

BLADE: Boosting LLM Decoding's Communication Efficiency in DRAM-based PIM

Yilong Zhao (Speaker), Fangxin Liu, Zongwu Wang, Mingjian Li, Mingxing Zhang, Chixiao Chen and Li Jiang

Shanghai Jiao Tong University, Shanghai Qi Zhi Institute,
Tsinghua University, Fudan University

ASP-DAC 2026

饮水思源 · 爱国荣校



🏛️ **Context:** Large Language Model (LLM) inference, **Prefill-Decoding Disaggregation**

- Prefill- \rightarrow GPU; Decoding \rightarrow DRAM-based **Processing-in-Memory (PIM)**

🏛️ **Challenges on Communication & Solutions:**

- Challenge 1: **High KV Cache Transpose Latency** During “GPU- \rightarrow PIM” KV Cache Transfer
- Solution 1: **Transpose on Transfer Method**
- Challenge 2: Static PIM Parallelism **Fails to Balance Communication & Computation**
- Solution 2: **Dynamical Parallelism Scaling Method**

🏛️ **Results:**

- Transpose on Transfer 1.54x; Dynamical Parallelism Scaling 2.09x

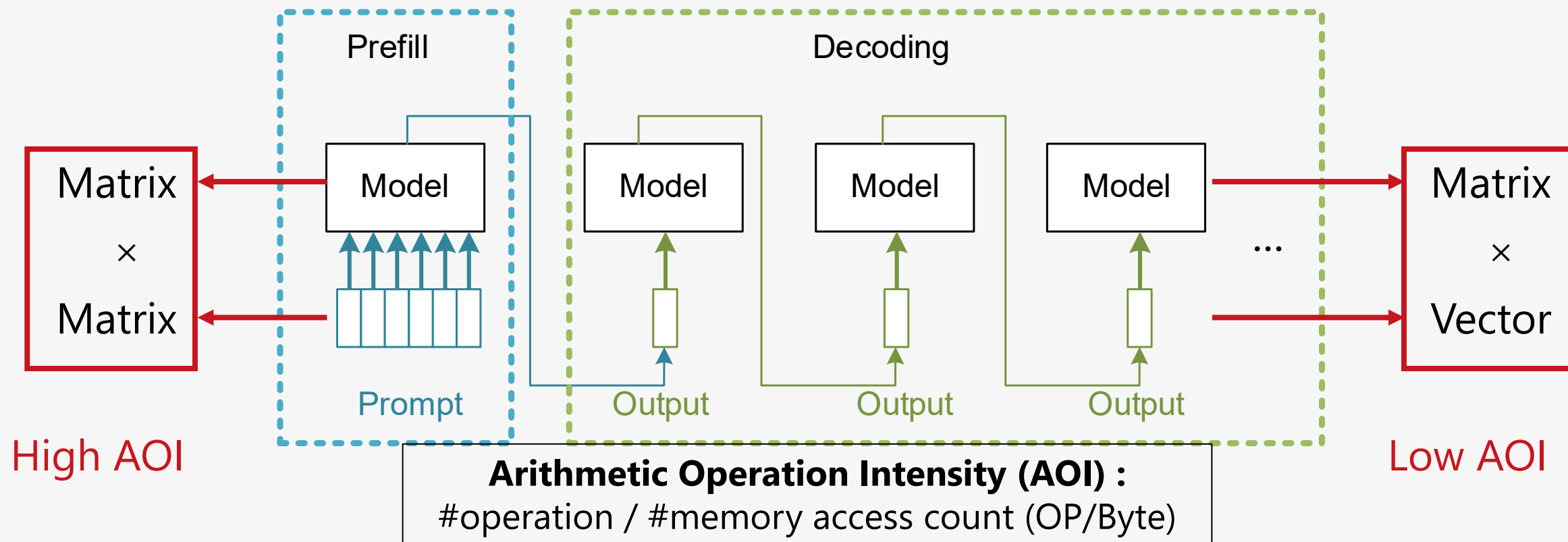


Background:

Large Language Models
PIM-based P-D Disaggregation

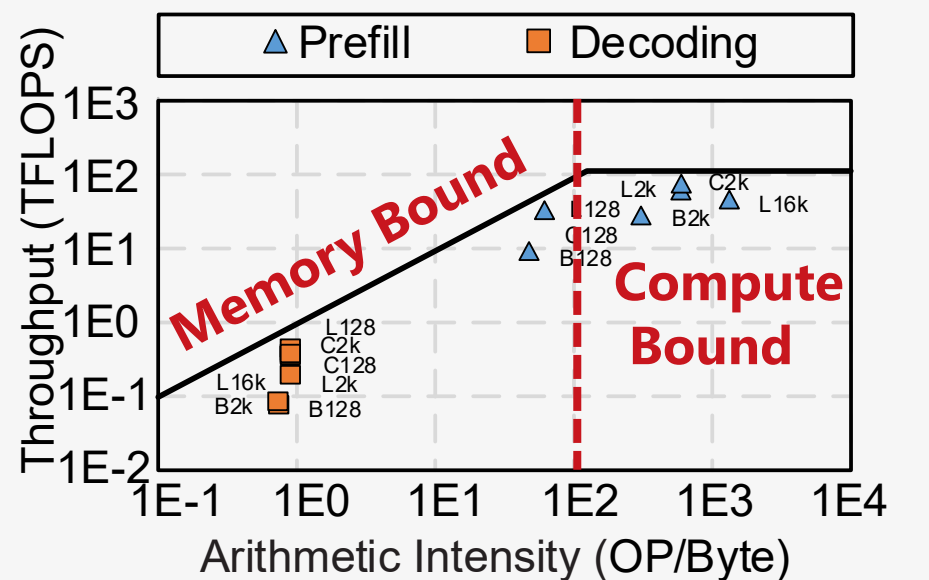
Prefill-Decoding Disaggregation

Two phases of Large Language Model (LLM) inference: **Prefill** & **Decoding**



AOI of Prefill / Decoding

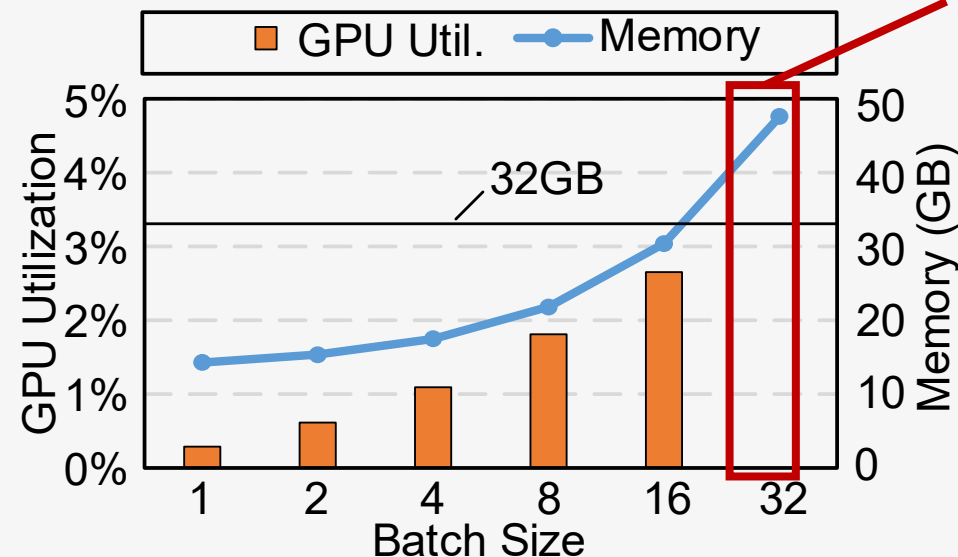
**Memory Footprint
Exceeds 32GB!
GPU Util still < 3%**



L:LLaMA-7B B:BLOOM-1b1 C:ChatGLM-9b

(a)

