

DistCA: Core Attention Disaggregation

Efficient Long-context LLM Training by Core Attention Disaggregation


Yonghao Zhuang^{*1}, **Junda Chen**^{*2}, Bo Pang², Yi Gu²,
Yimin Jiang³, Yibo Zhu³, Eric Xing¹, Hao Zhang²

¹ Carnegie Mellon University, ² UC San Diego, ³ Stepfun



*support for compute



 hao-ai-lab/DistCA

Blog: <https://haoailab.com/blogs/distca/>

Long-context LLM is the Norm

Model

| | |
|----------------------|---------------|
| Anthropic Opus 4.6 | 1M |
| OpenAI GPT-5 | 1M |
| Gemini 3.1 Pro | 2M |
| Qwen3-Next | 1M (extended) |
| Meta's Llama 4 Scout | 10M |



Gemini



Qwen



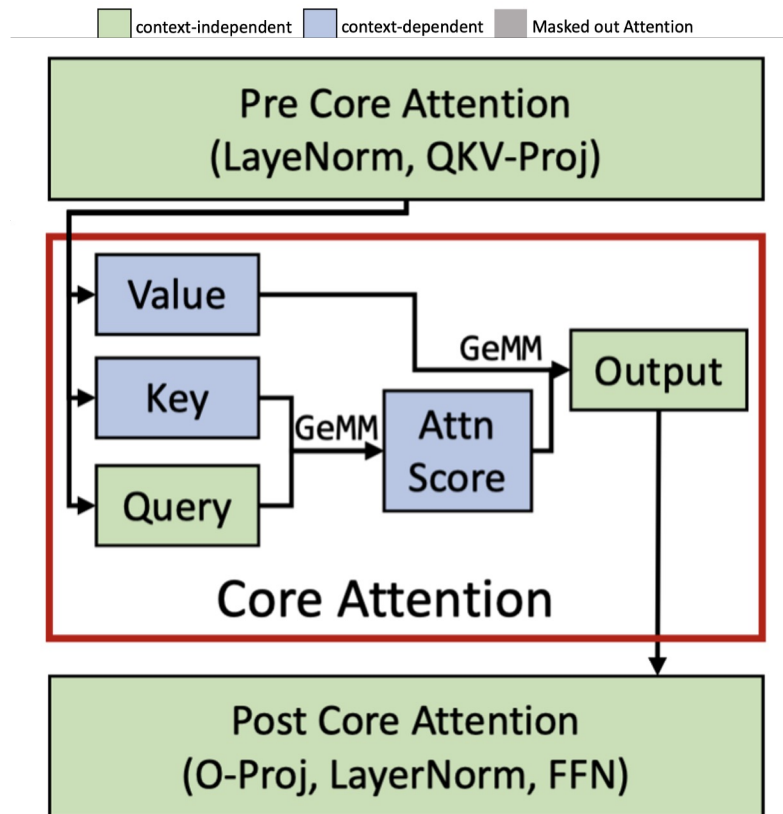
Meta

Applications



...

Background: Transformer Architecture



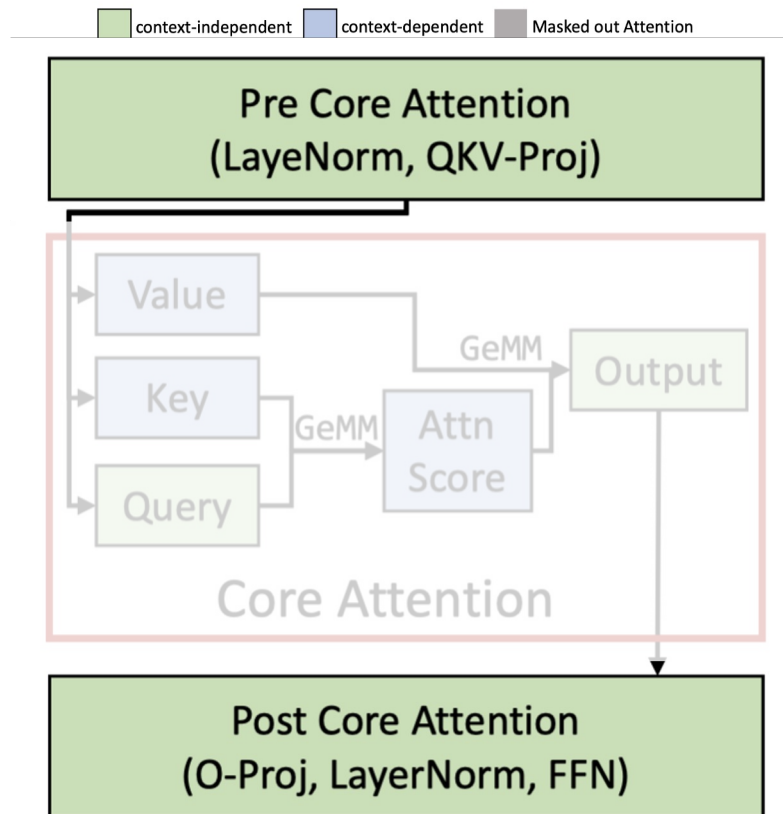
LLM Architecture

- Context Independent
 - nn.Linear
 - qkv/o-proj, FFN, Im_head, LayerNorm, RoPE
 - Stateful: model weight, activation
- Context dependent
 - Only Core Attention mechanism
 - The O = softmax(QK)V part
 - (e.g. Flash Attention kernel)
 - No Parameters. Stateless

Note:

- Core attention does not produce QKV

Background: Transformer Architecture



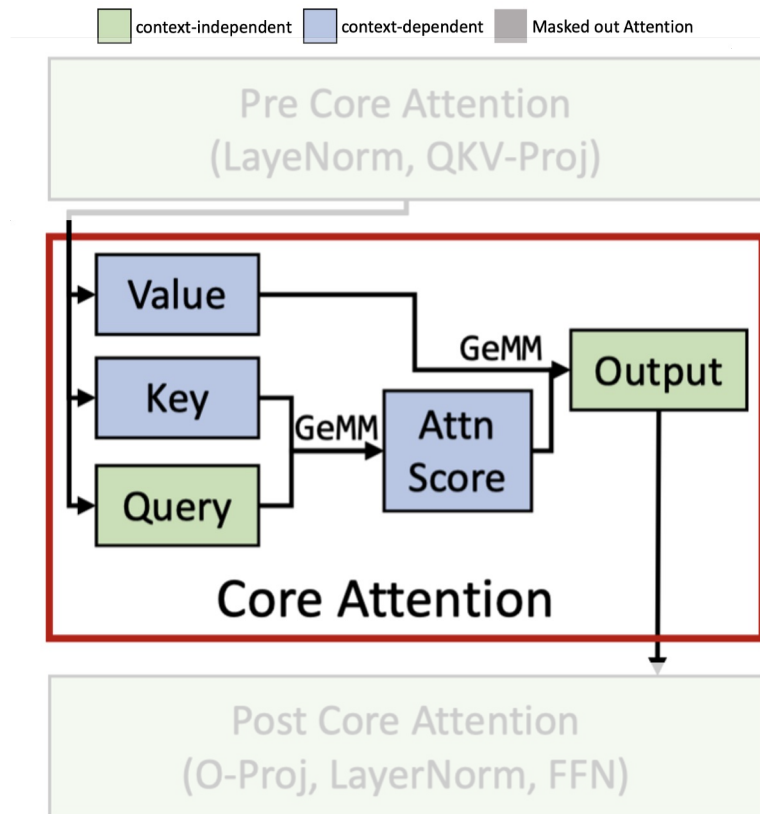
LLM Architecture

- **Context Independent (**Linear-ish parts**)**
 - nn.Linear
 - qkv/o-proj, FFN, Im_head, LayerNorm, RoPE
 - **Stateful: model weight, activation**
- **Context dependent**
 - Only Core Attention mechanism
 - The O = softmax(QK)V part
 - (e.g. Flash Attention kernel)
 - No Parameters. Stateless

Note:

- Core attention does not produce QKV

Background: Transformer Architecture



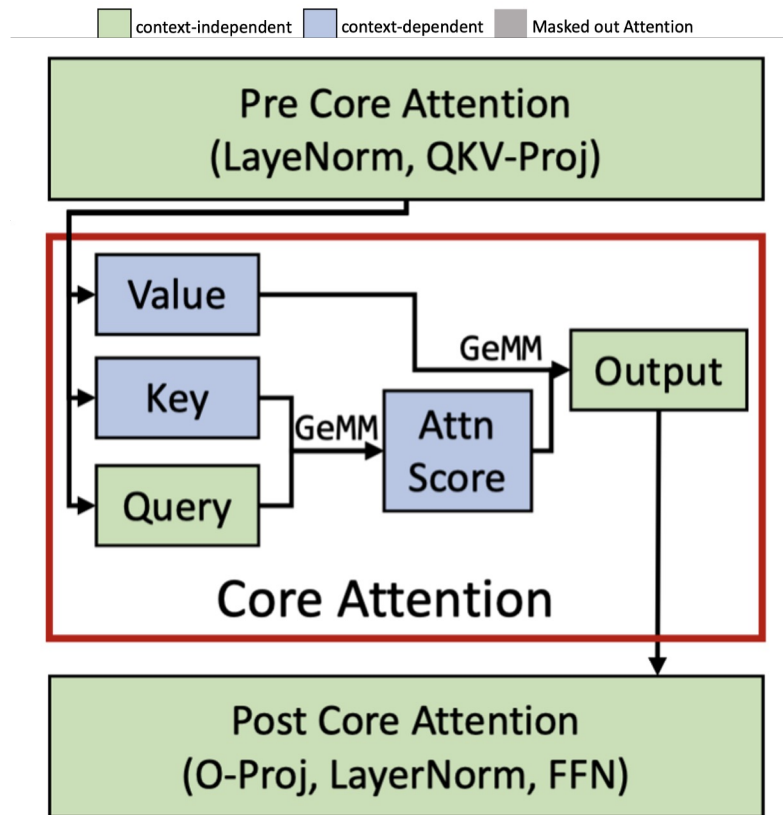
LLM Architecture




- Context Independent
 - nn.Linear
 - qkv/o-proj, FFN, Im_head, LayerNorm, RoPE
 - Stateful: model weight, activation
- Context dependent (**Core Attention**)
 - The O = softmax(QK)V part
 - (e.g. **Flash Attention** kernel)
 - No Parameters. Stateless.

Note:

- Core attention does not produce QKV

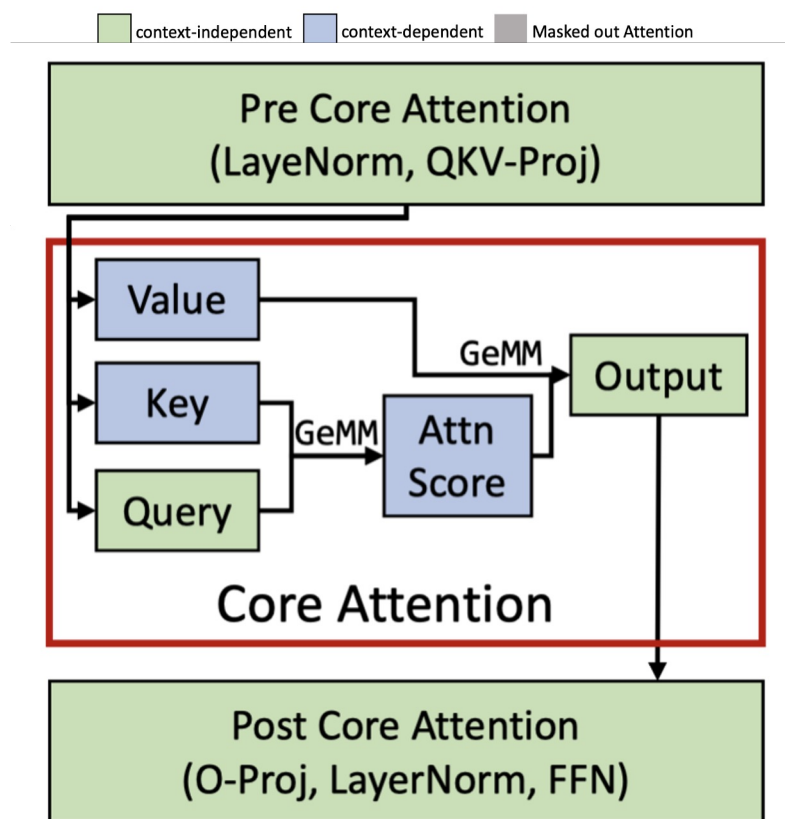
Background: Transformer Architecture (complexity)



| Layer | Memory | Compute |
|--|--|-------------|
|  CA | $\sim \text{const}$ (nearly constant) | $O(l^2)$ |
|  Linear | $O(l)$ | $O(l)$ |
|  MISC | $O(l)$ | ≈ 0 |

- CA:** core attention layer
- Linear:** FFN, qkvo-proj
- MISC:** layer norm, dropout, ...

