



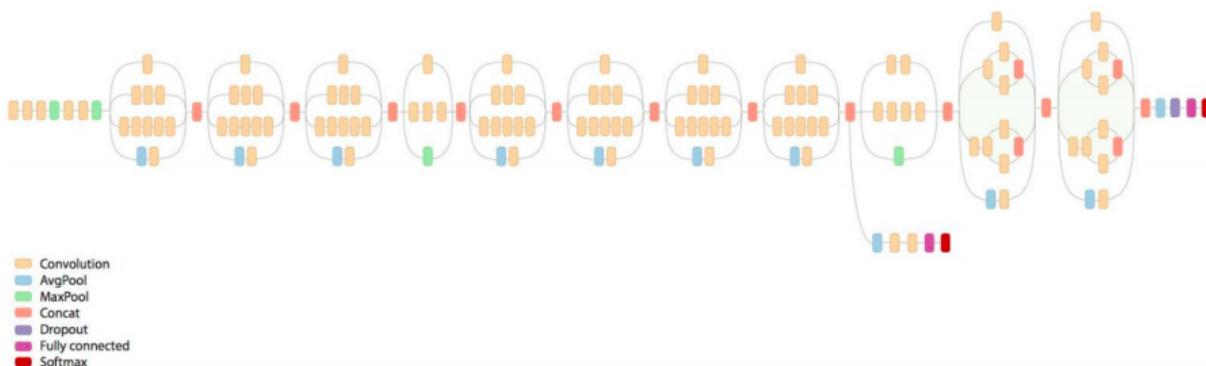
Distributed Learning - Data Parallelization

Amir H. Payberah
payberah@kth.se
2020-10-23



Training Deep Neural Networks

- ▶ Computationally intensive
- ▶ Time consuming



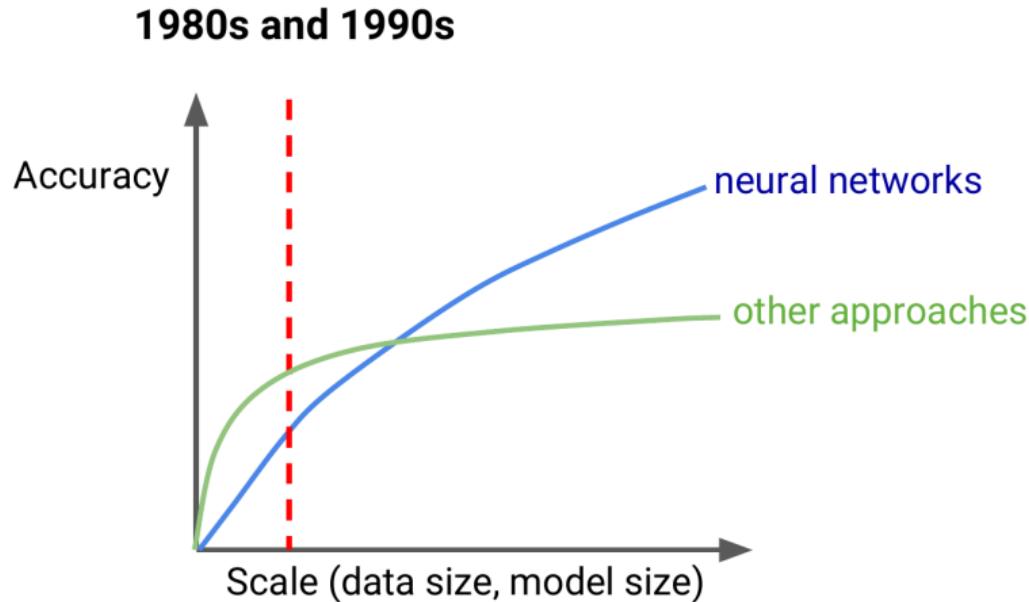
[<https://cloud.google.com/tpu/docs/images/inceptionv3onc--oview.png>]

Why?

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

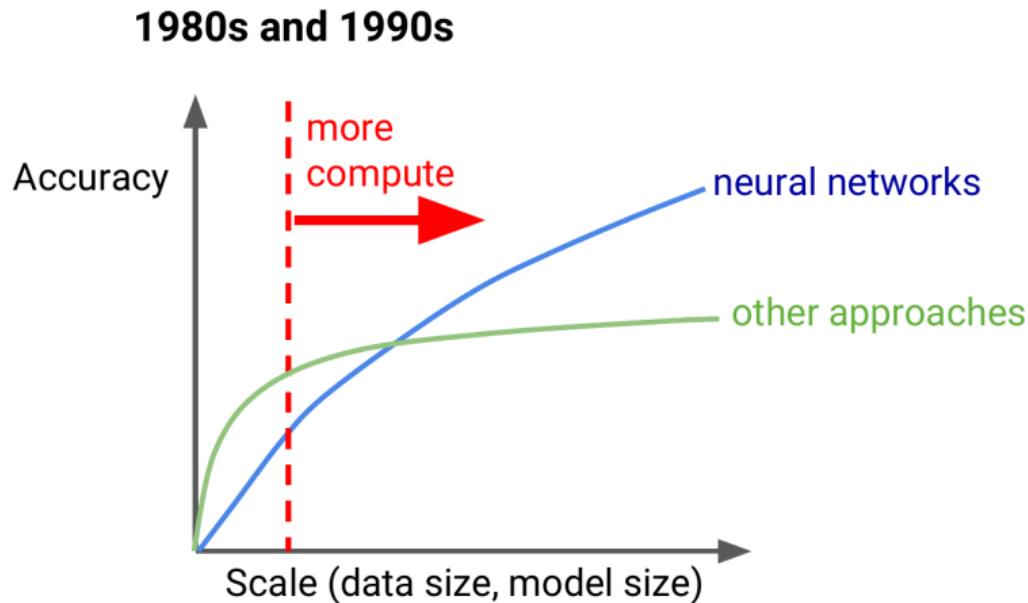


Accuracy vs. Data/Model Size



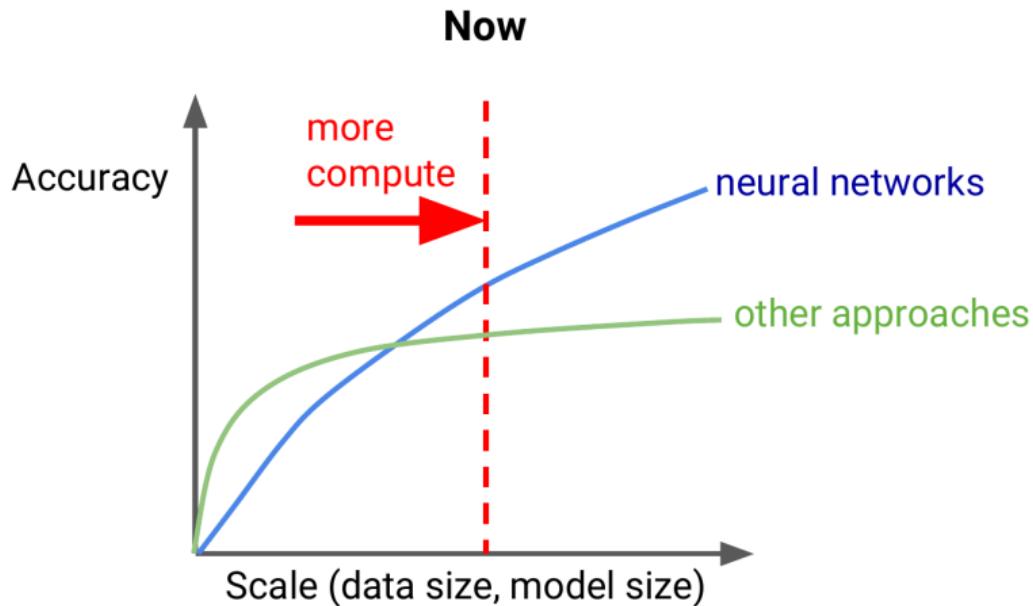
[Jeff Dean at AI Frontiers: Trends and Developments in Deep Learning Research]

Accuracy vs. Data/Model Size



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

Accuracy vs. Data/Model Size



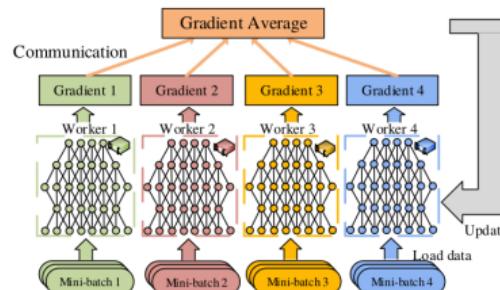
[Jeff Dean at AI Frontiers: Trends and Developments in Deep Learning Research]

Scalability



Data Parallelization (1/4)

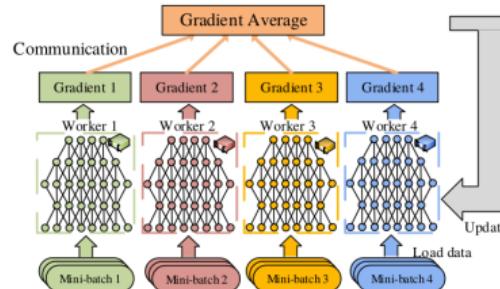
- ▶ Replicate a **whole model** on **every device**.



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

Data Parallelization (1/4)

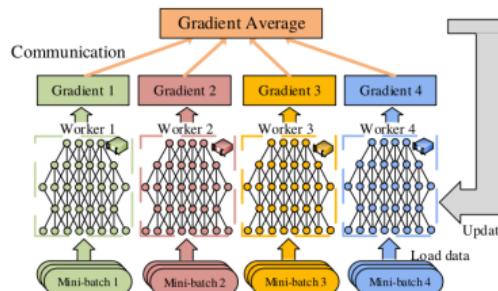
- ▶ Replicate a **whole model** on **every device**.
- ▶ Train **all replicas simultaneously**, using a **different mini-batch** for each.



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

Data Parallelization (2/4)

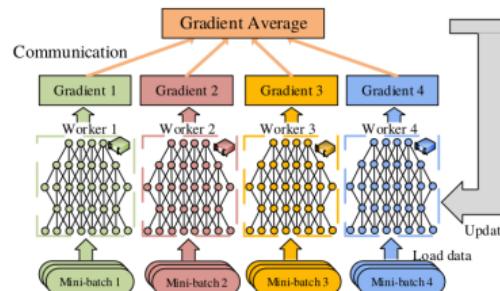
- ▶ k devices



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

Data Parallelization (2/4)

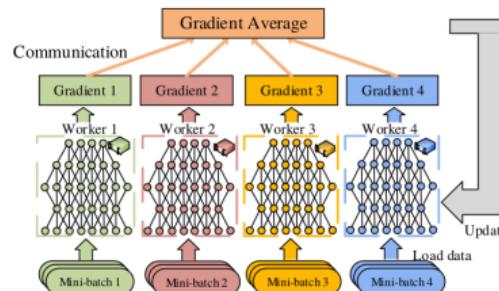
- ▶ k devices
- ▶ $J_i(\mathbf{w}) = \frac{1}{|\beta_i|} \sum_{x \in \beta_i} l(x, \mathbf{w}), \forall i = 1, 2, \dots, k$



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

Data Parallelization (2/4)

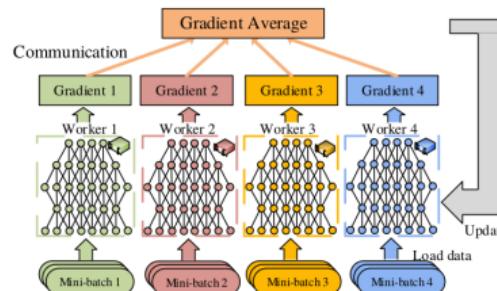
- ▶ k devices
- ▶ $J_i(\mathbf{w}) = \frac{1}{|\beta_i|} \sum_{\mathbf{x} \in \beta_i} l(\mathbf{x}, \mathbf{w}), \forall i = 1, 2, \dots, k$
- ▶ $G_i(\mathbf{w}, \beta_i) = \frac{1}{|\beta_i|} \sum_{\mathbf{x} \in \beta_i} \nabla l(\mathbf{w}, \mathbf{x})$



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

Data Parallelization (2/4)

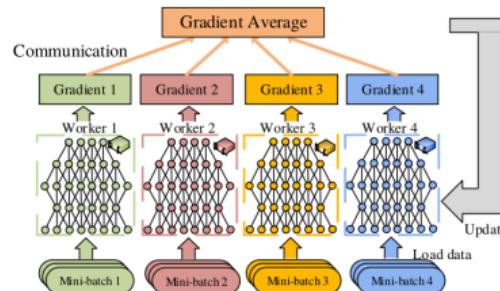
- ▶ k devices
- ▶ $J_i(\mathbf{w}) = \frac{1}{|\beta_i|} \sum_{\mathbf{x} \in \beta_i} l(\mathbf{x}, \mathbf{w}), \forall i = 1, 2, \dots, k$
- ▶ $G_i(\mathbf{w}, \beta_i) = \frac{1}{|\beta_i|} \sum_{\mathbf{x} \in \beta_i} \nabla l(\mathbf{w}, \mathbf{x})$
- ▶ $G_i(\mathbf{w}, \beta_i)$: the **local estimate** of the gradient of the loss function $\nabla J_i(\mathbf{w})$.