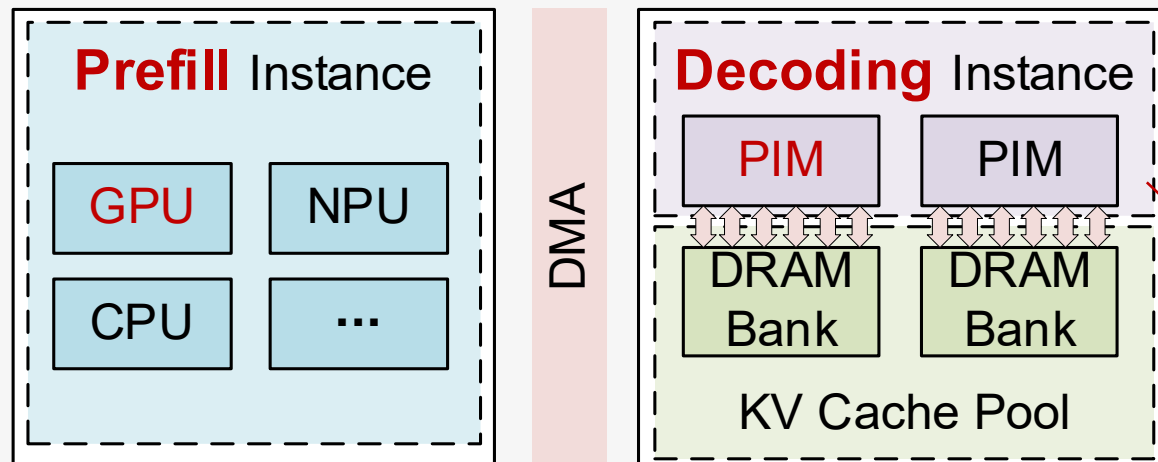
Prefill & Decoding on V100 GPU



(b)

Batching does not work for Decoding

PIM-based PD Disaggregation



Prefill: High AOI

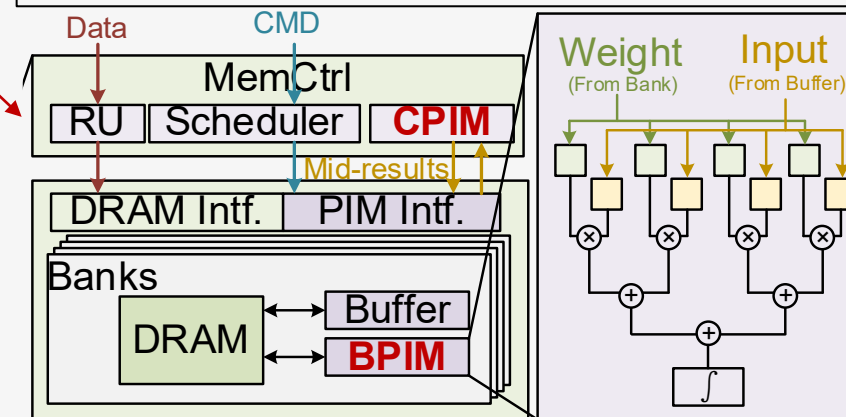
Decoding: Low AOI

GPU: High Computing Power

PIM: High Memory Bandwidth

Processing-in-Memory (PIM):

Integrate computational unit inside DRAM memory



Two-level PIM:

BPIM: In Bank for MAC

CPIM: In MemCtrl for mid-result reducing



Challenges & Solutions

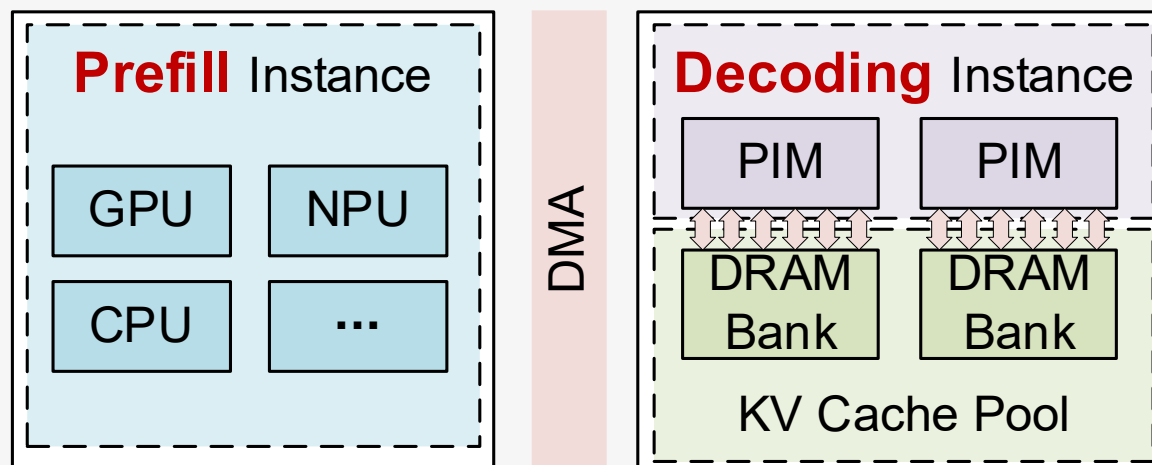
2 Challenges on Communication
Our Solution – BLADE Architecture

Communication Challenge 1 – V-Cache Transpose (1)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

During Key-Value Cache (KV Cache) Transfer



Step 1: In Prefill Instance
Generate **Output** & **KV Cache**

Step 3: In Decoding Instance
Consume **KV Cache**
and generate output

Step 2: Transfer **KV Cache**
to PIM side



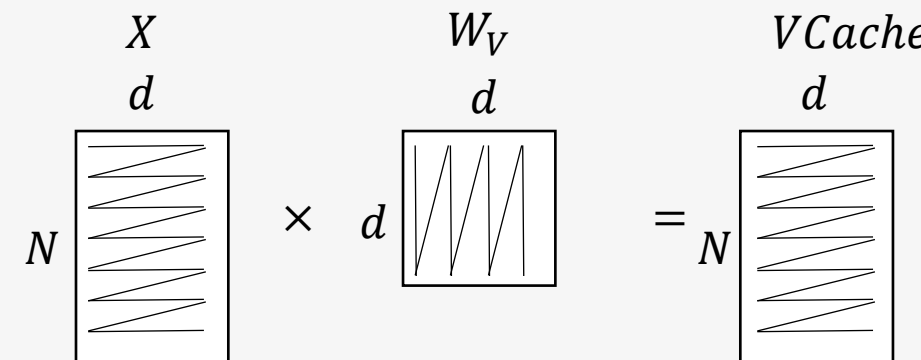
Communication Challenge 1 – V-Cache Transpose (2)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

⊙ (Prefill phase) Producing Value Cache:

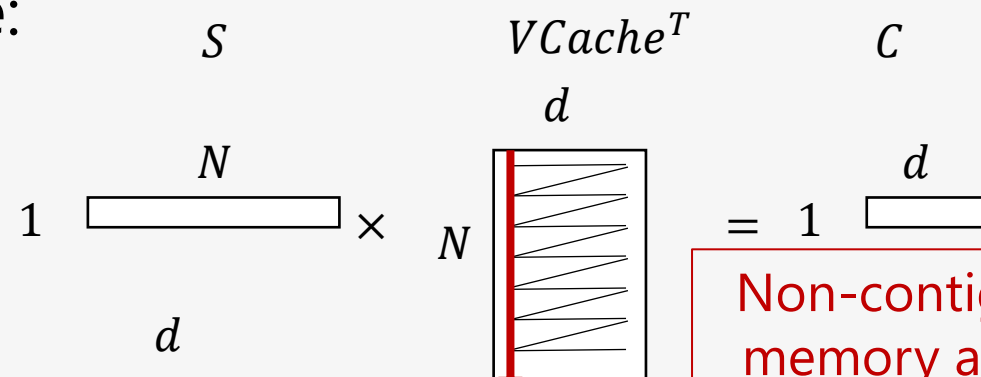
$$\begin{array}{c}
 \bullet \text{ } VCache = \begin{array}{c} X \\ (N \times d) \\ \swarrow \quad \searrow \\ \text{Token Num} \quad \text{Head Size} \end{array} \times \begin{array}{c} W_V \\ (d \times d) \end{array}
 \end{array}$$



Contiguous along
 d dimension

⊙ (Decoding phase) Consuming Value Cache:

$$\bullet \text{ } C = \begin{array}{c} S \\ (1 \times d) \end{array} \times \begin{array}{c} VCache^T \\ (N \times d) \end{array}$$



Preferred layout
(Require GPU Transpose!)