Background: Transformer Architecture (complexity)



| Layer | Memory | Compute |
|--------|--|-------------|
| CA | $\sim \text{const}$ (nearly constant) | $O(l^2)$ |
| Linear | $O(l)$ | $O(l)$ |
| MISC | $O(l)$ | ≈ 0 |

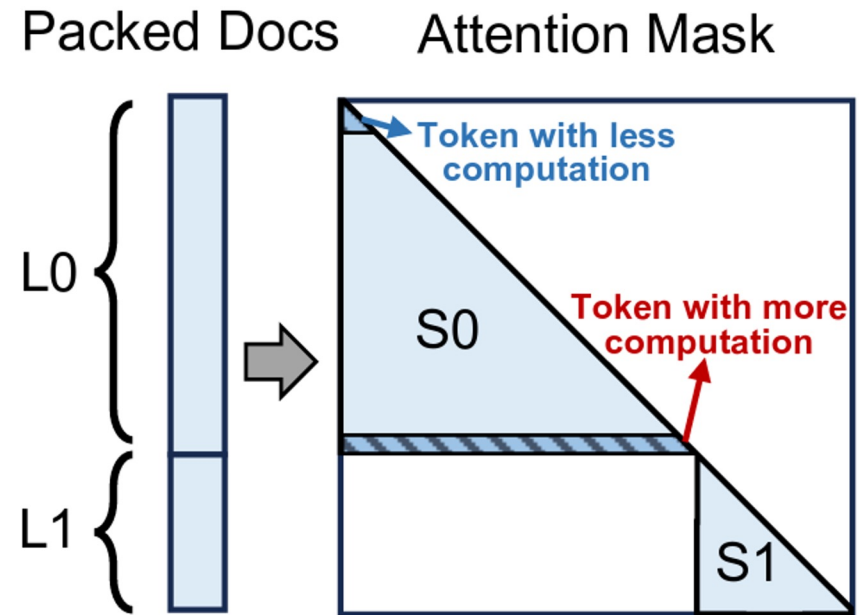
load balance

- CA:** core attention layer
- Linear:** FFN, qkvo-proj
- MISC:** layer norm, dropout, ...

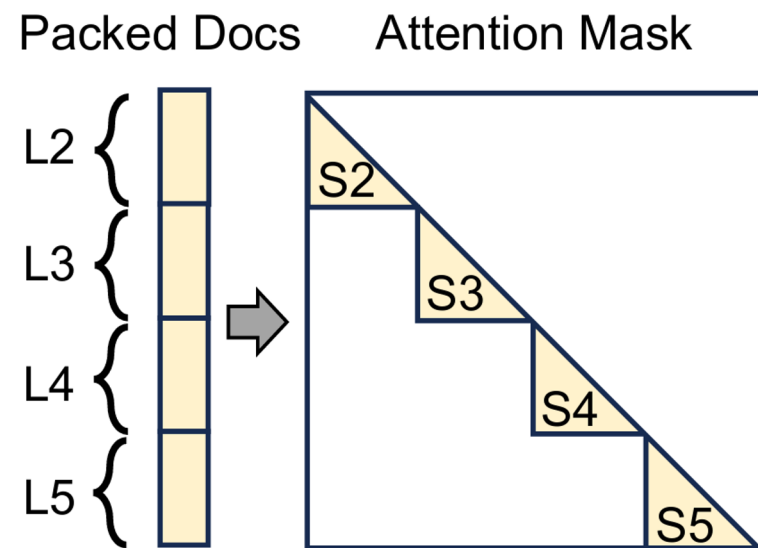
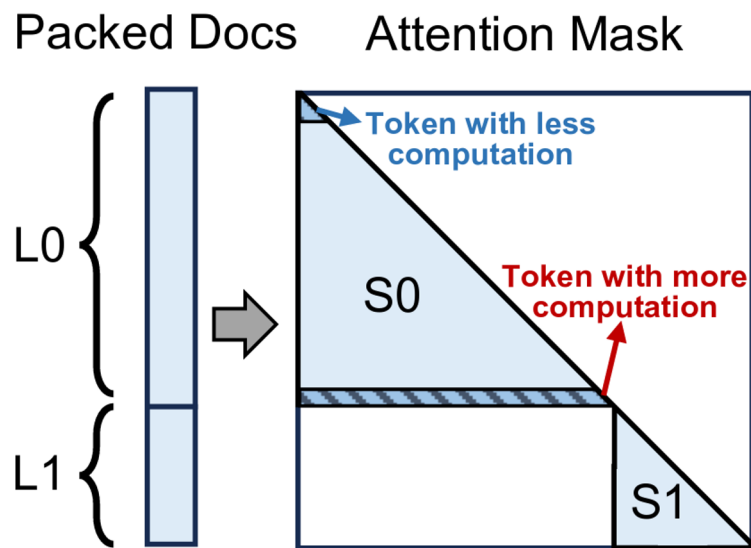
Problem: Imbalance

Data loader in LLM pretrain:

- Concat documents into multiple fixed-size chunks, each chunk is a data.
- Use a special attention mask to prevent attention computation between documents.



