# Samsung PIM/PNM for Transformer based AI

## : Energy Efficiency on PIM/PNM Cluster

Jin Hyun Kim, Yuhwan Ro, Jinin So, Sukhan Lee, Shin-haeng Kang, YeonGon Cho, Hyeonsu Kim, Byeongho Kim, Kyungsoo Kim ,Sangsoo Park, Jin-Seong Kim, Sanghoon Cha, Won-Jo Lee, Jin Jung, Jong-Geon Lee, Jieun Lee, JoonHo Song, Seungwon Lee,  Jeonghyeon Cho, Jaehoon Yu, and Kyomin Sohn
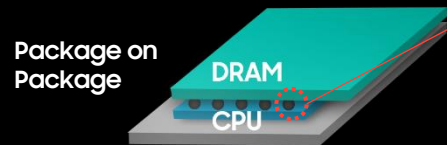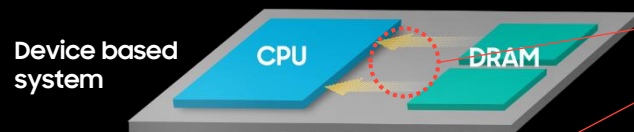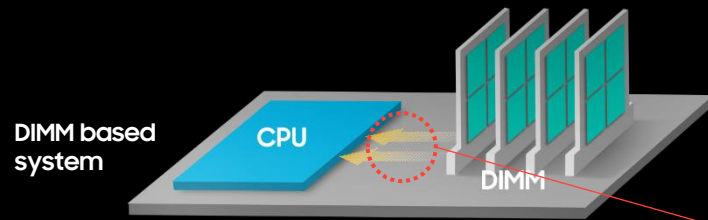
Samsung Semiconductor

# **Index**

# PIM/PNM on Memory Hierarchy and Energy Reduction

- Data movement consumes a lot of energy even for simple computation
- PIM/PNM technology can reduce energy consumption within a typical memory hierarchy
- PIM/PNM device for each layer must meet specific requirements: bandwidth(BW), power, capacity, etc.

## Energy Cost of Data Transfer

DRAM

Global Buffer

PE — PE
PE — ALU (RF / RF)

Fetch data to run a MAC operation here

### Normalized Energy Cost*

| | ALU | 1 x (Reference) |
| 0.5-1.0 kB | RF → ALU | 1 x |
| NoC: 200 – 1000 PEs | PE → ALU | 2 x |
| 100 – 500 kB | Buffer → ALU | 6 x |
| | DRAM → ALU | 200 x |

Source: * Y.-H. Eyeriss, 2016 ISCA

## Computing Coupled Memory Hierarchy

Higher Bandwidth

Capacity

[Memory requirement]

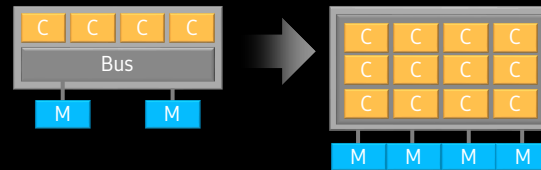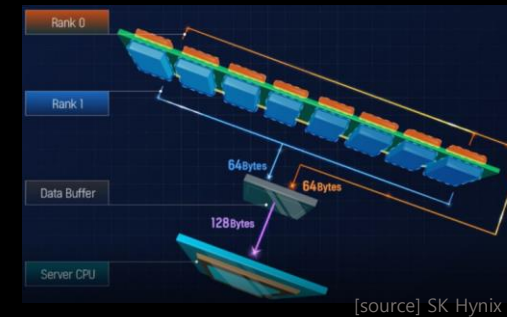| Reg | |
| (~256MB) LLC, **LLC-DRAM** | |
| (~24GB) HBM, **HBM-PIM** | Extreme BW and Power |
| (~256GB) DDRx, LPDDR, **LPDDR-PIM** | High capacity, BW for AI, Low power |
| (~512GB) CXL-DRAM, **CXL-PNM** | Extreme capacity, Moderate cost |
| (~8PB) Storage (SSD, **PBSSD, MS-SSD**) | |

# Traditional Approach to Overcome Memory Bottleneck

- While there are various methods to increase BW, it is difficult to achieve a dramatic increase
  - Limited by # of PCB wires, # of CPU ball, and thermal constraints

- Increasing # of the balls and PCB wires is physically and thermally bounded and is a expensive solution
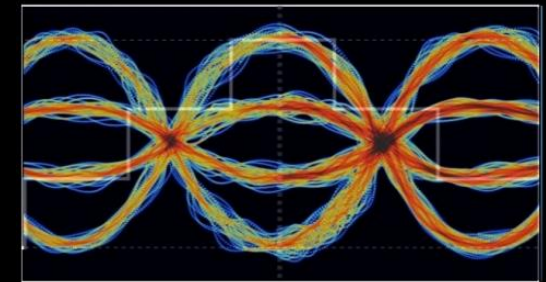  - MCR-DIMM, PAM3/4 signaling IO, 2K-IO or 3D stacking

**DIMM based system**

CPU
DIMM

**Device based system**

CPU
DRAM

**Package on Package**

DRAM
CPU

Memory BW is continuously requested due to AI application and multi-core host

C C C C
Bus
M M

C C C C
C C C C
C C C C
M M M M

**MCR-DIMM(Multiplexer Combined Ranks DIMM)**

Rank 0
Rank 1
Data Buffer
64Bytes
64Bytes
128Bytes
Server CPU

[source] SK Hynix

**PAM 3/4(Pulse Amplitude Modulation)**

[source] Micron

**2K-IO 2D or 3D stacking**

2D
3D

[source] Samsung
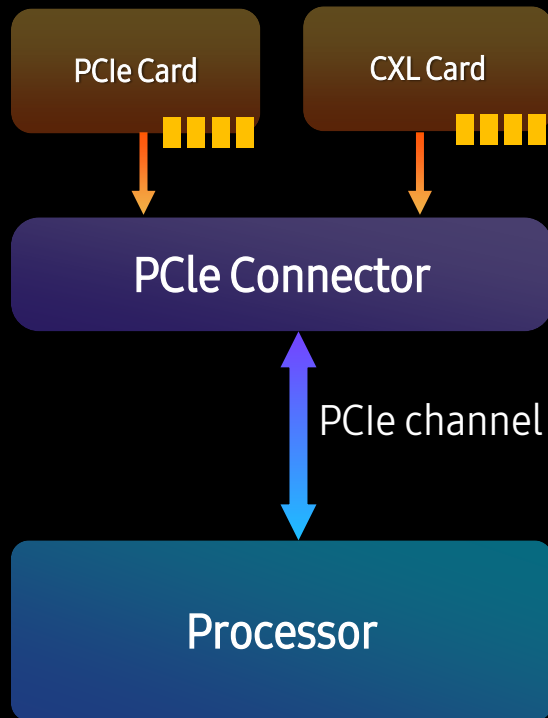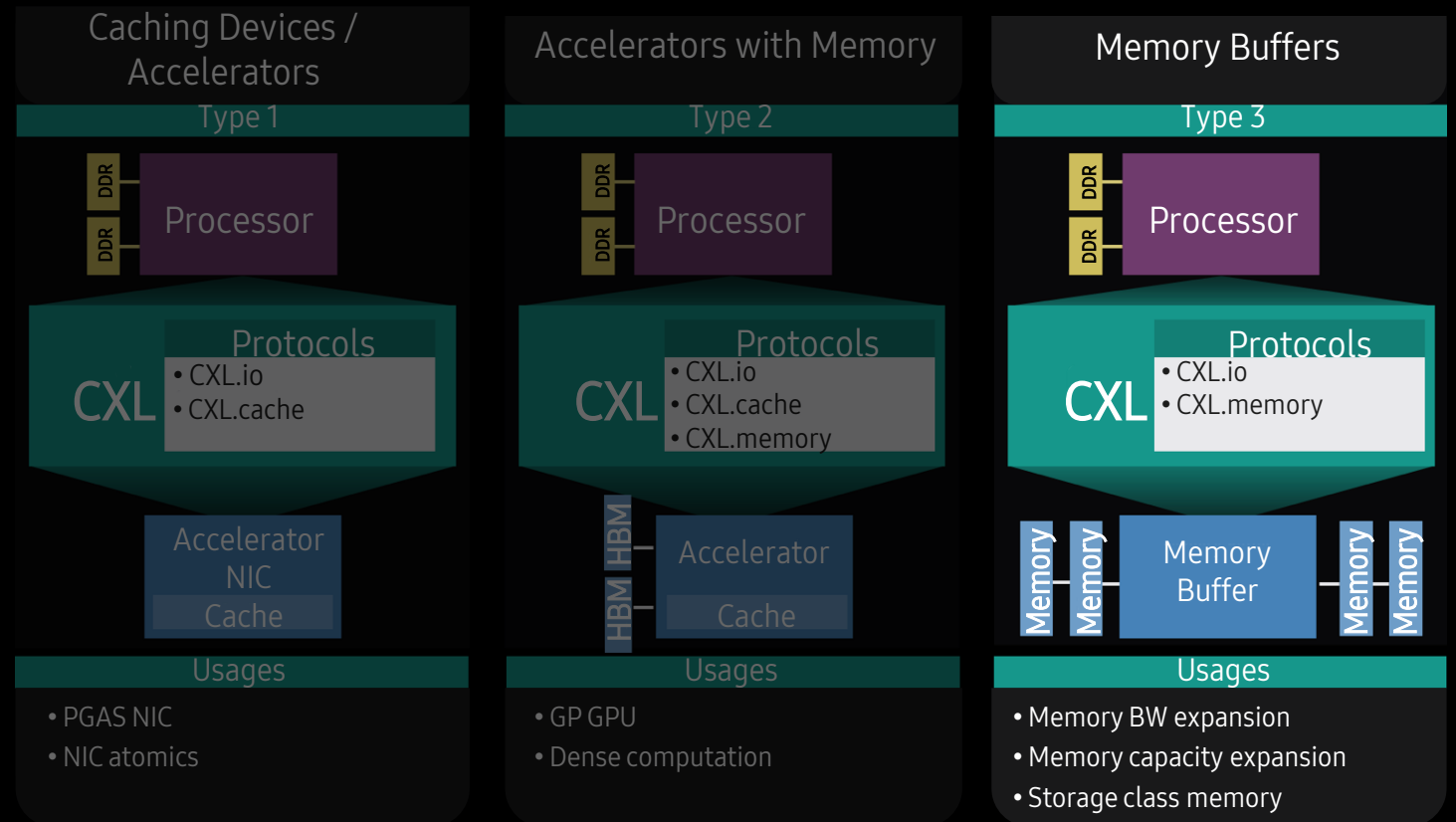
SAMSUNG

# CXL solution, Trend to Pay Attention to

- CXL is strong candidate for memory hierarchy to address performance and density
- Successful power-on of memory expander, SSD/pooling solutions are next big-thing

## CXL is a high performance, low latency protocol that leverages PCIe physical layer



## CXL is an open industry standard with broad industry support

### Caching Devices / Accelerators
Type 1

**CXL**
Protocols
- CXL.io
- CXL.cache

Processor (DDR)

