

Enabling Parallelism Hot Switching for Efficient Training of Large Language Models

Hao Ge*, Fangcheng Fu*, Haoyang Li, Xuanyu Wang, Sheng Lin, Yujie Wang, Xiaonan Nie, Hailin Zhang, Xupeng Miao(Purdue University), Bin Cui

PKU, Purdue University

December, 12th

- ① Background
- ② HotSPa
- ③ Experiments
- ④ Conclusions

① Background

② HotSPa

③ Experiments

④ Conclusions

Background-Transformers and Large Language Models

- Padding and Packing
 - In order to train sequences with different lengths together, techniques such as padding or packing are needed to preprocess the sequences

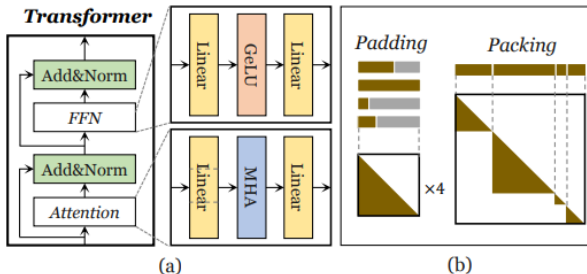


Figure 1. Illustration of (a) the architecture of Transformer layer; (b) sequence padding and packing.

Background- Deep Learning Training

- Computation Graph
- Mixed-precision Training
- Gradient Accumulation

Seq IDs and Lengths

①: 0.5K	⑧: 4.5K
②: 28K	⑨: 3.2K
③: 3.6K	⑩: 0.2K
④: 13.5K	⑪: 26K
⑤: 14.5K	⑫: 6.6K
⑥: 5K	⑬: 2K
⑦: 1.5K	⑭: 8K
⑧: 4K	⑮: 5K



Packed Seq and Lengths

packed(①⑦): 32K
packed(⑪⑤⑩): 31.5K
packed(④③②⑩): 31.8K
packed(⑭⑫⑮⑧⑨⑬⑥): 30.8K



Train with 4 grad acc steps

Figure 2. An example of training with gradient accumulation and sequence packing when the maximal supported length is 32K.

Background-Gradient Accumulation

- Data Parallelism (DP)
- Model Parallelism (MP)
- Sequence Parallelism (SP)
- Hybrid Parallelism

Data Parallelism (DP)

- Introduction
 - Each device maintains a full copy of model states execute locally and synchronize the model gradients globally.
- Need an extra round of all-gather communication
 - For example, when there are p micro-batches, the communication cost of ZeRO-3(FSDP) would be $p + \frac{1}{2}$ times of that of conventional DP.

Model Parallelism (MP)

- Tensor parallelism (TP)
 - Two parameters within one self-attention or FFN block are split in the column- and row-wise.
 - 4 all-reduce communication for one Transformer layer.
- Pipeline parallelism (PP)
 - A model is reckoned as a sequence of layers, and divided into multiple stages across devices.
 - P2P communication operations are needed to transfer the intermediate results.

Sequence Parallelism (SP)

A special form of DP

- SP splits the training samples (sequences) in the sequence dimension

Hybrid Parallelism

Combine DP, TP, and PP and researchers can tune the parallelism degree of each strategy for a given task to achieve better efficiency.

Motivation

- Two popular datasets for LLMs, CommonCrawl and Github, show significant skewness in the sequence lengths.
- This skewness leads to inefficiencies when applying static parallelism strategies in training.

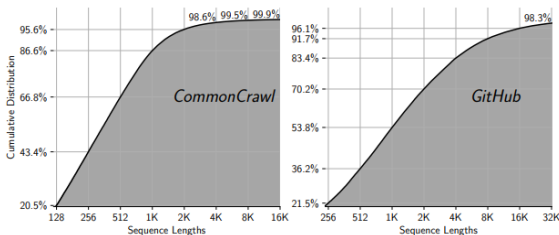


Figure 3. Cumulative distributions of sequence lengths.

