# Taming Load Balancing in Distributed LLM Training

Jiale Xu

# Table of Contents

- Background
- Prior Work
- Challenge
- *OSDI'25'* **WLB-LLM**
- Discussion on other load balancing issue in MLSys topic.

# Background: Distributed LLM Serving Paradigm-Memory Footprint
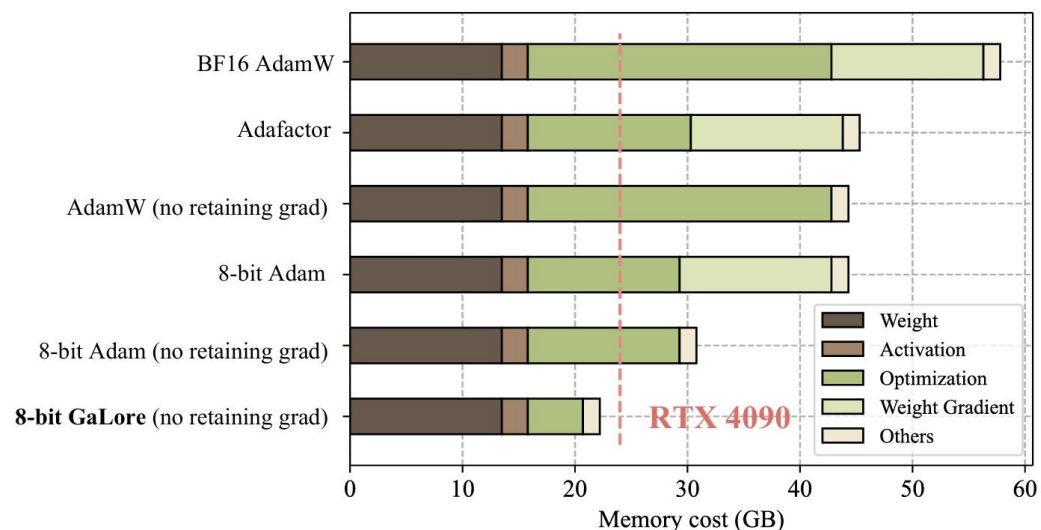


Figure 1: Estimated memory consumption of pre-training a LLaMA 7B model with a token batch size of 256 on a single device, without activation checkpointing and memory offloading[2]. Details refer to Section 5.5.
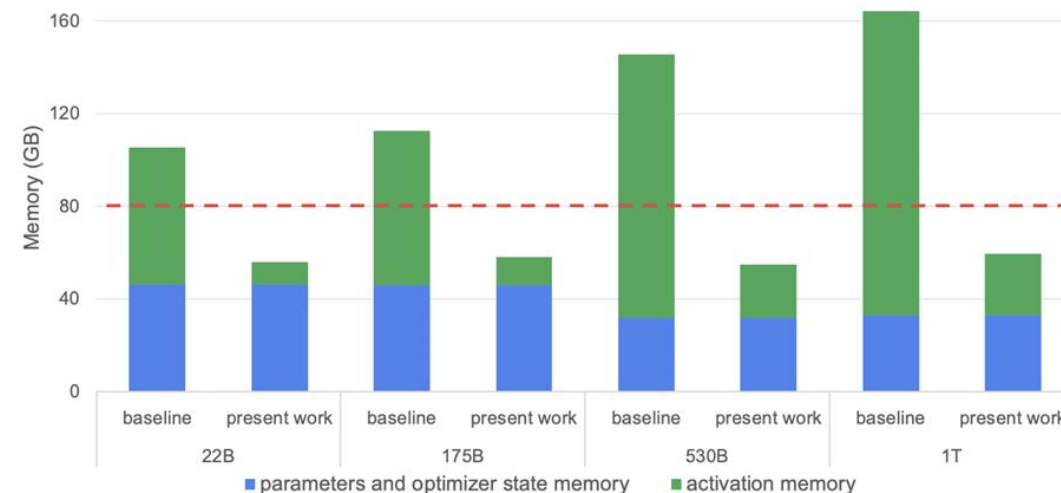


Figure 1: Parameters, optimizer state, and activations memory. The dashed red line represents the memory capacity of an NVIDIA A100 GPU. Present work reduces the activation memory required to fit the model. Details of the model configurations are provided in Table 3.

Model weight and Optimizer State is dominating when apply optimizations. ⟶ Activation becomes a bottleneck in memory footprint.

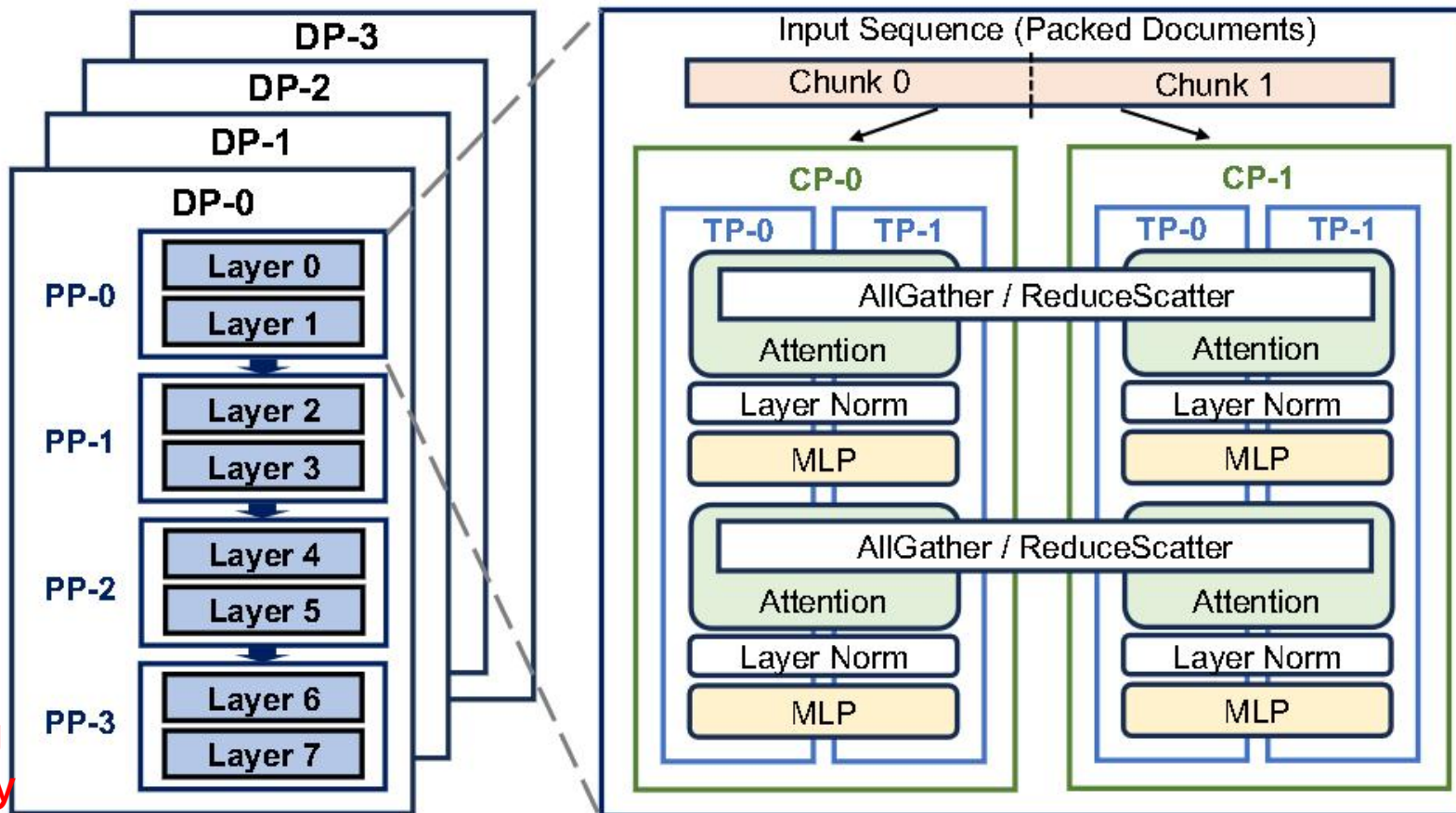After distributed training and model compression techniques.