Accelerator NIC
Cache

Usages
- PGAS NIC
- NIC atomics

### Accelerators with Memory
Type 2

**CXL**
Protocols
- CXL.io
- CXL.cache
- CXL.memory

Processor (DDR)

HBM Accelerator
Cache

Usages
- GP GPU
- Dense computation

### Memory Buffers
Type 3

**CXL**
Protocols
- CXL.io
- CXL.memory

Processor (DDR)

Memory Buffer

Usages
- Memory BW expansion
- Memory capacity expansion
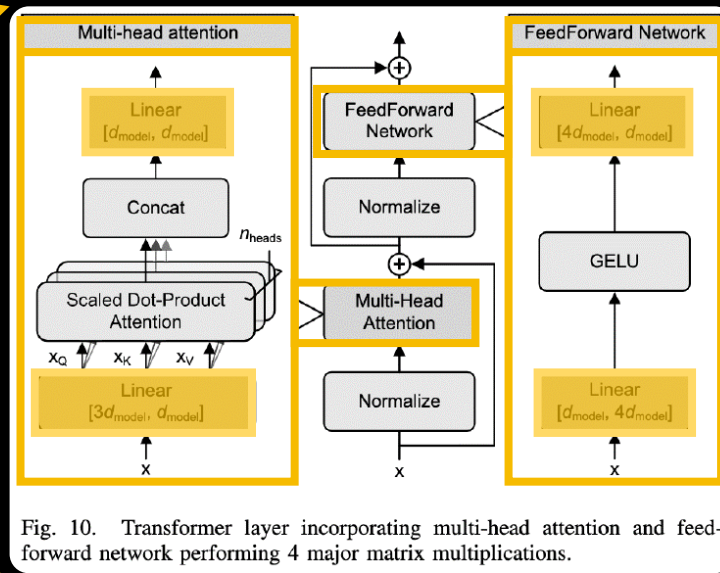- Storage class memory

SAMSUNG

# Bottleneck in GPT: Linear Layers in Generation Stage

- Main target: transformer decoders used in ChatGPT, GPT-3
  - Linear layers in multi-head attention(MHA) and feed-forward networks(FFN)

- Focus on memory-bottleneck in Generation stage
  - Generation Stage shows poor performance with GPU due to its memory-bound & sequential characteristic

**<Transformer>**

[source] Google Transformer

Fig. 10. Transformer layer incorporating multi-head attention and feed-forward network performing 4 major matrix multiplications.

**<Four major matrix multiplication block>**

Computation

Memory BW          Time
                   Inference Finish
Memory Load
& Store

**<Summarization>**

Computation

Memory BW          Time
                   Inference Finish
Memory Load
& Store

**< Generation >**

[source] Naver Clova

# GPT Profiling Result

- GPT workload consists of Summarization(computing-bound) and Generation(memory–bound)
- GEMV portion can be 60–80% of total generation latency, which are the target of PIM/PNM

Legend: GEMM$_{SUM}$ · GEMV$_{GEN}$ · VECTOR · GELU · SOFTMAX · RESIDUAL · ETC

**Number of Operations**

| 13.46 | 86.53 |

**Latency**

| 2.12 | 82.27 | 3.8 | 1.8 | 2.2 | 1.4 | 6.5 |

↔ **Matrix-Matrix/Vector Multiplication** | ↔ **Non-Linear Function**
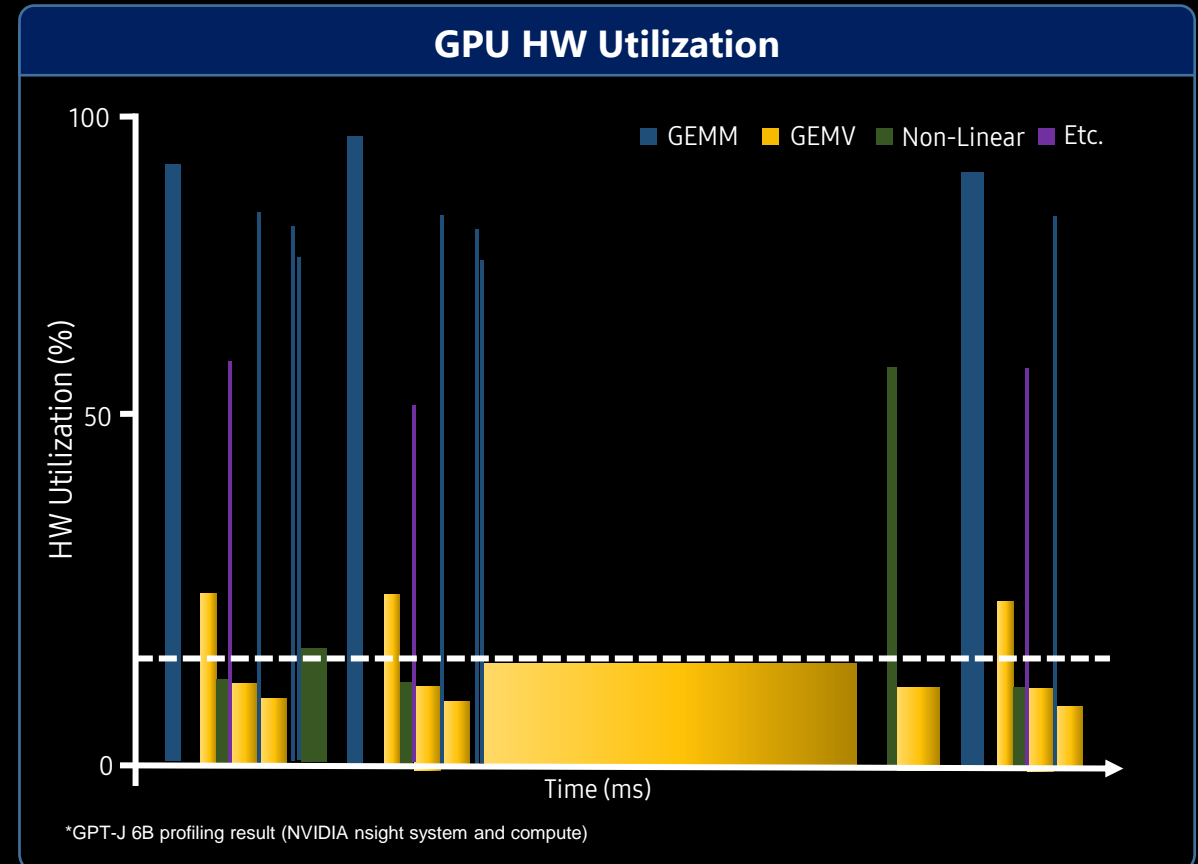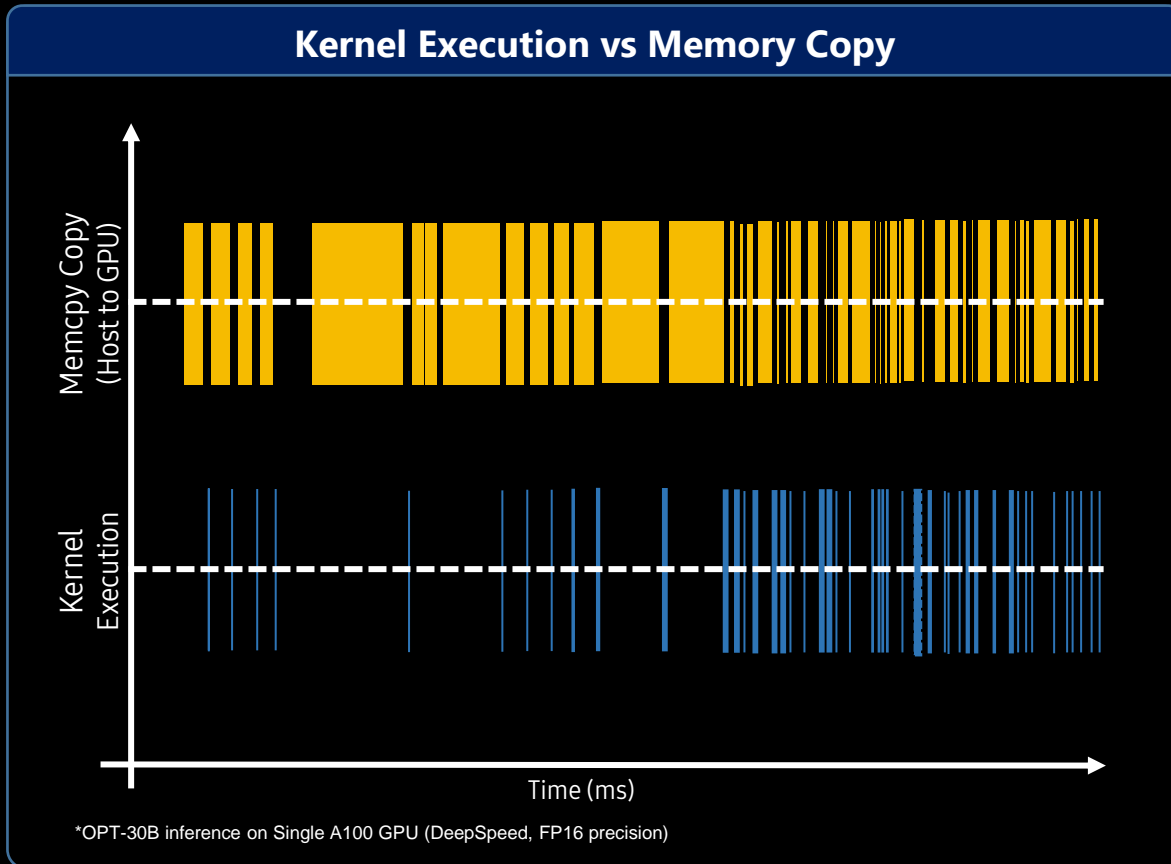
0% — 80% — 85% — 90% — 95% — 100%

*Profiling result is measured in A100 System (DeepSpeed + GPT-J 6B, FP16, Input/Output token:7/46)
GPT-j: Google JAX framework

https://community.openai.com/t/how-does-chatgpt-have-such-massive-token-limit/25738/6

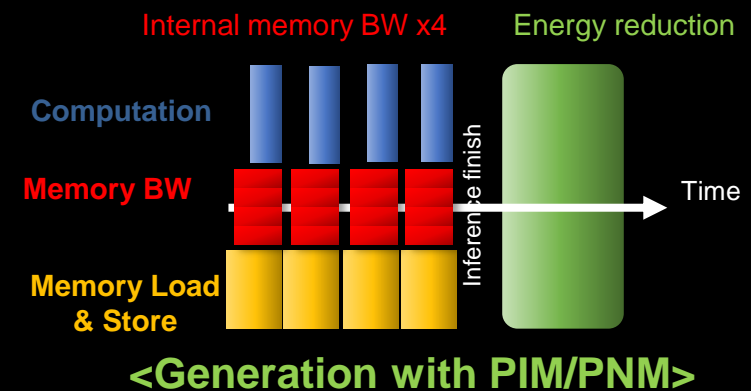| Stage | Computation | Latency |
|-------|-------------|---------|
| SUM | 78.95 GFLOP | 7.62 ms |
| GEN | 11.28 GFLOP | 6.58 ms |

