# PIMBA: A Processing-in-Memory Acceleration for Post-Transformer Large Language Model Serving

**Wonung Kim**          Yubin Lee          Yoonsung Kim          Jinwoo Hwang

Seongryong Oh          Jiyong Jung          Aziz Huseynov          Woong Gyu Park
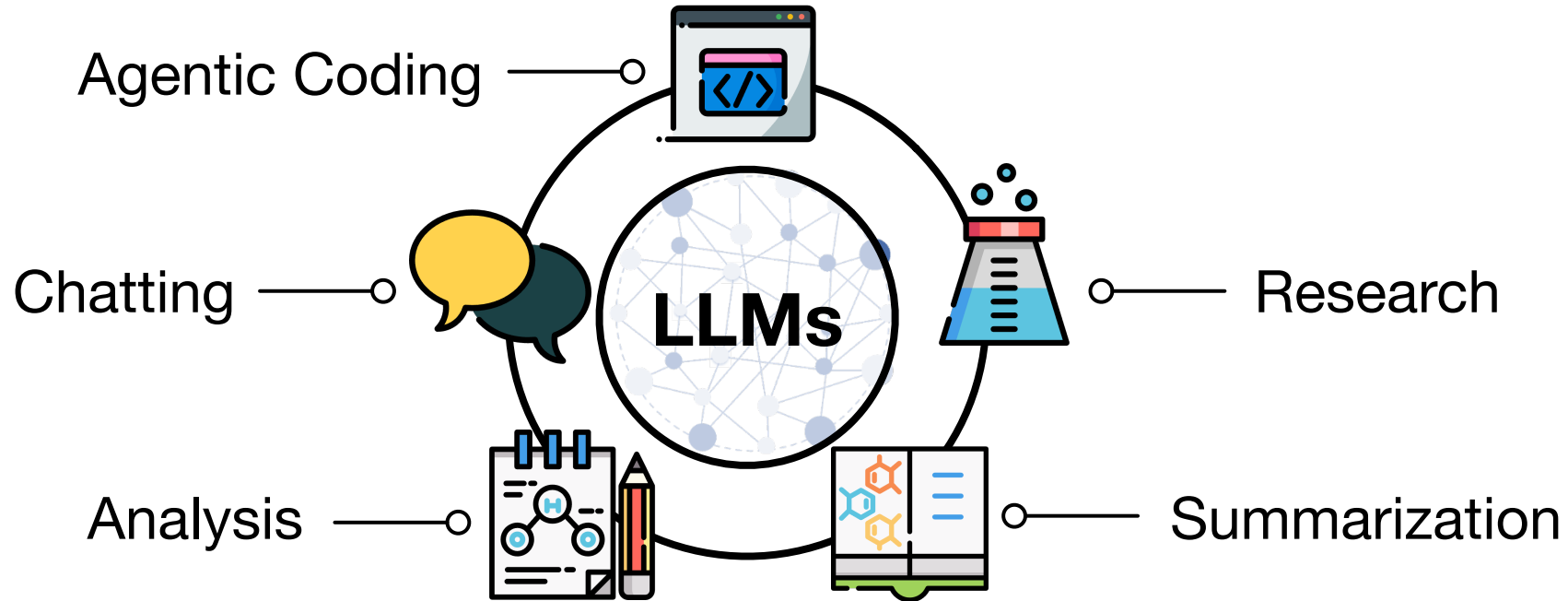
Chang Hyun Park[§]          Divya Mahajan[†]          Jongse Park

KAIST CASYS Lab          § Uppsala University          † Georgia Institute of Technology

