

Background Knowledge: Distributed ML

Mark Lekina Rorat Destin Niyomufasha

Dartmouth College

February 7, 2024

- ▶ **Distributed Machine Learning** involves training machine learning models across multiple computational resources by splitting the training process across machines, processors, or nodes, utilizing *data/model parallelism*, and employing coordination mechanisms such as *parameter servers* and *all-reduce* operations for effective synchronization.

Why do we need it?

- ▶ *Handling Large-Scale Data*: For datasets too large for a single machine.
- ▶ *Speeding Up Training*: Reduces training time through parallel processing.
- ▶ *Training More Complex Models*: Enables training of larger models not fitting in single GPU/CPU memory.
- ▶ *Resource Utilization*: Efficient use of computational resources across networks or clouds.
- ▶ *Geographical Distribution*: Trains on locally stored data, addressing privacy and bandwidth issues.
- ▶ *Fault tolerance*: Offers system resilience to node failures and scalable training capabilities.

Outline

- ▶ **Model vs. data parallelism**
- ▶ Centralized vs. decentralized optimization
- ▶ Synchronous vs. asynchronous scheduling
- ▶ Communication pattern used for exchanging parameters.

Model Parallelism

- ▶ Model parallelism involves partitioning a machine learning model into smaller parts or segments that can be processed in parallel.
- ▶ Simple implementation: encoder-decoder systems common in NLP, e.g., Transformer
- ▶ *Vertical partitioning*: applying splits between neural network layers.
- ▶ *Horizontal partitioning*: the layers themselves are partitioned.

Horizontal vs. Vertical partitioning

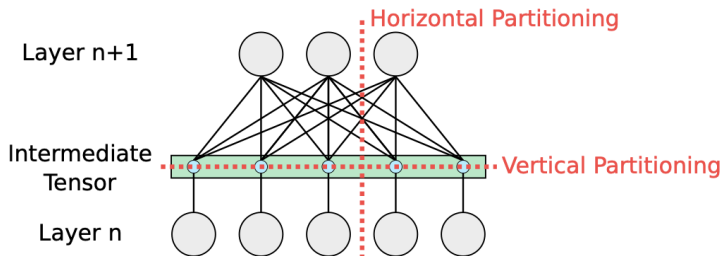


Figure 1: Horizontal vs. Vertical partitioning.¹

¹Source: <https://arxiv.org/pdf/2007.03970.pdf>

Challenges of Model Parallelism

- ▶ Requires meticulous synchronization of model segments to manage dependencies and ensure seamless data flow between partitions.
- ▶ High communication overhead, particularly for data-intensive exchanges between partitions.
- ▶ Scalability issues: as the number of partitions increases, the whole configuration needs to be restructured.

Data Parallelism

- ▶ Data parallelism involves splitting the dataset into smaller chunks and distributing these chunks across multiple processing units. Each unit then performs computations on its portion of the data, often using a copy of the same model.
- ▶ *Key idea:* the sum of per-parameter gradients computed using subsets of a mini-batch matches the per-parameter gradients for the entire input batch.

$$\frac{\partial L(x; w)}{\partial w} \approx \frac{\partial L(x^0; w)}{\partial w} + \dots + \frac{\partial L(x^n; w)}{\partial w}$$

- ▶ Most recent distributed ML systems prefer data parallelism over model parallelism due to its ease of scalability.

Outline

- ▶ Model vs. data parallelism
- ▶ **Centralized vs. decentralized optimization**
- ▶ Synchronous vs. asynchronous scheduling
- ▶ Communication pattern used for exchanging parameters.

Centralized optimization

- ▶ Involves architectures where the optimization process (*i.e.*, the updating of model parameters based on computed gradients) is coordinated through a central node or a set of central nodes known as the parameter server(s).
- ▶ *Parameter-server*: optimization instance(s) that update model parameters and send new model to worker nodes.
- ▶ *Workers*: perform backpropagation and send computed gradients to the parameter server.
- ▶ This approach ensures that all worker nodes are synchronized with the same global view of the model parameters.
- ▶ Parameter servers, however, can be a bottleneck and single point of failure.