# Utilization and Execution Time Breakdown

- Most of the execution time is spent for the memory copy from the host CPU memory to the GPU memory
- Utilization for performing GEMV operations (Generation stage) is seriously low, compared to GEMM
- As # of output tokens increases, GEMV operations dominate the inference time

## Kernel Execution vs Memory Copy

*OPT-30B inference on Single A100 GPU (DeepSpeed, FP16 precision)

## GPU HW Utilization

*GPT-J 6B profiling result (NVIDIA nsight system and compute)

# Acceleration by PIM/PNM on Generation stage

- Generation stage on GPT requires high capacity and bandwidth memory

- MHA and FFN can be fully offloaded to PIM/PNM, exploiting full bandwidth provided by PIM/PNM

- As a result, PIM/PNM can significantly reduce the time and energy spent on inference

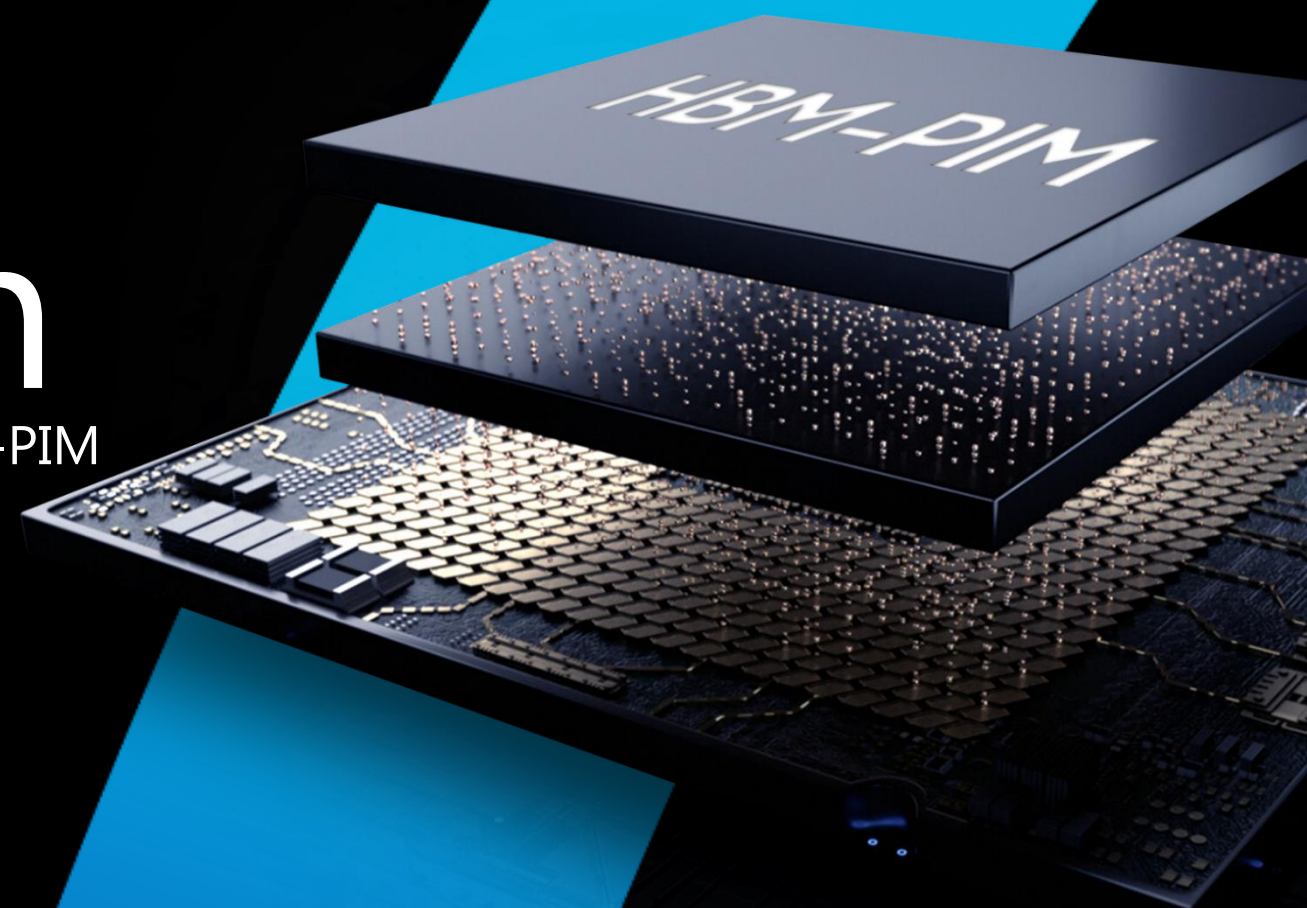**Summarization**          **Generation Stage**

GEMV Computation dominated

"Hello, my name"

**LM** → **LM** → **LM** → **LM** → ··· → **LM**

"is"    "James"   "Smith"   "and"    "."

0%        25%        50%        75%        100%

| 2.1% | GEMV accounts for 82.3% of total computation time (GPU) | 7.8% | 7.9% |

*GPT-J 6B latency breakdown on A100 GPU ■ GEMM ■ GEMV ■ Non-Linear ■ Etc.

Computation
Memory BW
Memory Load & Store
Inference finish → Time

**< Generation >**

Internal memory BW x4          Energy reduction

Computation
Memory BW
Memory Load & Store
Inference finish → Time
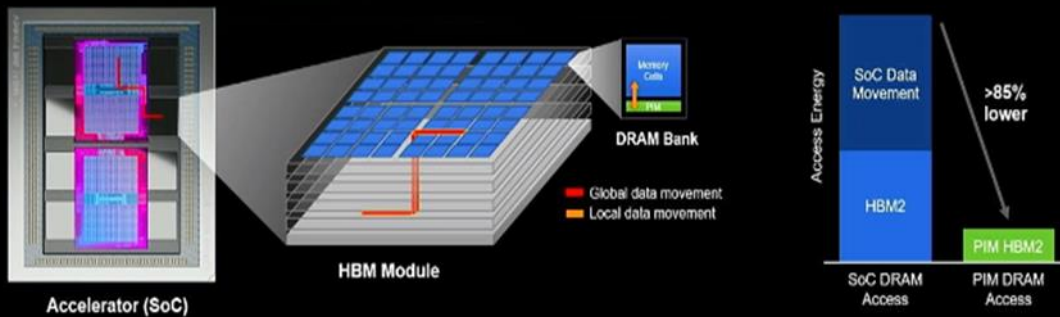
**<Generation with PIM/PNM>**

# PIM Solution
Redesigned to Advance AI : HBM-PIM / LPDDR-PIM

# Energy Advantage of PIM on Generative AI

- Since OpenAI focuses on developing new AI technologies and pushing the boundaries of what can be done with AI, it is likely that they will explore the use of PIM technology in the future.

- In ISSCC 2023, AMD mentioned
  - Key algorithmic kernels can be executed directly in memory, saving precious communication energy
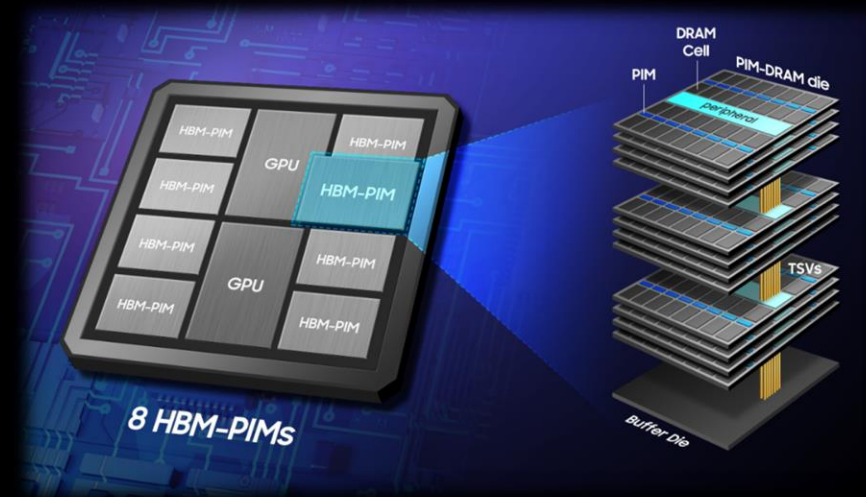  - PIM can reduce energy by 85% compared with conventional HBMs

## Processing-in-Memory



Key algorithmic kernels can be executed directly in memory, saving precious communication energy
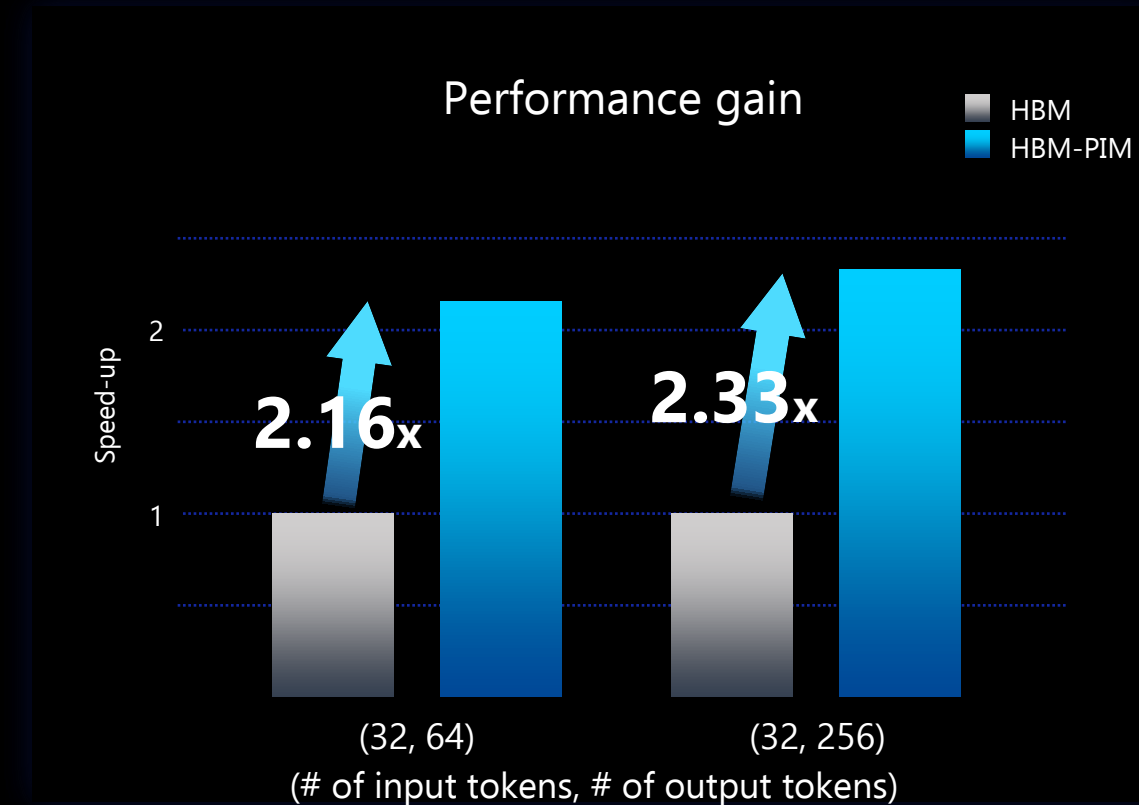
[source] AMD ISSCC