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

Data Parallelization (3/4)

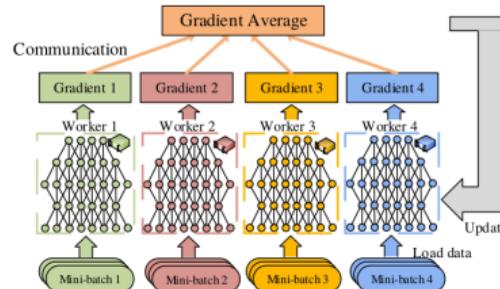
- ▶ Compute the gradients aggregation (e.g., mean of the gradients).
- ▶ $F(G_1, \dots, G_k) = \frac{1}{k} \sum_{i=1}^k G_i(\mathbf{w}, \beta_i)$



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

Data Parallelization (4/4)

- ▶ Update the model.
- ▶ $\mathbf{w} := \mathbf{w} - \eta F(\mathbf{G}_1, \dots, \mathbf{G}_k)$



[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**
- ▶ **Parallelism** of computations and communications



The Aggregation Algorithm



The Aggregation Algorithm

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



The Aggregation Algorithm

- ▶ How to aggregate gradients (compute the mean of the gradients)?
- ▶ Centralized - parameter server



The Aggregation Algorithm

- ▶ How to aggregate gradients (compute the mean of the gradients)?
- ▶ Centralized - parameter server
- ▶ Decentralized - all-reduce



The Aggregation Algorithm

- ▶ How to aggregate gradients (compute the mean of the gradients)?
- ▶ Centralized - parameter server
- ▶ Decentralized - all-reduce
- ▶ Decentralized - gossip

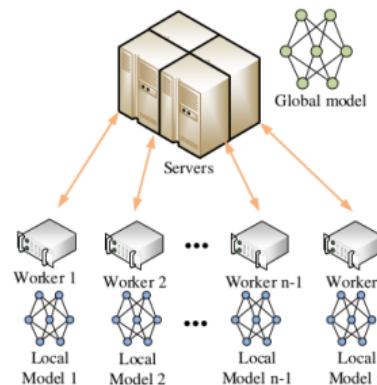


Aggregation - Centralized - Parameter Server

- ▶ Store the model parameters **outside of the workers**.

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]

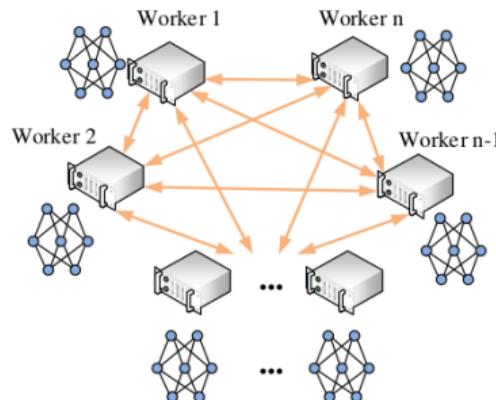


Aggregation - Distributed - All-Reduce

- ▶ Mirror all the model parameters across all workers (no PS).

Aggregation - Distributed - All-Reduce

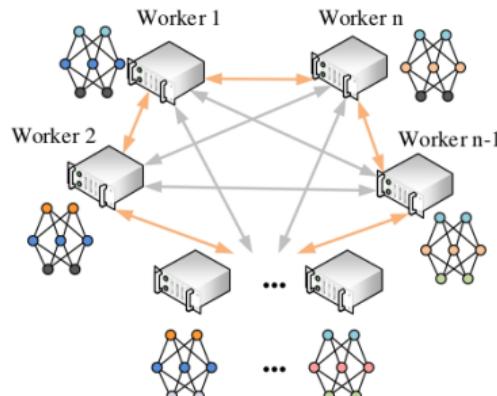
- ▶ Mirror all the model **parameters** across all workers (no PS).
- ▶ Workers **exchange** parameter updates **directly** via an **allreduce** operation.



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

Aggregation - Distributed - Gossip

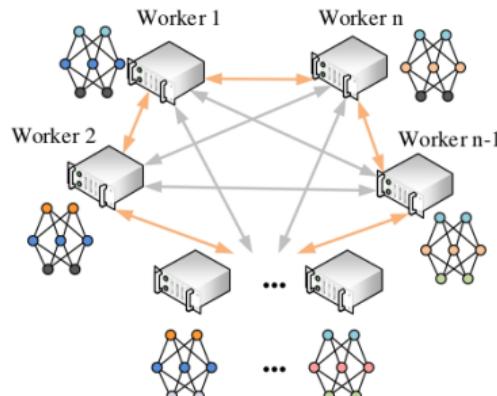
- ▶ No PS, and no global model.



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

Aggregation - Distributed - Gossip

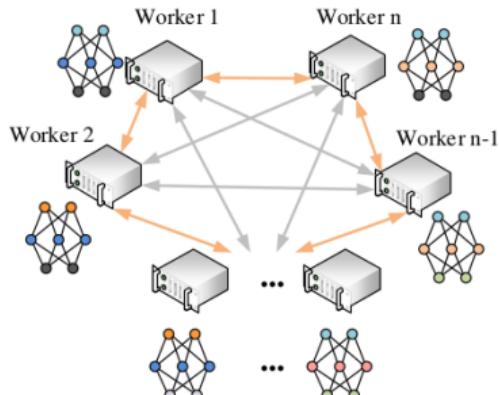
- ▶ No PS, and no global model.
- ▶ Every worker communicates updates with their neighbors.



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

Aggregation - Distributed - Gossip

- ▶ 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 and AllReduce (1/2)

- ▶ **Reduce**: reducing a **set of numbers** into a **smaller set of numbers** via a function.

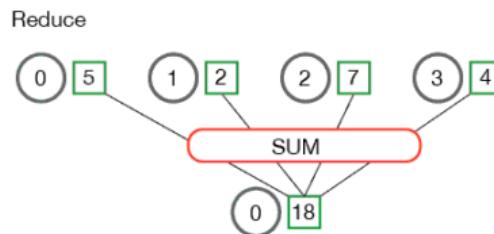


Reduce and AllReduce (1/2)

- ▶ **Reduce**: reducing a **set of numbers** into a **smaller set of numbers** via a function.
- ▶ E.g., `sum([1, 2, 3, 4, 5]) = 15`

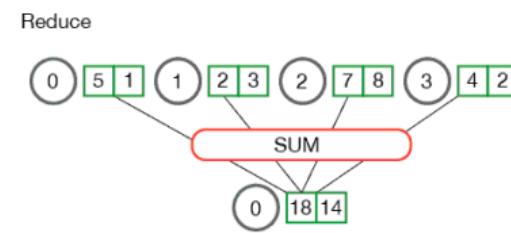
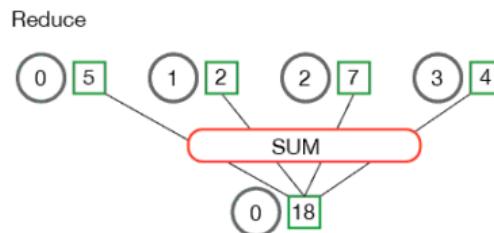
Reduce and AllReduce (1/2)

- ▶ **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**.



Reduce and AllReduce (1/2)

- ▶ **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>]

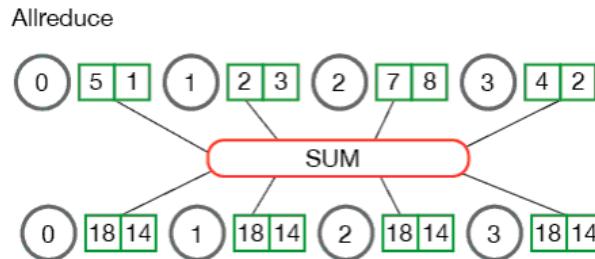


Reduce and AllReduce (2/2)

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

Reduce and AllReduce (2/2)

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



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

AllReduce Example

Initial state

Worker A
17 11 1 9

Worker B
5 13 23 14

Worker C
3 6 10 8

Worker D
12 7 2 12



After AllReduce operation

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>]

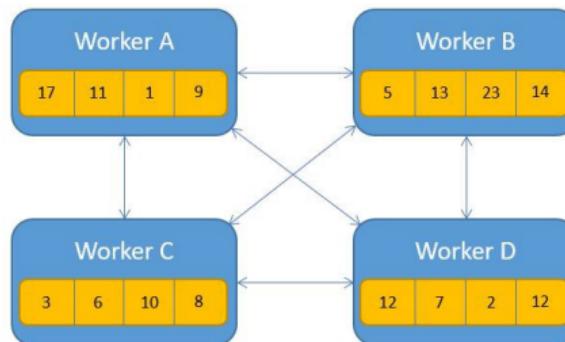


AllReduce Implementation

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

AllReduce Implementation - All-to-All AllReduce

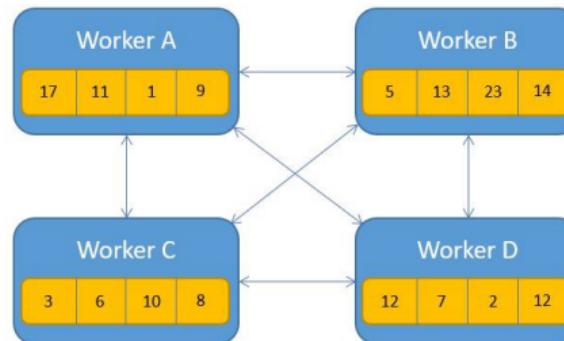
- ▶ Send the array of data to each other.
- ▶ Apply the reduction operation on each process.



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

AllReduce Implementation - All-to-All AllReduce

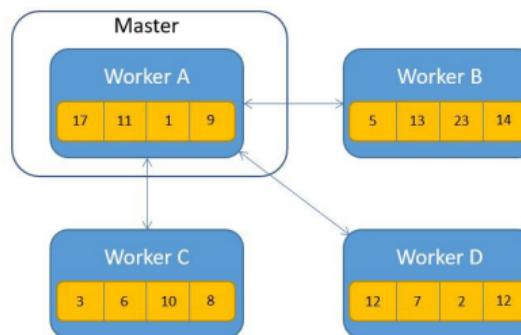
- ▶ Send the array of data to each other.
- ▶ Apply the reduction operation on each process.
- ▶ Too many unnecessary messages.



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

AllReduce Implementation - Master-Worker AllReduce

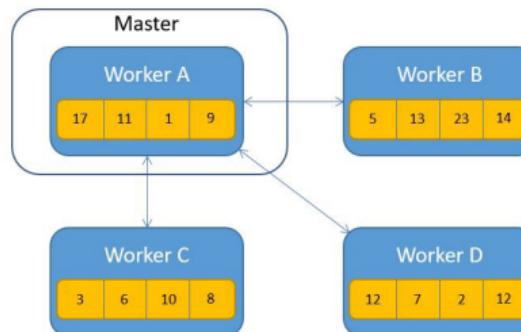
- ▶ Selecting one process as a **master**, gather all arrays into the master.
- ▶ Perform **reduction operations** locally in the **master**.
- ▶ Distribute the result to the **other processes**.



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

AllReduce Implementation - Master-Worker AllReduce

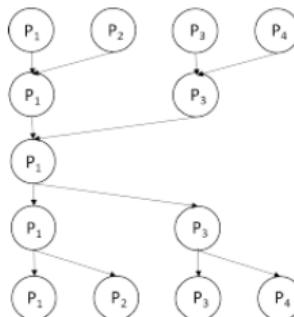
- ▶ Selecting one process as a **master**, gather all arrays into the master.
- ▶ Perform **reduction operations** locally in the **master**.
- ▶ Distribute the result to the **other processes**.
- ▶ The master becomes a **bottleneck** (**not scalable**).



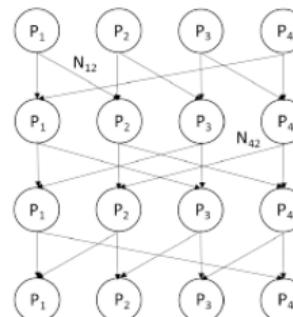
[<https://towardsdatascience.com/visual-intuition-on-ring-allreduce-for-distributed-deep-learning-d1f34b4911da>]

AllReduce Implementation - Other implementations

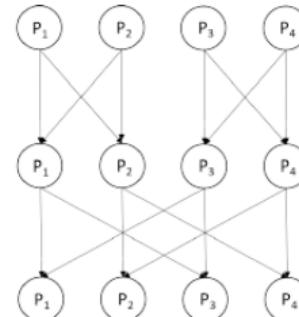
- ▶ Some try to minimize bandwidth.
- ▶ Some try to minimize latency.



(a) Tree AllReduce



(b) Round-robin AllReduce



(c) Butterfly AllReduce

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

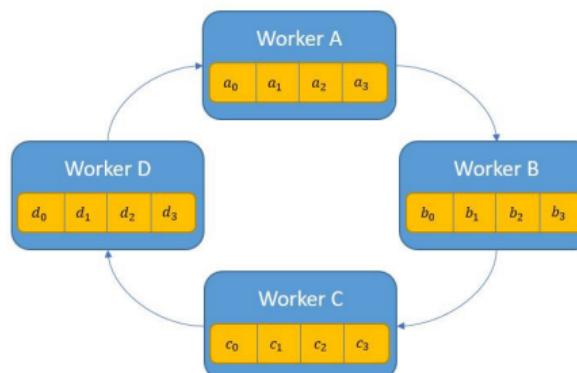


AllReduce Implementation - Ring-AllReduce (1/6)

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

AllReduce Implementation - Ring-AllReduce (2/6)

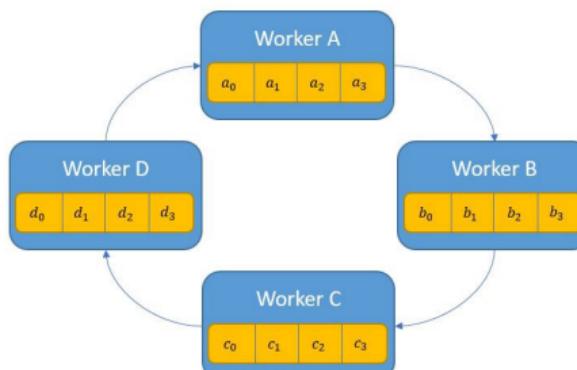
- In the **share-reduce** phase, each process p sends data to the process $(p+1) \% m$
 - m is the number of processes, and $\%$ is the modulo operator.