Problem: Imbalance



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

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

Problem: Imbalance at Scale

Data Parallel

Split data chunk

Pipeline Parallel

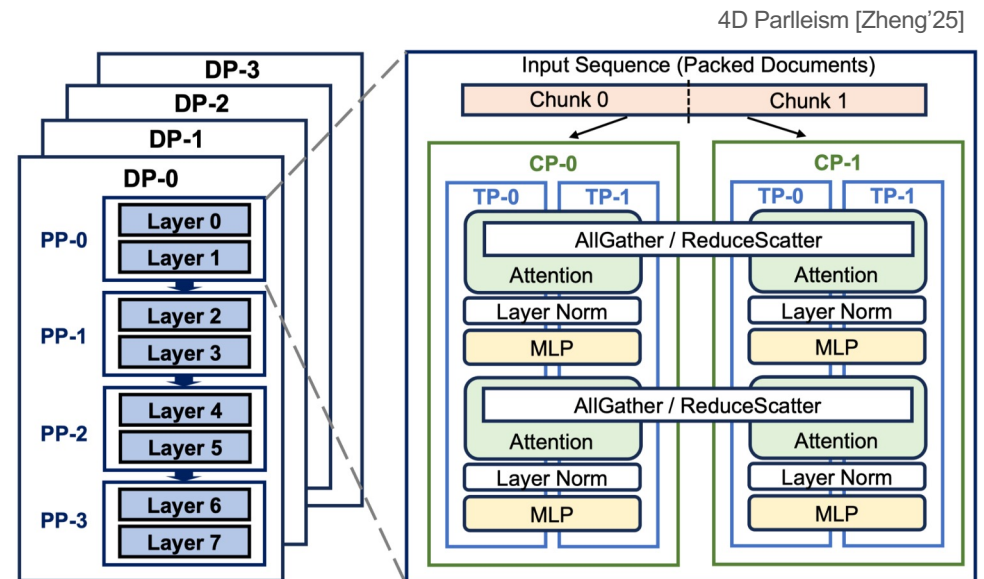
Split models weights

Tensor Parallel

Split attention heads & MLP

Context Parallel

Split sequence length dimension of input tensor



Problem: Imbalance at Scale

Data Parallel

Split data chunk

Pipeline Parallel

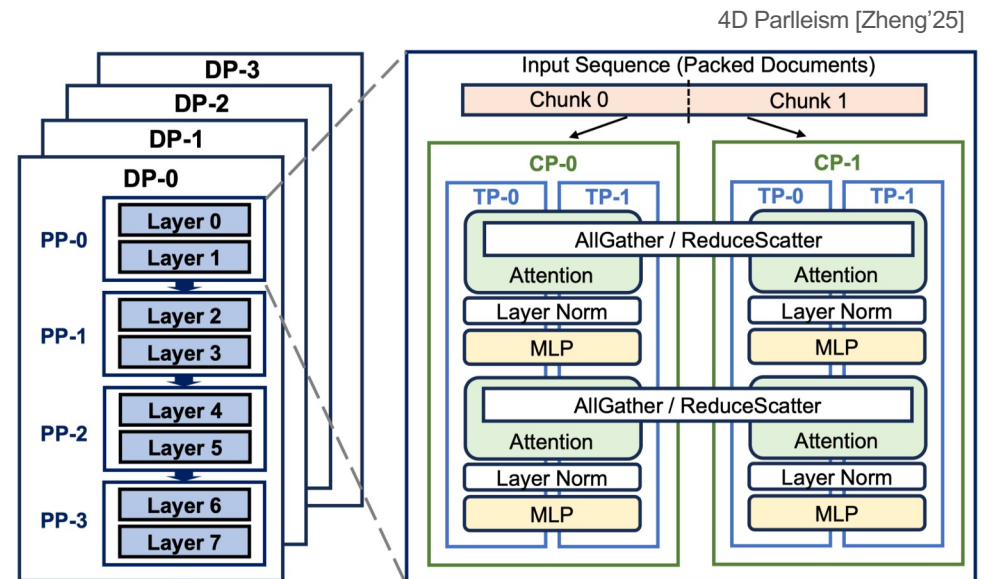
Split models weights

Tensor Parallel

Split attention heads & MLP

Context Parallel

Split sequence length dimension of input tensor

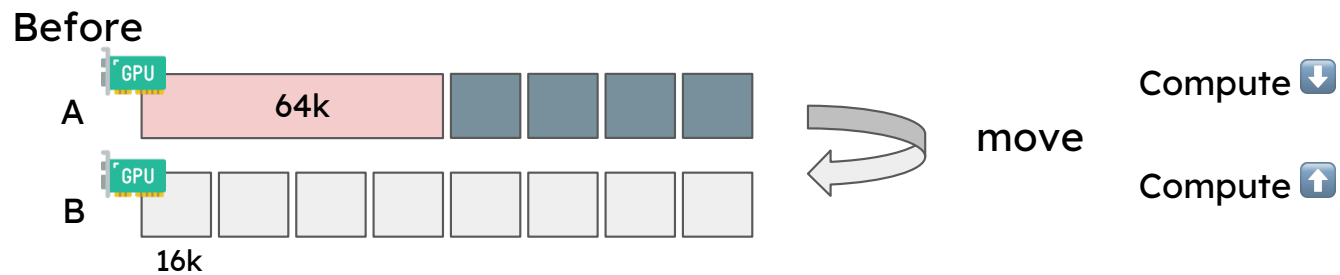


Problem: Imbalance at Scale

DP / PP

(1) Variable-len Data Chunk

- Move some documents from the heavy workload batch → light workload batch

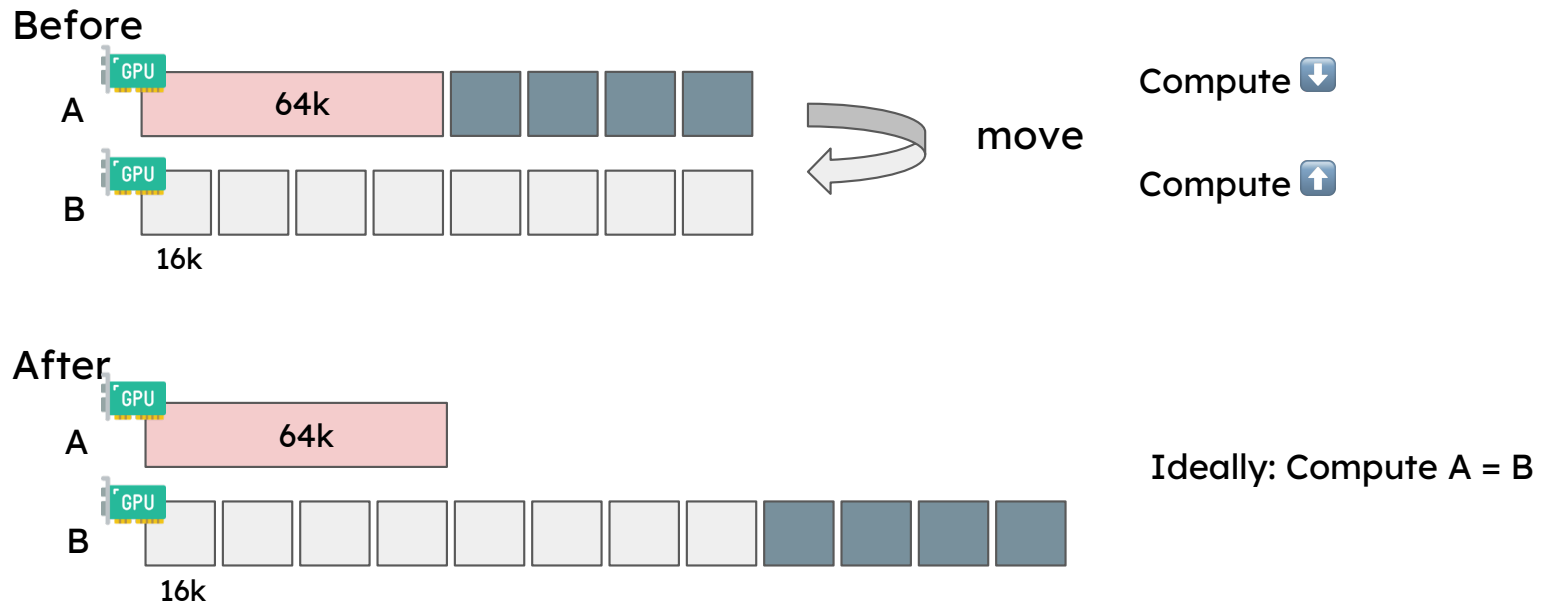


Problem: Imbalance at Scale

DP / PP

(1) Variable-len Data Chunk

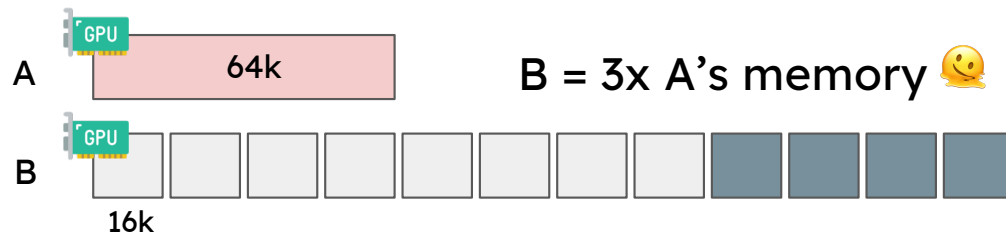
- Move some documents from the heavy workload batch → light workload batch



Problem: Imbalance at Scale

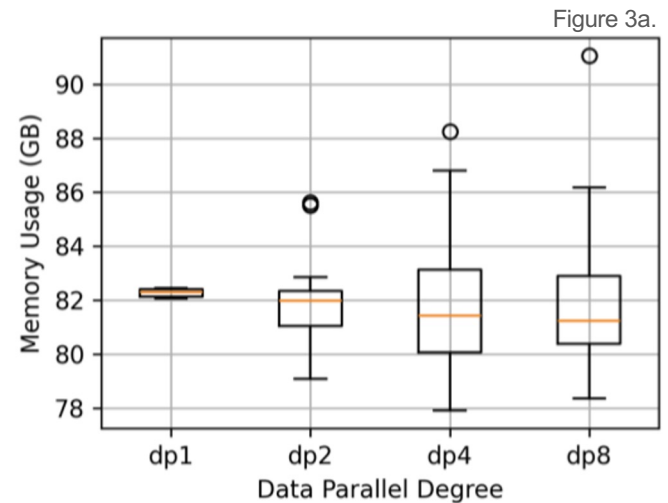
DP / PP

(1) Variable-len Data Chunk



Critical Drawbacks:

- Invitation to Memory imbalance
- Limit ability to balance with long sequence



Memory divergence can grow up to ~1.2x even in 8 nodes!

Problem: Imbalance at Scale

suffer from imbalance

Data Parallel

Split data chunk

Pipeline Parallel

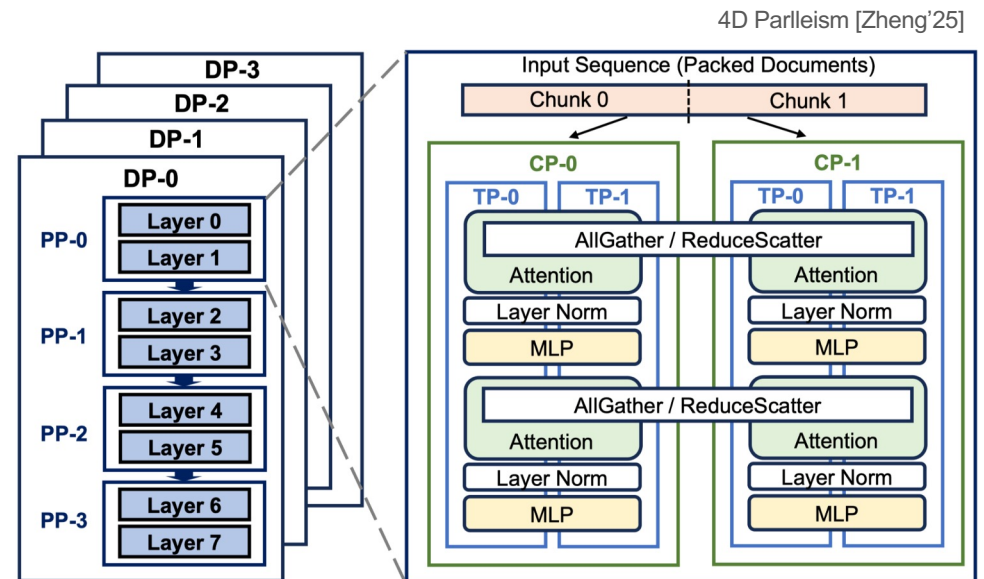
Split models weights

Tensor Parallel

Split attention heads & MLP

Context Parallel

Split sequence length dimension of input tensor
introduce overhead - does not scale well



Problem: Imbalance at Scale

suffer from imbalance

Data Parallel

Split data chunk

Pipeline Parallel

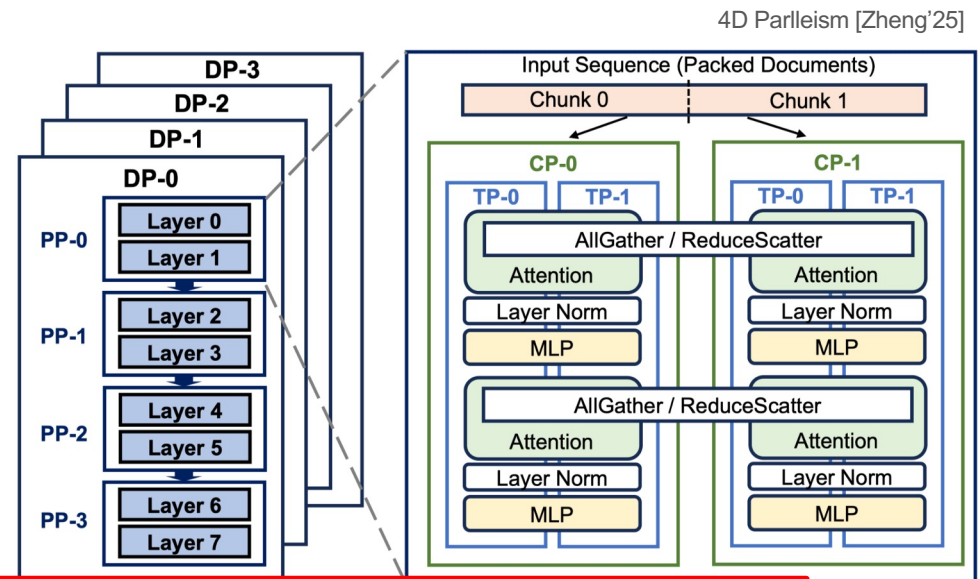
Split models weights

Tensor Parallel

Split attention heads & MLP

Context Parallel

Split sequence length dimension of input tensor

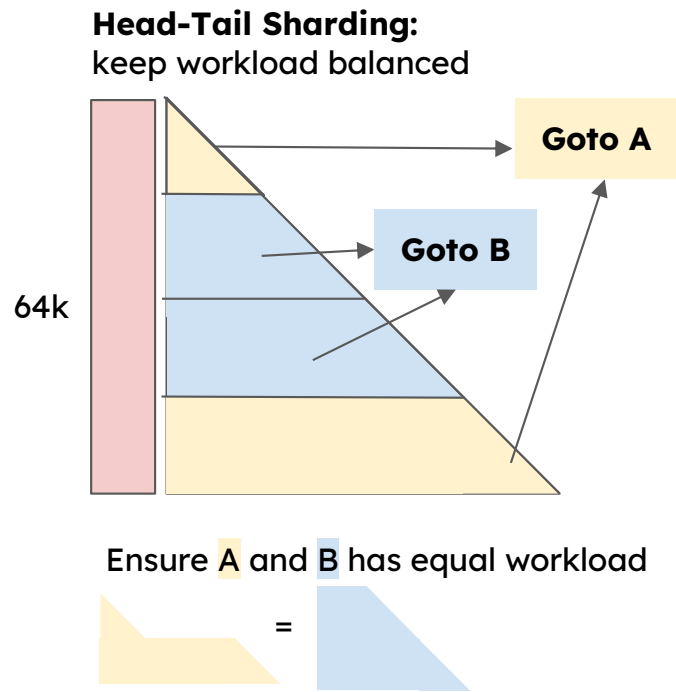
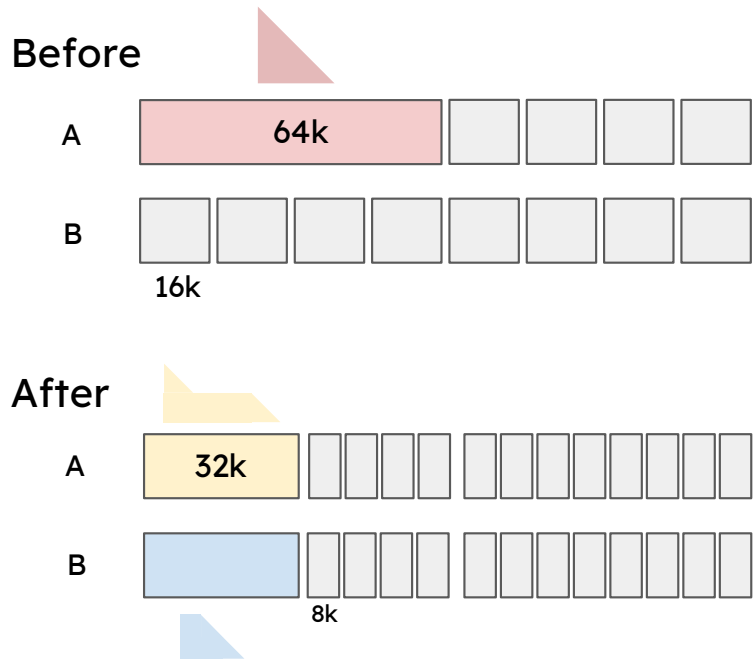