# Background: Distributed LLM Serving Paradigm-4D Parallel

**1. Data Parallel**
Enlarge input batch Size while duplicating model weight memory

**2. Pipeline Parallel**
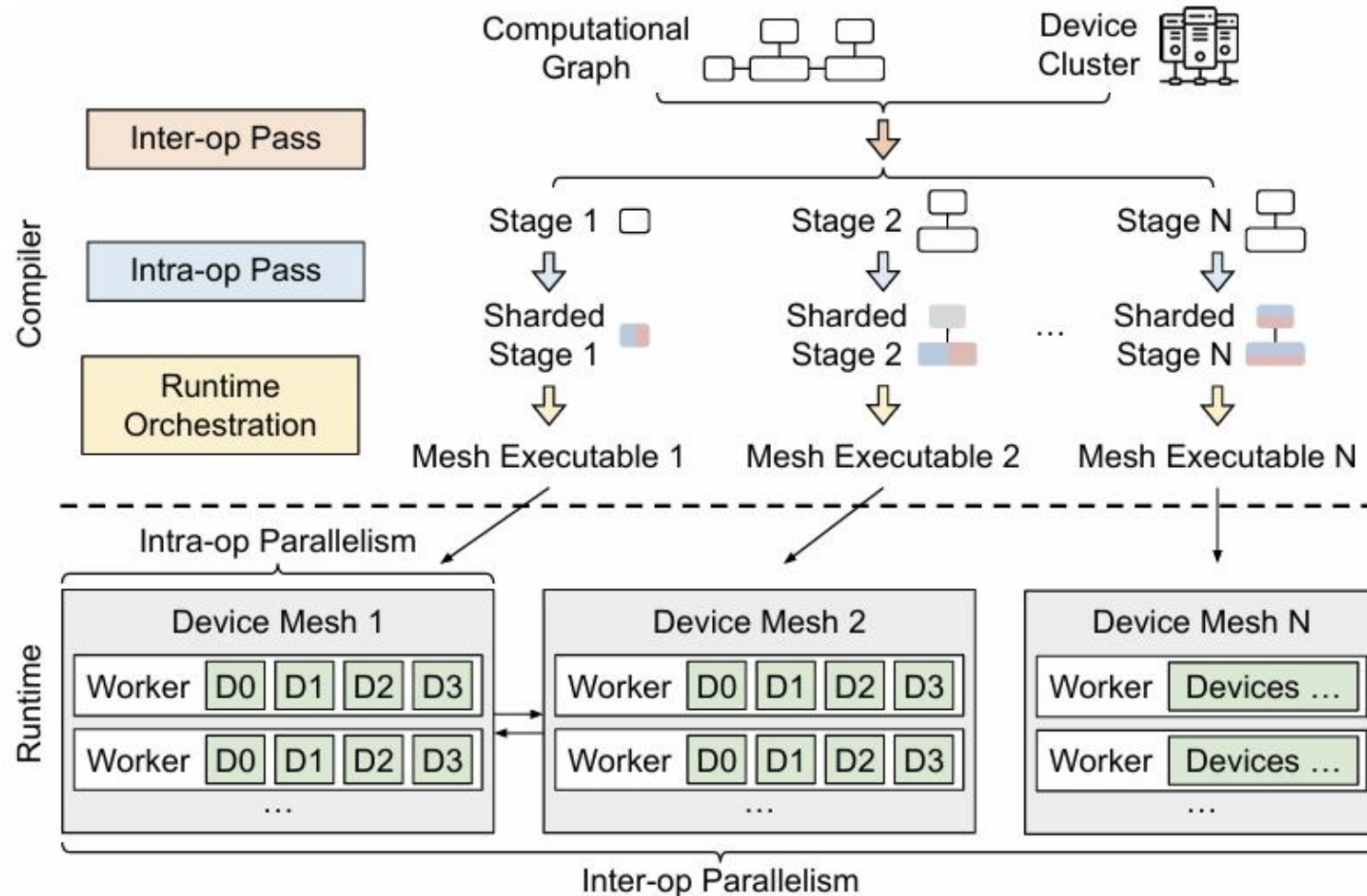Reducing model weight memory while undermining compute efficiency

**3. Context Parallel**
Reducing activation memory & increasing compute capability while duplicating model weight memory & communication overhead
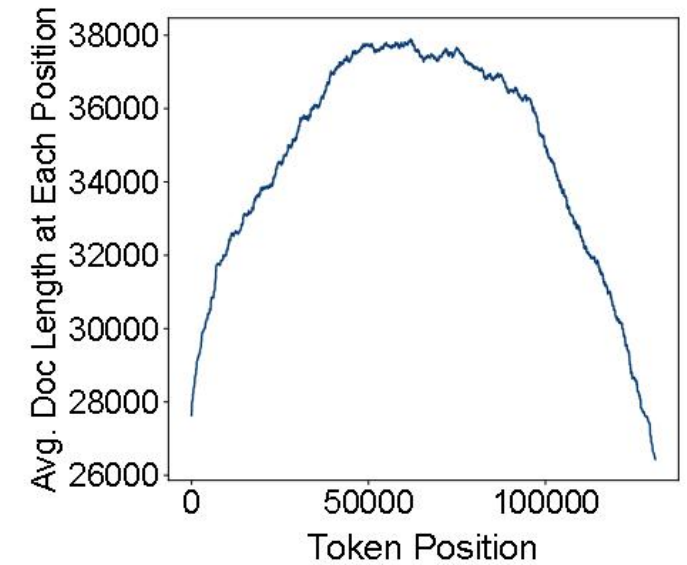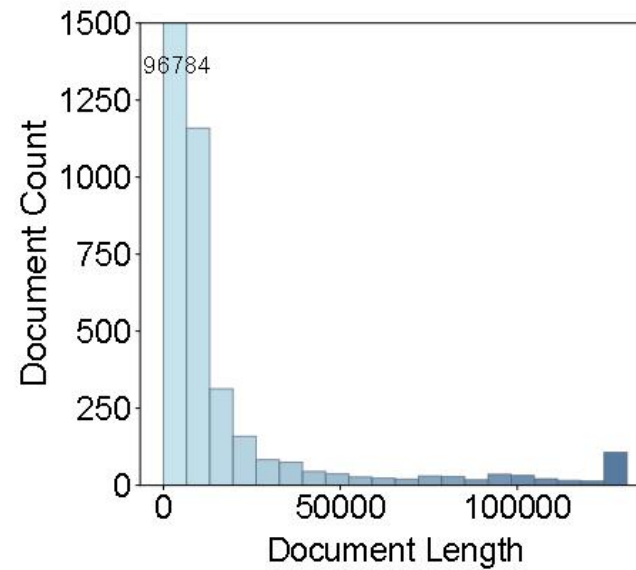
**4. Tensor Parallel**
Reducing overall memory & increasing compute capability while introducing communication overhead