## Future HBM-PIM concept



[source] Samsung, MemCon

# Generative AI on HBM-PIM

- Experimental setup: GPT-J (6B, 32 input tokens), single AMD MI100-PIM GPU
- About 2x greater system energy efficiency compared to the GPU with a normal HBM
- GPT can be accelerated by more than 2x over baseline



Energy efficiency

■ HBM
■ HBM-PIM

Relative energy efficiency

**1.79x**

**2.24x**

(32, 64)          (32, 256)

(# of input tokens, # of output tokens)

Performance gain

■ HBM
■ HBM-PIM

Speed-up

**2.16x**

**2.33x**

(32, 64)          (32, 256)

(# of input tokens, # of output tokens)

# Architecture of HBM-PIM Cluster

- Installed 96 AMD MI100 GPUs fabricated with HBM-PIM
- Accelerate large-scale workloads with high energy efficiency and low latency

HBM-PIM cluster

AMD MI100-PIM GPU

Server node

· Capacity : 24GB (4 cubes)
· PIM performance : 4.9 TFLOPS
· GPU performance : 184.6 TFLOPS (FP16)
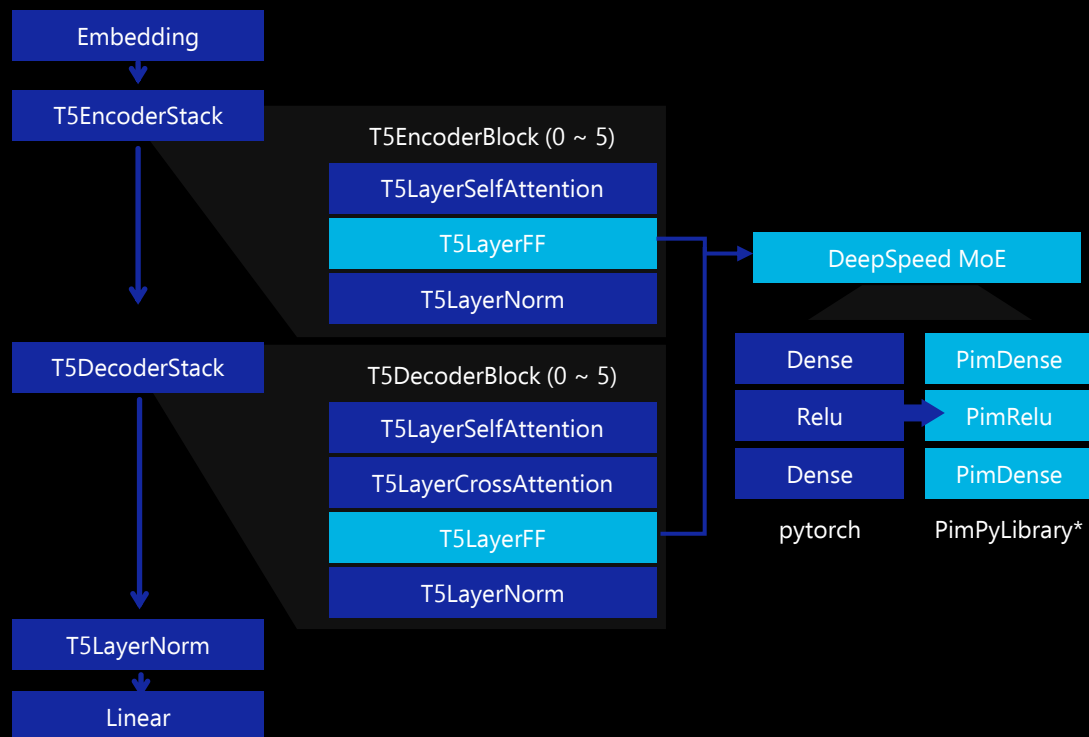
· 8 MI100-PIM GPUs per node

· Total 96 MI100-PIM GPUs in a cluster
. 12 nodes interconnected
  through 200G InfiniBand network
  (bi-section bandwidth : 1.2TB/s)
. Total memory capacity : 2.25TB
. Total PIM performance : 471.9 TFLOPS
. Total GPU performance : 17.7 PFLOPS (FP16)

# Workload : T5(Transformer)-Based MoE Model
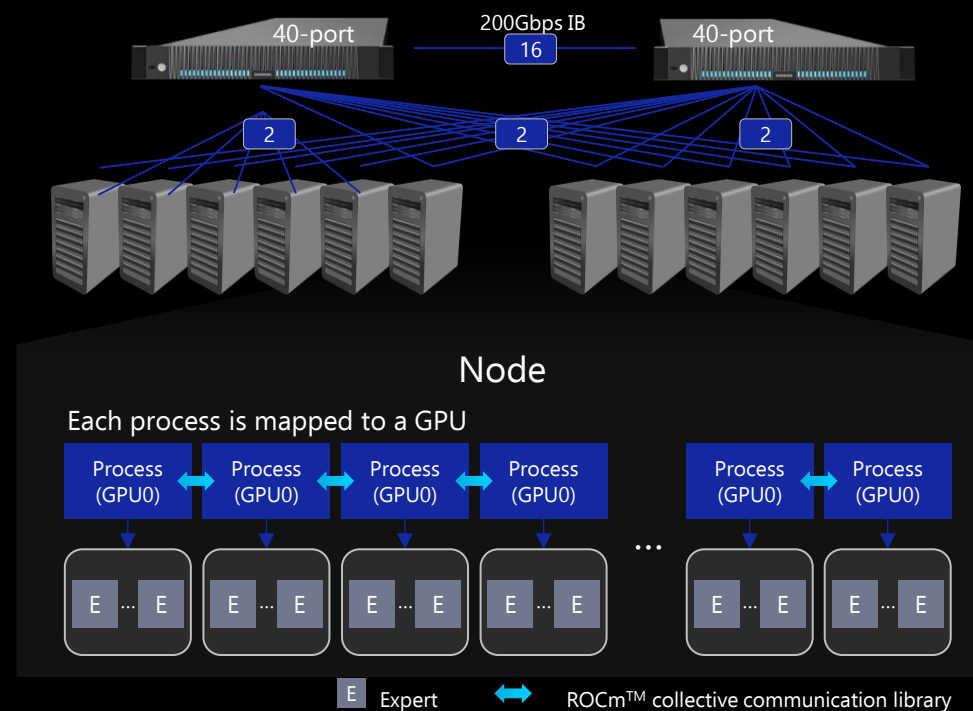(Mixture of Experts)

- DeepSpeed MoE replaces "T5LayerFF" layers to accelerate T5-large model with PIM
- The MoE layers are updated to use our PimPyLibrary* APIs



T5-MoE model architecture

MoE deployment on HBM-PIM cluster

* PimPyLibrary is a python library for providing PIM-enabled AI operators. PIM SDK provides not only PimPyLibrary but also full SW stack for utilizing PIM

# Energy Efficiency and Performance on MoE Model

- More than 3x greater system energy efficiency compared to normal GPU clusters
- Increases performance by more than 2x over baseline



Energy efficiency

HBM
HBM-PIM

Relative energy efficiency

3.61x

2.88x

32 MI100 GPUs    64 MI100 GPUs

Performance gain

HBM
HBM-PIM

Speed-up

2.71x

2.06x

32 MI100 GPUs    64 MI100 GPUs

* Acknowledgement: Jaeyoung Heo and professor Sungjoo Yoo (Seoul National University) provided the idea of PIM acceleration for this workload

# PIM S/W Stack for AI

- Support existing AI frameworks (e.g., PyTorch and TF) for users to utilize PIM functions
- PIM Runtime Library: Apply PIM and provide operator-level optimizations during PIM operation
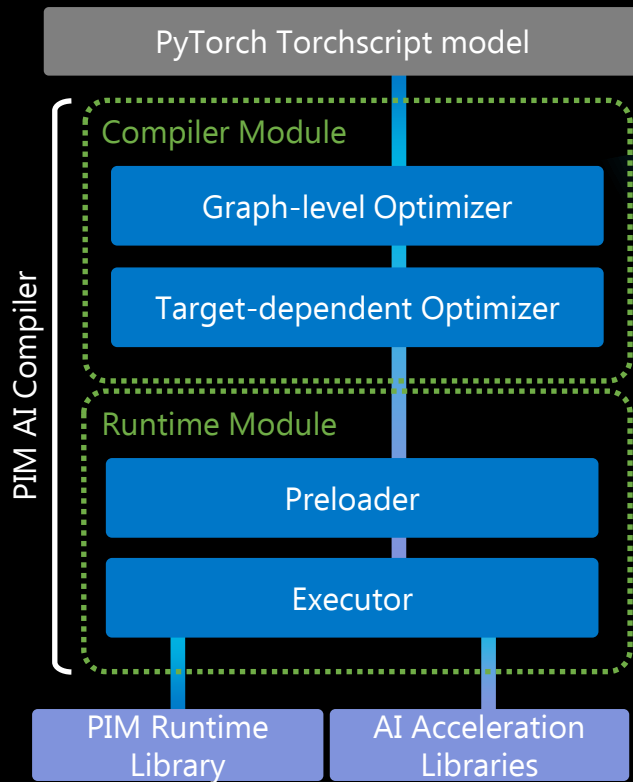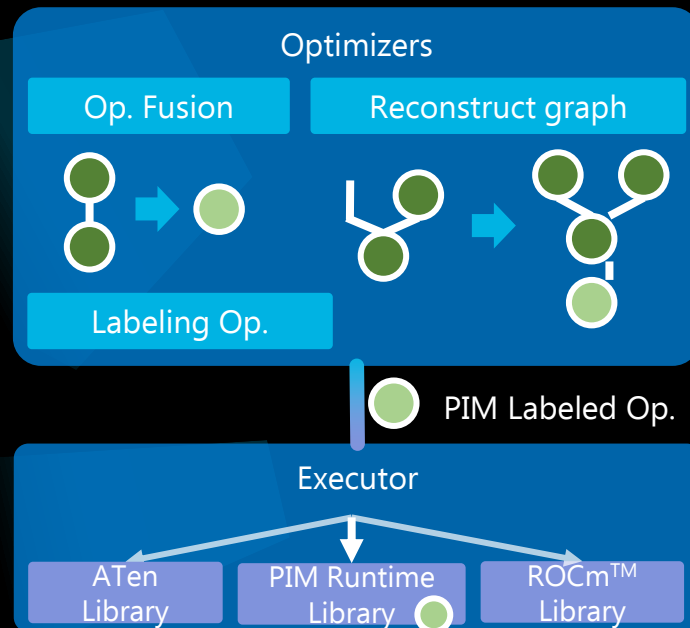- PIM AI Compiler: Provide graph-level optimizations during end-to-end execution

## Execution Pipeline

PyTorch Torchscript model

PIM AI Compiler

Compiler Module
- Graph-level Optimizer
- Target-dependent Optimizer

Runtime Module
- Preloader
- Executor

PIM Runtime Library | AI Acceleration Libraries

## Optimization Examples

- Operator Fusion is to mitigate kernel load overhead
- Reconstruct graph is to make possible to run on PIM

Optimizers
- Op. Fusion
- Reconstruct graph
- Labeling Op.