Non-contiguous
memory access:
Low DRAM
bandwidth

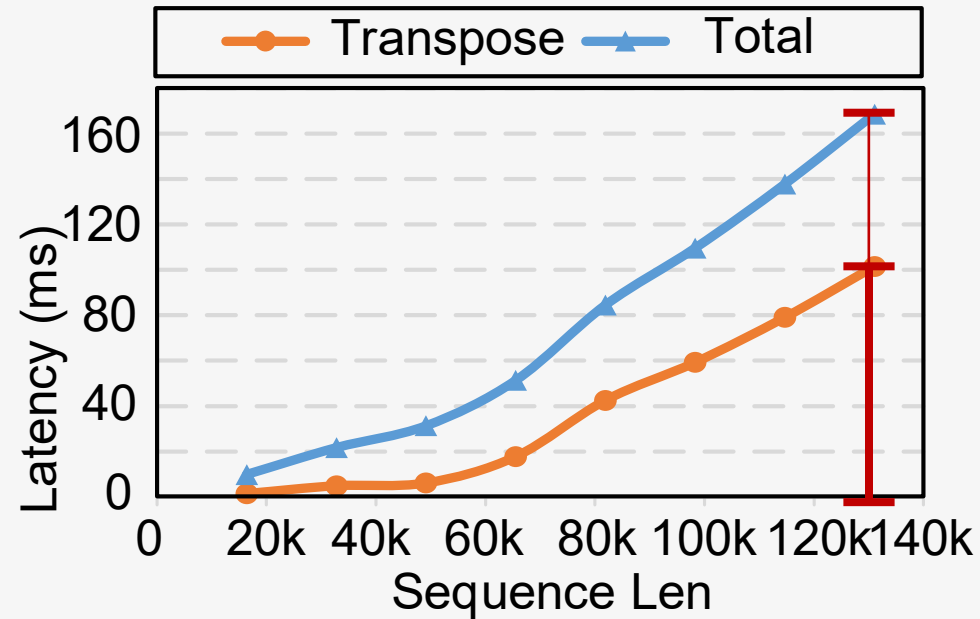


Communication Challenge 1 – V-Cache Transpose (3)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

⊗ Inefficiency of GPU Transpose:



**Over 78% of
KV Cache transfer time**

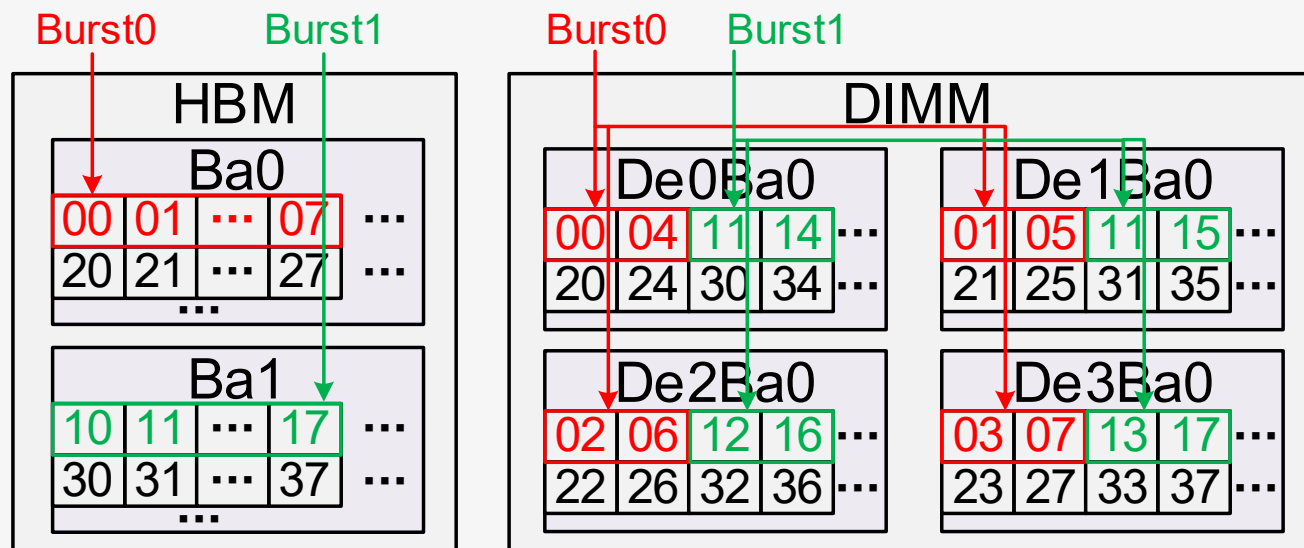


Solution to Challenge 1: Transpose-on-Transfer (1)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

Utilize the different access granularity between CPUs and PIM units:



HBM Burst:
Contiguous inside
one Bank

DIMM Burst:
Interleaved between Devices
(Device also called "Chip")

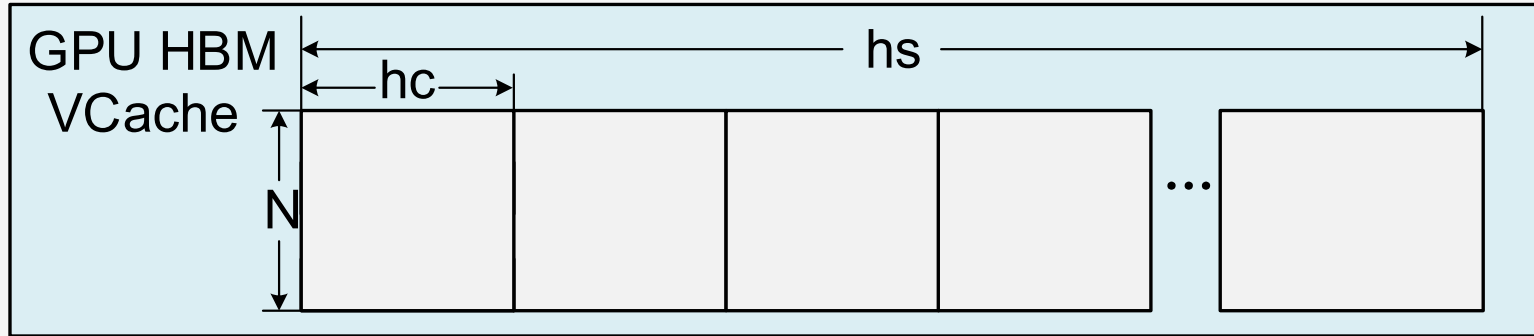


Solution for Challenge 1: Transpose-on-Transfer (2)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

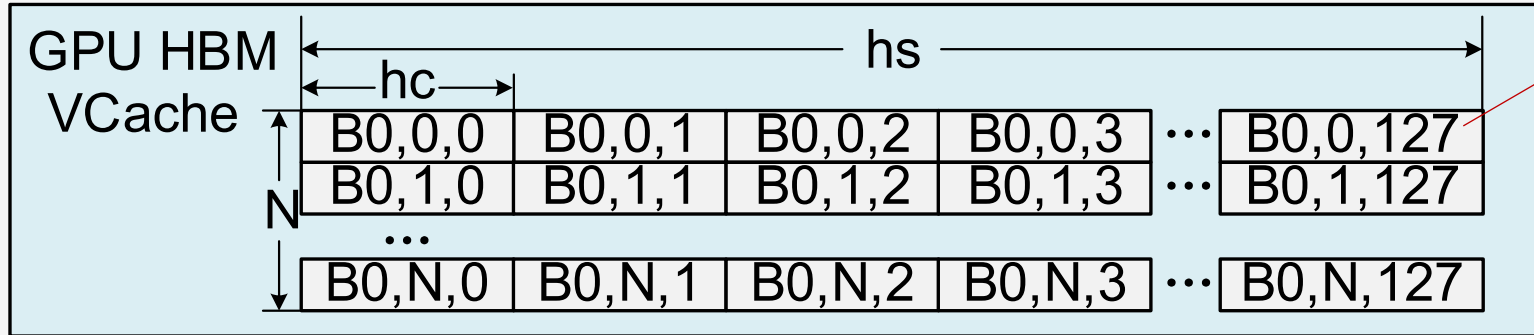
FP16 quant, $hs=128$, $hc=32$



Solution for Challenge 1: Transpose-on-Transfer (2)