# *Eternal Goal: Finding a more Efficient Distributed Training Plan*

# *Challenges*

Llama2-7b: 4K → Qwen-2.5: 1M



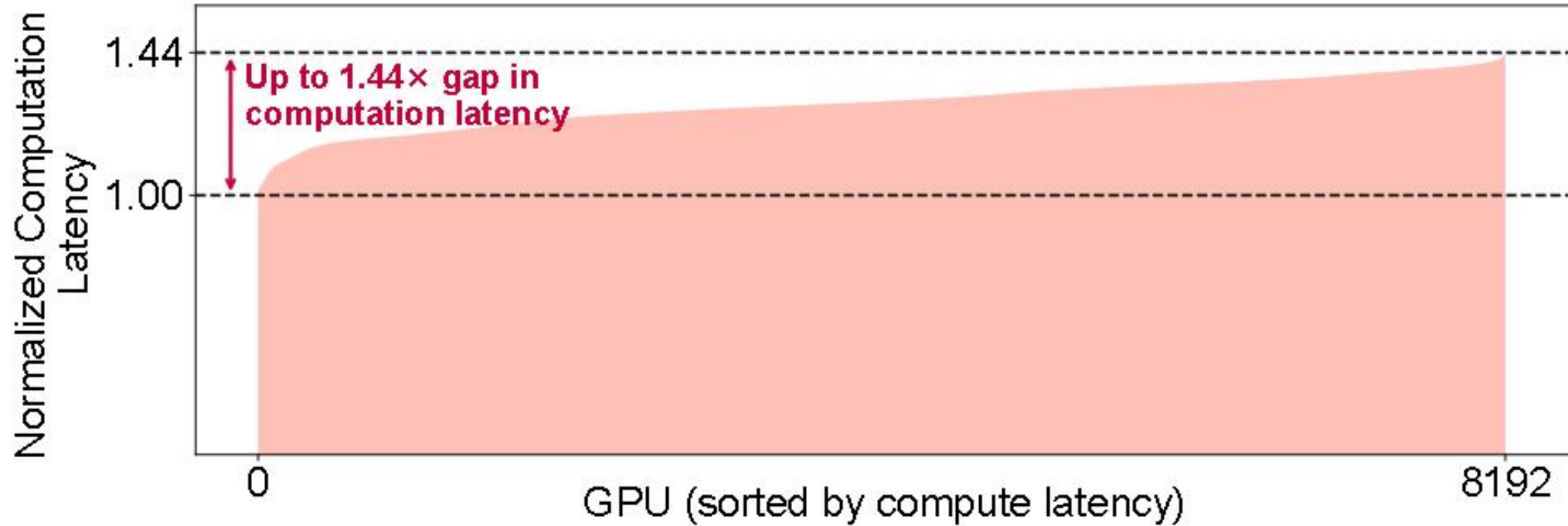**Growing context window size**

**Varied dataset length distribution**

# OSDI'25 WLB-LLM

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

Zheng Wang1,2, Anna Cai2, Xinfeng Xie2, Zaifeng Pan1, **Yue Guan**1,
Weiwei Chu2, Jie Wang2, Shikai Li2, Jianyu Huang2, Chris Cai2, Yuchen
Hao2, Yufei Ding1,2
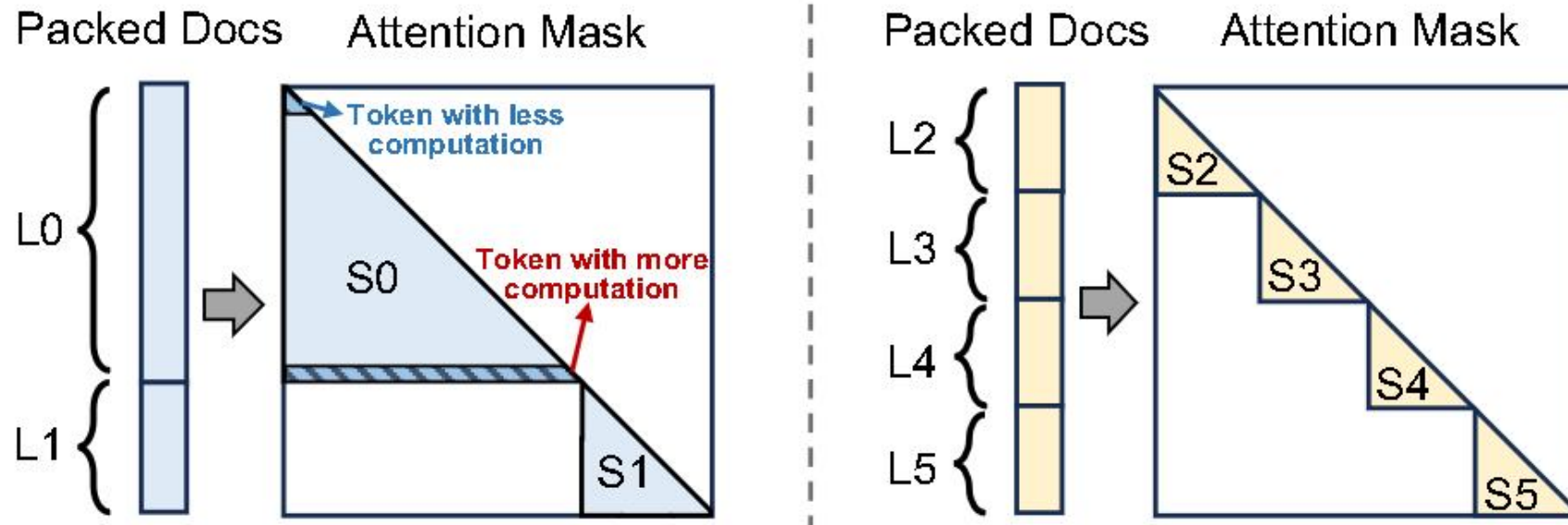University of California, San Diego1 Meta2

