# AttAcc! Unleashing the Power of PIM for Batched Transformer-based Generative Model Inference
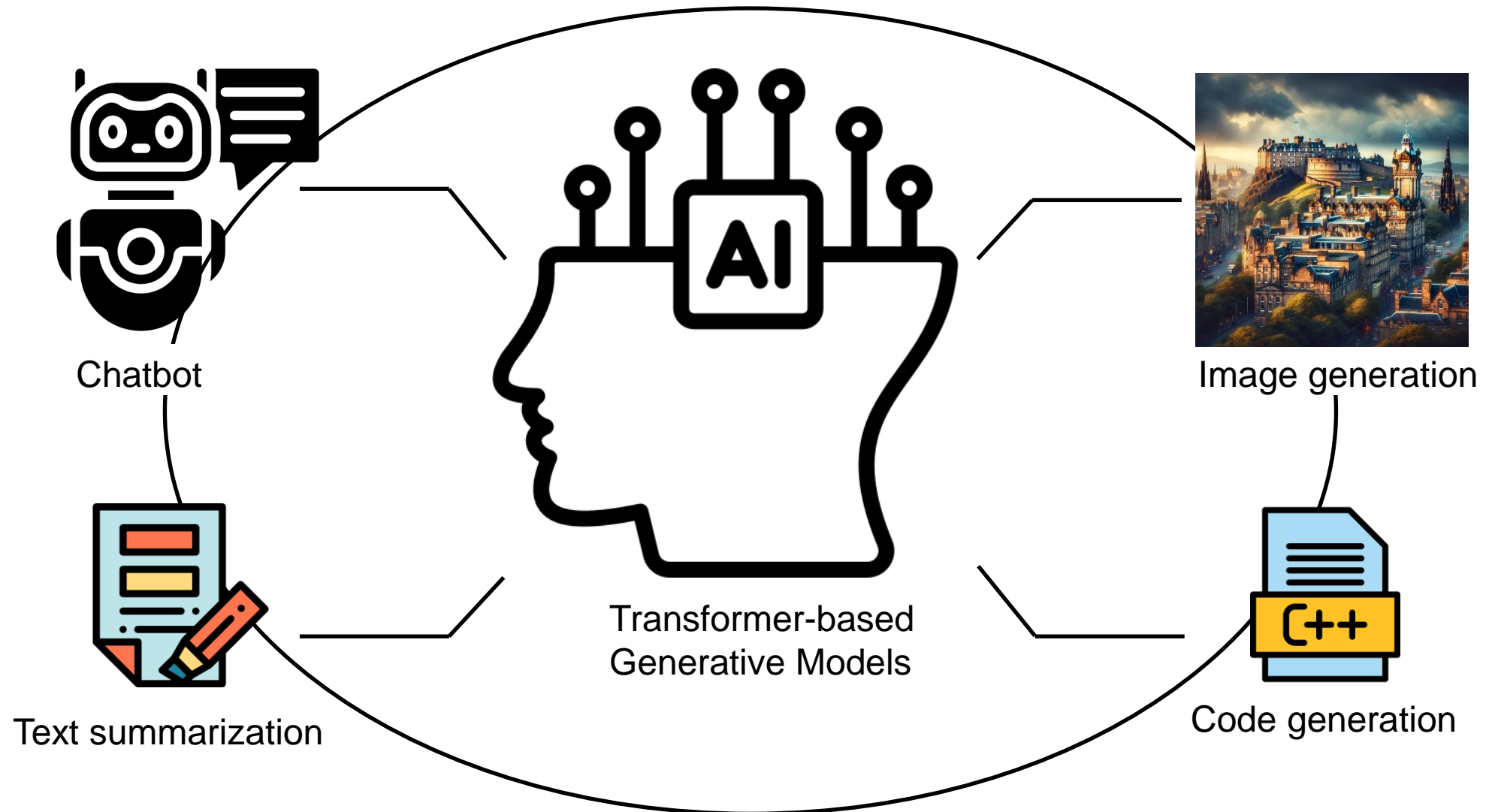
Jaehyun Park*[†], Jaewan Choi*[†], **Kwanhee Kyung**[†], Michael Jaemin Kim[†],

Yongsuk Kwon[†], Nam Sung Kim[‡], Jung Ho Ahn[†]

[†] Seoul National University, [‡] University of Illinois Urbana Champaign
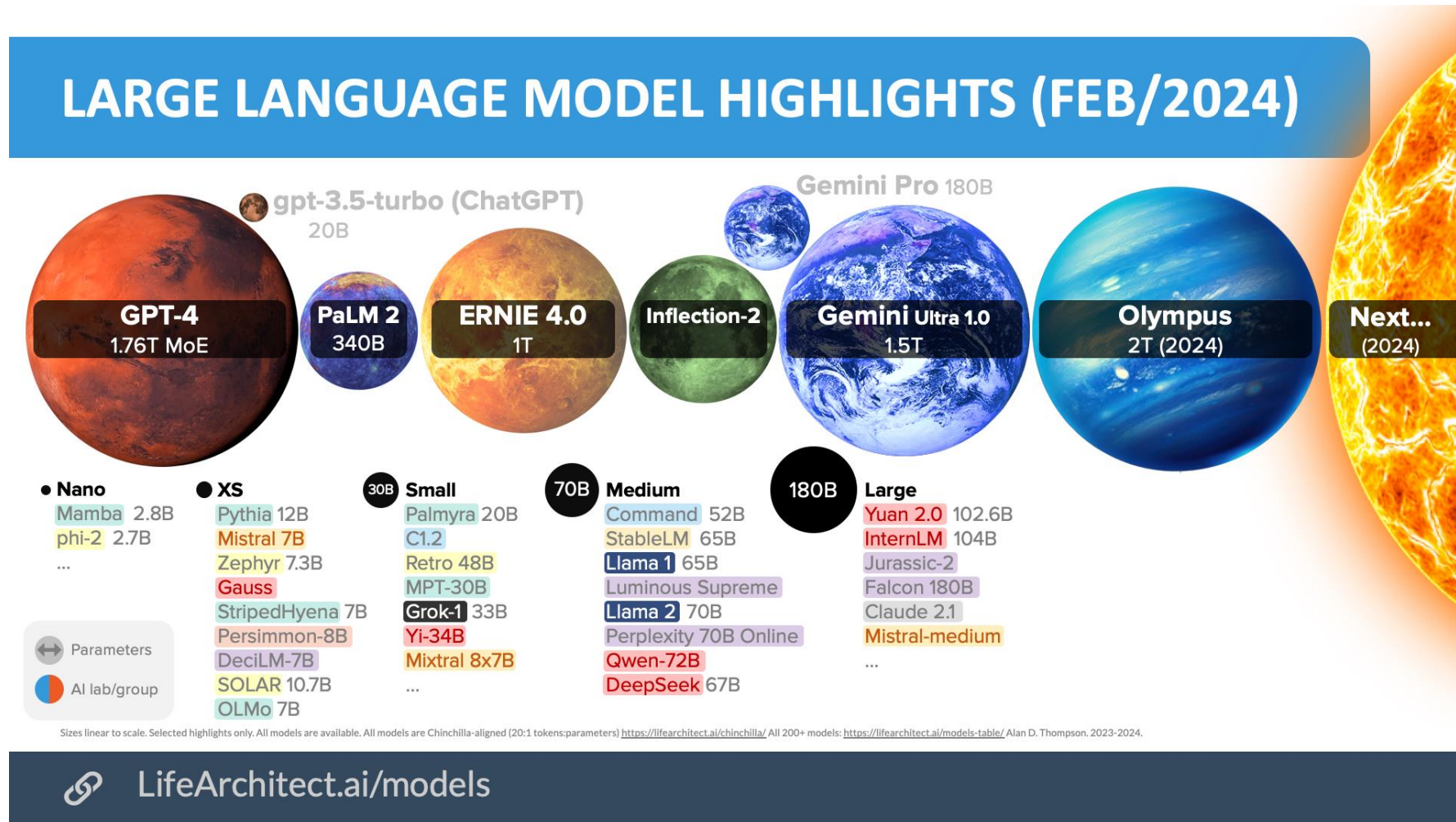
* Equally contributed

Presenter: Kwanhee Kyung (kwanhee.kyung@scale.snu.ac.kr)
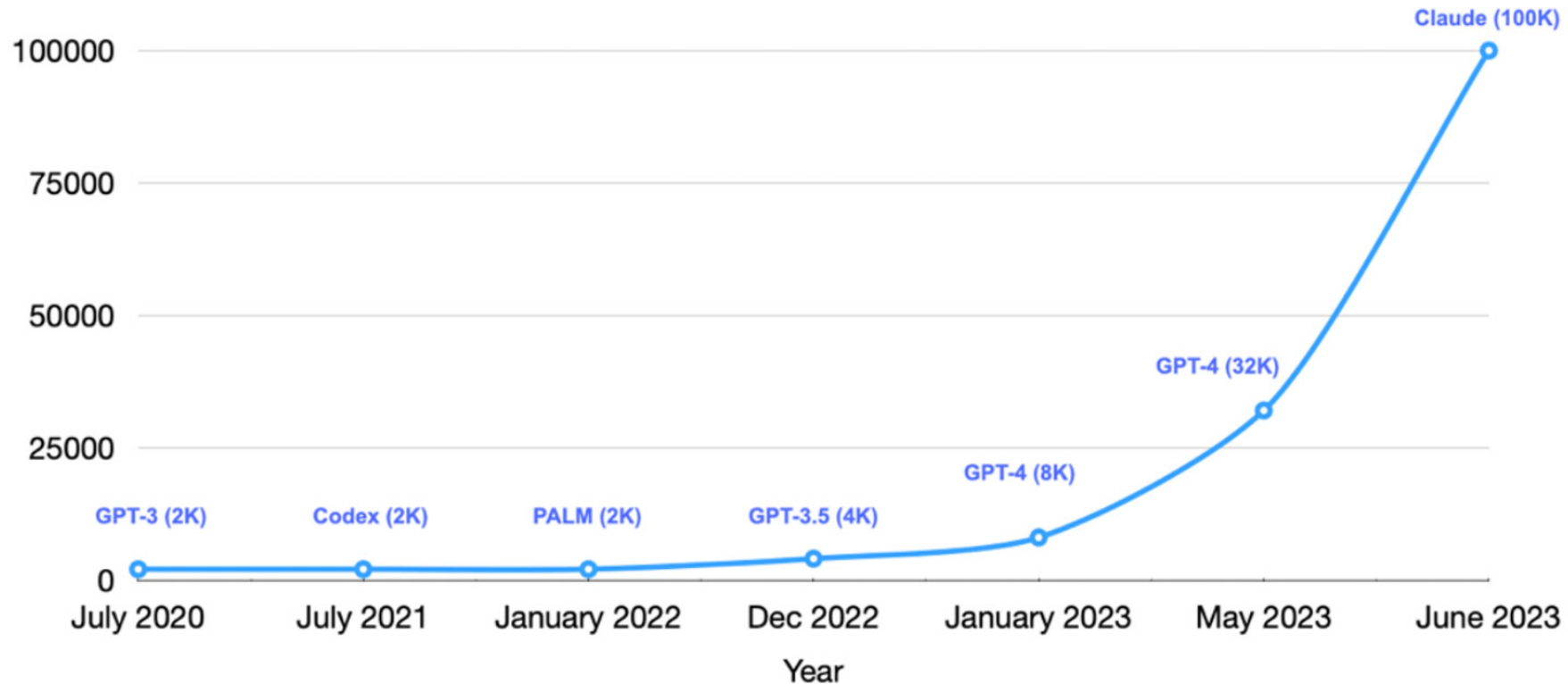
# Why Transformer-based Generative Model (TbGM) Inference?



Chatbot

Image generation

Text summarization

Transformer-based
Generative Models

Code generation

# Why Transformer-based Generative Model (TbGM) Inference?

- Model size



LARGE LANGUAGE MODEL HIGHLIGHTS (FEB/2024)

gpt-3.5-turbo (ChatGPT) 20B

Gemini Pro 180B

GPT-4 1.76T MoE

PaLM 2 340B

ERNIE 4.0 1T

Inflection-2

Gemini Ultra 1.0 1.5T

Olympus 2T (2024)

Next... (2024)

| ● Nano | ● XS | 30B Small | 70B Medium | 180B Large |
|---|---|---|---|---|
| Mamba 2.8B | Pythia 12B | Palmyra 20B | Command 52B | Yuan 2.0 102.6B |
| phi-2 2.7B | Mistral 7B | C1.2 | StableLM 65B | InternLM 104B |
| ... | Zephyr 7.3B | Retro 48B | Llama 1 65B | Jurassic-2 |
| | Gauss | MPT-30B | Luminous Supreme | Falcon 180B |
| | StripedHyena 7B | Grok-1 33B | Llama 2 70B | Claude 2.1 |
| | Persimmon-8B | Yi-34B | Perplexity 70B Online | Mistral-medium |
| | DeciLM-7B | Mixtral 8x7B | Qwen-72B | ... |
| | SOLAR 10.7B | ... | DeepSeek 67B | |
| | OLMo 7B | | | |

Parameters

AI lab/group

Sizes linear to scale. Selected highlights only. All models are available. All models are Chinchilla-aligned (20:1 tokens:parameters) https://lifearchitect.ai/chinchilla/ All 200+ models: https://lifearchitect.ai/models-table/ Alan D. Thompson. 2023-2024.