# Large Language Models Are Pervasive
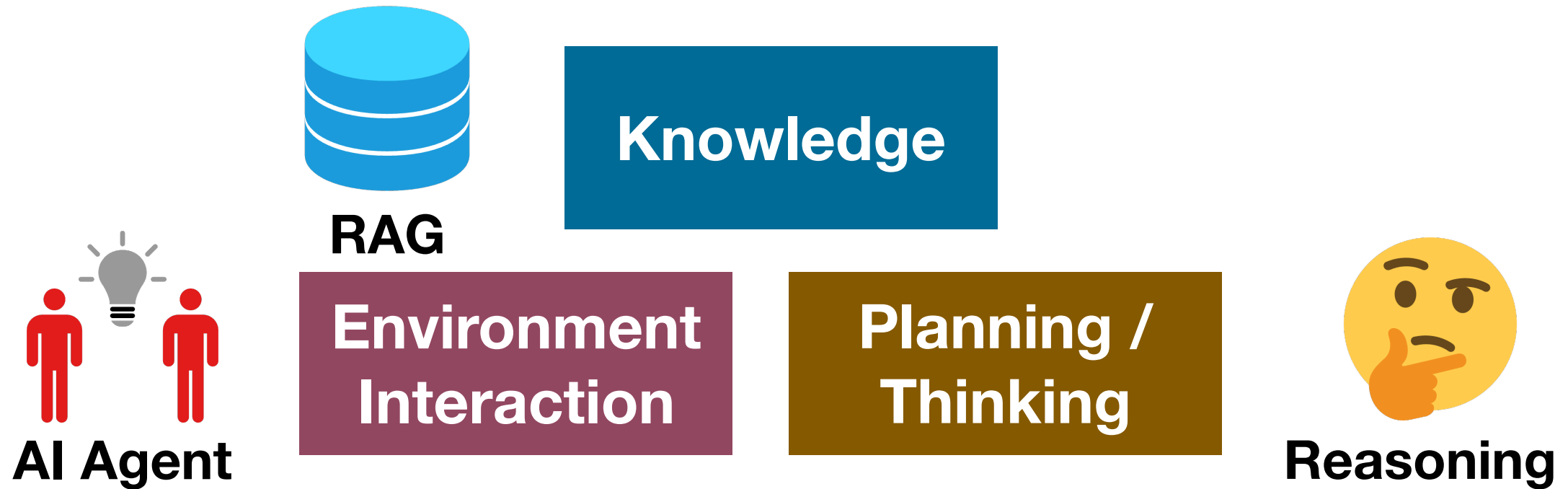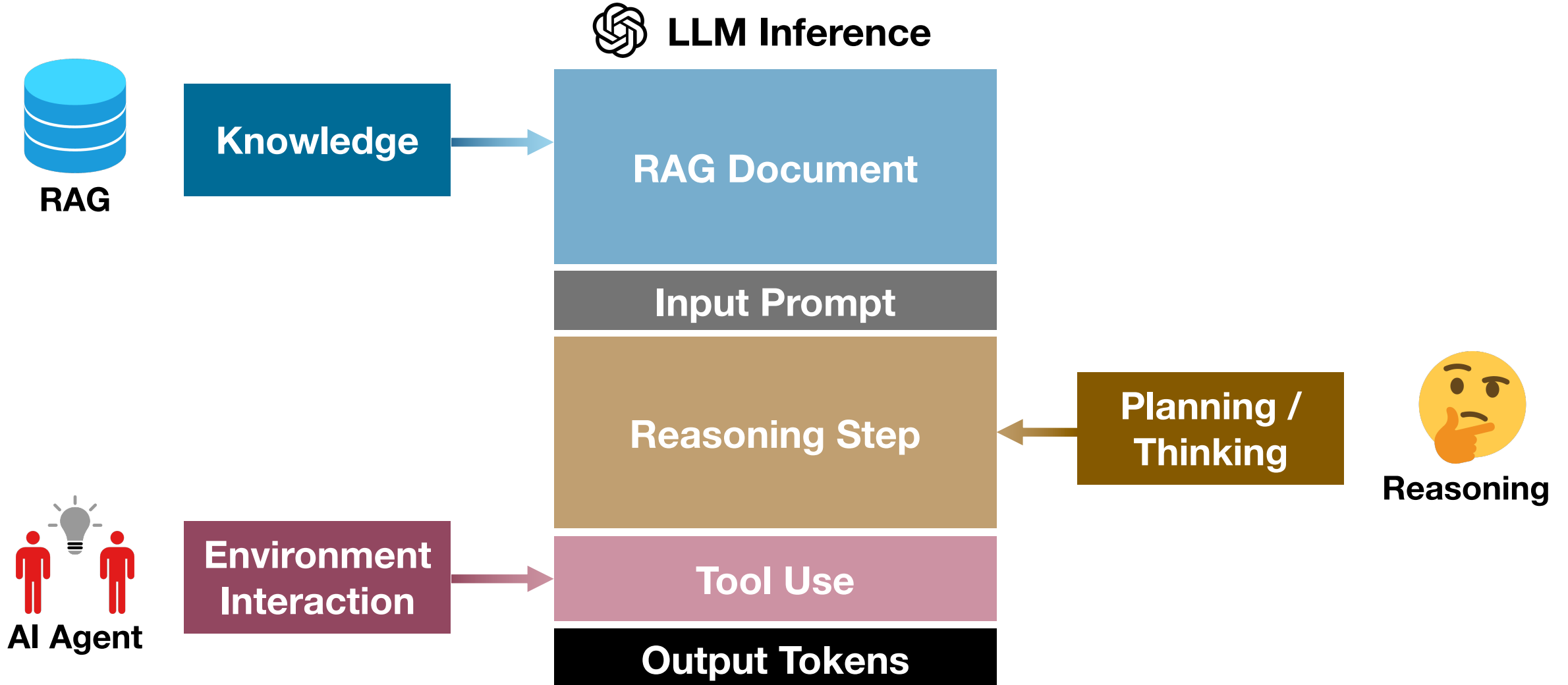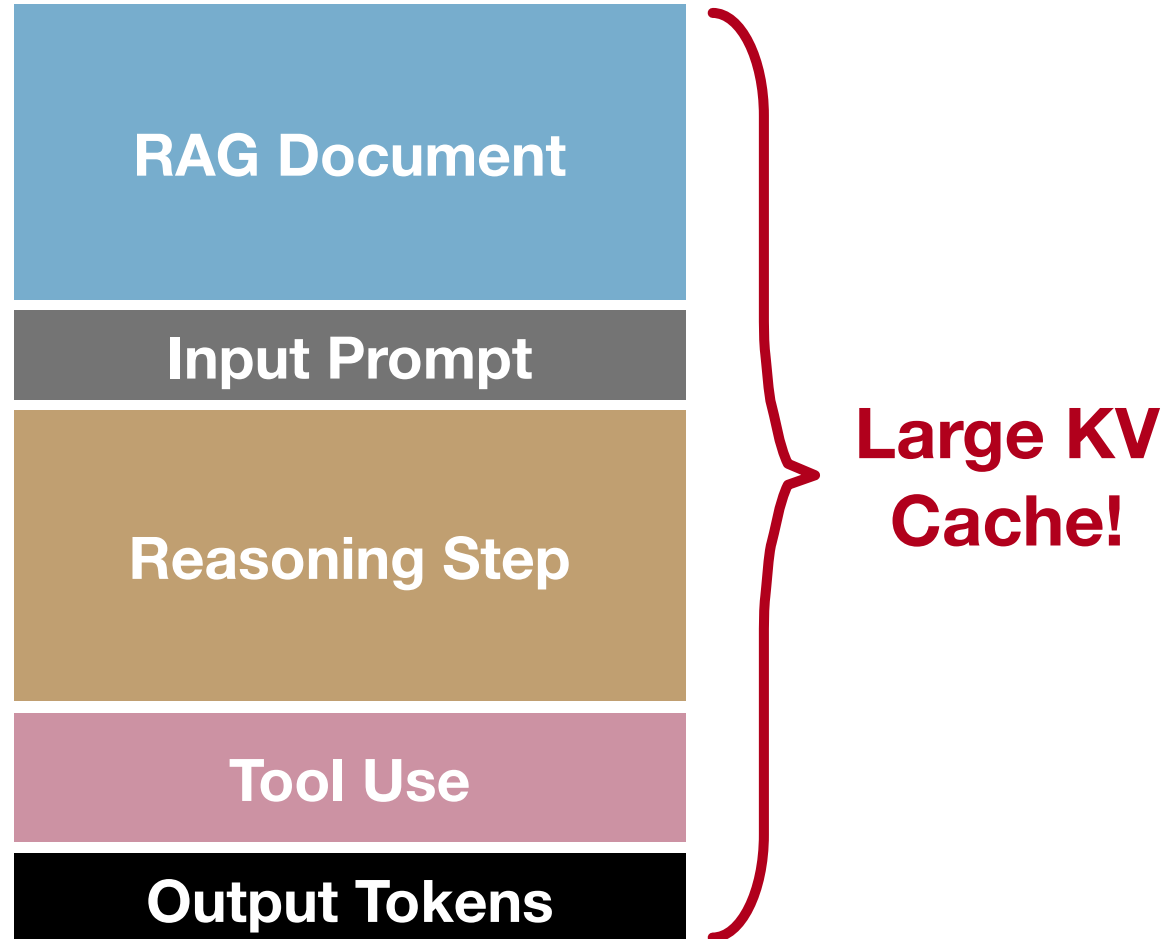
# Modern LLM Trends

# Modern LLM Trends

# Modern LLM Trends: Large KV Cache

# Modern LLM Trends: Large KV Cache

**LLM Inference**

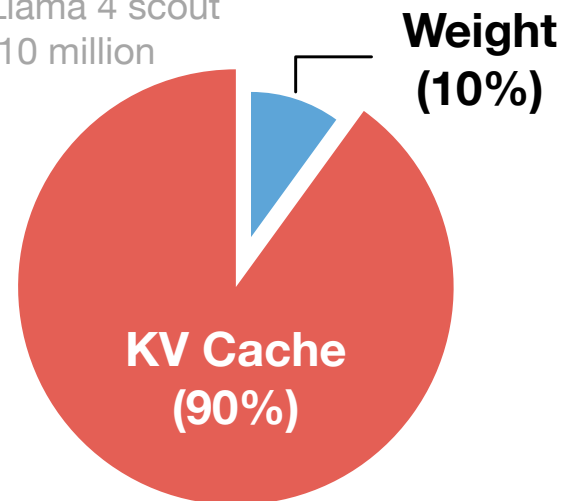| RAG Document |
|---|
| Input Prompt |
| Reasoning Step |
| Tool Use |
| Output Tokens |

| Model | Context Length |
|---|---|
| Llama 4 | 10M |
| Grok 4 | 2M |
| Gemini 2.5 | 1M |
| Claude 4 | 1M |