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

AllReduce Implementation - Ring-AllReduce (2/6)

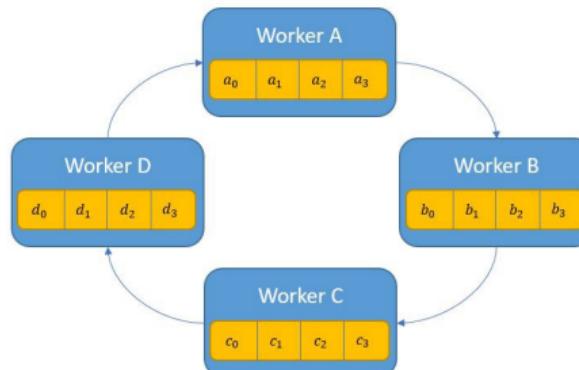
- In the **share-reduce** phase, each process p sends data to the process $(p+1) \% m$
 - m is the number of processes, and $\%$ is the modulo operator.
- The **array of data** on each process is divided to m chunks ($m=4$ here).



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

AllReduce Implementation - Ring-AllReduce (2/6)

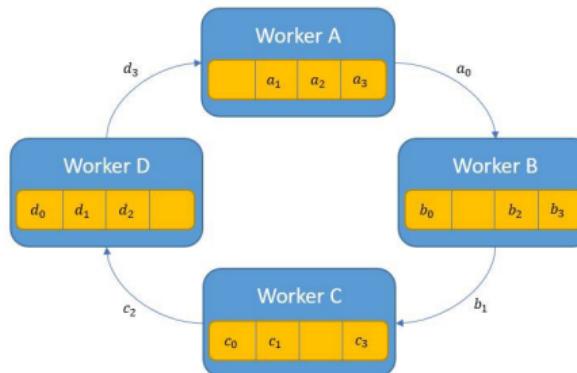
- In the **share-reduce** phase, each process p sends data to the process $(p+1) \% m$
 - m is the number of processes, and $\%$ is the modulo operator.
- The **array of data** on each process is divided to m chunks ($m=4$ here).
- Each one of these **chunks** will be **indexed** by i going forward.



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

AllReduce Implementation - Ring-AllReduce (3/6)

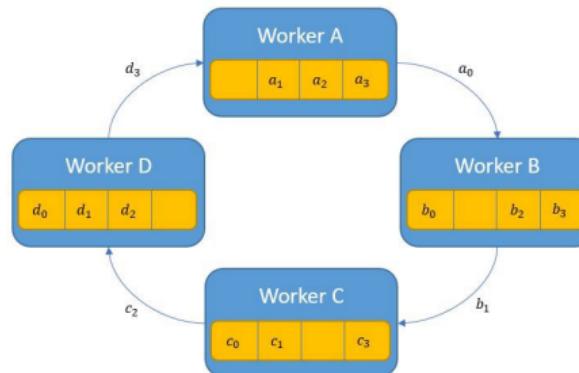
- In the first share-reduce step, process A sends a_0 to process B.



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

AllReduce Implementation - Ring-AllReduce (3/6)

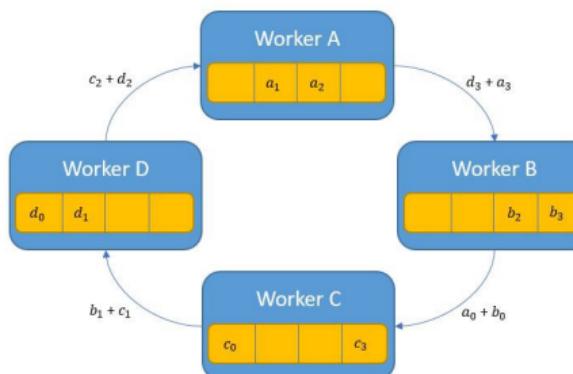
- ▶ In the **first share-reduce step**, process **A** sends **a₀** to process **B**.
- ▶ Process **B** sends **b₁** to process **C**, etc.



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

AllReduce Implementation - Ring-AllReduce (4/6)

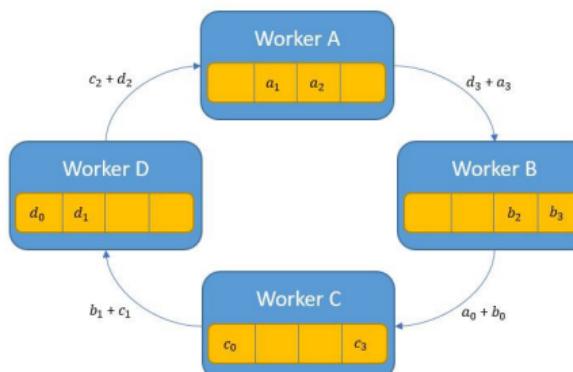
- When each process receives the data from the previous process, it applies the reduce operator (e.g., sum)



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

AllReduce Implementation - Ring-AllReduce (4/6)

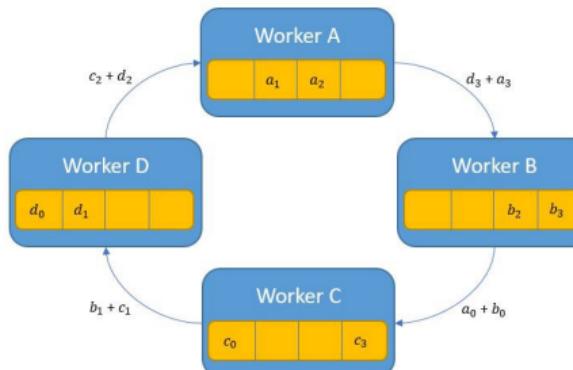
- When each process receives the data from the previous process, it applies the reduce operator (e.g., sum)
 - The reduce operator should be associative and commutative.



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

AllReduce Implementation - Ring-AllReduce (4/6)

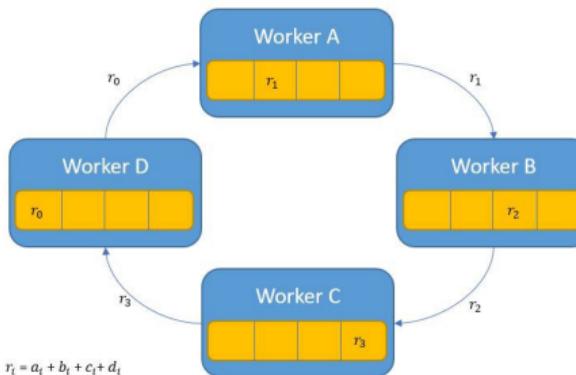
- ▶ When each process receives the data from the previous process, it applies the reduce operator (e.g., sum)
 - The reduce operator should be associative and commutative.
- ▶ It then proceeds to send it to the next process in the ring.



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

AllReduce Implementation - Ring-AllReduce (5/6)

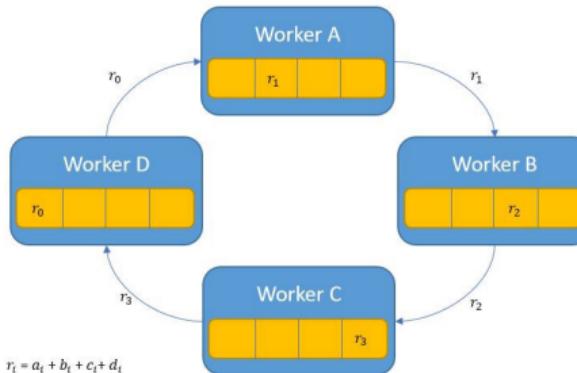
- The share-reduce phase finishes when each process holds the complete reduction of chunk i.



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

AllReduce Implementation - Ring-AllReduce (5/6)

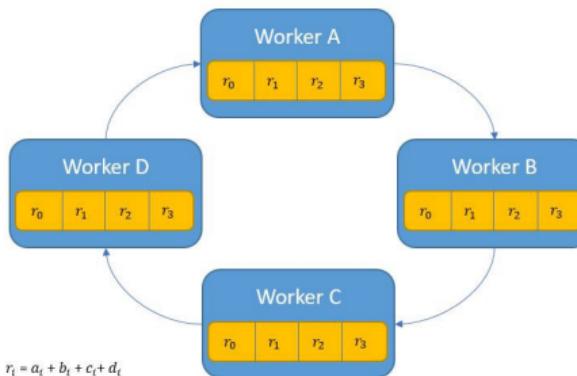
- ▶ The share-reduce phase finishes when each process holds the complete reduction of chunk i.
- ▶ At this point each process holds a part of the end result.



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

AllReduce Implementation - Ring-AllReduce (6/6)

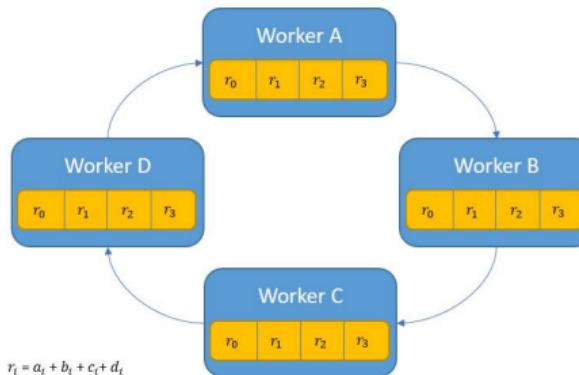
- The **share-only** step is the same process of sharing the data in a ring-like fashion without applying the reduce operation.



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

AllReduce Implementation - Ring-AllReduce (6/6)

- ▶ 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**.



[<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 vs. Ring-AllReduce

- ▶ N : number of elements, m : number of processes
- ▶ Master-Worker AllReduce



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.



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.



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 .



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



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.

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.

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



Synchronization

- ▶ 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**



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



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:



Reducing Synchronization Overhead

- ▶ Two directions for improvement:
 1. To **relax** the **synchronization** among all workers.



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.

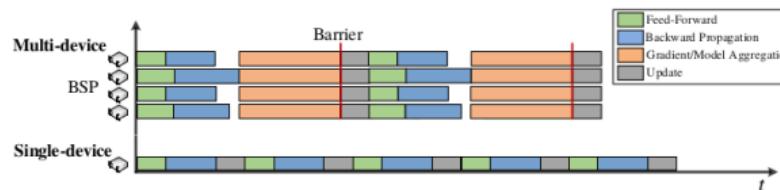


Communication Synchronization Models

- ▶ Synchronous
- ▶ Stale-synchronous
- ▶ Asynchronous
- ▶ Local SGD

Communication Synchronization - Synchronous

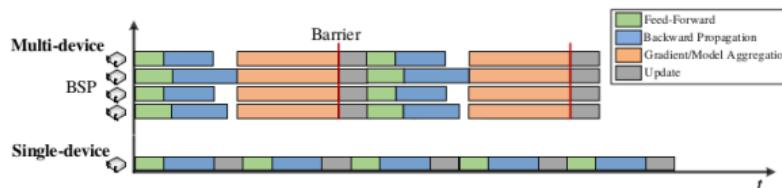
- ▶ After each **iteration**, the workers **synchronize** their parameter updates.



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

Communication Synchronization - Synchronous

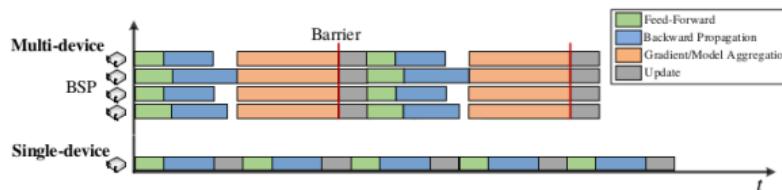
- ▶ After each **iteration**, the workers **synchronize** their parameter updates.
- ▶ Every worker must **wait** for **all workers** to **finish** the transmission of all parameters in the current iteration, before the **next training**.



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

Communication Synchronization - Synchronous

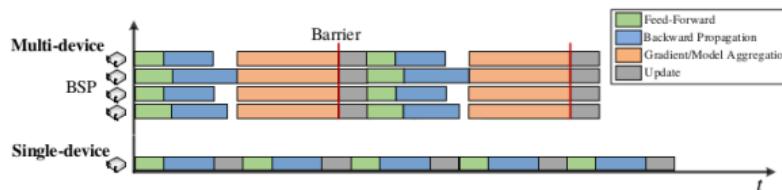
- ▶ After each **iteration**, the workers **synchronize** their parameter updates.
- ▶ Every worker must **wait** for **all workers** to **finish** the transmission of all parameters in the current iteration, before the **next training**.
- ▶ **Stragglers** can influence the overall system **throughput**.



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

Communication Synchronization - Synchronous

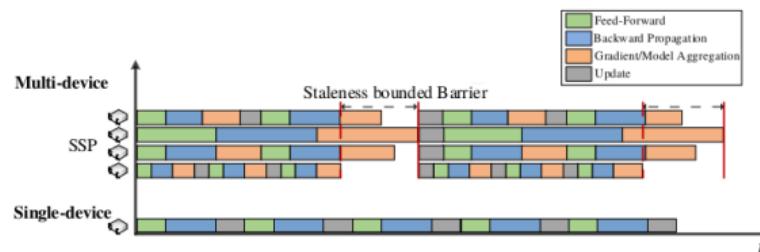
- ▶ After each **iteration**, the workers **synchronize** their parameter updates.
- ▶ Every worker must **wait** for **all workers** to **finish** the transmission of all parameters in the current iteration, before the **next training**.
- ▶ **Stragglers** can influence the overall system **throughput**.
- ▶ High **communication** cost that **limits** the system **scalability**.



