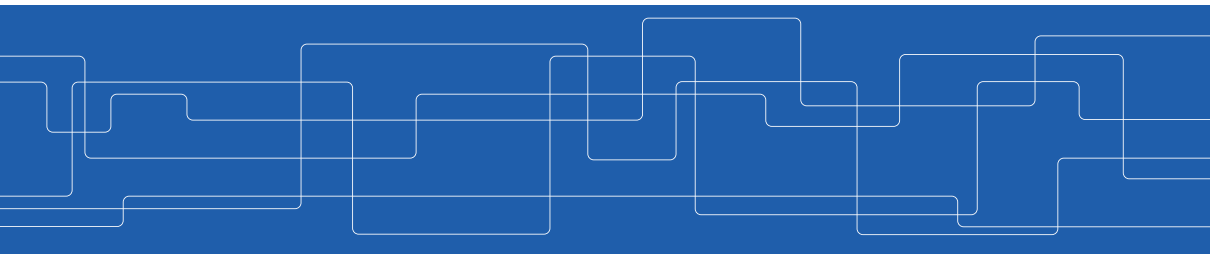# Distributed Deep Learning
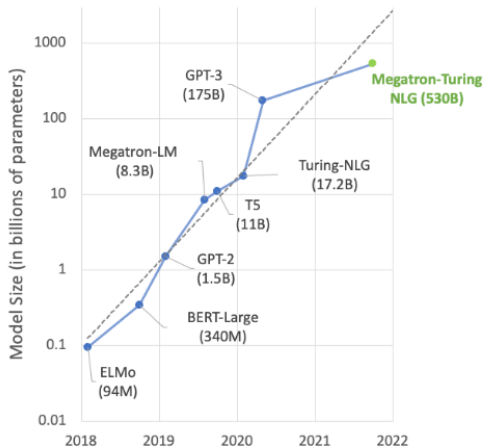
Slides by Amir H. Payberah and Jim Dowling

# The need for Distributed Training of DNNs
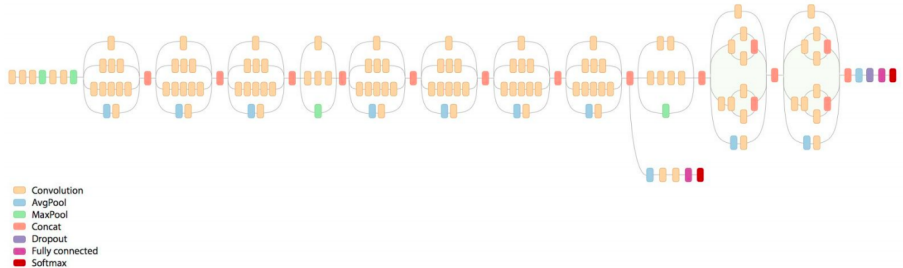
# Growth in the Size of Deep Neural Networks



Nvidia Link

- Computationally intensive
- Time consuming



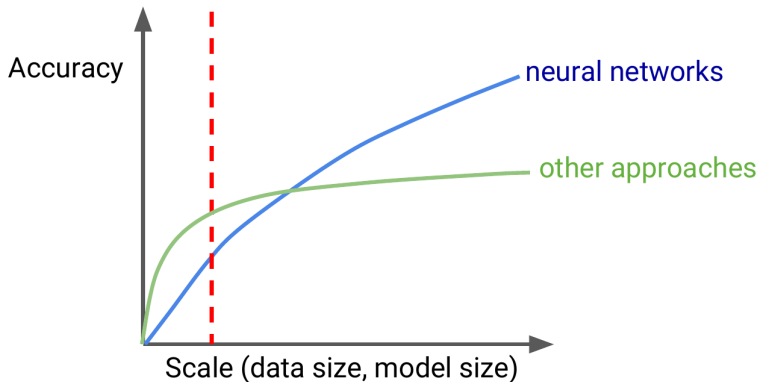[https://cloud.google.com/tpu/docs/images/inceptionv3onc--oview.png]