**Modern LLMs now support up to 10M context length**

Model: Llama 4 scout
Length: 10 million

Weight (10%)

KV Cache (90%)

**Even with GQA, MoE, KV Cache dominates!**

# Modern LLM Trends: Large KV Cache

LLM Inference

RAG Document

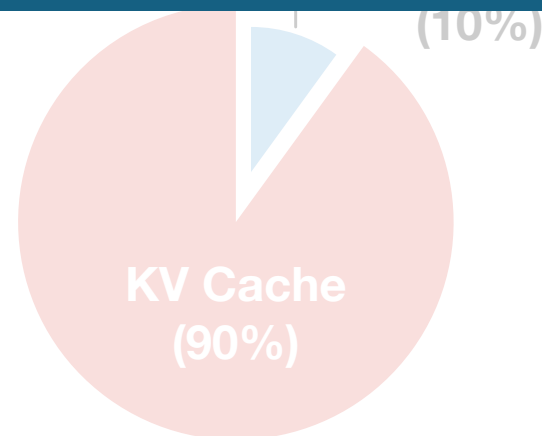| Model | Context Length |
|-------|----------------|
| Llama 4 | 10M |
| Grok 4 | 2M |

**Modern LLMs now support up to 10M**

Reasoning Step

Tool Use

Output Tokens

(10%)

KV Cache
(90%)

**Even with GQA, MoE, KV Cache dominates!**

It has become inevitable to confront
the KV cache memory bottleneck

KAIST School of Computing

CASYS | KAIST Computer Architecture & System Lab

# Post-Transformer Models

| Prompt | Please | add | some | test | cases | to | this | project |
|--------|--------|-----|------|------|-------|-----|------|---------|

Transformer

| KV #0 | KV #1 | KV #2 | KV #3 | KV #4 | KV #5 | KV #6 | KV #7 |
|-------|-------|-------|-------|-------|-------|-------|-------|

| Prompt | Please | add | some | test | cases | to | this | project |
|--------|--------|-----|------|------|-------|-----|------|---------|

Post-Transformer

Fixed State

Mamba [1]
GLA [3]

RetNet [2]

 [4]

 [5]

IBM [6]

[1] Mamba: Linear-Time Sequence Modeling with Selective State Spaces
[2] Retentive Network: A Successor to Transformer for Large Language Models
[3] Gated Linear Attention Transformers with Hardware-Efficient Training

[4] Nemotron-H: A Family of Accurate and Efficient Hybrid Mamba-Transformer Models
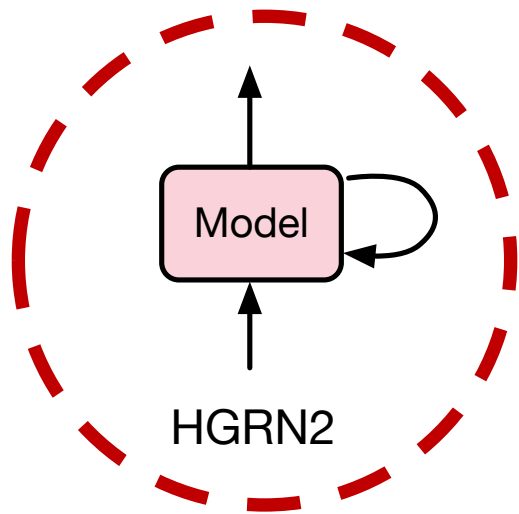[5] Samba: Simple Hybrid State Space Models for Efficient Unlimited Context Language Modeling
[6] IBM Granite 4.0: hyper-efficient, high performance hybrid models for enterprise
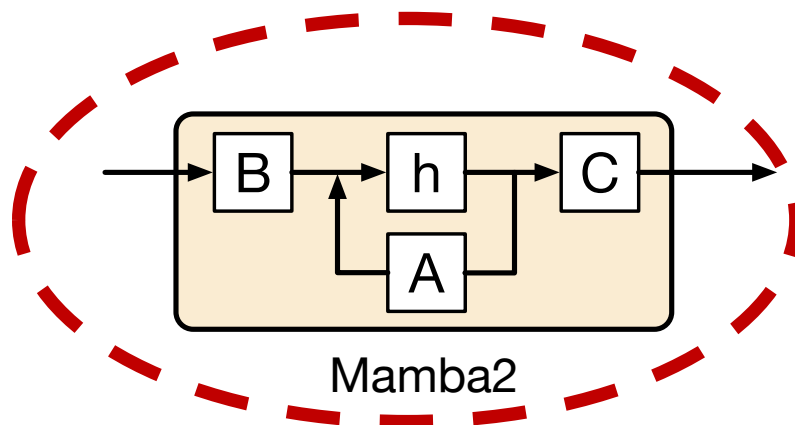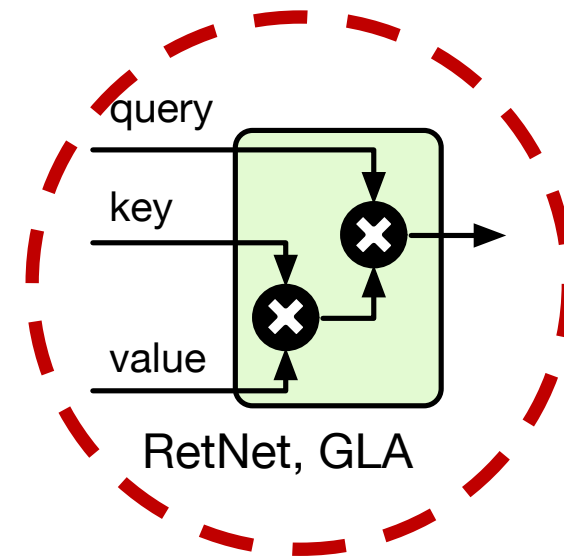
KAIST School of Computing