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

Communication Synchronization - Stale Synchronous (1/2)

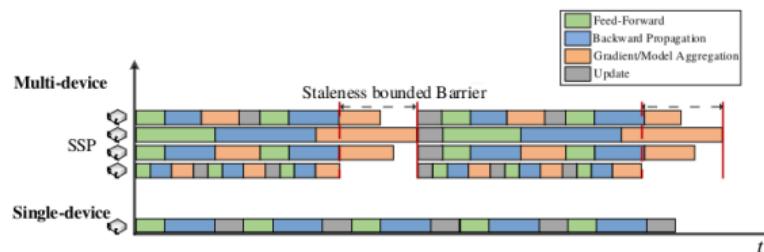
- ▶ Alleviate the straggler problem without losing synchronization.



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

Communication Synchronization - Stale Synchronous (1/2)

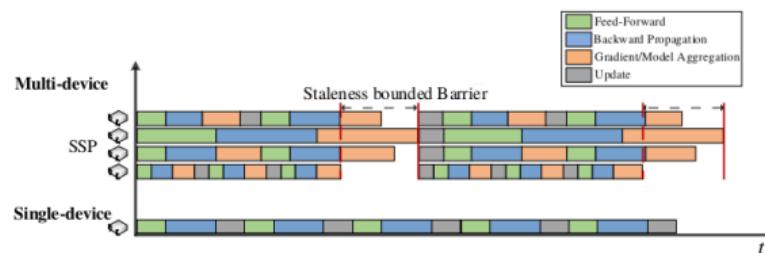
- ▶ Alleviate the straggler problem without losing synchronization.
- ▶ The faster workers to do **more updates** than the **slower workers** to reduce the waiting time of the faster workers.



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

Communication Synchronization - Stale Synchronous (1/2)

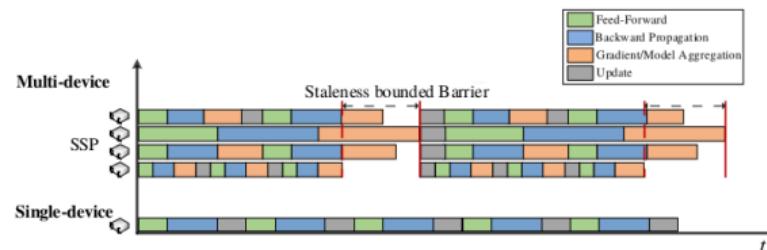
- ▶ Alleviate the straggler problem without losing synchronization.
- ▶ The faster workers to do **more updates** than the **slower workers** to reduce the waiting time of the faster workers.
- ▶ **Staleness bounded barrier** to limit the **iteration gap** between the fastest worker and the slowest worker.



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

Communication Synchronization - Stale Synchronous (2/2)

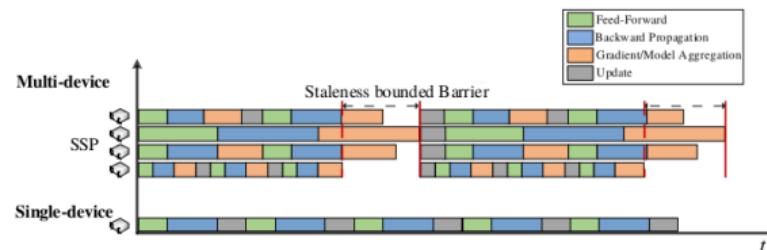
- ▶ For a maximum staleness bound s , the update formula of worker i at iteration $t+1$:
- ▶ $\mathbf{w}_{i,t+1} := \mathbf{w}_0 - \eta(\sum_{k=1}^t \sum_{j=1}^n G_{j,k} + \sum_{k=t-s}^t G_{i,k} + \sum_{(j,k) \in S_{i,t+1}} G_{j,k})$



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

Communication Synchronization - Stale Synchronous (2/2)

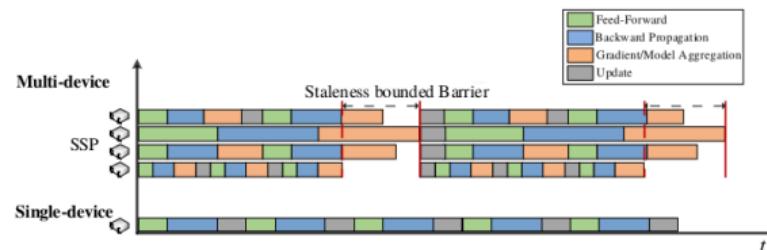
- ▶ For a maximum staleness bound s , the update formula of worker i at iteration $t+1$:
- ▶ $\mathbf{w}_{i,t+1} := \mathbf{w}_0 - \eta(\sum_{k=1}^t \sum_{j=1}^n G_{j,k} + \sum_{k=t-s}^t G_{i,k} + \sum_{(j,k) \in S_{i,t+1}} G_{j,k})$
- ▶ The update has three parts:



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

Communication Synchronization - Stale Synchronous (2/2)

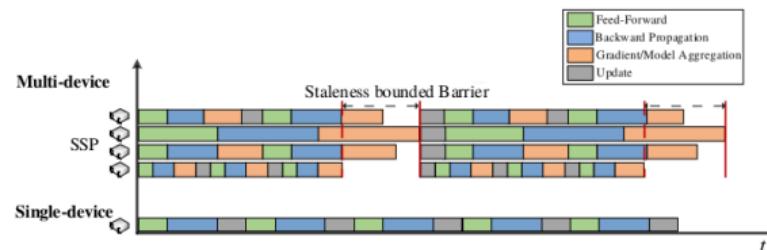
- ▶ For a maximum staleness bound s , the update formula of worker i at iteration $t+1$:
- ▶ $\mathbf{w}_{i,t+1} := \mathbf{w}_0 - \eta(\sum_{k=1}^t \sum_{j=1}^n G_{j,k} + \sum_{k=t-s}^t G_{i,k} + \sum_{(j,k) \in S_{i,t+1}} G_{j,k})$
- ▶ The update has three parts:
 1. Guaranteed pre-window updates from clock 1 to t over all workers.



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

Communication Synchronization - Stale Synchronous (2/2)

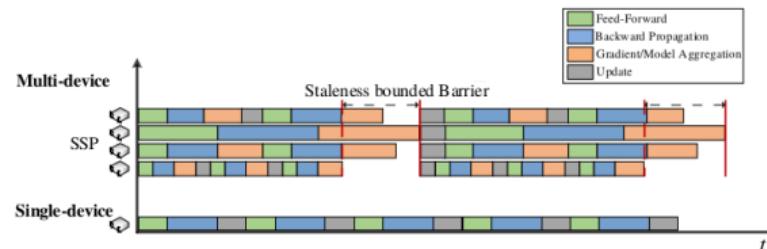
- ▶ For a maximum staleness bound s , the update formula of worker i at iteration $t+1$:
- ▶ $\mathbf{w}_{i,t+1} := \mathbf{w}_0 - \eta(\sum_{k=1}^t \sum_{j=1}^n G_{j,k} + \sum_{k=t-s}^t G_{i,k} + \sum_{(j,k) \in S_{i,t+1}} G_{j,k})$
- ▶ The update has three parts:
 1. Guaranteed pre-window updates from clock 1 to t over all workers.
 2. Guaranteed read-my-writes in-window updates made by the querying worker i .



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

Communication Synchronization - Stale Synchronous (2/2)

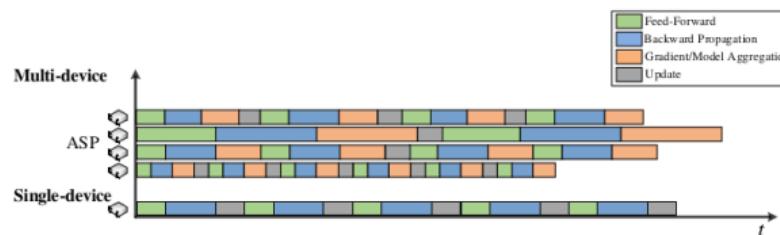
- ▶ For a maximum staleness bound s , the update formula of worker i at iteration $t+1$:
- ▶ $\mathbf{w}_{i,t+1} := \mathbf{w}_0 - \eta(\sum_{k=1}^t \sum_{j=1}^n G_{j,k} + \sum_{k=t-s}^t G_{i,k} + \sum_{(j,k) \in S_{i,t+1}} G_{j,k})$
- ▶ The update has three parts:
 1. Guaranteed pre-window updates from 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]

Communication Synchronization - Asynchronous (1/2)

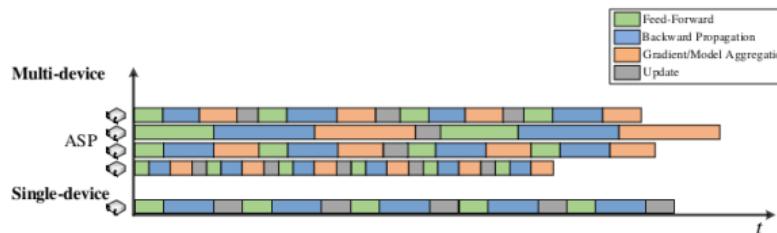
- ▶ It completely eliminates the synchronization.



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

Communication Synchronization - Asynchronous (1/2)

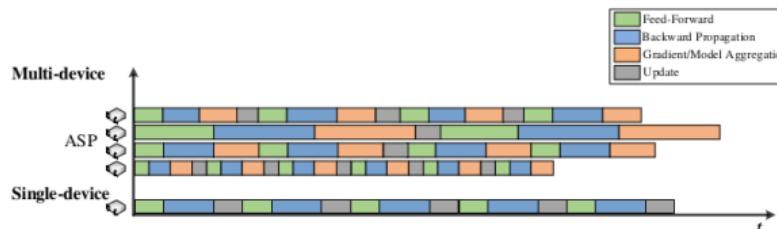
- ▶ It completely **eliminates** the synchronization.
- ▶ Each work **transmits its gradients** to the PS after it calculates the gradients.



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

Communication Synchronization - Asynchronous (1/2)

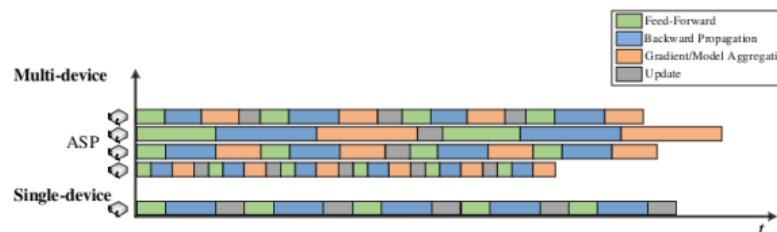
- ▶ It completely **eliminates** the synchronization.
- ▶ Each work **transmits its gradients** to the PS after it calculates the gradients.
- ▶ The PS updates the global model **without waiting** for the other workers.



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

Communication Synchronization - Asynchronous (2/2)

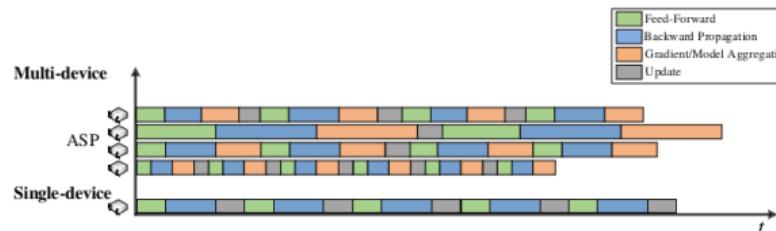
► $\mathbf{w}_{t+1} := \mathbf{w}_t - \eta \sum_{i=1}^n G_{i,t-\tau_{k,i}}$



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

Communication Synchronization - Asynchronous (2/2)

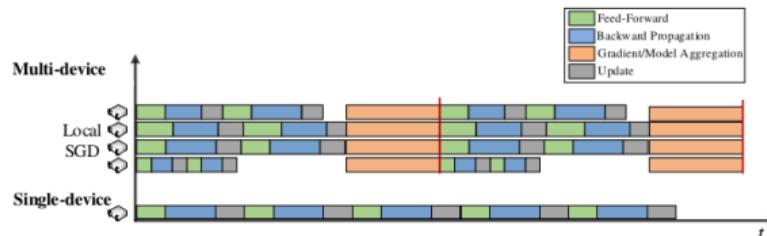
- ▶ $\mathbf{w}_{t+1} := \mathbf{w}_t - \eta \sum_{i=1}^n G_{i,t-\tau_{k,i}}$
- ▶ $\tau_{k,i}$ is the time delay between the moment when worker i calculates the gradient at the current iteration.



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

Communication Synchronization - Local SGD

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

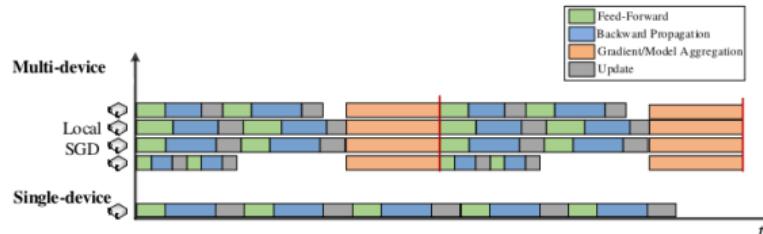


[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Communication Synchronization - Local SGD

- ▶ 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:

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



[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.



Communication Compression

- ▶ Reduce the communication traffic with **little impact** on the model convergence.
- ▶ Compress the exchanged gradients or models **before transmitting** across the network.