- **Massive** amount of training dataset
- **Large** number of parameters

**1980s and 1990s**



[Jeff Dean at AI Frontiers:  Trends and Developments in Deep Learning Research]

**1980s and 1990s**

Accuracy

more compute

neural networks

other approaches

Scale (data size, model size)

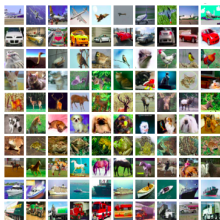[Jeff Dean at AI Frontiers: Trends and Developments in Deep Learning Research]

[Jeff Dean at AI Frontiers: Trends and Developments in Deep Learning Research]

# Fundamentals of Machine Learning

▶ E.g., tabular data, image, text, etc.

- E.g., linear models, neural networks, etc.
- $\hat{y} = f_w(\mathbf{x})$

# Loss function

- How good $\hat{y}$ is able to predict the expected outcome y.

- $J(\mathbf{w}) = \sum_{i=1}^{m} l(y_i, \hat{y}_i)$



- E.g., $J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^{m} (y_i - \hat{y}_i)^2$

# Objective

- Minimize the loss function

- $\arg\min_{\mathbf{w}} J(\mathbf{w})$

- $J(\mathbf{w}) = \sum_{i=1}^{m} l(y_i, \hat{y}_i)$

# Training

- $J(\mathbf{w}) = \sum_{i=1}^{m} l(y_i, \hat{y}_i)$

- Gradient descent, i.e., $\mathbf{w} := \mathbf{w} - \eta \nabla J(\mathbf{w})$



- Stochastic gradient descent, i.e., $\mathbf{w} := \mathbf{w} - \eta \tilde{g} J(\mathbf{w})$
  - $\tilde{g}$: gradient at a randomly chosen point.

- Mini-barch gradient descent, i.e., $\mathbf{w} := \mathbf{w} - \eta \tilde{g}_B J(\mathbf{w})$
  - $\tilde{g}$: gradient with respect to a set of $B$ randomly chosen points.

# Let's Scale the Learning

- Data parallelism

- Model parallelism

# Data Parallelism

# Data Parallelization (1/4)

▸ Replicate a whole model on every device.
▸ Train all replicas simultaneously, using a different mini-batch for each.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

# Data Parallelization (2/4)

- $k$ devices
- $J_j(\mathbf{w}) = \sum_{i=1}^{b_j} l(y_i, \hat{y}_i),\ \forall j = 1, 2, \cdots, k$
- $\tilde{g}_B J_j(\mathbf{w})$: gradient of $J_j(\mathbf{w})$ with respect to a set of $B$ randomly chosen points at device $j$.
- Compute $\tilde{g}_B J_j(\mathbf{w})$ on each device $j$.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

▶ Compute the mean of the gradients.

▶ $\tilde{g}_B J(\mathbf{w}) = \frac{1}{k} \sum_{j=1}^{k} \tilde{g}_B J_j(\mathbf{w})$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Update the model.
- $\mathbf{w} := \mathbf{w} - \eta\tilde{\mathbf{g}}_B J(\mathbf{w})$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

# Data Parallelization Design Issues

- The aggregation algorithm

- Communication synchronization and frequency

- Communication compression

# The Aggregation Algorithm

- How to aggregate gradients (compute the mean of the gradients)?

- Centralized - parameter server

- Decentralized - all-reduce

- Decentralized - gossip

# Aggregation - Centralized - Parameter Server

- Store the model parameters outside of the workers.

- Workers periodically report their computed parameters or parameter updates to a (set of) parameter server(s) (PSs).



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Mirror all the model parameters across all workers (no PS).
- Workers exchange parameter updates directly via an allreduce operation.



[Tang et al., Communication-Efficient Distributed Deep Learning:  A Comprehensive Survey, 2020]

- No PS, and no global model.
- Every worker communicates updates with their neighbors.
- The consistency of parameters across all workers only at the end of the algorithm.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Reduce: reducing a set of numbers into a smaller set of numbers via a function.
- E.g., sum([1, 2, 3, 4, 5]) = 15
- Reduce takes an array of input elements on each process and returns an array of output elements to the root process.