CASYS | KAIST Computer Architecture & System Lab

# We begin by analyzing
# how post-transformers operate
# to identify performance bottlenecks
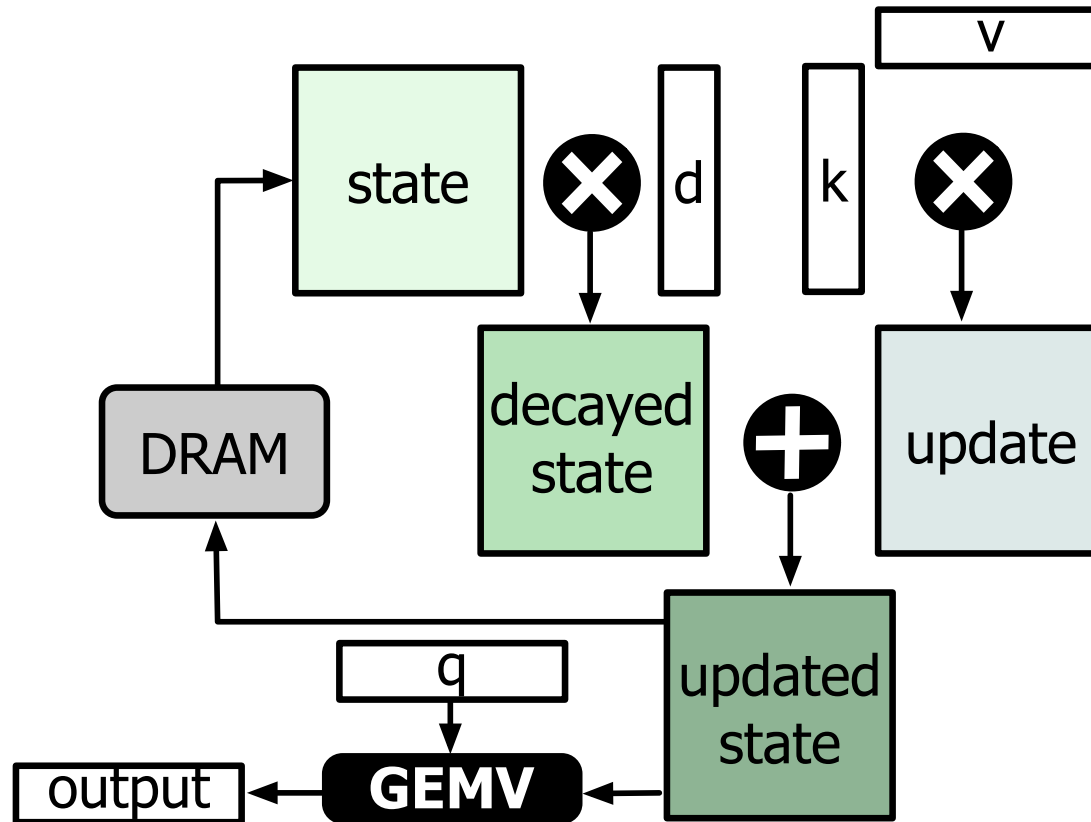
**Post-Trans formers**

HGRN2

RetNet, GLA

Mamba2

**But, unlike transformers, post-transformers exhibit diverse algorithmic forms**

We identified a common operator shared across these algorithms, which we call it as,
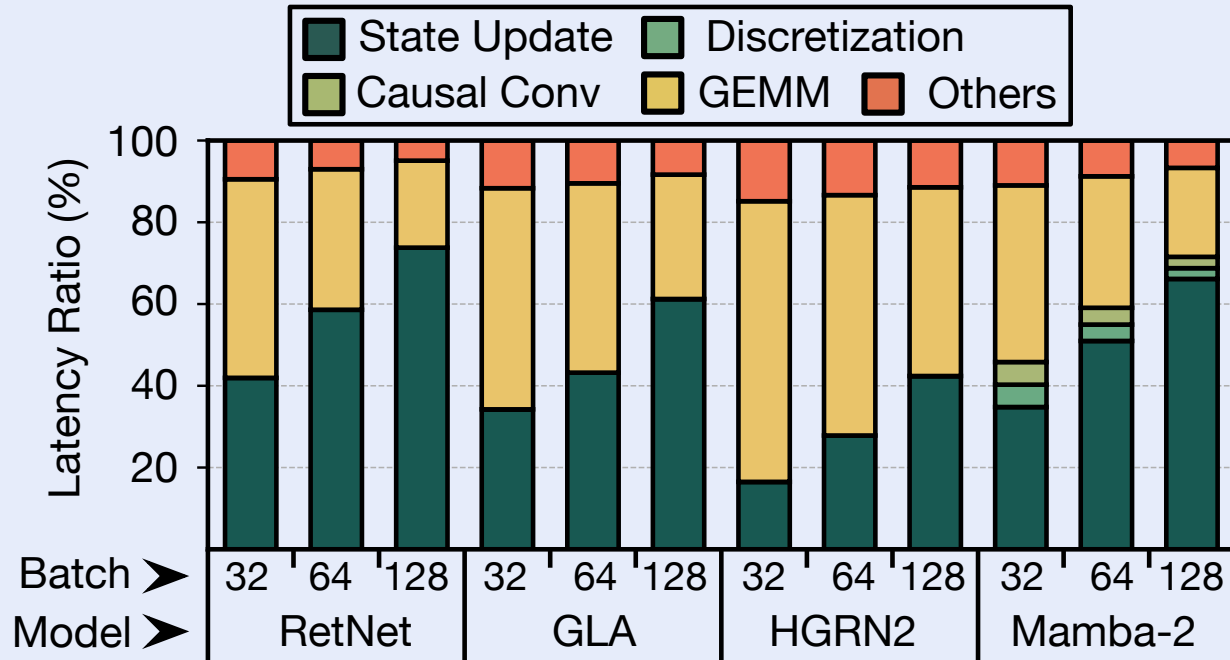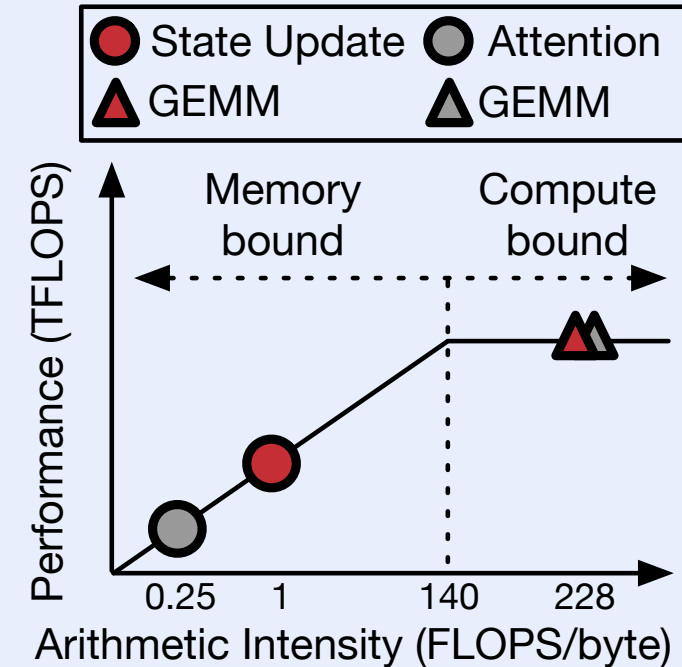
**State Update**