# *Finding 1:* Imbalanced computation latency within large-scale LLM training job.



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

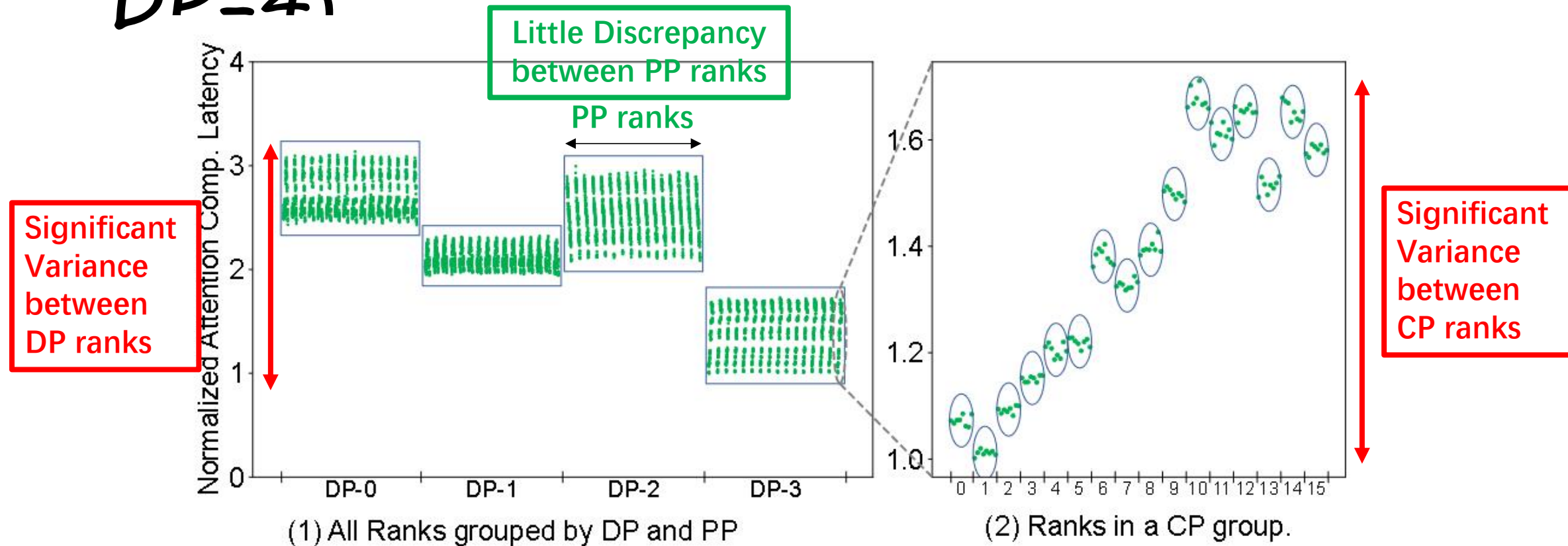# *Finding 2:* Imbalanced Reason: Attention Computation Nature



Document lengths: L0 + L1 = L2 + L3 + L4 + L5

Computation (triangle areas): S0 + S1 >> S2 + S3 + S4 + S5

(b) Reason of imbalance: input-dependent nature of attention computation and the varying input document length.
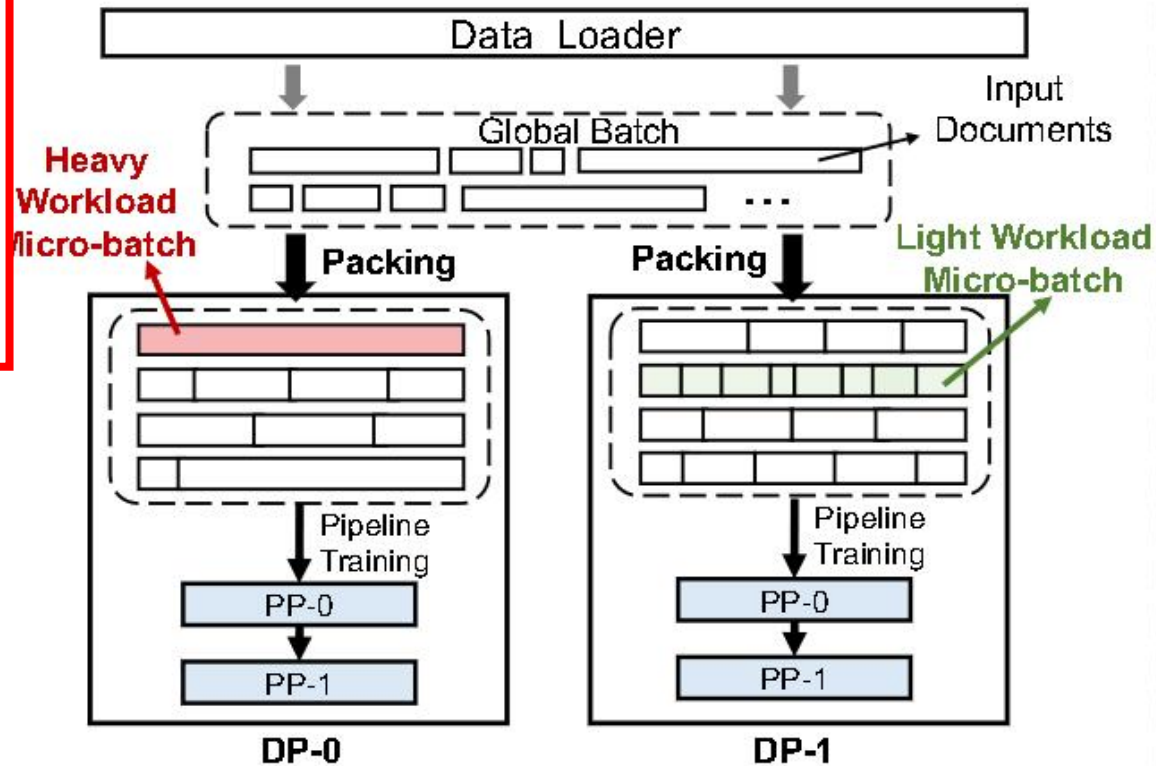
# Case Study: (TP=8, CP=16, PP=16, DP=4)



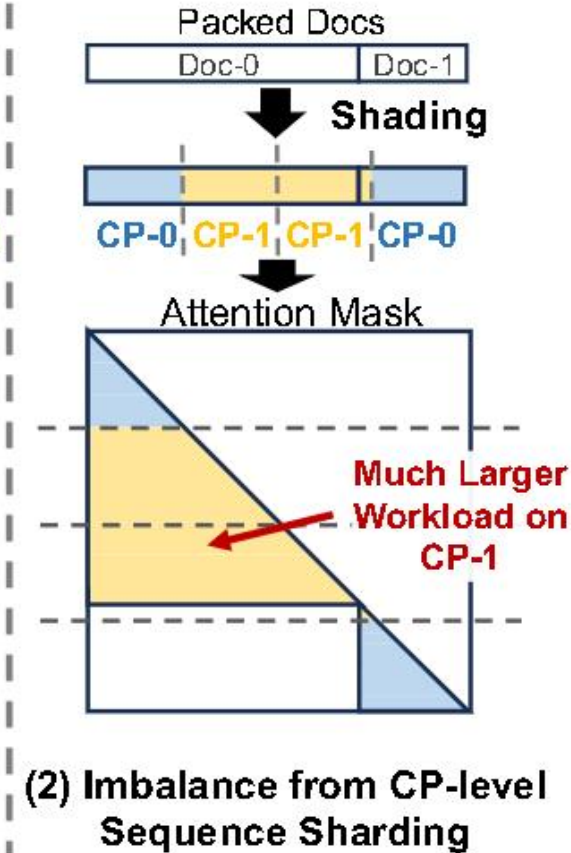(a) Imbalance Analysis (TP=8, CP=16, PP=16, DP=4): (1) Normalized computation latency (group by DP and PP); (2) Normalized computation latency in a CP group.