[https://mpitutorial.com/tutorials/mpi-reduce-and-allreduce]

▶ AllReduce stores reduced results across all processes rather than the root process.



[https://mpitutorial.com/tutorials/mpi-reduce-and-allreduce]

Initial state

After AllReduce operation



| Worker A | | | |
|---|---|---|---|
| 17 | 11 | 1 | 9 |

| Worker B | | | |
|---|---|---|---|
| 5 | 13 | 23 | 14 |

| Worker C | | | |
|---|---|---|---|
| 3 | 6 | 10 | 8 |

| Worker D | | | |
|---|---|---|---|
| 12 | 7 | 2 | 12 |

| Worker A | | | |
|---|---|---|---|
| 37 | 37 | 36 | 43 |

| Worker B | | | |
|---|---|---|---|
| 37 | 37 | 36 | 43 |

| Worker C | | | |
|---|---|---|---|
| 37 | 37 | 36 | 43 |

| Worker D | | | |
|---|---|---|---|
| 37 | 37 | 36 | 43 |

`[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]`

- All-to-all allreduce
- Master-worker allreduce
- Tree allreduce
- Round-robin allreduce
- Butterfly allreduce
- Ring allreduce

# AllReduce Implementation - All-to-All AllReduce

- **Send** the array of data to each other.

- Apply the reduction operation on each process.

- Too many unnecessary messages.



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- Selecting one process as a master, gather all arrays into the master.

- Perform reduction operations locally in the master.

- Distribute the result to the other processes.

- The master becomes a bottleneck (not scalable).



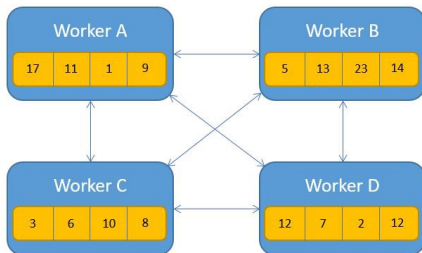[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- Some try to minimize bandwidth.
- Some try to minimize latency.



(a) Tree AllReduce     (b) Round-robin AllReduce     (c) Butterfly AllReduce

[Zhao H. et al., arXiv:1312.3020, 2013]

- The Ring-Allreduce has two phases:
  1. First, the share-reduce phase
  2. Then, the share-only phase

▶ In the share-reduce phase, each process `p` sends data to the process `(p+1)%m`
  • `m` is the number of processes, and `%` is the modulo operator.

▶ The array of data on each process is divided to `m` chunks (`m=4` here).

▶ Each one of these chunks will be indexed by `i` going forward.



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

▶ In the first share-reduce step, process A sends $a_0$ to process B.

▶ Process B sends $b_1$ to process C, etc.



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- When each process receives the data from the previous process, it applies the reduce operator (e.g., sum)
  - The reduce operator should be associative and commutative.
- It then proceeds to send it to the next process in the ring.
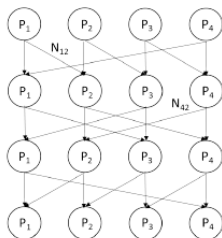


[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- The share-reduce phase finishes when each process holds the complete reduction of chunk i.

- At this point each process holds a part of the end result.



[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]

- The share-only step is the same process of sharing the data in a ring-like fashion without applying the reduce operation.

- This consolidates the result of each chunk in every process.