PIM Labeled Op.

Executor
- ATen Library
- PIM Runtime Library
- ROCm™ Library

- Runtime choose appropriate library per given tensor

## Effect of PIM with S/W optimizations

PyTorch ATen
PIM Runtime Library
PIM AI Compiler

PIM Runtime Library effect
PimAiCompiler effect

RNNT: 2.1x, 1.2x
GNMT: 1.8x, 1.2x

Relative Latency (1, 0.8, 0.6, 0.4, 0.2, 0)

# Promising Standard Programming Models

- **PIM-SYCL** accelerates upcoming HPC/AI applications on heterogeneous platform
  - SYCL supports CPU/GPU/NPU/FPGA in modern C++ template
- **PIM-OpenACC** is under development for legacy scientific applications
  - OpenACC enables incremental parallelization from C/Fortran serial code

### SYCL GEMV source code

```
// Buffer Allocation
buffer<sycl::half, 2> M{matrix, range<2>{N, N}};
buffer<sycl::half> X{in, range<1>{N}};
buffer<sycl::half> Y{out, range<1>{N}};

// Parallel Execution
q.submit([&](ext::samsung::pim_handler &ph) {
  ext::samsung::pim_accessor accM{M, ph, sycl::read_only};
  ext::samsung::pim_accessor accX{X, ph, sycl::read_only};
  ext::samsung::pim_accessor accY{Y, ph, sycl::write_only,
    sycl::property::no_init{}};

  ph.gemv(accY, accX, accM);
});
```

### SYCL GEMV performance

w/o PIM
w/ PIM

Execution time (ms)

Matrix Size (number of elements, 16-bit)

4096x1024  4096x2048  8192x1024  8192x2048  16384x102

# OneMCC S/W Standardization (To be)
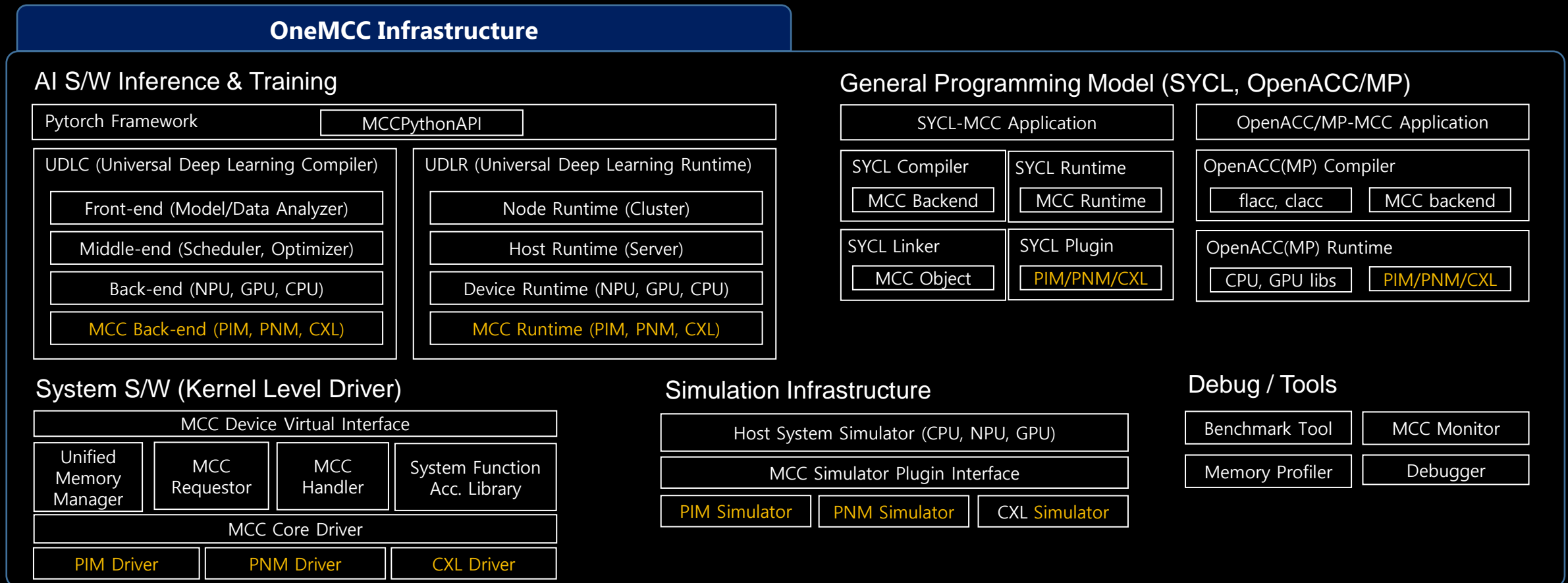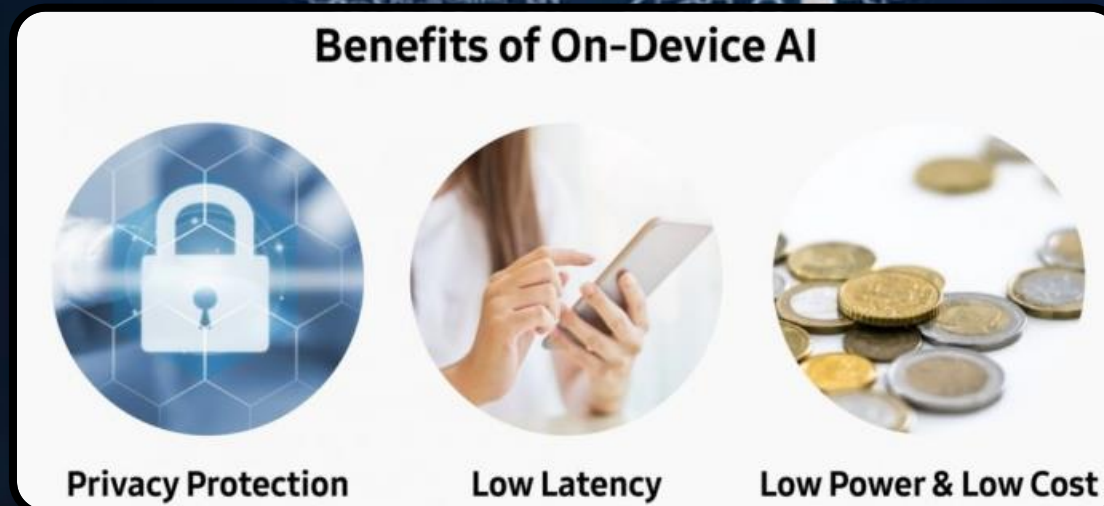
- OneMCC (Memory Coupled-Computing) is an open & standard S/W for PIM, PNM, CXL solutions
- Plan to provide standard programming model to support multi-architecture and domain
- Boost AI and HPC workloads with a variety of accelerators like CPUs, GPUs, and NPUs
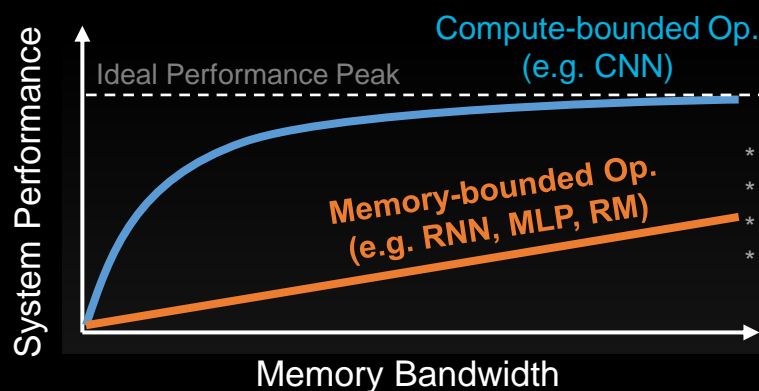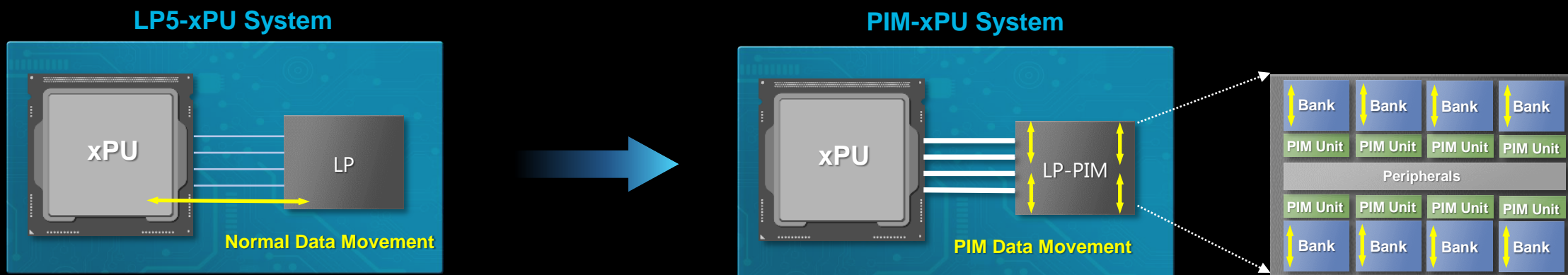
## OneMCC Infrastructure

### AI S/W Inference & Training

| Pytorch Framework | MCCPythonAPI |
|---|---|

| UDLC (Universal Deep Learning Compiler) | UDLR (Universal Deep Learning Runtime) |
|---|---|
| Front-end (Model/Data Analyzer) | Node Runtime (Cluster) |
| Middle-end (Scheduler, Optimizer) | Host Runtime (Server) |
| Back-end (NPU, GPU, CPU) | Device Runtime (NPU, GPU, CPU) |
| MCC Back-end (PIM, PNM, CXL) | MCC Runtime (PIM, PNM, CXL) |

### General Programming Model (SYCL, OpenACC/MP)

| SYCL-MCC Application | OpenACC/MP-MCC Application |
|---|---|

| SYCL Compiler | SYCL Runtime | OpenACC(MP) Compiler | |
|---|---|---|---|
| MCC Backend | MCC Runtime | flacc, clacc | MCC backend |

| SYCL Linker | SYCL Plugin | OpenACC(MP) Runtime | |
|---|---|---|---|
| MCC Object | PIM/PNM/CXL | CPU, GPU libs | PIM/PNM/CXL |

### System S/W (Kernel Level Driver)

| MCC Device Virtual Interface | | | |
|---|---|---|---|
| Unified Memory Manager | MCC Requestor | MCC Handler | System Function Acc. Library |

| MCC Core Driver | | |
|---|---|---|
| PIM Driver | PNM Driver | CXL Driver |

### Simulation Infrastructure

| Host System Simulator (CPU, NPU, GPU) | | |
|---|---|---|
| MCC Simulator Plugin Interface | | |
| PIM Simulator | PNM Simulator | CXL Simulator |

### Debug / Tools

| Benchmark Tool | MCC Monitor |
|---|---|
| Memory Profiler | Debugger |

# Processing-in-Memory for On-device Generative AI

- Expanding On-device AI Necessity:
  - Data center costs and power consumption are increasing due to the growing demand for cloud AI
  - Privacy concerns are rising as sensitive data is transmitted to the cloud for processing
  - Network connectivity is not always reliable or available, particularly in remote areas
- LPDDR-PIM  improves battery life by preventing memory over-provisioning just for bandwidth