# State Update Operation



❶ **Weight decay**

❷ **Outer product**

❸ **Update**

❹ **GEMV**

# Characterizing Post-Transformers



**Latency Breakdown**

- State Update
- Discretization
- Causal Conv
- GEMM
- Others

**Roofline Analysis**

- State Update
- Attention
- GEMM
- GEMM

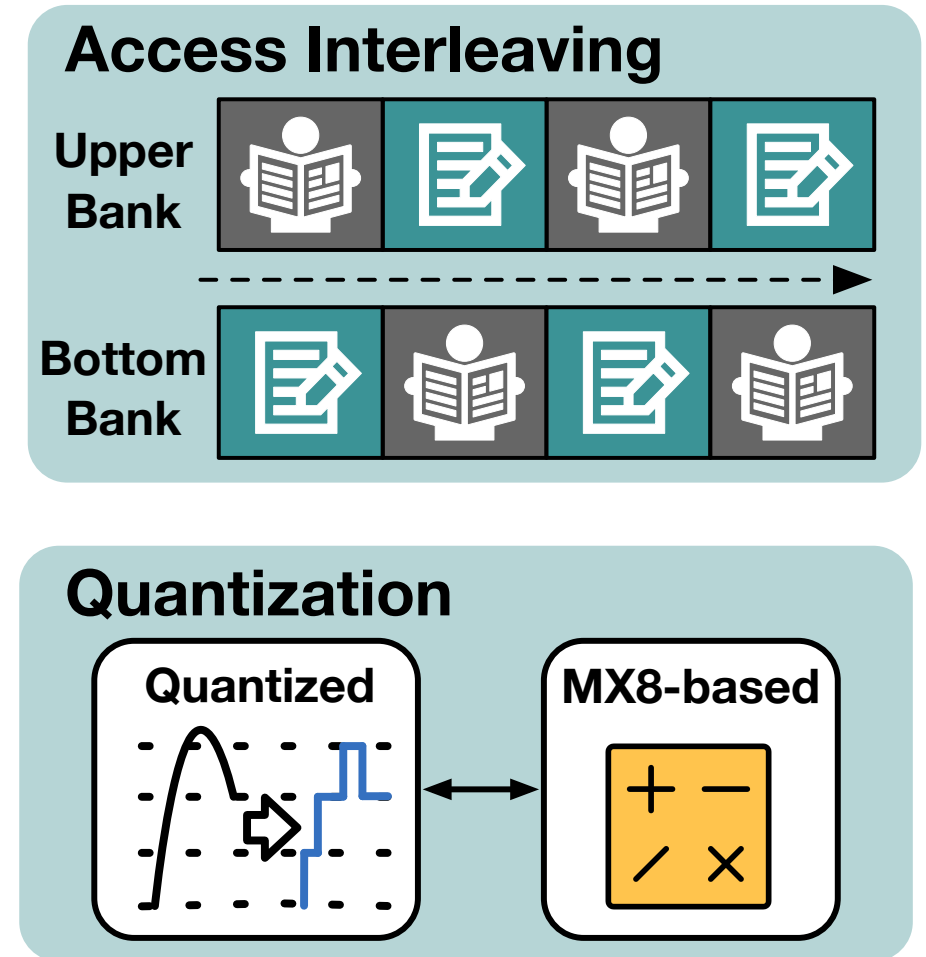- **State update operations dominate**
  - Due to lack of parameter reuse, state updates **cannot** be efficiently batched
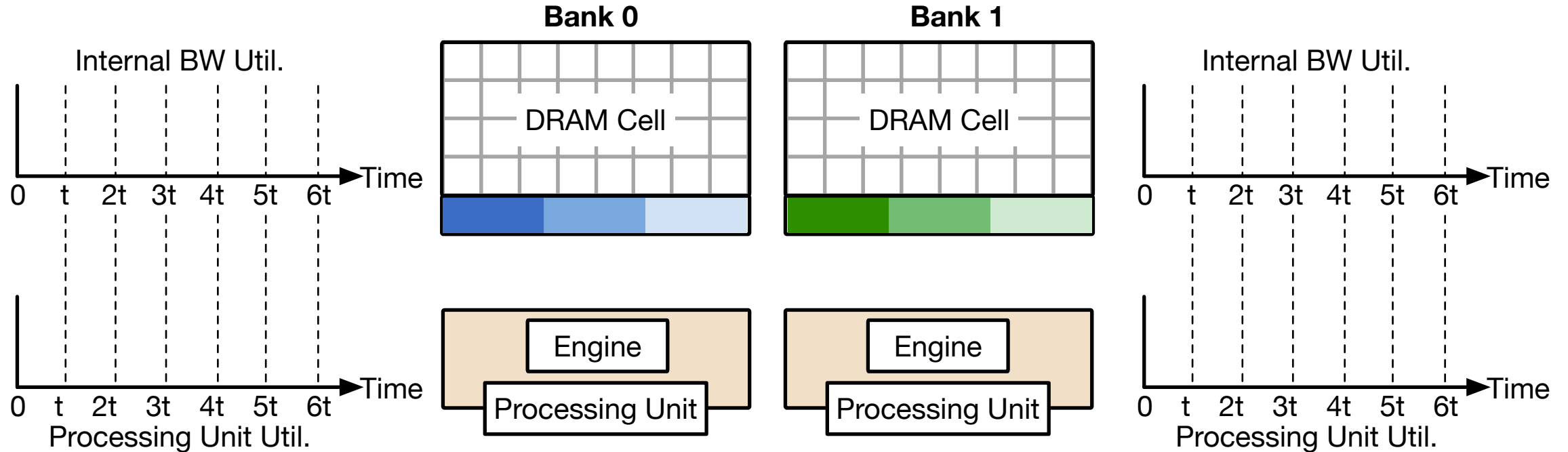  - Unlike GEMM, they have low arithmetic intensity, thus **memory-bound**

# Characterizing Post-Transformers

## Latency Breakdown

State Update  Discretization
Causal Conv  GEMM  Others

100

| Batch ➤ | 32 | 64 | 128 | 32 | 64 | 128 | 32 | 64 | 128 | 32 | 64 | 128 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model ➤ | | RetNet | | | GLA | | | HGRN2 | | | Mamba-2 | |

## Roofline Analysis

State Update  Attention
GEMM  GEMM

Memory  Compute

0.25  1  140  228
Arithmetic Intensity (FLOPS/byte)

**This memory bottleneck motivates us to develop PIMBA, which simultaneously leverages PIM and Quantization**

- **State update operations dominate**
  - Due to lack of parameter reuse, state updates **cannot** be efficiently batched
  - Unlike GEMM, they have low arithmetic intensity, thus **memory-bound**

KAIST School of Computing

CASYS | KAIST Computer Architecture & System Lab

# Pɪᴍʙᴀ Overview



**Diverse** operations needed!
We focus on optimizing area overhead

# State Updates in Existing PIM

# State Updates in Existing PIM

# State Updates in Existing PIM

# State Updates in Existing PIM

**Access Interleaving**

Quantization

# State Updates in Existing PIM

# State Updates in Existing PIM

Internal BW Util.

