

# WLB-LLM: Workload-Balanced 4D Parallelism for Large Language Model Training

**Zheng Wang**, Anna Cai, Xinfeng Xie, Zaifeng Pan, Yue Guan, Weiwei Chu,  
Jie Wang, Shikai Li, Jianyu Huang, Chris Cai, Yuchen Hao, Yufei Ding

**Presenter: Yi Liu**

UC San Diego

 **Meta**

 **usenix**  
THE ADVANCED  
COMPUTING SYSTEMS  
ASSOCIATION

 **OSDI 25**

# Author

## Research Direction:

- high-performance system design
- end-to-end optimization for deep learning

## Selected Publications:

- GMI-DRL(ATC'25)
- WLB-LLM(OSDI'25)
- FastTree(MLSys'25)



Zheng Wang

<https://zhengwang.info/>

# Outline

- Background
- Design
- Evaluation
- Thinking

# Background: Training LLM is costly

The scale of LLMs grow larger ➡ Training is costly

## Estimated training cost and compute of select AI models

Source: Epoch AI, 2024 | Chart: 2025 AI Index report

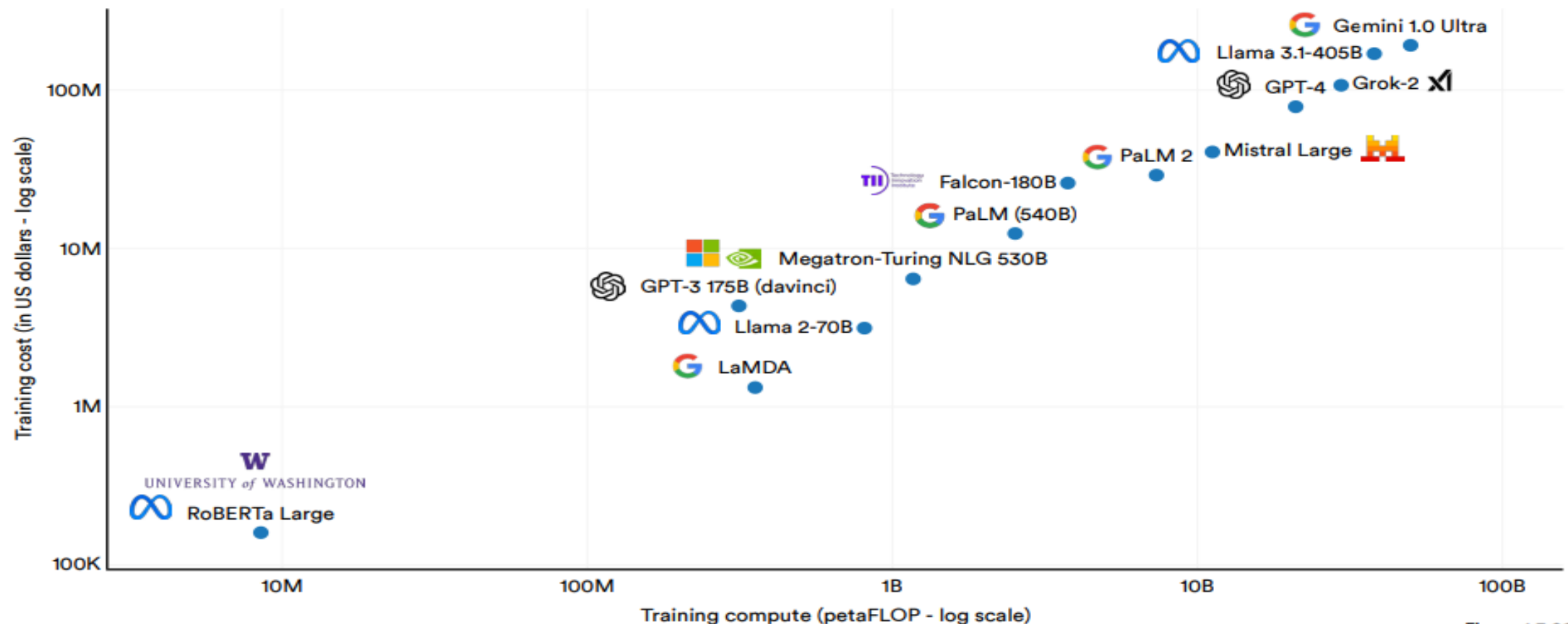
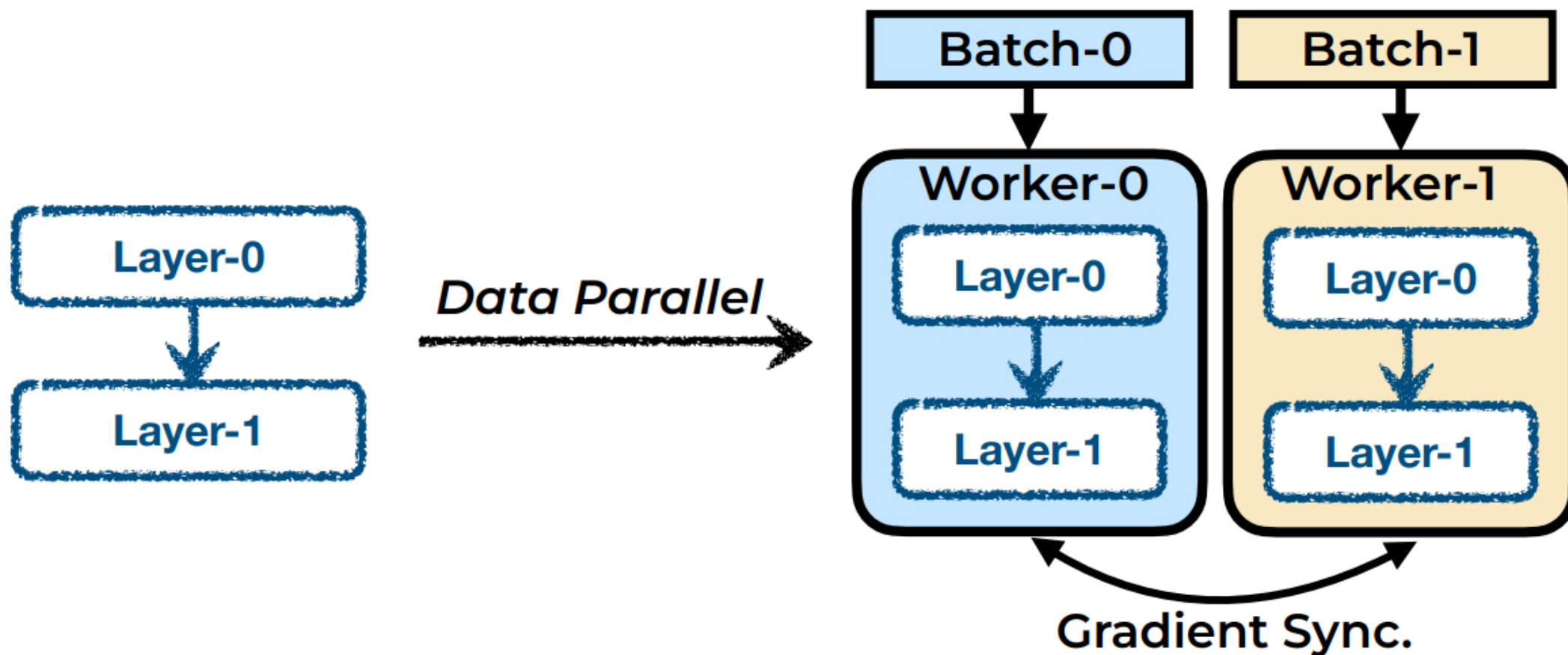