Problem: Imbalance at Scale

CP

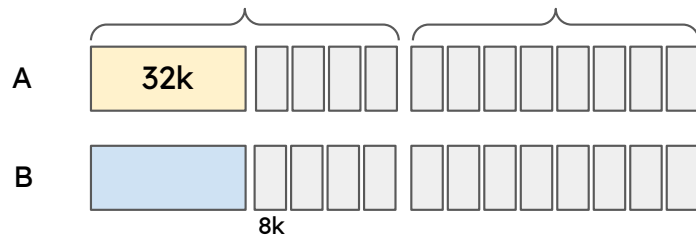
(2) Per-Doc Context Parallel

- (Head-Tail) Shard each document to different batch and make workload equal



Problem: Imbalance at Scale

(2) Per-Doc Context Parallel

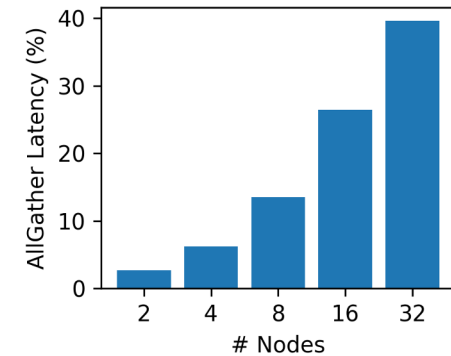


Drawbacks:

- Communication cost
 - allgather increases as the cp group size increases
- Memory footprint increases
 - All KV will store in one GPU (used for backward)

Per-Doc CP does not scale well!

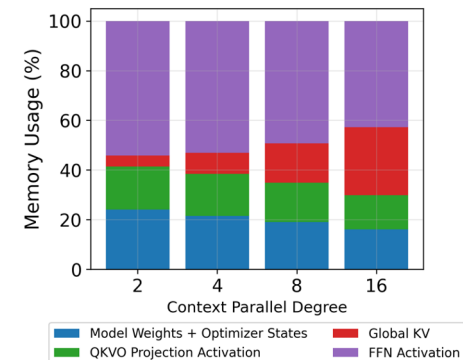
Figure 2a.



Latency

All Gather can reach to **~40%** latency (32 nodes)

Figure 2b.



Memory

All Gather buffer for KV can reach to **~20%** memory consumption (16 node)

Problem: Imbalance at Scale

suffer from imbalance

Data Parallel

Split data chunk

Pipeline Parallel

Split models weights

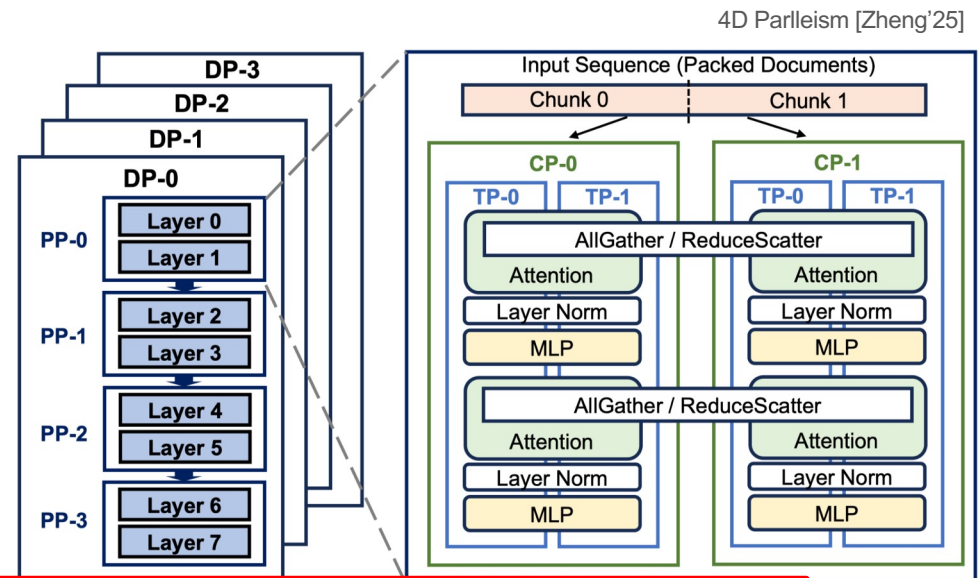
Tensor Parallel

Split attention heads & MLP

Context Parallel

Split sequence length dimension of input tensor

introduce overhead - does not scale well



Core Challenge: Load Balance

For a microbatch of documents D_i on GPU i with lengths $l(d)$ for document d :

Balancing memory (activation memory \propto total tokens):

$$\forall i, j: \sum_{d \in D_i} l(d) = \sum_{d \in D_j} l(d)$$

Balancing compute (attention is quadratic, rest is linear):

$$\forall i, j: \sum_{d \in D_i} [\alpha l^2(d) + \beta l(d)] = \sum_{d \in D_j} [\alpha l^2(d) + \beta l(d)]$$

Satisfying both constraints simultaneously is difficult in practice.

Observations

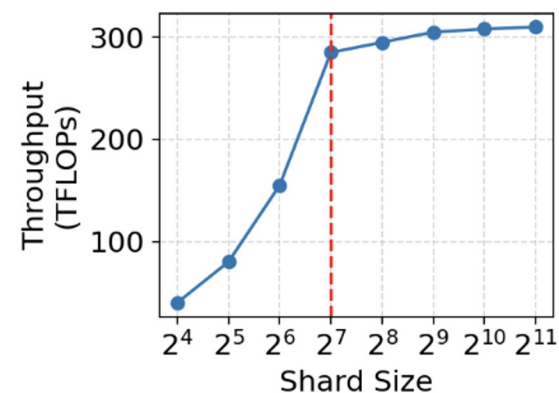
(1) CA is divisible.

FlashAttention compute CA in unit of 128-tokens

- i.e., each 128-token shard can be computed by a device, independent from devices of other shards.

128-token shards from different documents can be batched together



⇒ Shard attention into different pieces to equalize computation



Attention Throughput (TFLOPs)
Attention throughput can easily reach near peak throughput when shard size ≥ 128

Observations

- (1) CA is divisible.
- (2) CA is stateless.

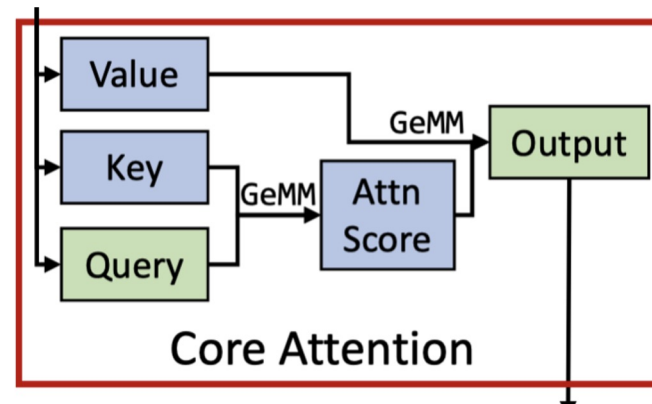
| Layer | Memory | Compute |
|--|--|----------|
|  CA | $\sim \text{const}$ (nearly constant) | $O(l^2)$ |
|  Linear | $O(l)$ stateless | $O(l)$ |

Complexity of CA vs Linear. CA memory footprint is small.

To let a device compute a shard's core attention, it only needs the input Query, Key, and Value.

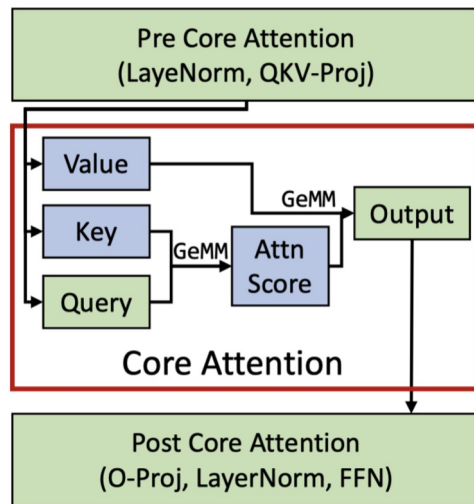
- No need for parameters
- No need to store intermediate activations

Result: small communication overhead



⇒ Sending small CA-related tensors via network is a manageable cost

DistCA: System Design



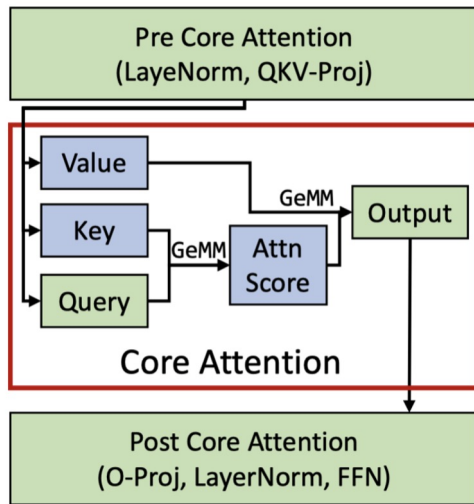
Coupled resource allocation
for both core-attn and others

- (1) CA is divisible.
- (2) CA is stateless.

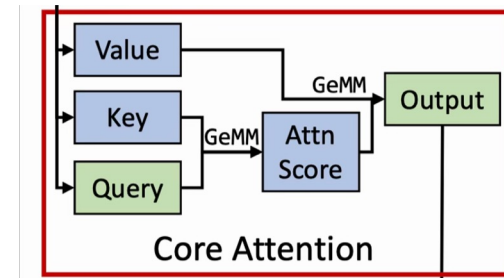
Can we balance
Core Attention independently
from other comps?

- Balance CA compute independently ?
- Balance MLP memory independently ?
- Without introducing more overhead ?