🔗 LifeArchitect.ai/models

# Why Transformer-based Generative Model (TbGM) Inference?

- Sequence length ($L$) supported by TbGM



https://hazyresearch.stanford.edu/blog/2023-03-27-long-learning

# TbGM Inference



Summarization
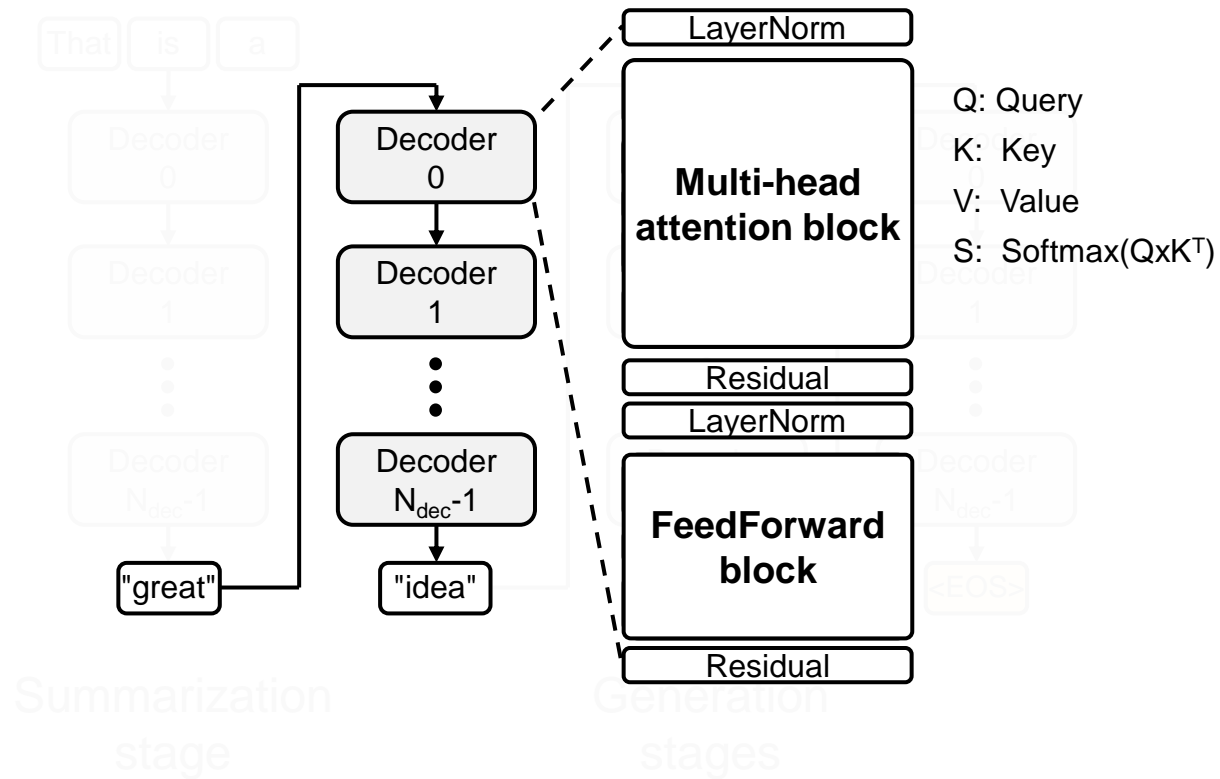(Sum) stage
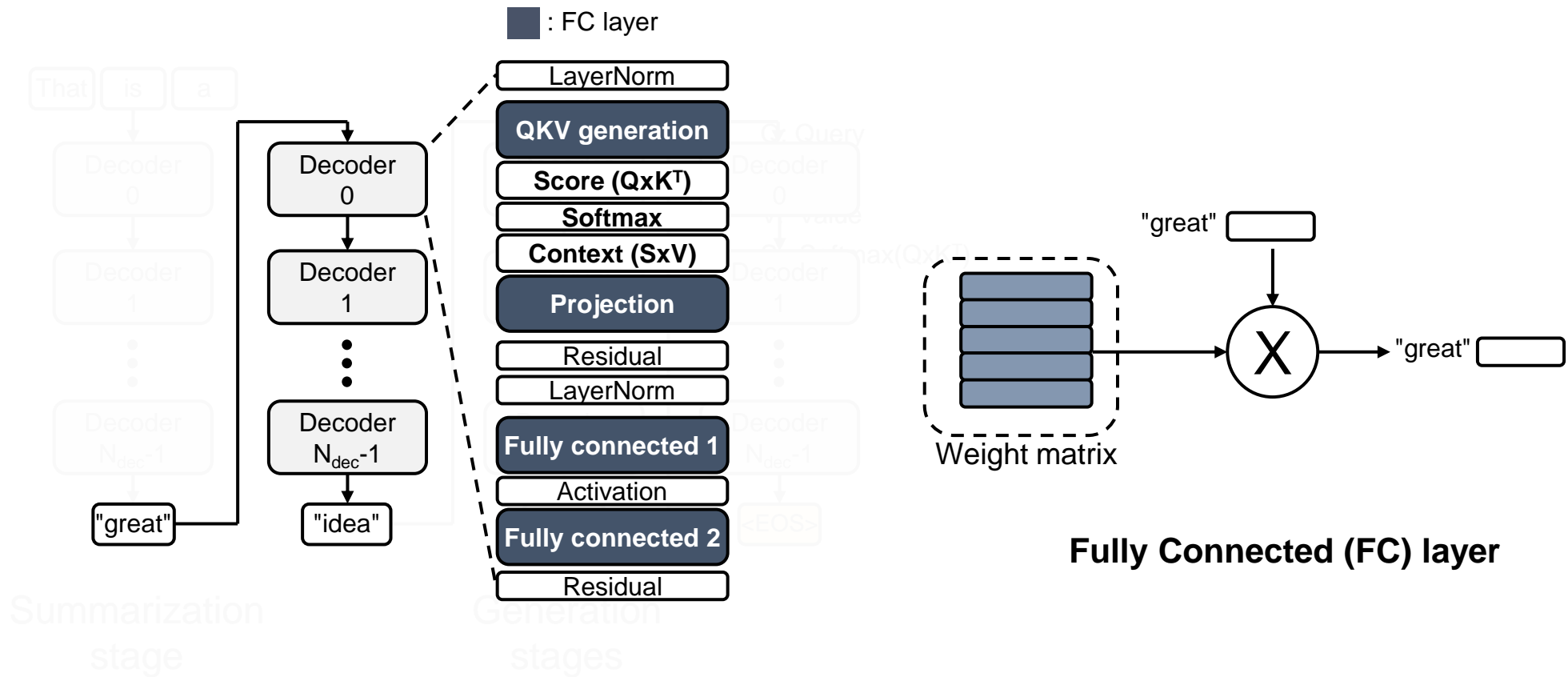(=Prefill phase)

Generation
(Gen) stages
(=Decode phase)

# Characteristics of the Gen Stage

# Characteristics of the Gen Stage

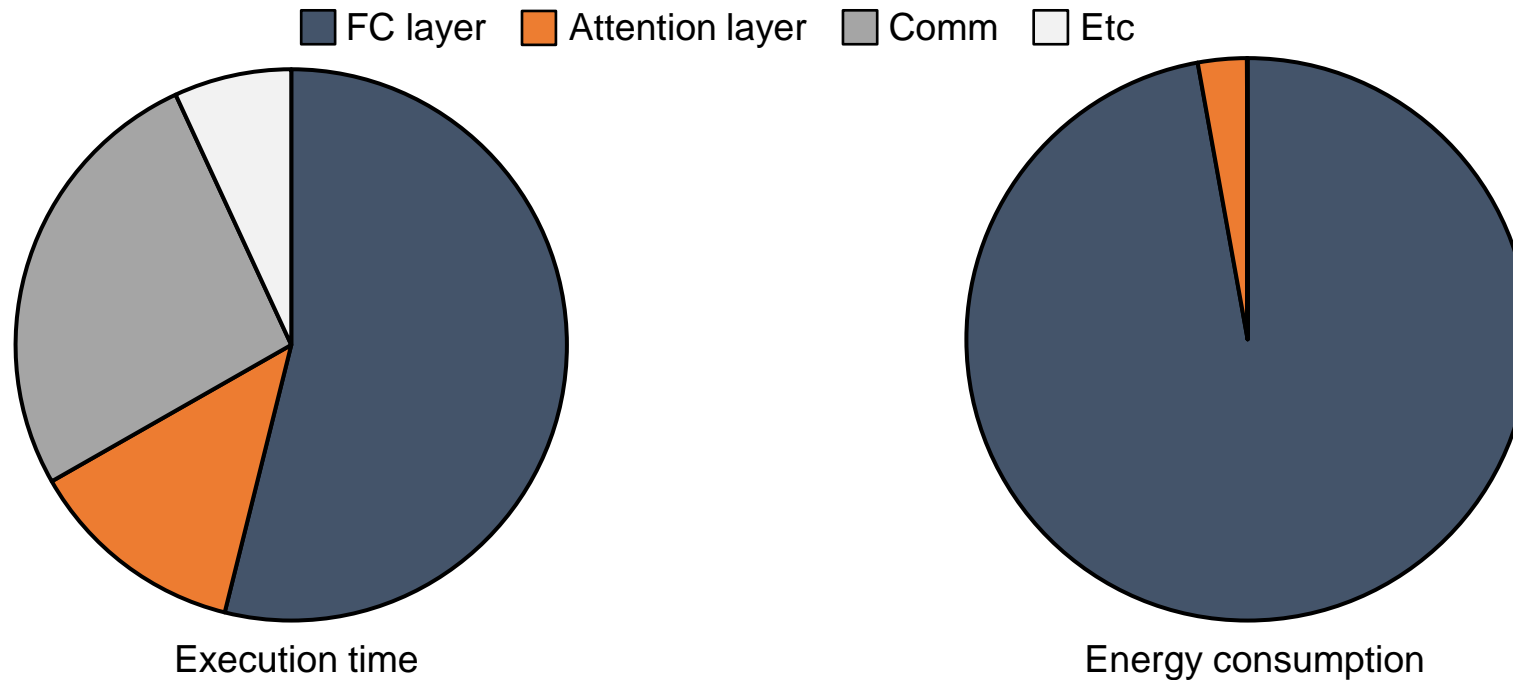- **FC layers** are all **general matrix-vector multiplications (GEMVs)**



■ : FC layer

LayerNorm
**QKV generation**
**Score (QxK$^T$)**
**Softmax**
**Context (SxV)**
**Projection**
Residual
LayerNorm
**Fully connected 1**
Activation
**Fully connected 2**
Residual

Decoder 0
Decoder 1
Decoder N$_{dec}$-1

"great"
"idea"

Weight matrix

"great"

X

"great"

**Fully Connected (FC) layer**

# Characteristics of the Gen Stage

- The **attention layer** also has **GEMVs** (GEMV$_{score}$ and GEMV$_{context}$)



Attention layer

# Prior TbGM Accelerators

- Many prior works [1,2,3] for TbGM accelerator focused on accelerating **FC layers**
  - FC layers account for a significant portion of execution time and energy consumption



Execution time

Energy consumption

Execution time and energy consumption breakdown of TbGM inference
(GPT3-175B on DGX-A100 with HBM3 , $L_{in}$=2048, $L_{out}$=128)

[1] S Lee et al., "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology," ISCA, 2021
[2] S Hong et al., "DFX: A Low-latency Multi-FPGA Appliance for Accelerating Transformer-based Text Generation," MICRO, 2022.
[3] D Kwon et al., "A 1ynm 1.25V 8Gb 16Gb/s/Pin GDDR6-Based Accelerator-in-Memory Supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep Learning Application," JSSC, 2023
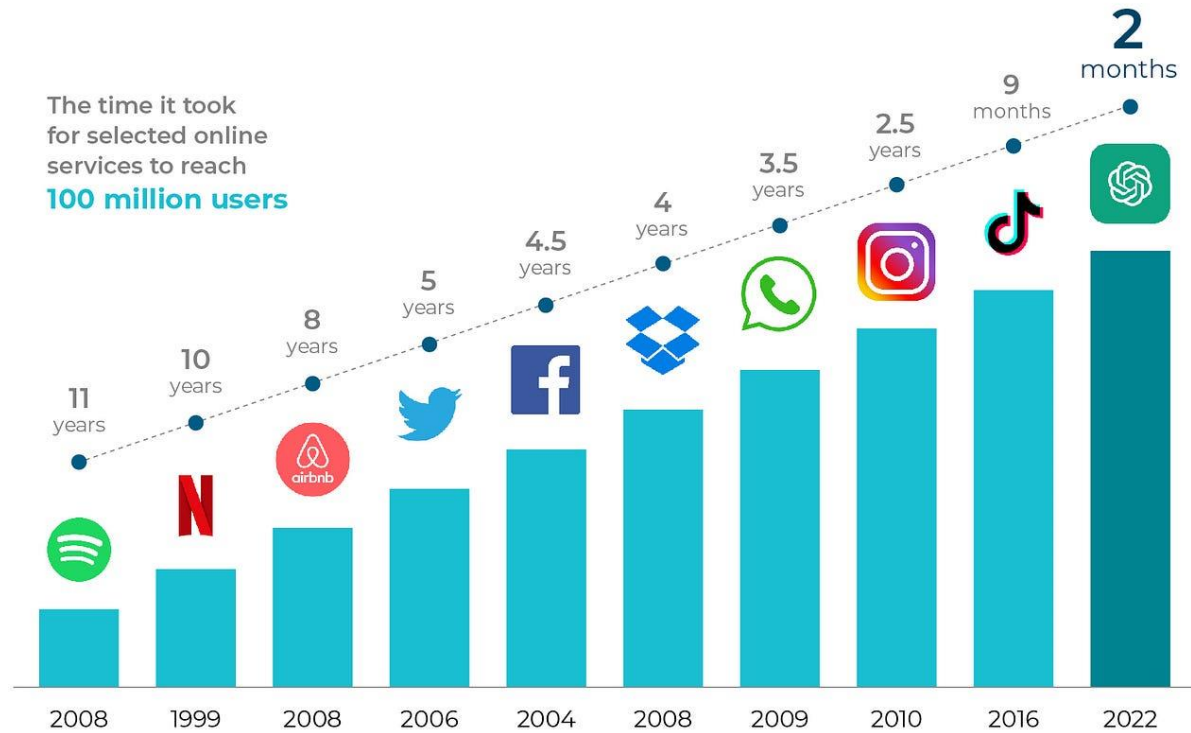
# Prior TbGM Accelerators