Figure 1.3.26

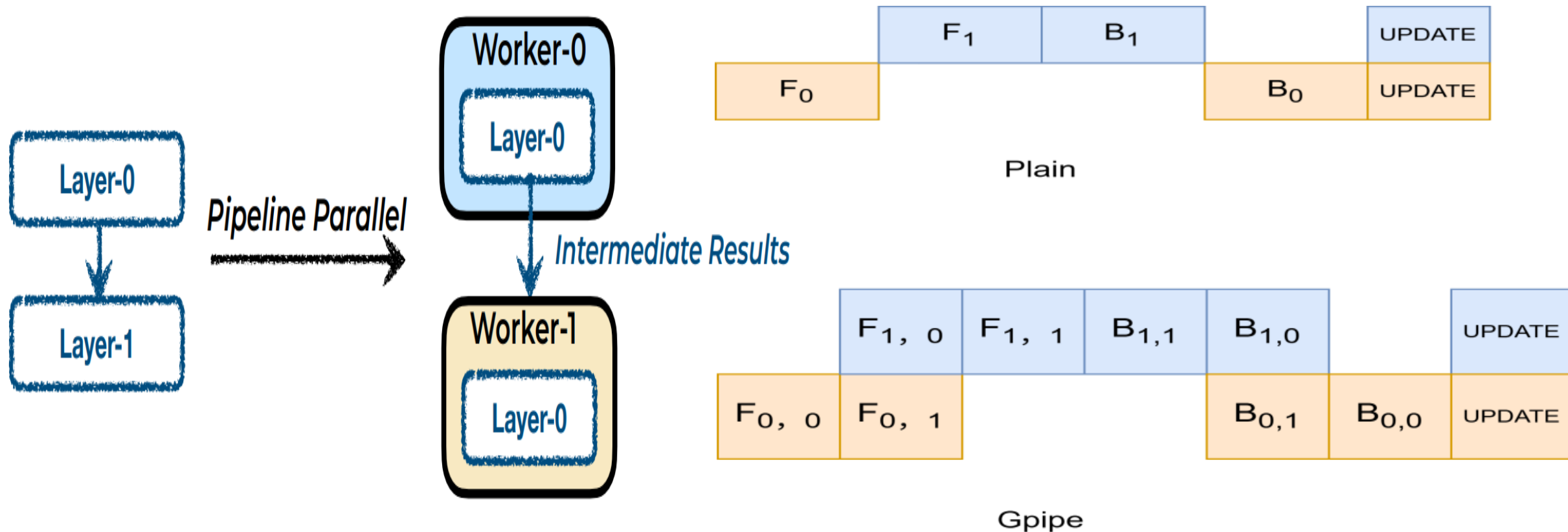
# Background: Data Parallelism

Duplicate models to different workers



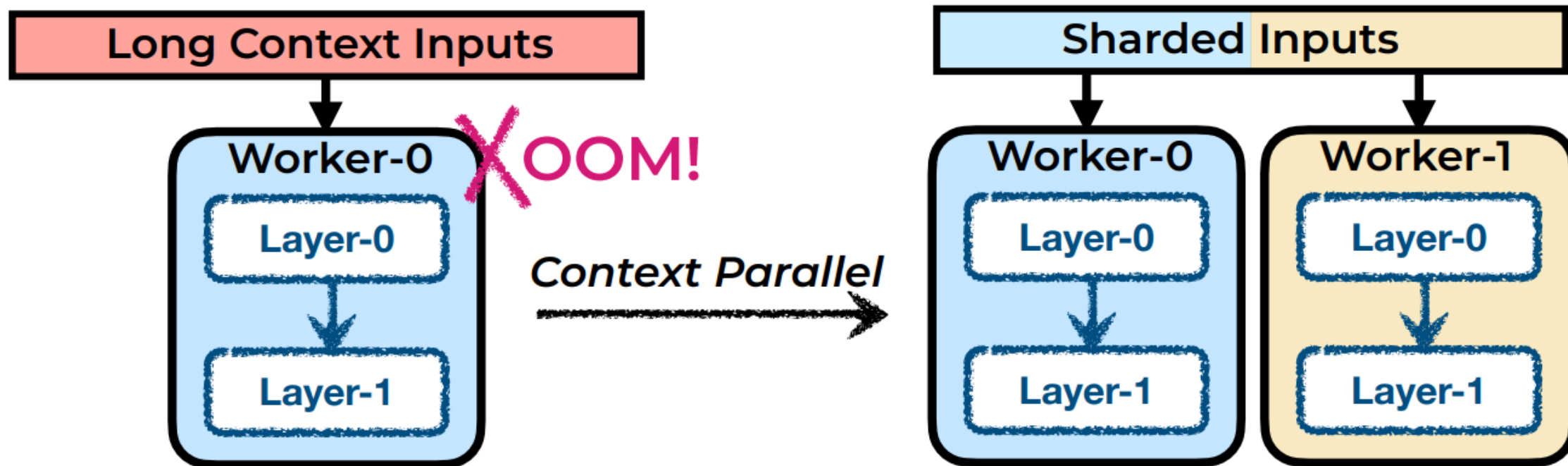
# Background: Pipeline Parallelism

Split model in a layer-wise manner



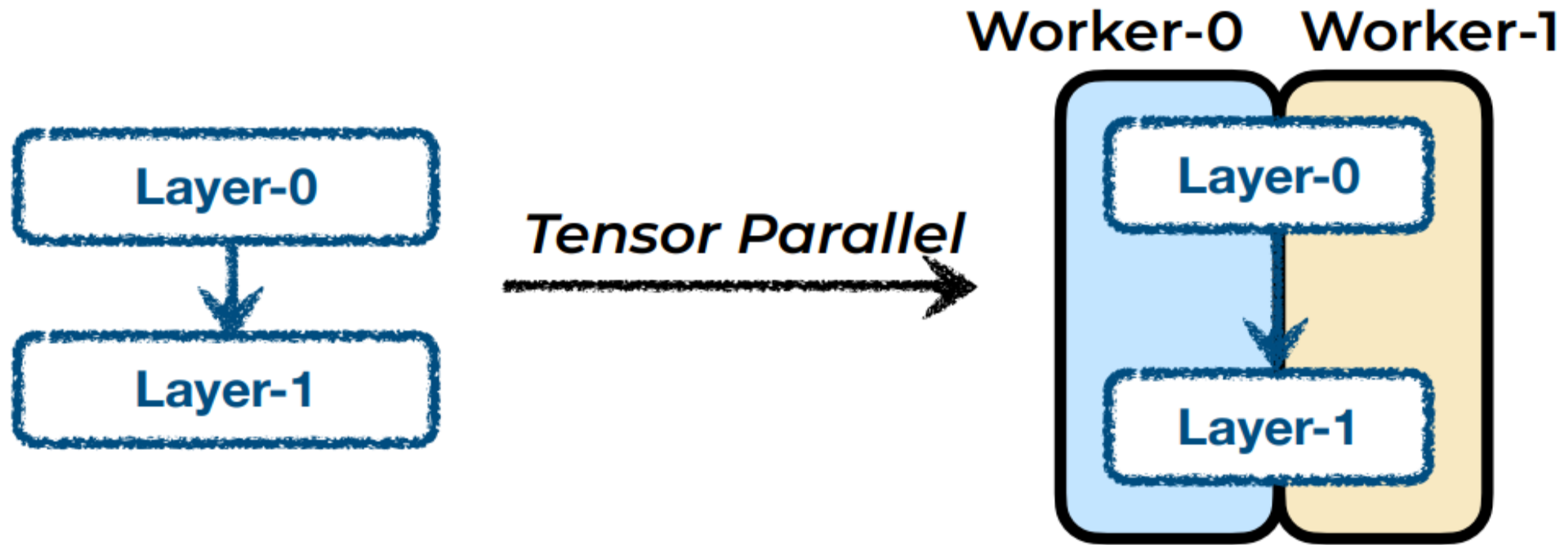
# Background: Context Parallelism

Shard the long-context input along sequence length dimension



# Background: Tensor Parallelism

Split model parameters within a layer



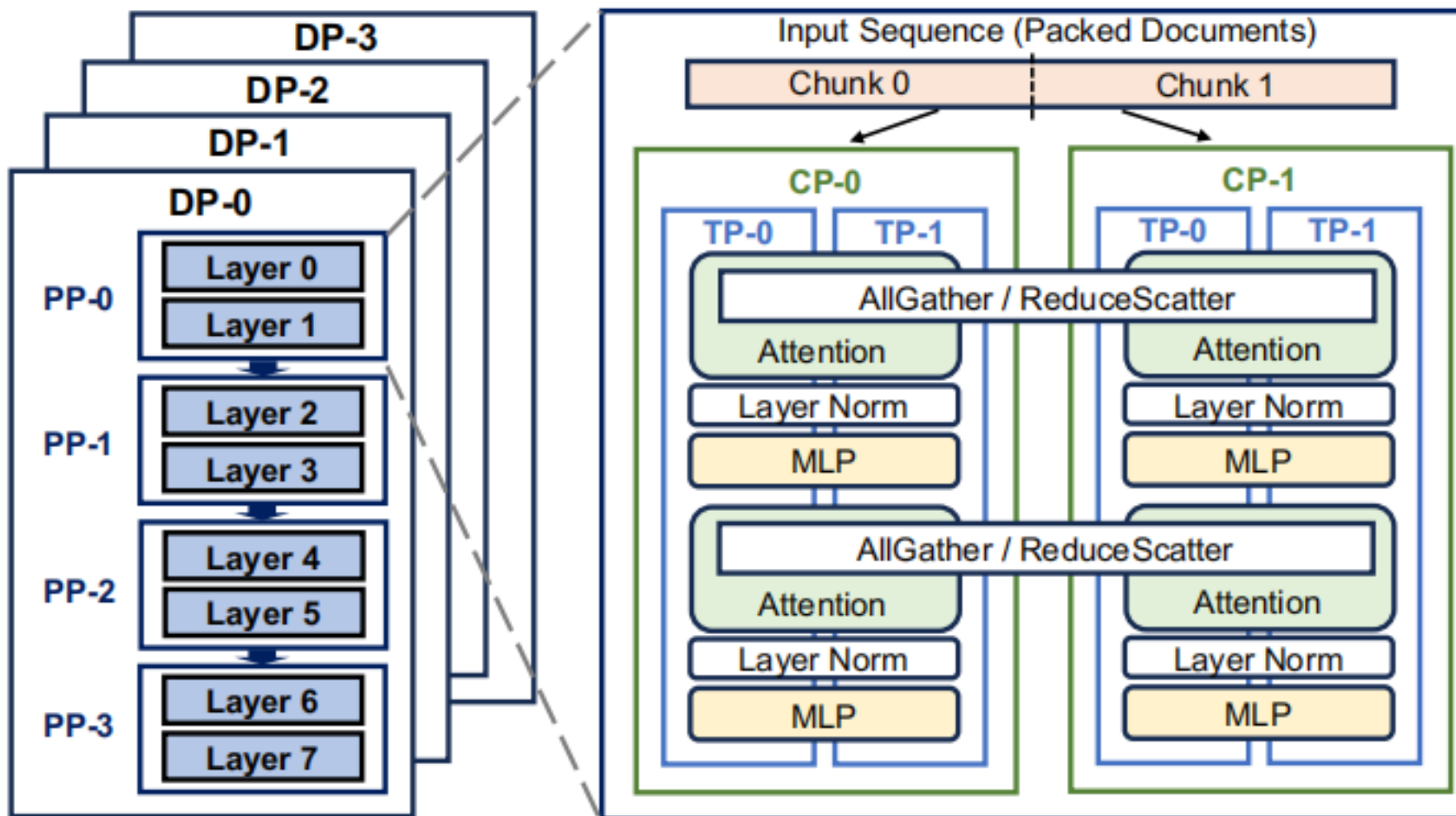


# Background: 4D parallelism

Training is costly