**Benefits of On-Device AI**

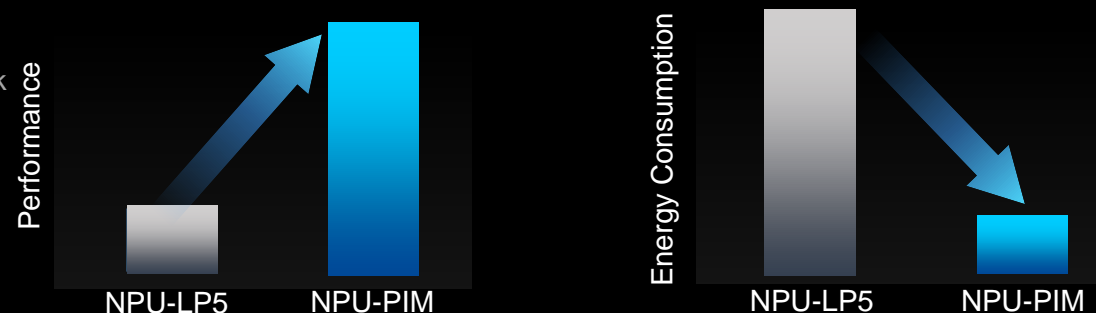**Privacy Protection**          **Low Latency**          **Low Power & Low Cost**

# LPDDR-PIM Concept

- Improve the 4.5x performance and save 72% of energy in the system with in-DRAM processing

  - Performance: Utilize up to 8× higher in-DRAM bandwidth by multi-bank parallel operation

  - Energy Efficiency: Reduce data movement energy by utilizing PIM unit

**LP5-xPU System**

xPU

LP

Normal Data Movement

**PIM-xPU System**

xPU

LP-PIM

PIM Data Movement

Bank  Bank  Bank  Bank

PIM Unit  PIM Unit  PIM Unit  PIM Unit

Peripherals

PIM Unit  PIM Unit  PIM Unit  PIM Unit

Bank  Bank  Bank  Bank

Ideal Performance Peak

Compute-bounded Op. (e.g. CNN)

Memory-bounded Op. (e.g. RNN, MLP, RM)

System Performance

Memory Bandwidth

* CNN: Convolutional Neural Network
* RNN: Recurrent Neural Network
* MLP: Multi-layer Perceptron
* RM: Recommendation Model
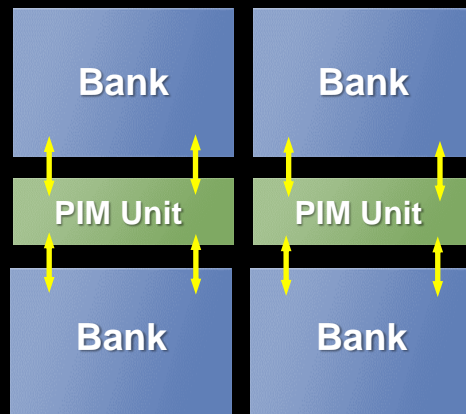
**System Performance and Energy Comparison**

Performance

NPU-LP5    NPU-PIM

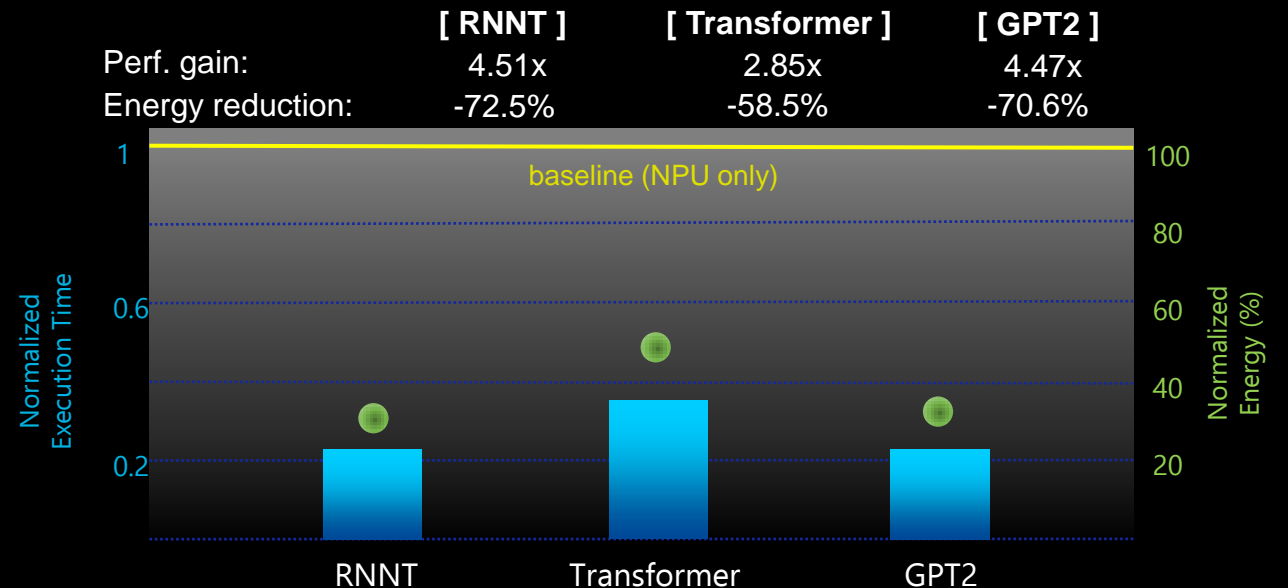Energy Consumption

NPU-LP5    NPU-PIM

# LPDDR-PIM Features

- Peak internal bandwidth: 102.4 GB/s

  - Using Bank-level parallelism, 8x bandwidth of the base LPDDR product
- Supporting native integer/floating point arithmetic and logical (and/or/...) operations
- Peak performance: 102.4 GFLOPS/s (FP16), 204.8 GOPS/s (INT8)
- Acceleration target: memory-bounded operations such as BLAS1 and BLAS2
- Samsung can support LPDDR-PIM simulator package to measure performance gain & energy reduction

BLAS1: Element-wise addition/multiplication
or layer normalization
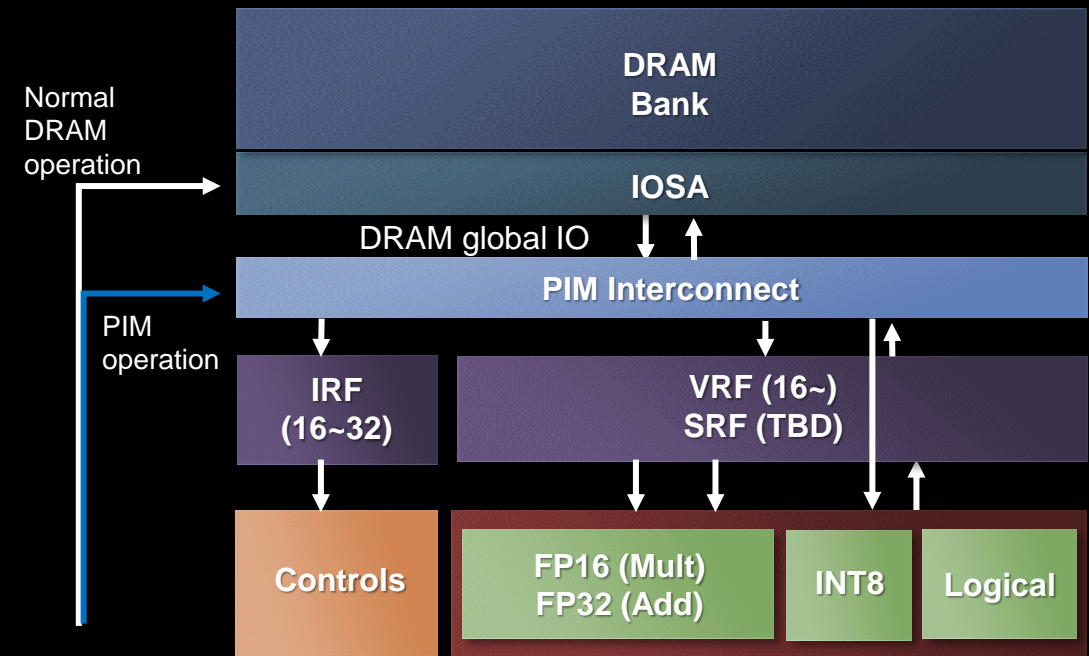BLAS2: vector-matrix multiplications

## Bank Parallelism

| Bank | Bank |
| :-: | :-: |
| PIM Unit | PIM Unit |
| Bank | Bank |

### End-to-end inference on AI application (simulation-based)

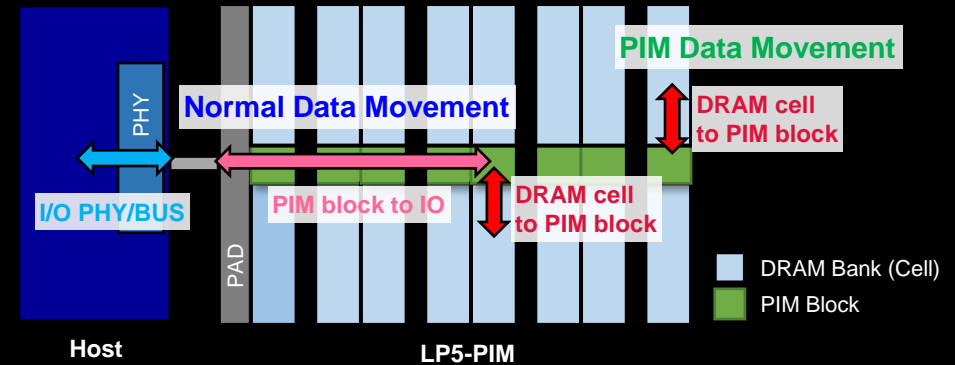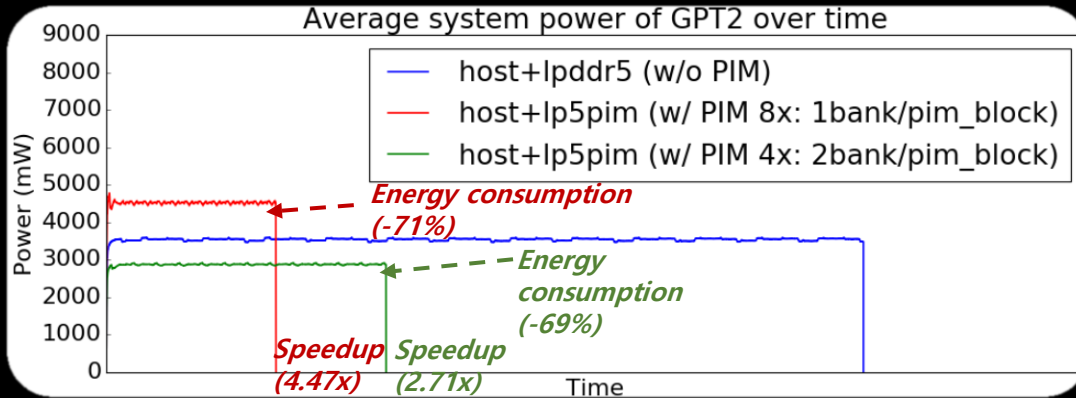| | [ RNNT ] | [ Transformer ] | [ GPT2 ] |
| --- | :-: | :-: | :-: |
| Perf. gain: | 4.51x | 2.85x | 4.47x |
| Energy reduction: | -72.5% | -58.5% | -70.6% |