0  t  2t  3t  4t  5t  6t  Time

Processing Unit Util.

**Bank 0**

DRAM Cell

**WRITE**

Engine

Processing Unit

**Bank 1**

DRAM Cell

**WRITE**

Engine

Processing Unit

Internal BW Util.

0  t  2t  3t  4t  5t  6t  Time

Processing Unit Util.

# State Updates in Existing PIM

**Unlike dot-product operations,
state update operations require both reads and writes**

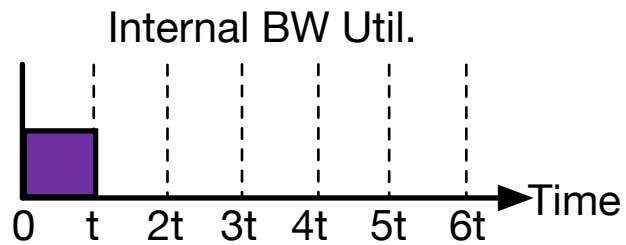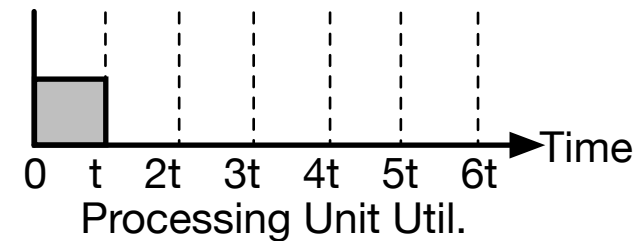# State Updates in Existing PIM

**This leads to underutilization of processing units during writes**
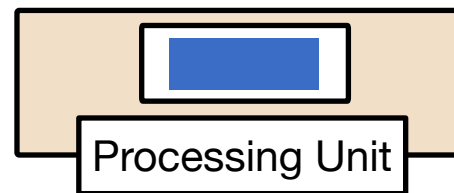
# Access Interleaving

# Access Interleaving

Internal BW Util.

Time
0  t  2t  3t  4t  5t  6t

**Bank 0**

DRAM Cell

**READ**

**Bank 1**

DRAM Cell

Internal BW Util.