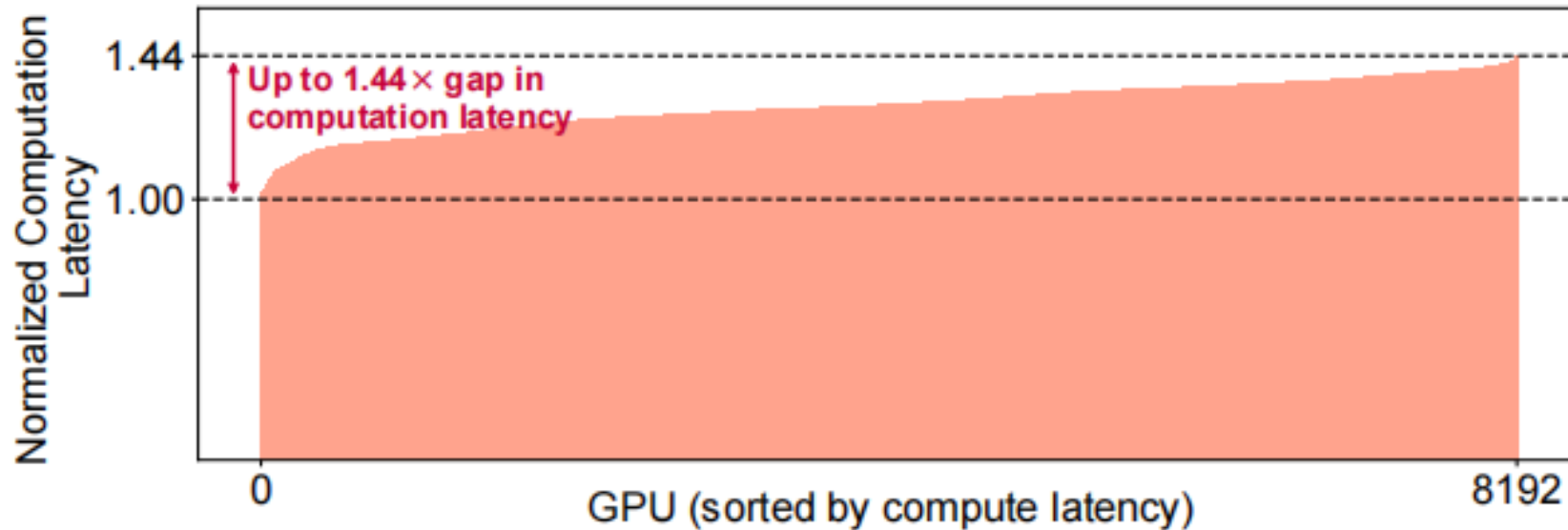


multi level parallelism



# Background: workload imbalance

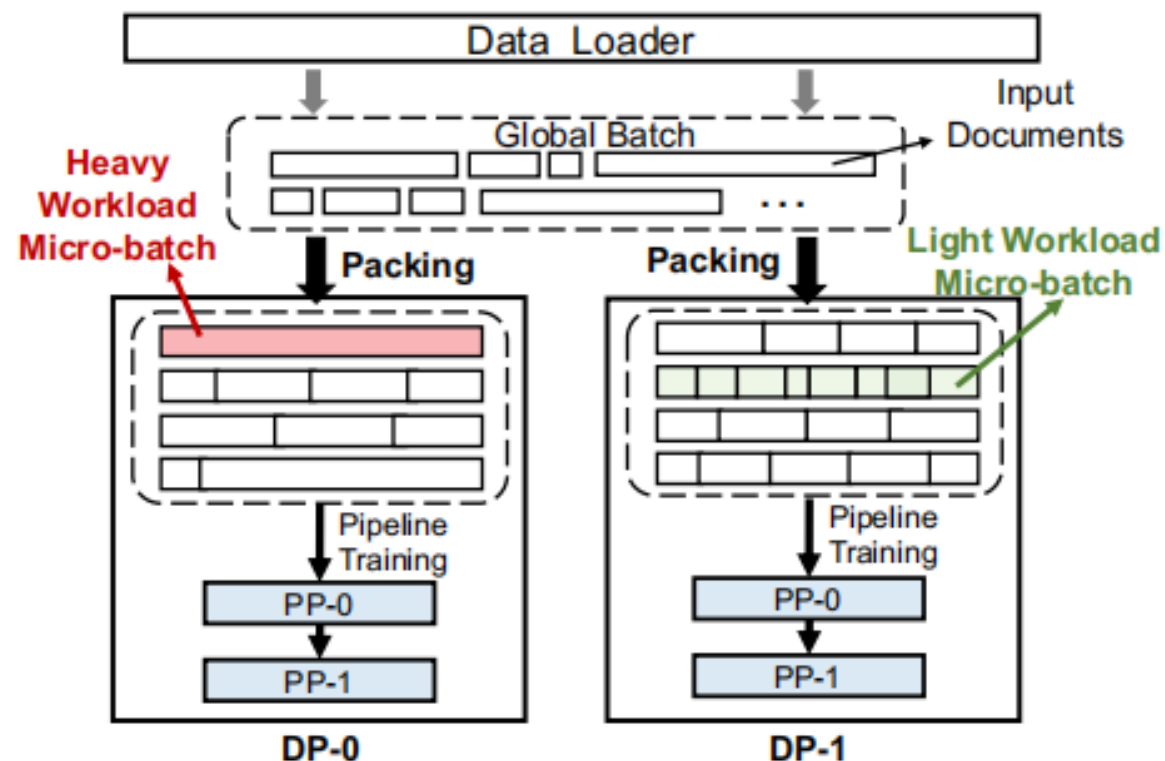
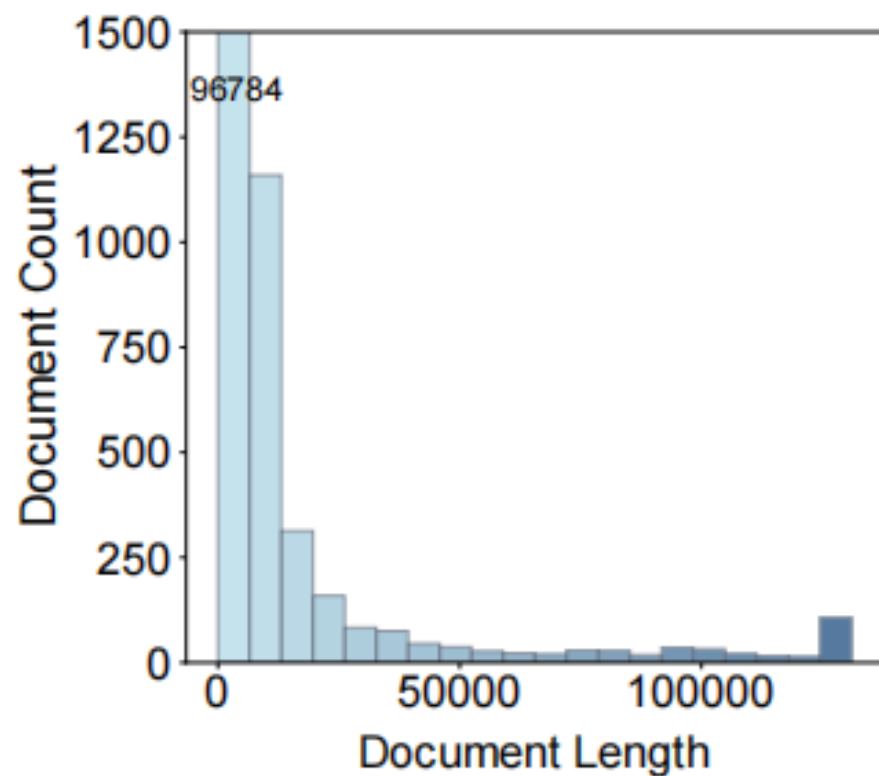
not fully utilizing GPUs → 1. the causes of workload imbalance  
2. how to resolve the imbalance



(a) Normalized computation latency in an 8K-GPU LLM training job.

# Background: PP level imbalance

varying input document length → imbalance in parallelism



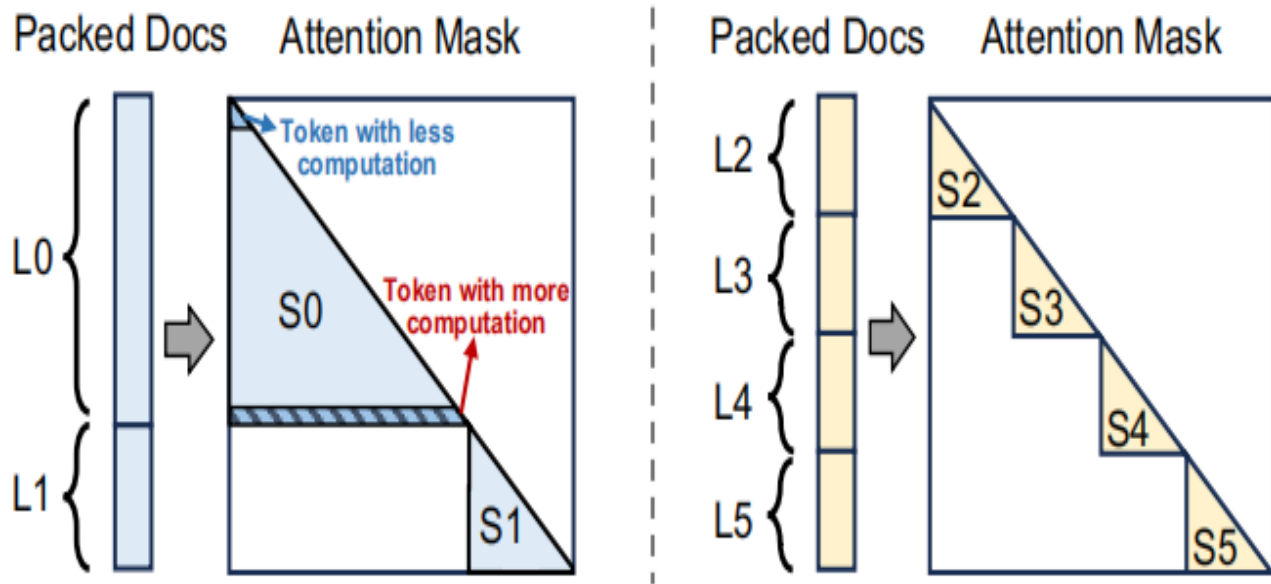
(1) Imbalance from PP-level Document Packing