$r_i = a_i + b_i + c_i + d_i$

[https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da]
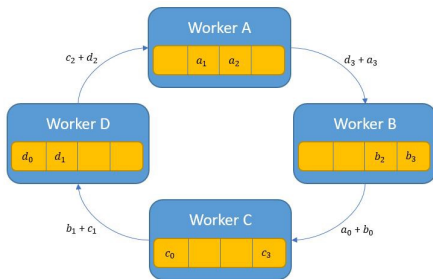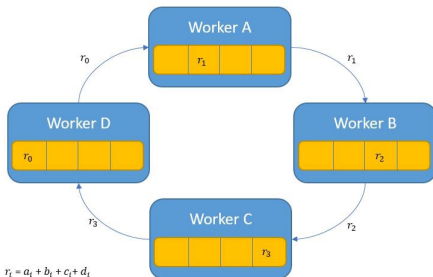
# Master-Worker AllReduce vs. Ring-AllReduce

- $N$: number of elements, $m$: number of processes

- Master-Worker AllReduce
  - First each process sends $N$ elements to the master: $N \times (m-1)$ messages.
  - Then the master sends the results back to the process: another $N \times (m-1)$ messages.
  - Total network traffic is $2(N \times (m-1))$, which is proportional to $m$.

- Ring-AllReduce
  - In the share-reduce step each process sends $\frac{N}{m}$ elements, and it does it $m-1$ times: $\frac{N}{m} \times (m-1)$ messages.
  - On the share-only step, each process sends the result for the chunk it calculated: another $\frac{N}{m} \times (m-1)$ messages.
  - Total network traffic is $2(\frac{N}{m} \times (m-1))$.

# Communication Synchronization and Frequency

- When to synchronize the parameters among the parallel workers?

# Communication Synchronization (1/2)

- Synchronizing the model replicas in data-parallel training requires communication
  - between workers, in allreduce
  - between workers and parameter servers, in the centralized architecture

- The communication synchronization decides how frequently all local models are synchronized with others.

# Communication Synchronization (2/2)

- It will influence:
  - The communication traffic
  - The performance
  - The convergence of model training

- There is a trade-off between the communication traffic and the convergence.

# Reducing Synchronization Overhead

- Two directions for improvement:

  1. To relax the synchronization among all workers.

  2. The frequency of communication can be reduced by more computation in one iteration.

- Synchronous

- Stale-synchronous

- Asynchronous

- Local SGD

- After each iteration, the workers synchronize their parameter updates.

- Every worker must wait for all workers to finish the transmission of all parameters in the current iteration, before the next training.

- Stragglers can influence the overall system throughput.

- High communication cost that limits the system scalability.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Alleviate the straggler problem without losing synchronization.

- The faster workers to do more updates than the slower workers to reduce the waiting time of the faster workers.

- Staleness bounded barrier to limit the iteration gap between the fastest worker and the slowest worker.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- For a maximum staleness bound $s$, the update formula of worker $i$ at iteration $t+1$:

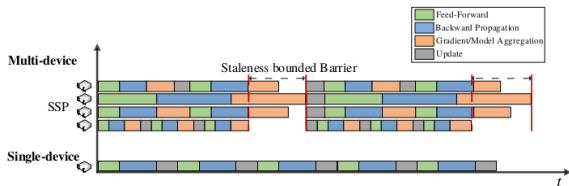- $\mathbf{w}_{i,t+1} := \mathbf{w}_0 - \eta(\sum_{k=1}^{t} \sum_{j=1}^{n} G_{j,k} + \sum_{k=t-s}^{t} G_{i,k} + \sum_{(j,k) \in S_{i,t+1}} G_{j,k})$

- The update has three parts:
  1. Guaranteed pre-window updates from clock $1$ to $t$ over all workers.
  2. Guaranteed read-my-writes in-window updates made by the querying worker $i$.
  3. Best-effort in-window updates. $S_{i,t+1}$ is some subset of the updates from other workers during period $[t-s]$.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- It completely eliminates the synchronization.

- Each work transmits its gradients to the PS after it calculates the gradients.

- The PS updates the global model without waiting for the other workers.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

▶ $\mathbf{w}_{t+1} := \mathbf{w}_t - \eta \sum_{i=1}^{n} G_{i,t-\tau_{k,i}}$

▶ $\tau_{k,i}$ is the time delay between the moment when worker `i` calculates the gradient at the current iteration.