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

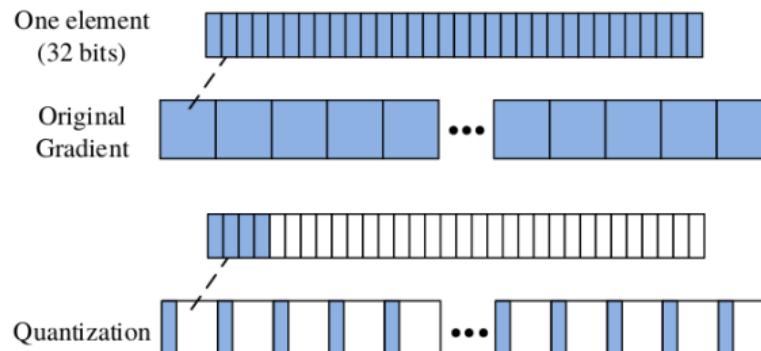


Communication Compression

- ▶ Reduce the communication traffic with **little impact** on the model convergence.
- ▶ Compress the exchanged gradients or models **before transmitting** across the network.
- ▶ Quantization
- ▶ Sparsification

Communication Compression - Quantization

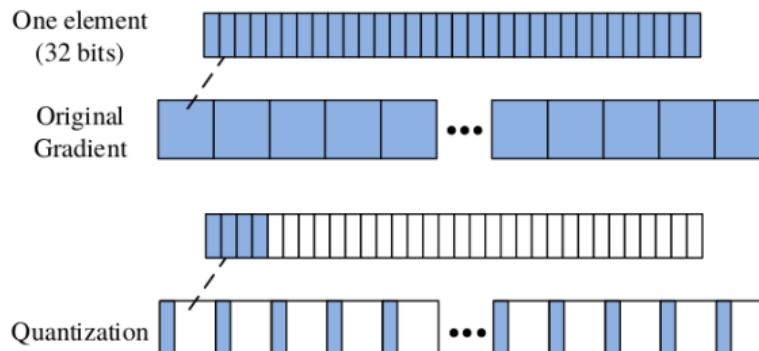
- ▶ Using lower bits to represent the data.



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

Communication Compression - Quantization

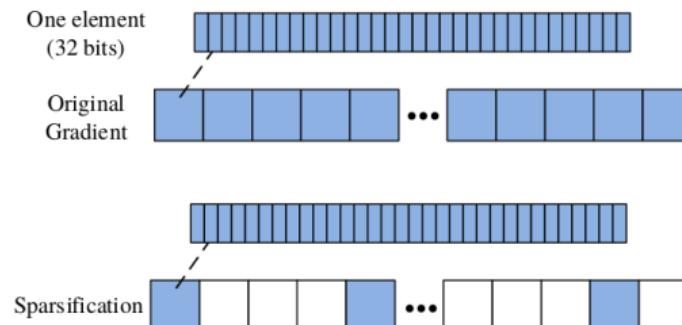
- ▶ Using lower bits to represent the data.
- ▶ The gradients are of low precision.



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

Communication Compression - Sparsification

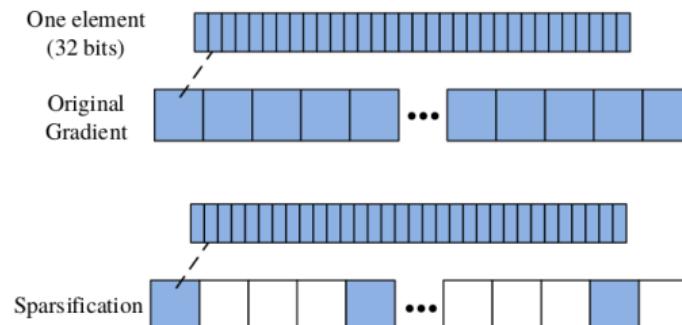
- ▶ Reducing the **number of elements** that are transmitted at each iteration.



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

Communication Compression - Sparsification

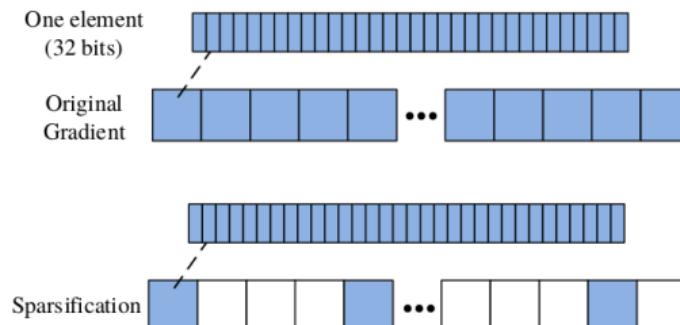
- ▶ Reducing the **number of elements** that are transmitted at each iteration.
- ▶ Only **significant gradients** are required to **update the model parameter** to guarantee the convergence of the training.



[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]



Parallelism of Computations and Communications



Parallelism of Computations and Communications (1/3)

- ▶ The layer-wise structure of deep models makes it possible to **parallels** the communication and computing tasks.

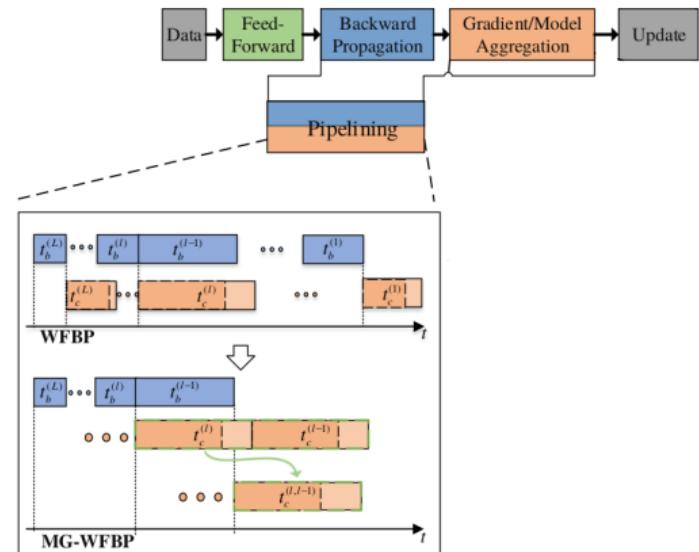


Parallelism of Computations and Communications (1/3)

- ▶ The layer-wise structure of deep models makes it possible to **parallelize** the communication and computing tasks.
- ▶ **Optimizing** the order of computation and communication such that the communication cost can be **minimized**

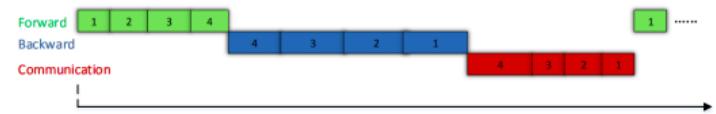
Parallelism of Computations and Communications (2/3)

- ▶ Wait-free backward propagation (WFBP)
- ▶ Merged-gradient WFBP (MG-WFBP)

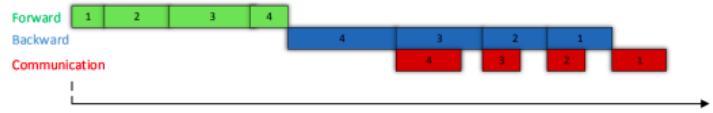
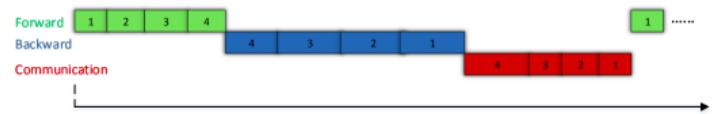


[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Parallelism of Computations and Communications (3/3)

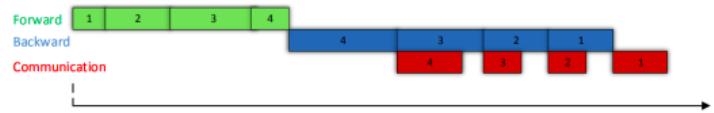
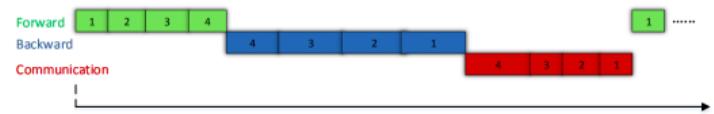


Parallelism of Computations and Communications (3/3)

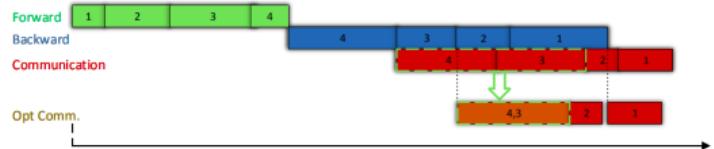


Wait-free backward propagation (WFBP)

Parallelism of Computations and Communications (3/3)



Wait-free backward propagation (WFBP)



Merged-gradient WFBP (MG-WFBP)

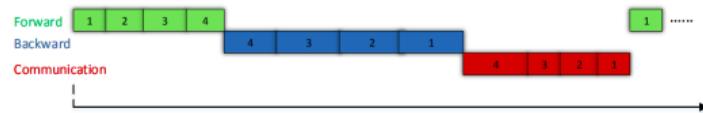
[shi et al., MG-WFBP: Efficient Data Communication for Distributed Synchronous SGD Algorithms, 2018]



TicTac: Accelerating Distributed Deep Learning with Communication Scheduling

Computation vs. Communication

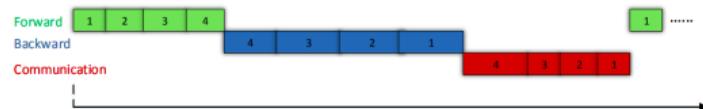
- ▶ The **iteration time** in deep learning systems depends on the time taken by
 1. Computation
 2. Communication
 3. The **overlap** between the two



[shi et al., MG-WFBP: Efficient Data Communication for Distributed Synchronous SGD Algorithms, 2018]

Computation vs. Communication

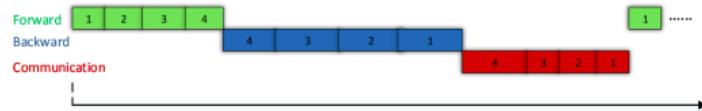
- ▶ The **iteration time** in deep learning systems depends on the time taken by
 1. Computation
 2. Communication
 3. The **overlap** between the two
- ▶ When **workers receive the parameters from the PS at the beginning of each iteration**, all parameters **are not used simultaneously**.



[shi et al., MG-WFBP: Efficient Data Communication for Distributed Synchronous SGD Algorithms, 2018]

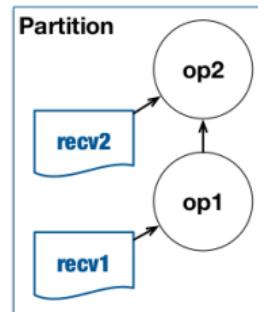
Computation vs. Communication

- ▶ The **iteration time** in deep learning systems depends on the time taken by
 1. Computation
 2. Communication
 3. The **overlap** between the two
- ▶ When **workers** receive the parameters from the PS at the beginning of each iteration, all parameters **are not used simultaneously**.
- ▶ Identifying the best **schedule** of parameter transfers is critical for **reducing** the **blocking on computation**.

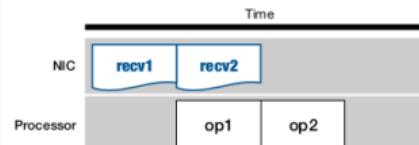


[shi et al., MG-WFBP: Efficient Data Communication for Distributed Synchronous SGD Algorithms, 2018]

Good vs. Bad Execution Order



(a) Toy Computational Graph



(b) Good Execution Order



(c) Bad Execution Order

[Hashemi et al., TicTac: Accelerating Distributed Deep Learning with Communication Scheduling, 2019]



High GPU Utilization

- ▶ High GPU utilization can be achieved in two ways:



High GPU Utilization

- ▶ High GPU utilization can be achieved in two ways:
 1. When total communication time is less than or equal to the computation time.



High GPU Utilization

- ▶ High GPU utilization can be achieved in two ways:
 1. When total communication time is less than or equal to the computation time.
 2. With efficient overlap of communication and computation.



High GPU Utilization

- ▶ High GPU utilization can be achieved in two ways:
 1. When total communication time is less than or equal to the computation time.
 2. With efficient overlap of communication and computation.
- ▶ Techniques improve GPU utilization:



High GPU Utilization

- ▶ High GPU utilization can be achieved in two ways:
 1. When total communication time is less than or equal to the computation time.
 2. With efficient overlap of communication and computation.
- ▶ Techniques improve GPU utilization:
 - Increasing computation time



High GPU Utilization

- ▶ High GPU utilization can be achieved in two ways:
 1. When total communication time is less than or equal to the computation time.
 2. With efficient overlap of communication and computation.
- ▶ Techniques improve GPU utilization:
 - Increasing computation time
 - Decreasing communication time



High GPU Utilization

- ▶ High GPU utilization can be achieved in two ways:
 1. When total communication time is less than or equal to the computation time.
 2. With efficient overlap of communication and computation.
- ▶ Techniques improve GPU utilization:
 - Increasing computation time
 - Decreasing communication time
 - Better interleaving of computation and communication



Overlap Coefficient (1/2)

- Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$



Overlap Coefficient (1/2)

- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$
- ▶ T: the actual iteration time

Overlap Coefficient (1/2)

- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$
- ▶ T: the **actual iteration time**
- ▶ N: the **communication time**



Overlap Coefficient (1/2)

- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$
- ▶ T: the **actual iteration time**
- ▶ N: the **communication time**
- ▶ C: the **computation time**



Overlap Coefficient (1/2)

- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$
- ▶ T: the **actual iteration time**
- ▶ N: the **communication time**
- ▶ C: the **computation time**
- ▶ **N + C** is the iteration time when there is **no overlap**



Overlap Coefficient (2/2)

- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$

Overlap Coefficient (2/2)

- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$
- ▶ The maximum overlap possible is given by $\min(N,C)$, which is achieved when the smaller quantity completely overlaps with the large quantity.