# Background: reason of imbalance

varing computation



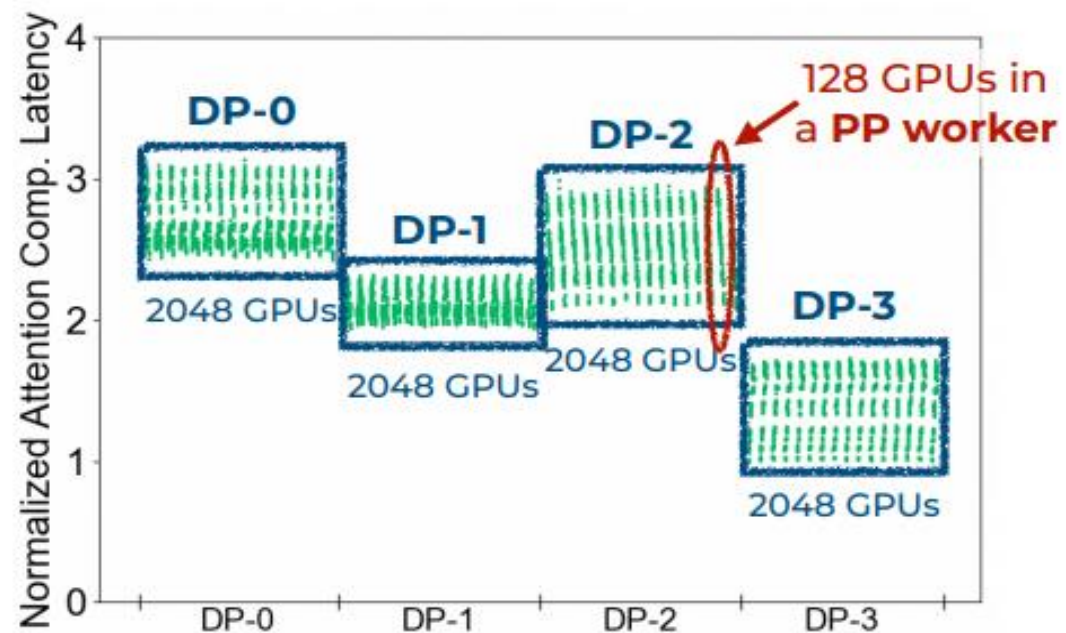
imbalance across PP workers



Document lengths:  $L_0 + L_1 = L_2 + L_3 + L_4 + L_5$

Computation (triangle areas):  $S_0 + S_1 \gg S_2 + S_3 + S_4 + S_5$

**Normalized Attention Computation Latency**  
(8K GPU: DP-4, PP-16, CP-16, TP-8)

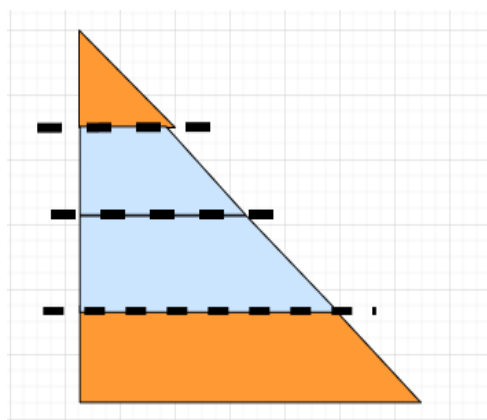
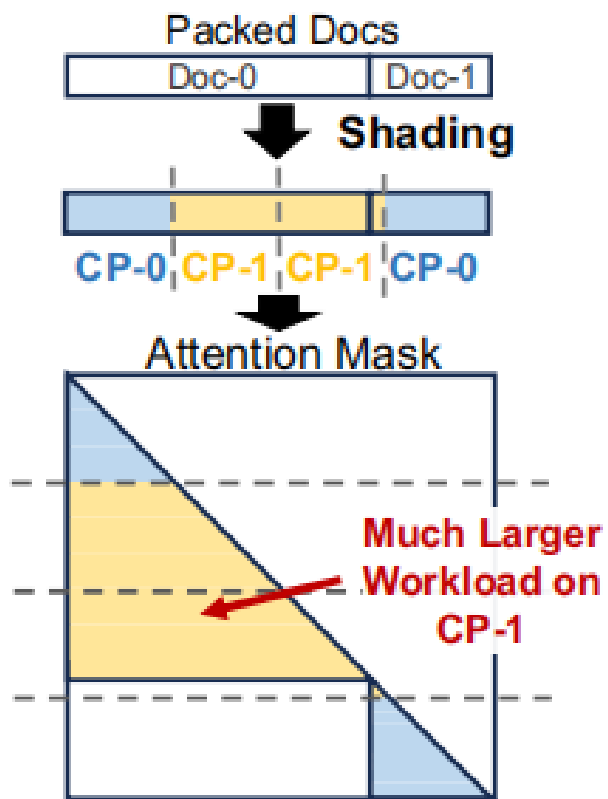