DistCA: System Design



Coupled resource allocation
for both core-attn and others

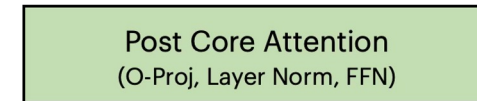
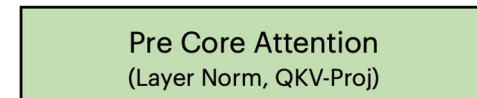
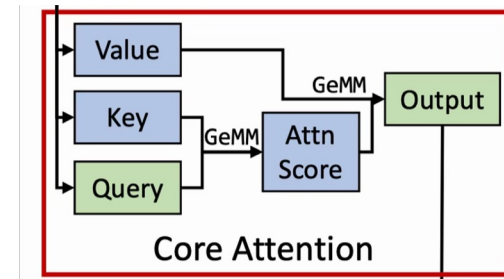


Decoupled resource allocation
disaggregate CA from others

DistCA: System Design

We independently consider how to:

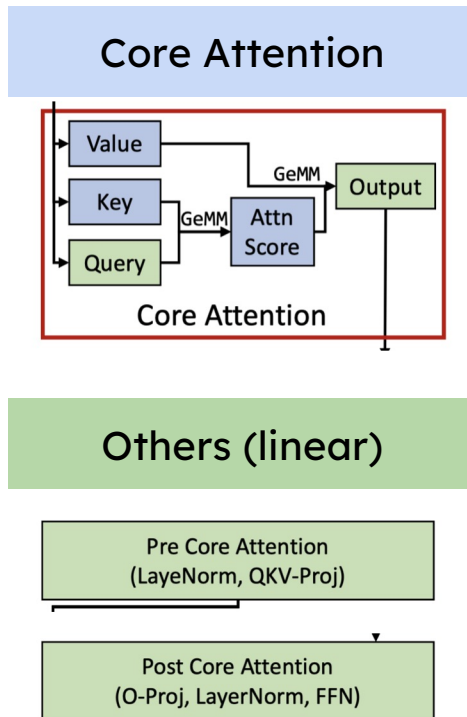
- **Balance Core Attention,**
 - ~ No memory constraint
 - $O(n^2)$ compute
- **Balance Other Parts**
 - Context-Independent
 - $O(n)$ compute and memory cost
- **Disaggregate:**
 - Communicate between 2 stages
 - Reduce network cost (prove to be small)



Decoupled resource allocation
disaggregate CA from others

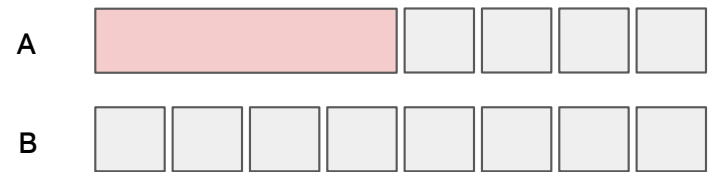
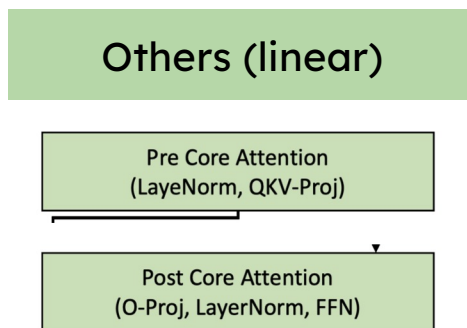
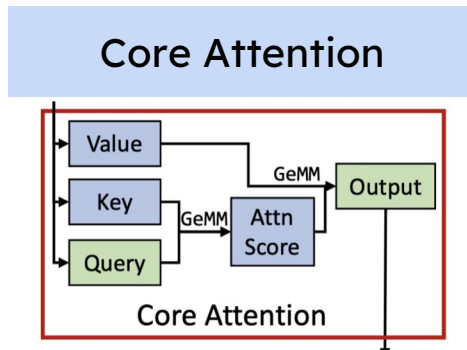
DistCA

(1) Disaggregate + Scheduler to balance CA



DistCA

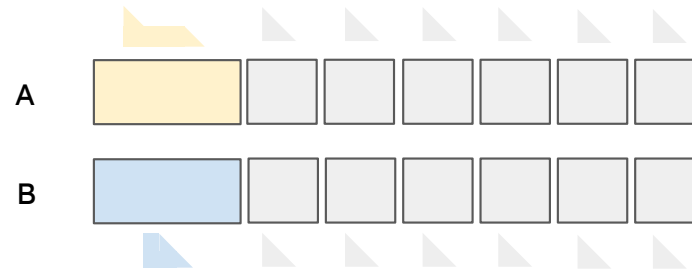
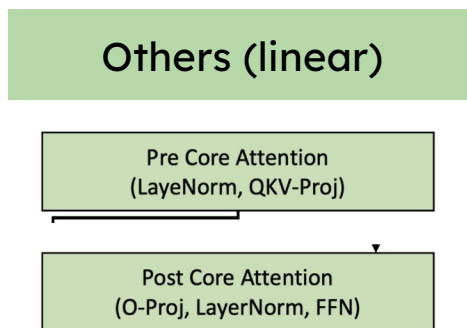
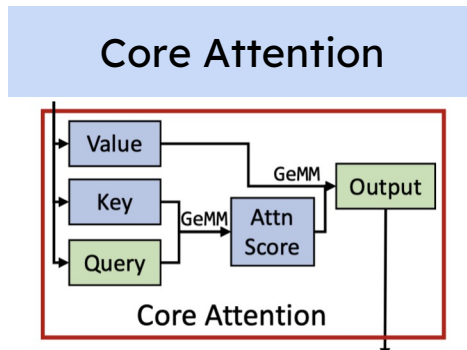
(1) Disaggregate + Scheduler to balance CA




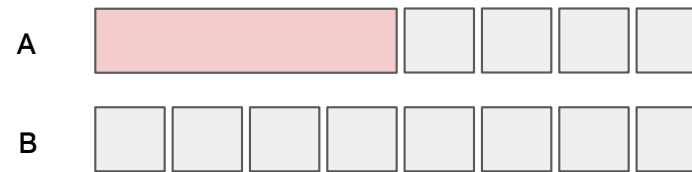
- ✓ Compute balanced
- ✓ Memory balanced

DistCA

(1) Disaggregate + Scheduler to balance CA



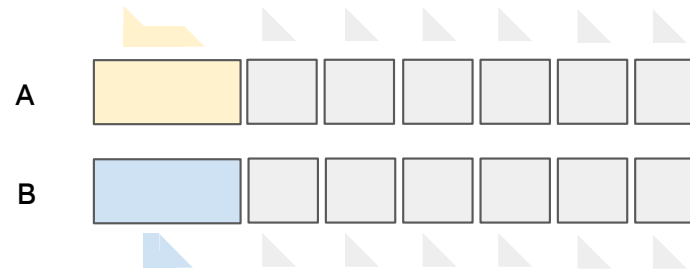
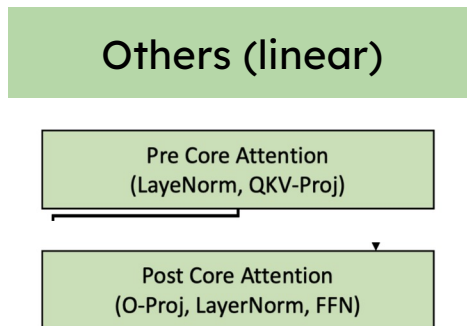
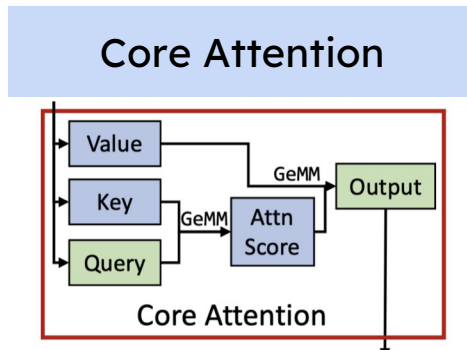
 Scheduler
• Balance compute



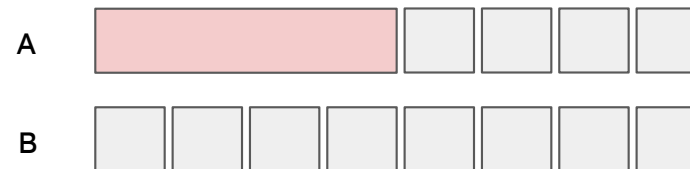
✓ Compute balanced
✓ Memory balanced

DistCA

(1) Disaggregate + Scheduler to balance CA



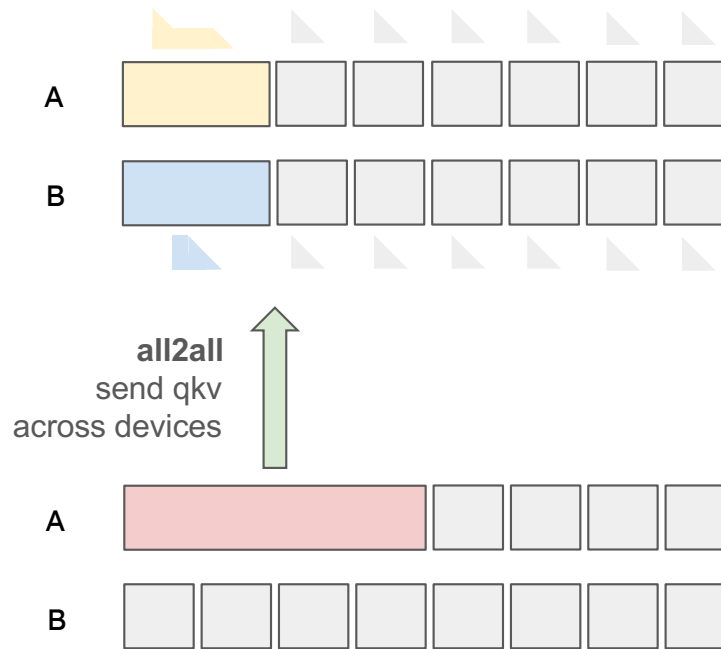
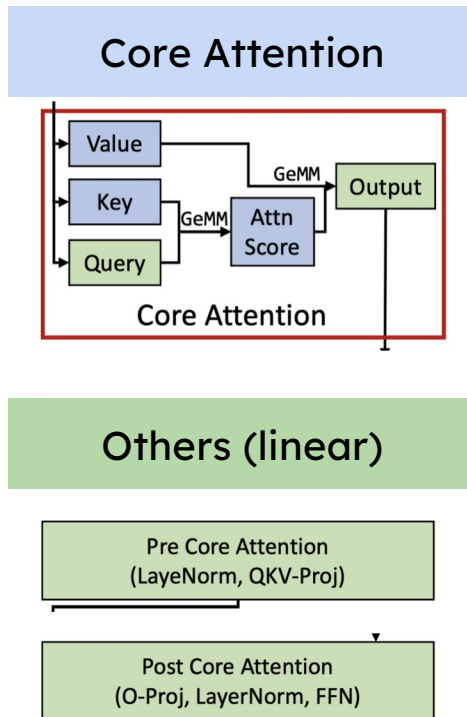
- Scheduler
- Balance compute
 - ✓ Compute balanced




- ✓ Compute balanced
- ✓ Memory balanced

DistCA

(1) Disaggregate + Scheduler to balance CA

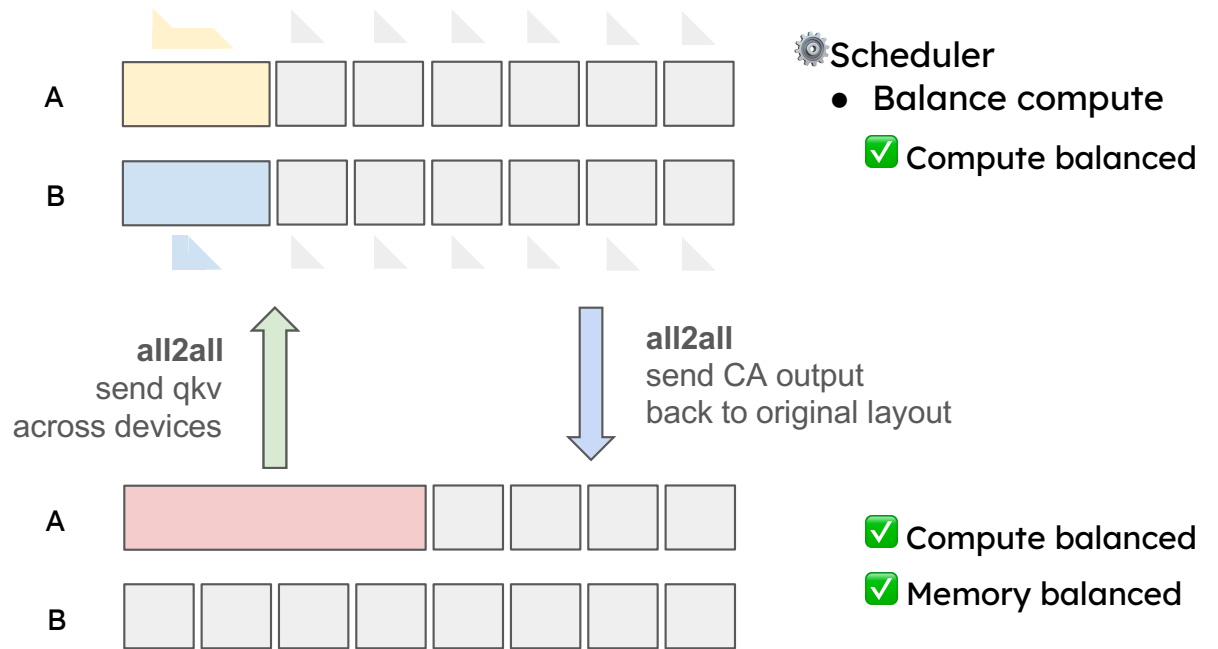
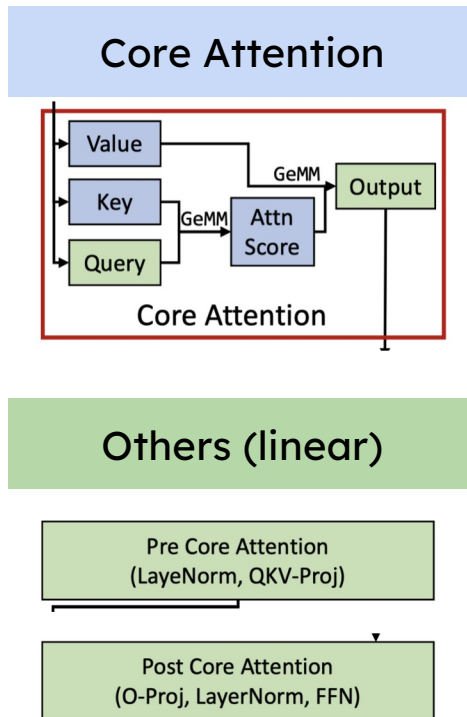


