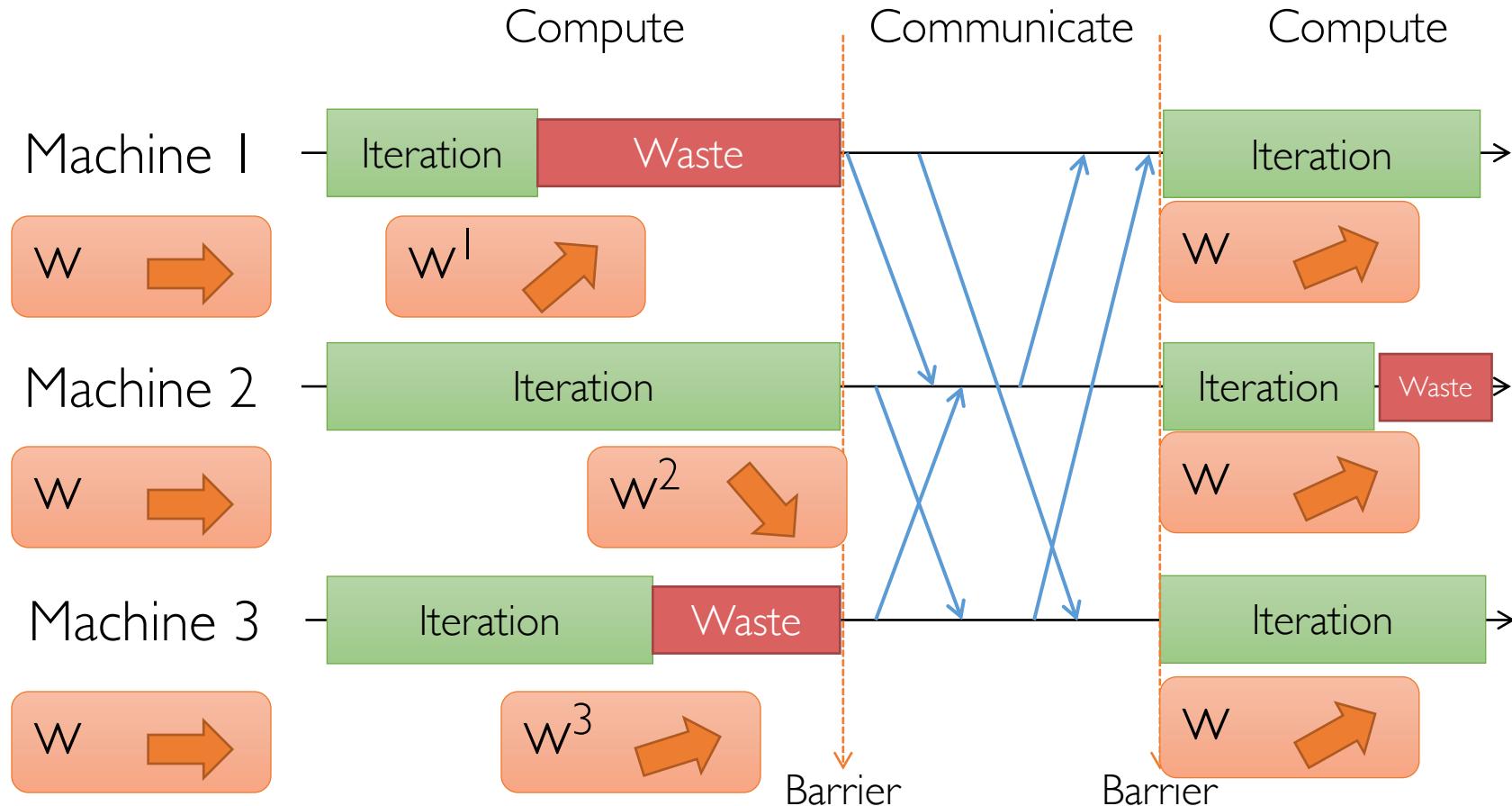
Model Parallelism: Break the model and distribute processing of every layer to multiple PEs

$$w^1 = w^0 - \frac{\alpha}{B} \sum_{i=1}^B \frac{\partial \mathcal{J}(w^0)}{\partial w}$$

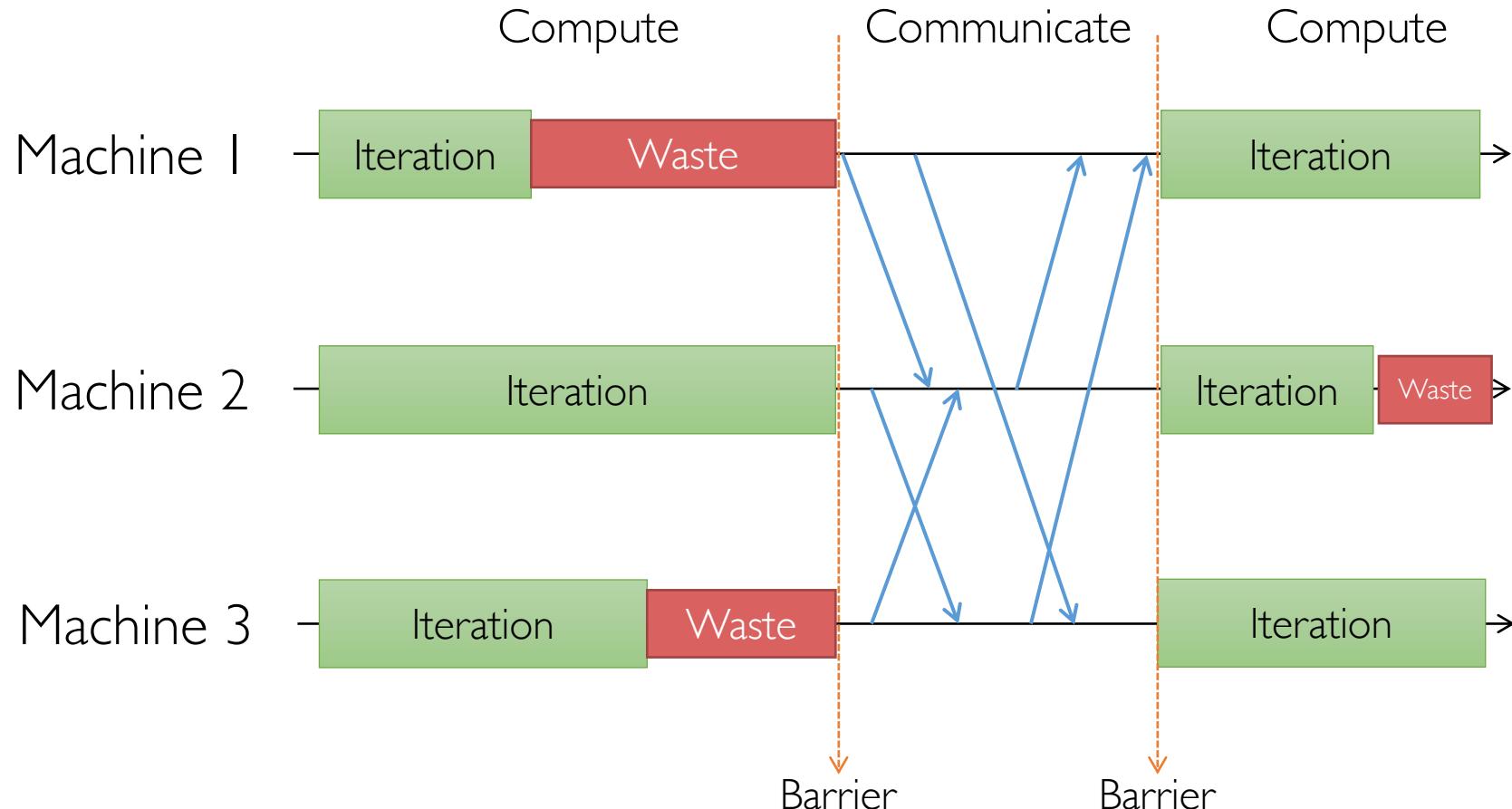
For either approach it is also possible to use **synchronous** or **asynchronous** updates

$$w^1 = w^0 - \frac{\alpha}{B} \sum_{i=1}^B \frac{\partial \mathcal{J}(w^0)}{\partial w}$$

Bulk Synchronous Parallel (BSP) Execution

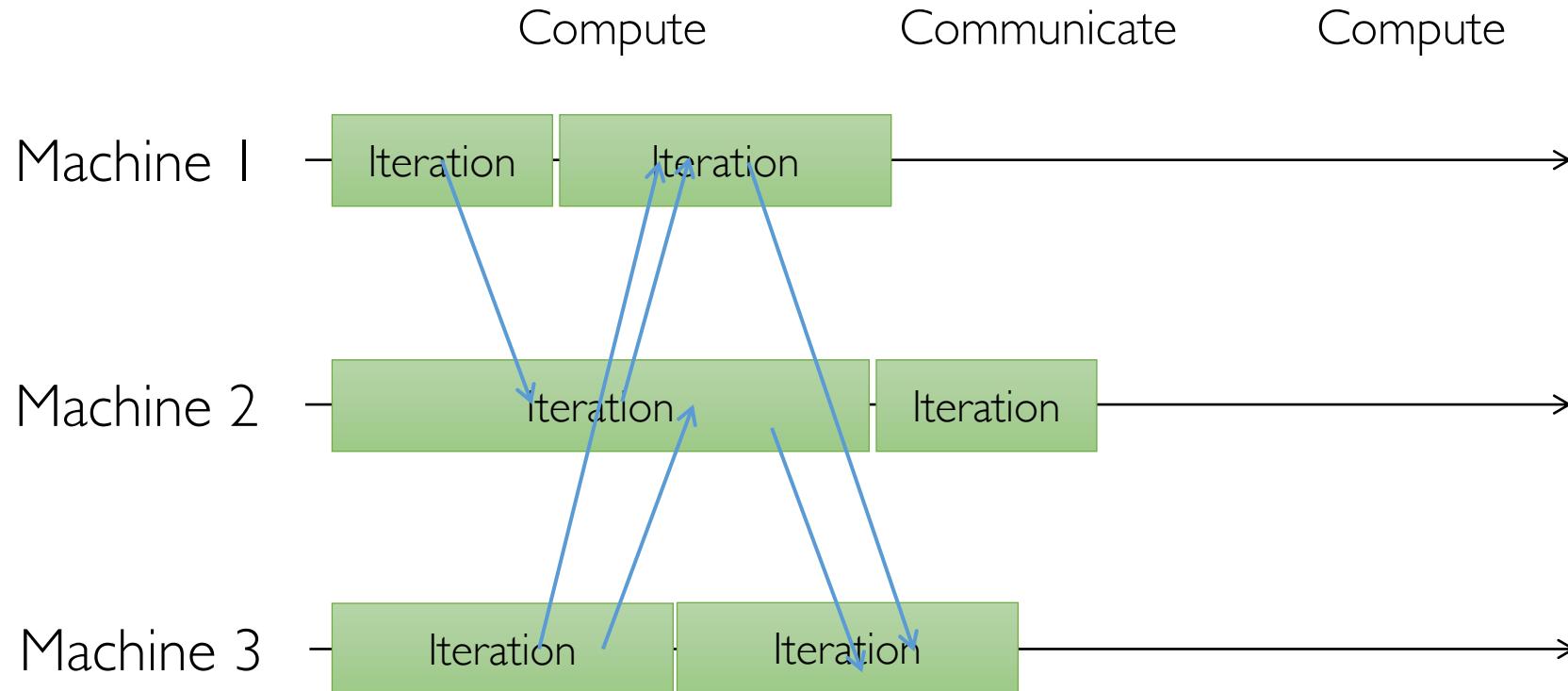


Bulk Synchronous Parallel (BSP) Execution



Enable more frequent coordination on parameter values

Asynchronous Execution

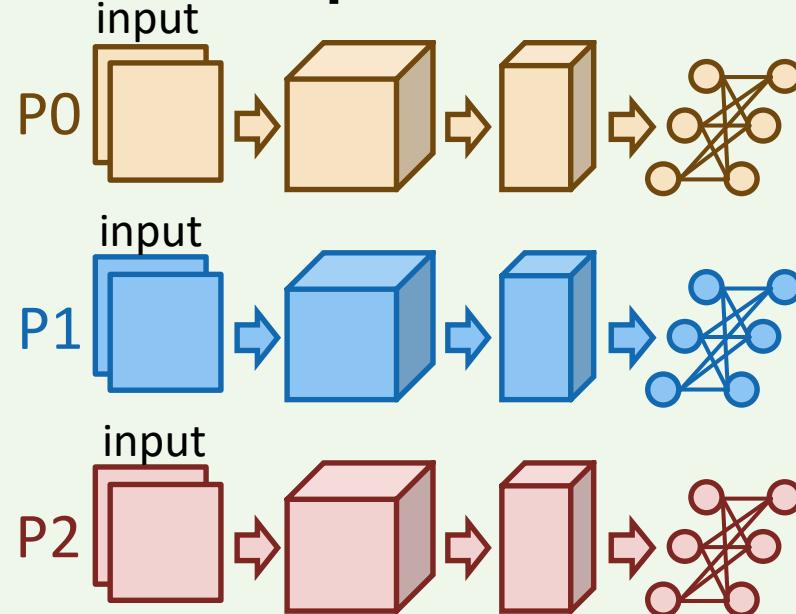


Enable more frequent coordination on parameter values, but often results in generalization loss. Today we will only focus on synchronous training.

Synchronous Data Parallel

Parallel and distributed training

Data parallelism



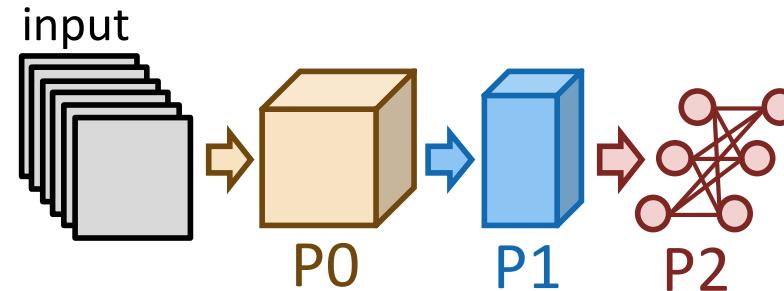
Pros:

- a. Easy to realize

Cons:

- a. Not work for large models
- b. High allreduce overhead

Pipeline parallelism



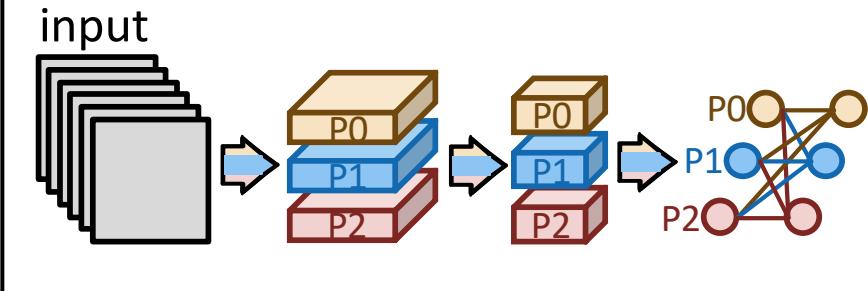
Pros:

- a. Make large model training feasible
- b. No collective, only P2P

Cons:

- a. Bubbles in pipeline
- b. Removing bubbles leads to stale weights

Model parallelism



Pros:

- a. Make large model training feasible

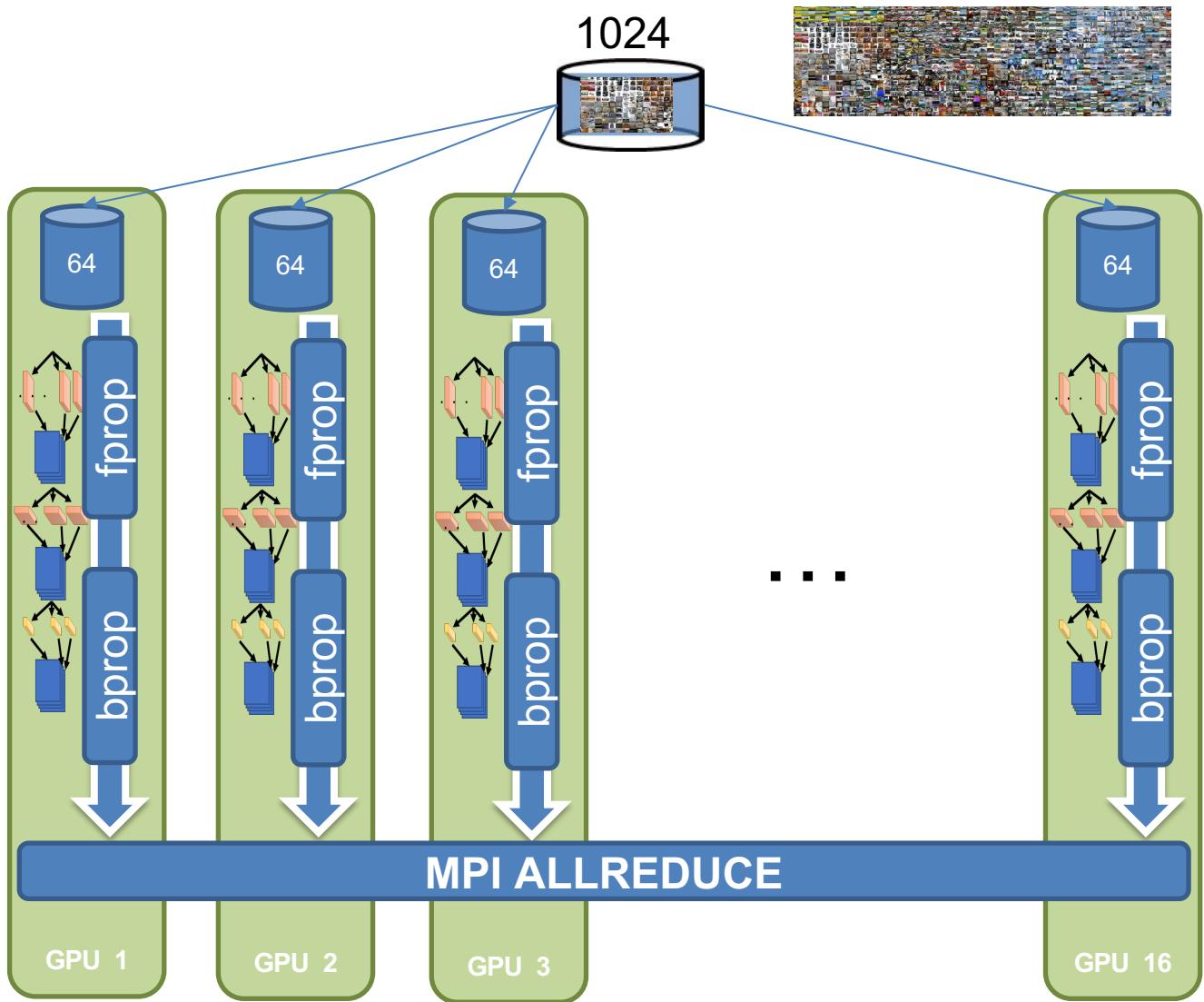
Cons:

- b. Communication for each operator (or each layer)

Synchronous Data Parallelism

- Compute the entire model on each processor
- Distribute the batch evenly across each processor:
 - 1024 batch distributed over 16 PEs: 64 images per GPU
- Communicate gradient updates through **allreduce**

$$w^1 = w^0 - \frac{\alpha}{B} \sum_{i=1}^B \frac{\partial \mathcal{J}(w^0)}{\partial w}$$



All Reduce

$$w^1 = w^0 - \frac{\alpha}{B} \sum_{i=1}^B \frac{\partial \mathcal{J}(w^0)}{\partial w}$$

$$a_1 = \sum_{i=1}^{B/4} \frac{\partial \mathcal{J}}{\partial w}$$

GPU 1

$$b_1 = \sum_{i=B/4}^{2B/4} \frac{\partial \mathcal{J}}{\partial w}$$

GPU 2

$$c_1 = \sum_{i=2B/4}^{3B/4} \frac{\partial \mathcal{J}}{\partial w}$$

GPU 3

$$d_1 = \sum_{i=3B/4}^B \frac{\partial \mathcal{J}}{\partial w}$$

GPU 4

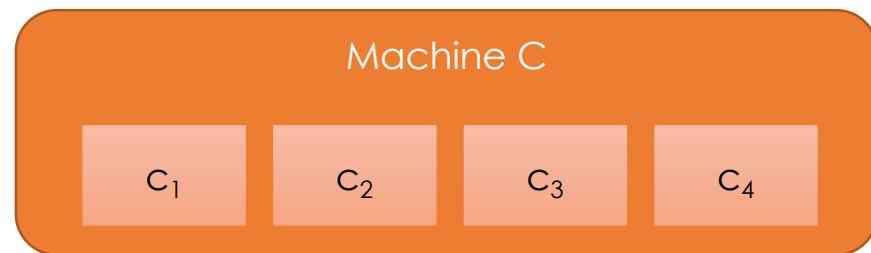
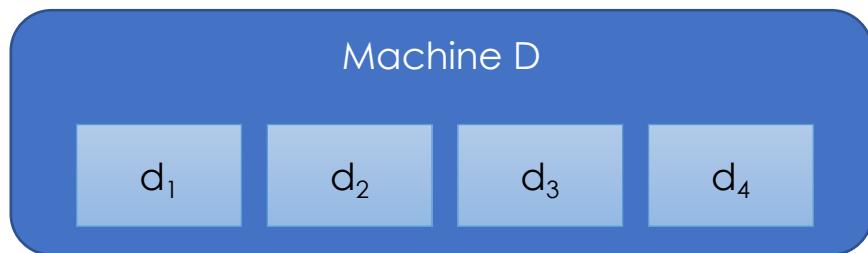
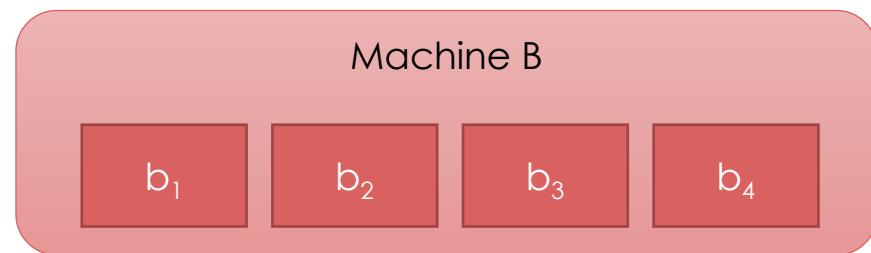
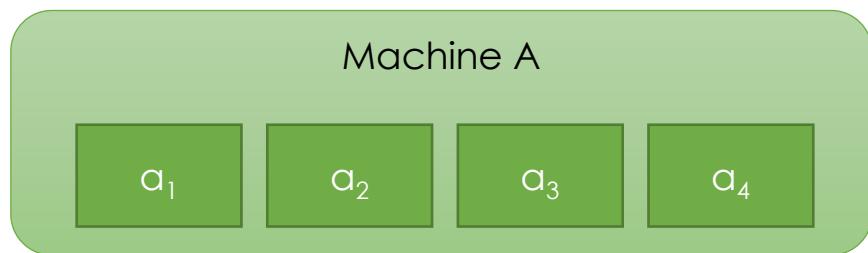
MPI ALLREDUCE

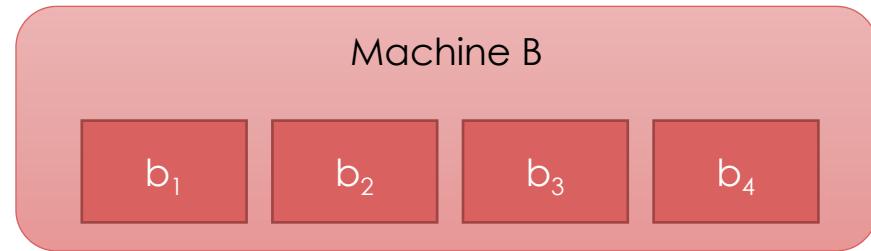
$$\sum_{i=1}^B \frac{\partial \mathcal{J}}{\partial w} = a_1 + b_1 + c_1 + d_1$$

All Reduce

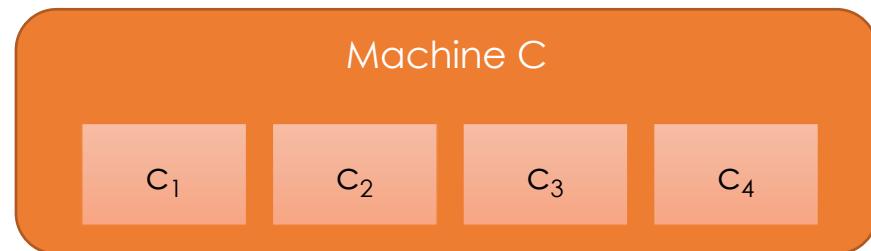
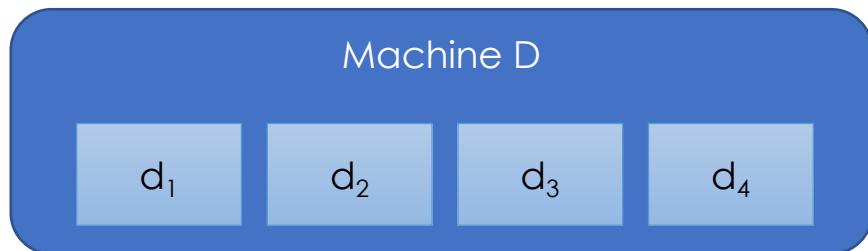
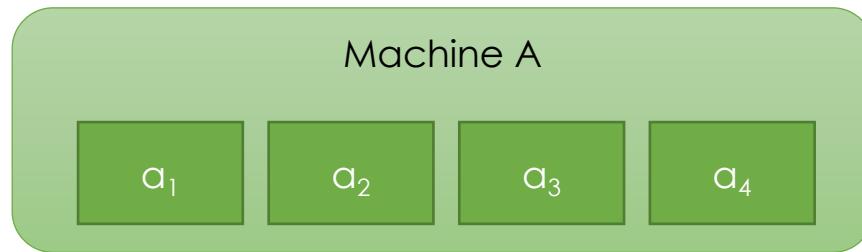
There are many different all reduce algorithms, each with their own trade offs.

For simplicity, assume our model has 4 layers, and is trained on $P=4$ machines



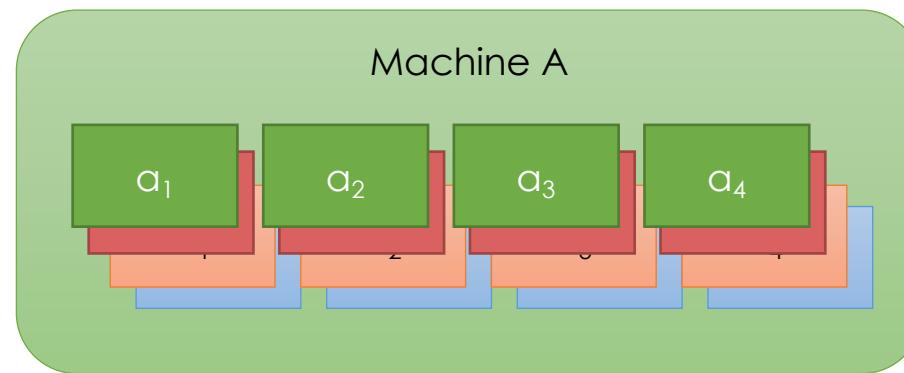


Parameter Server (Single Master All-Reduce)



Parameter Server

Machine B



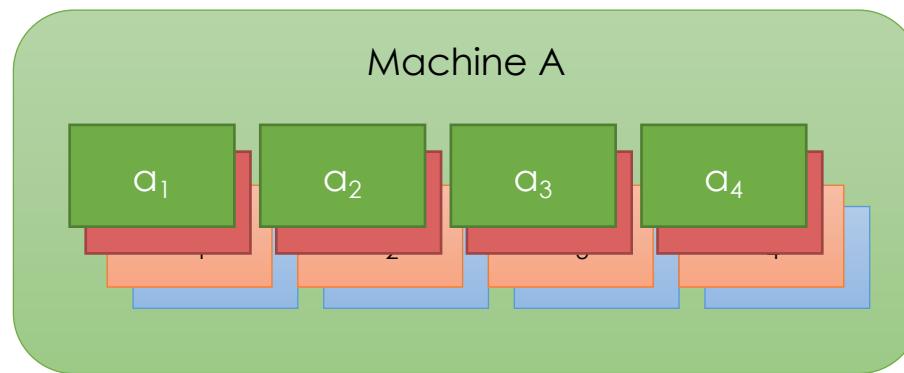
Sends **(P-1) * N** Data
➤ **P** Machines
➤ **N** Parameters

Machine D

Machine C

Parameter Server

Machine B



Sends **(P-1) * N** Data
➤ **P** Machines
➤ **N** Parameters

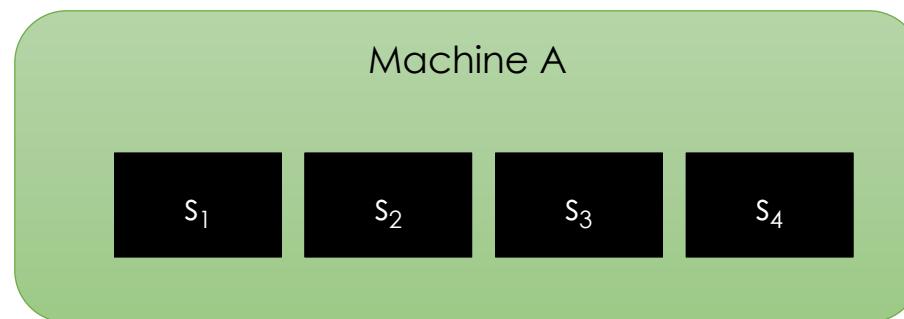
$$s_i = a_i + b_i + c_i + d_i$$

Machine D

Machine C

Parameter Server

Machine B

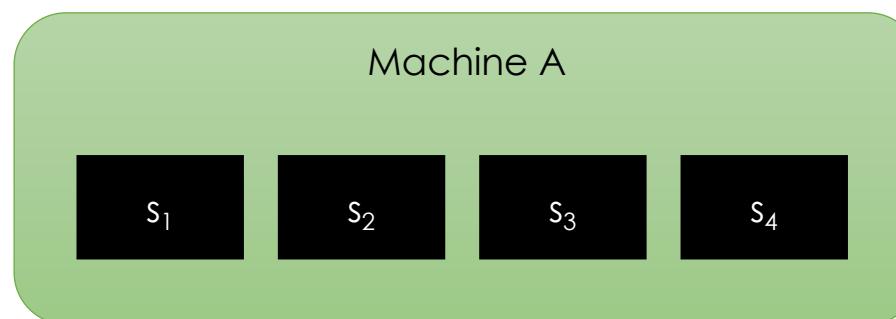
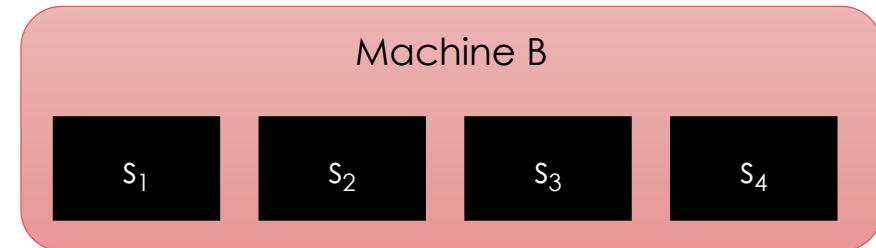


Communicate **(P-1) * N** Data
➤ **P** Machines
➤ **N** Parameters

$$s_i = a_i + b_i + c_i + d_i$$



Parameter Server



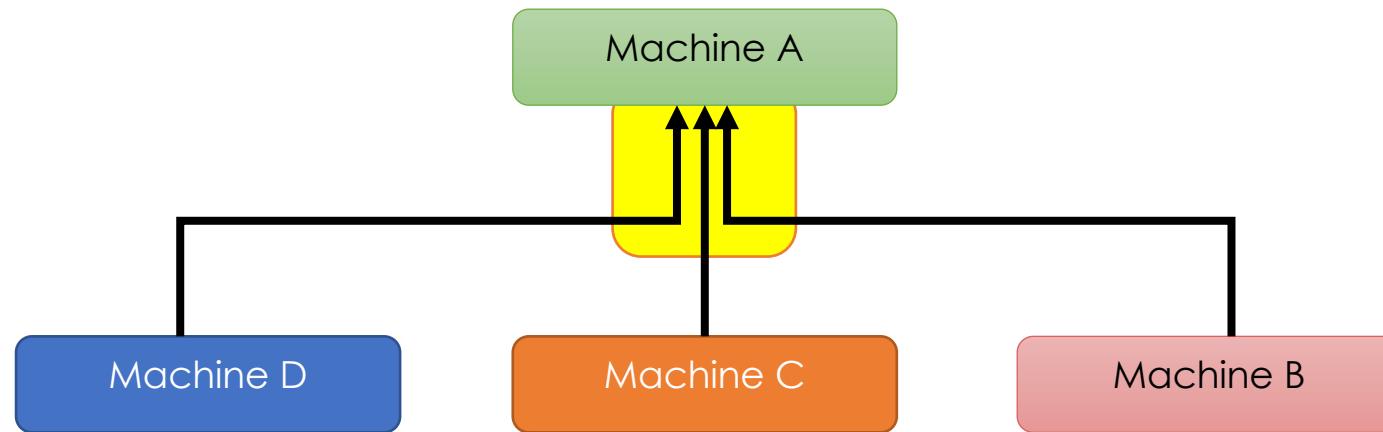
*2
Communicate $(P-1) * N$ Data
➤ P Machines
➤ N Parameters

$$s_i = a_i + b_i + c_i + d_i$$



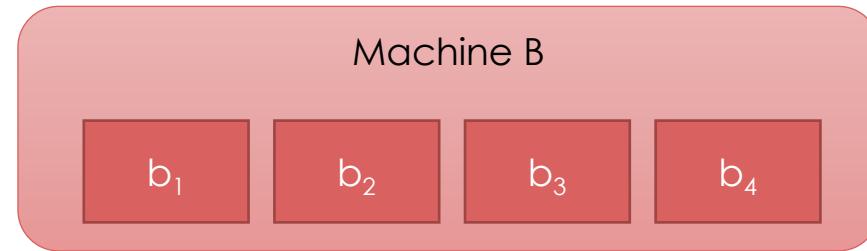
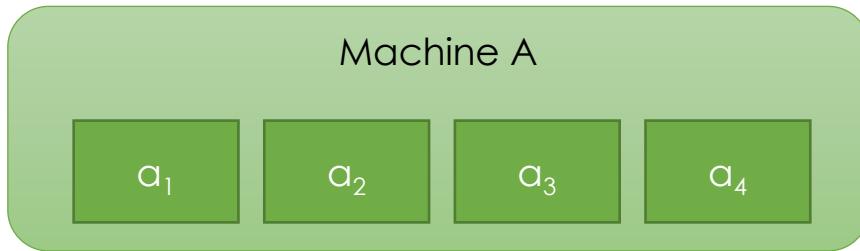
Parameter Server

Comm $(P-1) * N$ Data
➤ P Machines
➤ N Parameters

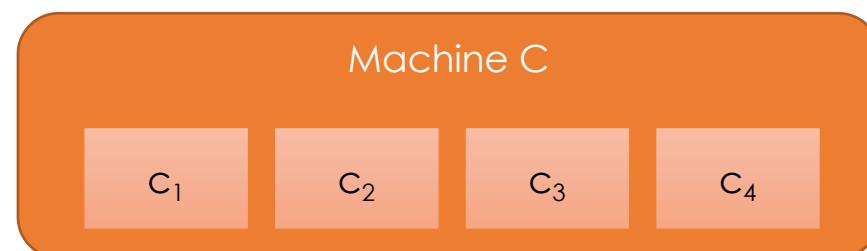
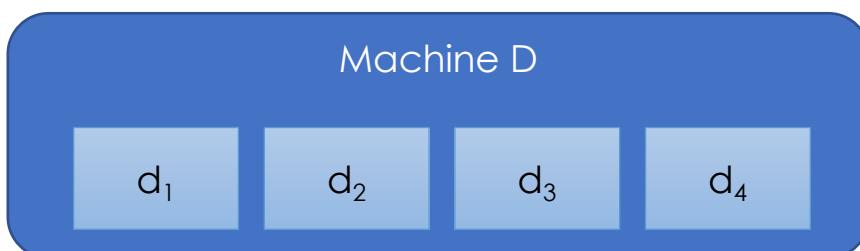


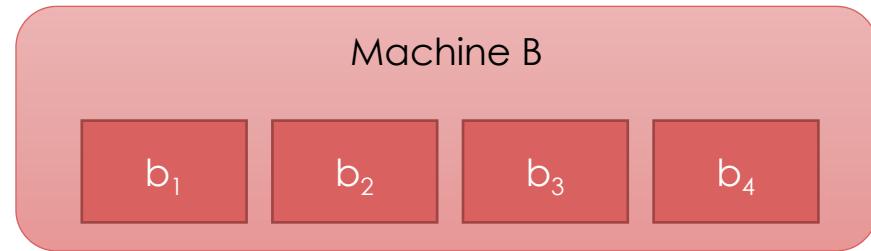
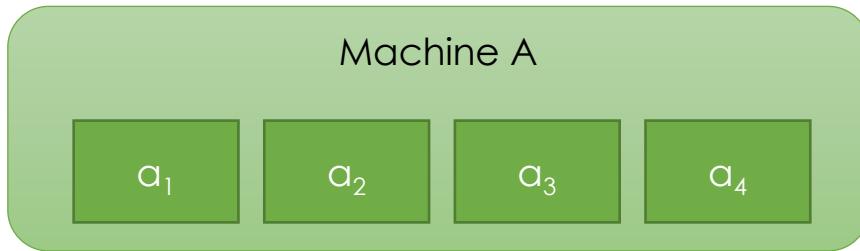
Issues?