Time
0  t  2t  3t  4t  5t  6t
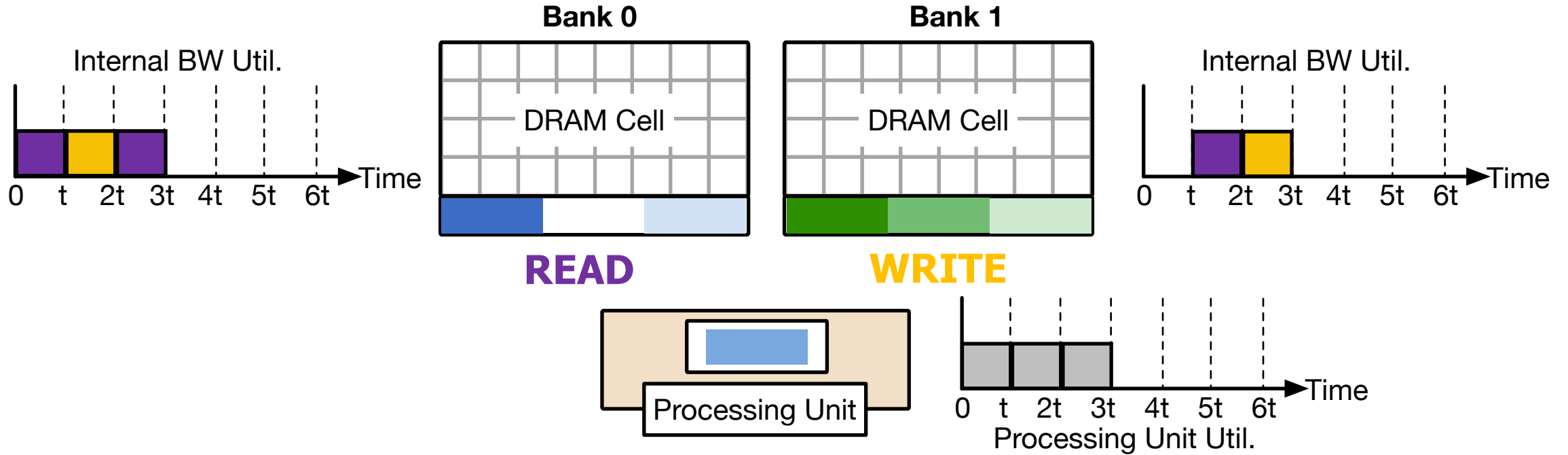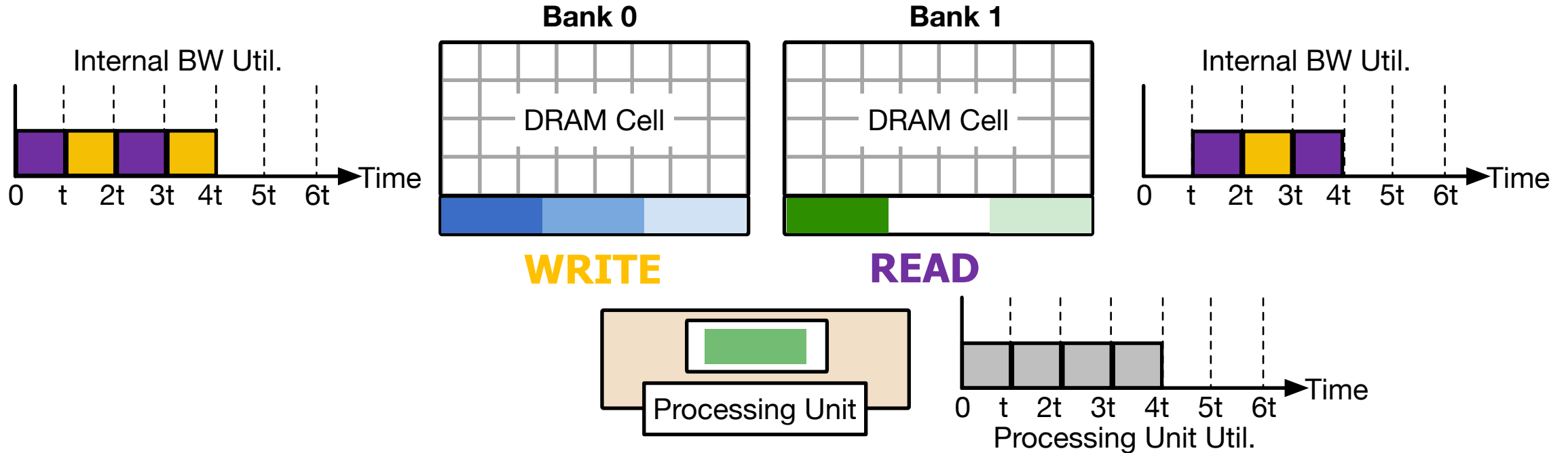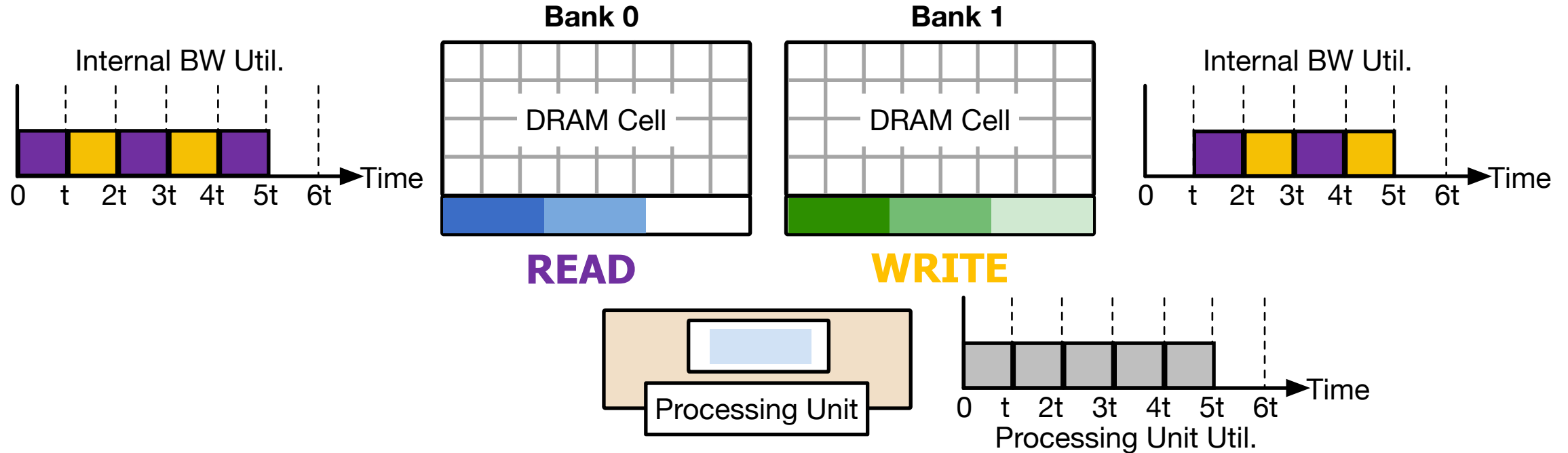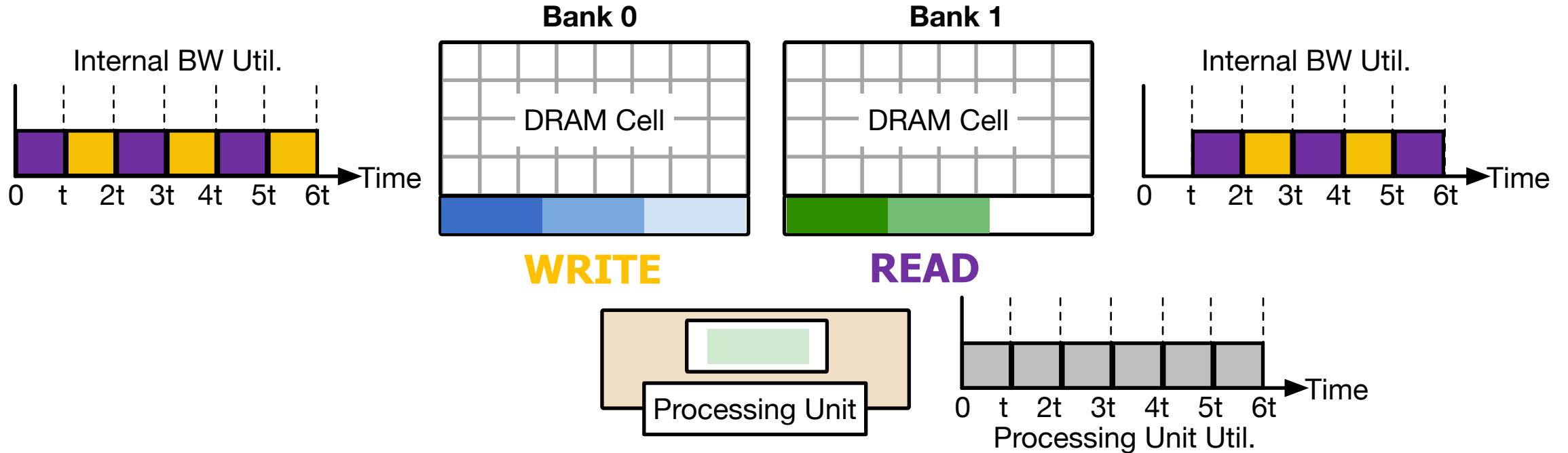
Processing Unit
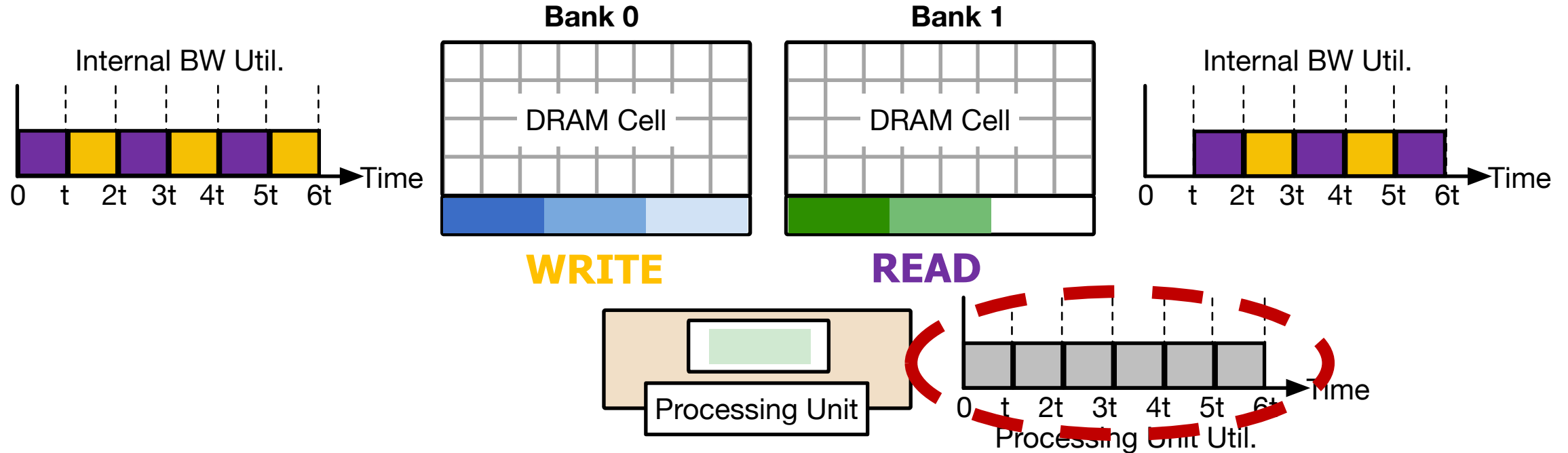
Time
0  t  2t  3t  4t  5t  6t
Processing Unit Util.