Centralized optimization

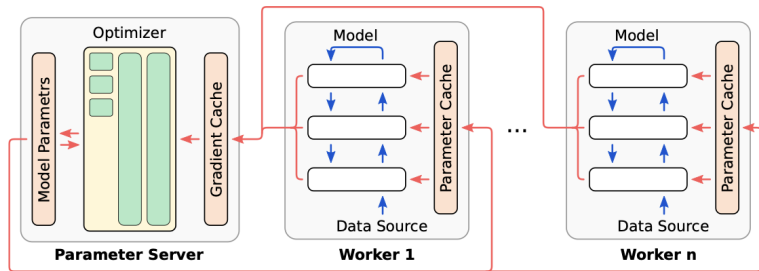


Figure 2: Workers evaluate the model to generate gradients (blue). The parameter server consumes them to update the model (red).³

³Source: <https://arxiv.org/pdf/2007.03970.pdf>

Decentralized optimization

- ▶ Refers to the approach where the optimization process is spread across multiple nodes without a central coordinator like a parameter server.
- ▶ Each node directly communicates with one or more other nodes in the network to exchange information (e.g., gradients, parameters) and update its local model.
- ▶ Significantly reduces bottlenecks associated with centralized approaches and improve fault tolerance, as there's no single point of failure.
- ▶ However, nodes may have different views of the model at a given time, which can lead to convergence issues.

Decentralized optimization

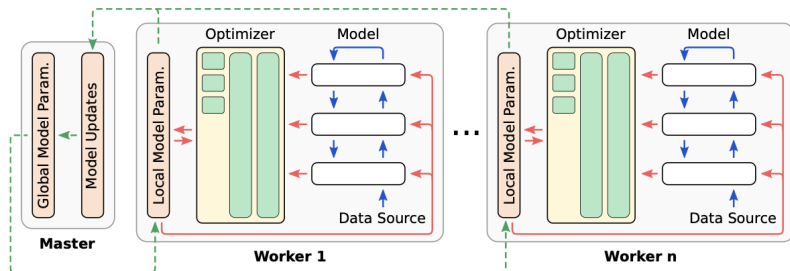


Figure 3: The master node generates the next global model by combining local model replicas (green) that were trained in isolation by the workers (red).⁵

⁵Source: <https://arxiv.org/pdf/2007.03970.pdf>

Exploration vs. exploitation phase

- ▶ Multiple independent workers concurrently try to solve a similar but not exactly the same problem.
- ▶ *Exploration phase*: workers iteratively evaluate the loss function using different mini-batches and independently update their local models.
- ▶ *Exploitation phase*: workers share their model updates with the master node, which merges the updates to distill adjustments that work better on average across the investigated portion of the training dataset, and shared the new parameters with the workers.

Exploration vs. exploitation phase

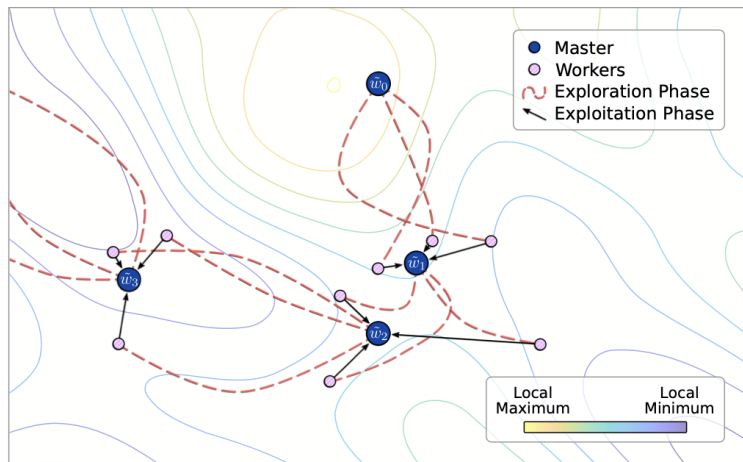


Figure 4: Workers start with the global model, then explore independently before combining their findings to refine the global model. ⁷

Outline

- ▶ Model vs. data parallelism
- ▶ Centralized vs. decentralized optimization
- ▶ **Synchronous vs. asynchronous scheduling**
- ▶ Communication pattern used for exchanging parameters.

Synchronous Scheduling

- ▶ *Simultaneous updates*: All worker nodes must complete their tasks before the system proceeds to the next step to ensure that updates are based on the same model version.
- ▶ *Consistency and stability*: Ensures consistency in model updates, as all nodes synchronize at the end of each iteration, leading to potentially more stable convergence behaviors.
- ▶ *Potential for idle time*: Can lead to inefficiencies due to resource under-utilization if some nodes are slower than others, as faster nodes must wait for stragglers to complete before proceeding.

Asynchronous Scheduling

- ▶ *Independent updates*: Nodes update the shared model as soon as they complete their tasks, without waiting for other nodes. Updates are applied as soon as they are available, potentially speeding up the overall process.
- ▶ *Efficiency and scalability*: Maximizes resource utilization and efficiency by eliminating idle time, making it well-suited for environments with heterogeneous computing resources.
- ▶ *Risk of inconsistency*: Can lead to inconsistencies in the model state, as updates may be applied based on stale parameters, potentially causing slower convergence or divergence if not properly managed.