- Many prior works [1,2,3] for TbGM accelerator focused on accelerating **FC layers**
  - FC layers account for a significant portion of execution time and energy consumption

■ FC layer  ■ Attention layer  ■ Comm  □ Etc

**However, a key assumption so far is that the batch size is 1**

Execution time

Energy consumption

Execution time and energy consumption breakdown of TbGM inference
(GPT3-175B on DGX-A100 with HBM3 , $L_{in}$=2048, $L_{out}$=128)

[1] S Lee et al., "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology," ISCA, 2021
[2] S Hong et al., "DFX: A Low-latency Multi-FPGA Appliance for Accelerating Transformer-based Text Generation," MICRO, 2022.
[3] D Kwon et al., "A 1ynm 1.25V 8Gb 16Gb/s/Pin GDDR6-Based Accelerator-in-Memory Supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep Learning Application," JSSC, 2023

# Outline

- We discover that **the attention layer**
  - becomes more important in batched TbGM inference
  - poses several challenges in conventional systems

- To address these challenges, we propose
  - processing-in-memory (PIM)-based accelerator (*AttAcc*) for the attention layer
  - heterogeneous system with *AttAcc* and xPU for end-to-end TbGM inference
  - optimizations that improve utilization of the heterogeneous system

# Why Large Batch Size Matters

- Ensure sufficient requests from increased TbGM inference usage



Chat-GPT sprints to 100 million users

The time it took for selected online services to reach 100 million users

| Service | Time | Year |
| --- | --- | --- |
| Spotify | 11 years | 2008 |
| Netflix | 10 years | 1999 |
| airbnb | 8 years | 2008 |
| Twitter | 5 years | 2006 |
| Facebook | 4.5 years | 2004 |
| Dropbox | 4 years | 2008 |
| WhatsApp | 3.5 years | 2009 |
| Instagram | 2.5 years | 2010 |
| TikTok | 9 months | 2016 |
| ChatGPT | 2 months | 2022 |

Source: World of Statistics

# Why Large Batch Size Matters

- Ensure sufficient requests from increased TbGM inference usage

- Batching technique for TbGM inference [1] enables high throughput and energy efficiency



Throughput and energy consumption per output token of TbGM inference (GPT-3 175B, DGX A100 with HBM3)

[1] G Yu et al., "Orca: A Distributed Serving System for Transformer-Based Generative Models," OSDI, 2022

# What If Batch Size Increases?

- **FC layers** become more **compute-intensive**
    - Weight matrices are shared across different requests



: FC layers in a Gen stage

**Weights** reused across multiple requests

# What If Batch Size Increases?

- **The attention layer** is still **memory-intensive**
  - The attention layer has unique **KV matrices per request**
  - The arithmetic intensity **remains nearly 1** regardless of the batch size



**Key and Value (KV) matrices per request**

# Challenges of Large Batch Size on Conventional Systems



Throughput per output token of TbGM inference
(GPT-3 175B, DGX A100 with HBM3)

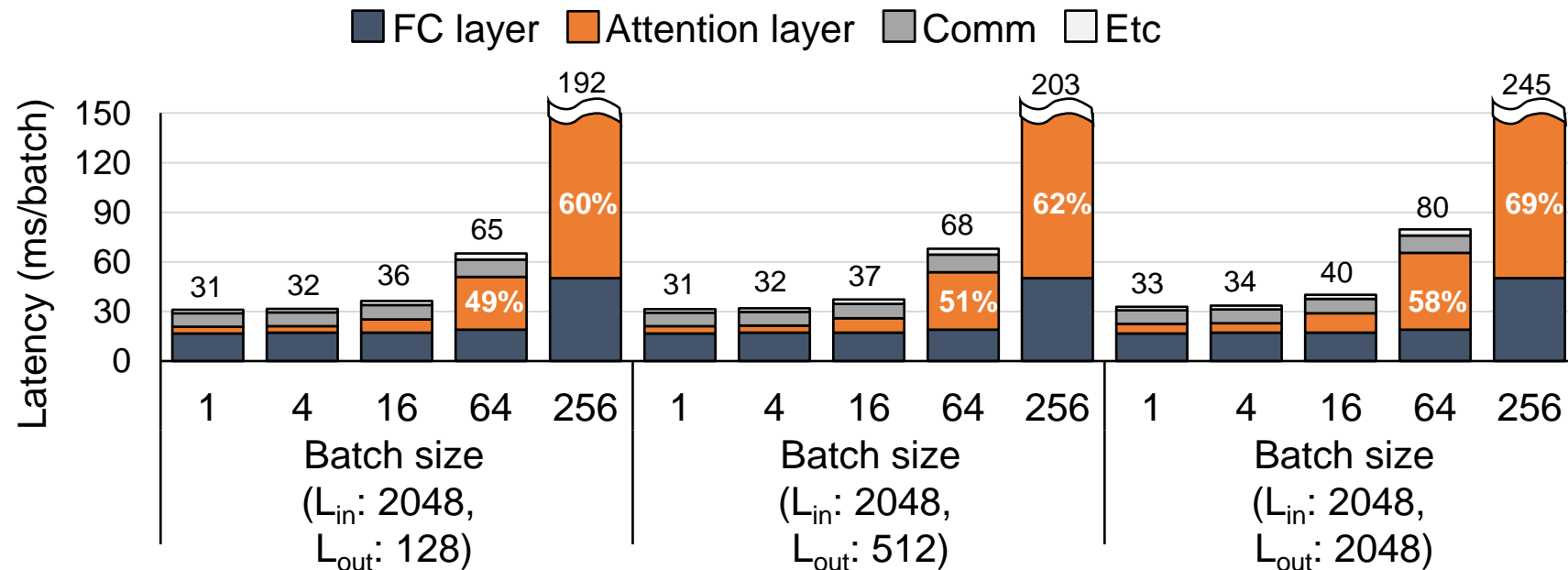# Challenges of Large Batch Size on Conventional Systems

- **Large memory capacity requirement** from KV matrices
  - KV matrices require more memory capacity in proportion to batch size.



Throughput and required memory capacity per output token of TbGM inference
(GPT-3 175B, DGX A100 with HBM3)

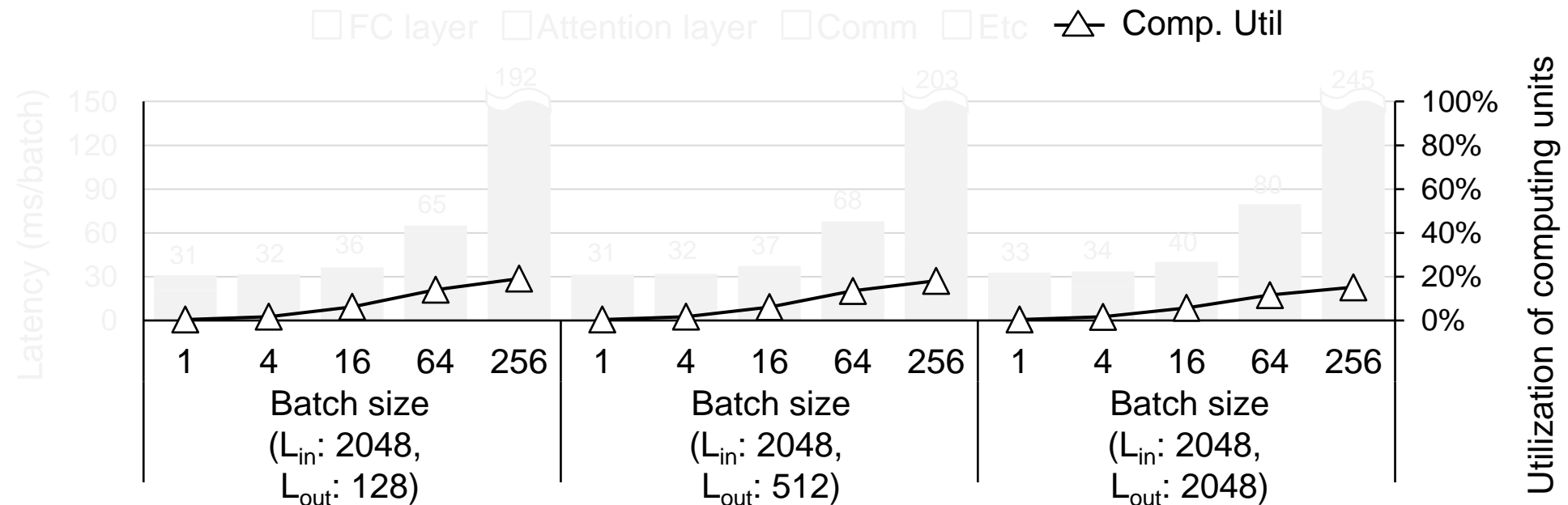# Challenges of Large Batch Size on Conventional Systems

- **Large memory capacity requirement** from KV matrices
  - KV matrices require more memory capacity in proportion to batch size.



Throughput and required memory capacity per output token of TbGM inference
(GPT-3 175B, DGX A100 with HBM3)

# Challenges of Large Batch Size on Conventional Systems

- **Long latency** of the attention layer
  - The latency of the attention layer increases linearly with batch size.
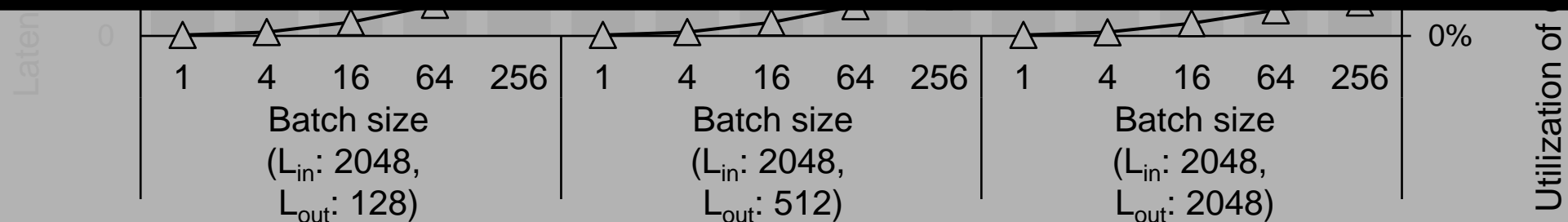  - It can limit batch sizes under service level objectives (SLOs).



The Gen stage time breakdown and compute utilization
(GPT-3 175B, DGX-unlimited memory capacity)

# Challenges of Large Batch Size on Conventional Systems

- **Low utilization of computing units** from attention layer



The Gen stage time breakdown and compute utilization
(GPT-3 175B, DGX-unlimited memory capacity)

# Challenges of Large Batch Size on Conventional Systems

- **Low utilization of computing units** from attention layer
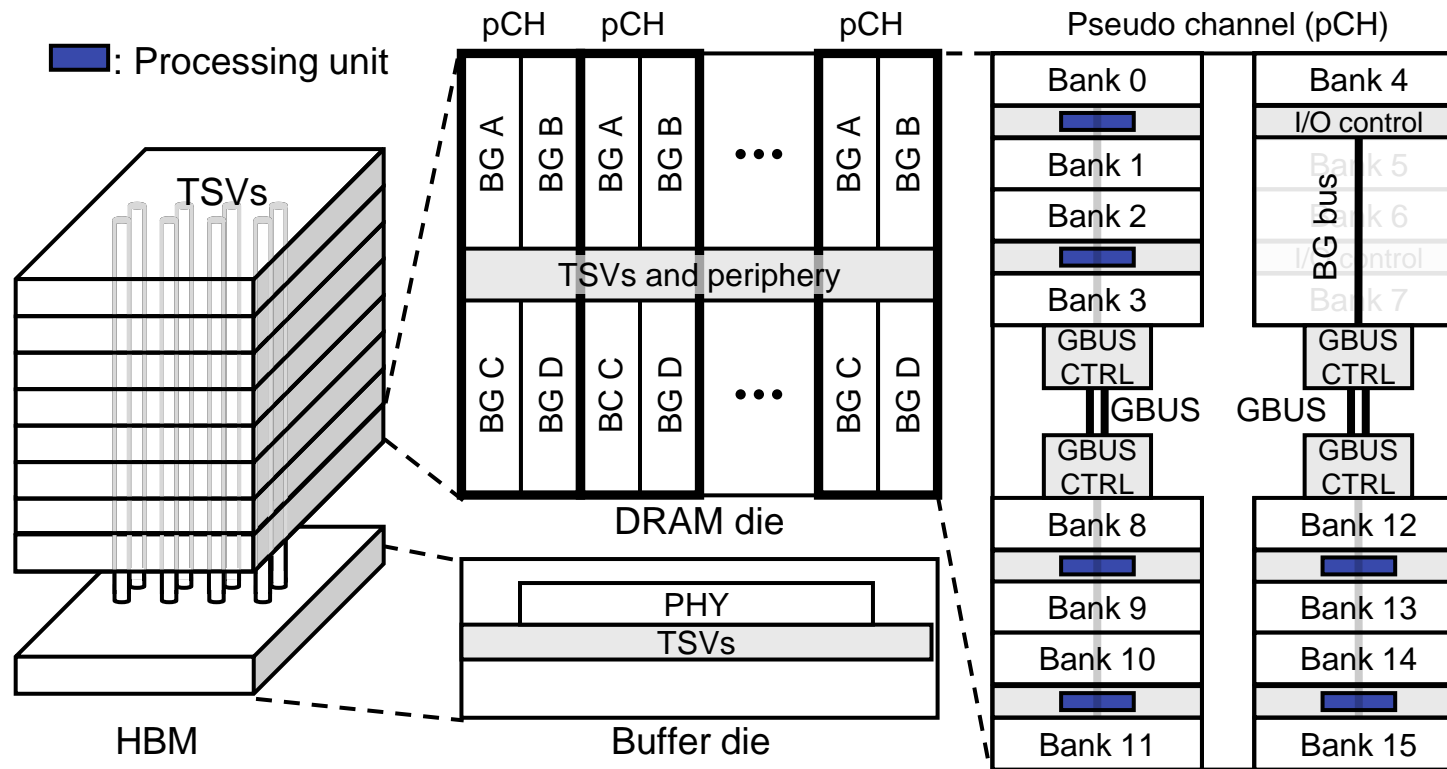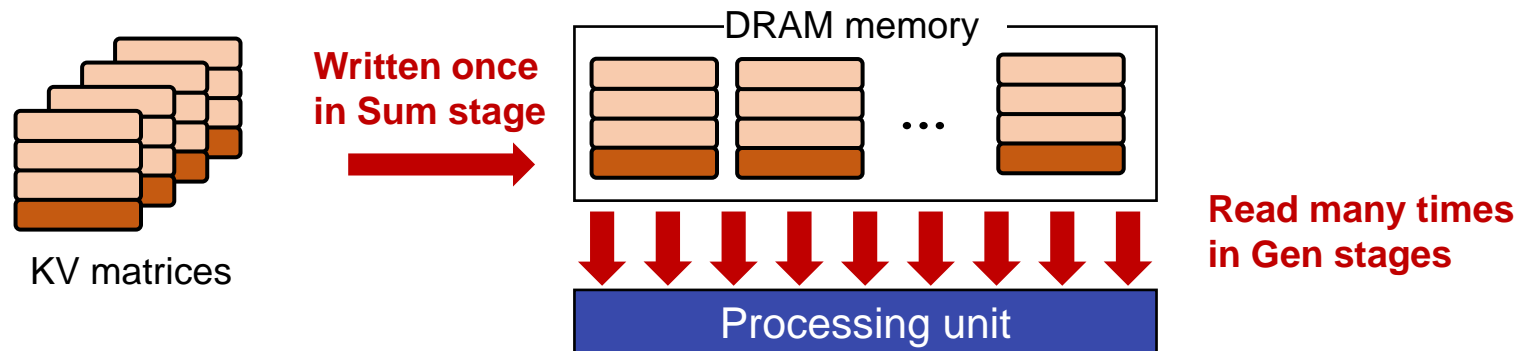
☐ FC layer ☐ Attention layer ☐ Comm ☐ Etc ─△─ Comp. Util

