# UM-PIM: DRAM-based PIM with Uniform & Shared Memory Space

**Yilong Zhao,** Mingyu Gao, Fangxin Liu*, Yiwei Hu, Zongwu Wang, Han Lin, Ji Li, He Xian, Hanlin Dong, Tao Yang, Naifeng Jing, Xiaoyao Liang, Li Jiang*

Shanghai Jiao Tong University
2024/7/2

Outline

Background: Process-in-memory and Memory Interleaving

UM-PIM: DRAM-based PIM with Uniform & Shared Memory Space

Evaluation

Outline

Background: Process-in-memory and Memory Interleaving

UM-PIM: DRAM-based PIM with Uniform & Shared Memory Space

Evaluation

# Background: the "Memory Wall"

⌘ According to Computation / Memory

- Compute bound: Matrix Multiplication

- Memory bound: Element-wise (dropout，masking)
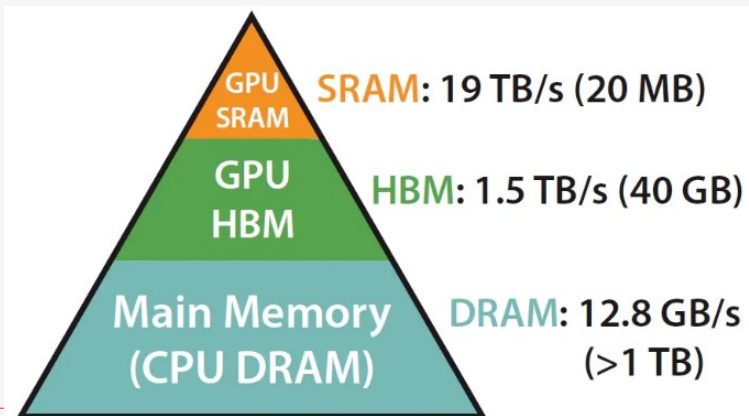  Reduce (Layer nom, Sum)

⌘ **Memory bottleneck becomes severe with the emergence of large models**

**Attention on GPT2**



**Memory increases with Token linearly**

**Llama2-13B Model**

| Token len | 4K | 16K | 128K |
|---|---|---|---|
| FP Memory | 3G | 12G | 100G |
| BP Memory | 10G | 39G | 320G |

**Accelerate with DRAM?**



GPU SRAM — SRAM: 19 TB/s (20 MB)