# Case Study: (TP=8, CP=16, PP=16, DP=4)

**DP ranks have different micro batches**



(b) Document packing at PP level and sequence sharding at CP level.
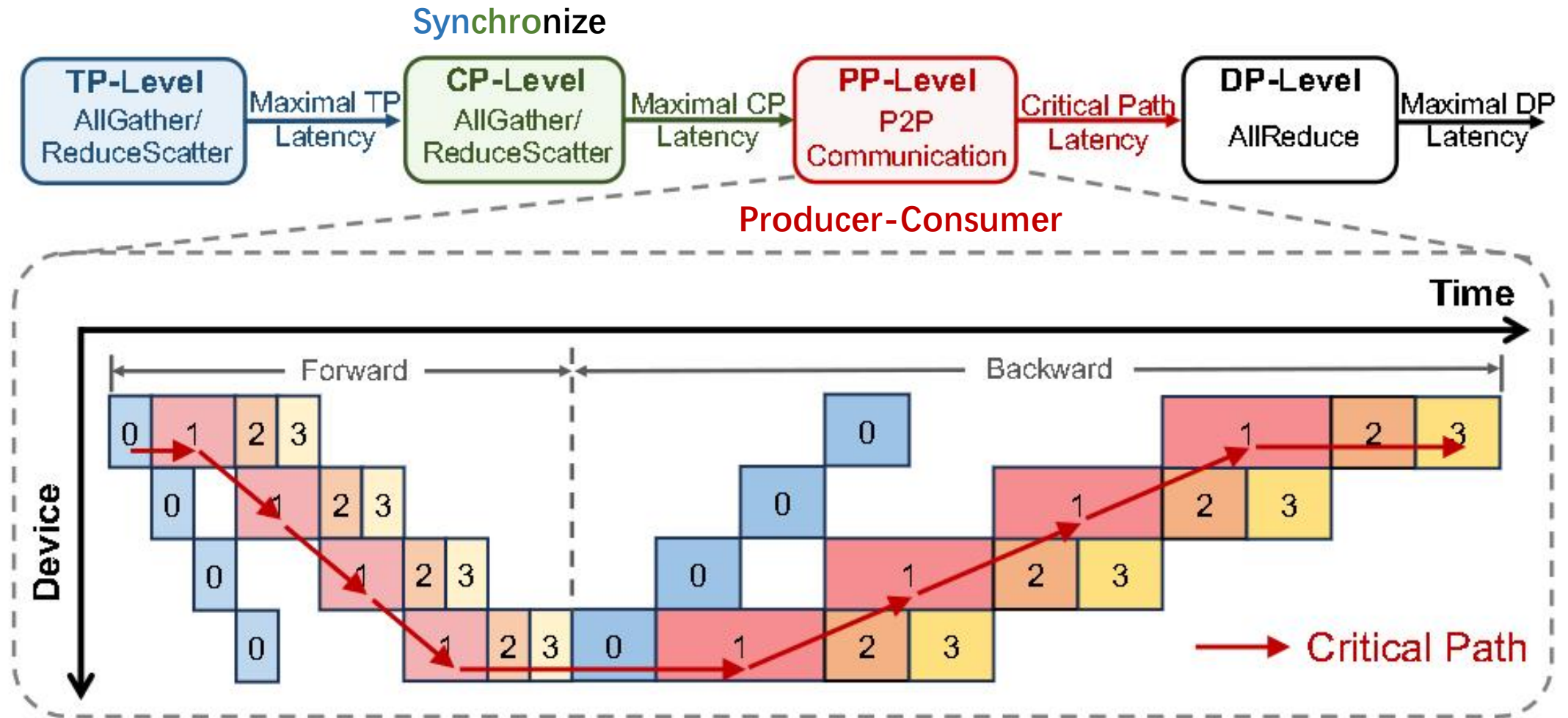
(1) Imbalance from PP-level Document Packing

(2) Imbalance from CP-level Sequence Sharding

**CP ranks have attention load balancing issue**

**PP ranks share the same micro-batch**

# Latency Propagation Chain

# Two Optimizations

**Better PP packing**

**Better Document shard balancing**



(1) Imbalance from PP-level Document Packing

(2) Imbalance from CP-level Sequence Sharding

(b) Document packing at PP level and sequence sharding at CP level.

# Better PP Packing



$$\text{minimize} \quad \max(\sum_{i=1}^{N}(W_a(x_{ij} \cdot d_i) + W_l(x_{ij} \cdot d_i))),$$

$$j = 1, \cdots, M$$

$$\text{subject to} \quad \sum_{j=1}^{M} x_{ij} = 1, \qquad i = 1, \cdots, N \qquad (2)$$

$$\sum_{i=1}^{N} x_{ij} \cdot d_i \leq L_{max}, \qquad j = 1, \cdots, M$$

$$x_{ij} \in \{0, 1\}$$

# Better PP Packing

# *Better Document shard balancing*

# Better Document shard balancing



Less efficient computation
Better load balancing

# Discussion On Kernel Load Balancing



Figure 10: Attention kernel performance profiling: (Left) Attention forward latency; (Right) Achieved TFLOPs of the attention forward kernel.



**FlashInfer with load balancing logic**

# Evaluation – Setup

| 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) |

Table 1: Model and 4D parallelism configurations.

# Evaluation – E2E

Takeaway #1.  Larger Model, Less improve. Reason: Larger communication

Figure 12: Training performance speedups of *WLB-LLM* and *Fixed-4D* over *Plain-4D* across various configurations.

Figure 13: Performance breakdown of *WLB-LLM* on the 7B model with a 128K context window.

Figure 14: Speedups of *WLB-LLM* on the 7B model across context window sizes.

Takeaway #2.  PP-Level Load balancing is more important.

Takeaway #3. Larger Context lead to more speedup Reason: larger context window raises the likelihood of outlier documents appearing

# Related Work(s)

SOSP'24 Enabling Parallelism Hot Switching for Ecient Training of Large Language Models

ASPLOS'25 FlexSP: Accelerating Large Language Model Training via Flexible Sequence Parallelism

|  | Targeted Parallelism |
| --- | --- |
| Hot-switching | TP |
| FlexSP | SP |
| WLB-LLM | PP+SP |

# Comments

Pros:
1.  Direct and strong motivation and clear design writing.
2.  Comprehensive evaluation with decent speedup, sensible analysis, and progressive breakdown.

Cons:
1.  Not a very impressive/novel idea.

# Two ways to achieve load balancing

1.  **Request Scheduling**
    Dispatching request evenly
    to different workers.


2.   **Resource Re-orchestration**
Re-allocate unused resources.