**We propose Processing-in-Memory (PIM)-based accelerator for the attention layer**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 4 | 16 | 64 | 256 | 1 | 4 | 16 | 64 | 256 | 0% |

Batch size
($L_{in}$: 2048,
$L_{out}$: 128)

Batch size
($L_{in}$: 2048,
$L_{out}$: 512)

Batch size
($L_{in}$: 2048,
$L_{out}$: 2048)

The Gen stage time breakdown and compute utilization
(GPT-3 175B, DGX-unlimited memory capacity)

# Processing-In-Memory (PIM)

- PIM exploits abundant internal bandwidth to processing units (PUs) closer to the memory.
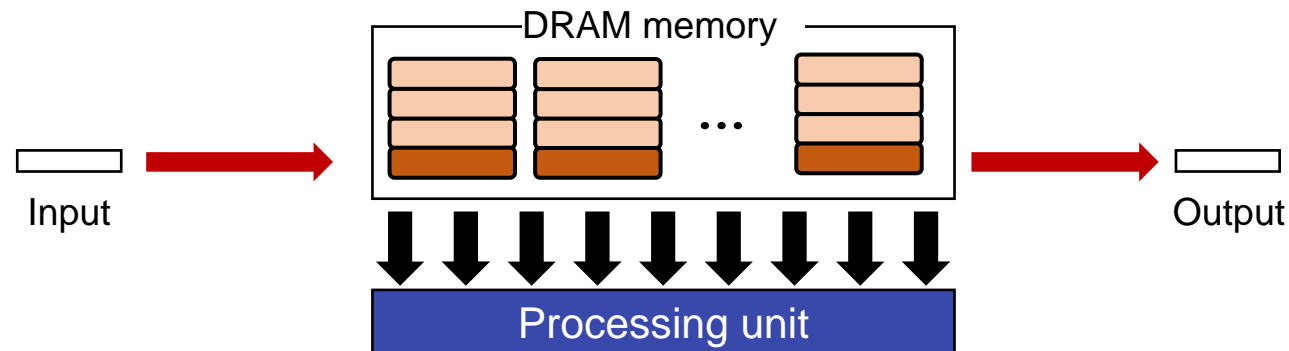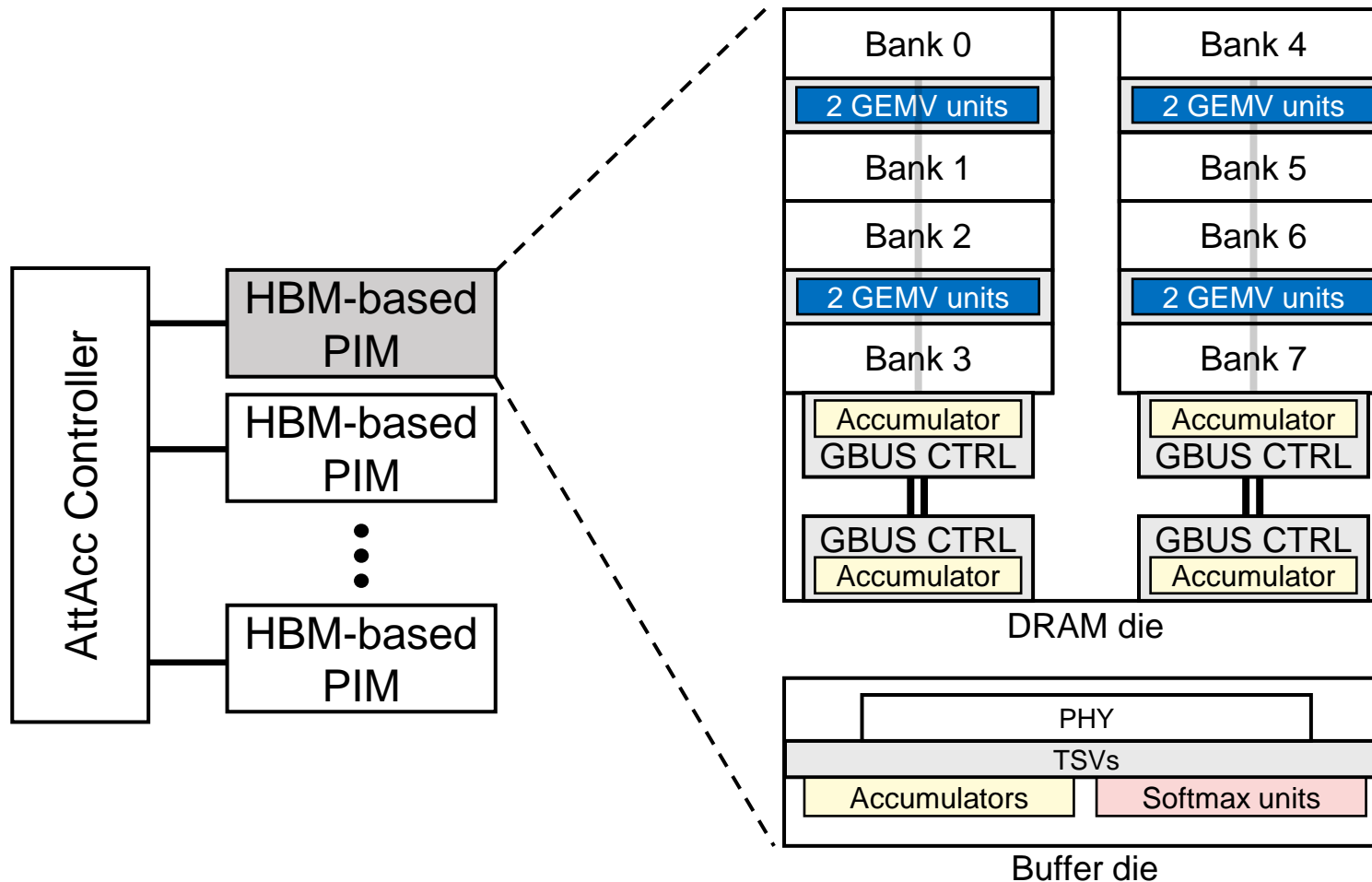
# Why Processing-In-Memory (PIM) for the Attention Layer?

- High memory bandwidth requirement
  - The attention layer is **memory-intensive GEMV**
  - The size of the KV matrices is **too large to be cached**

- Relatively low external bandwidth requirement
  - KV matrices are **written once** in the Sum stage and **read many** times in Gen stages

# Why Processing-In-Memory (PIM) for the Attention Layer?

- High memory bandwidth requirement
  - The attention layer is **memory-intensive GEMV**
  - The size of the KV matrices is **too large to be cached**

- Relatively low external bandwidth requirement
  - KV matrices are **written once** in the Sum stage and **read many** times in Gen stages
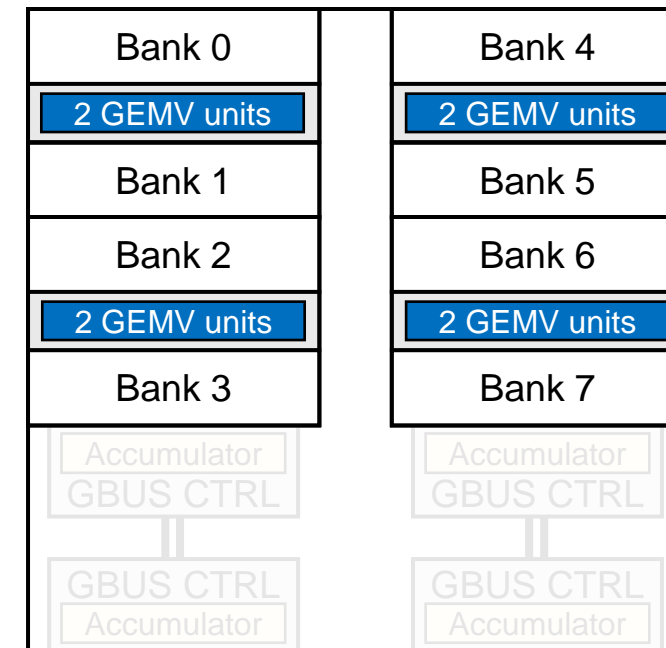  - The input and output of the attention layer are vectors that are **much smaller than KV matrices**.

# AttAcc: PIM-based Attention Accelerator

- We propose **AttAcc**, which consists of HBM-based PIMs and a controller

- HBM-based PIM has
    - GEMV units
    - Softmax unit
    - Accumulators

# AttAcc: GEMV Unit

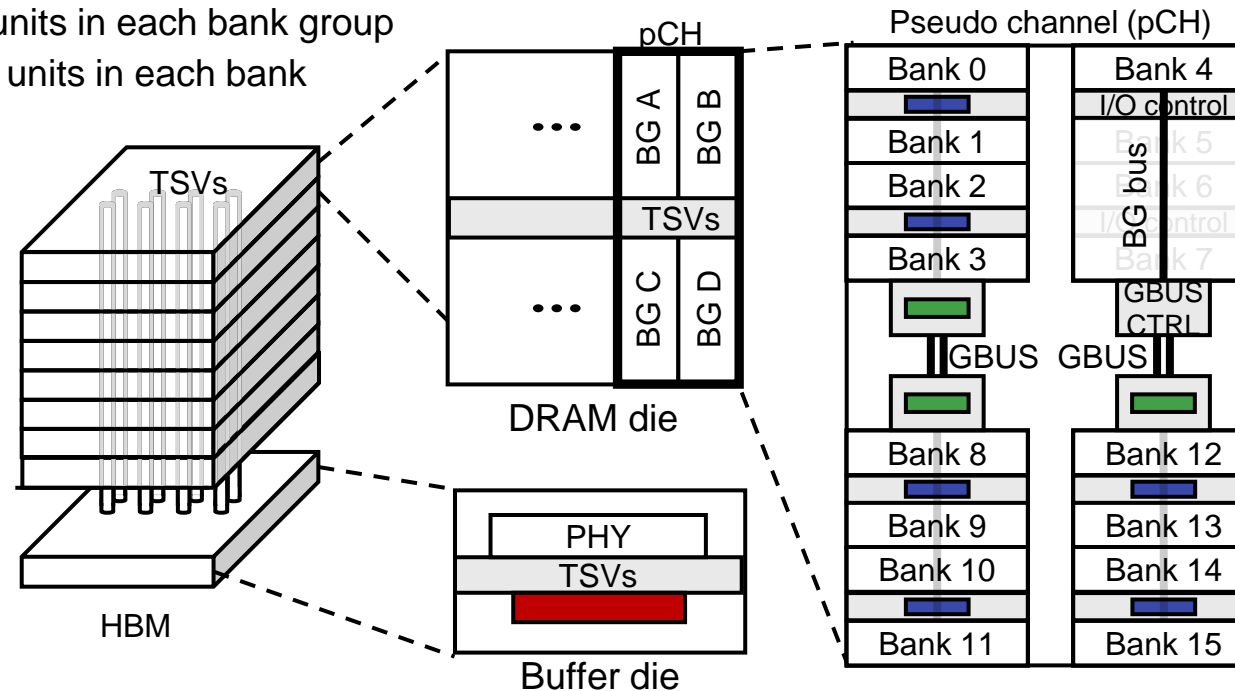- Placed on each bank similar to Samsung HBM-PIM [1] and Hynix AiM [2]

| Bank 0 | Bank 4 |
|---|---|
| **2 GEMV units** | **2 GEMV units** |
| Bank 1 | Bank 5 |
| Bank 2 | Bank 6 |
| **2 GEMV units** | **2 GEMV units** |
| Bank 3 | Bank 7 |
| Accumulator | Accumulator |
| GBUS CTRL | GBUS CTRL |
| GBUS CTRL | GBUS CTRL |
| Accumulator | Accumulator |