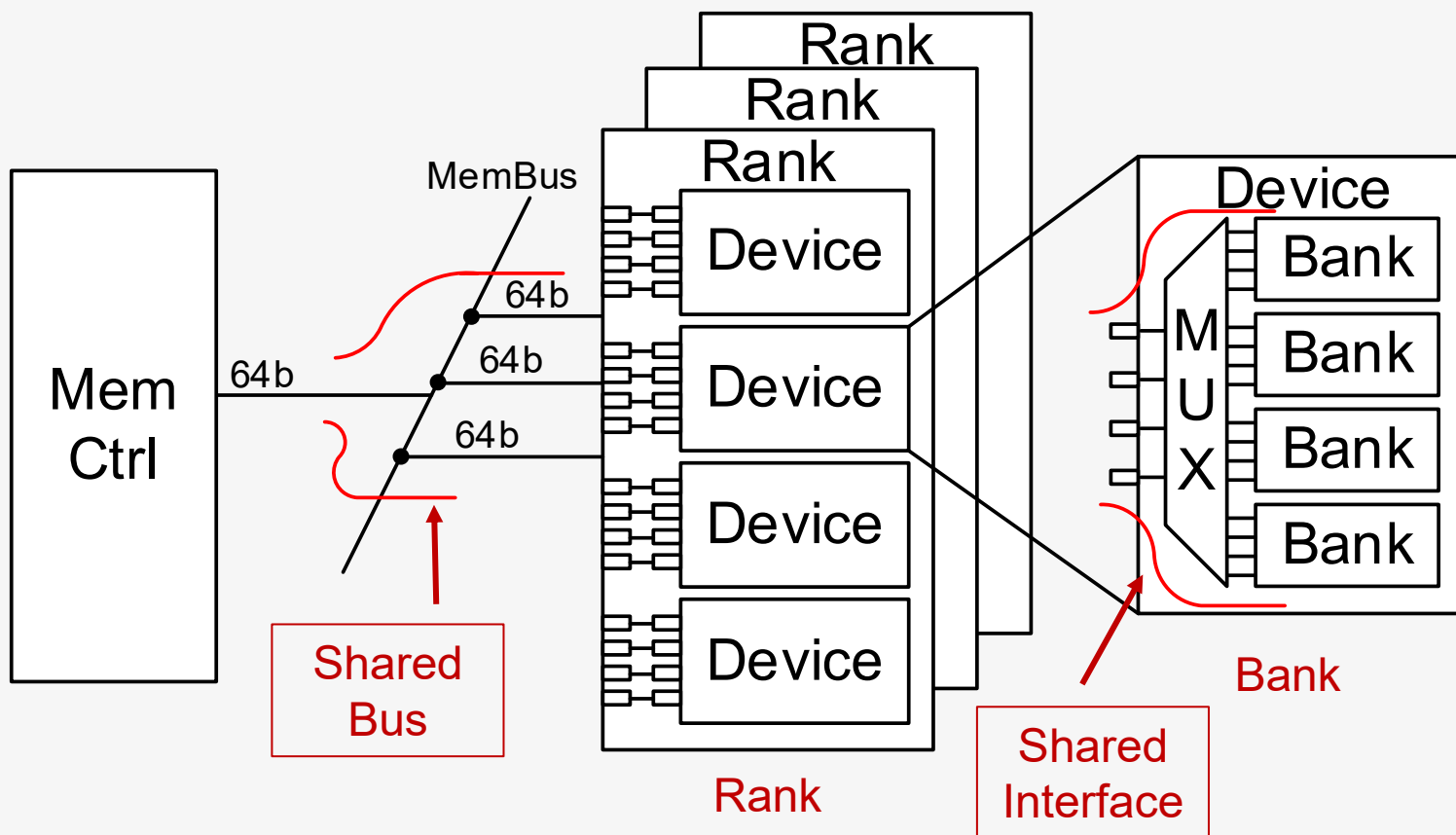GPU HBM — HBM: 1.5 TB/s (40 GB)

Main Memory (CPU DRAM) — DRAM: 12.8 GB/s (>1 TB)

# Background: the "Memory Wall" in DRAM

- DRAM ranks/banks can only be accessed sequentially, and shares data line!