# Background:CP level imbalance

sequence sharding

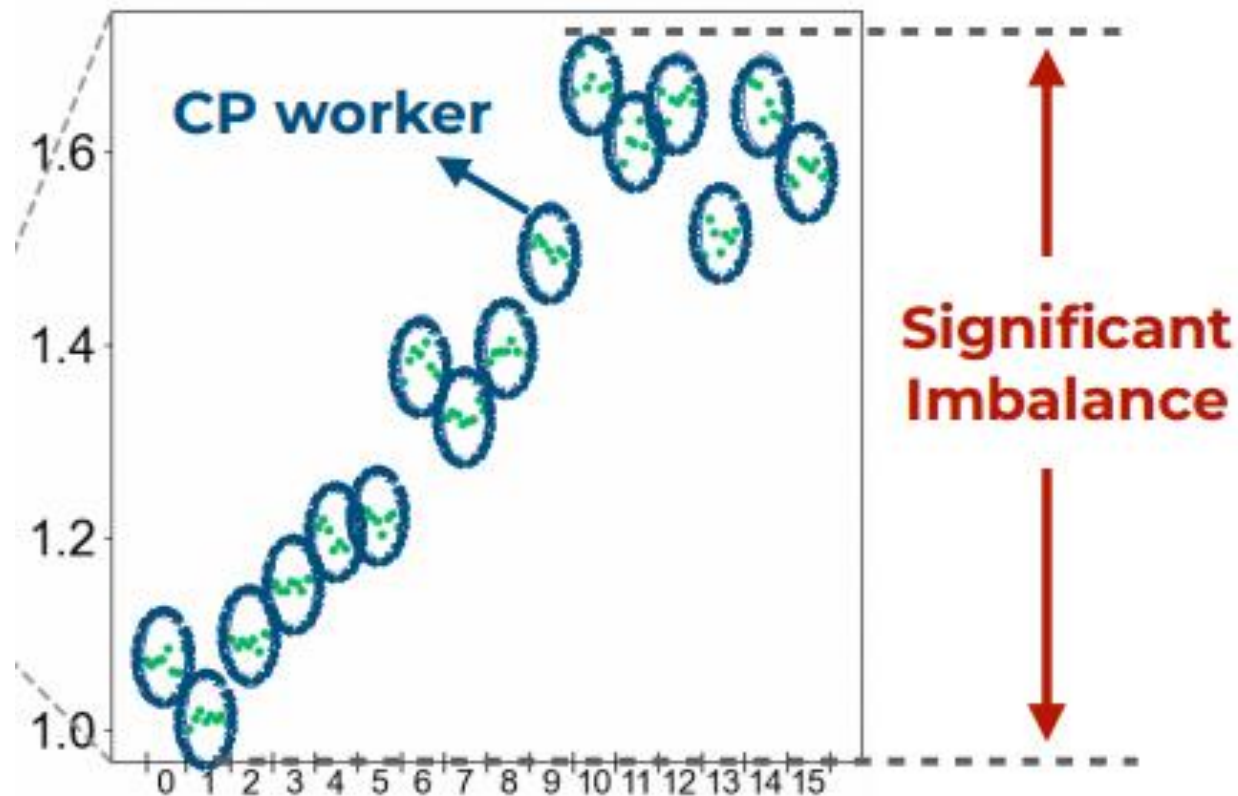


imbalance across CP workers



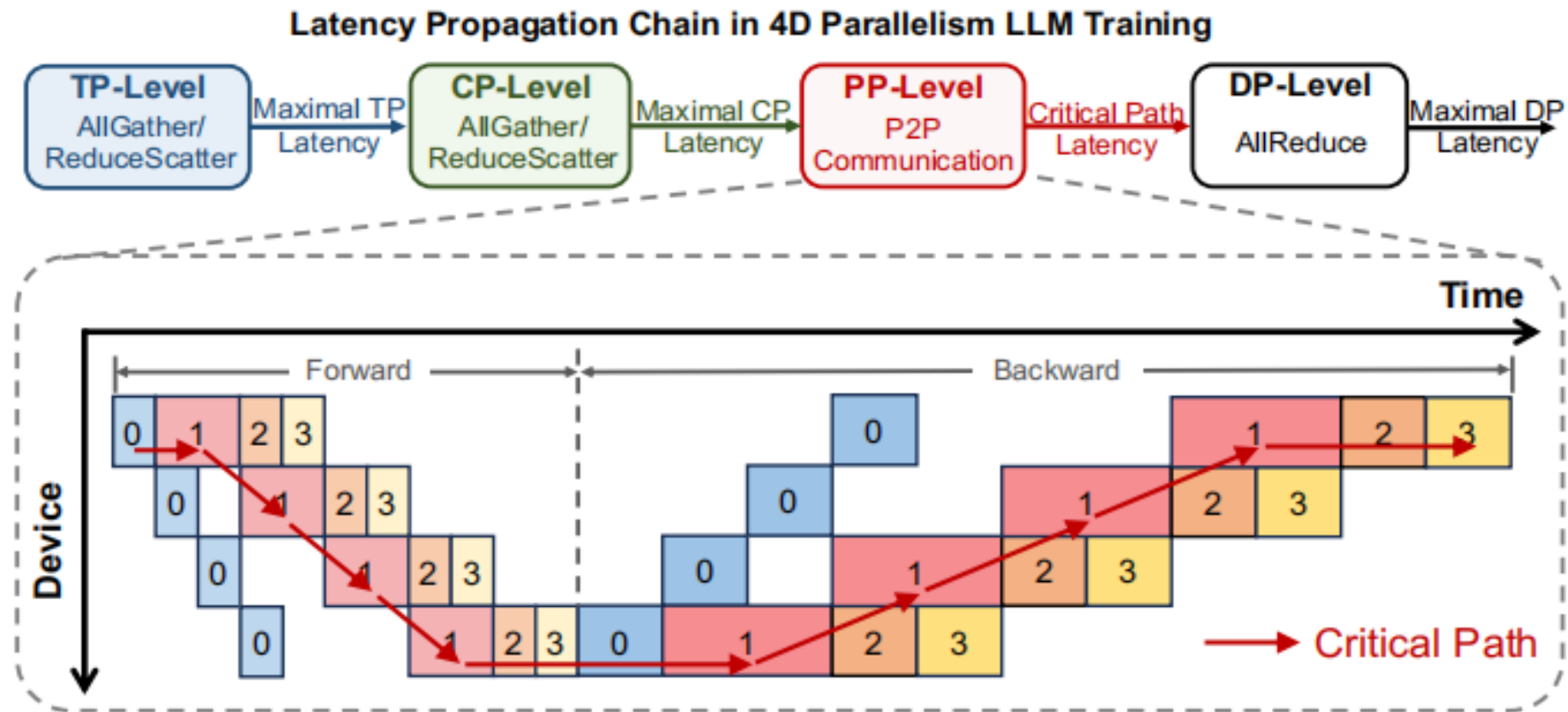
(2) Imbalance from CP-level  
Sequence Sharding

Zoom in to a PP worker (16 CP workers)



# Background: influence of imbalance

imbalance will be accumulated ➡ end-to-end training latency



# Design:WLB-LLM

- PP-level → 1.variance-length packing  
2.heuristic outlier document delay optimization
- CP-level → fine-grained and adaptive sharding optimization

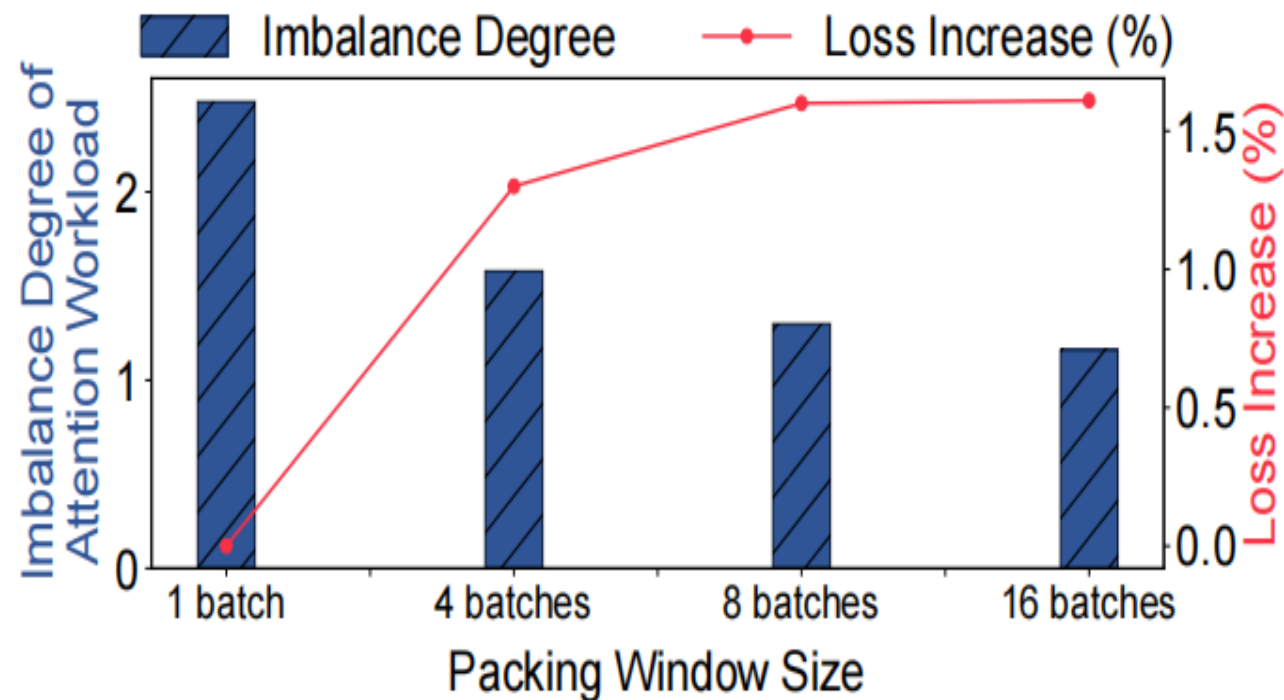
# Design:Fixed-length

Fixed-length packing



Loss increase

$$\begin{aligned} &\text{minimize} && \max\left(\sum_{i=1}^N x_{ij} \cdot d_i^2\right), j = 1, \dots, M \\ &\text{subject to} && \sum_{j=1}^M x_{ij} = 1, \quad i = 1, \dots, N \\ & && \sum_{i=1}^N x_{ij} \cdot d_i \leq L, \quad j = 1, \dots, M \\ & && x_{ij} \in \{0, 1\} \end{aligned}$$