DRAM die

PHY

TSVs

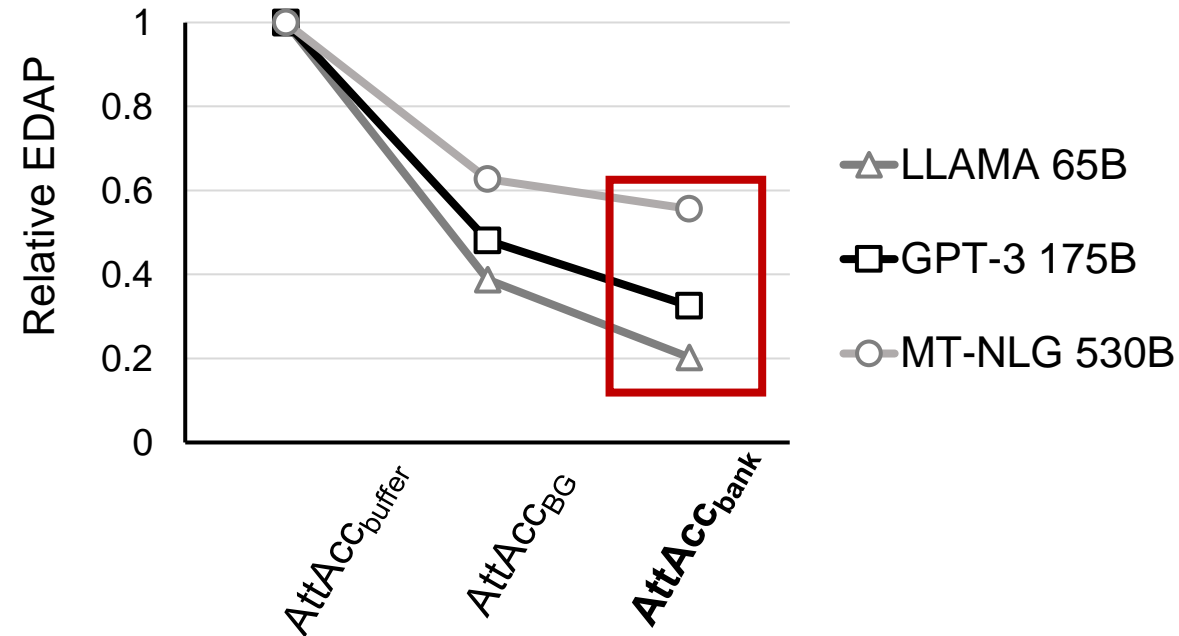Accumulators    Softmax units

Buffer die

[1] S Lee et al., "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology," ISCA, 2021
[2] D Kwon et al., "A 1ynm 1.25V 8Gb 16Gb/s/Pin GDDR6-Based Accelerator-in-Memory Supporting 1TFLOPS MAC Operation a...                 2023

# AttAcc: GEMV Unit

- Placed on each bank similar to Samsung HBM-PIM [1] and Hynix AiM [2]
  - *AttAcc$_{Buffer}$ vs AttAcc$_{BG}$ vs AttAcc$_{Bank}$*



AttAcc$_{buffer}$: GEMV units in each pCH of buffer die
AttAcc$_{BG}$: GEMV units in each bank group
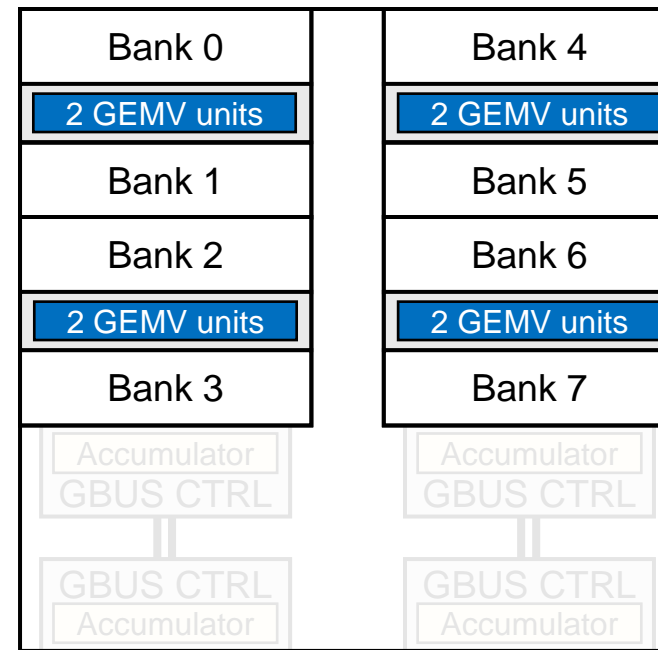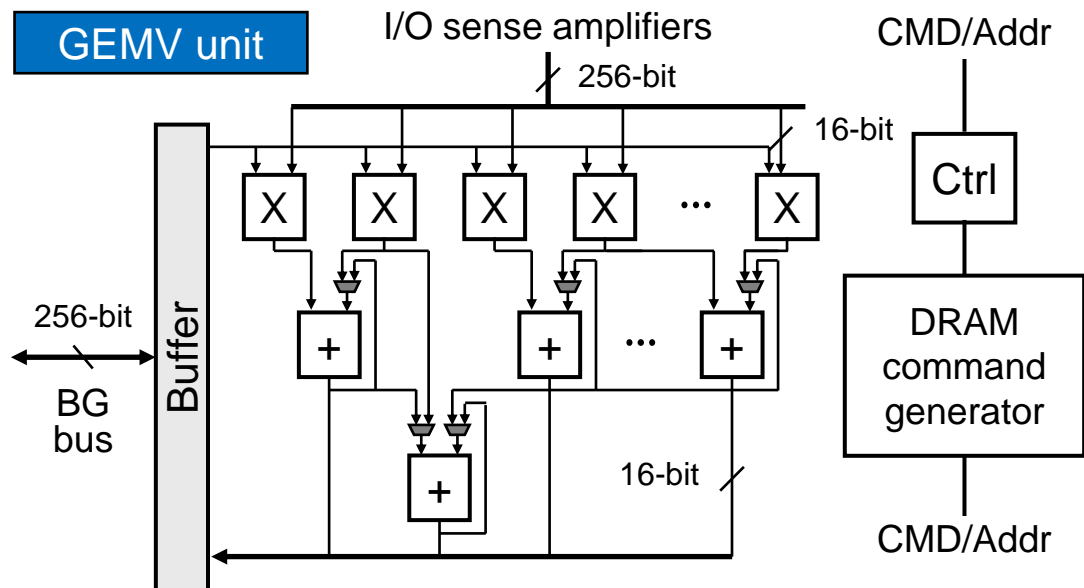AttAcc$_{bank}$: GEMV units in each bank

[1] S Lee et al., "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology," ISCA, 2021
[2] D Kwon et al., "A 1ynm 1.25V 8Gb 16Gb/s/Pin GDDR6-Based Accelerator-in-Memory Supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep Learning Application," JSSC, 2023

# AttAcc: GEMV Unit

- Placed on each bank similar to Samsung HBM-PIM [1] and Hynix AiM [2]
  - *AttAcc*$_{Buffer}$ *vs AttAcc*$_{BG}$ *vs AttAcc*$_{Bank}$

[1] S Lee et al., "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology," ISCA, 2021
[2] D Kwon et al., "A 1ynm 1.25V 8Gb 16Gb/s/Pin GDDR6-Based Accelerator-in-Memory Supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep Learning Application," JSSC, 2023

# AttAcc: GEMV Unit

- Placed on each bank similar to Samsung HBM-PIM [1] and Hynix AiM [2]
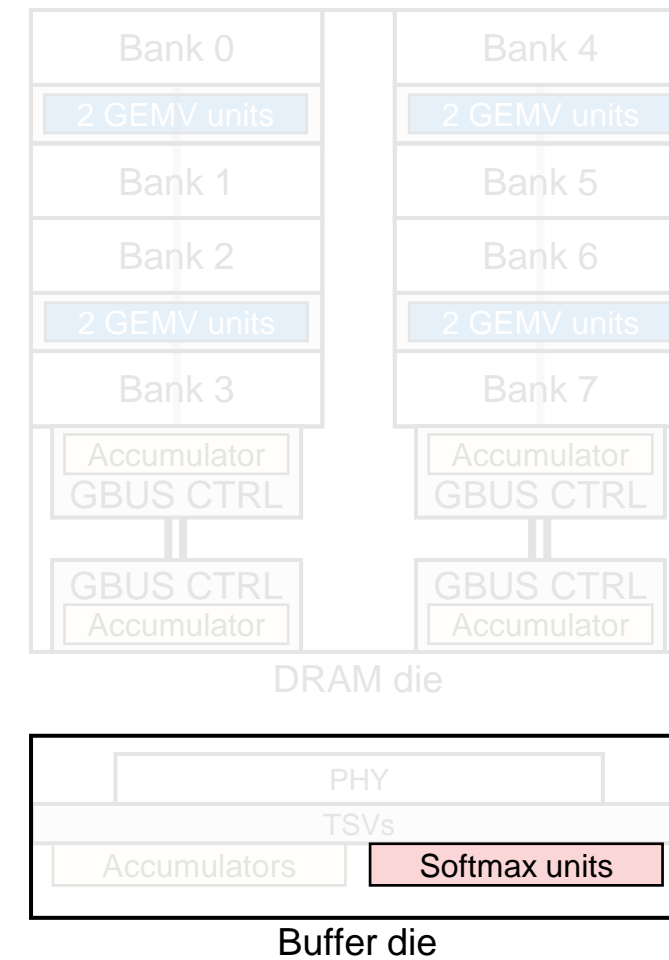- FP16 multipliers, FP16 adders, buffer for input vectors, and control unit.

# AttAcc: Softmax Unit
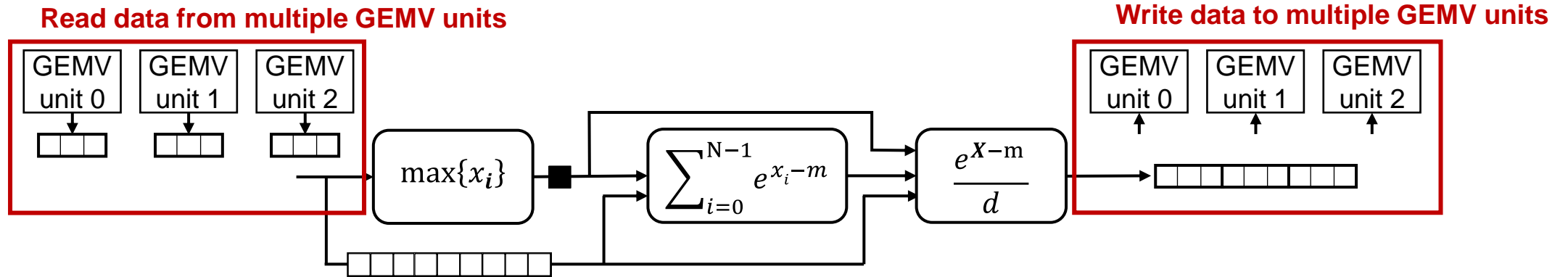
- Placed on buffer die

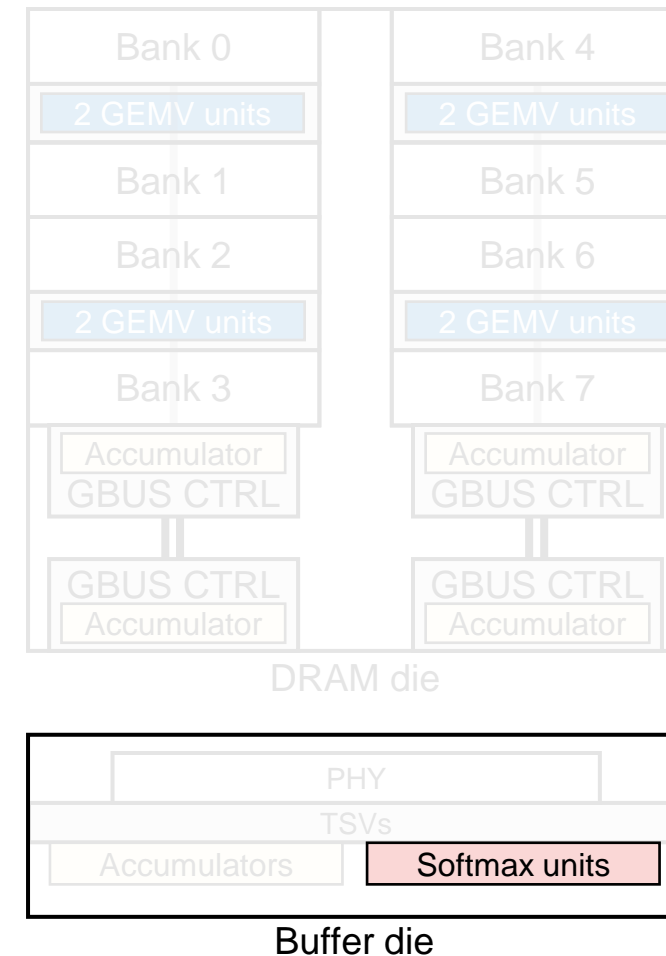| Bank 0 | Bank 4 |
|---|---|
| 2 GEMV units | 2 GEMV units |
| Bank 1 | Bank 5 |
| Bank 2 | Bank 6 |
| 2 GEMV units | 2 GEMV units |
| Bank 3 | Bank 7 |

| Accumulator | Accumulator |
|---|---|
| GBUS CTRL | GBUS CTRL |

| GBUS CTRL | GBUS CTRL |
|---|---|
| Accumulator | Accumulator |

DRAM die

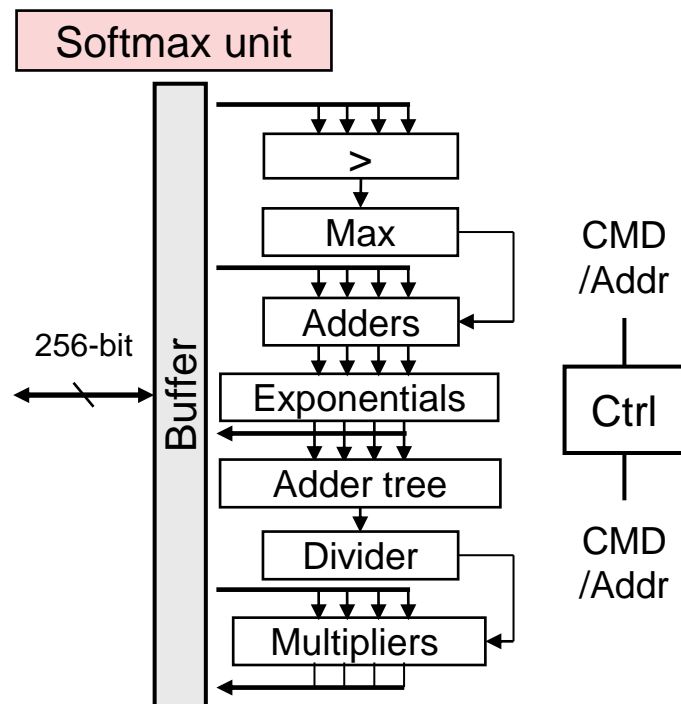| PHY | |
|---|---|
| TSVs | |
| Accumulators | Softmax units |