① Background

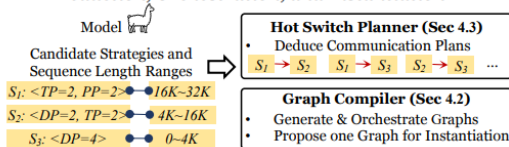
② HotSPa

③ Experiments

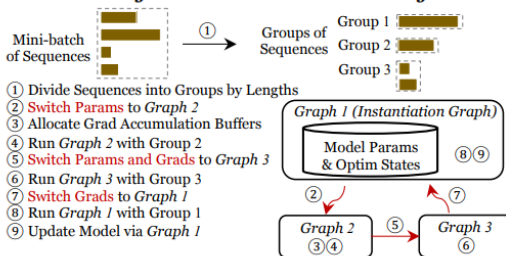
④ Conclusions

Overview

Deduction, Orchestration, and Instantiation



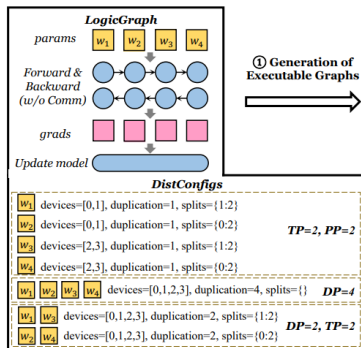
Training with Parallelism Hot Switching



Graph Compilation

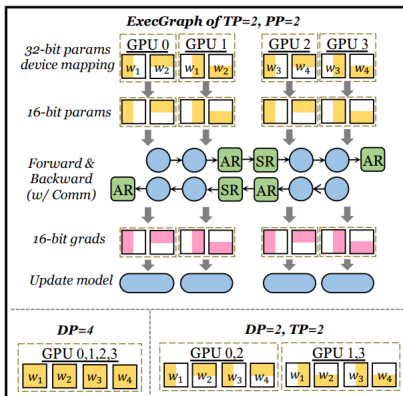
- Generation of Logical Graph
- Generation of Executable Graphs
- Orchestration of Executable Graphs

Logical Graph



Generate *DistConfigs* based on candidate parallelism strategies, the assigned devices, the number of duplications, and a map to indicate how a multi-dimensional parameter is split.

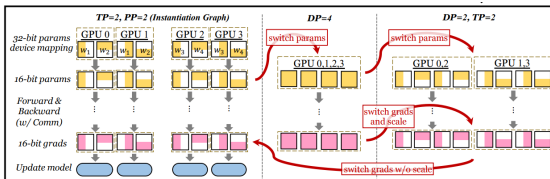
Generation of Executable Graphs



Generate *ExecGraphs* considering three types of insertions:

- For model parameters, the type casting operators are inserted.
- Inside the forward and backward propagation, communication operators.
- For model gradients, accumulation operators are inserted.

Orchestration of Executable Graphs



- Firstly, we need to identify the candidate ExecGraphs that minimize the memory occupation.
- Secondly, the ExecGraphs are re-ordered for better efficiency.
- Thirdly, prune the type casting and model update operators of the other ExecGraphs like cast model parameter to the 16-bit version.

Model Hot Switching

Two key characteristics:

- Intra-node communication is preferable than inter-node: In typical GPU clusters, GPUs within a node are connected by NVLink, which has a higher communication bandwidth than Infini-band or Ethernet.
- GPU connectivities are full-duplex: it is reasonable to minimize the maximum sending volume of all devices.

Algorithm 1: Our heuristic hot switching planner.

```

1 Initialize hot switching plan  $\mathcal{P} = \{\}$ ;
2 Initialize intra- and inter-node communication
   volume  $V_i^{(inter)}, V_i^{(intra)}$  as 0 for each device  $i$ ;
3 foreach model parameter/gradient slice do
4   Determine the owner (source) devices  $S$ ;
5   Determine the target (destination) devices  $D$ ;
6   foreach  $dst$  in  $D$  do
7     if  $dst \notin S$  then
8       Partition  $S$  into  $S^{(intra)}, S^{(inter)}$ ;
9       if  $S_i^{(intra)}$  is not empty then
10         $src \leftarrow \arg \min_i \{V_i^{(intra)} | i \in S^{(intra)}\}$ ;
11         $V_{src}^{(intra)} \leftarrow V_{src}^{(intra)} + \text{sizeof}(\text{slice})$ ;
12      else
13         $src \leftarrow \arg \min_i \{V_i^{(inter)} | i \in S^{(inter)}\}$ ;
14         $V_{src}^{(inter)} \leftarrow V_{src}^{(inter)} + \text{sizeof}(\text{slice})$ ;
15         $\mathcal{P} \leftarrow \mathcal{P} \cup (\text{slice}, src, dst)$ ;
16 return  $\mathcal{P}$ ;