# Design: Variable-length Packing

Fixed context length

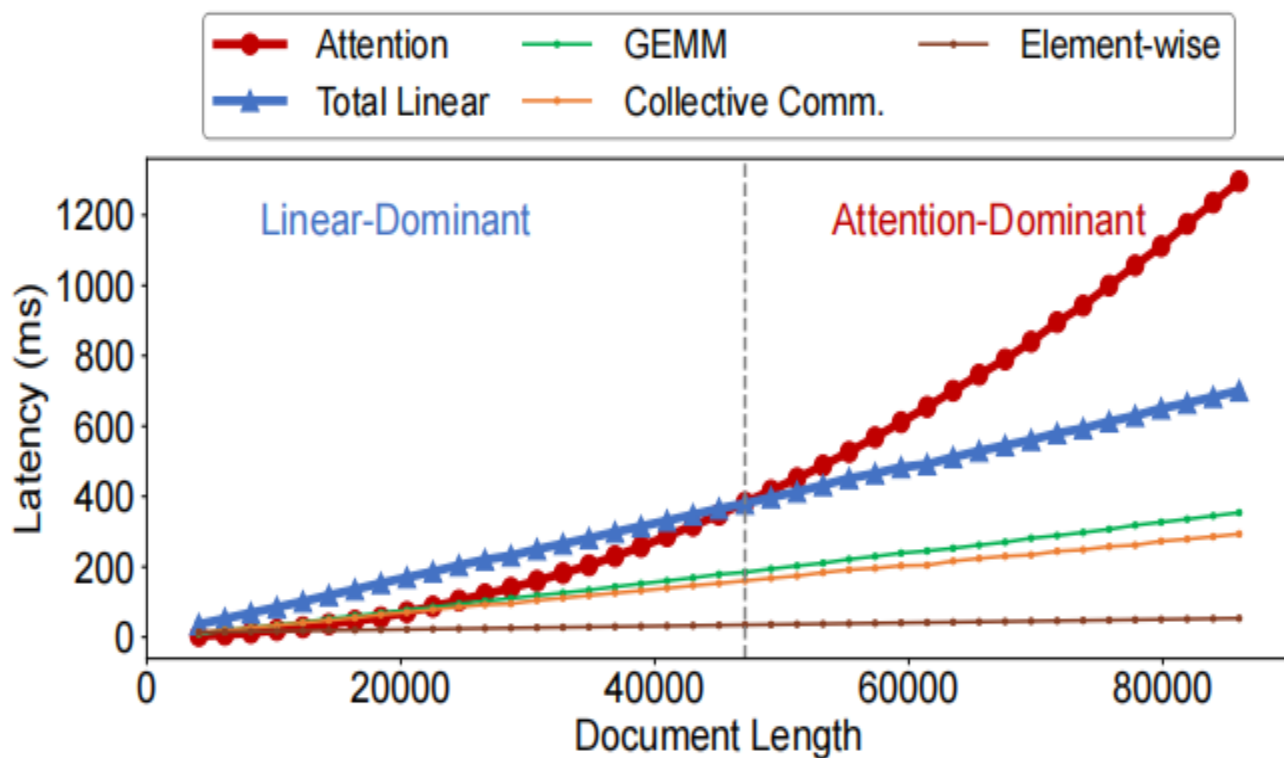


variable length

balance attention workload



balance attention and linear



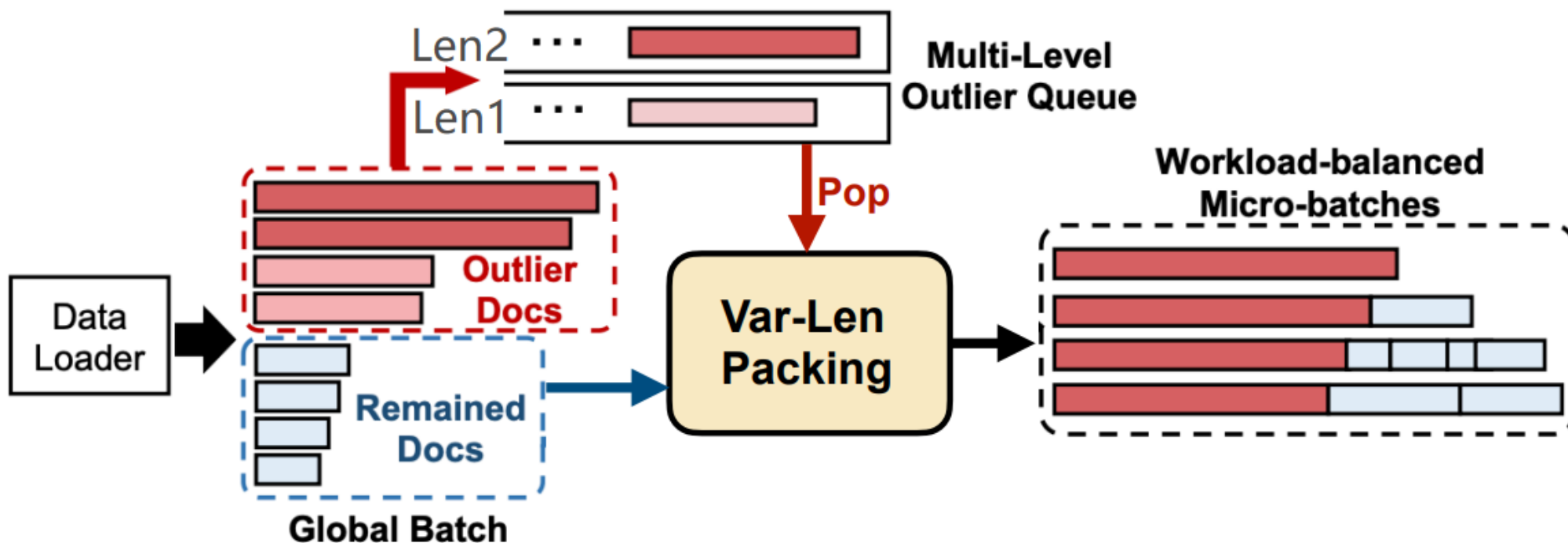
$$\begin{aligned} \text{minimize} \quad & \max \left( \sum_{i=1}^N (W_a(x_{ij} \cdot d_i) + W_l(x_{ij} \cdot d_i)) \right), \\ & j = 1, \dots, M \\ \text{subject to} \quad & \sum_{j=1}^M x_{ij} = 1, & i = 1, \dots, N \\ & \sum_{i=1}^N x_{ij} \cdot d_i \leq L_{\max}, & j = 1, \dots, M \\ & x_{ij} \in \{0, 1\} \end{aligned}$$

# Design: Outlier document delay

short document insufficient



delay outlier document



# Design:Heuristic Packing Algorithm

ILP solver takes a long time → heuristic packing algorithm

## Algorithm 1: Heuristic Var-length Packing Algorithm

```
input :Dataloader:  $D$ , Waiting Queues:  $Q$ ,  
        #Micro-Batch per Iteration:  $N$ ,  
        Sequence Length Upper bound:  $L_{max}$   
output :Packed Micro-Batches for Training:  $B$   
1  Remained_Doc = [];  
2  for Cur_Batch in  $D$  do  
3      Doc_Set = Remained_Doc;  
4      for Doc in Cur_Batch do  
5          /* Delay the execution of outlier documents. */  
6          if Doc.Is_Outlier() then  
7              | Q.Add(Doc);  
8          else  
9              | Doc_Set.Push(Doc);  
10         end  
11         for  $q$  in  $Q$  do  
12             if  $\text{len}(q) \geq N$  then  
13                 /* Pop outlier documents for the current batch. */  
14                 | Doc_Set.Push(q.Pop(N));  
15             end  
16         end  
17         /* Sort the documents in descending order by length. */  
18         Doc_Set.Sort_by_Length();
```

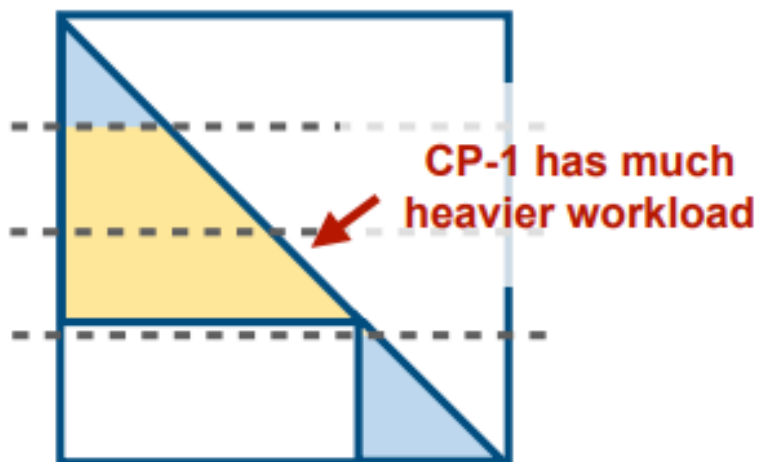
```
17     /* Start packing. */  
18     New_Batch = Create_Batch( $N$ );  
19     for Doc in Doc_Set do  
20         /* Get micro-batches with minimum workload/length. */  
21         W_idx = New_Batch.Get_Min_Workload();  
22         L_idx = New_Batch.Get_Min_Length();  
23         if New_Batch[W_idx].Len() + Doc.Len() <  $L_{max}$  then  
24             | New_Batch[W_idx].Push(Doc);  
25         else  
26             if New_Batch[L_idx].Len() + Doc.Len() <  $L_{max}$  then  
27                 | New_Batch[L_idx].Push(Doc);  
28             else  
29                 | Remained_Doc.Push(Doc);  
30             end  
31         end  
32     end  
33     B.Push(New_Batch);
```

# Design:per-document sharding

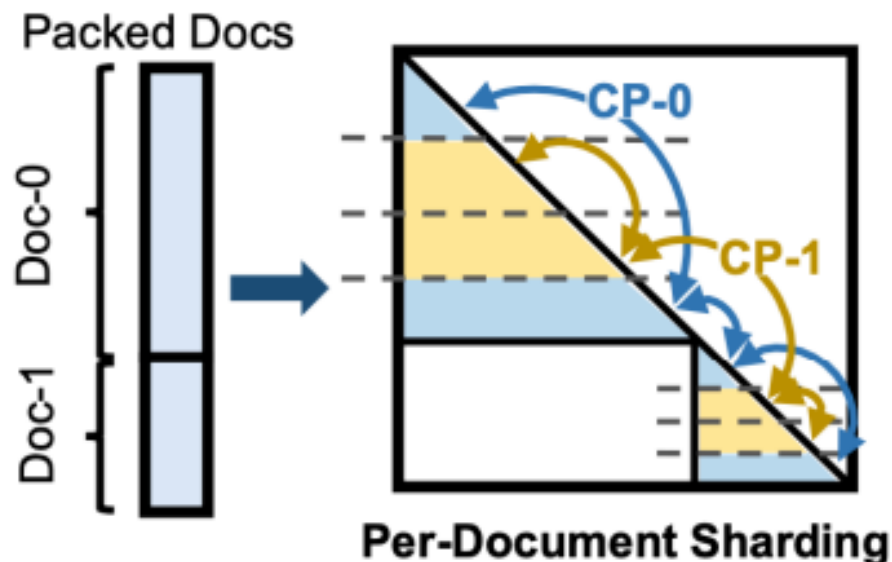
per-sequence sharding



per-document sharding



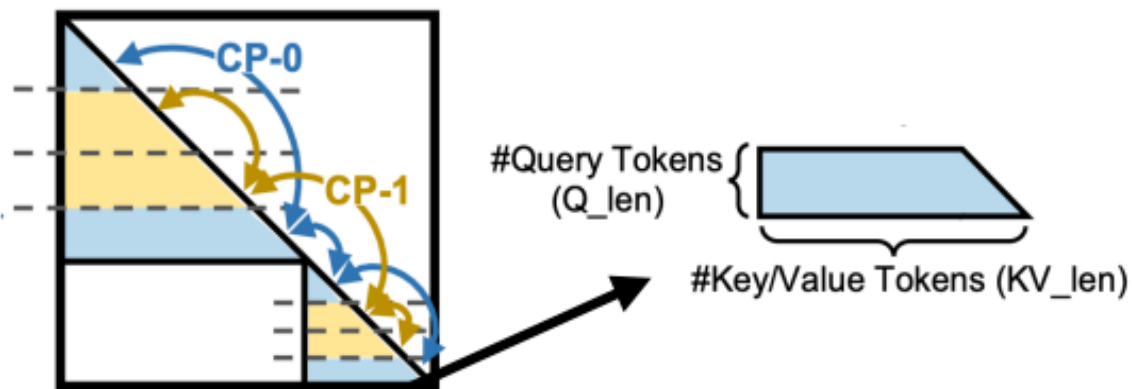
**Per-Sequence Sharding**



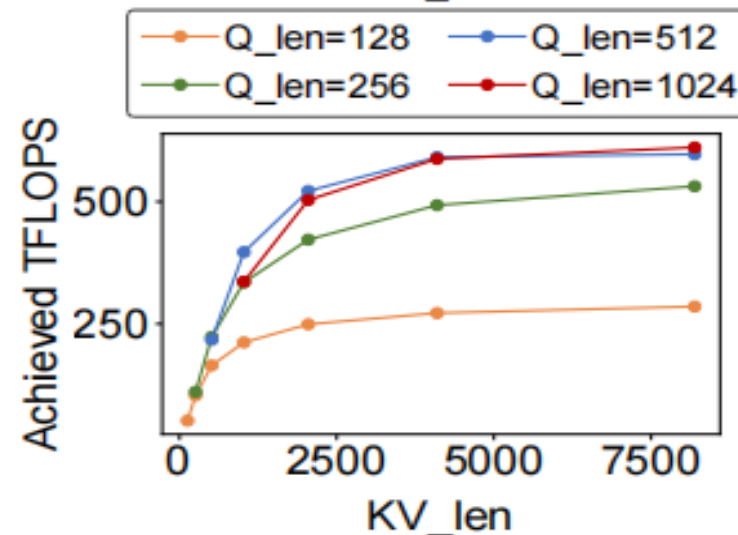
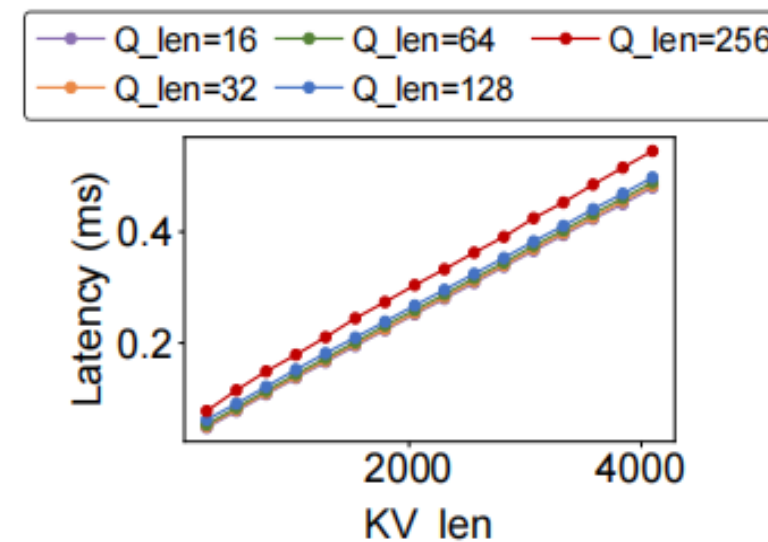
**Per-Document Sharding**