Rank

MemBus

Mem Ctrl

64b

64b

64b

64b

Device

Device

Device

Device

Rank

Shared Bus

Device
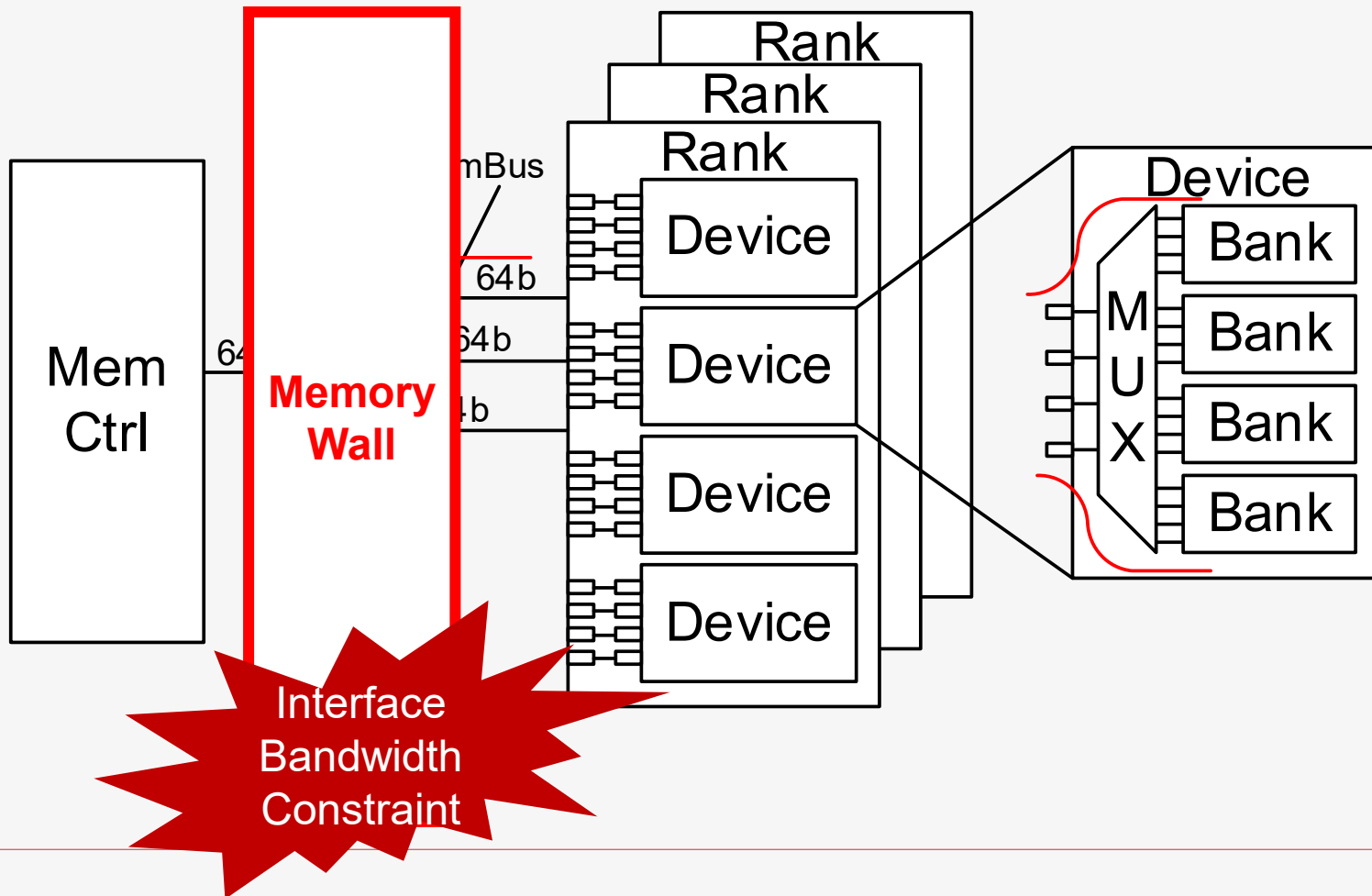
M U X

Bank

Bank

Bank

Bank

Bank

Shared Interface

Ranks (Banks) share memory interface

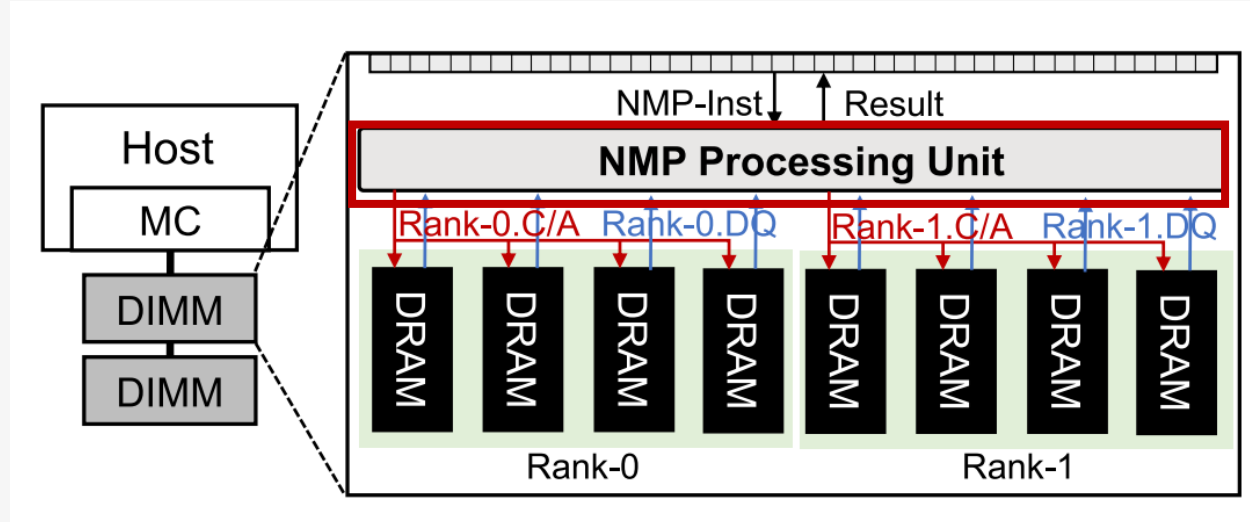Only one Rank (Bank) can be activated at the same time!