Buffer die

# AttAcc: Softmax Unit

- Placed on buffer die
  - Communication with multiple GEMV units is required
  - Complex processing units and requirement for large SRAM buffers for intermediate vectors
  - Placing softmax unit on DRAM die is overkill.
    - For GPT-3 175B, the FLOPs of softmax is **50 times smaller** than GEMVs in the attention layer
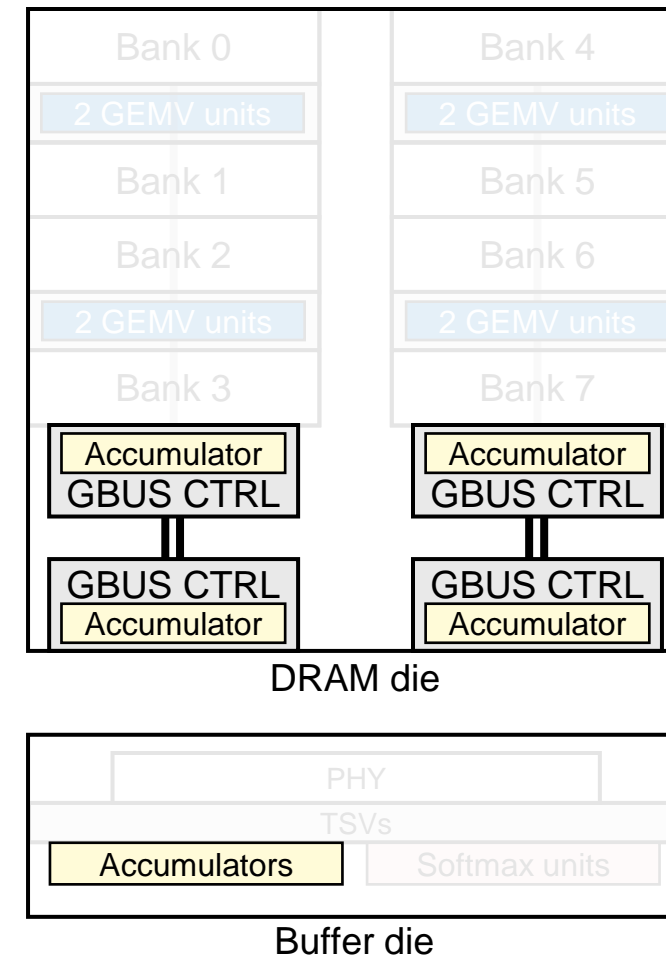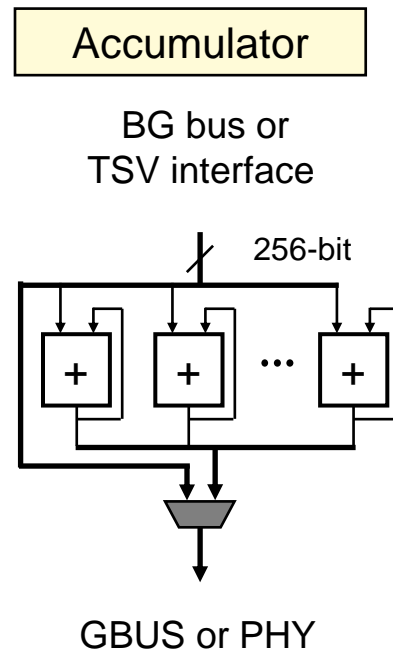
# AttAcc: Softmax Unit

- Placed on buffer die

- Processing units such as exponents, multipliers, and adders supporting FP32

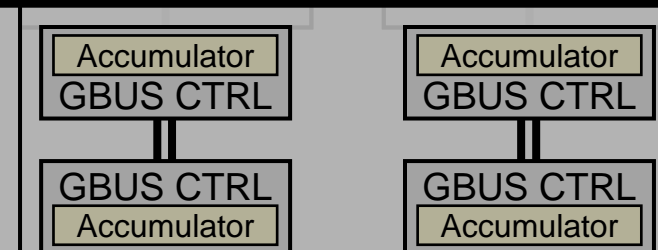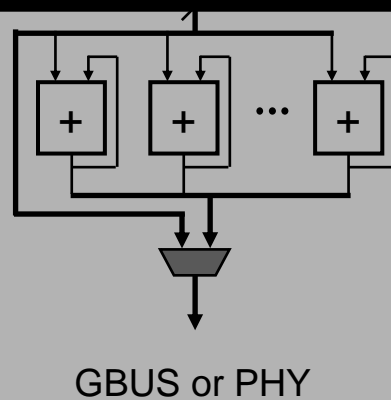- Buffer for intermediate vectors and control unit.

# AttAcc: Accumulator

- Placed hierarchically between the GEMV and the softmax units
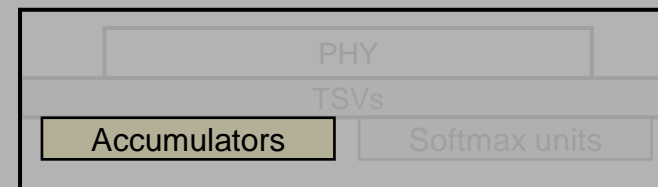- Supports the reduction of partial results from different GEMV units



Accumulator

BG bus or
TSV interface

256-bit

GBUS or PHY



| Bank 0 | Bank 4 |
|---|---|
| 2 GEMV units | 2 GEMV units |
| Bank 1 | Bank 5 |
| Bank 2 | Bank 6 |
| 2 GEMV units | 2 GEMV units |
| Bank 3 | Bank 7 |

Accumulator
GBUS CTRL

Accumulator
GBUS CTRL

GBUS CTRL
Accumulator

GBUS CTRL
Accumulator

DRAM die

PHY

TSVs

Accumulators          Softmax units

Buffer die

# AttAcc: Accumulator

- Placed hierarchically between the GEMV and the softmax units

- Supports the reduction of partial results from different GEMV units

Bank 0    Bank 4



**Please refer to the full paper for more detailed design exploration and data mapping.**

| + | + | ... | + |

GBUS or PHY

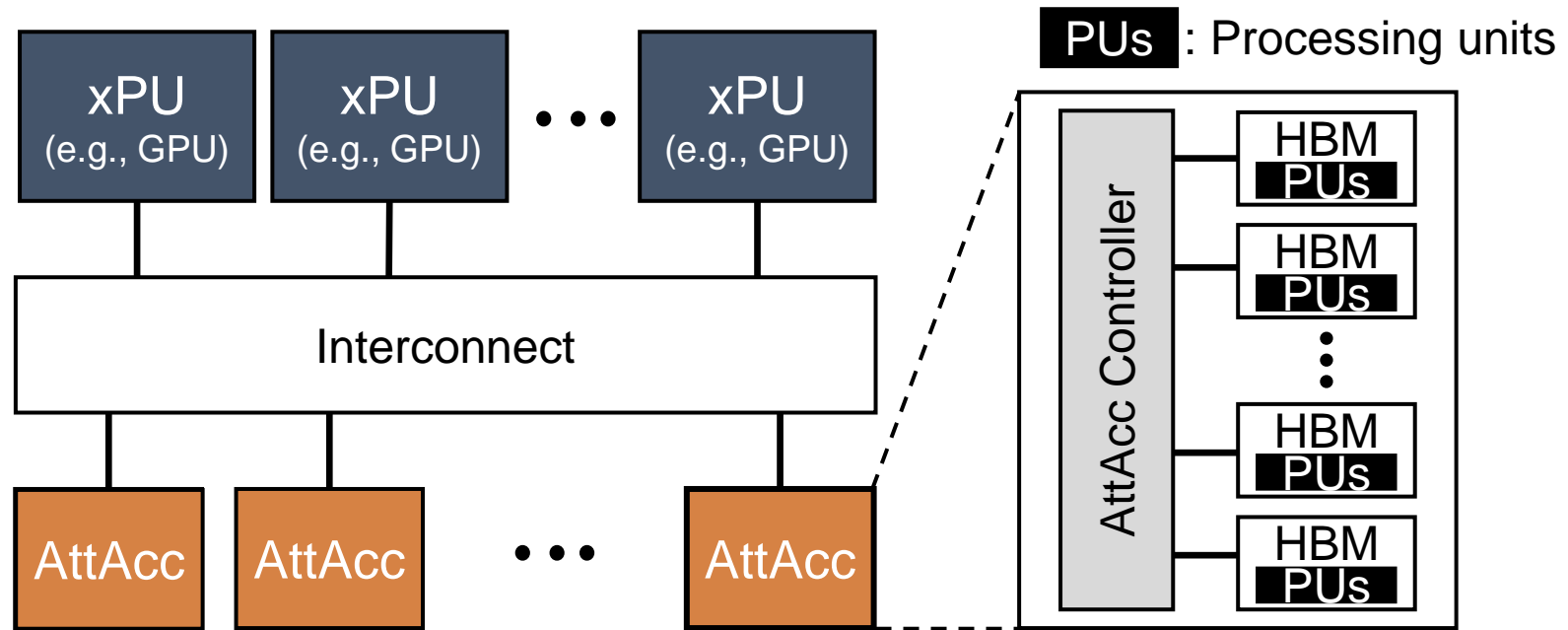| Accumulator | | Accumulator |
| GBUS CTRL | | GBUS CTRL |
| GBUS CTRL | | GBUS CTRL |
| Accumulator | | Accumulator |