FP16 quant, $hs=128$, $hc=32$



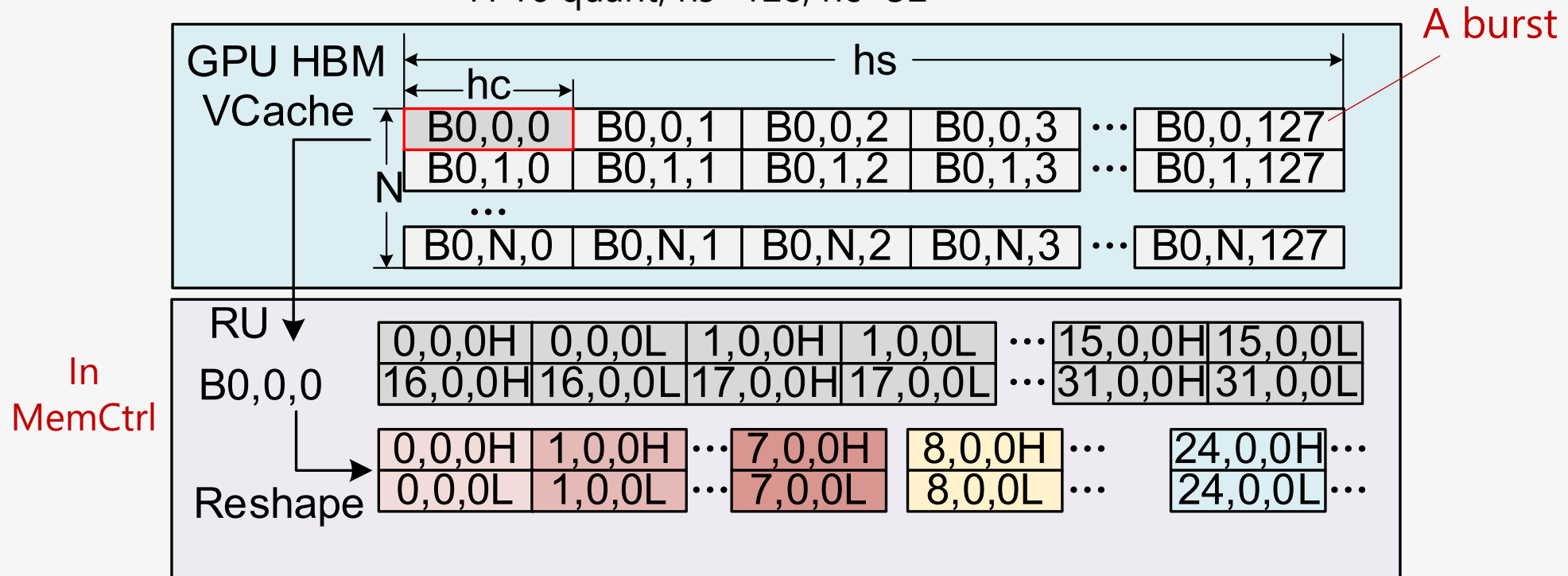
A burst

Solution for Challenge 1: Transpose-on-Transfer (2)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

FP16 quant, hs=128, hc=32

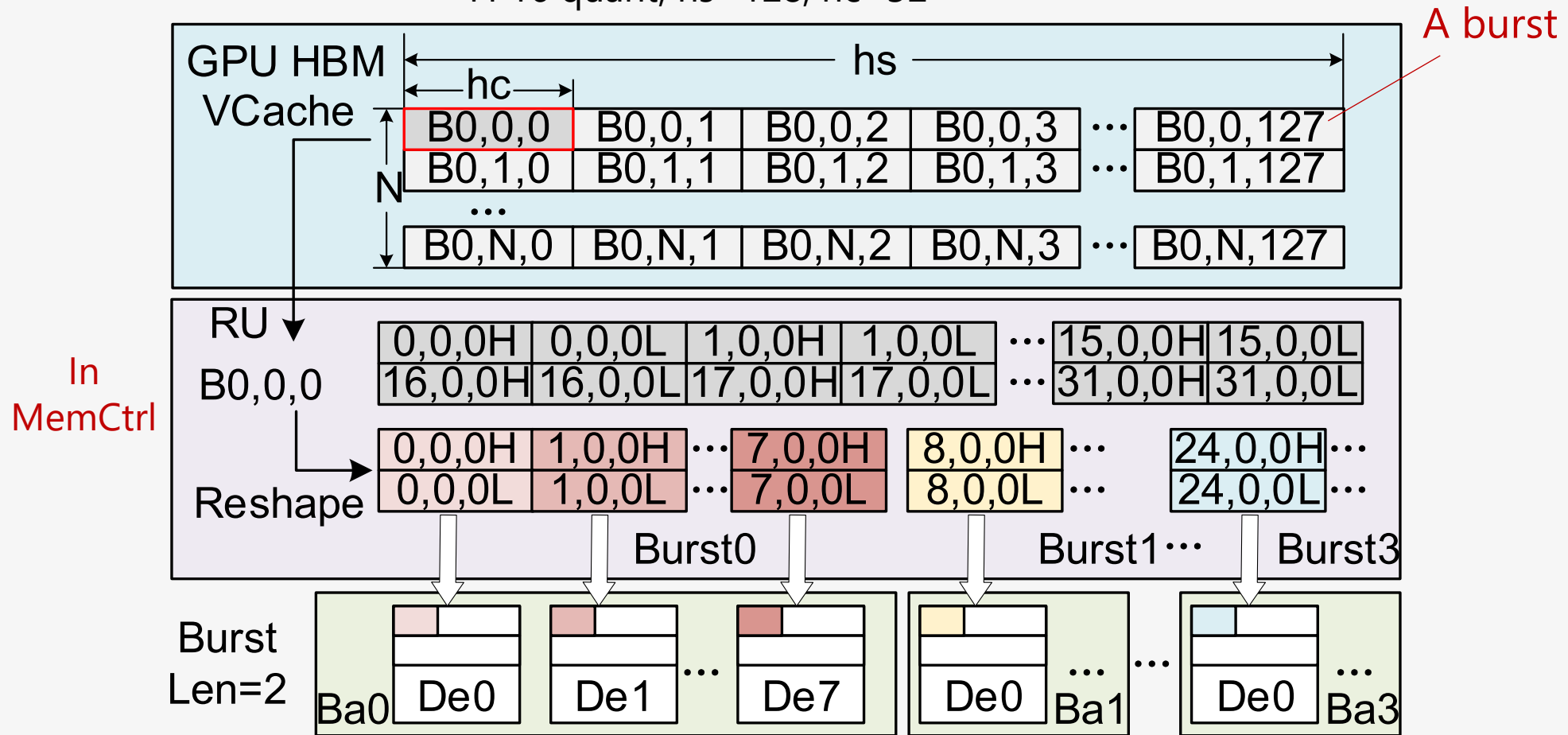


Solution for Challenge 1: Transpose-on-Transfer (2)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