```

① Background

② HotSPa

③ Experiments

④ Conclusions

Case Studies

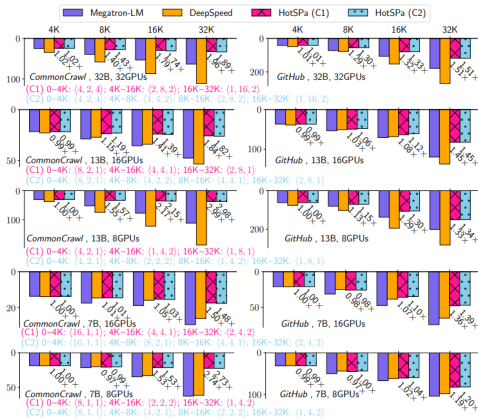


Figure 7. End-to-end evaluation (measured in seconds per mini-batch). We present two hot switching combinations for each experiments and the corresponding speedups compared with baselines (format: "seq_len_range: (DP, TP, PP)").

Case Studies

Table 3. Case studies of training with different combinations of parallelism strategies (32B, 32GPUs, $s=32K$). We dissect the running time for different groups of sequence lengths. "Others" includes the time cost of hot switching in one mini-batch.

Breakdown	CommonCrawl (time in seconds)				GitHub (time in seconds)			
	Static	C1	C2	C3	Static	C1	C2	C3
0~1K	42.1	22.7 (2.40×)	22.7 (2.40×)	12.8 (3.28×)	34.9	26.3 (2.78×)	26.3 (2.78×)	8.6 (4.05×)
1K~4K	12.5			8.2 (1.52×)	38.4			17.1 (2.24×)
4K~8K	3.5			2.7 (1.29×)	25.0			15.4 (1.62×)
8K~16K	2.9	5.4 (1.18×)	2.2 (1.31×)	2.2 (1.31×)	27.3	37.4 (1.40×)	20.2 (1.35×)	20.2 (1.35×)
16K~32K	2.4		2.4 (1.00×)	2.4 (1.00×)	53.0		53.0 (1.00×)	53.0 (1.00×)
Others	-	1.8	3.5	4.9	-	1.8	3.5	4.9
Total	63.4	32.3 (1.96×)	33.5 (1.89×)	33.2 (1.90×)	178.6	118.5 (1.50×)	118.4 (1.50×)	119.2 (1.49×)

(Static) 0~32K: (1, 16, 2) (C1) 0~4K: (4, 2, 4); 4K~16K: (2, 8, 2); 16K~32K: (1, 16, 2)

(C2) 0~4K: (4, 2, 4); 4K~8K: (4, 4, 2); 8K~16K: (2, 8, 2); 16K~32K: (1, 16, 2)

(C3) 0~1K: (8, 4, 1); 1K~4K: (4, 2, 4); 4K~8K: (4, 4, 2); 8K~16K: (2, 8, 2); 16K~32K: (1, 16, 2)

Model Hot Switching

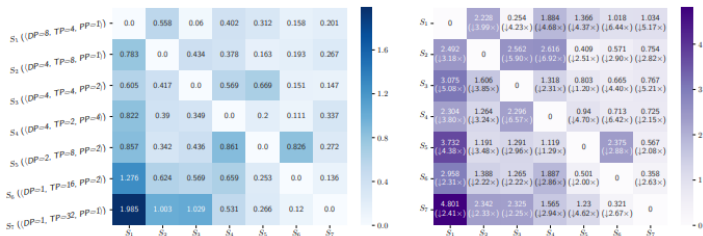


Figure 8. Time cost (in seconds) of switching between parallelism strategies (LLaMA2-32B, 32GPUs). The value in the i -th row and j -th column indicates switching from S_i to S_j . (Left: time cost of our work. Right: time cost without the optimizations in Section 4.3.)

- 1 Background
- 2 HotSPa
- 3 Experiments
- 4 Conclusions**

Future Work

Limitations:

- Reliance on user-provided combinations of strategies
- When the number of gradient accumulation steps is large, the hot switching overhead can be amortize
- Many parallelism strategies that are not incorporated

Thanks!