# Hot-switching (SOSP'24)



**Deduction, Orchestration, and Instantiation**
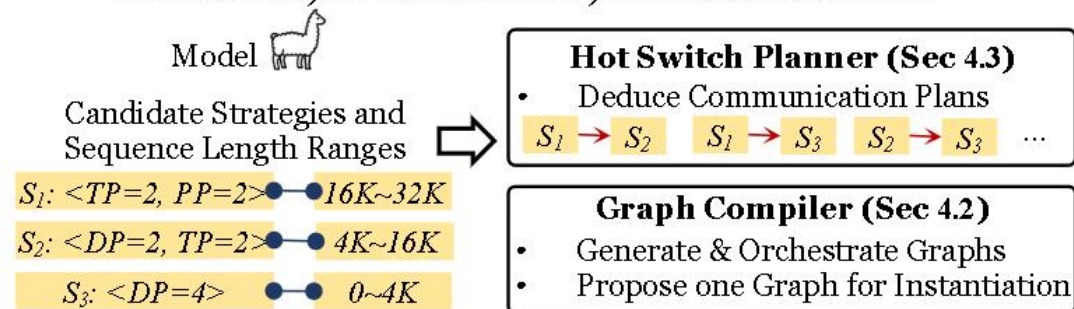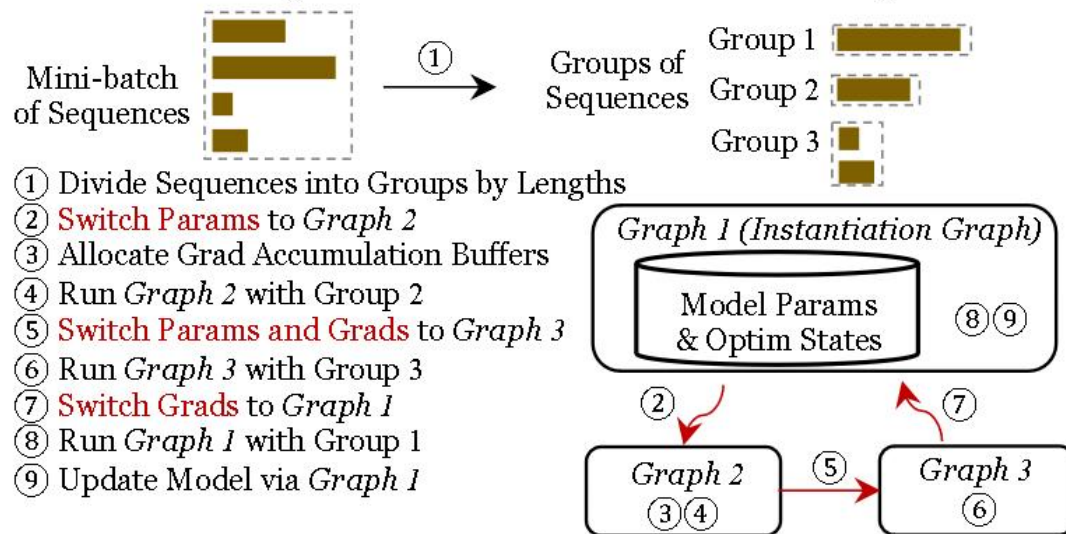
Model

Candidate Strategies and
Sequence Length Ranges

$S_1$: <TP=2, PP=2> — 16K~32K

$S_2$: <DP=2, TP=2> — 4K~16K

$S_3$: <DP=4> — 0~4K

**Hot Switch Planner (Sec 4.3)**
- Deduce Communication Plans

$S_1 \rightarrow S_2$   $S_1 \rightarrow S_3$   $S_2 \rightarrow S_3$ ...

**Graph Compiler (Sec 4.2)**
- Generate & Orchestrate Graphs
- Propose one Graph for Instantiation

**Training with Parallelism Hot Switching**

Mini-batch of Sequences

① Groups of Sequences

Group 1
Group 2
Group 3

① Divide Sequences into Groups by Lengths
② Switch Params to *Graph 2*
③ Allocate Grad Accumulation Buffers
④ Run *Graph 2* with Group 2
⑤ Switch Params and Grads to *Graph 3*
⑥ Run *Graph 3* with Group 3
⑦ Switch Grads to *Graph 1*
⑧ Run *Graph 1* with Group 1
⑨ Update Model via *Graph 1*

*Graph 1 (Instantiation Graph)*

Model Params
& Optim States   ⑧⑨

② ⑦

*Graph 2*
③④

⑤

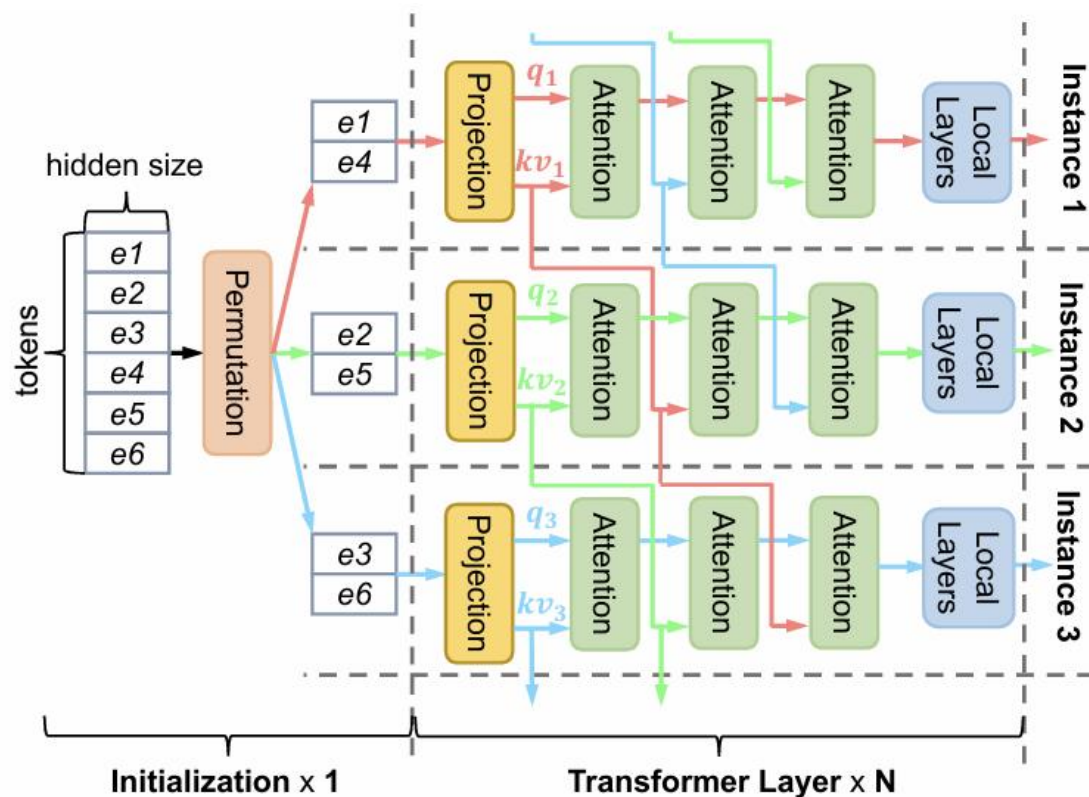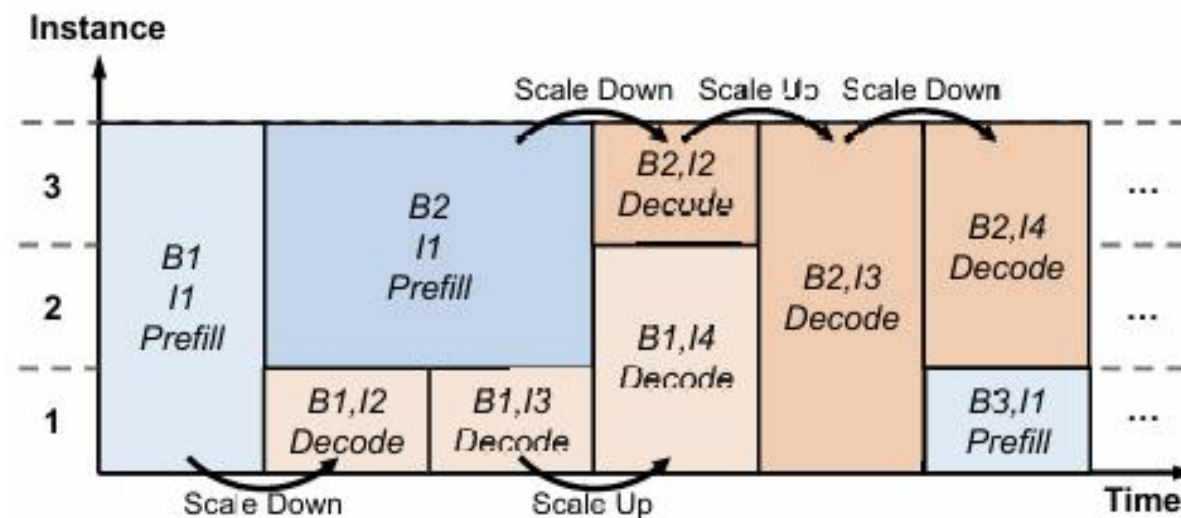*Graph 3*
⑥

