

## A NOTATIONS

Throughout the discussion of modeling the training cost, we introduce the notations in Table 1.

**Table 1: Table of notations.**

Symbol	Description
$d$	GPU device.
$\mathbf{D}$	Set of $N$ GPU devices.
$\mathbf{C}$	Tensor core of $\mathbf{D}$ .
$\mathbf{M}$	Memory limit of $\mathbf{D}$ .
$m_d$	GPU memory bandwidth of device $d$ .
$c_d$	Tensor core computation power of device $d$ .
$\mathbf{A}$	Communication matrix between devices describing the latency.
$\mathbf{B}$	Communication matrix between devices describing the bandwidth.
$\alpha_{d,d'}$	Latency between devices $d$ and $d'$ .
$\beta_{d,d'}$	Bandwidth between devices $d$ and $d'$ .
$L$	Total number of layers in the model.
$H$	Size of the hidden dimension in a transformer block.
$B_{\text{type}}$	Byte size for computational precision.
$S$	Sequence length in a transformer block.
$B_{mb}$	Micro-batch size.
$n_{mb}$	Number of micro-batches.
$D_{pp}^i$	Pipeline parallel degree for the $i$ -th pipeline.
$D_{dp}$	Data parallel degree.
$\mathbf{d}_{i,j}$	Set of GPUs serves the $j$ -th stage in the $i$ -th pipeline.
$ \mathbf{d}_{i,j} $	Tensor model parallel degree of the $j$ -th stage in the $i$ -th pipeline.
$\mathbf{D}_{dp}^k$	Set of GPUs serves the $k$ -th transformer layer in data parallelism.
$l_{i,j}$	Number of layers on the $i$ -th pipeline $j$ -th stage.
$s_{i,j}$	Starting index of transformer layers on the $i$ -th pipeline $j$ -th stage.
$\sigma$	Assignment of devices in $\mathbf{D}$ .

## B COST ESTIMATION

In this section, we model the COMM-COST, COMP-COST, and MEM-CUMSUM step by step. First we model cost for each transformer layer, and then model the end-to-end cost for the model as follows:

### Modeling Cost Layer-wisely

- **Tensor model parallel communication cost.** Suppose activation recompute is enabled. A transformer layer is running over a set of GPU  $\mathbf{d}_{i,j}$ , the tensor model parallel communication cost for a micro-batch can be estimated by:

$$\text{COMM-TP-LAYER}(\mathbf{d}_{i,j}) = 12 \cdot \max_{d \in \mathbf{d}_{i,j}} \sum_{d' \in \mathbf{d}_{i,j} - \{d\}} \left( \alpha_{d,d'} + \frac{B_{mb}SHB_{\text{type}}}{|\mathbf{d}_{i,j}| \beta_{d,d'}} \right) \quad (1)$$

- **Data parallel communication cost.** The data parallel communication cost for a transformer layer can be estimated by:

$$\text{COMM-DP-LAYER}(\mathbf{D}_{dp}^k) = 2 \cdot \max_{d \in \mathbf{D}_{dp}^k} \sum_{d' \in \mathbf{D}_{dp}^k - \{d\}} \left( \alpha_{d,d'} + \frac{12H^2 B_{\text{type}}}{|\mathbf{D}_{dp}^k| \beta_{d,d'}} \right) \quad (2)$$

- **Pipeline parallel communication cost.** Notice that pipeline parallel communication only happens when two layers are on different stages, e.g., the  $j$ -th stage and  $j+1$ -th stage in the  $i$ -th pipeline. It can be treated as two steps: the  $j$ -th stage send to the  $j+1$ -th stage in the forward pass and the  $j+1$ -th stage broadcast the information to every GPU in this stage (Or in the backward pass, the  $j$ -th stage recv from the  $j+1$ -th stage and then broadcast). Define  $\text{COMM-PP-HOP}(\mathbf{d}_{i,D_{pp}^i}, \mathbf{d}_{i,D_{pp}^i+1}) = 0$  for convenience. The communication cost for a micro-batch can be estimated by:

$$\begin{aligned} \text{COMM-PP-HOP}(\mathbf{d}_{i,j}, \mathbf{d}_{i,j+1}) = & 2 \cdot \min_{d \in \mathbf{d}_{i,j}, d' \in \mathbf{d}_{i,j+1}} \left( \alpha_{d,d'} + \frac{B_{mb}SHB_{\text{type}}}{\beta_{d,d'}} \right) \\ & + \sum_{d'' \in \mathbf{d}_{i,j+1} - d'} \left( \alpha_{d',d''} + \frac{B_{mb}SHB_{\text{type}}}{|\mathbf{d}_{i,j+1}| \beta_{d',d''}} \right) \end{aligned} \quad (3)$$

- **Computation cost for a micro-batch.** Assume tensor model parallelism is always running over the same type of device  $d$  for any layer, and activation recomputation is enabled. The computation cost can be estimated by:

$$\text{COMP-TP-LAYER}(\mathbf{d}_{i,j}) = \frac{96B_{mb} \cdot SH^2 \left(1 + \frac{S}{6H}\right)}{c_d |\mathbf{d}_{i,j}|} \quad (4)$$

### Modeling Cost for Each Parallel Strategy

- **Data parallelism cost.** Different pipeline stages synchronize gradient simultaneously. The Data parallelism cost is bounded by the slowest pipeline stage, which can be estimated as:

$$\text{COMM-DP} = \max_{i,j} \left[ \sum_{k=s_{i,j}}^{l_{i,j}} \text{COMM-DP-LAYER}(\mathbf{D}_{dp}^k) \right] \quad (5)$$

- **Pipeline and tensor model parallelism cost.** For  $i$ -th pipeline to execute, the cost consists of the computation and communication cost for each stage (indexed by  $j$ ):

$$\begin{aligned} \text{STAGE}(\mathbf{d}_{i,j}) = \\ \sum_{k=1}^{l_{i,j}} [\text{COMP-TP-LAYER}(\mathbf{d}_{i,j}) + \text{COMM-TP-LAYER}(\mathbf{d}_{i,j})] \end{aligned} \quad (6)$$

Notice that the slowest stage bounds the pipeline parallel stage. Thus, we formulate the pipeline and tensor model parallelism cost as below:

$$\begin{aligned} \text{PIPELINE-TIME}(i) = \\ \sum_{j=1}^{D_{pp}^i} (\text{STAGE}(\mathbf{d}_{i,j}) + \text{COMM-PP-HOP}(\mathbf{d}_{i,j}, \mathbf{d}_{i,j+1})) \\ + (n_{mb} - 1) \cdot \max_{j=2, \dots, D_{pp}^i} (\text{STAGE}(\mathbf{d}_{i,j}) \\ + \text{COMM-PP-HOP}(\mathbf{d}_{i,j}, \mathbf{d}_{i,j+1})) \end{aligned} \quad (7)$$

**Modeling End-to-end time:** One iteration time is determined by the slowest pipeline and the data parallel cost, which can be estimated as follows:

$$\begin{aligned} \text{COMM-COST}(\sigma) + \text{COMP-COST}(\sigma) \\ = \max_{i=1, \dots, D_{dp}} \text{PIPELINE-TIME}(i) + \text{COMM-DP} \end{aligned} \quad (8)$$

**Modeling Memory Cost:** Suppose full activation recompute and naive data parallelism are applied. The memory cost of parameters and activations can be estimated as below:

$$\text{MEM-CUMSUM}(\sigma) = \frac{48H^2 B_{type}}{|\mathbf{d}_{i,j}|} + B_{mb} SH B_{type} \quad (9)$$