DRAM die

PHY

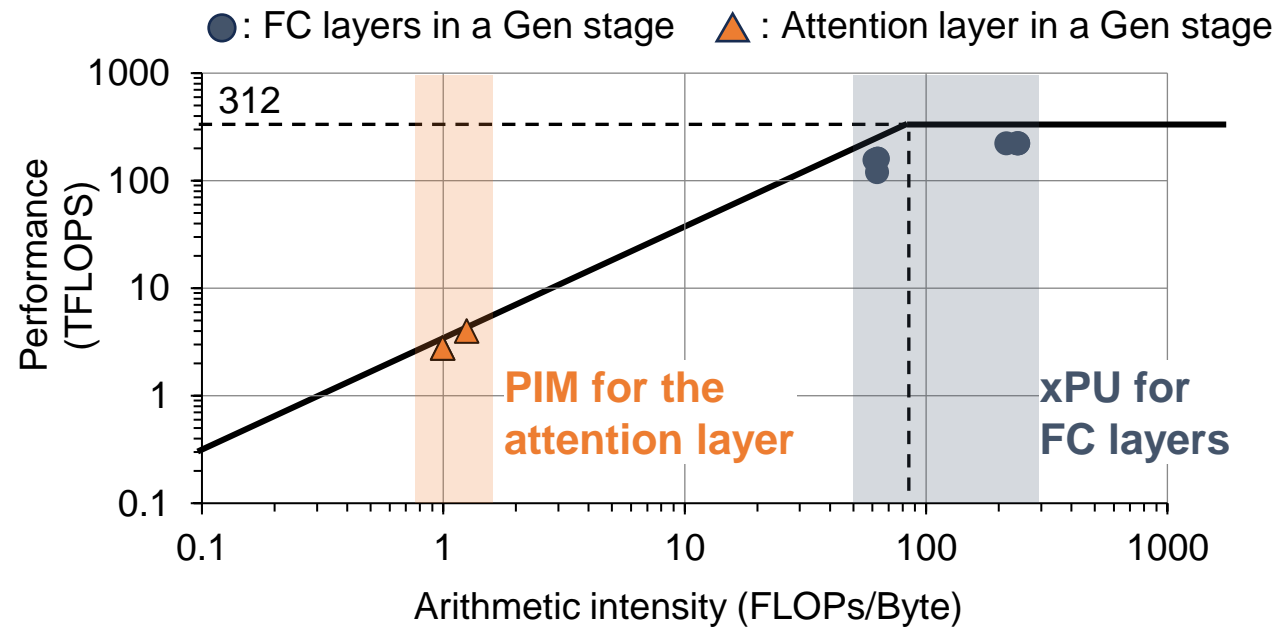TSVs

| Accumulators | Softmax units |

Buffer die

# Heterogeneous System with *AttAcc*



Heterogeneous system with xPUs and *AttAcc*s

# Heterogeneous System with *AttAcc*

- High performance
  - High computing power of **xPU for batched FC layers** with high FLOPs/Byte
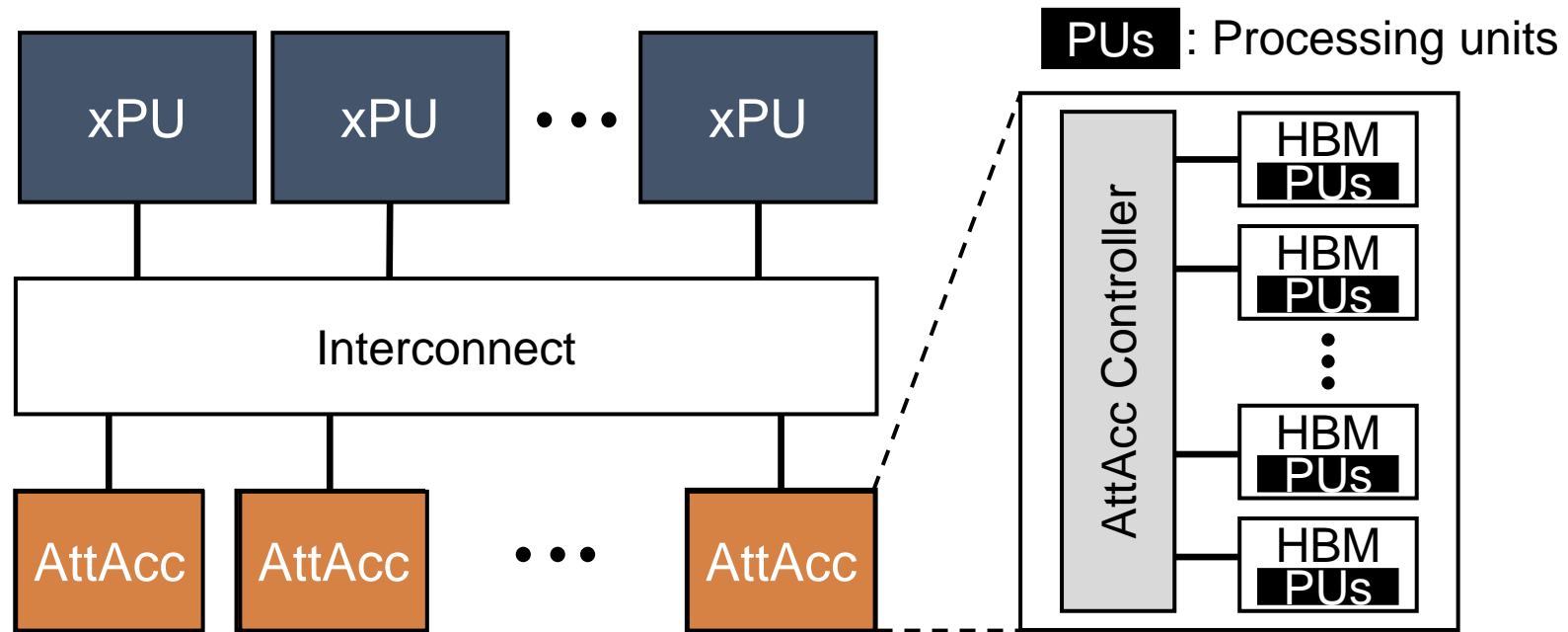  - Amplified memory bandwidth of **PIM for the attention layer** with low FLOPs/Byte



Roofline model of the DGX-A100 with HBM3

# Heterogeneous System with *AttAcc*

- High performance
  - High computing power of **xPU for batched FC layers** with high FLOPs/Byte
  - Amplified memory bandwidth of **PIM for the attention layer** with low FLOPs/Byte

- High energy efficiency
  - Leveraging **high reusability of weights** of batched FC layers through on-chip caches in xPU
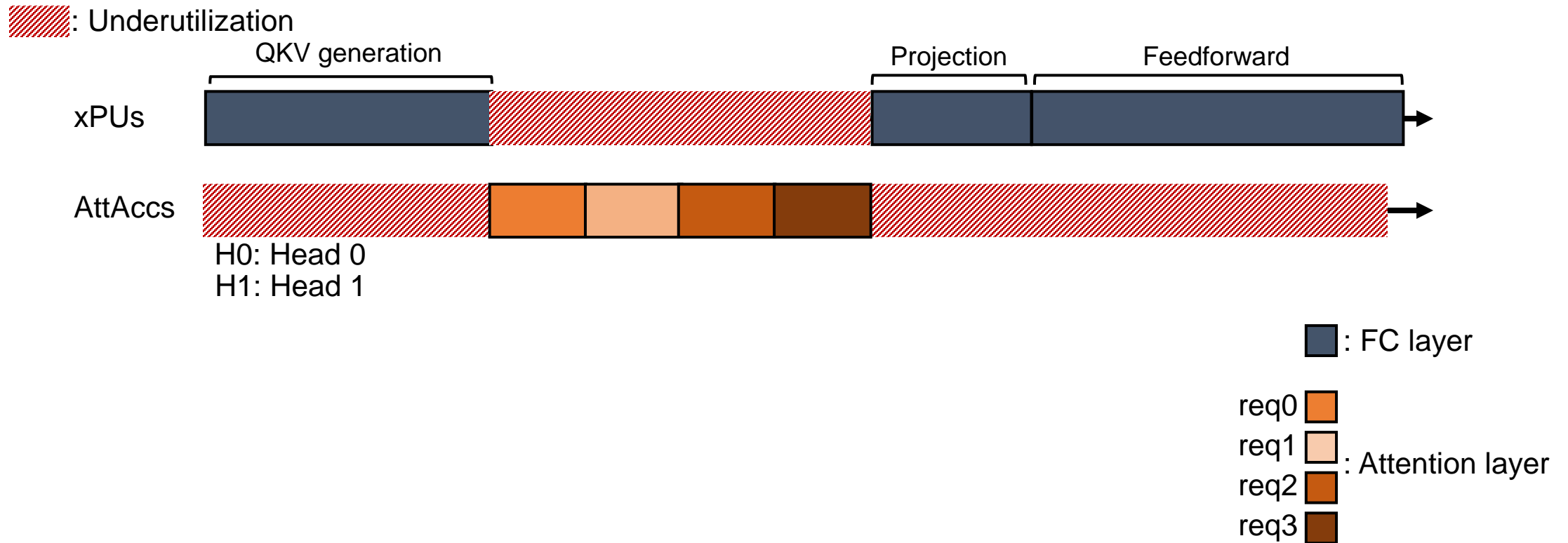  - Benefit from **short data transfer** in PIM for the attention layer

# Heterogeneous System with *AttAcc*

- Proposed system consists of multiple xPUs (e.g., GPU, TPU) and attention accelerators *AttAcc*s
  - Batched FC layers on multiple xPUs
  - Attention layers on *AttAcc*s

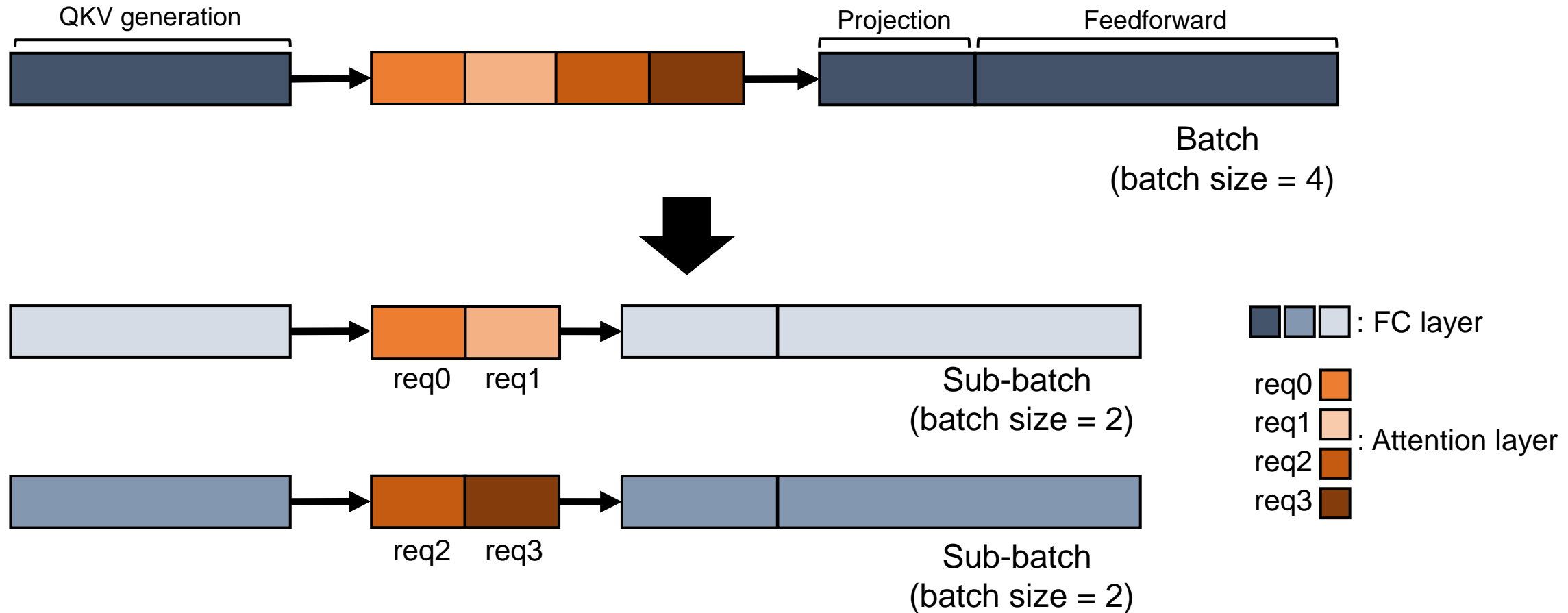- *AttAcc*s and xPUs can be connected via an interface such as NVLINK, PCIe, and CXL.
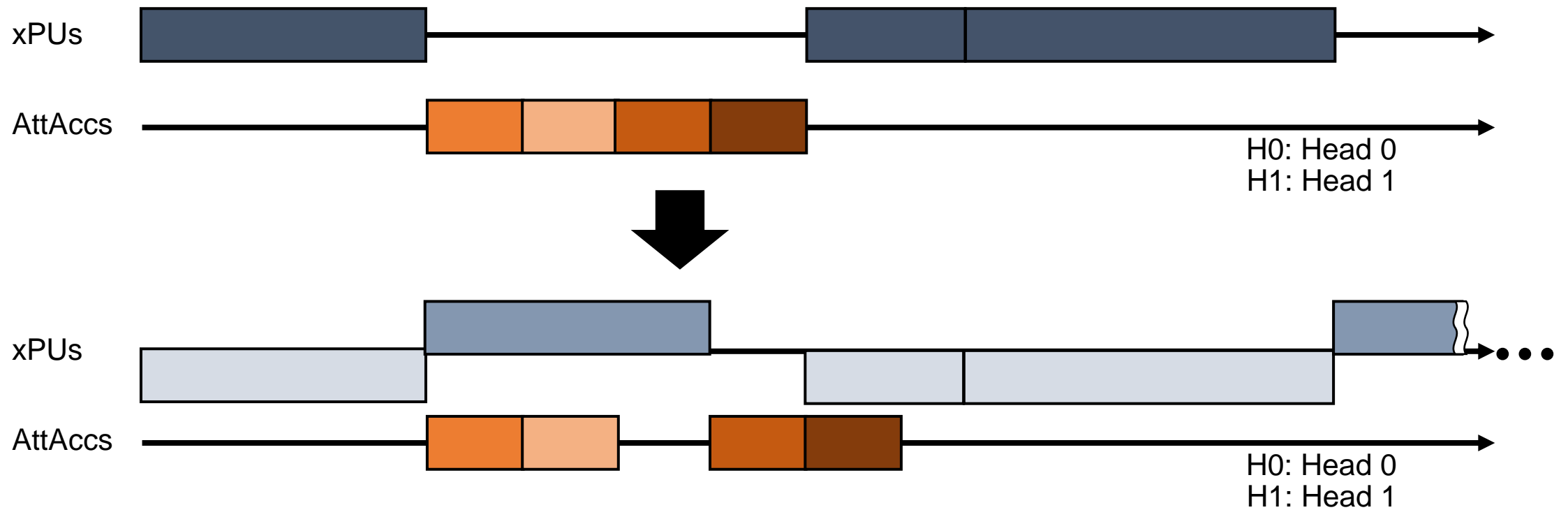


Heterogeneous system with xPUs and *AttAcc*s

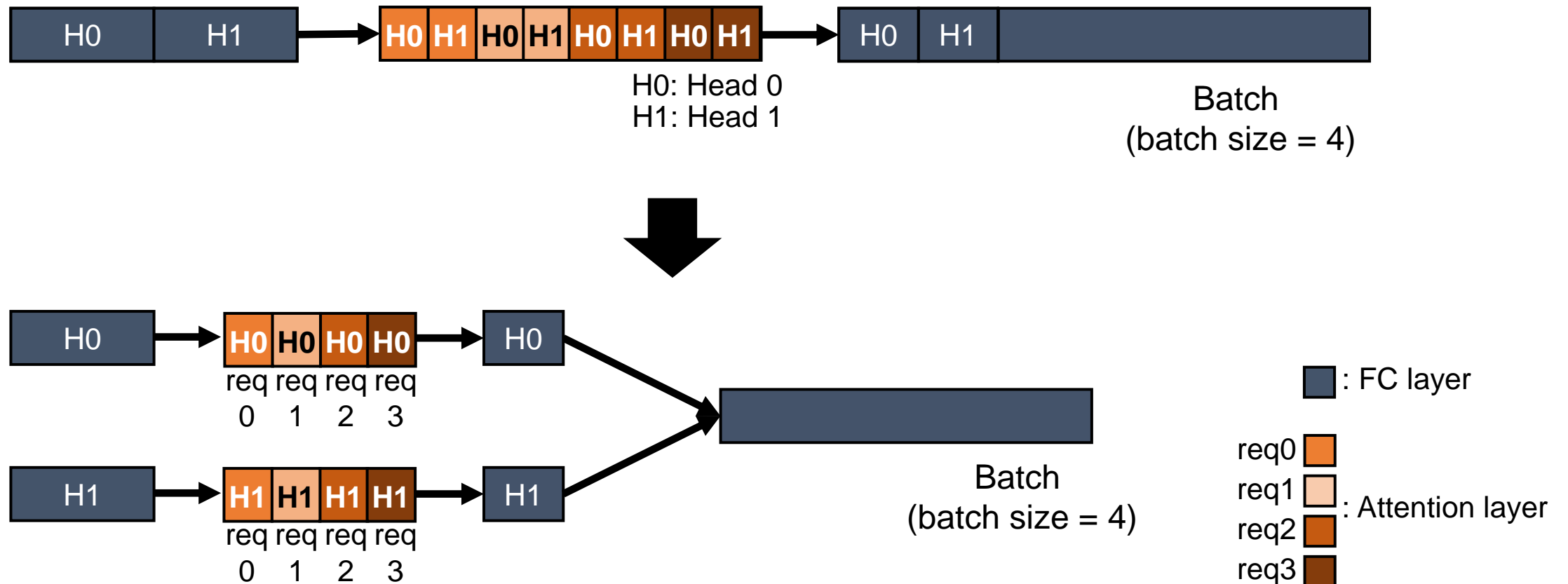# Execution Flow of the Heterogeneous System