- High **fan-in** on Machine A
- **$(P-1) * N$ Bandwidth** for Machine A



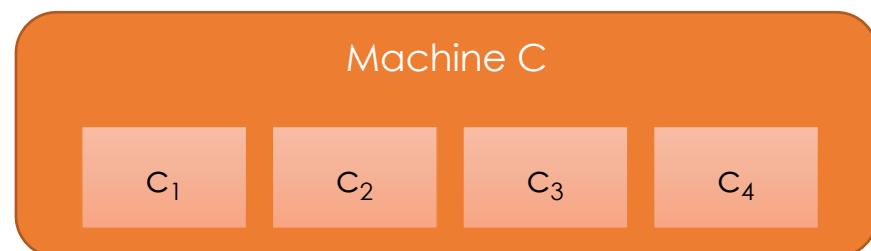
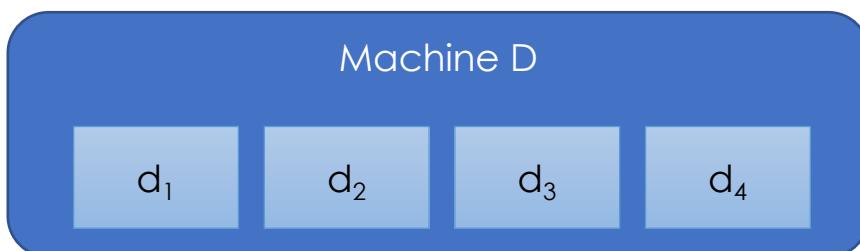
Parameter Server All Reduce

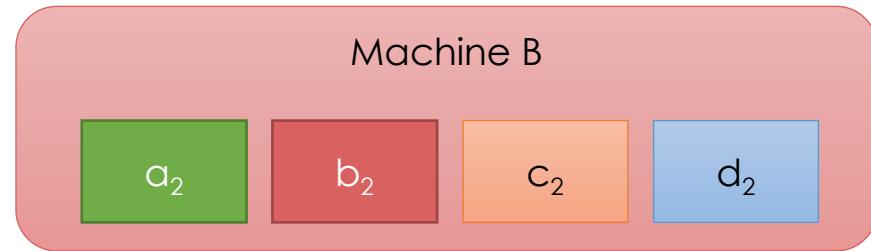
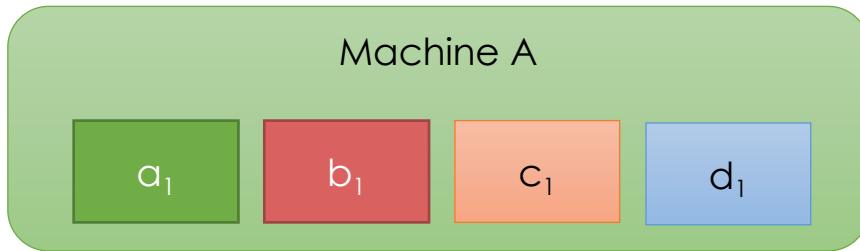




Send each entry to parameter server for that entry.

- Key 1 → A
- Key 2 → B
- Key 3 → C
- Key 4 → D

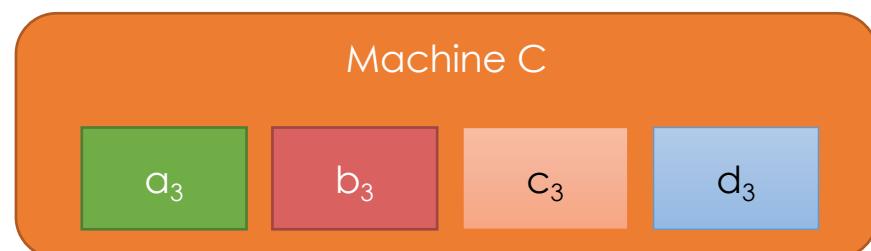
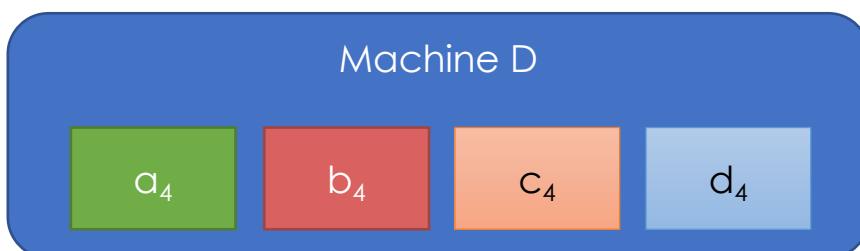




Each machine sends N/P data to all other machines.

(P-1) * N/P

- P Machines
- N Parameters



Machine A

s_1

Machine B

s_2

Compute local sum on each machine

$$s_i = a_i + b_i + c_i + d_i$$

Machine D

s_4

Machine C

s_3

Machine A

s_1

Machine B

s_2

Each machine broadcasts* the sum (N/P data size) to all other machines.

(P-1) * N/P

- **P** Machines
- **N** Parameters

Machine D

s_4

Machine C

s_3

* Technically All Gather based on MPI communication definition

Machine A



Machine B



Total Communication per machine:

$2 * (P-1) * N/P$ (roughly independent of P)

- **P** Machines
- **N** Parameters

Machine D

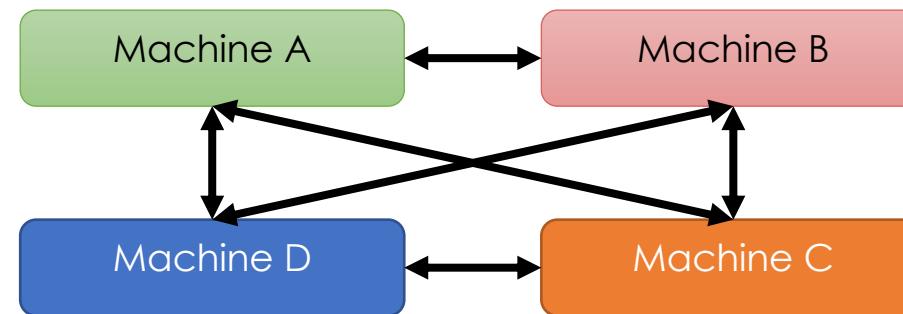


Machine C

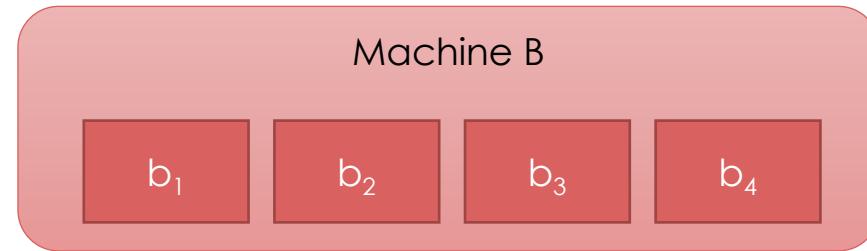
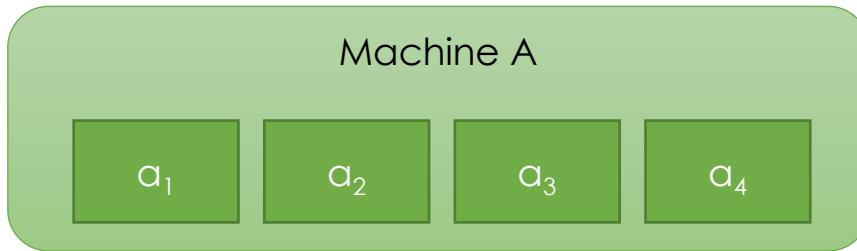


Parameter Server All-Reduce

- Same amount of total data transmitted as before, but spread evenly across all machines instead of just one

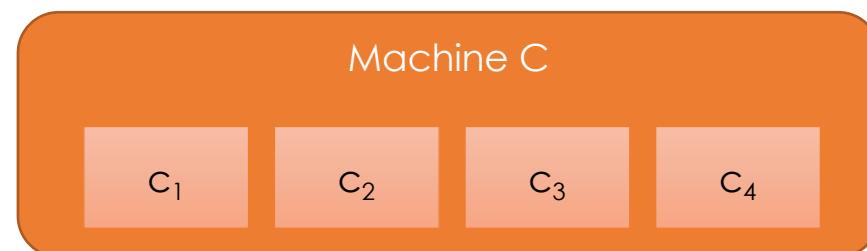
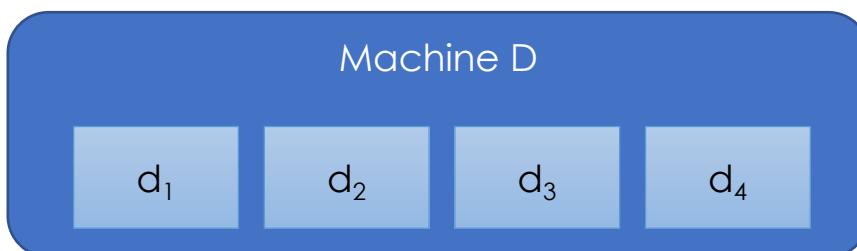


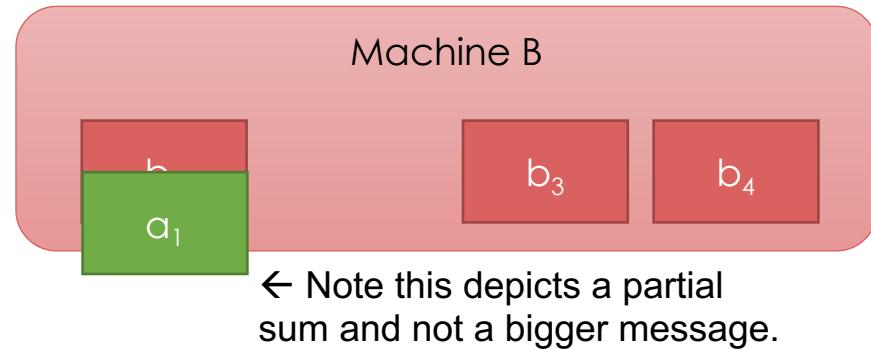
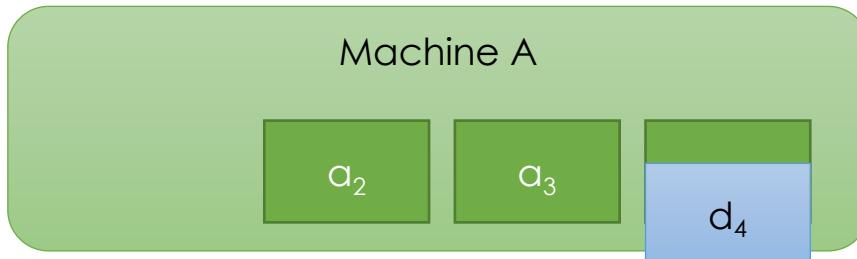
- Same **high fan-in** ($P-1$)
- **Reduced** Inbound Bandwidth = $2*(P-1)N/P$
 - Previously $2*(P-1)*N$ for the parameter server



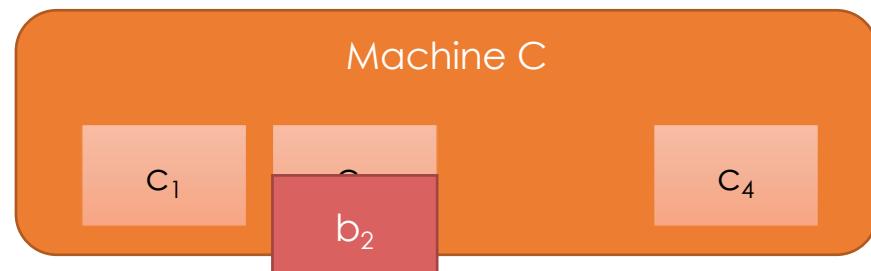
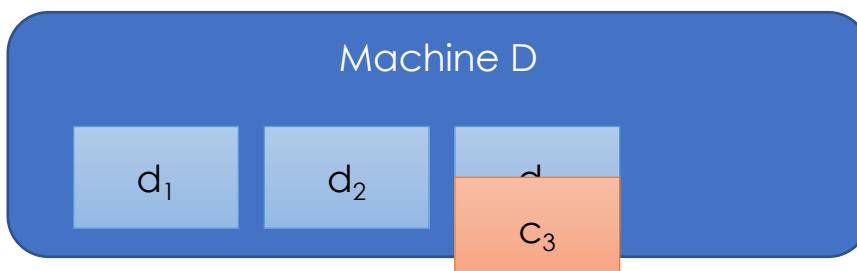
Ring All Reduce

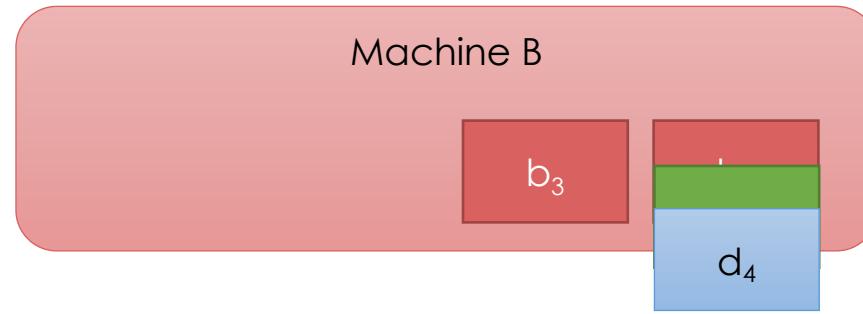
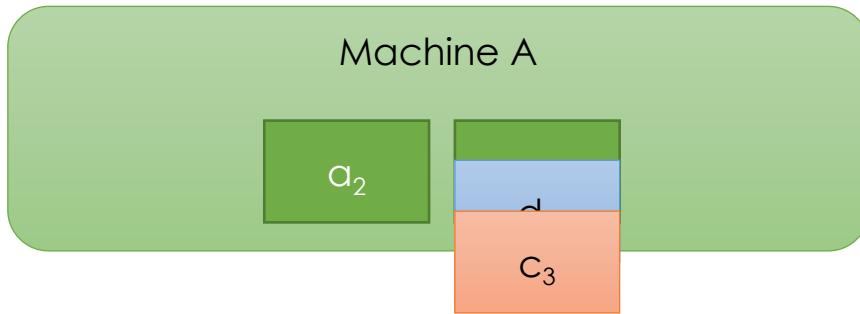
Send messages in a ring to reduce fan-in.



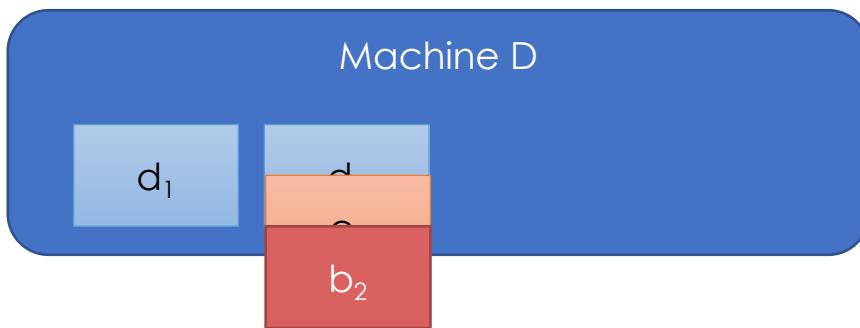


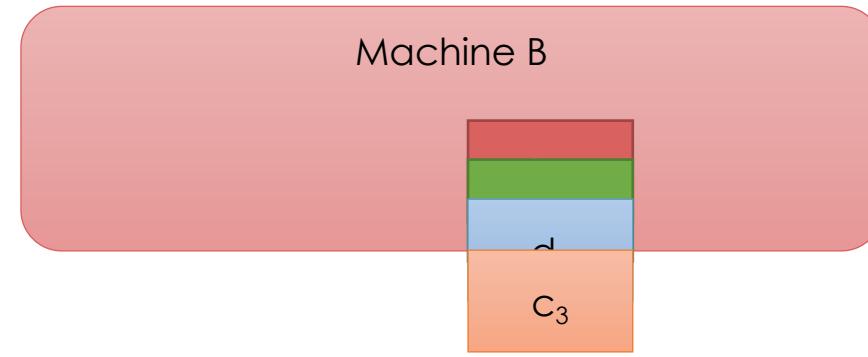
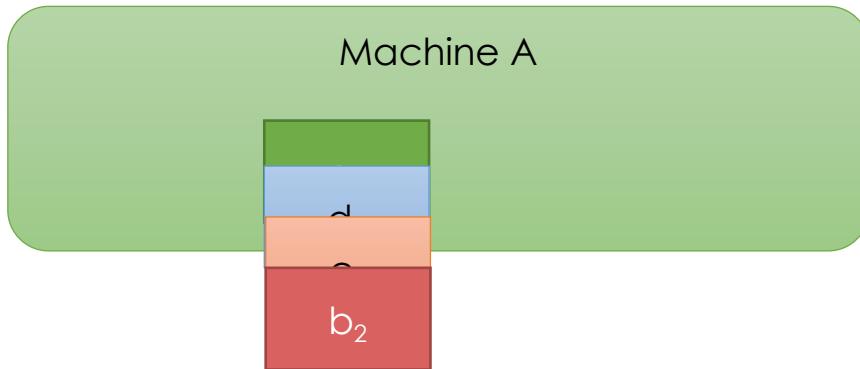
Ring All Reduce





Ring All Reduce





Ring All Reduce



Machine A

s_2

Machine B

s_3

Ring All Reduce

Each machine sends N/P data to next machine each of $(p-1)$ rounds:

$(P-1) * N/P$ (doesn't depend on P!)