# Access Interleaving

# Access Interleaving

# Access Interleaving

# Access Interleaving

# Access Interleaving

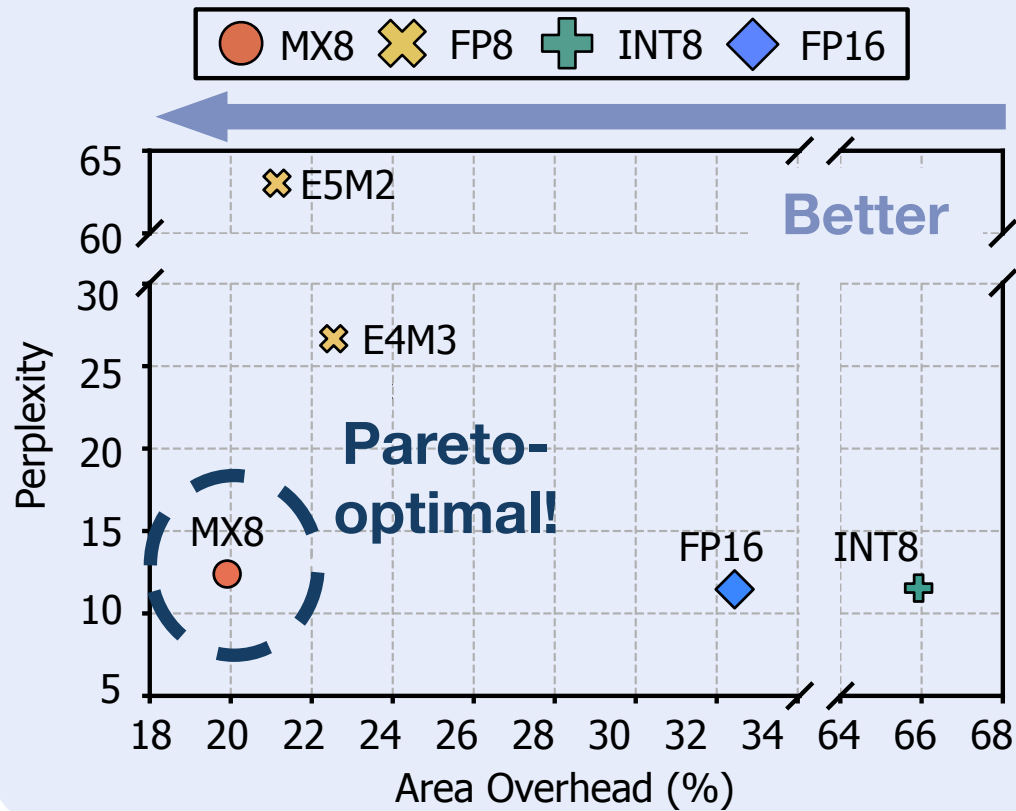**The processing units are now fully utilized,
which in turn reduces the area overhead by half
while maintaining the same throughput**

# MX-based Processing Element

- **MX operations require only simple addition, multiplication, and logic operations**

# Experimental Methodology

- **Baselines**
  - **GPU**: NVIDIA A100 GPUs
  - GPU w/ Quantization (**GPU+Q**): A100 + 8-bit quantized state
  - GPU w/ HBM-PIM (**GPU+PIM**): A100 + Samsung HBM-PIM

- **Models**
  - RetNet, GLA, HGRN2, Mamba2, Zamba2
  - small-scale (2.7B, 7B)
  - large-scale (70B)

- **Simulation**
  - GPU: extends AttAcc system simulator
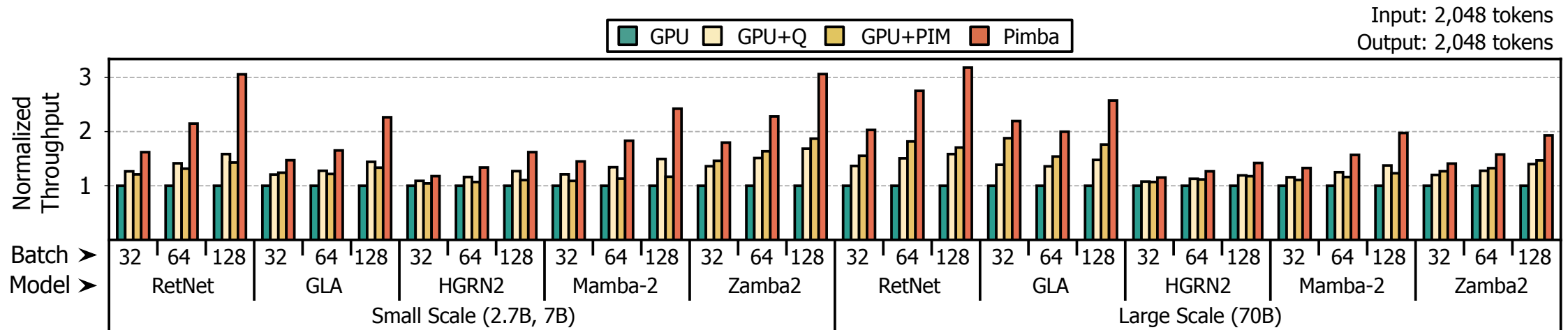  - PIM: extends cycle-accurate Ramulator2

# Throughput Results



- **PɪMBA** achieves **14.6×** faster state update operations compared to GPU
- **PɪMBA** delivers **up to 4.1×** higher decoding throughput compared to GPU

# More Results on Paper

- Accuracy evaluation

- Performance improvements on attention-based transformers

- Decode phase latency breakdown

- Energy consumption

- RTL area and power overhead

- Comparison with existing PIM-based LLM serving system

- General adoption of PIMBA

# Conclusion

- **PIMBA**
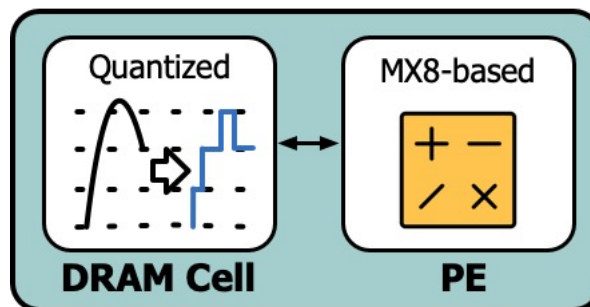  - An efficient PIM-based post-transformer acceleration solution

- **Contributions**
  - We conduct first comprehensive study of post-transformer LLMs
  - We analyze unique characterizations of post-transformer LLMs
  - We propose novel access interleaving strategy and quantization-based PIM

## Access Interleaving



Upper Bank

Bottom Bank

Time

## Quantization



Quantized — MX8-based

DRAM Cell — PE

## Code