FP16 quant, hs=128, hc=32

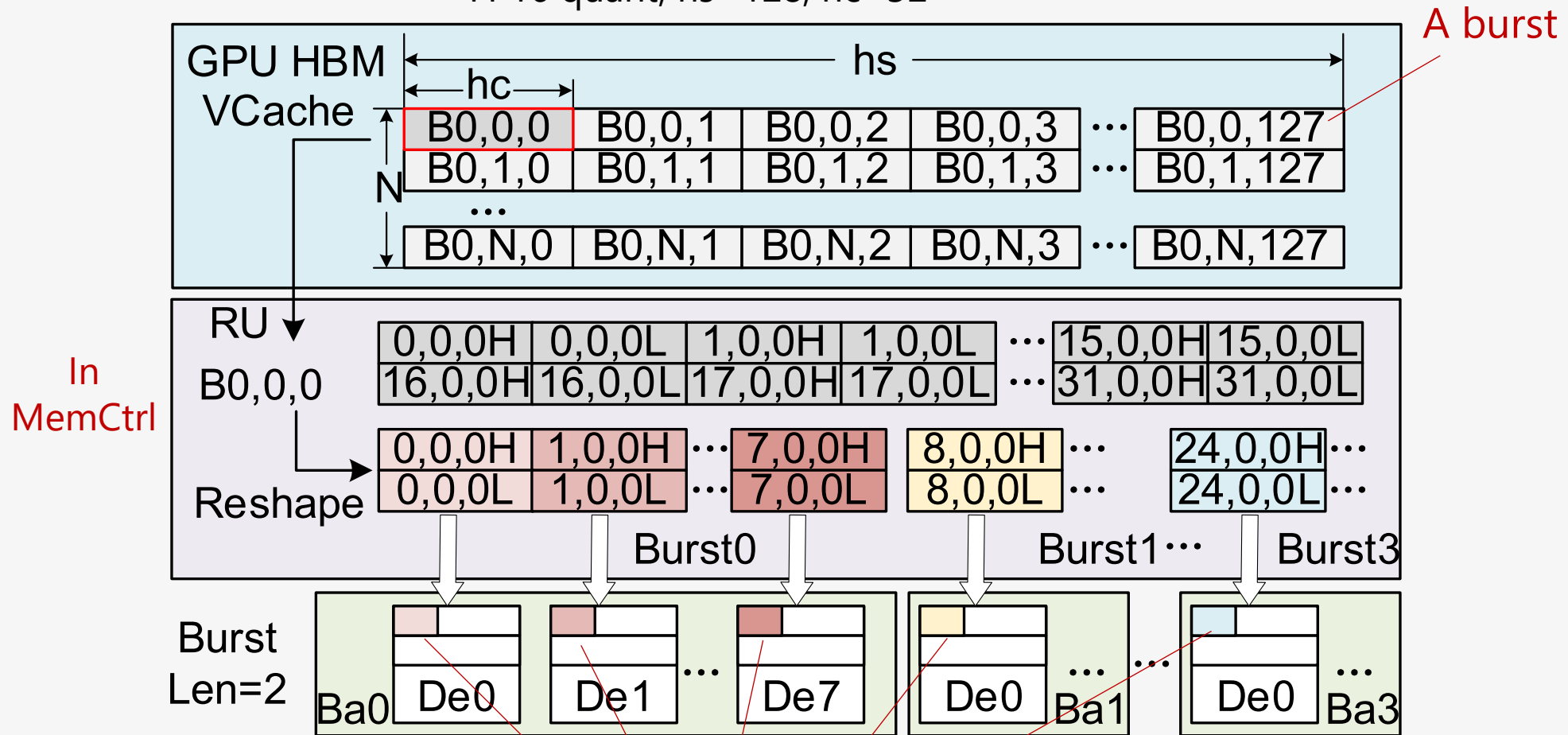


Solution for Challenge 1: Transpose-on-Transfer (2)



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

FP16 quant, hs=128, hc=32



Adjacent Elements on different Devices



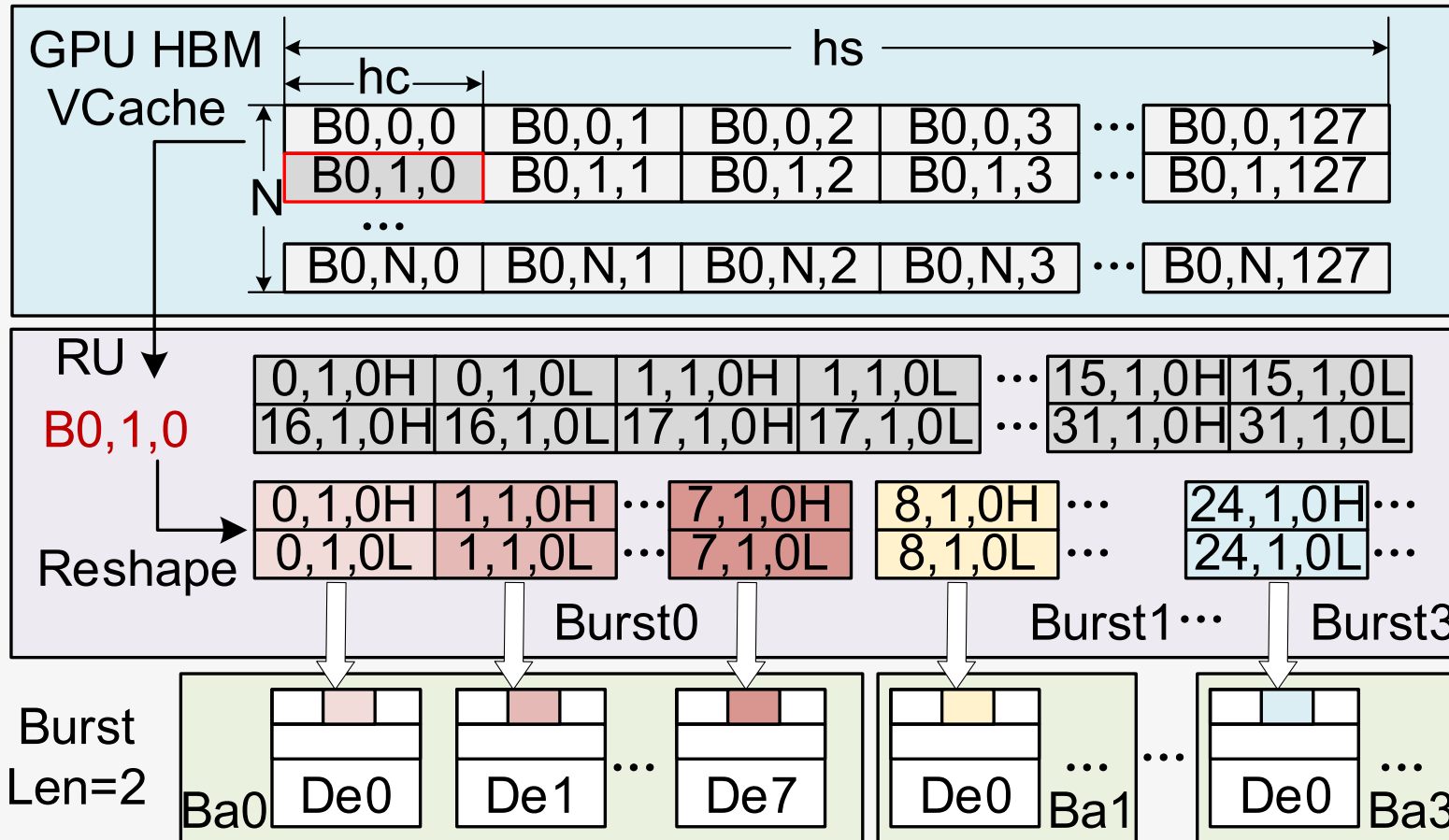
Solution for Challenge 1: Transpose-on-Transfer (2)