➤ Fan-in Per Round:

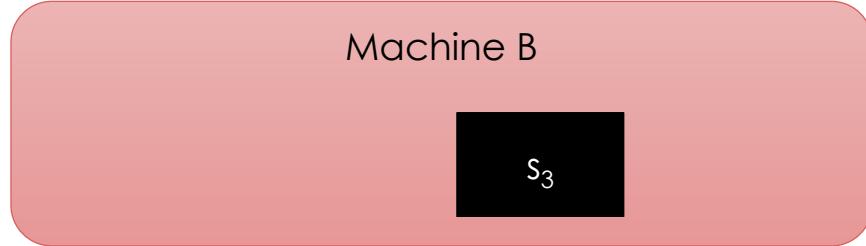
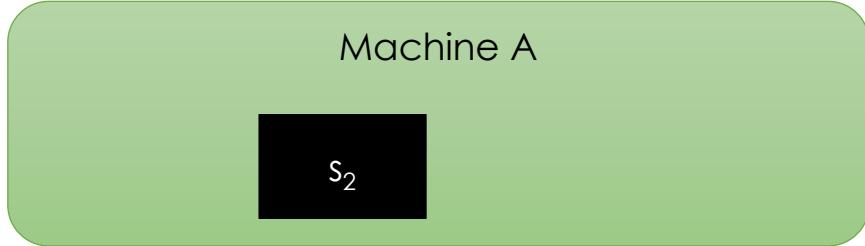
➤ 1 (doesn't depend on P)

Machine D

s_1

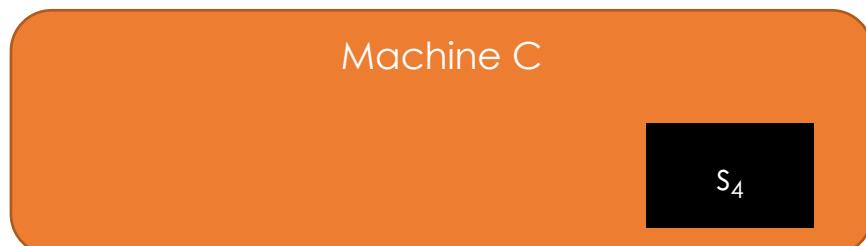
Machine C

s_4



Ring All Reduce

Broadcast stage* repeats process sending messages forwarding sums (same communication costs).



* Technically All Gather based on MPI communication definition

Machine A

s_1 s_2

Machine B

s_2 s_3

Ring All Reduce

Machine D

s_1 s_4

Machine C

s_3 s_4

Machine A

s_1 s_2 s_4

Machine B

s_1 s_2 s_3

Ring All Reduce

Machine D

s_1 s_3 s_4

Machine C

s_2 s_3 s_4

Machine A

S₁ S₂ S₃ S₄

Machine B

S₁ S₂ S₃ S₄

Ring All Reduce

Machine D

S₁ S₂ S₃ S₄

Machine C

S₁ S₂ S₃ S₄

Machine A

s_1 s_2 s_3 s_4

A horizontal row of four black rectangular boxes, each containing a white letter 's' followed by a subscript number from 1 to 4. They are arranged side-by-side within a light green rounded rectangular frame.

Machine B

s_1 s_2 s_3 s_4

A horizontal row of four black rectangular boxes, each containing a white letter 's' followed by a subscript number from 1 to 4. They are arranged side-by-side within a light red rounded rectangular frame.

Ring All Reduce

Machine D

s_1 s_2 s_3 s_4

A horizontal row of four black rectangular boxes, each containing a white letter 's' followed by a subscript number from 1 to 4. They are arranged side-by-side within a light blue rounded rectangular frame.

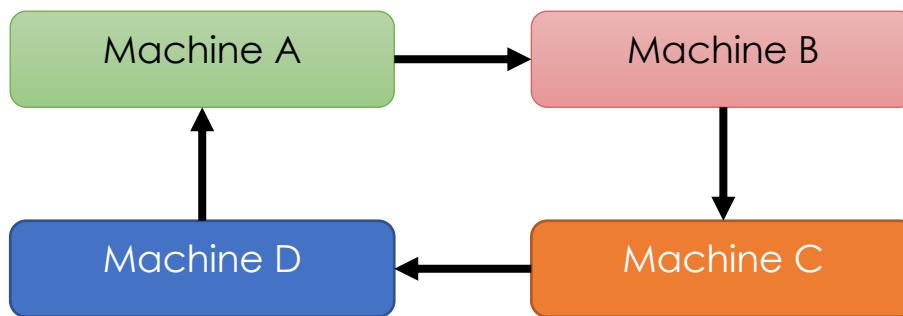
Machine C

s_1 s_2 s_3 s_4

A horizontal row of four black rectangular boxes, each containing a white letter 's' followed by a subscript number from 1 to 4. They are arranged side-by-side within a light orange rounded rectangular frame.

Ring All-Reduce

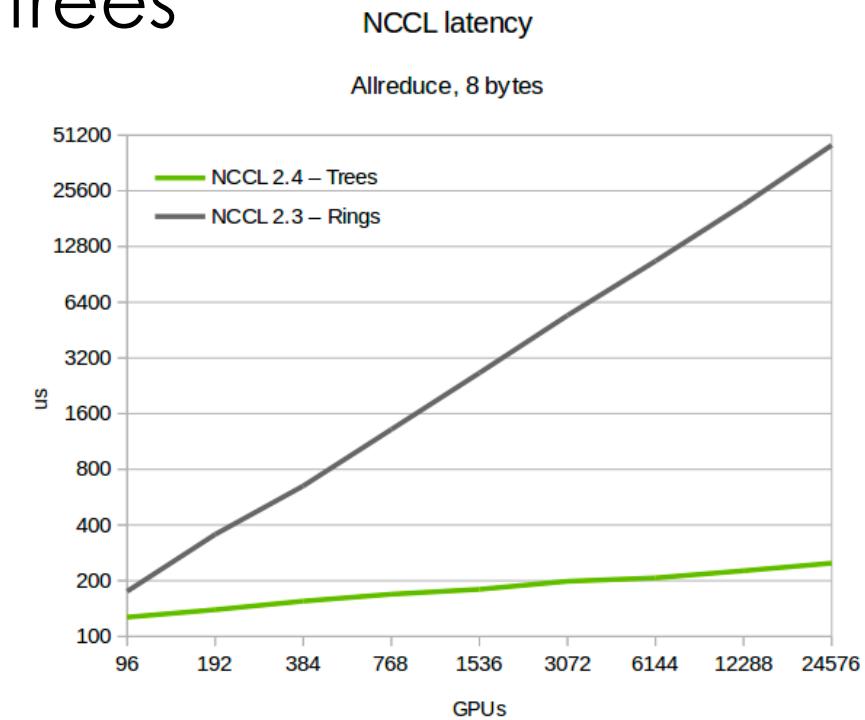
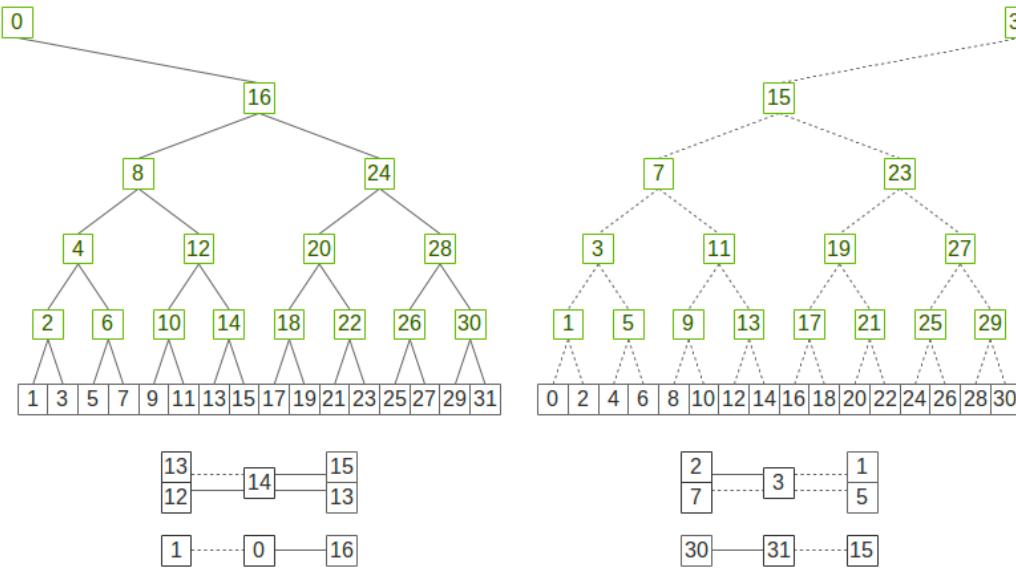
- Simplified communication topology with low fan-in



- Overall communication
 - Same total communication: $2*(P-1)*N$, but evenly distributed
 - Each Machine communicates $2*(P-1)N/P$ (almost independent of P)
 - Fan-in is constant (doesn't depend on P)
- Issue: Number of communication rounds ($P-1$)

Double Binary Tree All-Reduce

- Two overlaid binary reduction trees

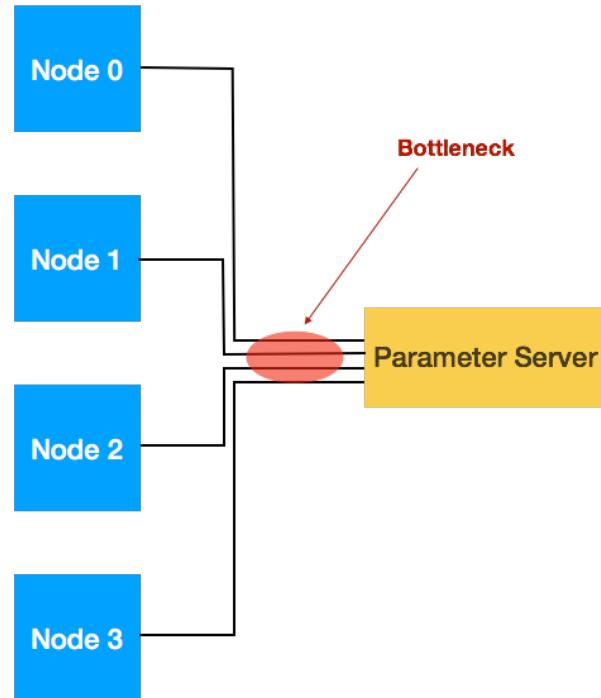


- Double the fan-in → $\log(p)$ rounds of communication
 - Currently used on Summit super-computer and latest NCCL

Complexity Summary

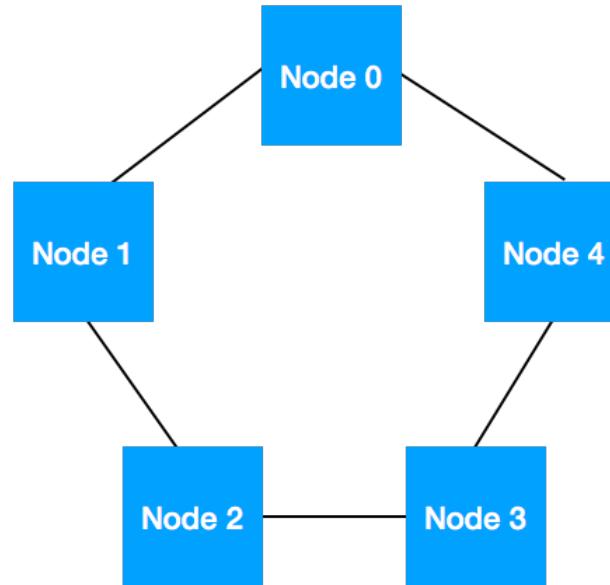
$$T_{comm} = (\alpha + PN\beta)$$

α latency
 β bandwidth
 N message size
 P #processes



Parameter Server

$$T_{comm} = 2((P - 1)\alpha + \frac{P - 1}{P}N\beta)$$



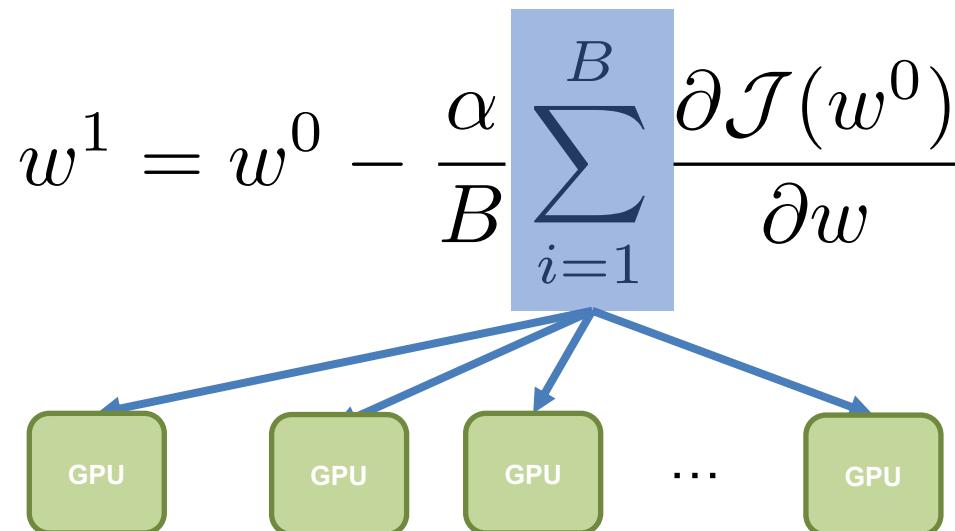
Ring All-reduce

Great Reference: T. Rajeev, R. Rabenseifner, and W. Gropp. "Optimization of collective communication operations in MPICH." *The International Journal of High Performance Computing Applications*, 2005.

Data Parallel Training Complexity Analysis

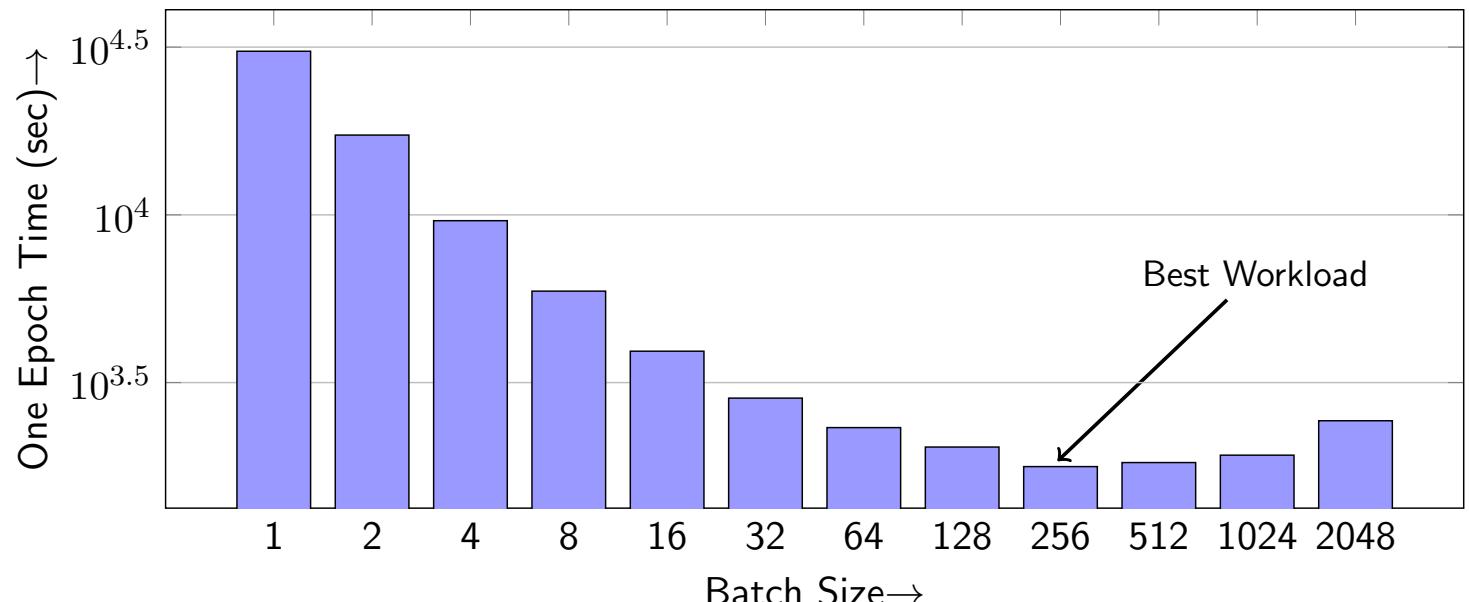
- Question: Comm time of ring allreduce is independent of the number of processors. So what limits scalability?

$$T_{comm}(\text{batch}) = 2 \sum_{i=0}^L \left(\alpha(P-1) + \beta \frac{P-1}{P} |W_i| \right)$$



Limits of Data Parallel Scaling

- The maximum limit of processors that you can use is $P=B$
- But this often leads to very low utilization of the hardware and would not yield any speed up