Overlap Coefficient (2/2)

- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$
- ▶ The maximum overlap possible is given by $\min(N,C)$, which is achieved when the smaller quantity completely overlaps with the large quantity.
- ▶ GPU utilization: $U = \frac{C}{T} = \frac{C}{N+C-\alpha\min(N,C)} = \frac{1}{1+\rho-\alpha\min(\rho,1)}$
- ▶ $\rho = \frac{N}{C}$: the communication/computation ratio



Scheduling Algorithm

- ▶ Prioritize transfers that speed up the critical path in the DAG, by reducing blocking on computation caused by parameter transfers.



Scheduling Algorithm

- ▶ Prioritize transfers that speed up the critical path in the DAG, by reducing blocking on computation caused by parameter transfers.
- ▶ TIC: Timing-Independent Communication scheduling



Scheduling Algorithm

- ▶ Prioritize transfers that speed up the critical path in the DAG, by reducing blocking on computation caused by parameter transfers.
- ▶ TIC: Timing-Independent Communication scheduling
- ▶ TAC: Timing-Aware Communication scheduling



Scheduling Algorithm

- ▶ Prioritize transfers that speed up the critical path in the DAG, by reducing blocking on computation caused by parameter transfers.
- ▶ TIC: Timing-Independent Communication scheduling
 - Prioritize those transfers that reduces blocking on network transfers.
- ▶ TAC: Timing-Aware Communication scheduling



Scheduling Algorithm

- ▶ Prioritize transfers that speed up the critical path in the DAG, by reducing blocking on computation caused by parameter transfers.
- ▶ TIC: Timing-Independent Communication scheduling
 - Prioritize those transfers that reduces blocking on network transfers.
- ▶ TAC: Timing-Aware Communication scheduling
 - Prioritize those transfers that reduces the blocking of computation.



TIC

- ▶ Prioritize those transfers that reduces blocking on network transfers.



TIC

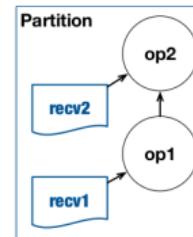
- ▶ Prioritize those transfers that reduces blocking on network transfers.
- ▶ Prioritize based only on vertex dependencies in the DAG.



- ▶ Prioritize those transfers that reduces blocking on network transfers.
- ▶ Prioritize based only on vertex dependencies in the DAG.
- ▶ Higher priorities are given to transfers that are least blocking on computation.

- ▶ Prioritize those transfers that reduces blocking on network transfers.
- ▶ Prioritize based only on vertex dependencies in the DAG.
- ▶ Higher priorities are given to transfers that are least blocking on computation.
- ▶ Ignore the ops time, and use the number of communication ops instead.

- ▶ Prioritize those transfers that reduces blocking on network transfers.
- ▶ Prioritize based only on vertex dependencies in the DAG.
- ▶ Higher priorities are given to transfers that are least blocking on computation.
- ▶ Ignore the ops time, and use the number of communication ops instead.
- ▶ E.g., $op_1.M = \text{Time}(recv_1)$ and $op_2.M = \text{Time}(recv_1) + \text{Time}(recv_2)$.





TAC

- ▶ Prioritize those transfers that reduces the blocking of computation.

- ▶ Prioritize those transfers that reduces the blocking of computation.
- ▶ Prioritize transfers that maximize α by using information on:

- ▶ Prioritize those transfers that reduces the blocking of computation.
- ▶ Prioritize transfers that maximize α by using information on:
 1. Vertex dependencies among ops specified by the computational DAG.



- ▶ Prioritize those transfers that reduces the blocking of computation.
- ▶ Prioritize transfers that maximize α by using information on:
 1. Vertex dependencies among ops specified by the computational DAG.
 2. Execution time of each op estimated with time oracle.

- ▶ Prioritize those transfers that reduces the blocking of computation.
- ▶ Prioritize transfers that maximize α by using information on:
 1. Vertex dependencies among ops specified by the computational DAG.
 2. Execution time of each op estimated with time oracle.
- ▶ To achieve this goal, the algorithm focuses on two cases:
 1. Any communication and computation overlapping?
 2. If no, choose one which eliminates the computation block sooner.



CARAMEL: Accelerating Decentralized Distributed Deep Learning with Computation Scheduling



CARAMEL

- ▶ Decentralized aggregation (no PS)



CARAMEL

- ▶ Decentralized aggregation (no PS)
- ▶ Improve efficiency of decentralized DNN training



CARAMEL

- ▶ Decentralized aggregation (no PS)
- ▶ Improve efficiency of decentralized DNN training
- ▶ In terms of iteration time and GPU utilization



CARAMEL

- ▶ Decentralized aggregation (no PS)
- ▶ Improve efficiency of decentralized DNN training
- ▶ In terms of iteration time and GPU utilization
- ▶ CARAMEL achieves this goal through:



CARAMEL

- ▶ Decentralized aggregation (no PS)
- ▶ Improve efficiency of decentralized DNN training
- ▶ In terms of iteration time and GPU utilization
- ▶ CARAMEL achieves this goal through:
 1. Computation scheduling that expands the feasible window of transfer for each parameter (transfer boundaries)



CARAMEL

- ▶ Decentralized aggregation (no PS)
- ▶ Improve efficiency of decentralized DNN training
- ▶ In terms of iteration time and GPU utilization
- ▶ CARAMEL achieves this goal through:
 1. Computation scheduling that expands the feasible window of transfer for each parameter (transfer boundaries)
 2. Network optimizations that smoothen the load



Computation Scheduling (1/2)

- ▶ In decentralized aggregation, all workers should have the parameter available for aggregation before the transfer can be initiated.



Computation Scheduling (1/2)

- ▶ In decentralized aggregation, all workers should have the parameter available for aggregation before the transfer can be initiated.
- ▶ There are multiple feasible orders for executing operations in a DAG.



Computation Scheduling (1/2)

- ▶ In decentralized aggregation, all workers should have the parameter available for aggregation before the transfer can be initiated.
- ▶ There are multiple feasible orders for executing operations in a DAG.
- ▶ The parameters may become available at different workers in varying orders.



Computation Scheduling (1/2)

- ▶ In **decentralized aggregation**, all **workers** should have the parameter available for aggregation before the **transfer** can be initiated.
- ▶ There are **multiple feasible orders** for executing operations in a DAG.
- ▶ The **parameters** may become available at different workers in **varying orders**.
- ▶ The **transfer boundaries** of a parameter represent the **window** when a **parameter** can be aggregated **without blocking** computation.



Computation Scheduling (1/2)

- ▶ In **decentralized aggregation**, all **workers** should have the parameter available for aggregation before the **transfer** can be initiated.
- ▶ There are **multiple feasible orders** for executing operations in a DAG.
- ▶ The **parameters** may become available at different workers in **varying orders**.
- ▶ The **transfer boundaries** of a parameter represent the **window** when a **parameter** can be aggregated **without blocking** computation.
- ▶ The **start boundary** is determined by the **completion** of the computation operation that **updates the parameter**.

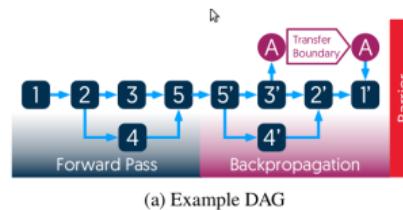


Computation Scheduling (1/2)

- ▶ In **decentralized aggregation**, all **workers** should have the parameter available for aggregation before the **transfer** can be initiated.
- ▶ There are **multiple feasible orders** for executing operations in a DAG.
- ▶ The **parameters** may become available at different workers in **varying orders**.
- ▶ The **transfer boundaries** of a parameter represent the **window** when a **parameter** can be aggregated **without blocking** computation.
- ▶ The **start boundary** is determined by the **completion** of the computation operation that **updates the parameter**.
- ▶ The **end boundary** is the computation operation that **reads the parameter**.

Computation Scheduling (2/2)

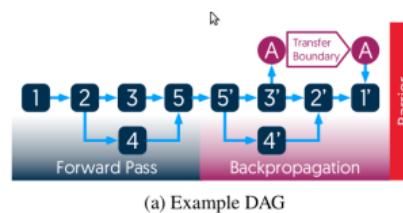
- ▶ CARAMEL expands these boundaries through **scheduling optimizations** of the **computation DAG**, where it



[Hashemi et al., CARAMEL: Accelerating Decentralized Distributed Deep Learning with Model-Aware Scheduling, 2020]

Computation Scheduling (2/2)

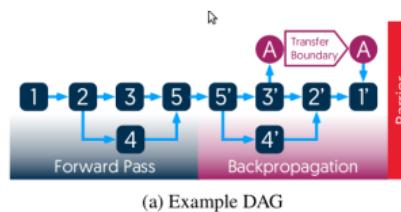
- ▶ CARAMEL expands these boundaries through **scheduling optimizations** of the **computation DAG**, where it
 1. Moves the **start boundaries earlier**.



[Hashemi et al., CARAMEL: Accelerating Decentralized Distributed Deep Learning with Model-Aware Scheduling, 2020]

Computation Scheduling (2/2)

- CARAMEL expands these boundaries through **scheduling optimizations** of the **computation DAG**, where it
 1. Moves the **start boundaries earlier**.
 2. Pushes the **end boundary** by postponing the execution of some computation operations to the **forward pass of next iteration**.



(a) Example DAG



(b) Best Schedule



(c) Worst Schedule

[Hashemi et al., CARAMEL: Accelerating Decentralized Distributed Deep Learning with Model-Aware Scheduling, 2020]



Network Optimization

- ▶ Optimizations for **smoothening the network load** include:



Network Optimization

- ▶ Optimizations for **smoothening the network load** include:
 1. **Batching** of **small parameters** to reduce the network overhead.



Network Optimization

- ▶ Optimizations for **smoothening the network load** include:
 1. Batching of **small parameters** to reduce the network overhead.
 2. Adaptive **splitting and pipelining of parameters** to accelerate aggregation of large data.



Defining the Environment

- ▶ T: the **actual iteration** time
- ▶ N: the **communication** time
- ▶ C: the **computation** time



Defining the Environment

- ▶ T : the **actual iteration** time
- ▶ N : the **communication** time
- ▶ C : the **computation** time
- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$

Defining the Environment

- ▶ T : the **actual iteration** time
- ▶ N : the **communication** time
- ▶ C : the **computation** time
- ▶ Overlap coefficient: $\alpha = \frac{N+C-T}{\min(N,C)}$
- ▶ GPU utilization: $U = \frac{C}{T} = \frac{C}{N+C-\alpha\min(N,C)} = \frac{1}{1+\rho-\alpha\min(\rho,1)}$
- ▶ $\rho = \frac{N}{C}$: the communication/computation ratio



CARAMEL Algorithm

- ▶ Dataflow DAG Optimizer
- ▶ Network Transfer Scheduler
- ▶ Parameter Batcher
- ▶ Adaptive Depth Enforcer



Dataflow DAG Optimizer

- ▶ Stage 1: Determining the **best** order.



Dataflow DAG Optimizer

- ▶ Stage 1: Determining the **best** order.
- ▶ Stage 2: Enforcing the **best** order.



Dataflow DAG Optimizer

- ▶ Stage 1: Determining the best order.
 - Increasing the overlap coefficient α by prioritizing the computations that activates the communication operations as early as possible.
- ▶ Stage 2: Enforcing the best order.



Dataflow DAG Optimizer

- ▶ Stage 1: Determining the **best order**.
 - Increasing the overlap coefficient α by prioritizing the **computations** that activates the **communication** operations **as early as possible**.
- ▶ Stage 2: Enforcing the **best order**.
 - Iteratively activate parameters in the **best order** chosen in the previous stage.

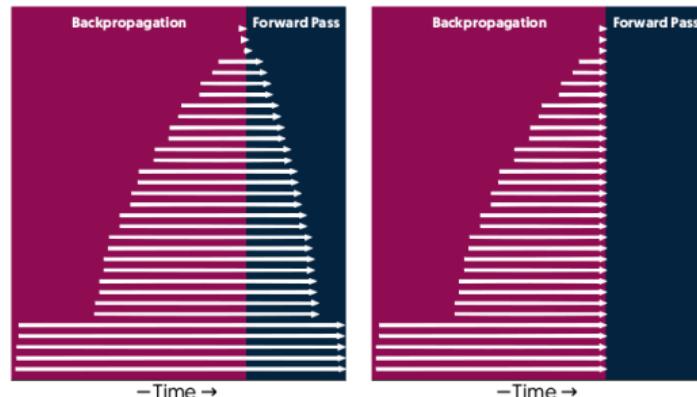


Dataflow DAG Optimizer

- ▶ Stage 1: Determining the **best order**.
 - Increasing the overlap coefficient α by prioritizing the **computations** that activates the **communication** operations **as early as possible**.
- ▶ Stage 2: Enforcing the **best order**.
 - Iteratively activate parameters in the **best order** chosen in the previous stage.
 - Ensuring that at each given time, **only ops needed for the target parameter update** can be executed.

Network Transfer Scheduler

- ▶ Increasing the overlap coefficient α by scheduling parameter transfers efficiently.
- ▶ Transfers are scheduled in both backward pass and forward pass.



(a) Parameter Server

(b) Decentralized aggregation

[Hashemi et al., CARAMEL: Accelerating Decentralized Distributed Deep Learning with Model-Aware Scheduling, 2020]



Parameter Batcher

- ▶ Small parameters incur large overhead.



Parameter Batcher

- ▶ Small parameters incur large overhead.
- ▶ Combining small parameters in to groups.



Parameter Batcher

- ▶ Small parameters incur large overhead.
- ▶ Combining small parameters in to groups.
- ▶ Parameters larger than a certain threshold are transferred without batching.



Adaptive Depth Enforcer

- ▶ Two stages in decentralized algorithms: **transferring** and **aggregating** data across nodes.