-  Scheduler
- Balance compute
 - ✓ Compute balanced

- ✓ Compute balanced
- ✓ Memory balanced

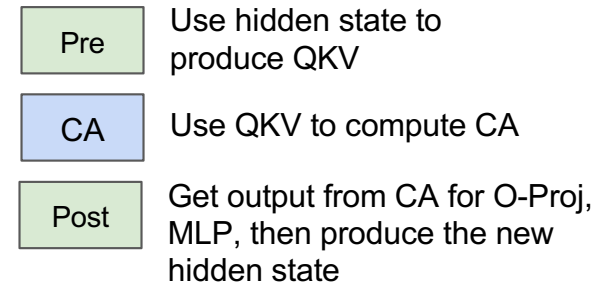
DistCA

(1) Disaggregate + Scheduler to balance CA

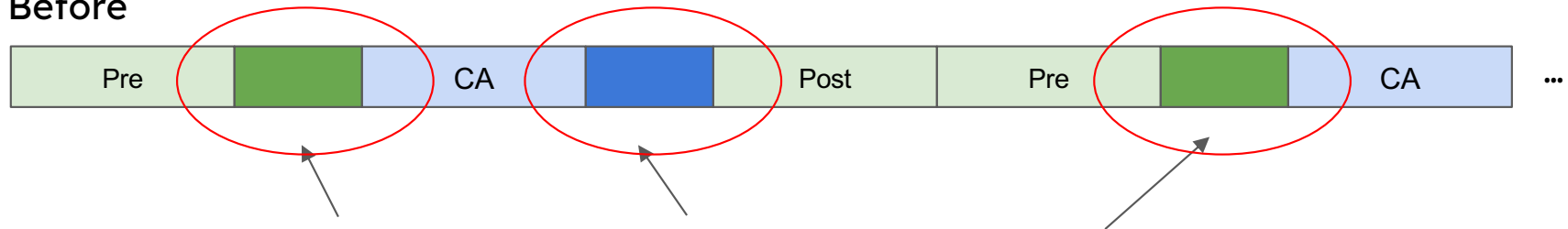


DistCA

- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost



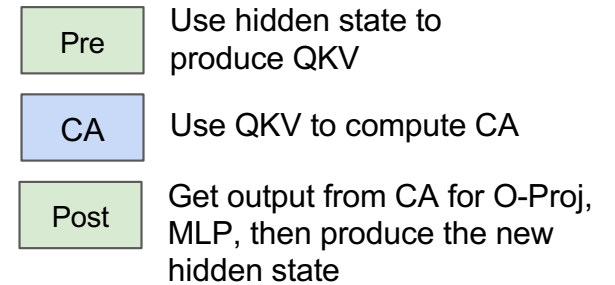
Before



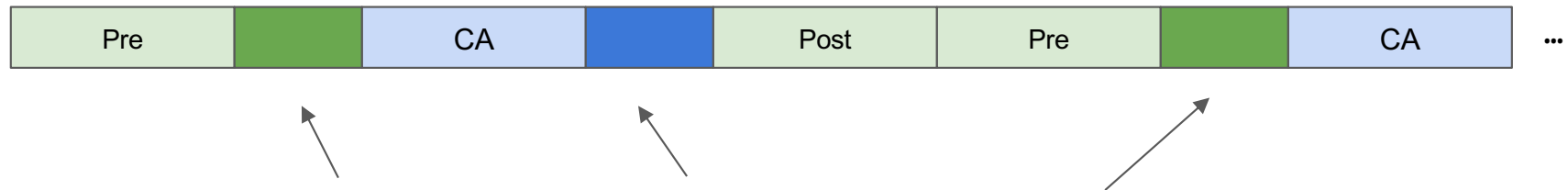
Communication overhead to send QKV / CA output across devices

DistCA

- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost



Before

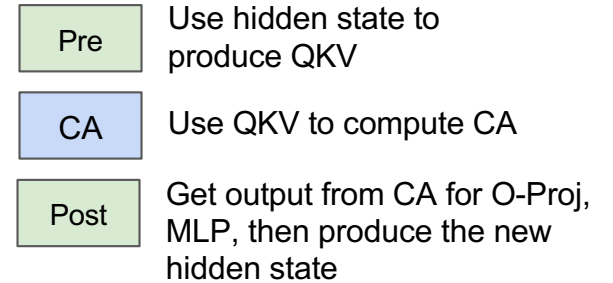


Communication overhead to send QKV / CA output across devices

💡 LLM training we have multiple micro-batches (usually not just one)

DistCA

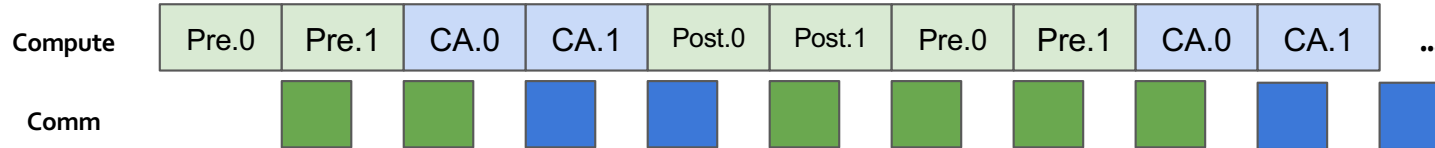
- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost



Before



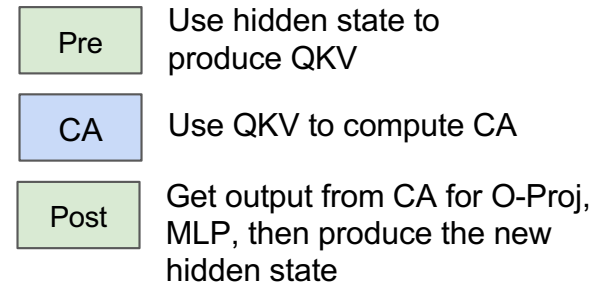
After



Overlap computation with communication (ping-pong scheduling 2 microbatch)

DistCA

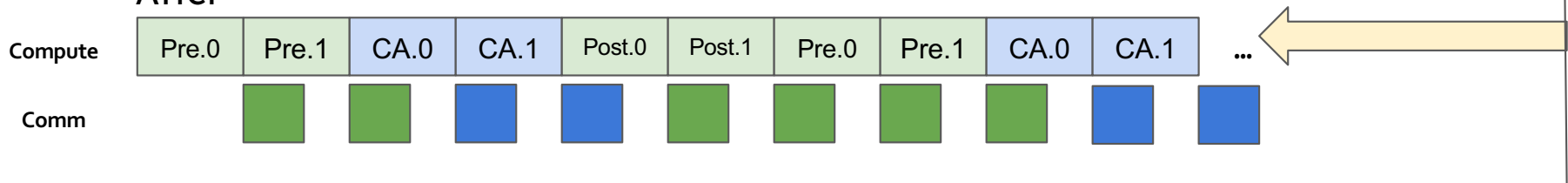
- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost



Before



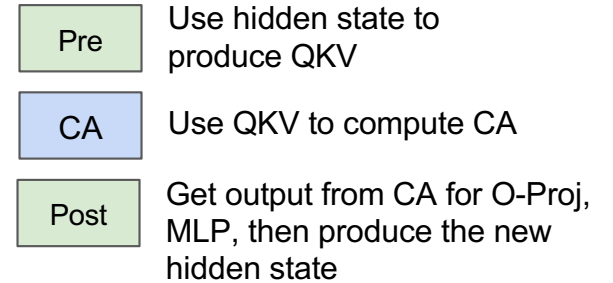
After



Overlap computation with communication (ping-pong scheduling 2 microbatch)

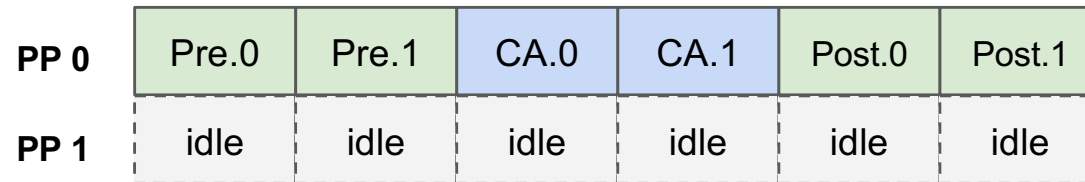
DistCA

- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost
- (3) Exploit Pipeline Bubbles



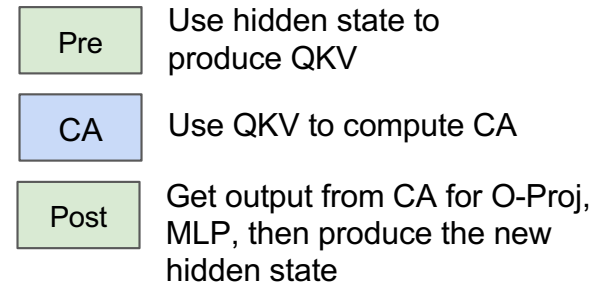
(1) Utilize Idle GPU across PP Ranks

✗ Other systems: Idle across PP ranks



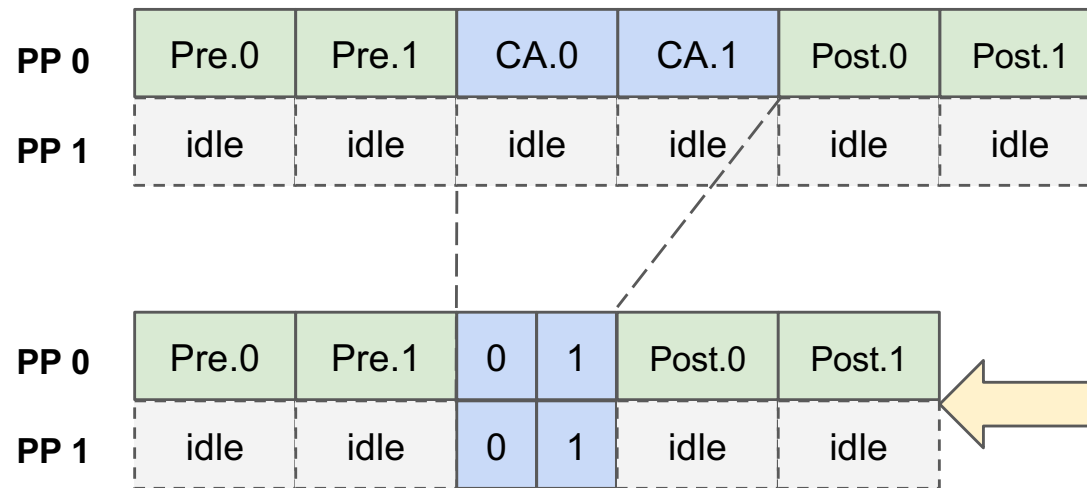
DistCA

- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost
- (3) Exploit Pipeline Bubbles