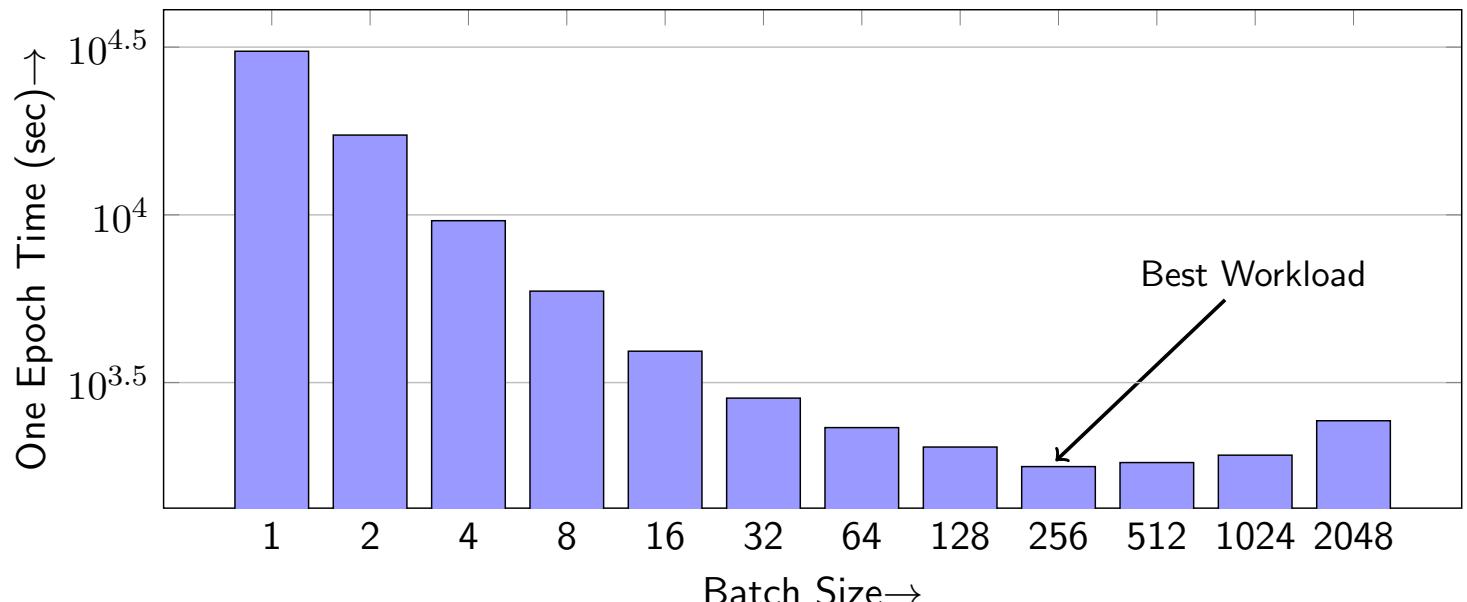


One epoch training time of AlexNet computed on an Intel KNL system

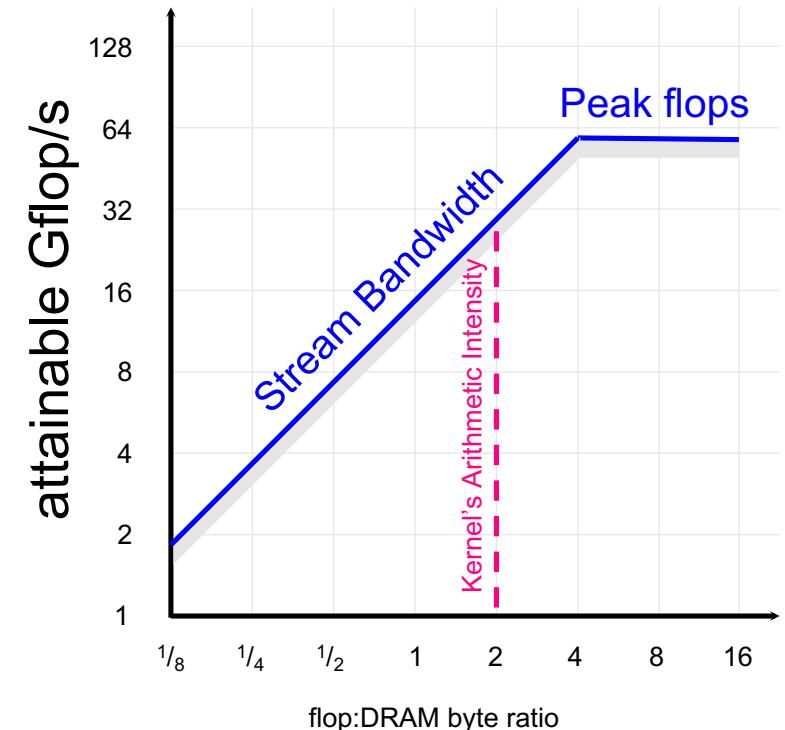
- Why does this happen?
- Remember roofline model?

Limits of Data Parallel Scaling

- The maximum limit of processors that you can use is $P=B$
- But this often leads to very low utilization of the hardware and would not yield any speed up



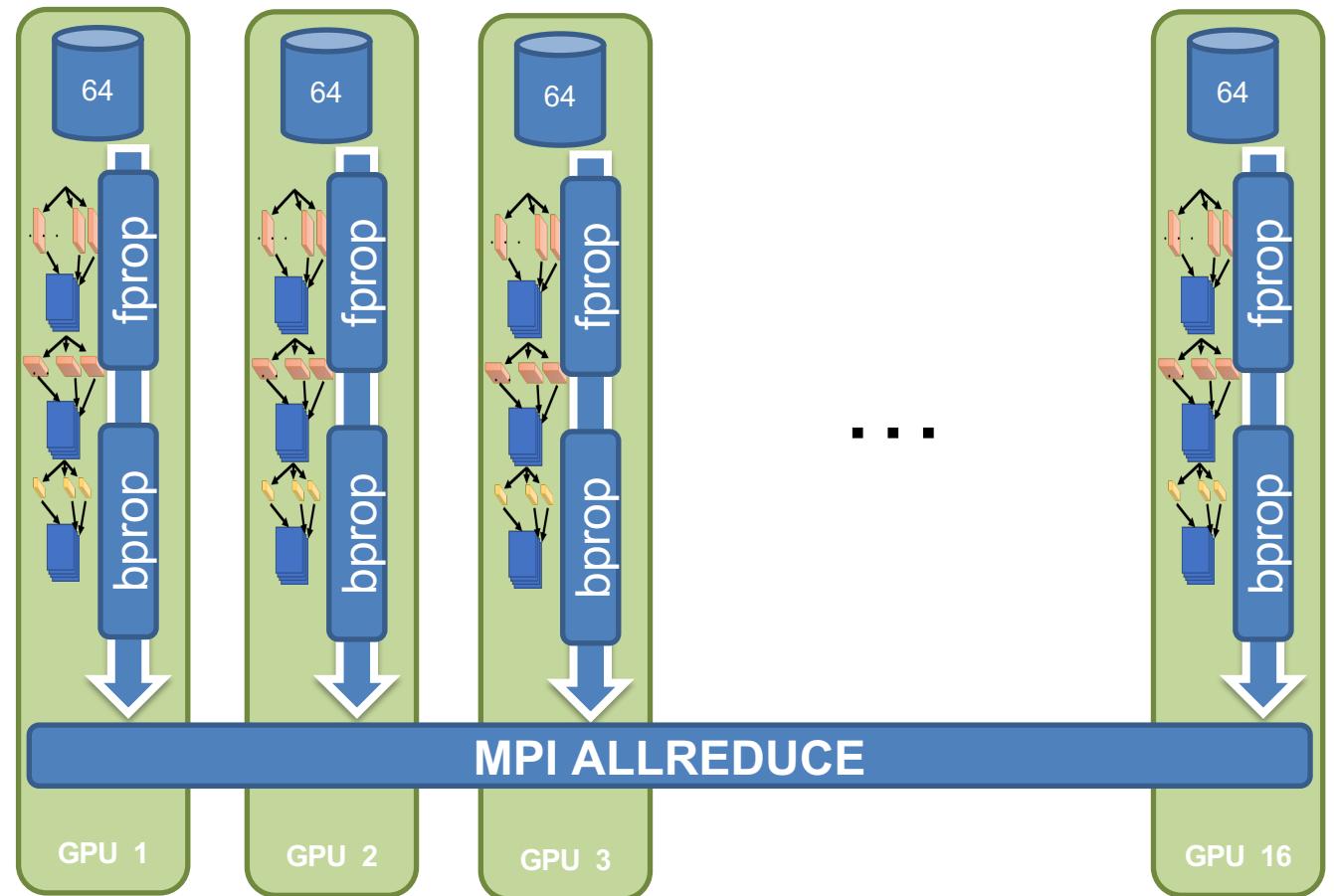
One epoch training time of AlexNet computed on an Intel KNL system



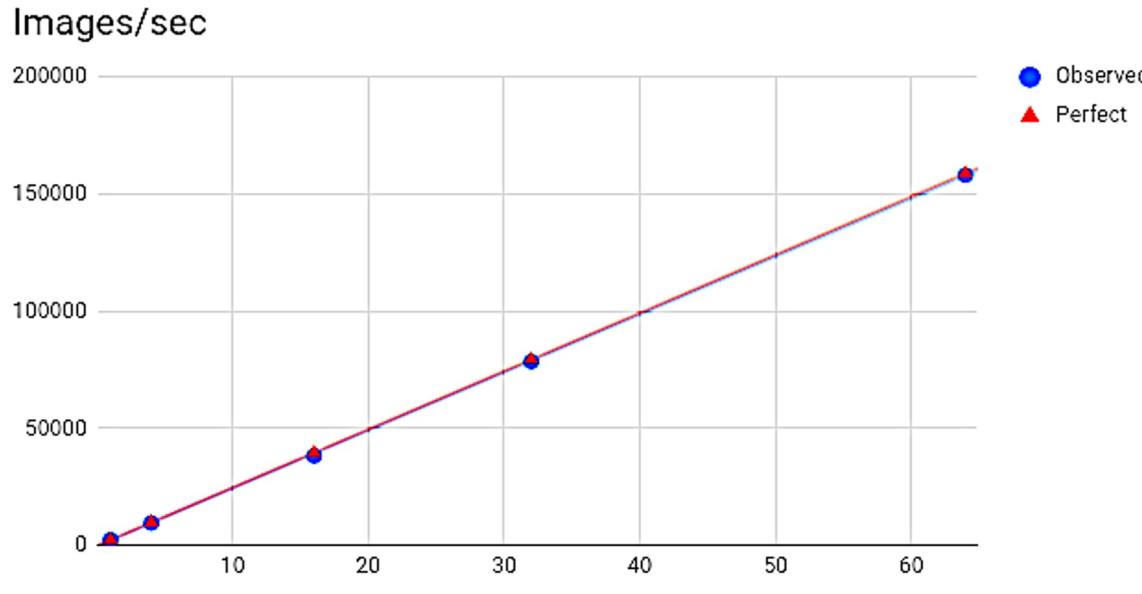
Scaling Data Parallel Training

If we want to keep scaling synchronous SGD then we have to keep **increasing** the batch size.

1024



Naively increasing Batch size leads to perfect results but ...



$$\left[\frac{\text{“Learning”}}{\text{Second}} \right] = \left[\frac{\text{“Learning”}}{\text{Record}} \right] \times \left[\frac{\text{Record}}{\text{Second}} \right]$$

Convergence Machine Learning Property

Throughput System Property

Bigger isn't Always Better

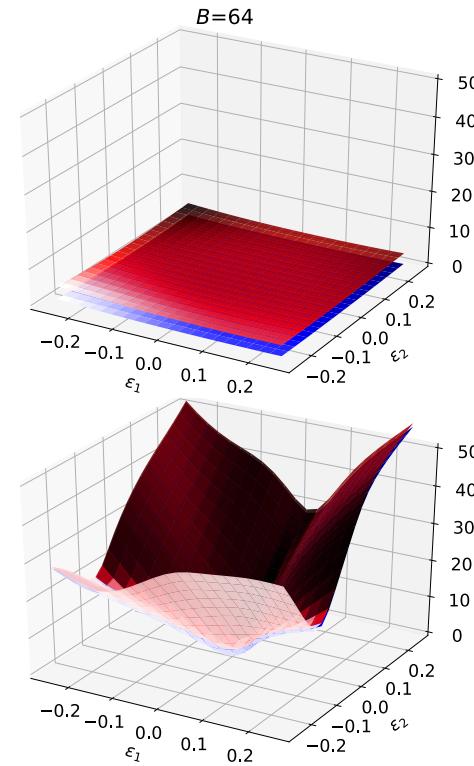
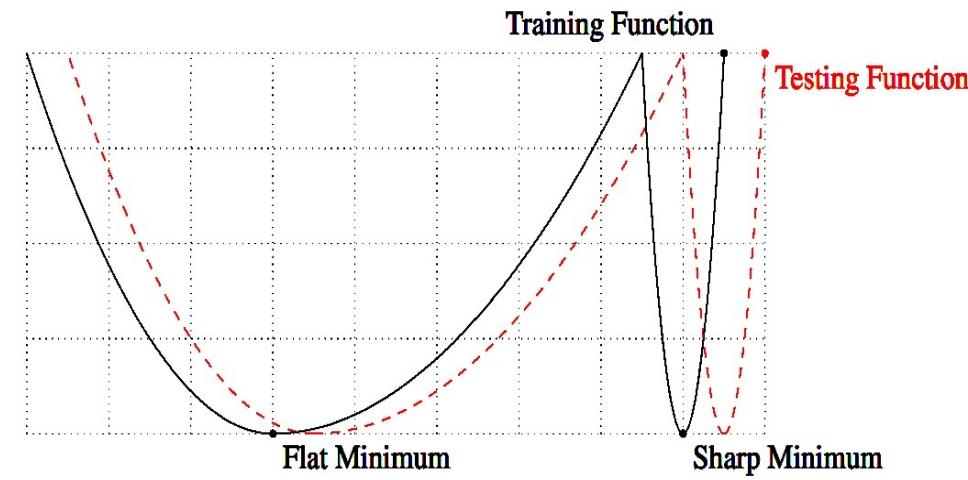
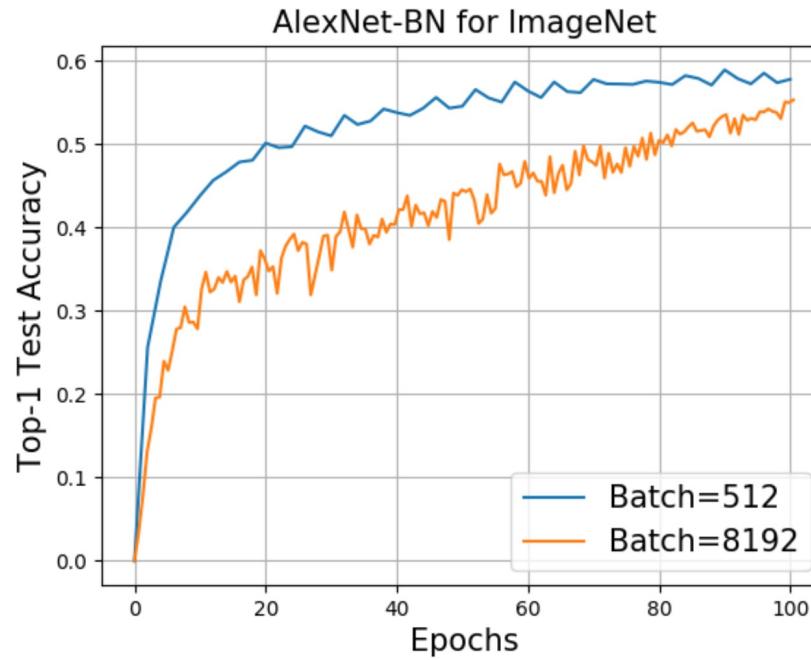
- Motivation for larger batch sizes
 - More opportunities for parallelism → but is it useful?
 - Recall ($1/n$ variance reduction):

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

- Is a variance reduction helpful?
 - Only if it lets you take bigger steps (move faster)
 - Does it affect the final prediction accuracy?

Problems with Large Batch Training

- Larger Batch leads to **sub-optimal generalization**
- A common belief is that large batch training gets attracted to “**sharp minimas**”

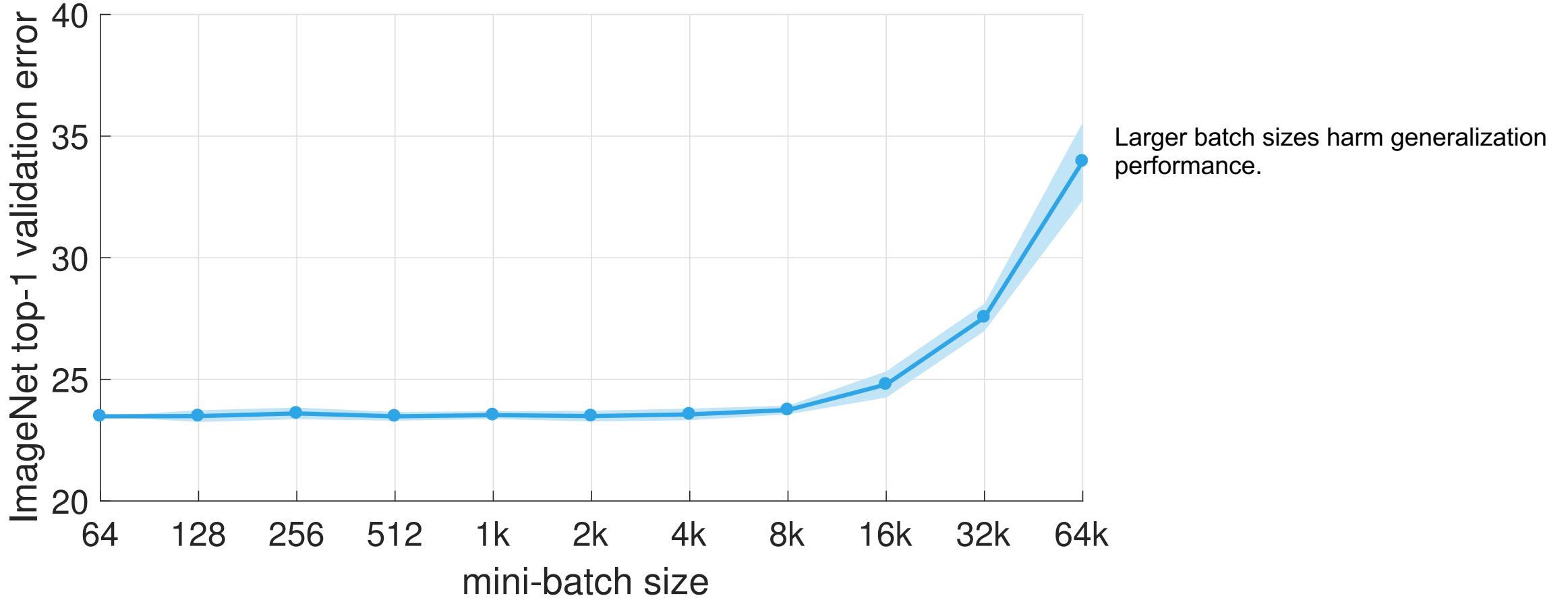


Keskar et al., On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima, ICLR’16.

Z. Yao, A. Gholami, Q. Lei, K. Keutzer, M. Mahoney. Hessian-based Analysis of Large Batch Training and Robustness to Adversaries, NeurIPS’18.

Ginsburg, Boris, Igor Gitman, and Yang You. "Large Batch Training of Convolutional Networks with LARS." arXiv:1708.03888, 2018.