# Hot-switching (SOSP'24)

# LoongServe (SOSP'24)



**Elastic Sequence Parallelism**



**Runtime Example**
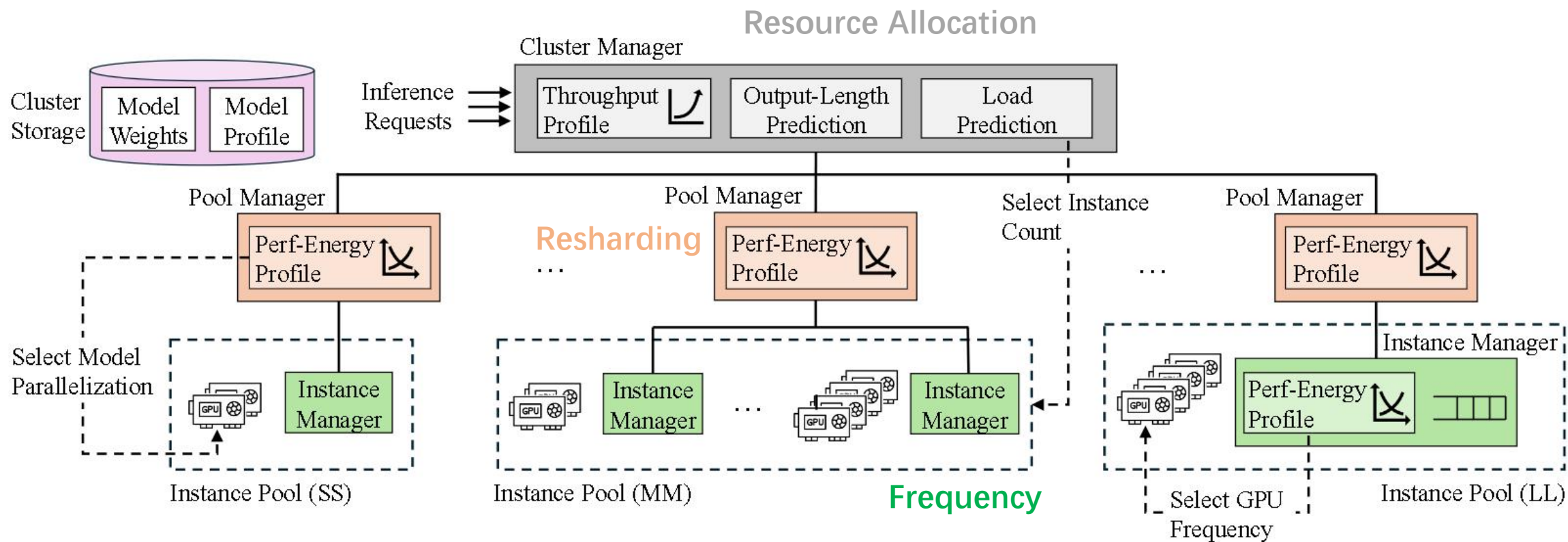
# DynamoLLM (HPCA'24 Best Paper)



Fig. 4: DynamoLLM architecture: a hierarchy of controllers with cluster resources split into per request-type pools.
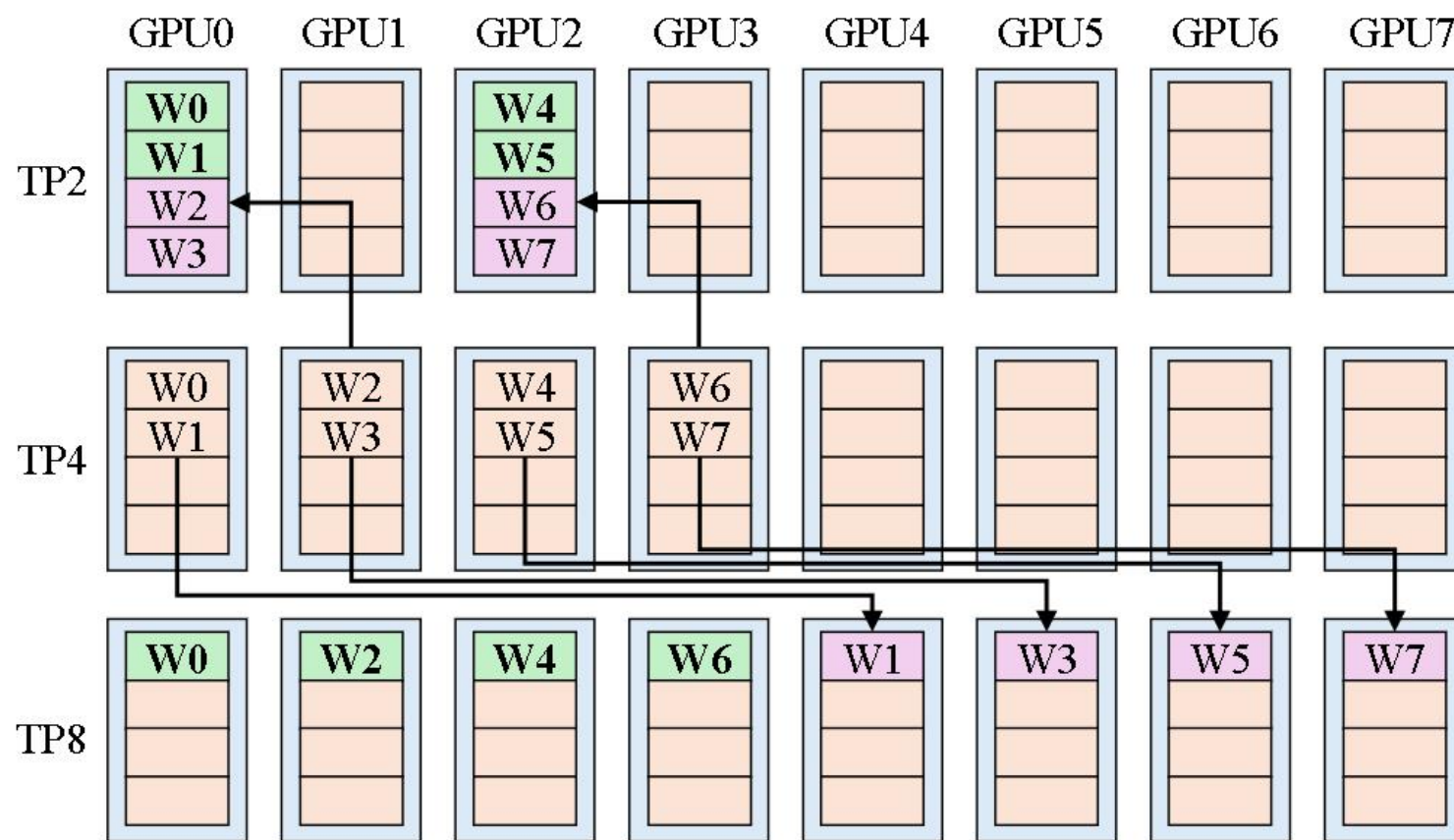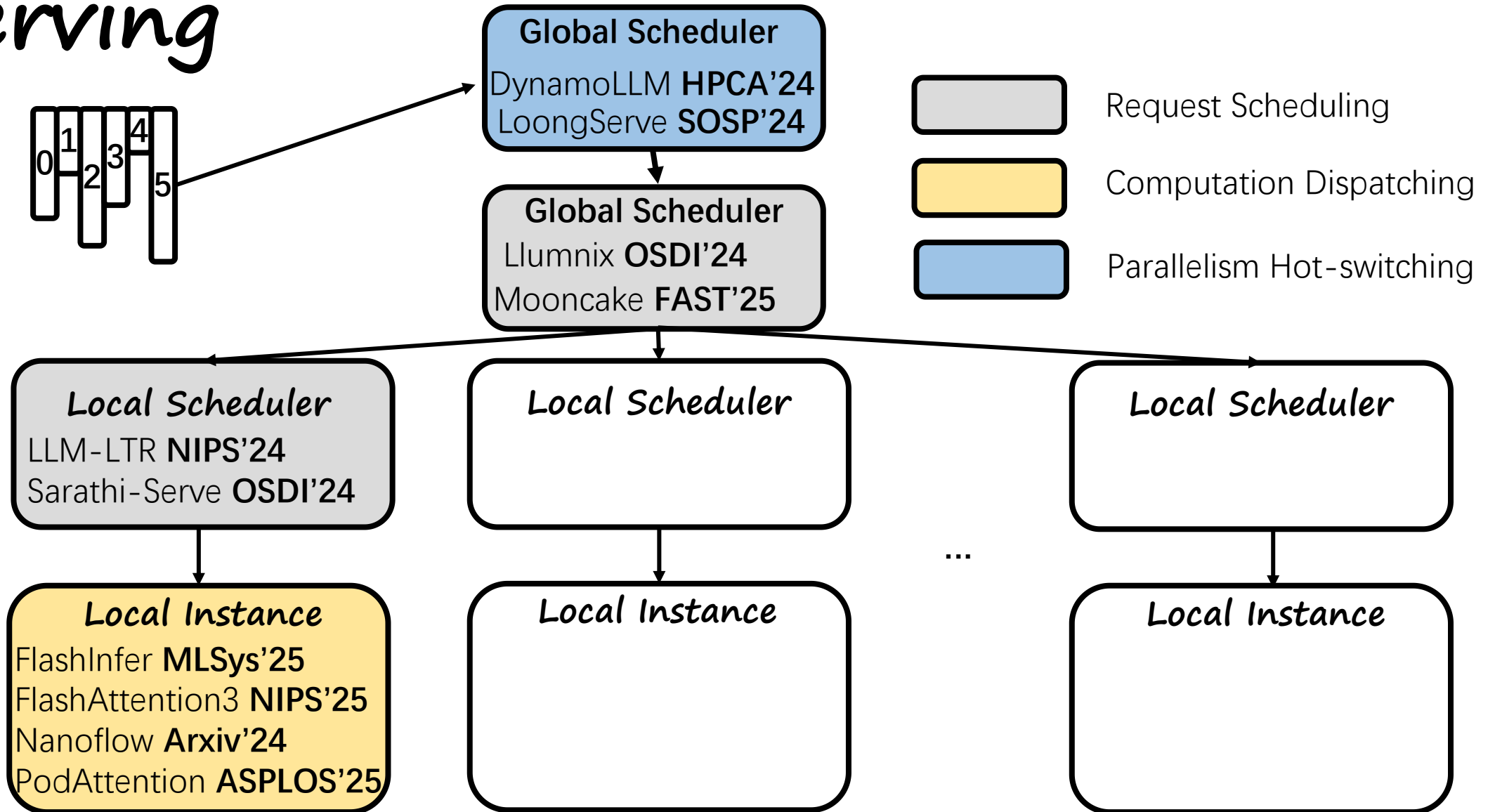
# DynamoLLM (HPCA'24 Best Paper)



Fig. 5: Example of re-sharding a TP4 model to TP2/TP8.

# Discussion: Load Balancing in LLM Serving

**Global Scheduler**

DynamoLLM **HPCA'24**
LoongServe **SOSP'24**

**Global Scheduler**

Llumnix **OSDI'24**
Mooncake **FAST'25**

*Local Scheduler*

LLM-LTR **NIPS'24**
Sarathi-Serve **OSDI'24**

*Local Instance*

FlashInfer **MLSys'25**
FlashAttention3 **NIPS'25**
Nanoflow **Arxiv'24**
PodAttention **ASPLOS'25**

*Local Scheduler*

*Local Instance*

*Local Scheduler*

*Local Instance*

...

Request Scheduling

Computation Dispatching

Parallelism Hot-switching

# Root Cause: Request Length Discrepancy

|  | Batching | Computation | Objective |  |
|---|---|---|---|---|
| LLM Training | Static | Prefill | Throughput |  |
| LLM Serving | Continuous | Prefill + Decode | SLO |  |