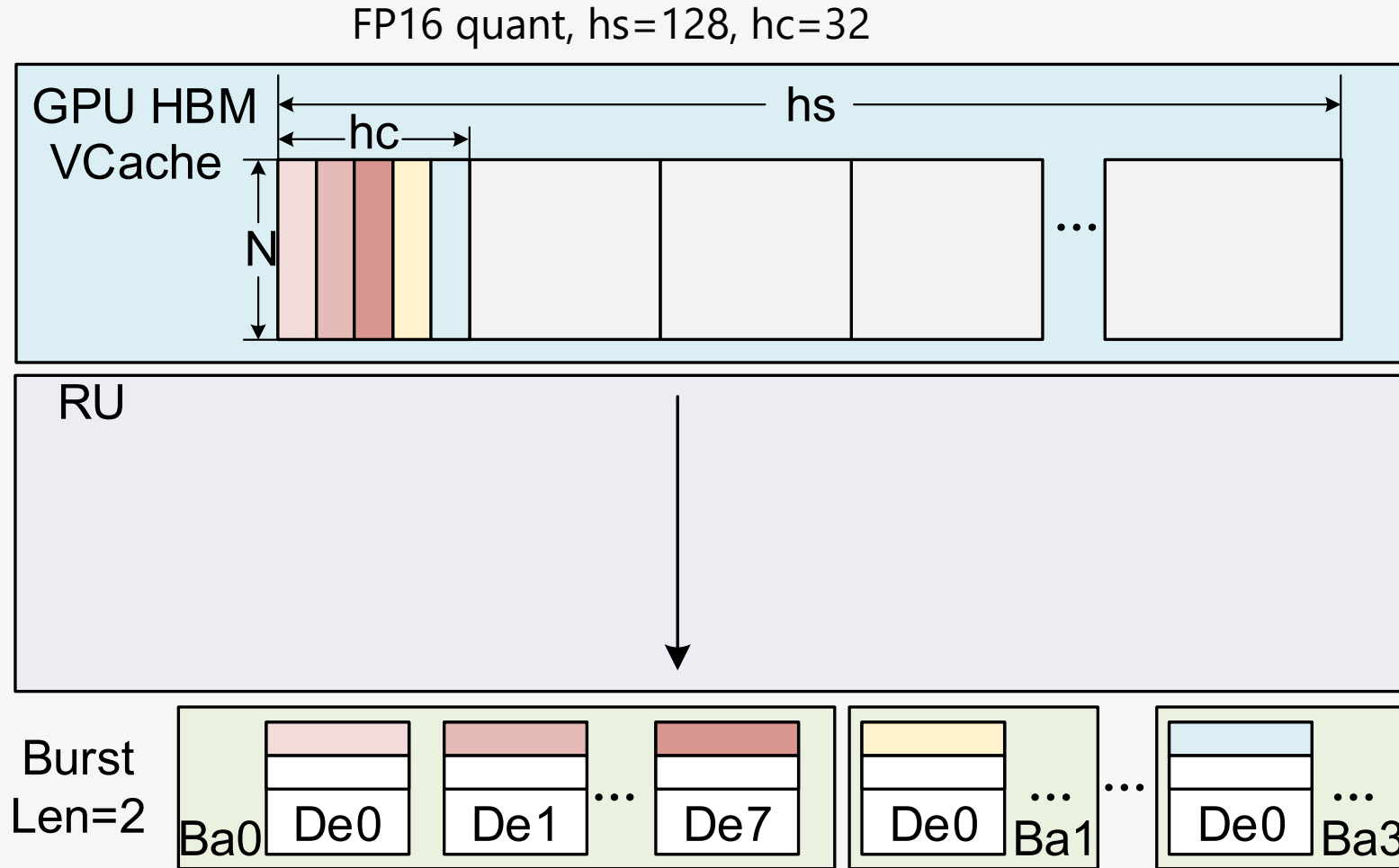
上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

FP16 quant, $hs=128$, $hc=32$

Next Burst



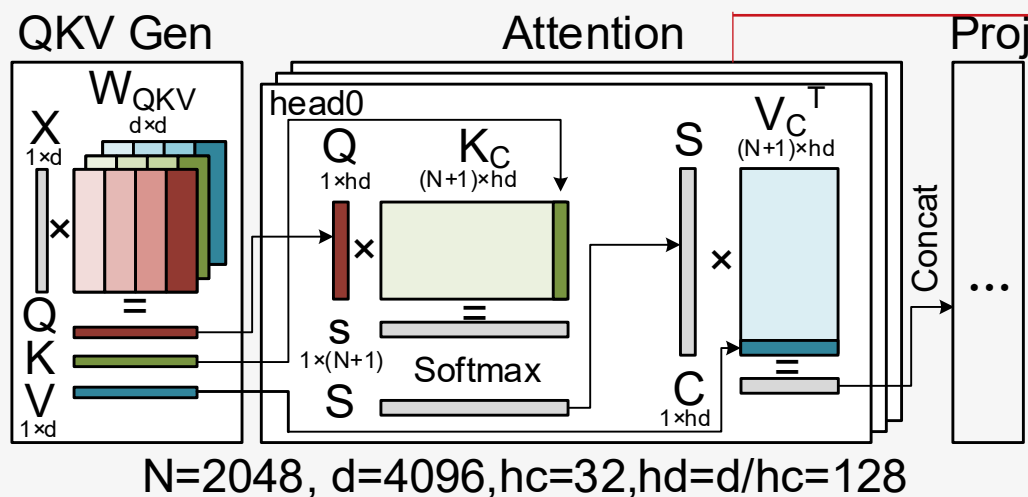
Solution for Challenge 1: Transpose-on-Transfer (2)



Transposed!

Communication Challenge 2 – Data Movement

Attention Layer Task Division on GPU



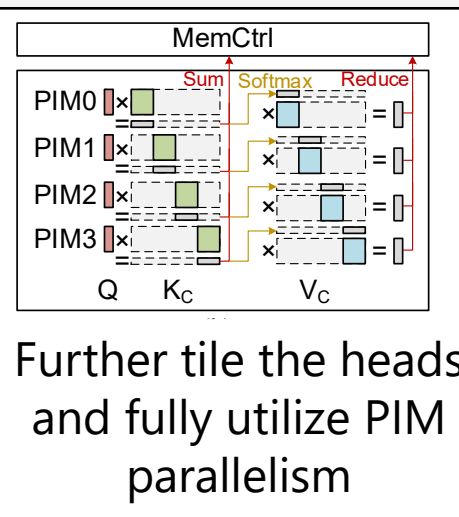
Divide Task based on **Heads**

16~32 Heads

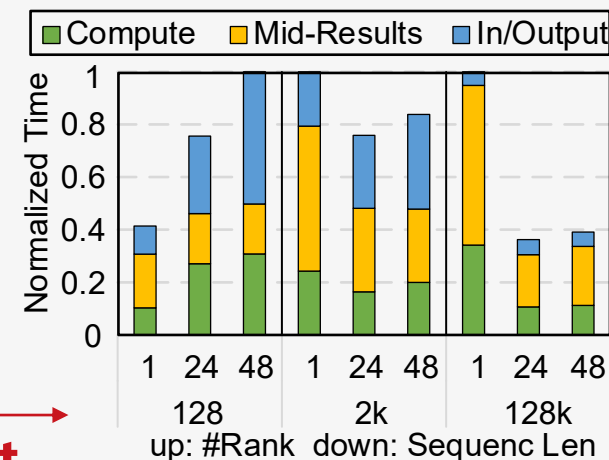
PIM Units Parallelism:

- A 8 Channel*4 Rank DDR5 Memory System:
 - 32 (Ranks) * 16 (Banks) -> 2048 (PIM Units)

Under-utilized!