# LPDDR-PIM Architecture

- PIM Unit is placed between every 1 bank (maximum performance) or 2 banks (moderate area overhead)
- PIM Unit: 256-bit SIMD FPU and registers (~640 bytes per PIM block)
  - Supporting operations: FP16 multiplication, FP32 accumulation, int8 arithmetic, etc.
  - PIM registers : Instruction (IRF), Vector (VRF), and Scalar (SRF)
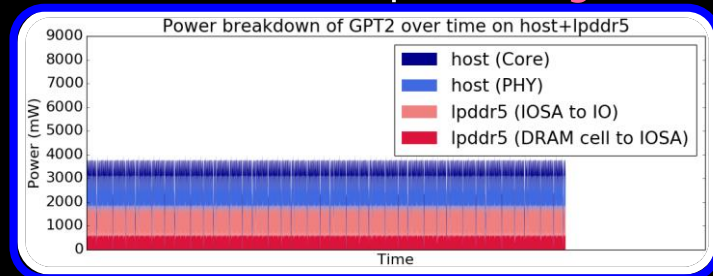


PAD
Bank

DRAM Bank (Cell)

PIM Block

PIM Support Modules

*Not scaled to actual physical silicon area

Normal DRAM operation

PIM operation

DRAM Bank

IOSA

DRAM global IO

PIM Interconnect

IRF (16~32)

VRF (16~)
SRF (TBD)

Controls

FP16 (Mult)
FP32 (Add)

INT8

Logical

# LPDDR-PIM System Performance/Power Analysis
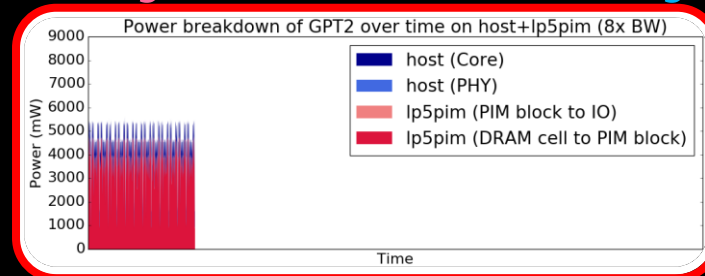
- Evaluate the performance and power consumption of GPT2
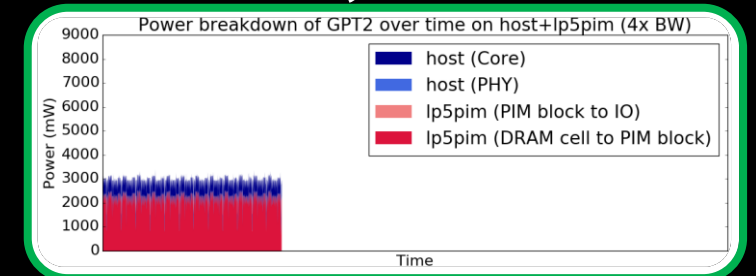  - LP5-PIM improves energy efficiency by shorter execution time



- Power consumption of DRAM internal component (red) increases proportionally
- Power consumption of global I/O bus (light red) and I/O PHYs (light blue) considerably decreases



[host + lpddr5]

[host + lp5pim (8x: 1 bank/PIM block)]

[host + lp5pim (4x: 2 bank/PIM block)]

* The power ratio is 17.9% (DRAM cell to IOSA), 31.3% (IOSA to IO), 31.1% (PHY), and 19.7% (core) for baseline.
* The power ratios are 85.3% (DRAM cell to PIM block), 14.6% (core), and 0.1% (etc) for lp5-pim(8x), and 76.9% (DRAM cell to PIM block), 23% (core), and 0.1% (etc) for lp5-pim(8x), respectively.

* Simulation experiment uses 4 memory (lp5/pim) channels
* lp5pim (DRAM cell to PIM block) includes ACT, PRE, IDLE, and REF
* The result of simulation is that more than 99% of memory traffic decreases by PIM
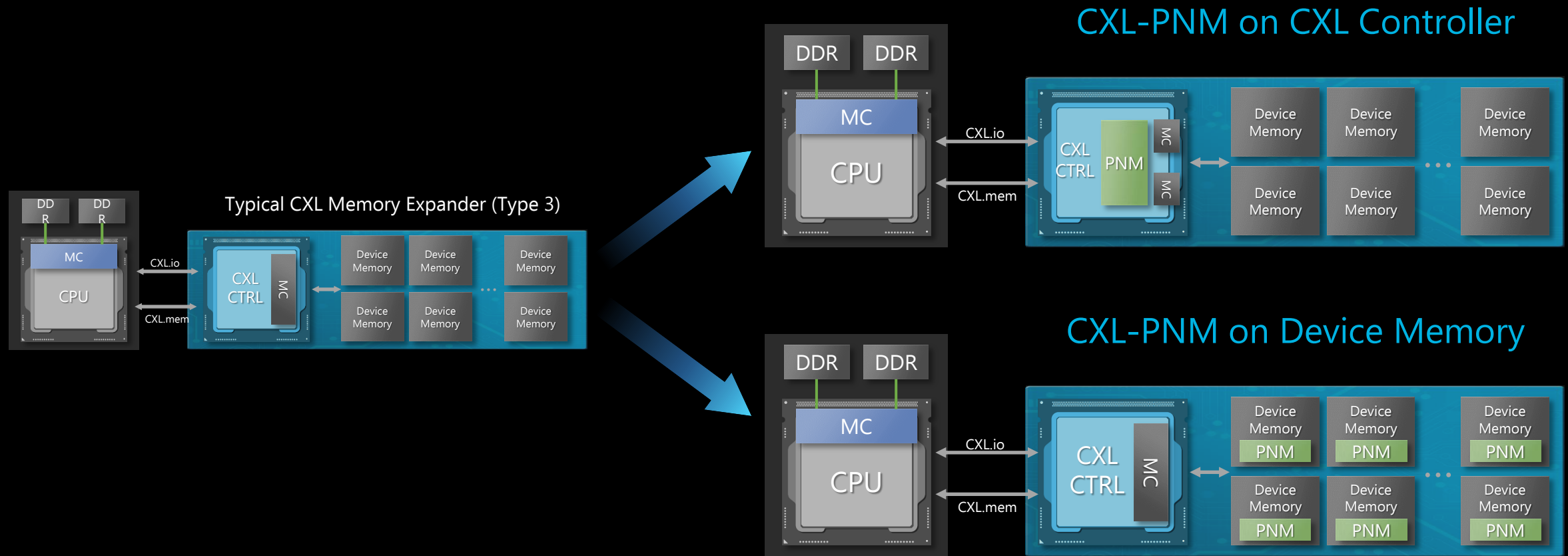* host (core) includes processing and IDLE power

SAMSUNG

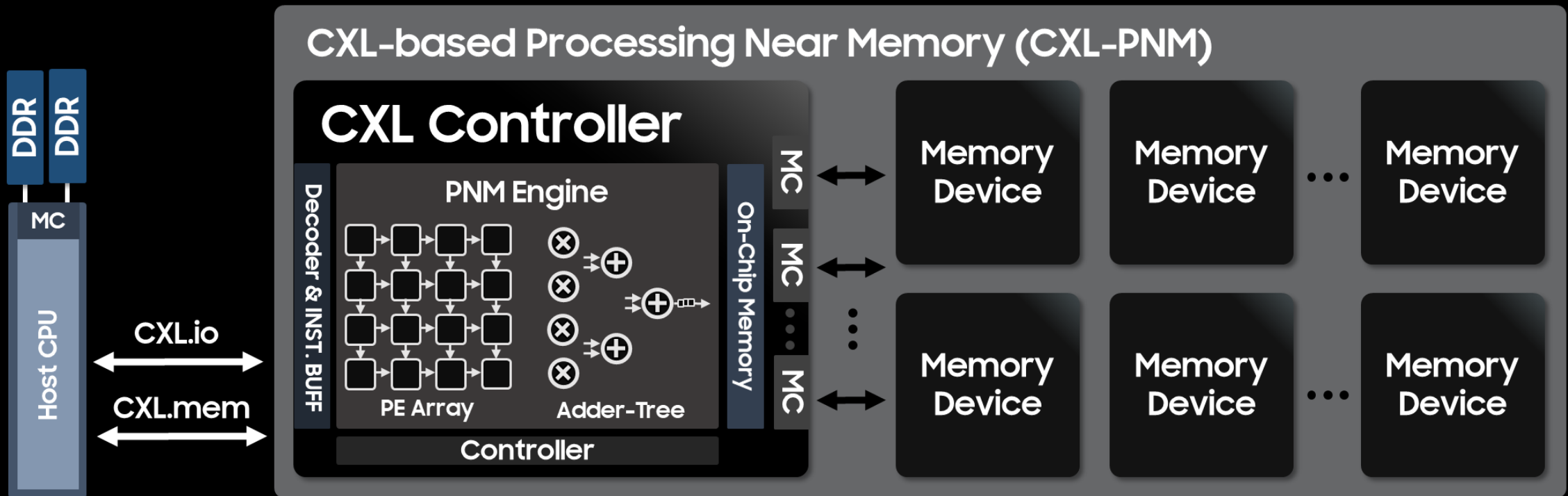CXL-PNM

Industry's 1st CXL-PNM (Processing-near-Memory)

# CXL-PNM Architecture

- A CXL-based Processing-near-Memory (PNM) Solution
- Two types of CXL-PNM: on CXL controller and on device memory



CXL-PNM on CXL Controller

Typical CXL Memory Expander (Type 3)

CXL-PNM on Device Memory
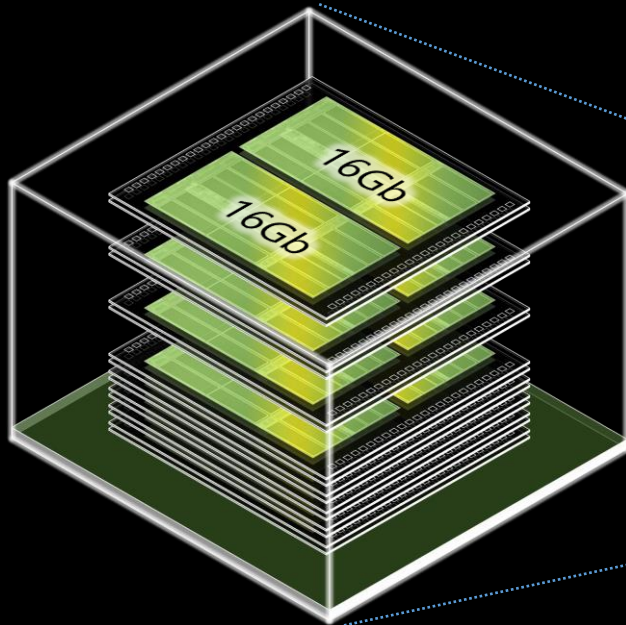
# CXL-PNM Architecture for GPT