: Underutilization

QKV generation

Projection    Feedforward

xPUs

AttAccs

H0: Head 0
H1: Head 1

: FC layer

req0
req1
req2          : Attention layer
req3

39

# Naïve Approach: Batch-level Pipelining



QKV generation

Projection    Feedforward

Batch
(batch size = 4)

req0   req1

Sub-batch
(batch size = 2)

req2   req3

Sub-batch
(batch size = 2)

: FC layer

req0
req1           : Attention layer
req2
req3

# Naïve Approach: Batch-level Pipelining



xPUs
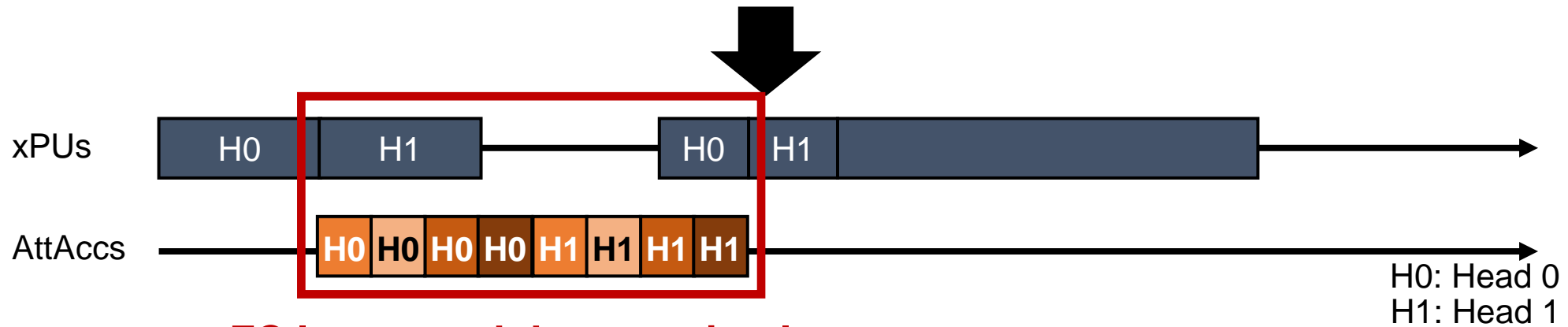
AttAccs

H0: Head 0
H1: Head 1

xPUs

AttAccs

H0: Head 0
H1: Head 1

# Head-level Pipelining

- FC layers that precede or follow the attention layer can be divided into heads.



H0: Head 0
H1: Head 1

Batch
(batch size = 4)

: FC layer

req0
req1    : Attention layer
req2
req3

# Head-level Pipelining



FC layers and the attention layer can be overlapped

# FeedForward Co-processing



Offload parts
of FC layers
to AttAcc

H0: Head 0
H1: Head 1

xPUs

AttAccs

# FeedForward Co-processing



xPUs

H0  H1                    H0  H1

AttAccs

H0 H0 H0 H0 H1 H1 H1 H1

H0: Head 0
H1: Head 1

xPUs

H0  H1          H0  H1

AttAccs

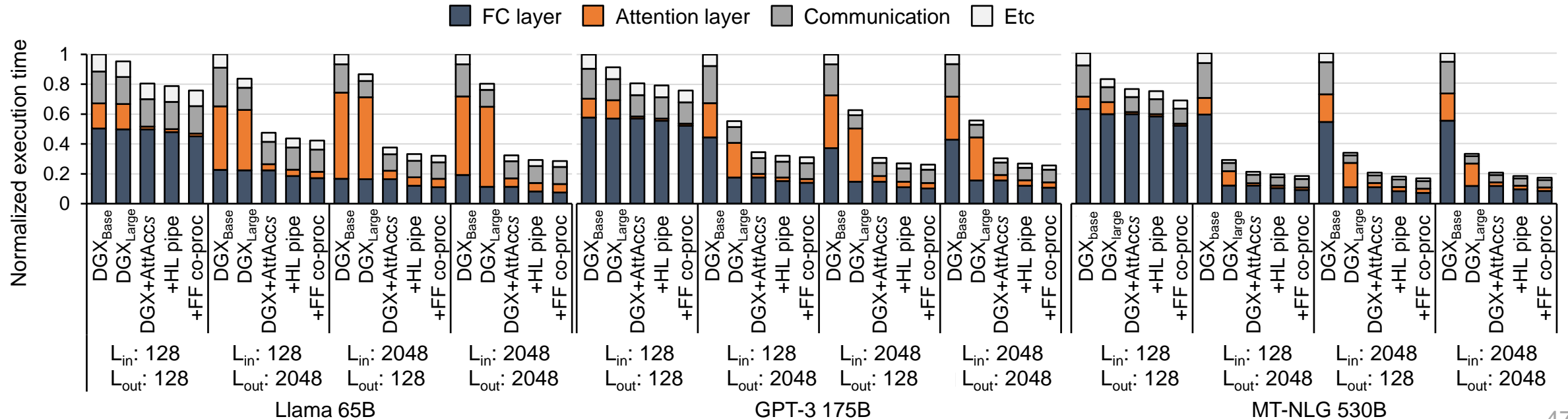H0 H0 H0 H0 H1 H1 H1 H1

H0: Head 0
H1: Head 1

# Experimental Setup

- Performance
  - Ramulator2 [1] and in-house simulator to evaluate **AttAcc** and DGX, respectively

- Energy and area
  - RTL synthesis for compute units and CACTI for buffer
  - The area overhead of **AttAcc**s is 10.84% of a HBM.
    - = Scaling the area to DRAM process for units in DRAM die

- Target model
  - Various size of TbGMs: Llama 65B, GPT-3 175B, and MT-NLG 530B

- Comparison
  - $DGX_{Base}$: DGX-A100 having 40 HBM stacks
  - $DGX_{Large}$: DGX-A100 having 80 HBM stacks
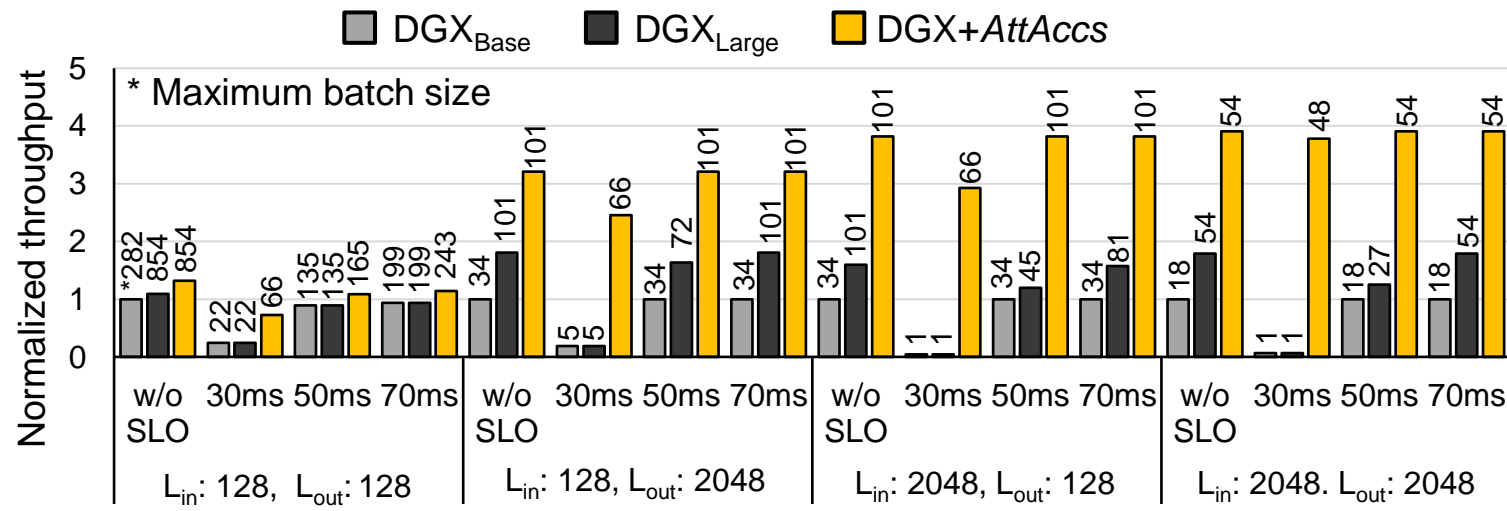  - DGX+**AttAcc**: $DGX_{Base}$ + 8 **AttAcc**s with 5 HBM stacks each

[1] H Luo et al., "Ramulator 2.0: A Modern, Modular, and Extensible DRAM Simulator," arXiv, 2023.

# Evaluation (Performance)

- *DGX+**AttAcc**s* outperforms $DGX_{Base}$ and even $DGX_{Large}$ up to by 5.93x and 2.81x, repectively
  - 4.84x and 2.48x from ***AttAcc***
  - 1.15x from head-level pipelining
  - 1.10x from feedforward co-processing
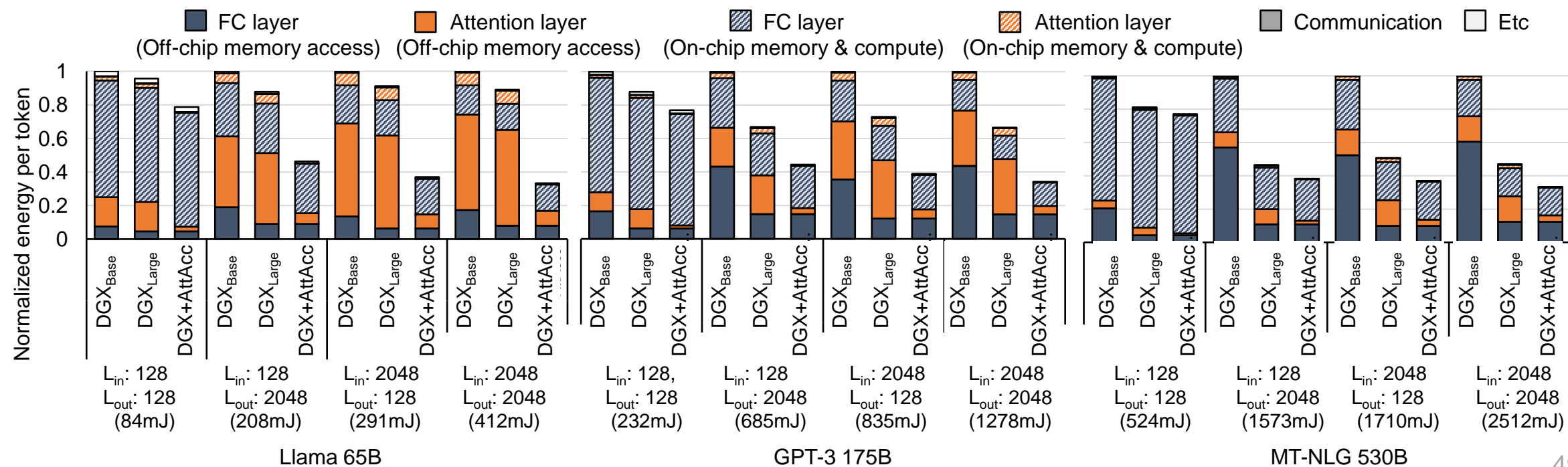
# Evaluation (Performance)

- *DGX+**AttAcc**s* achieves further throughput improvement under SLO constraint
  - Performance improvement from relieving the batch size constraints caused by SLO



Normalized throughput of GPT-3 175B inference for various SLOs

# Evaluation (Energy Efficiency)

- Energy consumption of *DGX+**AttAcc**s* compared to $DGX_{Base}$ ($DGX_{Large}$) is reduced by up to
  - 66.7% (62.6%) for Llama 65B
  - 65.9% (48.8%) for GPT-3 175B
  - 66.8% (29.1%) for MT-NLG 530B

# Conclusion

- We discovered that **the attention layer** poses a constraint on the batch size in conventional systems (e.g., DGX) due to the **long latency** and **memory capacity requirements**.

- We proposed a **heterogeneous system** (*DGX + AttAccs*) with the conventional system for the batched FC layer and *AttAccs* for the attention layer, leveraging PIM architecture.

- We explored GEMV unit placement and data mapping in the PIM architecture and proposed efficient pipelining and co-processing optimizations to improve system utilization.

- *DGX+AttAccs* achieved higher throughput (up to 2.81×) and energy efficiency (up to 2.67×) compared to the monolithic GPU system.

**Thank you!**

**Question?**