# Design: per-document sharding

sharding balance vs. kernel efficiency



**Fine-grained sharding may leads to small document shards**

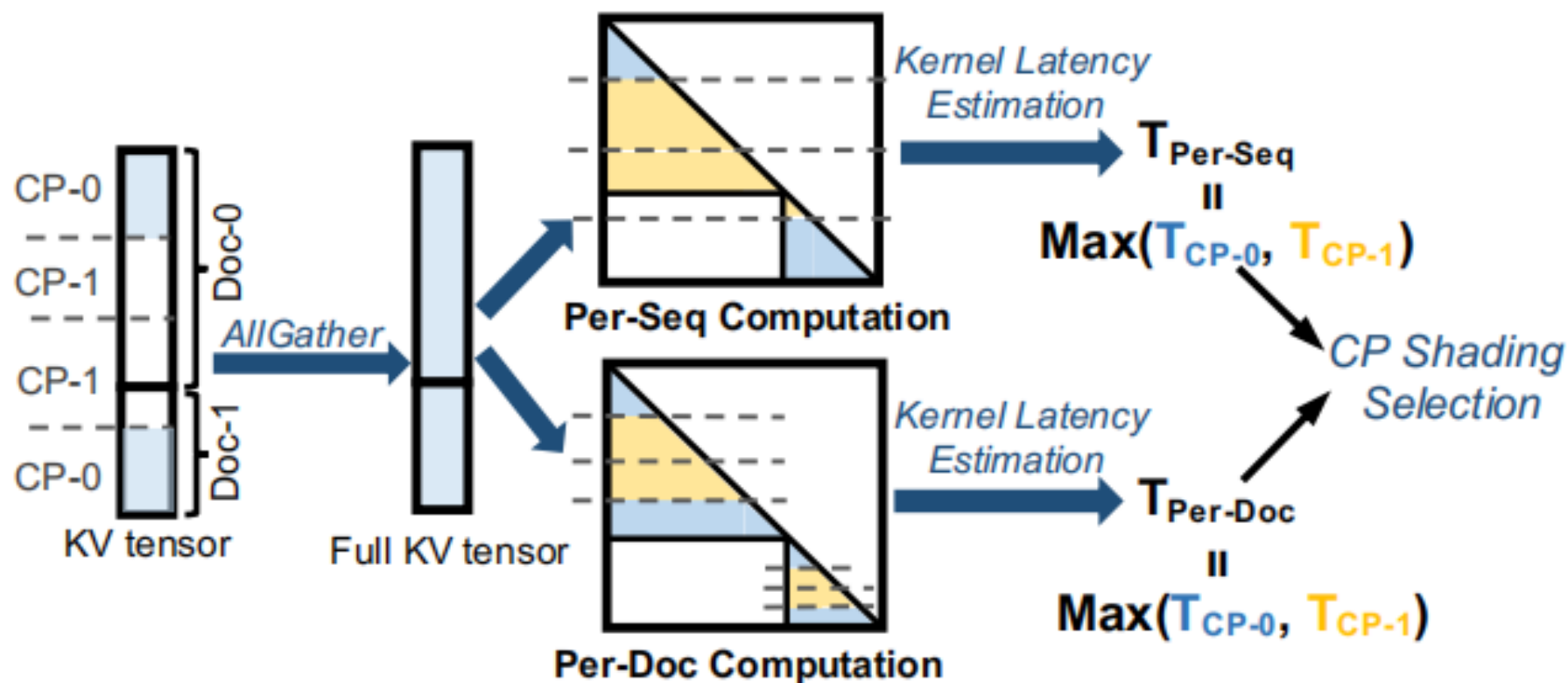


# Design:adaptive sharding

trade-off



adaptive sharding



# Evaluation: setting

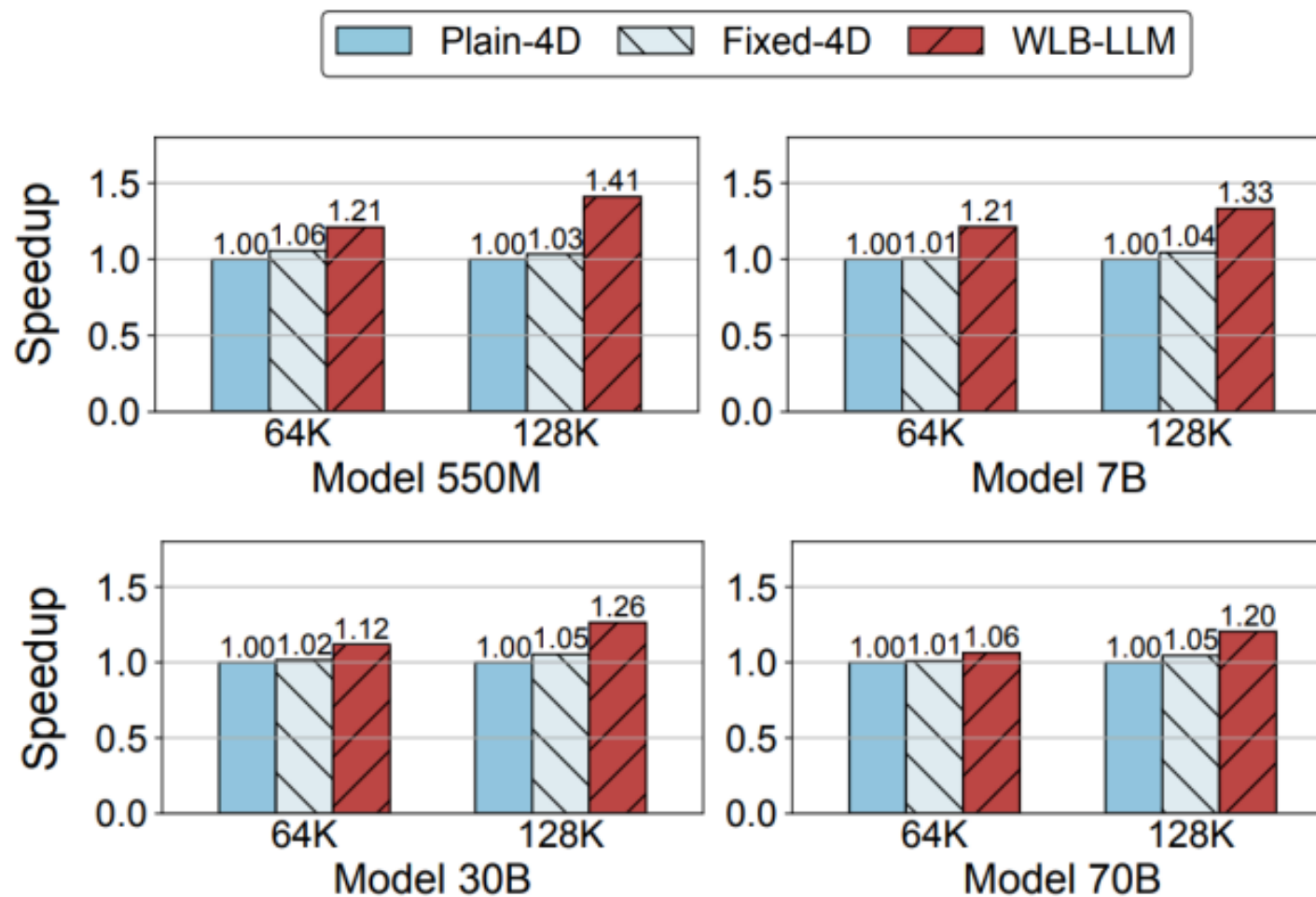
model and parallelism config

| Model Size | Context Window | #GPU | 4D Parallelism Configs (TP, CP, PP, DP) |
|------------|----------------|------|---|
| 550M       | 64K            | 32   | (2, 2, 4, 2)                            |
|            | 128K           | 32   | (2, 4, 4, 1)                            |
| 7B         | 64K            | 32   | (4, 2, 4, 1)                            |
|            | 128K           | 64   | (8, 2, 4, 1)                            |
| 30B        | 64K            | 64   | (8, 2, 4, 1)                            |
|            | 128K           | 128  | (8, 4, 4, 1)                            |
| 70B        | 64K            | 256  | (16, 4, 4, 1)                           |
|            | 128K           | 256  | (16, 4, 4, 1)                           |