- Heterogeneous compute unit (PE Array and Adder-Tree) on PNM Engine
  - Adder-tree designed to perform GEMV operation (Generation stage)
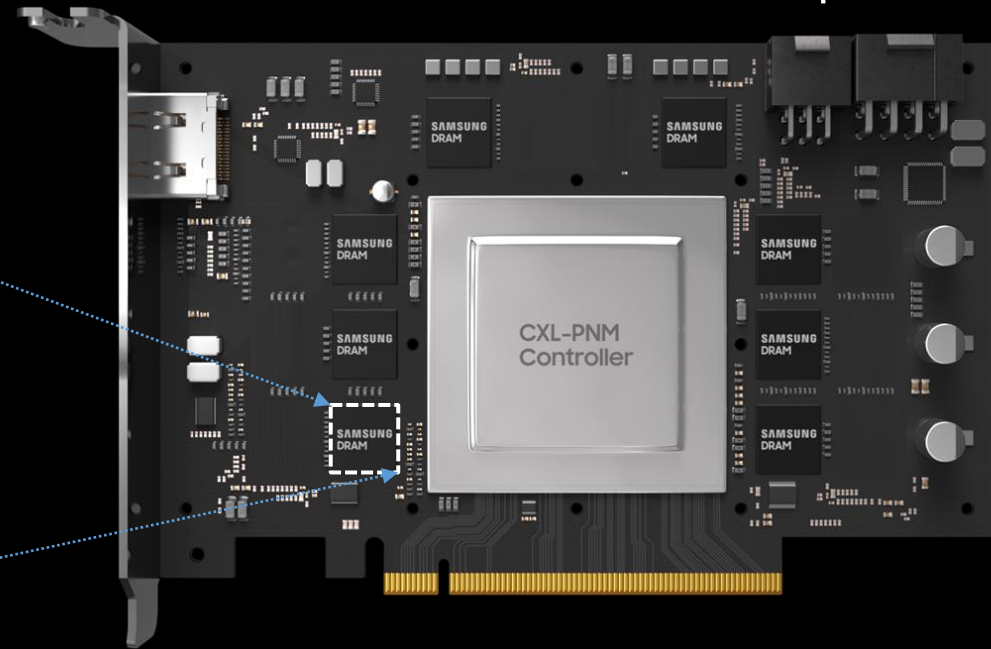  - PE Array for acceleration of GEMM operation (Summarization stage)

# 512GB 1.1TB/s CXL-PNM Concept

- CXL-PNM is able to be used for a wider range of systems including AI/ML accelerators
- Compared to other solutions, it can give unique trade-off among capacity, bandwidth, and power
- CXL-PNM can provide 512GB capacity and 1.1TB/s bandwidth
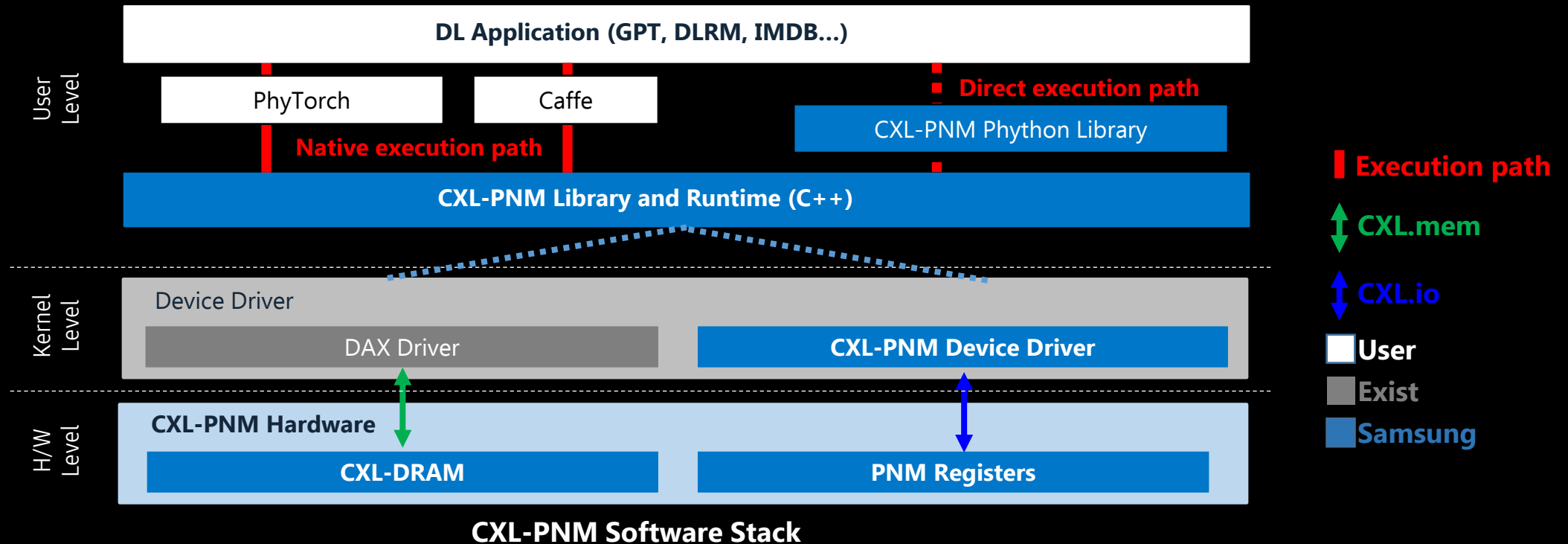
512Gb DRAM Package

CXL-PNM Module Concept

# CXL-PNM Software Stack

- CXL-PNM software stack for users to seamlessly and transparently utilizes the CXL-PNM platform
- CXL-PNM software stack includes user-level of library, runtime and kernel level of device driver
- CXL-PNM software stack supports two execution paths
  - Native execution path – Automatically offload PNM operation without modification of application source code
  - Direct execution path – Explicitly call PNM operations on the user application



**CXL-PNM Software Stack**

# Energy and Throughput Comparison for LLM

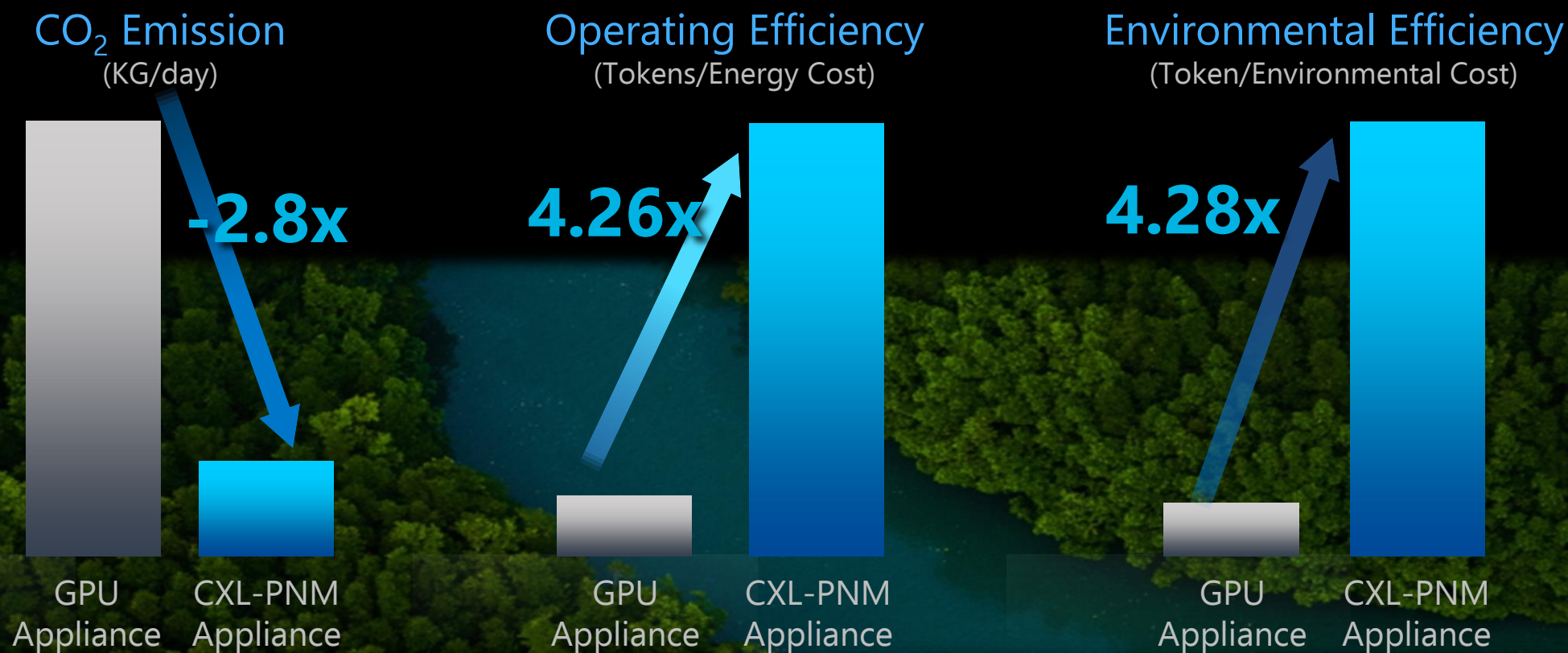- End-to-end performance evaluation of GPU and CXL-PNM with the same number of devices
  - CXL-PNM gives 2.9× higher energy efficiency, only 10.8% lower throughput, compared to A-GPU
  - However, CXL-PNM with large capacity can accelerate large-scale LLMs w/o any communication overhead
  - Multiple CXL-PNM can offer 4.4× higher energy efficiency and 53% higher throughput than multiple GPUs

## Energy Efficiency
(Token/Energy)

**2.89x**

**4.4x**

| A-GPU (1 device) | CXL-PNM (1 device) | A-GPU (8 devices) | CXL-PNM (8 devices) |

## Throughput
(Token/Second)

**-0.1x**

**1.53x**

| A-GPU (1 device) | CXL-PNM (1 device) | A-GPU (8 devices) | CXL-PNM (8 devices) |

• Single Device (OPT-13B), 8 Device (OPT-66B), input token(64) output token(1024)

# CO$_2$ Emission Reduction by CXL-PNM

- Operating/Environmental (Energy/CO$_2$) cost of GPU/CXL-PNM appliance with eight devices
  - GPU appliance is 2.8× more expensive than CXL-PNM appliance for the environmental cost
  - Operating efficiency of CXL-PNM appliance reduces the amount of CO$_2$ emission
  - CXL-PNM appliance is 4.3× more efficient than that of GPU appliance



CO$_2$ Emission
(KG/day)

-2.8x

GPU Appliance     CXL-PNM Appliance

Operating Efficiency
(Tokens/Energy Cost)

4.26x

GPU Appliance     CXL-PNM Appliance

Environmental Efficiency
(Token/Environmental Cost)

4.28x

GPU Appliance     CXL-PNM Appliance

- GPU Appliance (A-GPU server System with 8 A-GPU, FasterTransformer), CXL-PNM Appliance (8 CXL-PNM Card)

# Samsung AI memory