Adaptive Depth Enforcer

- ▶ Two stages in decentralized algorithms: **transferring** and **aggregating** data across nodes.
 - In each step, data is transferred on the network, and is sent to application to be reduced, before the result is sent again over the network.



Adaptive Depth Enforcer

- ▶ Two stages in decentralized algorithms: **transferring** and **aggregating** data across nodes.
 - In each step, data is transferred on the network, and is sent to application to be reduced, before the result is sent again over the network.
- ▶ This process **reduces the network utilization** since the network is **not utilized** during the **reduction** at the application layer.



Adaptive Depth Enforcer

- ▶ Two stages in decentralized algorithms: **transferring** and **aggregating** data across nodes.
 - In each step, data is transferred on the network, and is sent to application to be reduced, before the result is sent again over the network.
- ▶ This process **reduces the network utilization** since the network is **not utilized** during the **reduction** at the application layer.
- ▶ Chunk (break) the data in to a few pieces, and transfer each **chunk** independently in parallel.



Adaptive Depth Enforcer

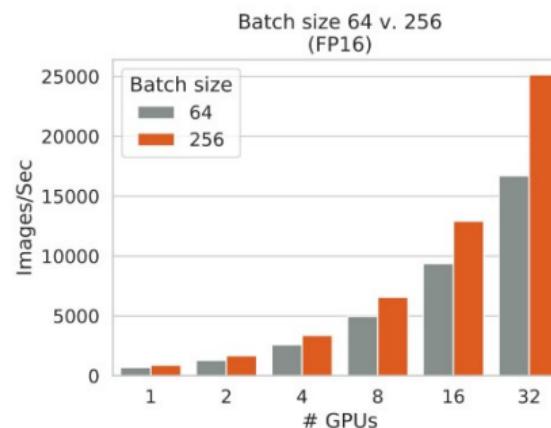
- ▶ Two stages in decentralized algorithms: **transferring** and **aggregating** data across nodes.
 - In each step, data is transferred on the network, and is sent to application to be reduced, before the result is sent again over the network.
- ▶ This process **reduces the network utilization** since the network is **not utilized** during the **reduction** at the application layer.
- ▶ Chunk (break) the data in to a few pieces, and transfer each **chunk** independently in parallel.
- ▶ While one chunk is being reduced on the CPU, another chunk can be sent over the network: this enables **pipelining** of **network transfer and application-level processing** across various chunks.



Distributed SGD and Batch Size

Batch Size vs. Number of GPUs

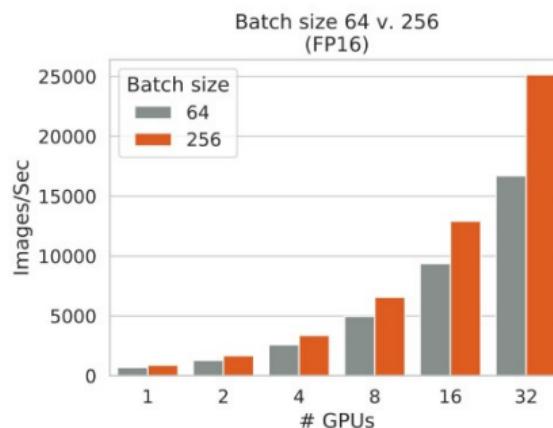
► $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$



[<https://medium.com/@emwatz/lessons-for-improving-training-performance-part-1-b5efd0f0dcea>]

Batch Size vs. Number of GPUs

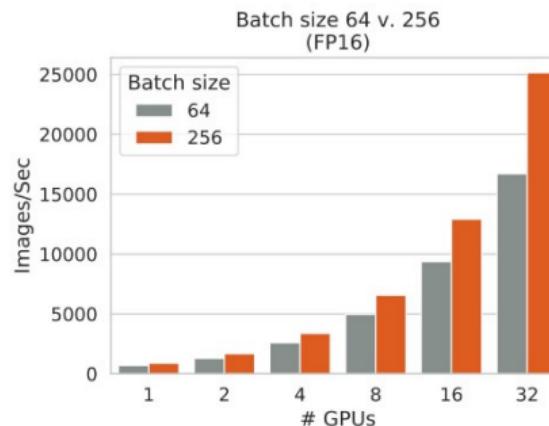
- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$
- ▶ The more samples processed during each batch, the faster a training job will complete.



[<https://medium.com/@emwatz/lessons-for-improving-training-performance-part-1-b5efd0f0dcea>]

Batch Size vs. Number of GPUs

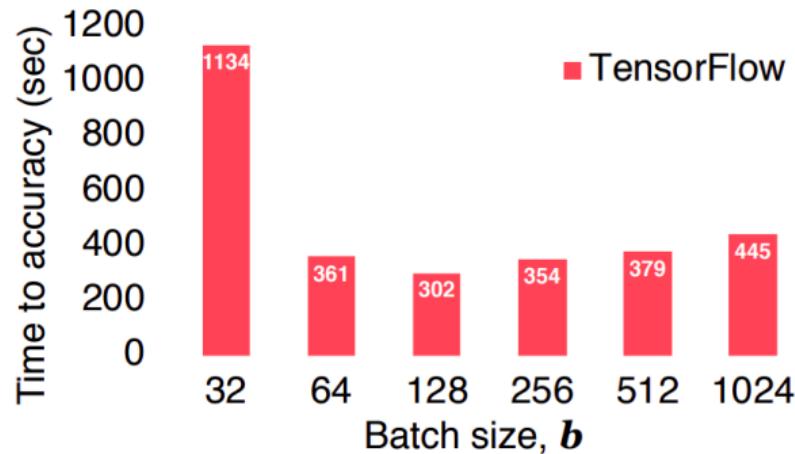
- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$
- ▶ The more samples processed during each batch, the faster a training job will complete.
- ▶ E.g., ImageNet + ResNet-50



[<https://medium.com/@emwatz/lessons-for-improving-training-performance-part-1-b5efd0f0dcea>]

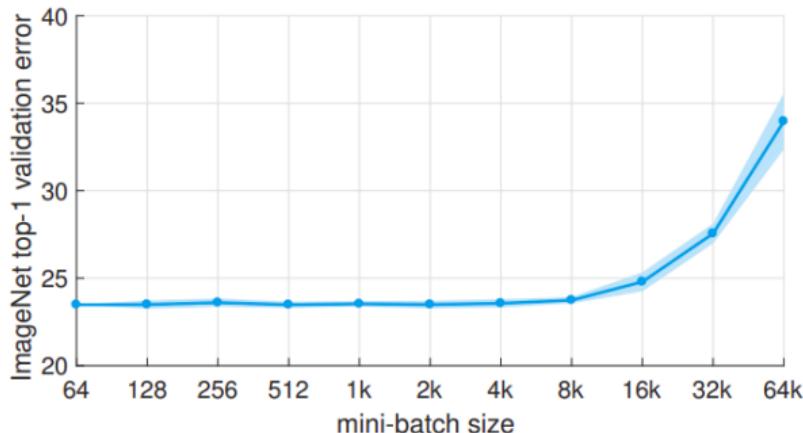
Batch Size vs. Time to Accuracy

- ▶ ResNet-32 on Titan X GPU



[Peter Pietzuch - Imperial College London]

Batch Size vs. Validation Error



[Goyal et al., Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour, 2018]



Improve the Validation Error



Improve the Validation Error

- ▶ Scaling learning rate
- ▶ Batch normalization
- ▶ Label smoothing
- ▶ Momentum



Scaling Learning Rate

► $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$.



Scaling Learning Rate

- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$.
- ▶ **Linear scaling:** multiply the **learning rate** by **k**, when the **mini batch size** is multiplied by **k**.



Scaling Learning Rate

- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$.
- ▶ **Linear scaling**: multiply the **learning rate** by **k**, when the **mini batch size** is multiplied by **k**.
- ▶ **Constant warmup**: start with a **small learning rate** for few epochs, and then increase the learning rate to **k times learning rate**.

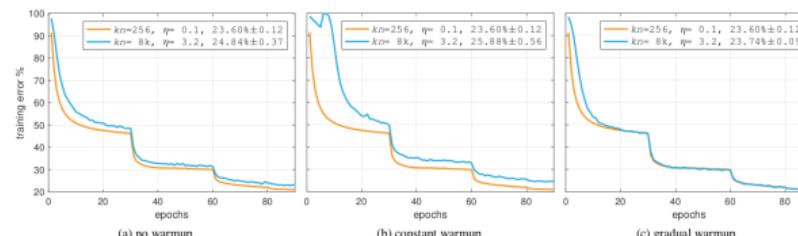


Scaling Learning Rate

- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$.
- ▶ **Linear scaling**: multiply the **learning rate** by **k**, when the **mini batch size** is multiplied by **k**.
- ▶ **Constant warmup**: start with a **small learning rate** for few epochs, and then increase the learning rate to **k times learning rate**.
- ▶ **Gradual warmup**: start with a **small learning rate**, and then gradually increase it by a constant for each epoch till it reaches **k times learning rate**.

Scaling Learning Rate

- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l(\mathbf{x}, \mathbf{w})$.
- ▶ **Linear scaling**: multiply the **learning rate** by k , when the **mini batch size** is multiplied by k .
- ▶ **Constant warmup**: start with a **small learning rate** for few epochs, and then increase the learning rate to k times learning rate.
- ▶ **Gradual warmup**: start with a **small learning rate**, and then gradually increase it by a constant for each epoch till it reaches k times learning rate.



[Goyal et al., Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour, 2018]



Batch Normalization (1/2)

- ▶ Changes in **minibatch size** change the underlying **loss function** being optimized.



Batch Normalization (1/2)

- ▶ Changes in **minibatch size** change the underlying **loss function** being optimized.
- ▶ **Batch Normalization** computes statistics along the **minibatch dimension**.



Batch Normalization (1/2)

- ▶ Changes in minibatch size change the underlying loss function being optimized.
- ▶ Batch Normalization computes statistics along the minibatch dimension.

$$\mu_\beta = \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \mathbf{x}$$

$$\sigma_\beta^2 = \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} (\mathbf{x} - \mu_\beta)^2$$



Batch Normalization (2/2)

- ▶ Zero-centering and normalizing the inputs, then scaling and shifting the result.

$$\hat{\mathbf{x}} = \frac{\mathbf{x} - \mu_\beta}{\sqrt{\sigma_\beta^2 + \epsilon}}$$
$$\mathbf{z} = \alpha \hat{\mathbf{x}} + \gamma$$

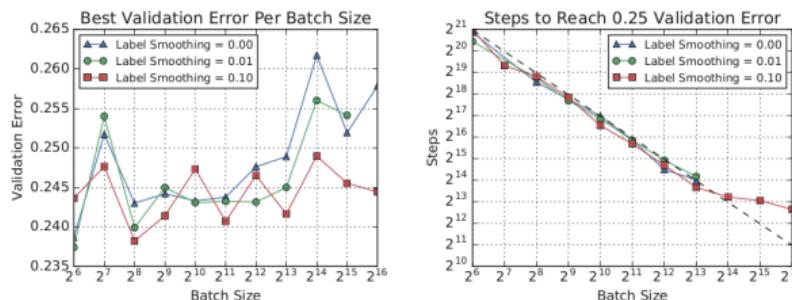
- ▶ $\hat{\mathbf{x}}$: the zero-centered and normalized input.
- ▶ \mathbf{z} : the output of the BN operation, which is a scaled and shifted version of the inputs.
- ▶ α : the scaling parameter vector for the layer.
- ▶ γ : the shifting parameter (offset) vector for the layer.
- ▶ ϵ : a tiny number to avoid division by zero.

Label Smoothing

- ▶ A generalization technique.
- ▶ Replaces one-hot encoded label vector \mathbf{y}_{hot} with a mixture of \mathbf{y}_{hot} and the uniform distribution.

$$\mathbf{y}_{\text{ls}} = (1 - \alpha)\mathbf{y}_{\text{hot}} + \alpha/\mathbf{K}$$

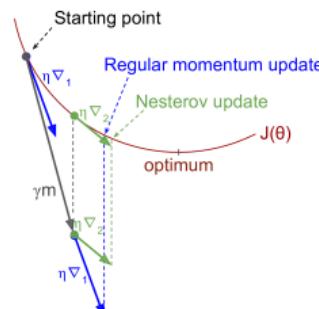
- ▶ K is the number of label classes, and α is a hyperparameter.