[Tang et al., Communication-Efficient Distributed Deep Learning:  A Comprehensive Survey, 2020]

▶ All workers run several iterations, and then averages all local models into the newest global model.

▶ If $\mathcal{I}_T$ represents the synchronization timestamps, then:

$$\mathbf{w}_{i,t+1} = \begin{cases} \mathbf{w}_{i,t} - \eta G_{i,t} & \text{if } t+1 \notin \mathcal{I}_T \\ \mathbf{w}_{i,t} - \eta \frac{1}{n} \sum_{i=1}^{n} G_{i,t} & \text{if } t+1 \in \mathcal{I}_T \end{cases}$$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]
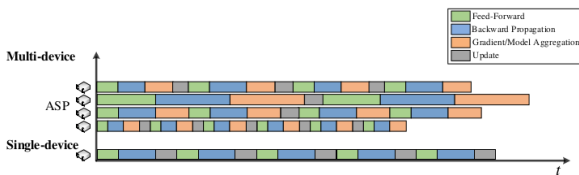
# Communication Compression

# Communication Compression

- ▶ Reduce the communication traffic with little impact on the model convergence.

- ▶ Compress the exchanged gradients or models before transmitting across the network.

- ▶ Quantization

- ▶ Sparsification

- Useing lower bits to represent the data.

- The gradients are of low precision.



One element
(32 bits)

Original
Gradient

Quantization
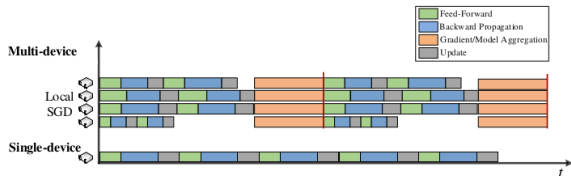
[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

# Communication Compression - Sparsification

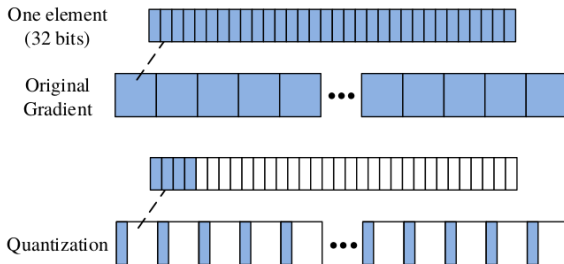- Reducing the number of elements that are transmitted at each iteration.

- Only significant gradients are required to update the model parameter to guarantee the convergence of the training.

- E.g., the zero-valued elements are no need to transmit.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

# Model Parallelism

# Model Parallelization

▶ The model is split across multiple devices.

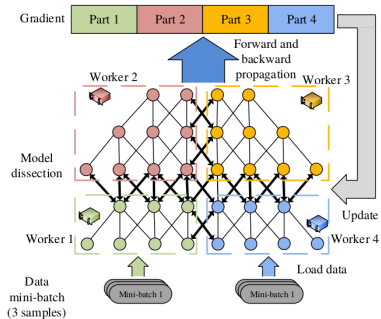▶ Depends on the architecture of the NN.



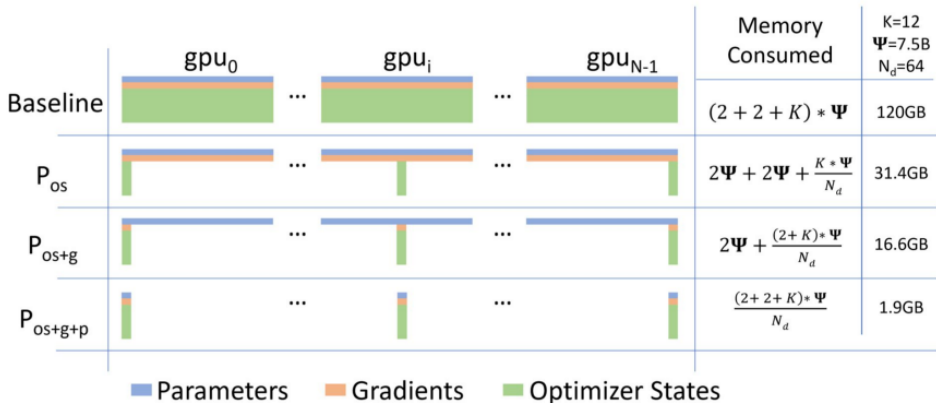[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