(1) Utilize Idle GPU across PP Ranks

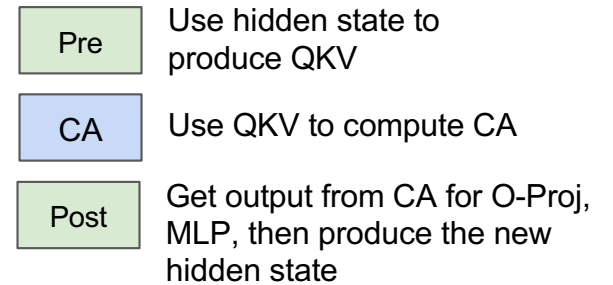
✗ Other systems: Idle across PP ranks



✓ DistCA: idle GPU share the CA workload

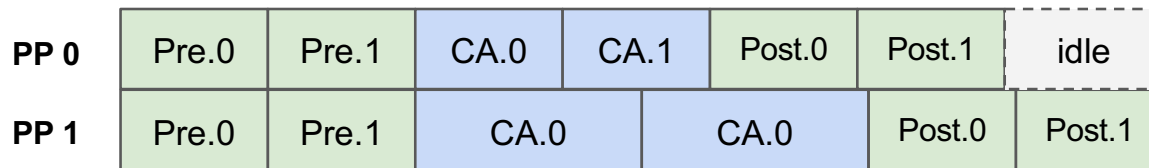
DistCA

- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost
- (3) Exploit Pipeline Bubbles



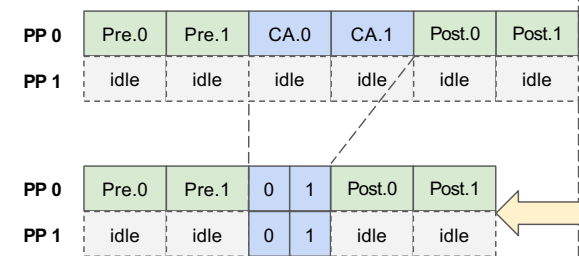
(2) Eliminate imbalance CA across PP

✗ Other systems: CA-imbalance causing bubble



(1) Utilize Idle GPU across PP Ranks

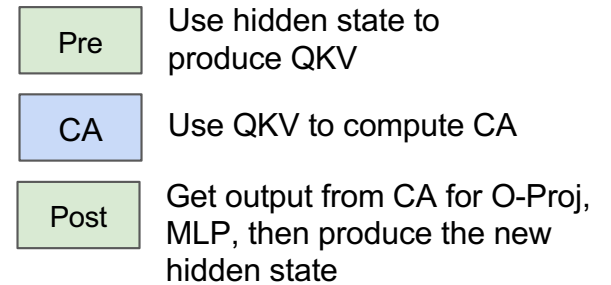
✗ Other systems: Idle across PP ranks



✔ DistCA: idle GPU share the CA workload

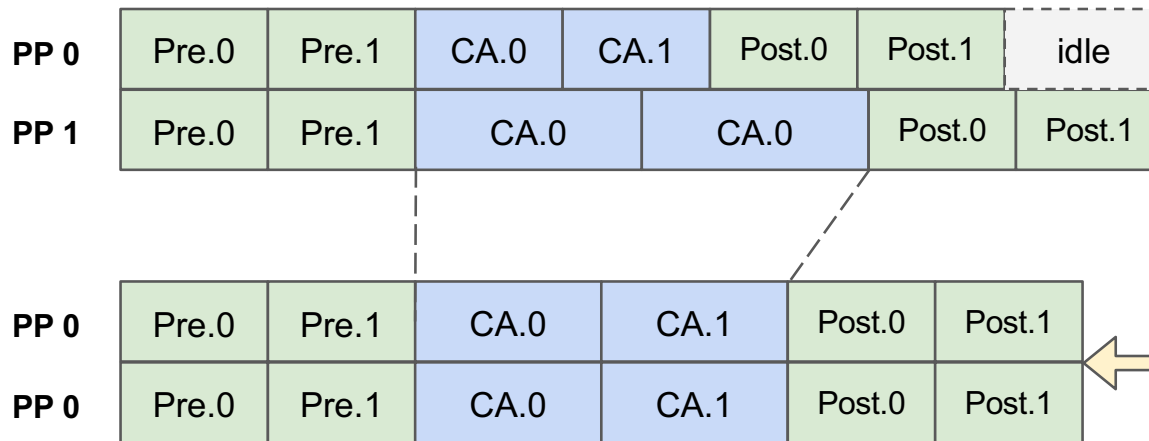
DistCA

- (1) Disaggregate + Scheduler to balance CA
- (2) Ping-Pong Schedule: hide communication cost
- (3) Exploit Pipeline Bubbles



(2) Eliminate imbalance CA across PP

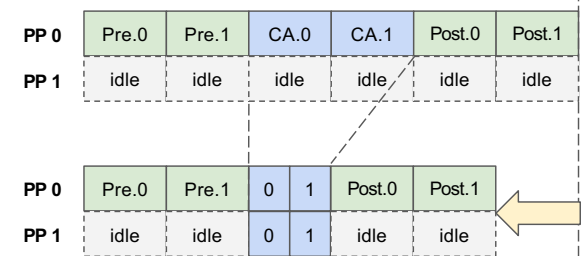
✗ Other systems: CA-imbalance causing bubble



✓ DistCA: CA and non-CA is balanced

(1) Utilize Idle GPU across PP Ranks

✗ Other systems: Idle across PP ranks

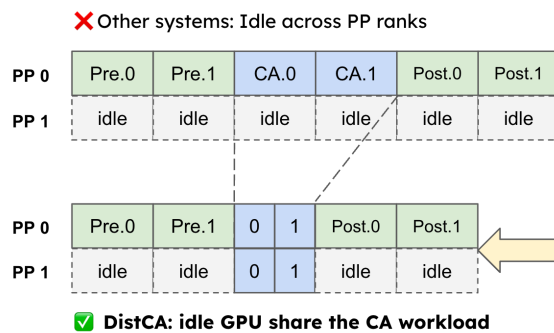


✓ DistCA: idle GPU share the CA workload

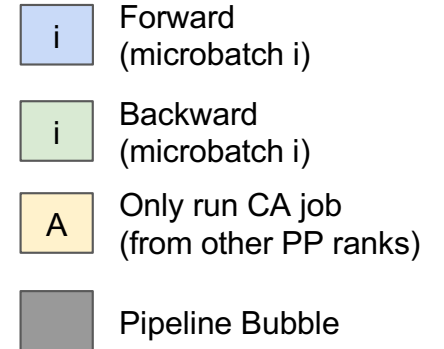
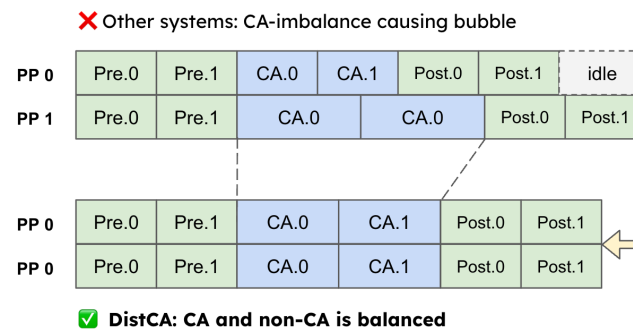
DistCA

(3) Exploit Pipeline Bubbles

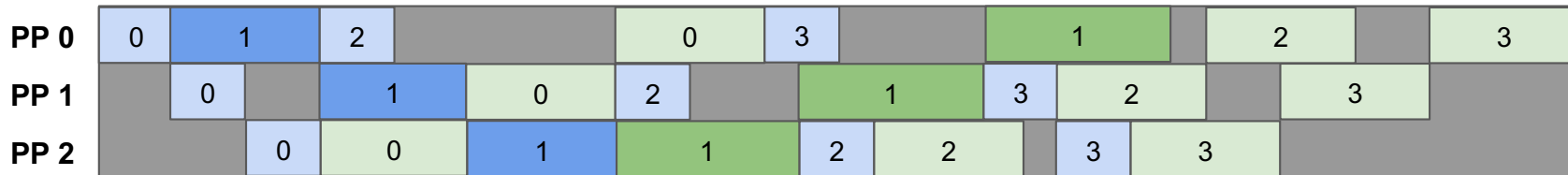
(1) Utilize Idle GPU across PP Ranks



(2) Eliminate imbalance CA across PP



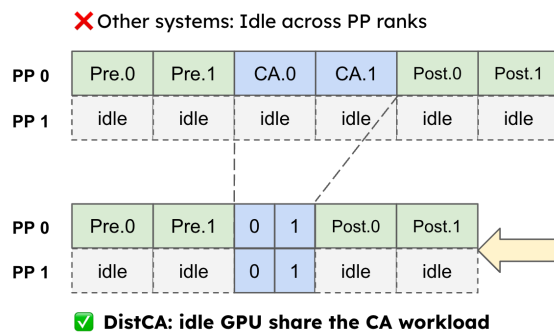
✗ Other systems



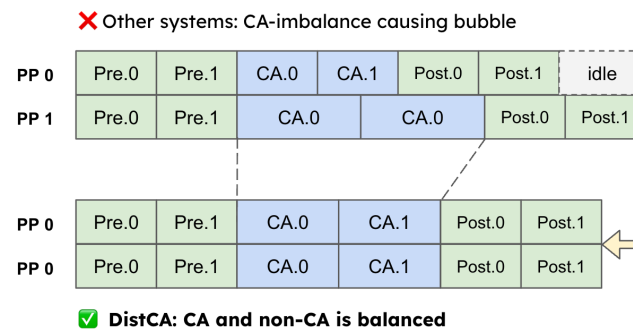
DistCA

(3) Exploit Pipeline Bubbles

(1) Utilize Idle GPU across PP Ranks

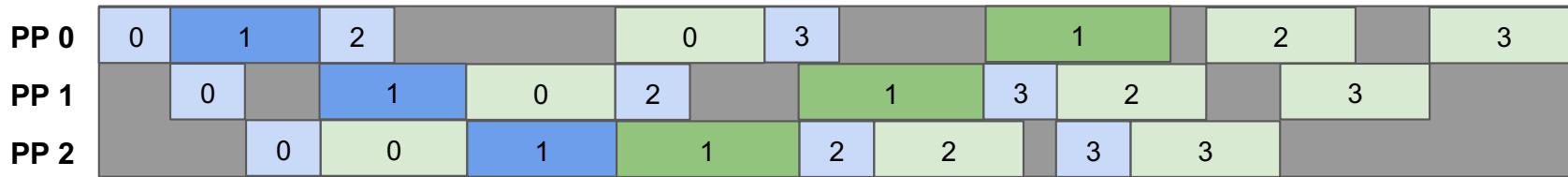


(2) Eliminate imbalance CA across PP

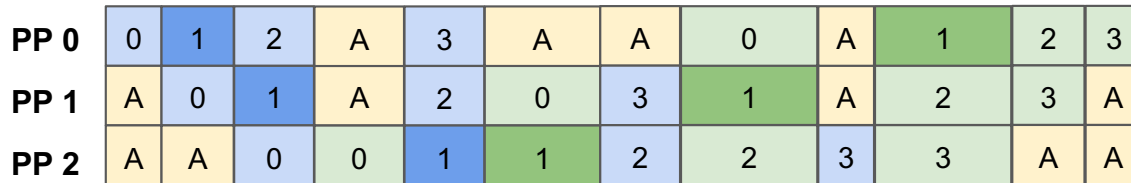


- i Forward (microbatch i)
- i Backward (microbatch i)
- A Only run CA job (from other PP ranks)
- Pipeline Bubble