Leads to heavier mid-result movement overhead



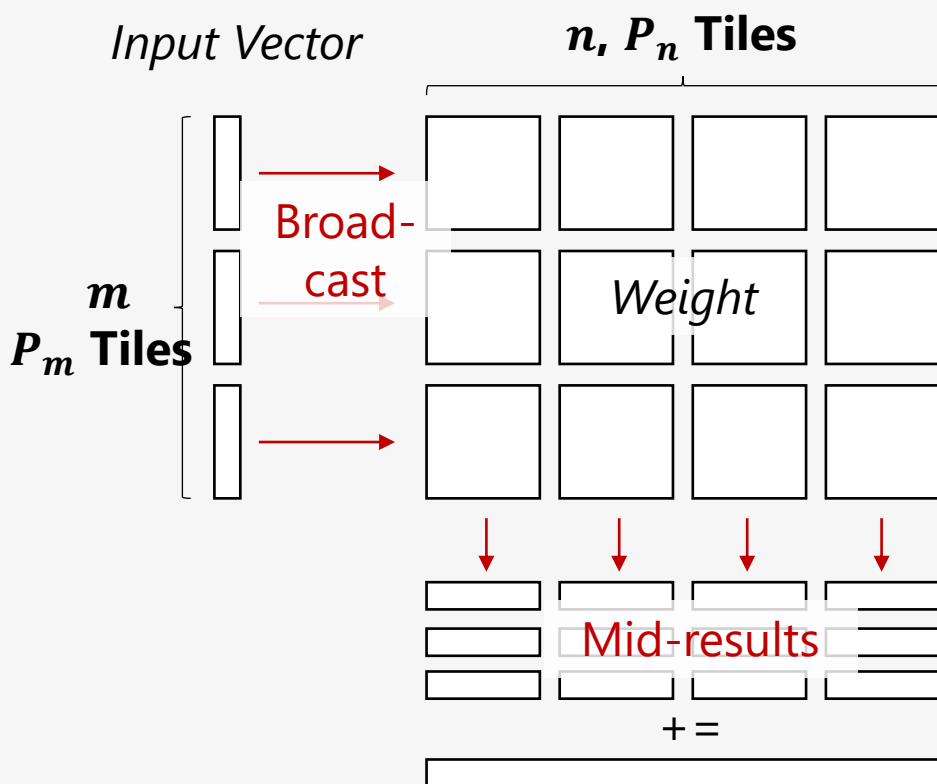
Solution to Challenge 2: Dynamical Parallelism Scaling



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

Estimate the PIM Computing & Communication Overhead

Linear Layer:



PIM Computation Time:

$$T_{comp} = \frac{mn}{P_m P_n \cdot BDW_{PIM}}$$

Communication Time:



$$T_{comm} = \frac{mP_n + nP_m + n_c P_m}{n_c \cdot BDW_{DRAM}}$$

Bandwidth
(n_c : channel num)

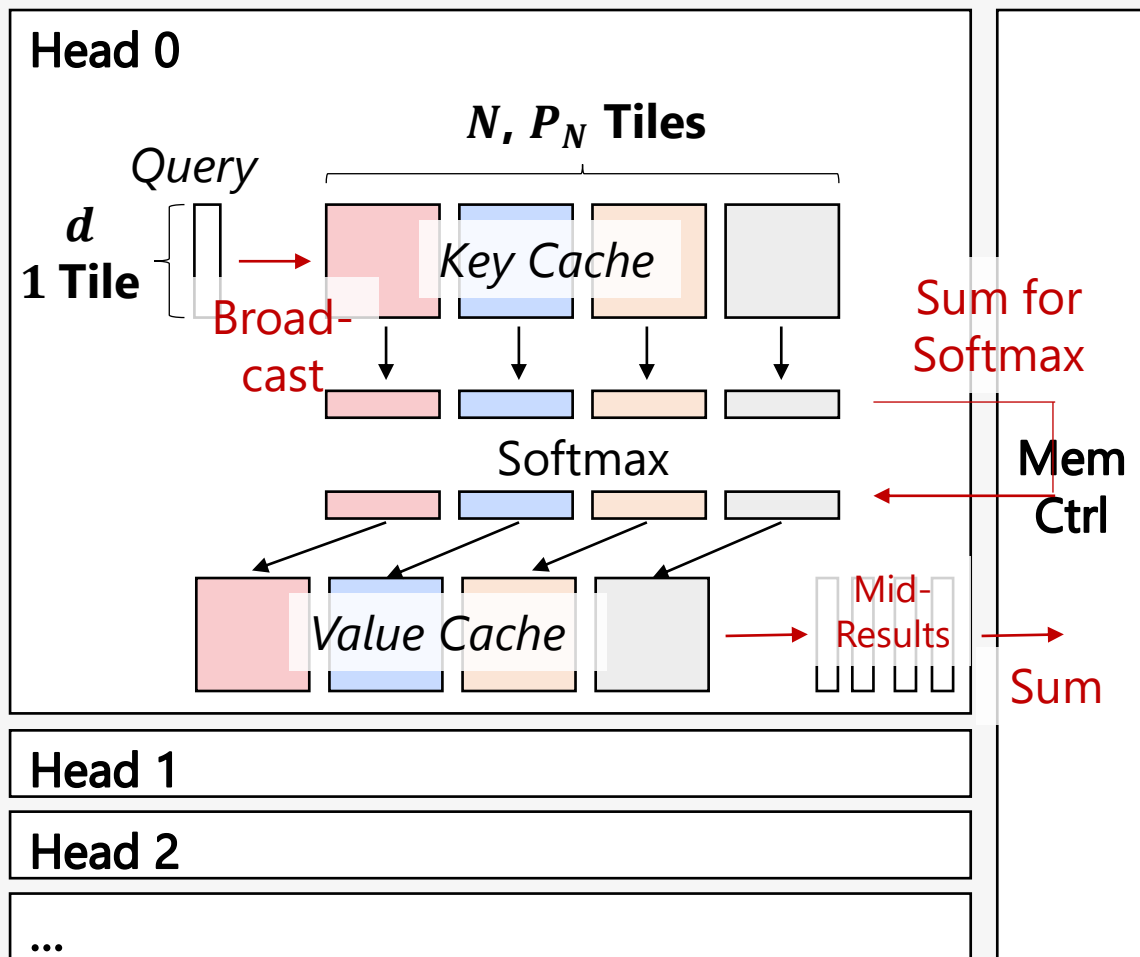


Solution to Challenge 2: Dynamical Parallelism Scaling



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

Attention Layer:



PIM Computation Time:

$$T_{comp} = \frac{2dN}{P_N \cdot BDW_{PIM}}$$

Communication Time:

$$T_{comm} = \frac{dP_N + 2P_N + dP_N + n_c d}{n_c \cdot \underset{\substack{\text{Bandwidth} \\ (n_c: \text{channel num})}}{BDW_{DRAM}}}$$

Broadcast **Softmax Sum** **Mid-Result (Intra-Channel)** **Mid-Result (Inter-Channel)**



Solution to Challenge 2: Dynamical Parallelism Scaling



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

⊗ PIM Computation Time: $T_{comp} = \frac{2dN}{P_N \cdot BDW_{PIM}}$

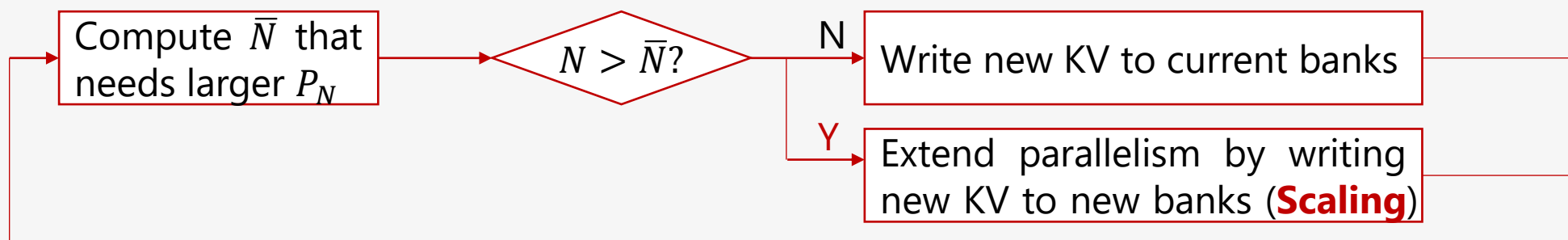


**N (#tokens in KV Cache)
Grows during decoding!**

⊗ Communication Time: $T_{comm} = \frac{dP_N + 2P_N + dP_N + n_c d}{n_c \cdot BDW_{DRAM}}$