Baseline

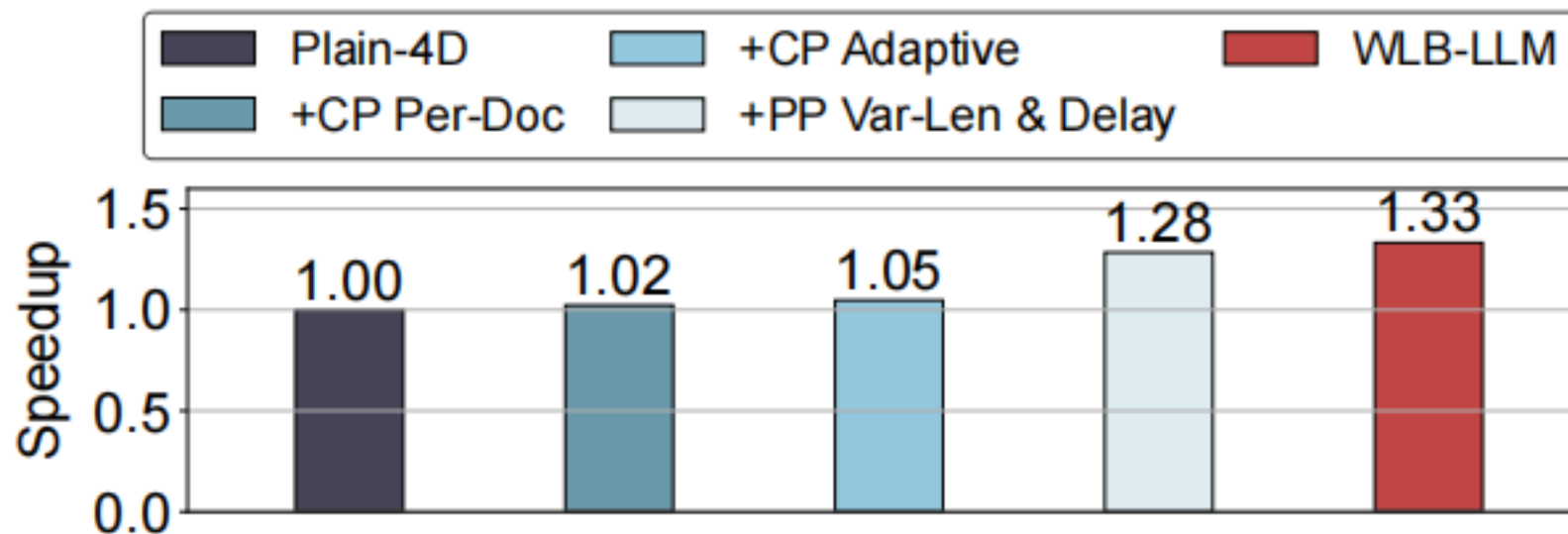
- plain-4D:
  - no optimizing packing
  - per-seq sharding
- fixed-4D:
  - fixed-len packing
  - better sharding

# Evaluation





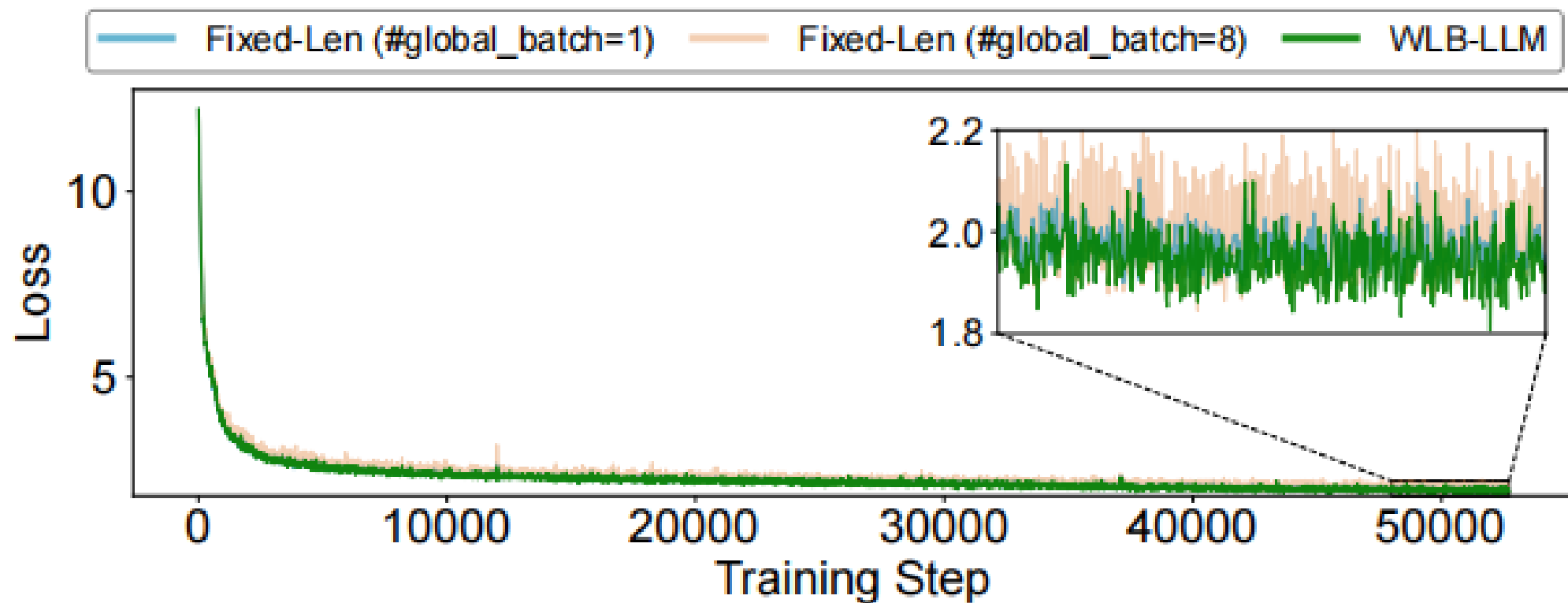
# Evaluation: optimization analysis



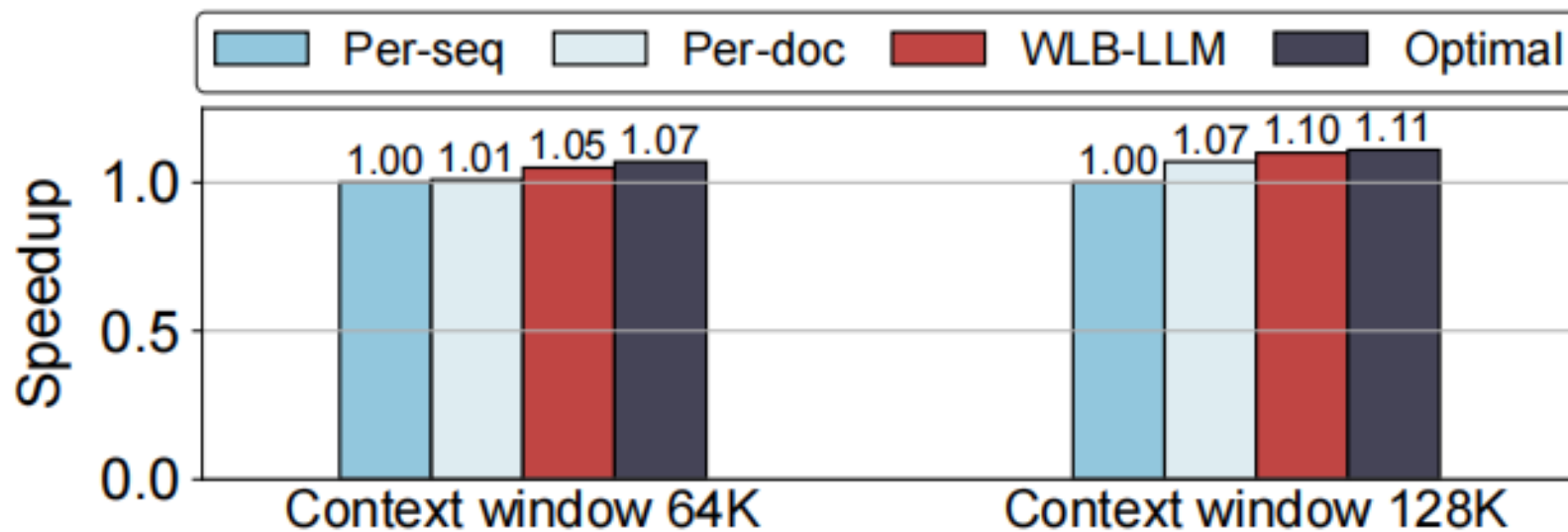
# Evaluation: packing

| Packing Method          |                 | Imbalance Degree | Packing Overhead (ms) |
|-------------------------|-----------------|------------------|-----------------------|
| Method                  | Config          |                  |                       |
| <i>Original Packing</i> | /               | 1.44             | 0                     |
| <i>Fixed-Len Greedy</i> | #global batch=1 | 1.41             | 4                     |
|                         | #global batch=2 | 1.22             | 5                     |
|                         | #global batch=4 | 1.11             | 5                     |
|                         | #global batch=8 | 1.08             | 5                     |
| <i>Fixed-Len Solver</i> | #global batch=1 | 1.40             | 467                   |
|                         | #global batch=2 | 1.18             | 1488                  |
|                         | #global batch=4 | 1.09             | 25313                 |
| <i>WLB-LLM</i>          | #queue=1        | 1.24             | 8                     |
|                         | #queue=2        | 1.05             | 20                    |
|                         | #queue=3        | 1.05             | 23                    |

# Evaluation: model convergence



# Evaluation: sharding



# Thinking

- 文章有什么疑问  
在考虑本文方法对模型效果的影响时，其基线对比的是用固定长度 packing 方法的效果，而不是和不加任何优化的模型训练作对比，其影响程度未知
- 在此之上：  
文中的 packing 算法的针对异常文档的队列的长度阈值是预设的，是否可以根据数据的分布自适应的确定长度阈值
- 于此平行：  
本文是优化训练过程中的负载不均衡提升性能，对于推理服务系统能否使用相同方法，优化推理延迟

Thanks