Generalization Gap Problem



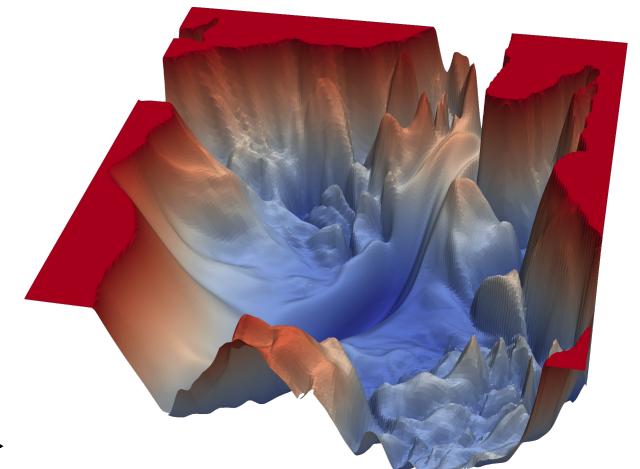
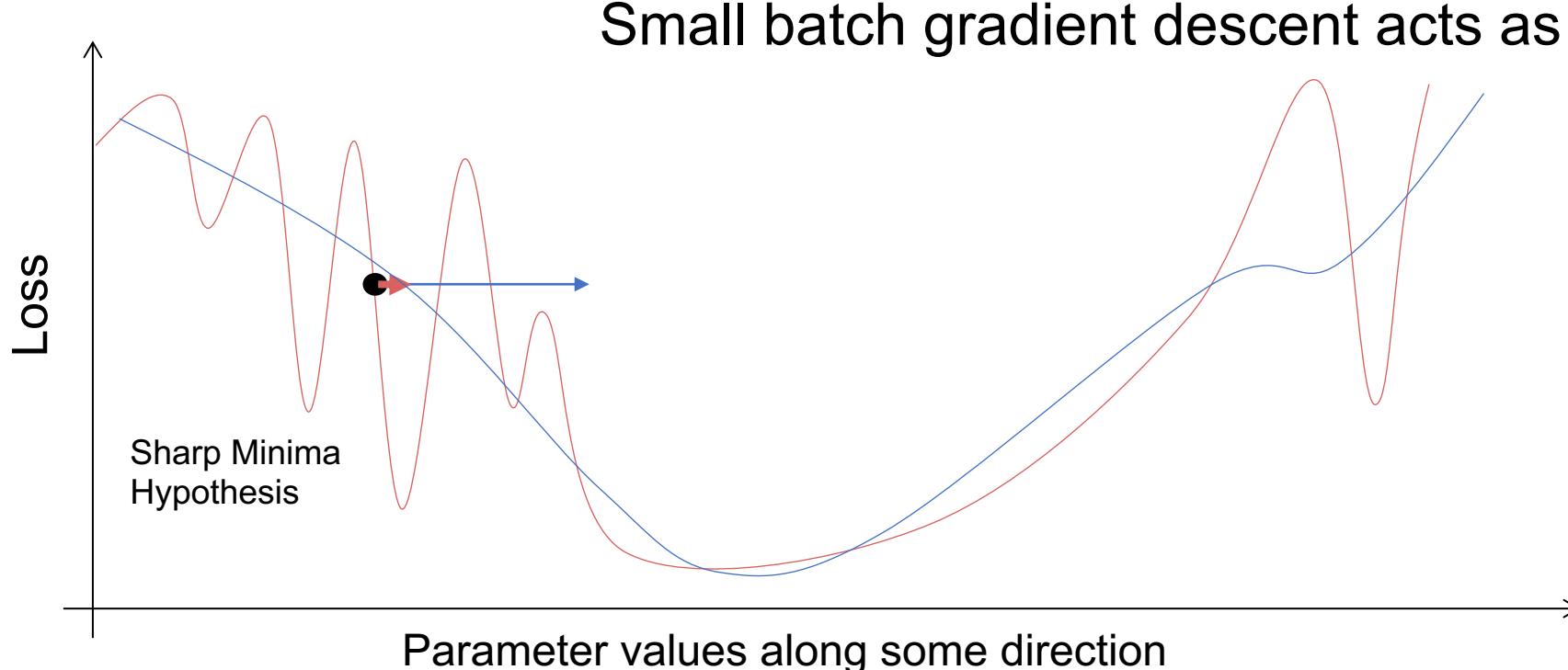
Why? Large Batch Reduces Noise and may Get Trapped in Local Minima

Objective function

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N l(x_i, y_i, \theta)$$

Update rule

$$\theta_{t+1} = \theta_t - \eta_t \frac{1}{|B|} \sum_{(x,y) \in B} \nabla_{\theta} l(x, y, \theta_t)$$



Active Research 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...

Data Parallelism Summary

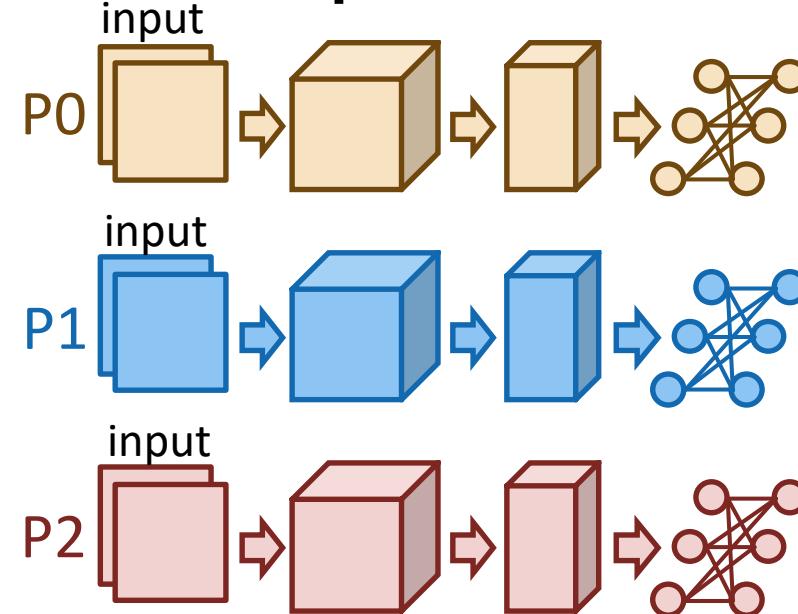
- An efficient parallel training method where the comm time is independent of processors with ring allreduce
- Very easy to implement. Only requires allreduce operation before updating parameters
- Very challenging to scale. Using large batch training is not an option as it hurts generalization performance.
 - Existing solutions often require a lot of tuning (outside of ResNet-50 on ImageNet)
- Does not work for large models such as GPT-3 which are too large to fit in one GPU
- Processes are never idle

Pipeline Parallelism

Really a form of model parallelism

Parallel and distributed training

Data parallelism



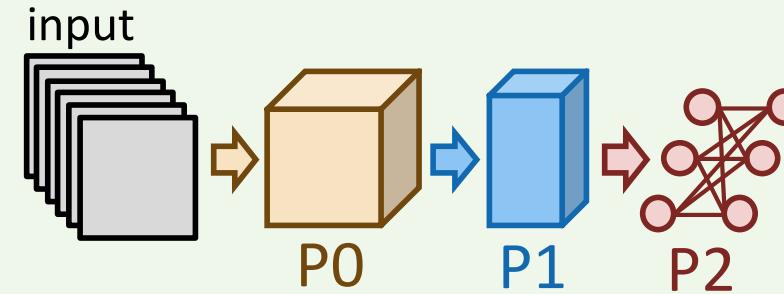
Pros:

- a. Easy to realize

Cons:

- a. Not work for large models
- b. High allreduce overhead

Pipeline parallelism



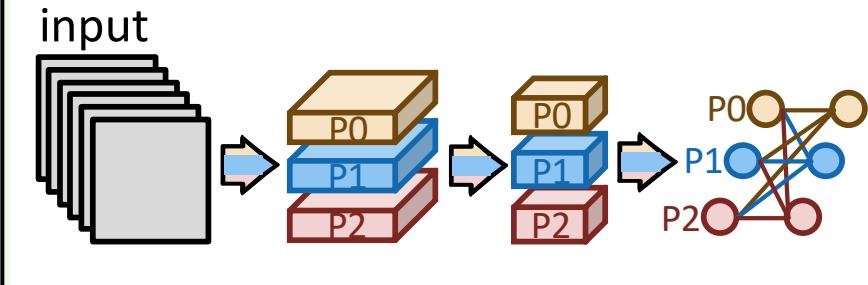
Pros:

- a. Make large model training feasible
- b. No collective, only P2P

Cons:

- a. Bubbles in pipeline
- b. Removing bubbles leads to stale weights

Model parallelism



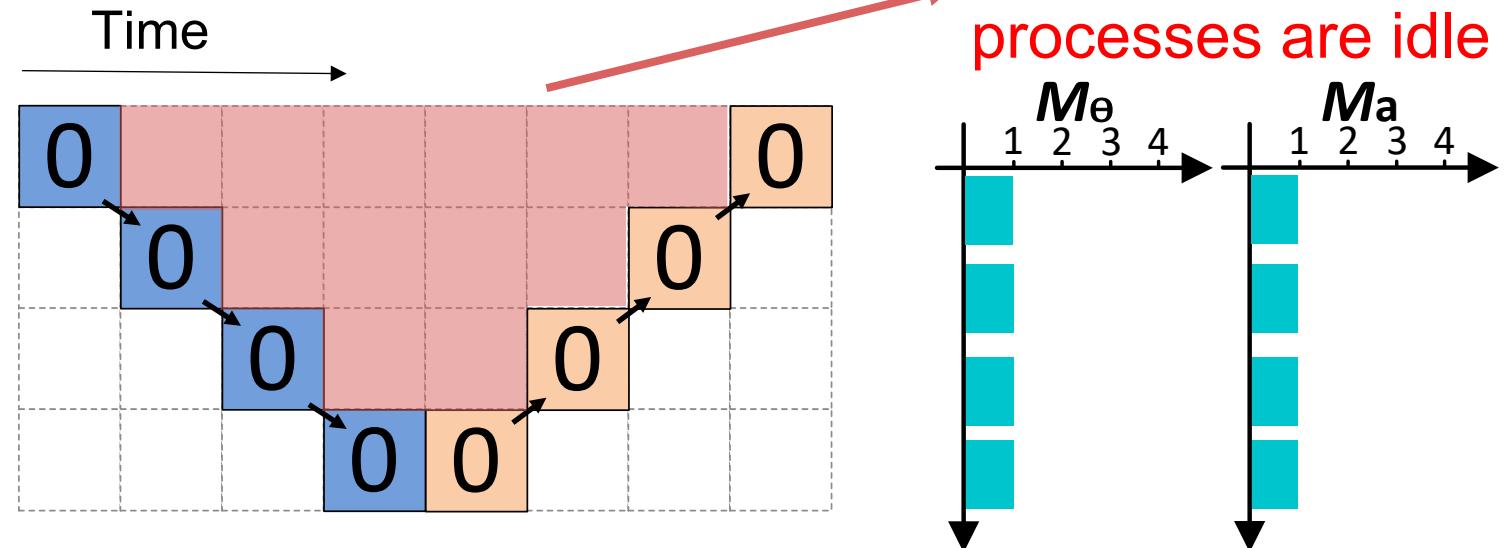
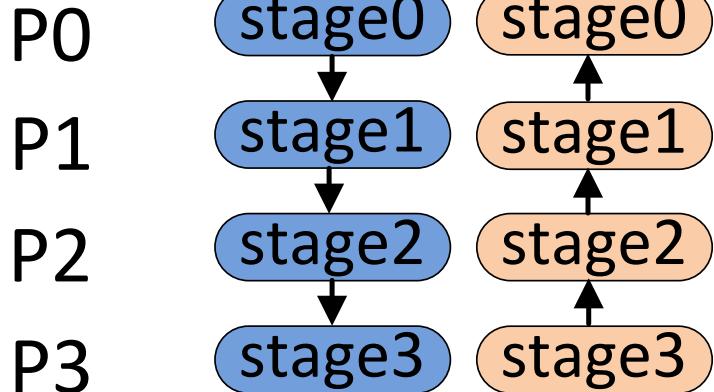
Pros:

- a. Make large model training feasible

Cons:

- b. Communication for each operator (or each layer)

Pipeline Parallelism



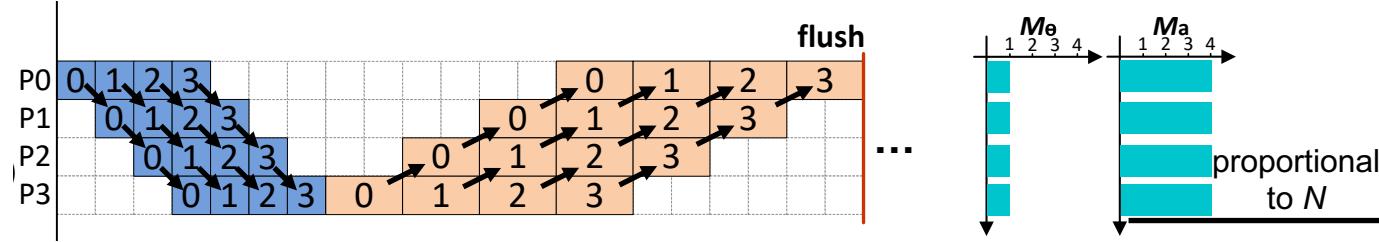
□ Bubble

■ ■ Forward and backward passes of
model replica0 for micro-batch x

M_e Memory consumption for the weights

M_a Memory consumption for the activations

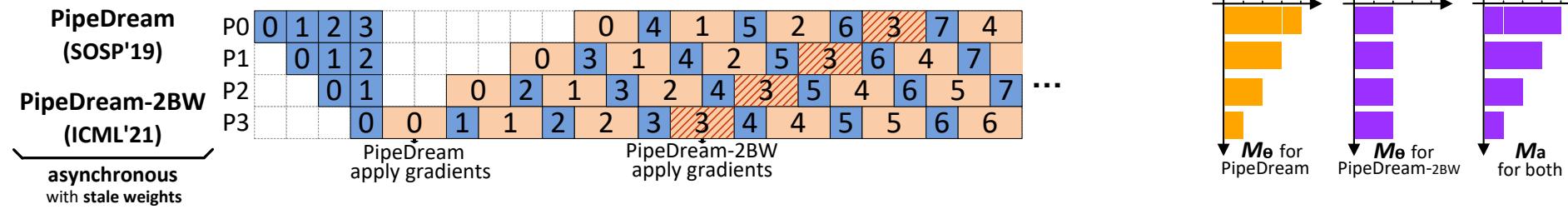
GPipe [NeurIPS'19]: Reduce Bubble with Micro-Batching



- GPipe reduces the bubble size by breaking the batch size into smaller pieces to reduce the idle time of the processes
- Pro: Reduces bubble size in an easy to implement manner
- Con: Significantly increases activation memory

 Bubble
x Forward and backward passes of
model replica 0 for micro-batch x
 M_e Memory consumption for the weights
 M_a Memory consumption for the activations

PipeDream[SOSP'19]: Use Async Updates to remove Bubble



- Pipedream uses asynchronous training: Avoid any idling by always doing a forward/backward pass irrespective of stale gradients/weights
- Pro: No bubble
- Con: As with other async methods this does affect model accuracy and convergence, and as such has not been adopted in industry.

Asynchronous Methods

- General advice: Training methods that adversely affect generalization are not adopted, unless there is a 10x speed improvement.
- Otherwise, there are so many moving parts that can go wrong in training NNs, that most often practitioners stay away from async methods unless absolutely necessary
 - For example training very large rec systems.

Pipeline Parallelism Summary

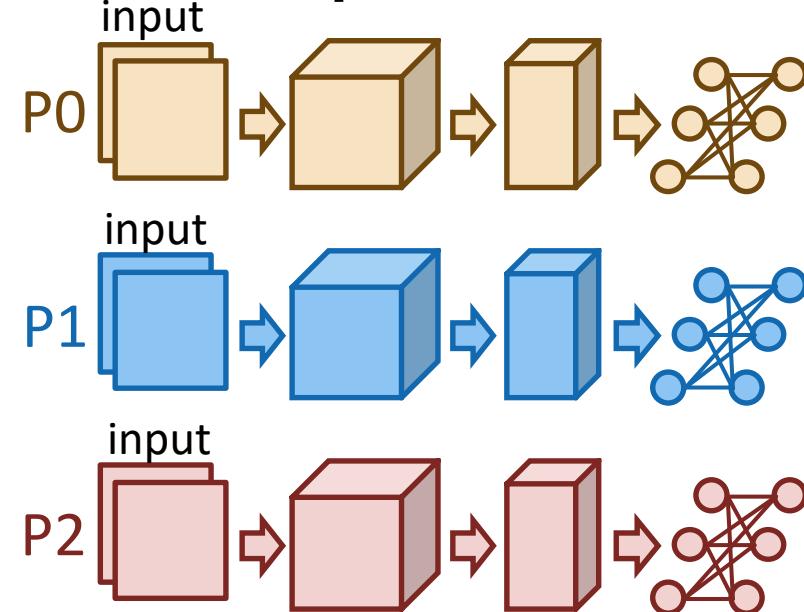
- Slightly more involved algorithm than data parallel method but with the advantage of only requiring point to point communication
- Ideal for large scale training to thousands of processes where point-to-point communication is much cheaper than collective operations such as allreduce or all-gather
- Requires special handling of bubble that results in idle processes

Model Parallelism

AKA Operator Parallelism

Parallel and distributed training

Data parallelism



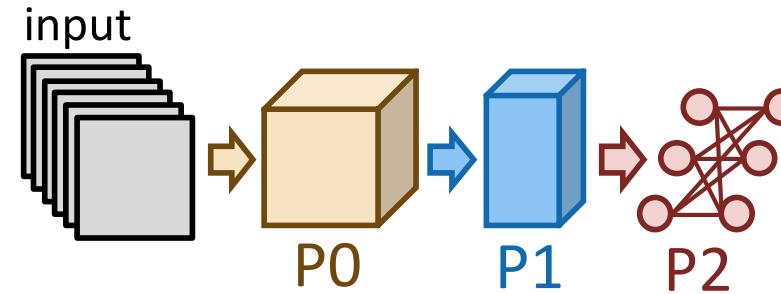
Pros:

- a. Easy to realize

Cons:

- a. Not work for large models
- b. High allreduce overhead

Pipeline parallelism



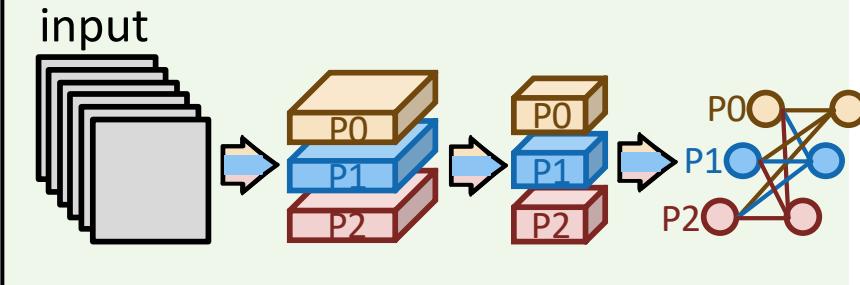
Pros:

- a. Make large model training feasible
- b. No collective, only P2P

Cons:

- a. Bubbles in pipeline
- b. Removing bubbles leads to stale weights

Model parallelism



Pros:

- a. Make large model training feasible

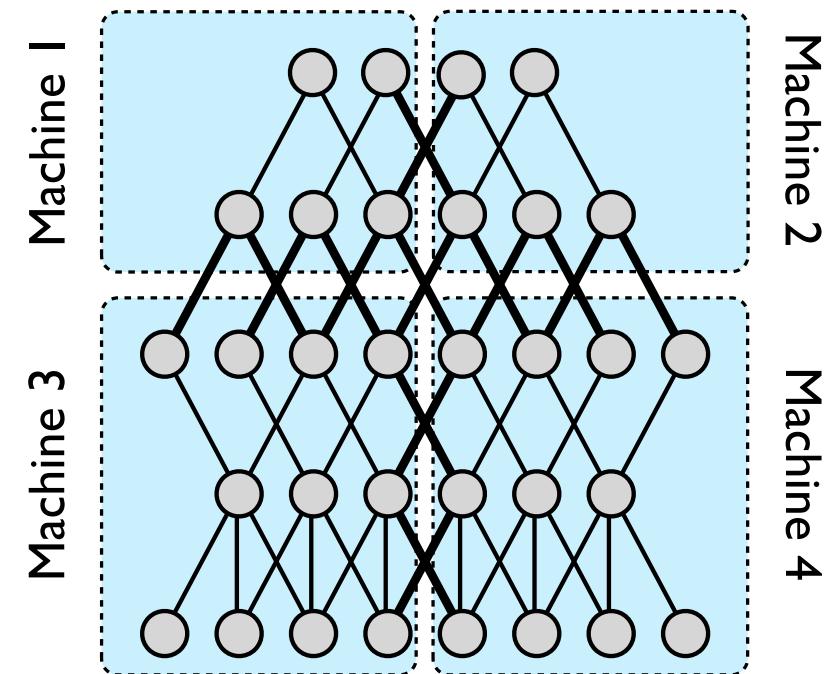
Cons:

- b. Communication for each operator (or each layer)

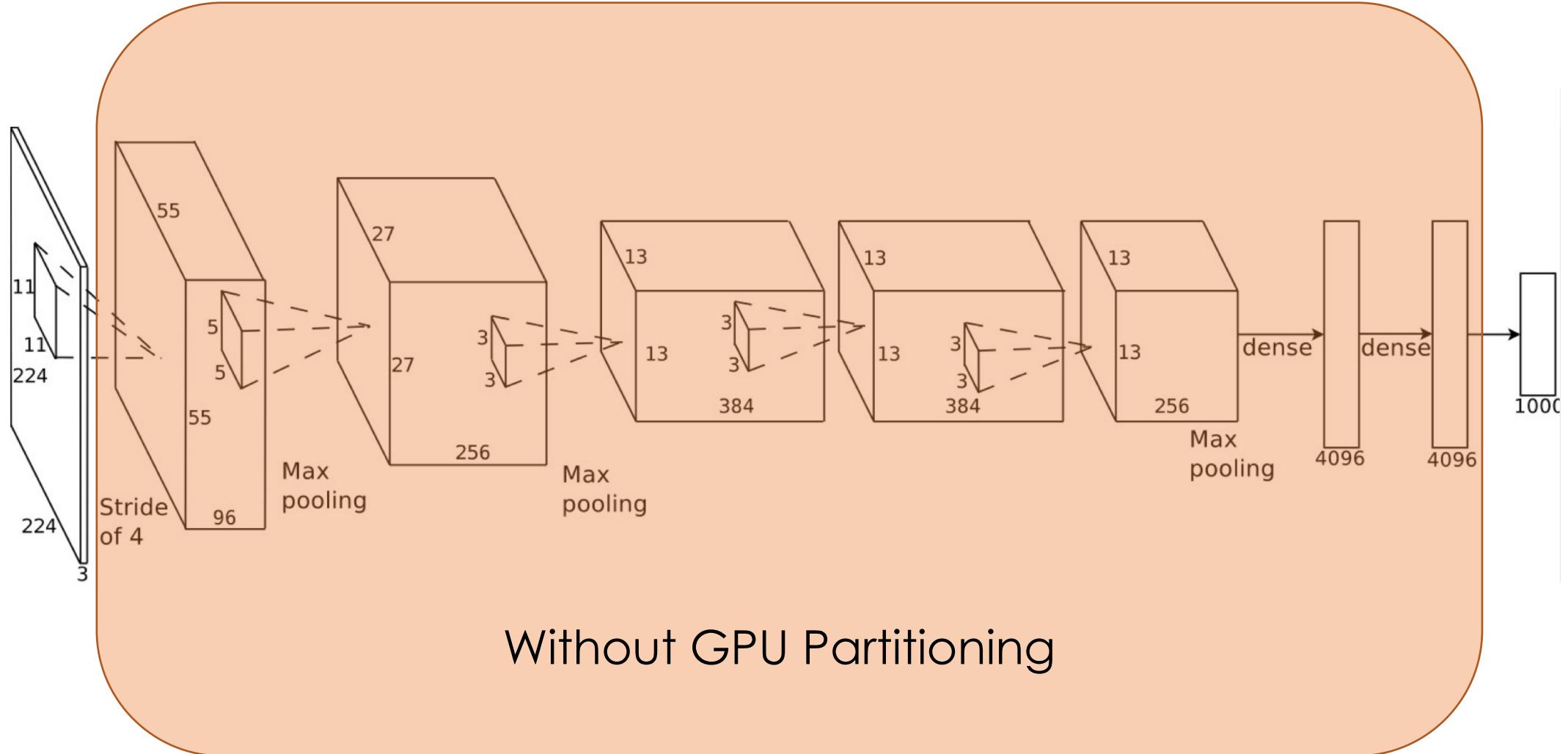
Model Parallelism

Divide the model across machines and replicate the data.

- Supports large models and activations
- Requires communication within single evaluation
- How to best divide a model?
 - Split across layers
 - Only one set of layers active a time → poor work balance
 - This is basically pipeline parallelism
 - Split individual layers
 - which dimension?
 - Weights or spatial → depends on operation

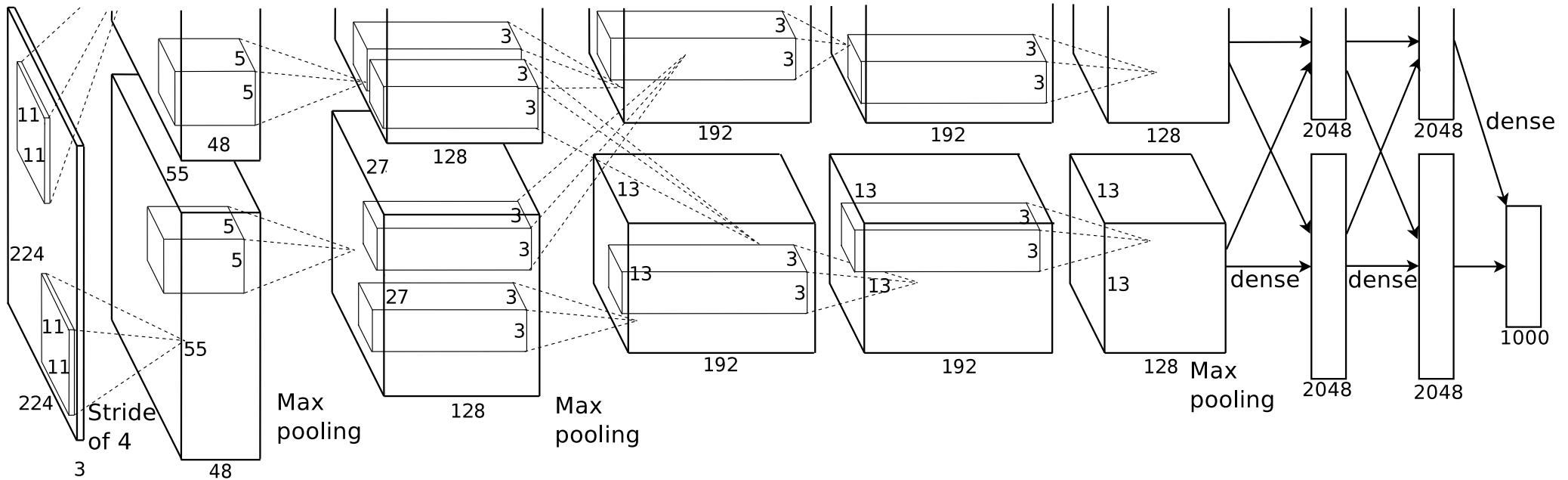


The AlexNet Architecture



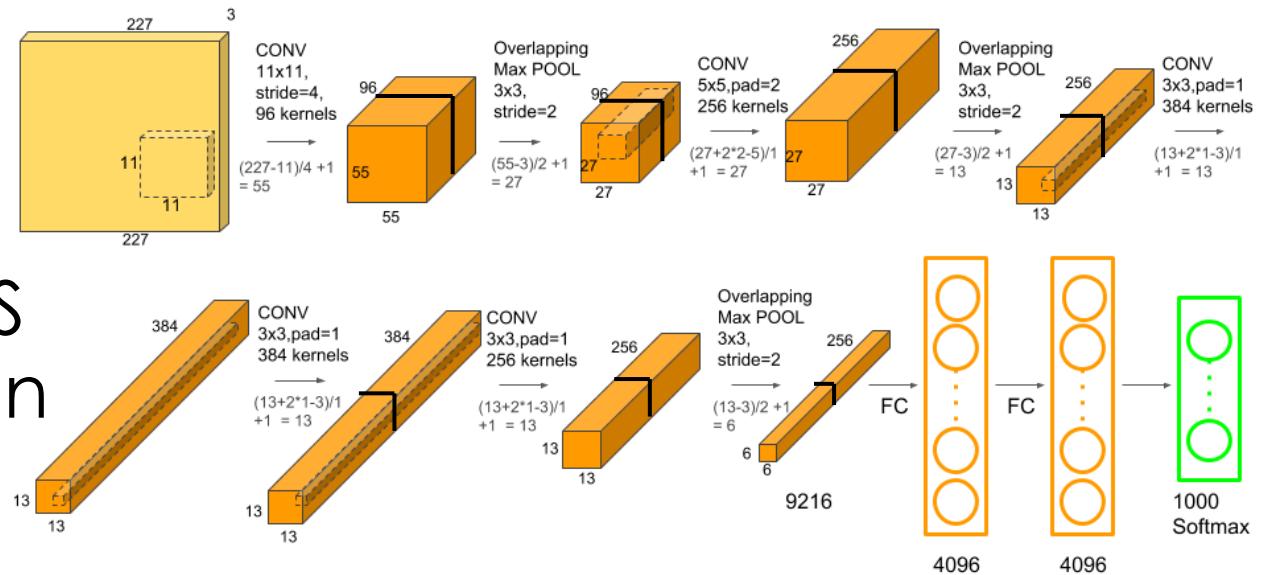
The Actual AlexNet Architecture

From the paper “ImageNet Classification with Deep Convolutional Neural Networks”



Training on Multiple GPUs

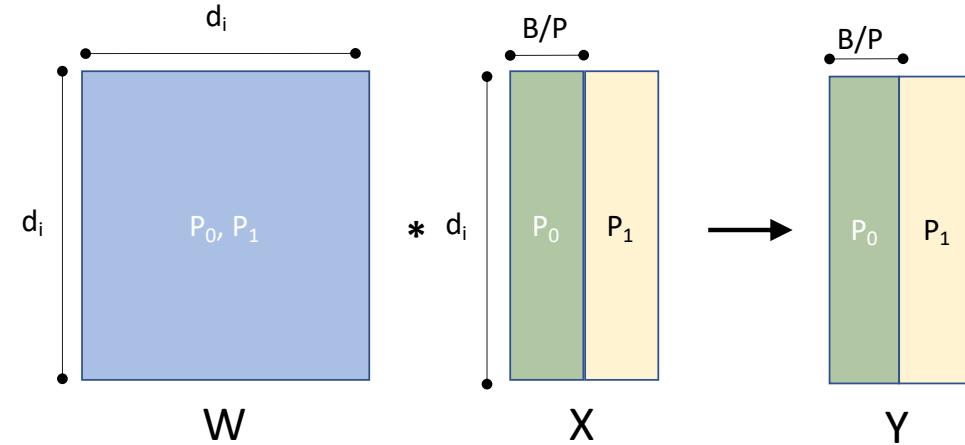
- Limited by GPU **memory** using Nvidia GTX 580 (3GB RAM)
 - 60M Parameters ~ **240 MB**
 - Need to cache activation maps for backpropagation
 - Batch size = 128
 - $128 * (227*227*3 + 55*55*96*2 + 96*27*27*2 + 256*27*27*2 + 256*13*13*2 + 13*13*384*2 + 256*13*13 + 6*6*256 + 4096 + 4096 + 1000) * 4 \text{ Bytes} \sim 782\text{MB Activations}$
 - That is assuming no overhead and single precision values
- Tuned splitting across GPUS to balance communication and computation



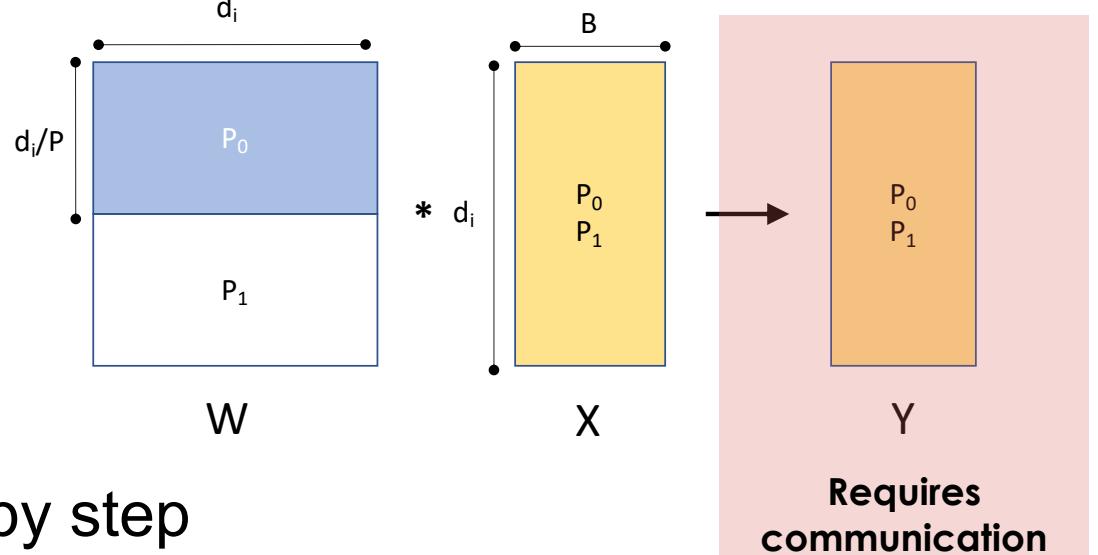
Model Parallelism: Comm Analysis

It helps to think of the operations in matrix form. Consider an FC layer

Data Parallelism: Partition input across different Processors (batch dimension)



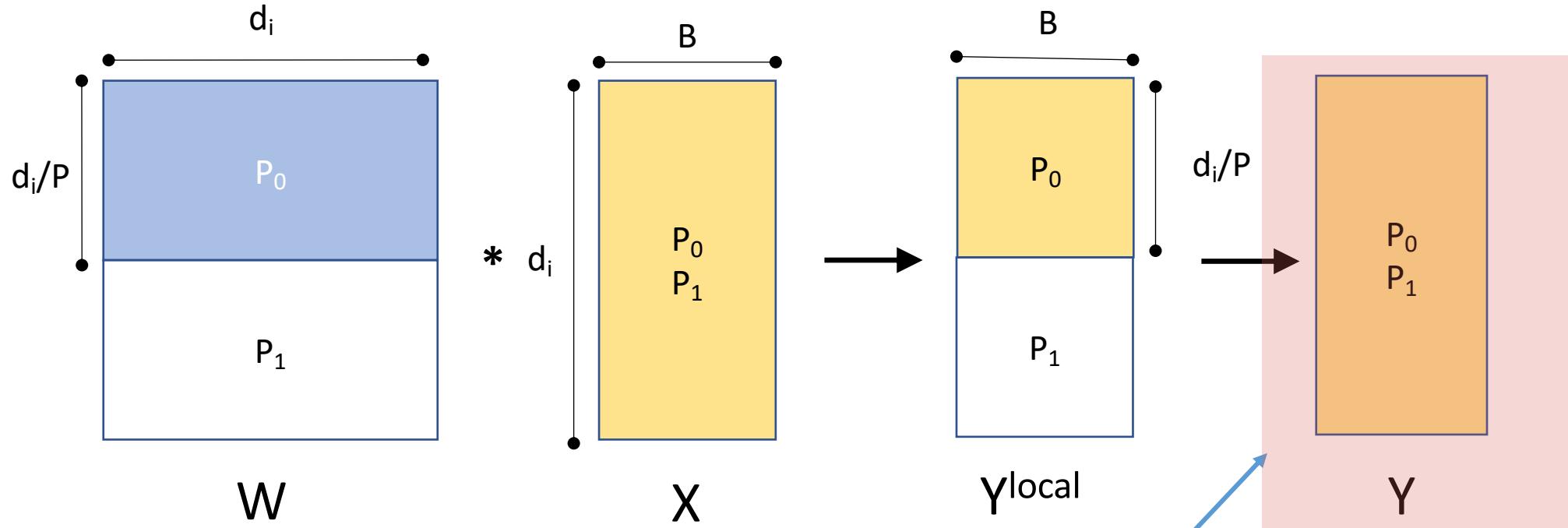
Model Parallelism: Partition weights across different Processes (W dimension)



Let's discuss the communication details, step by step

Requires communication

Comm Analysis: Forward Pass



- Requires an all gather communication so that all processes get each others activation data
- Same cost as all reduce without the 2x factor