**Larger N needs larger P_N to balance
computation & communication**

Dynamical Parallelism Scaling:





Experiments

Experimental Setup

⊗ Benchmarks: **ChatGLM-9B** (GHA, G=2), **LLaMA-7B**, **BLOOM-1B1**

⊗ GPU Baselines:

- **GPU_OFF**: Decoding with GPU, KV-Cache offloaded to CPU memory
- **GPU_NO**: Decoding with GPU, KV-Cache always in GPU memory
 - (Faster, but support shorter KV-Cache)

⊗ PIM Baselines:

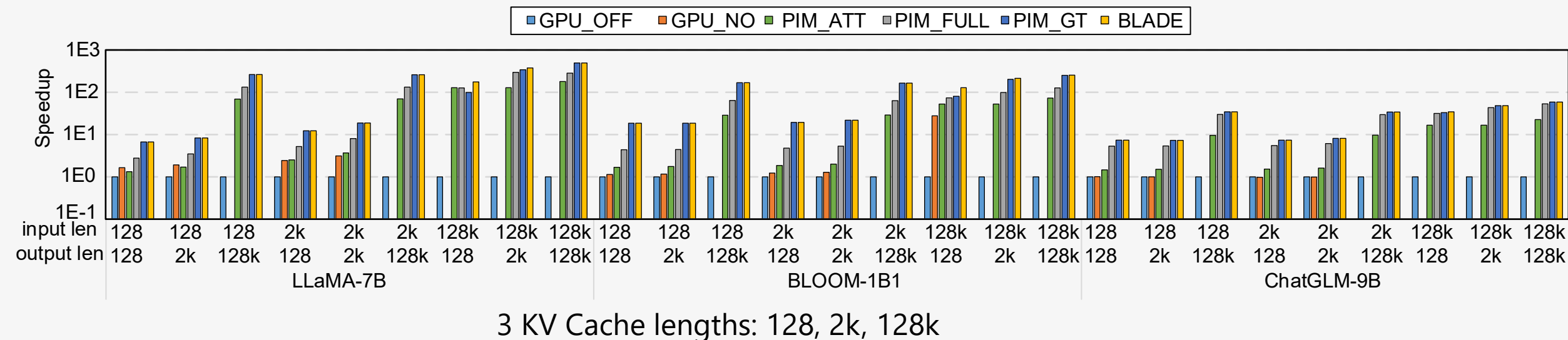
- **PIM_ATT**[ASPLOS'24]: Linear layer compute on GPU
- **PIM_GT**: V-cache transposed by GPU (Without Tech1 **Transpose-on-Transfer**)
- **PIM_FULL**: Always exploit full PIM parallelism (Without Tech2 **Dynamical Parallelism Scaling**)

⊗ Simulation:

- DRAM model: Ramulator2
- Circuits: Synopsys DC + TSMC 90nm

System Configuration	
GPU	Nvidia V100, 112TFLOPS
DRAM	DDR4, 12 Channel*4 Ranks
PIM	32GOPS, 64 per Rank

Results (1) – Overall Performance



BLADE:

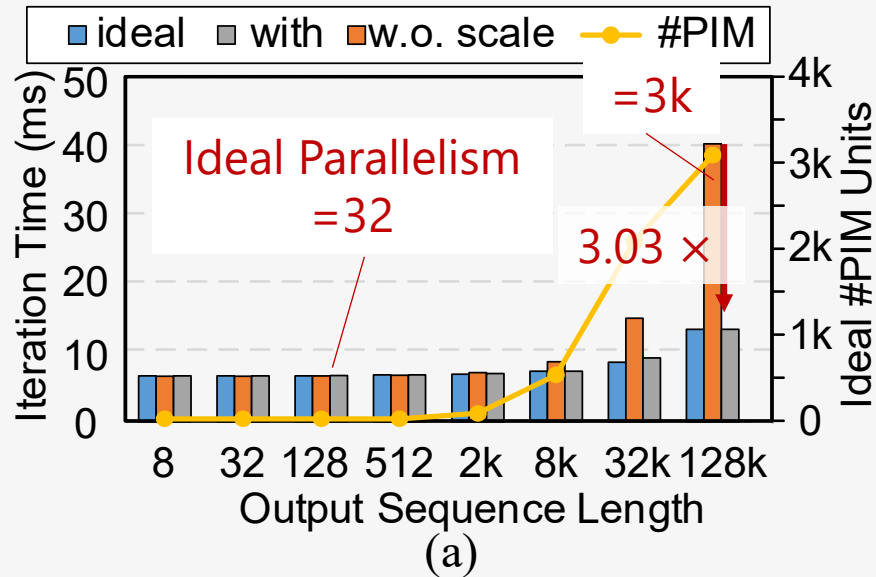
$4.93 \times$ to **PIM_ATT** (Linear Layers are on GPU)

$105.7 \times$ to **GPU_OFF**

$2.09 \times$ to **PIM_FULL** (w.o. Tech 2 Dynamical Parallelism Scaling)

$1.54 \times$ to **PIM_GT** (w.o. Tech 1 Transpose on Transfer)

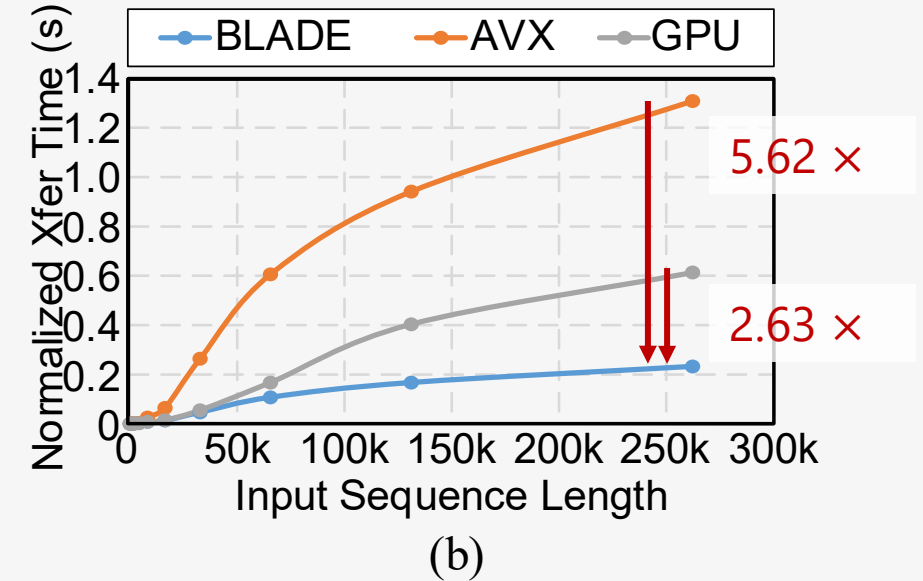
Results (2)



Dynamical Parallelism Scaling:

Baseline: Always use the initial parallelism

3.03 × for 128k KV Cache



Transpose-on-Transfer:

2.63 × to **GPU**-based Transpose
5.62 × to **CPU**-based Transpose (AVX Instructions)



Context: Large Language Model (LLM) inference, **Prefill-Decoding Disaggregation**

- Prefill- \rightarrow GPU; Decoding \rightarrow DRAM-based **Processing-in-Memory (PIM)**

Challenges on Communication & Solutions:

- Challenge 1: **High KV Cache Transpose Latency** During “GPU- \rightarrow PIM” KV Cache Transfer
- Solution 1: **Transpose on Transfer Method**
- Challenge 2: Static PIM Parallelism **Fails to Balance Communication & Computation**
- Solution 2: **Dynamical Parallelism Scaling Method**

Results:

- Transpose on Transfer 1.54x; Dynamical Parallelism Scaling 2.09x



上海交通大學

SHANGHAI JIAO TONG UNIVERSITY

Thank You!

飲水思源 愛國榮校