✗ Other systems



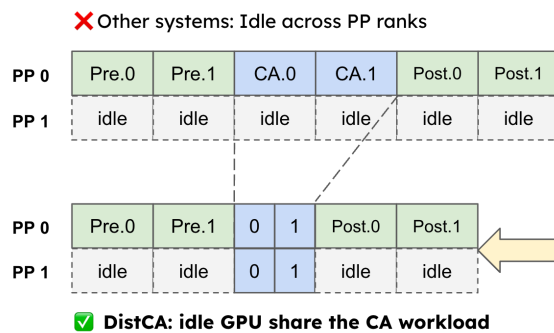
✓ DistCA



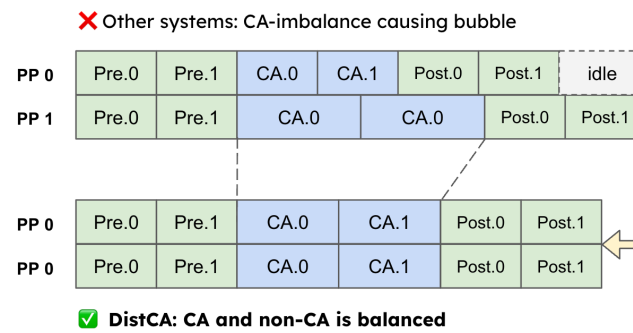
DistCA

(3) Exploit Pipeline Bubbles

(1) Utilize Idle GPU across PP Ranks

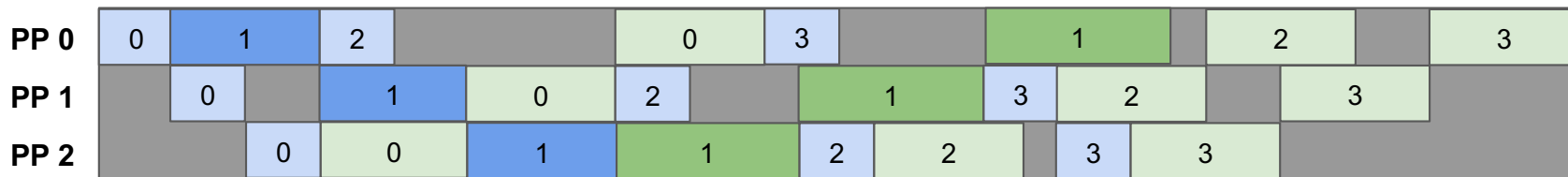


(2) Eliminate imbalance CA across PP

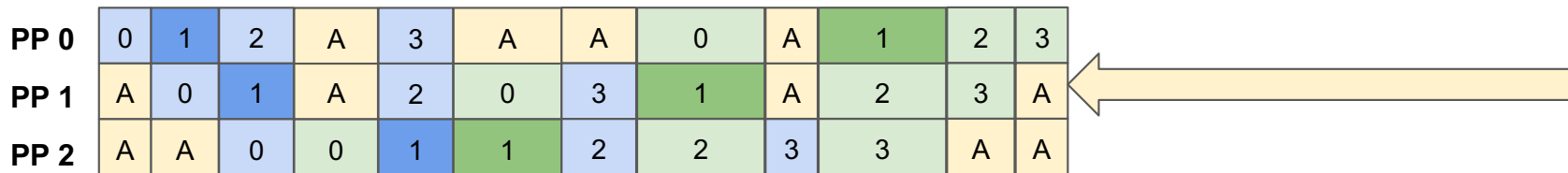


- i Forward (microbatch i)
- i Backward (microbatch i)
- A Only run CA job (from other PP ranks)
- Pipeline Bubble

✗ Other systems



✓ DistCA



DistCA: Evaluation

- 8~64x H200 Node (up to 512 GPUs)
- NVLink for intra-node comm, IB for inter-node comm
- Model: Llama-7B, 34B
- SeqLen: 128k, 256k, 512k (384k for 4D parallel)
- Baselines: Megatron-LM, WLB-LLM, FlexSP

Table 3. 3D Training Configurations.

| Model | MaxDocLen | Batch Size | #GPU |
|-----------|-----------|------------|--------------|
| Llama-8B | 128K | 8, 16, 32 | 64, 128, 256 |
| Llama-8B | 256K | 4, 8, 16 | 64, 128, 256 |
| Llama-8B | 512K | 2, 4, 8 | 64, 128, 256 |
| Llama-34B | 128K | 4, 8, 16 | 64, 128, 256 |
| Llama-34B | 256K | 2, 4, 8 | 64, 128, 256 |
| Llama-34B | 512K | 2, 4, 8 | 64, 128, 256 |

Table 4. 4D Parallel Training Configurations.

| Model | MaxDocLen | Batch Size | #GPU |
|-----------|-----------|-------------|---------------|
| Llama-8B | 128K | 32, 64, 128 | 64, 128, 256 |
| Llama-8B | 256K | 16, 32, 32 | 64, 128, 256 |
| Llama-8B | 512K | 8, 8, 16 | 64, 128, 256 |
| Llama-34B | 128K | 32, 64, 128 | 128, 256, 512 |
| Llama-34B | 256K | 16, 32, 32 | 128, 256, 512 |
| Llama-34B | 384K | 8, 8, 16 | 128, 256, 512 |



DistCA: Core Attention Disaggregation

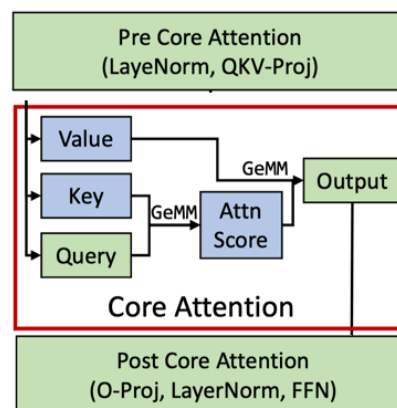
Efficient Long-context LLM Training by Core Attention Disaggregation

Yonghao Zhuang*¹, **Junda Chen***², Bo Pang², Yi Gu², Yimin Jiang³, Yibo Zhu³, Eric Xing¹, Hao Zhang²

¹ Carnegie Mellon University, ² UC San Diego, ³ Stepfun



DistCA



Existing training systems
colocates
Linear & Quadratic Compute

- LLM training with long context suffer from **workload imbalance**
- **DistCA** disaggregates CA from others \Rightarrow Balance compute and memory
- Achieves up to **1.9x** over Megatron-LM, **1.35x** speedup over WLB-LLM
- Scales better at larger cluster sizes and larger context length

Email: junda@ucsd.edu

Backup Slides

Experiment

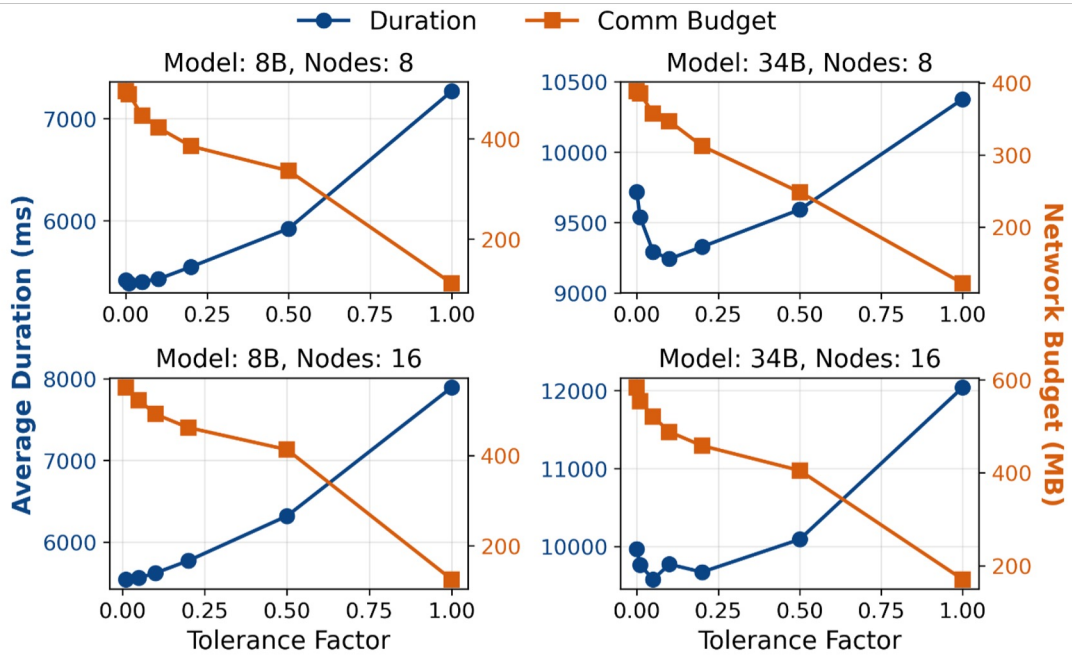


Figure 12. Impact of the compute imbalance tolerance factor.

Experiment

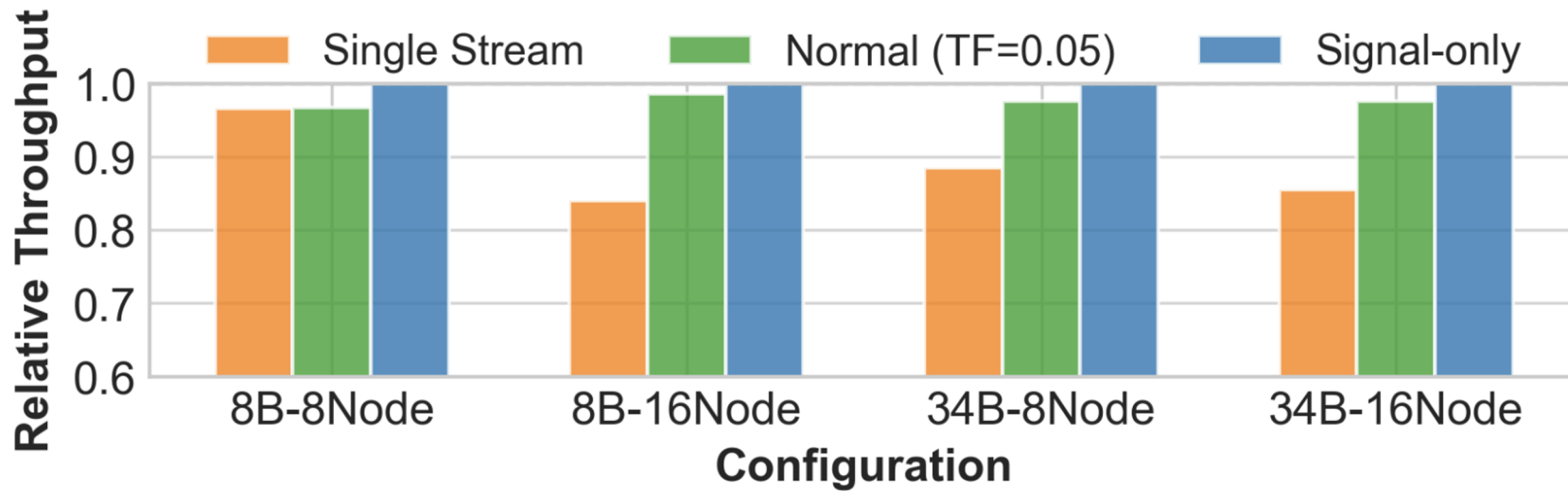


Figure 11. Throughput for different communication patterns.