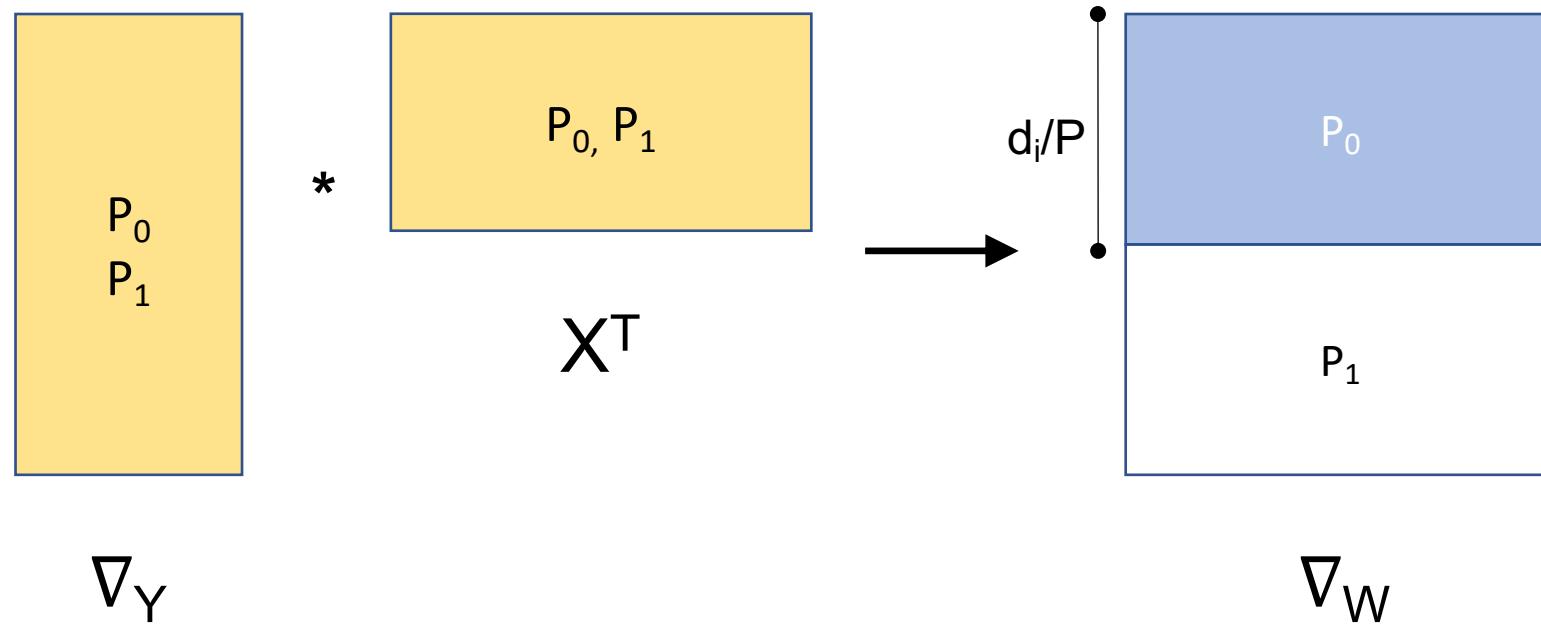
$$\sum_{i=1}^L \left(\beta(P-1) \frac{Bd_i}{P} \right)$$

* Ignoring latency term for notational simplicity

$$\nabla_Y * X^T = \nabla_W$$

$$W^T * \nabla_Y = \nabla_X$$

Backward Pass: Weights

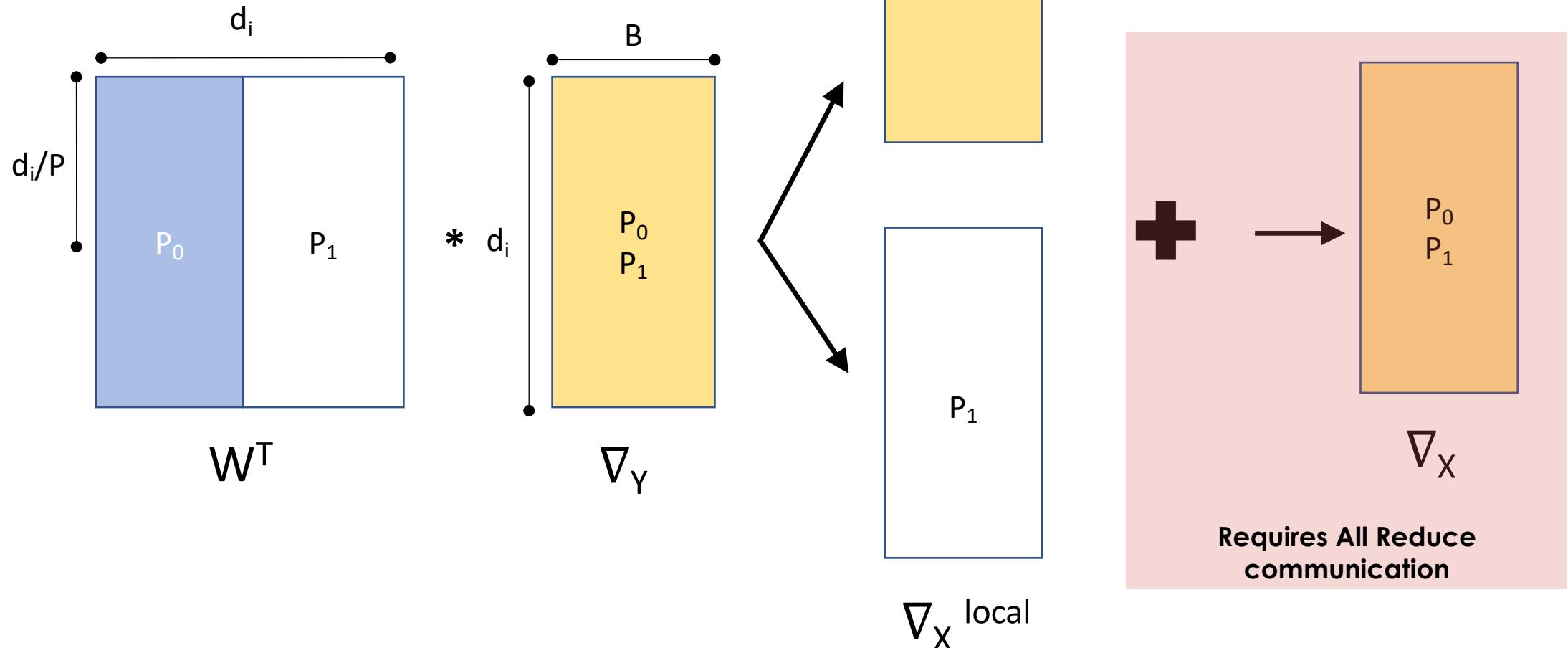


- No communication needed as every processor only needs the gradient of its own parameters
- This makes model parallelism very effective for cases where the model size is large

Backward Pass: Inputs

$$\nabla_Y * X^T = \nabla_W$$

$$W^T * \nabla_Y = \nabla_X$$



- Aggregating activation delta requires an allreduce operation

$$2 \sum_{i=2}^L \left(\beta(P-1) \frac{Bd_i}{P} \right)$$

Comm Complexity Analysis

In Model Parallelism we need two forms of communication:

1. All Gather operation so that all processors get all the activations
2. All reduce operation for backpropagating activation gradients

$$T_{comm}(\text{model}) = \sum_{i=1}^L \left(\beta(P-1) \frac{Bd_i}{P} \right) + 2 \sum_{i=2}^L \left(\beta(P-1) \frac{Bd_i}{P} \right)$$

All Gather All Reduce

Model vs Data Parallelism?

- When does it make sense to use Model vs Data Parallelism?

$$T_{comm}(\text{model}) = \sum_{i=1}^L \left(\beta(P-1) \frac{Bd_i}{P} \right) + 2 \sum_{i=2}^L \left(\beta(P-1) \frac{Bd_i}{P} \right)$$

$$T_{comm}(\text{data}) = \sum_{i=1}^L \left(\beta(P-1) \frac{d_i^2}{P} \right)$$

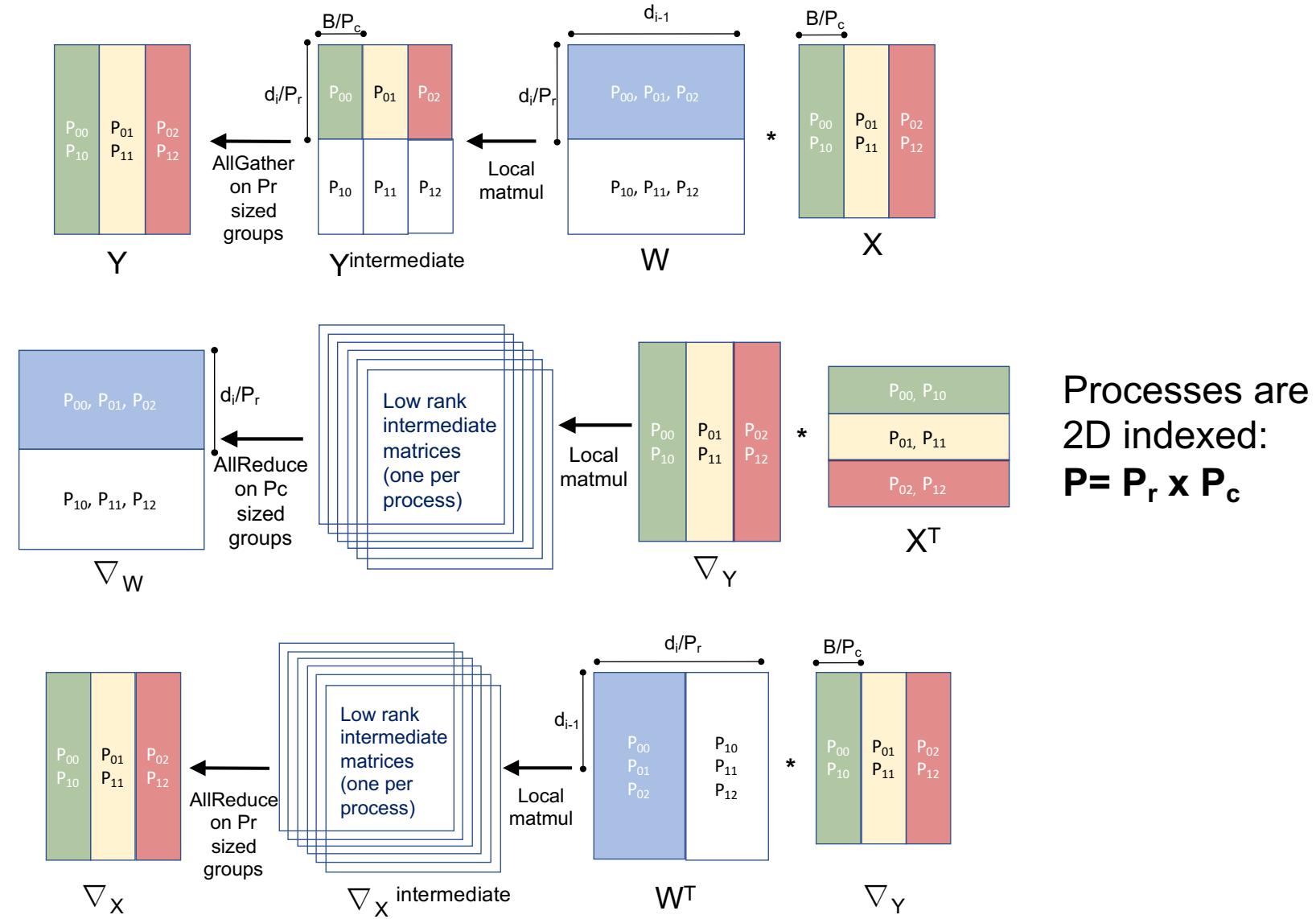
- Model parallelism reduces the quadratic complexity of d_i .
- It is useful for layers with very large weights $d_i \gg 1$
- It makes sense to use an integrated/hybrid data and model parallelism

Model Parallelism Summary

- Has **better comm complexity for large FC layers** than Data parallel approach
- Makes training large models feasible by breaking it into smaller parts
- However, requires **blocking collective communication** during **both** forward pass (all gather), as well as backwards pass (all reduce)
- Slightly **harder to implement** than data/pipeline parallel

Integrated Model and Data Parallelism

For a linear graph we can find the optimal hybrid method for analyzing the communication complexity, coupled with hardware utilization [1]



General Hybrid Methods

For a general computational graph we need to decide on:

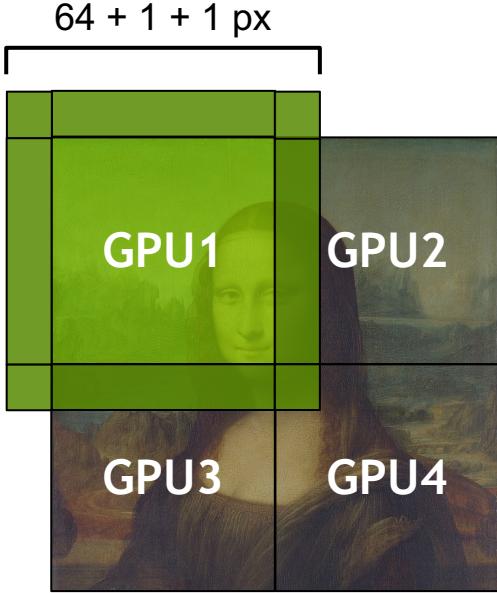
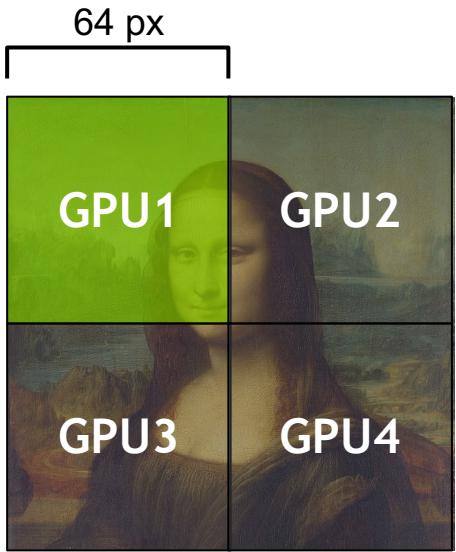
- How many processes to assign for DP
- Which axes to break the model: operator vs pipeline
- How to efficiently map the GPUs to the resulting execution graph
- ...

For a general non-linear graph this leads to a combinatorically large search space

Spatial Parallelism

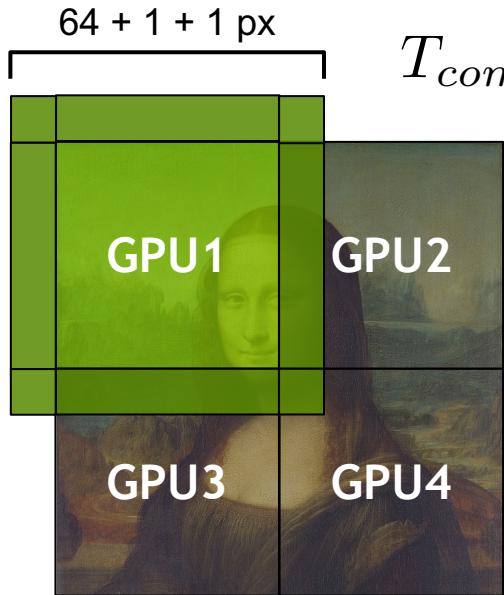
Spatial Parallel Training

- The general idea is to break the input into smaller pieces and distribute the work among different processors
 - Need to exchange boundary points for spatial convolutions



$$\begin{aligned} T_{comm}(\text{domain}) &= \sum_{i=1}^L (\alpha + \beta B X_W^i X_C^i k_h^i / 2) \\ &+ \sum_{i=1}^L (\alpha + \beta B Y_W^i Y_C^i k_w^i / 2) \\ &+ 2 \sum_{i=1}^L \left(\alpha \log(P) + \beta \frac{P-1}{P} |W_i| \right) \end{aligned}$$

Communication Complexity



$$T_{comm}(\text{domain}) = \sum_{i=0}^L (\alpha + \beta BX_W^i X_C^i k_h^i / 2) + \sum_{i=0}^L (\alpha + \beta BY_W^i Y_C^i k_w^i / 2) + 2 \sum_{i=0}^L \left(\alpha \log(P) + \beta \frac{P-1}{P} |W_i| \right)$$

Exchanging horizontal pixels

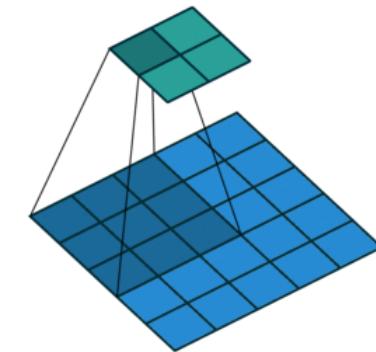
Exchanging vertical pixels

All reduce Cost
(same as before)

Useful for High Resolution Training

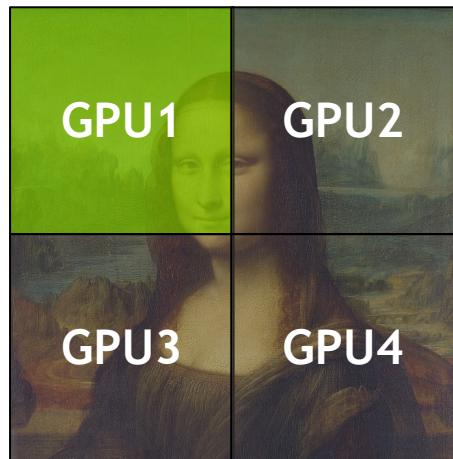
- Domain parallel scaling on V100 GPUs
 - 3x3 Conv, Batch=32, Channel=64

Resolution	Gpus	Fwd. wall-clock	Bwd. wall-clock
128×128	1	2.56 ms (1.0×)	6.63 ms (1.0×)
	2	1.52 ms (1.7×)	3.50 ms (1.9×)
	4	1.23 ms (2.1×)	2.33 ms (2.8×)
256×256	1	10.02 ms (1.0×)	26.81 ms (1.0×)
	2	5.34 ms (1.9×)	11.79 ms (2.3×)
	4	3.11 ms (3.2×)	6.96 ms (3.9×)
512×512	1	45.15 ms (1.0×)	126.11 ms (1.0×)
	2	20.18 ms (2.2×)	60.15 ms (2.1×)
	4	10.65 ms (4.2×)	26.76 ms (4.7×)



Spatial Parallelism Summary

- A little harder to implement since you need to exchange the boundary points
- Only effective for high resolution input data
 - Limits the number of processors that can be effectively utilized



Acknowledgments

Many slides from

Prof. Kurt Keutzer, Suresh Krishna, Prof. Patterson, Michael Pellauer (Nvidia), Prof. Sophia Shao, Naveen Kumar (Google), Shigang Li, Pallas Group