Decentralized Asynchronous Scheduling

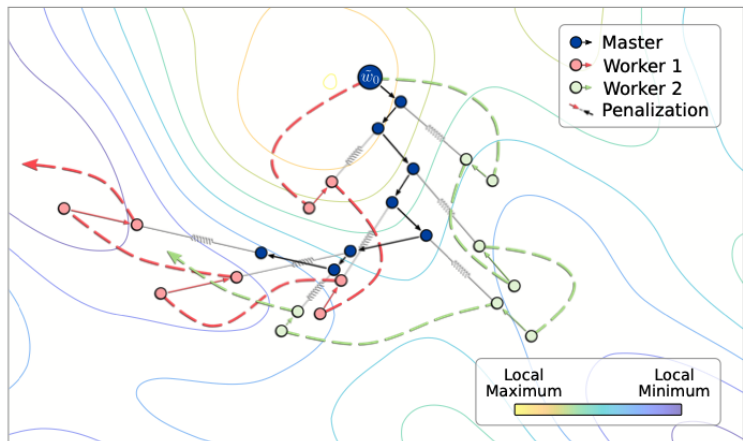


Figure 5: Workers explore locally, guided by a master node that steers them towards a shared optimum through a dance of attraction and penalties.⁹

⁹Source: <https://arxiv.org/pdf/2007.03970.pdf>

Bounded Asynchronous Scheduling

- ▶ *Delayed updates within bounds:* Allows for updates from nodes to be applied asynchronously, but within a certain bounded delay, ensuring that no update is based on information that is too outdated.
- ▶ *Balance between efficiency and consistency:* Aims to strike a balance between the efficiency of asynchronous methods and the consistency of synchronous methods, potentially reducing idle time while avoiding significant divergence in model states.
- ▶ *Complexity in implementation:* Implementing and maintaining a bounded delay can add complexity to the system, requiring mechanisms to track and manage the age of updates.

Outline

- ▶ Model vs. data parallelism
- ▶ Centralized vs. decentralized optimization
- ▶ Synchronous vs. asynchronous scheduling
- ▶ **Communication pattern used for exchanging parameters**

Communication patterns used for exchanging parameters

- ▶ Parameter Server
- ▶ All-Reduce
- ▶ Ring All-Reduce
- ▶ Gossip Protocols

Parameter Server

- ▶ As described earlier, one or more nodes hold global model parameters, whereas workers send gradients and receive updates from the parameter server.
- ▶ Scales to many workers and supports both synchronous and asynchronous updates.
- ▶ Potential bottleneck at parameter server.

All-Reduce

- ▶ Collective computation of an operation on data from all nodes. In the context of distributed ML, it aggregates gradients from all nodes to compute the global gradient.
- ▶ Ideal for synchronous data-parallel training, where gradients need global aggregation.

Ring All-Reduce

- ▶ A variant of the All-Reduce pattern where nodes form a logical ring and exchange data in a circular manner.
- ▶ Nodes in a ring exchange data with a neighbor, perform partial reductions and propagate the updates.
- ▶ Improves efficiency, since it reduces communication overhead and evenly distributes load to workers, making it ideal for large-scale training.

Ring All-Reduce

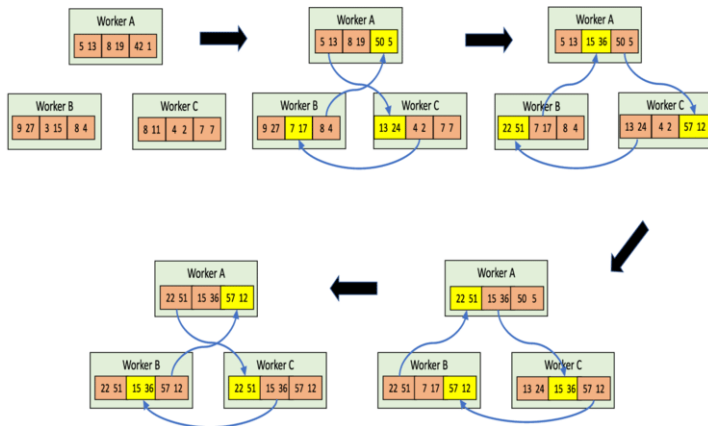


Figure 6: Nodes in a ring exchange gradients/parameters, performing on-the-fly reductions, until each has the final result.¹¹

¹¹Source: <https://doi.org/10.1007/s42514-019-00018-4>

Gossip Protocols

- ▶ Nodes randomly exchange parameters/gradients for decentralized convergence.
- ▶ Resilient to node failures, with good scalability.
- ▶ Ideal for decentralized training where strict order or central coordination is challenging.
- ▶ Convergence, however, may be slower or less predictable.