- Memory requirement from model states ($\Psi$ := # parameters)
  - Parameters (fp16): $2\Psi$
  - Gradients (fp16): $2\Psi$
  - Optimizer states: e.g. Adam $12\Psi$
    i. Parameters (fp32): $4\Psi$
    ii. Momentum (fp32): $4\Psi$
    iii. Variance (fp32): $4\Psi$

- Approach: partition each of them to all DP processes

| | gpu$_0$ | | gpu$_i$ | | gpu$_{N-1}$ | Memory Consumed | K=12 $\Psi$=7.5B $N_d$=64 |
|---|---|---|---|---|---|---|---|
| Baseline | | ... | | ... | | $(2 + 2 + K) * \Psi$ | 120GB |
| P$_{os}$ | | ... | | ... | | $2\Psi + 2\Psi + \frac{K * \Psi}{N_d}$ | 31.4GB |
| P$_{os+g}$ | | ... | | ... | | $2\Psi + \frac{(2 + K) * \Psi}{N_d}$ | 16.6GB |
| P$_{os+g+p}$ | | ... | | ... | | $\frac{(2 + 2 + K) * \Psi}{N_d}$ | 1.9GB |

■ Parameters  ■ Gradients  ■ Optimizer States

- Communication analysis ($\Psi$ := # parameters)
  - Baseline DP: one all-reduce, $2\Psi$
  - $P_{os+g}$: $2\Psi$
    i. Scatter-reduce on gradients: $\Psi$
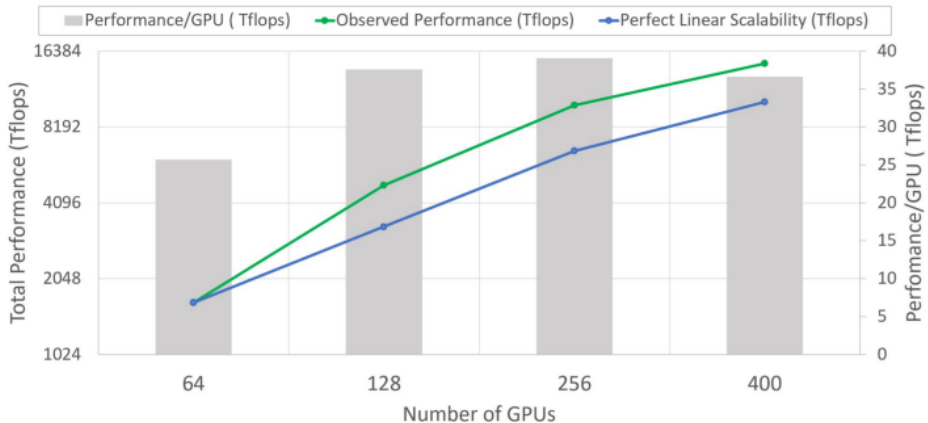    ii. All-gather on updated parameters: $\Psi$

- Communication analysis ($\Psi$ := # parameters)
  - Baseline DP: one all-reduce, $2\Psi$
  - $P_{os+g}$: $2\Psi$
    i. Scatter-reduce on gradients: $\Psi$
    ii. All-gather on updated parameters: $\Psi$
  - $P_{os+g+p}$: $3\Psi$ (1.5x communication)
    i. All-gather on parameters for forward: $\Psi$
    ii. All-gather on parameters for backward: $\Psi$
    iii. Scatter-reduce on gradients: $\Psi$

# Summary

# Summary

- Scalability matters

- Parallelization

- Data Parallelization
  - Parameter server vs. AllReduce
  - Synchronized vs. asynchronized

- Model Parallelization
  - DeepSpeed-Zero

Thanks!