[Shallue et al., Measuring the Effects of Data Parallelism on Neural Network Training, 2019]

Momentum (1/3)

- ▶ Regular gradient descent optimization: $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla J(\mathbf{w})$



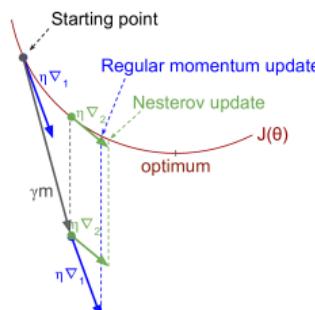
[Aurélien Géron, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2019]

Momentum (1/3)

- ▶ Regular gradient descent optimization: $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla J(\mathbf{w})$
- ▶ At each iteration, **momentum optimization** adds the **local gradient** to the **momentum vector \mathbf{m}** .

$$\mathbf{m} \leftarrow \beta \mathbf{m} + \eta \nabla J(\mathbf{w})$$

$$\mathbf{w} \leftarrow \mathbf{w} - \mathbf{m}$$

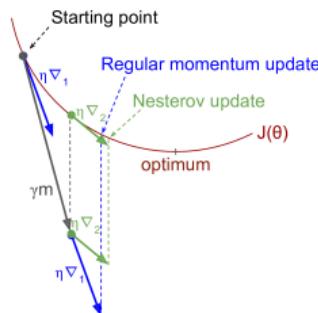


[Aurélien Géron, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2019]

Momentum (2/3)

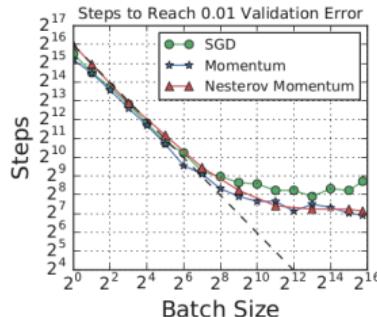
- ▶ **Nesterov momentum** measure the gradient of the cost function **slightly ahead** in the direction of the momentum.

$$\mathbf{m} = \beta \mathbf{m} + \eta \nabla J(\mathbf{w} + \beta \mathbf{m})$$
$$\mathbf{w} \leftarrow \mathbf{w} - \mathbf{m}$$



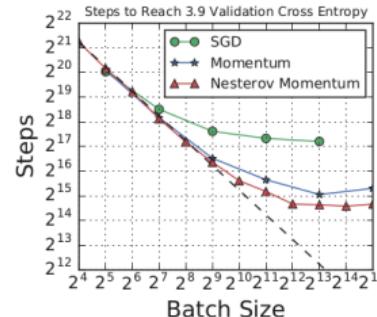
[Aurélien Géron, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2019]

Momentum (3/3)

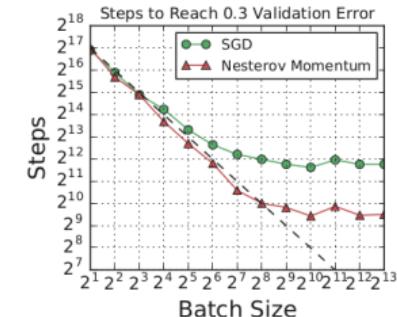


(a) Simple CNN on MNIST

[Shallue et al., Measuring the Effects of Data Parallelism on Neural Network Training, 2019]



(b) Transformer Shallow on LM1B



(c) ResNet-8 on CIFAR-10



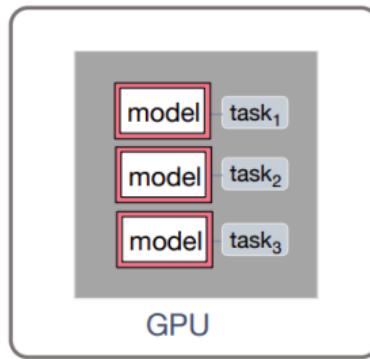
CROSSBOW: Scaling Deep Learning with Small Batch Sizes on Multi-GPU Servers



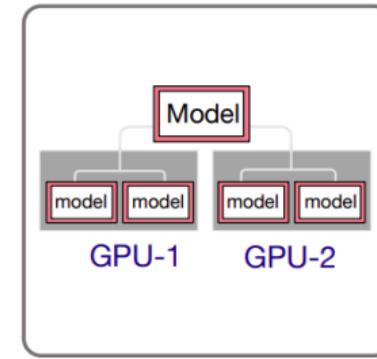
- ▶ How to design a deep learning system that scales training with multiple GPUs, even when the preferred batch size is small?

Crossbow

(1) How to increase efficiency with small batches?



(2) How to synchronise model replicas?



[Peter Pietzuch - Imperial College London]

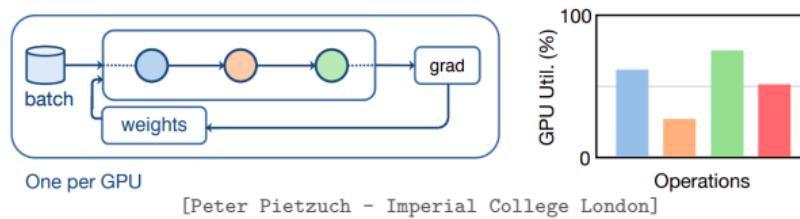


Problem: Small Batches

- ▶ Small batch sizes underutilise GPUs.

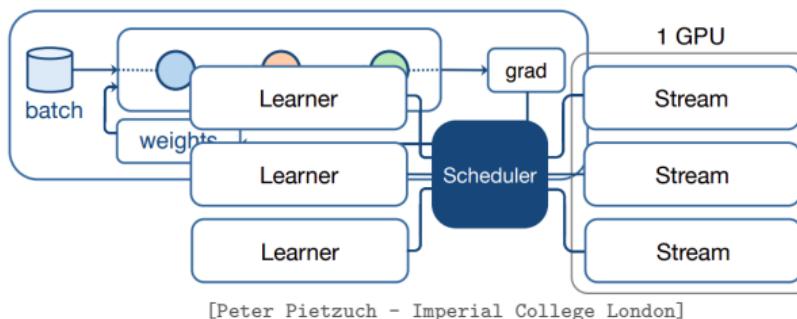
Problem: Small Batches

- ▶ **Small** batch sizes **underutilise** GPUs.
- ▶ **One batch** per GPU: **not enough data** and instruction parallelism for every operator.



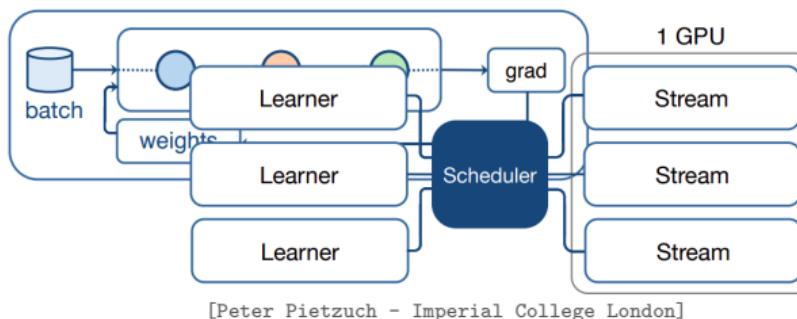
Idea: Multiple Replicas Per GPU

- ▶ Train **multiple model replicas** per GPU.
- ▶ A **learner** is an entity that trains a **single model replica** **independently** with a given batch size.



Idea: Multiple Replicas Per GPU

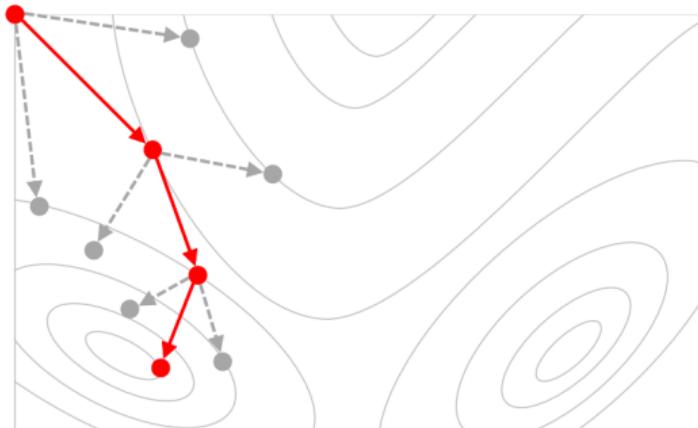
- ▶ Train **multiple model replicas** per GPU.
- ▶ A **learner** is an entity that trains a **single model replica** **independently** with a given batch size.



- ▶ But, now we must **synchronise** a **large number** of **model replicas**.

Problem: Similar Starting Point

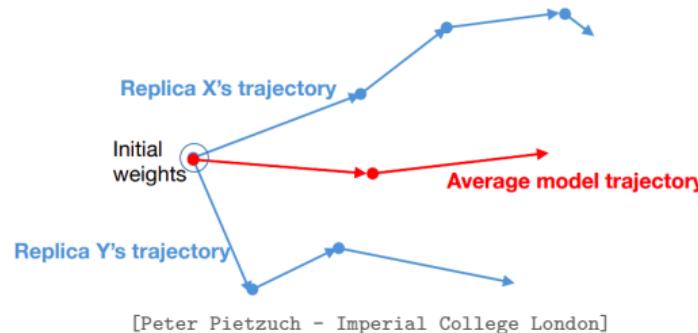
- ▶ All learners always **start** from the **same point**.
- ▶ Limited **exploration** of parameter space.



[Peter Pietzuch - Imperial College London]

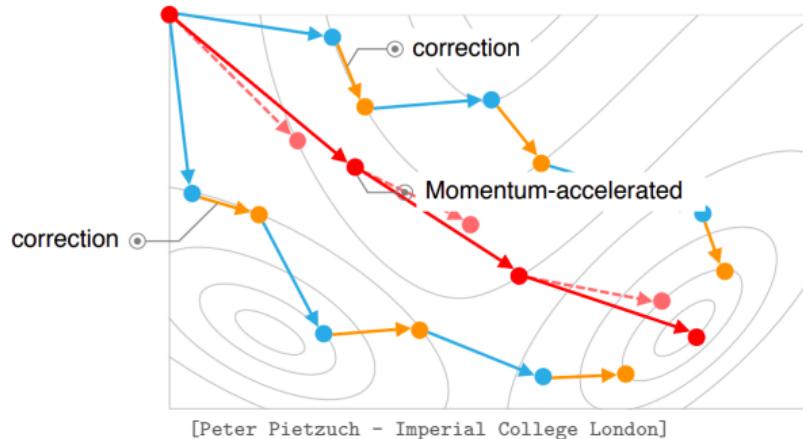
Idea: Independent Replicas

- ▶ Maintain **independent** model **replicas**.
- ▶ **Increased exploration** of space through parallelism.
- ▶ Each model replica uses **small batch size**.



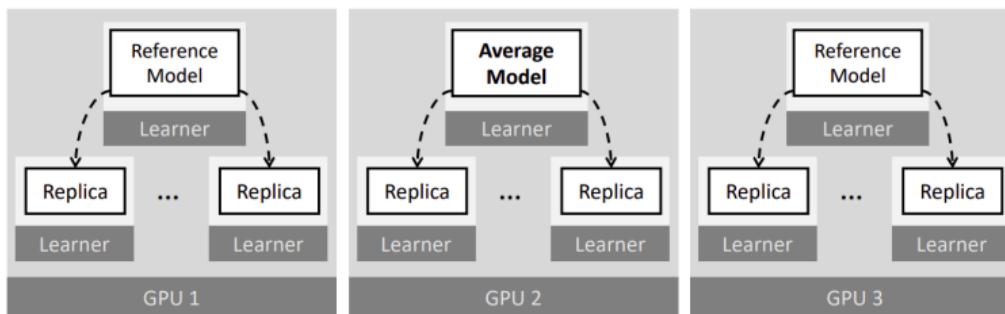
Crossbow: Synchronous Model Averaging

- ▶ Allow learners to diverge, but correct trajectories based on average model.
- ▶ Accelerate average model trajectory with **momentum** to find minima faster.



GPUs with Synchronous Model Averaging

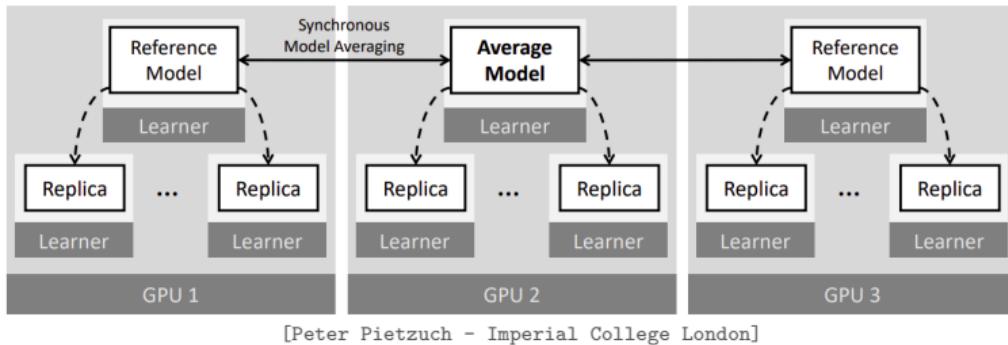
- ▶ Synchronously apply corrections to **model replicas**.



[Peter Pietzuch - Imperial College London]

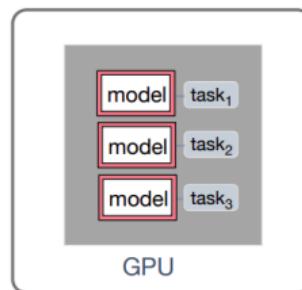
GPUs with Synchronous Model Averaging

- ▶ Ensures **consistent view** of **average model**.
- ▶ Takes **GPU bandwidth** into account during synchronisation.



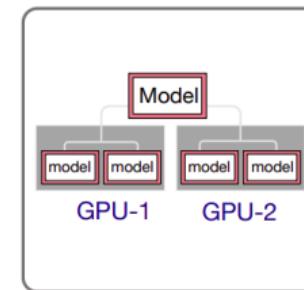
Crossbow

(1) How to increase efficiency with small batches?



Train multiple
model replicas
per GPU

**(2) How to synchronise
model replicas?**



Use synchronous
model averaging

[Peter Pietzuch - Imperial College London]



Summary



Summary

- ▶ Data-parallel
- ▶ The aggregation algorithm
- ▶ Communication synchronization
- ▶ Communication compression
- ▶ Parallelism of computations and communications
- ▶ TicTac and Caramel
- ▶ Batch Size



Reference

- ▶ Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020
- ▶ Hashemi et al., TicTac: Accelerating Distributed Deep Learning with Communication Scheduling, 2019
- ▶ Hashemi et al., CARAMEL: Accelerating Decentralized Distributed Deep Learning with Model-Aware Scheduling, 2020
- ▶ P. Goyal et al., Accurate, large minibatch sgd: Training imagenet in 1 hour, 2017
- ▶ C. Shallue et al., Measuring the effects of data parallelism on neural network training, 2018
- ▶ A. Koliousis et al. CROSSBOW: scaling deep learning with small batch sizes on multi-gpu servers, 2019



Questions?