# Background: the Memory Wall

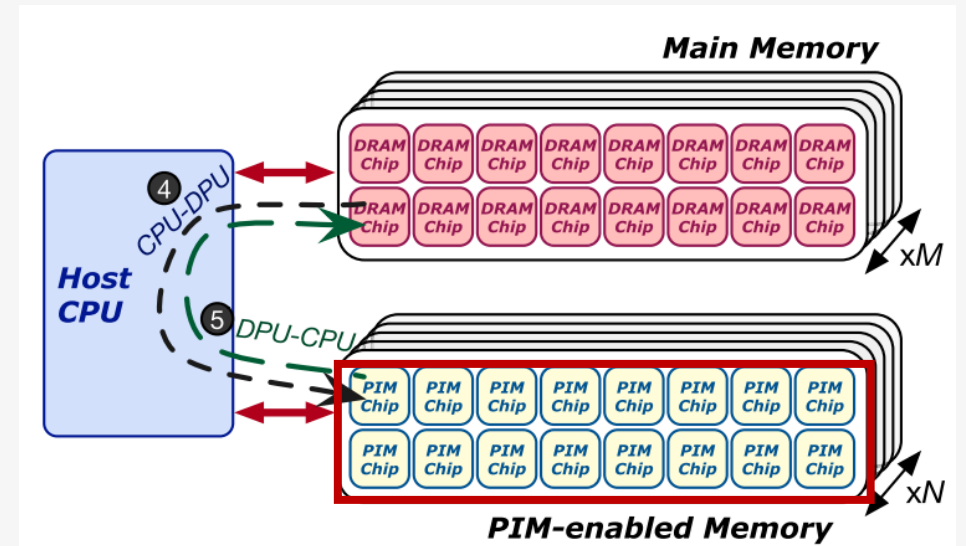- The INTERFACE BANDWIDTH CONSTRAINT makes a **Memory Wall** Between CPU and DRAM!

# Background: Process in Memory (PIM)

- Integrate computing units (**PIM units**) within DRAM to better utilize internal bandwidth
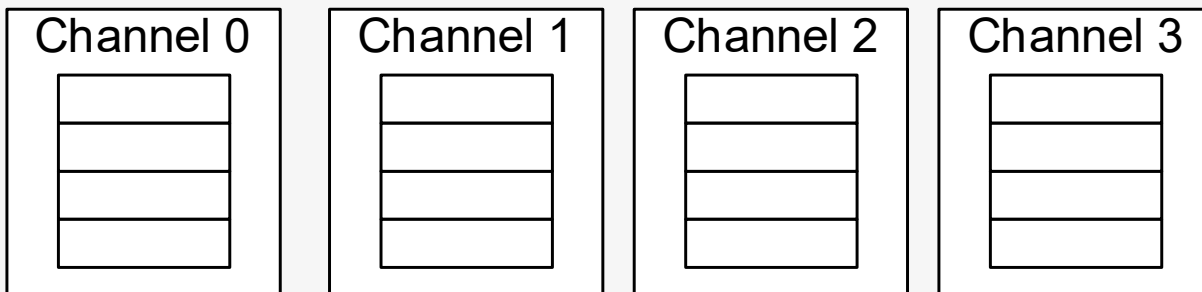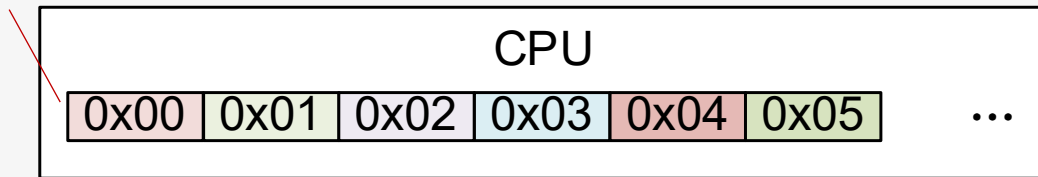


AxDIMM: in DIMM



UPMEM, AiM: In Bank

# Memory Interleaving

To increase CPU bandwidth:

- Adjacent data block stored in different channels
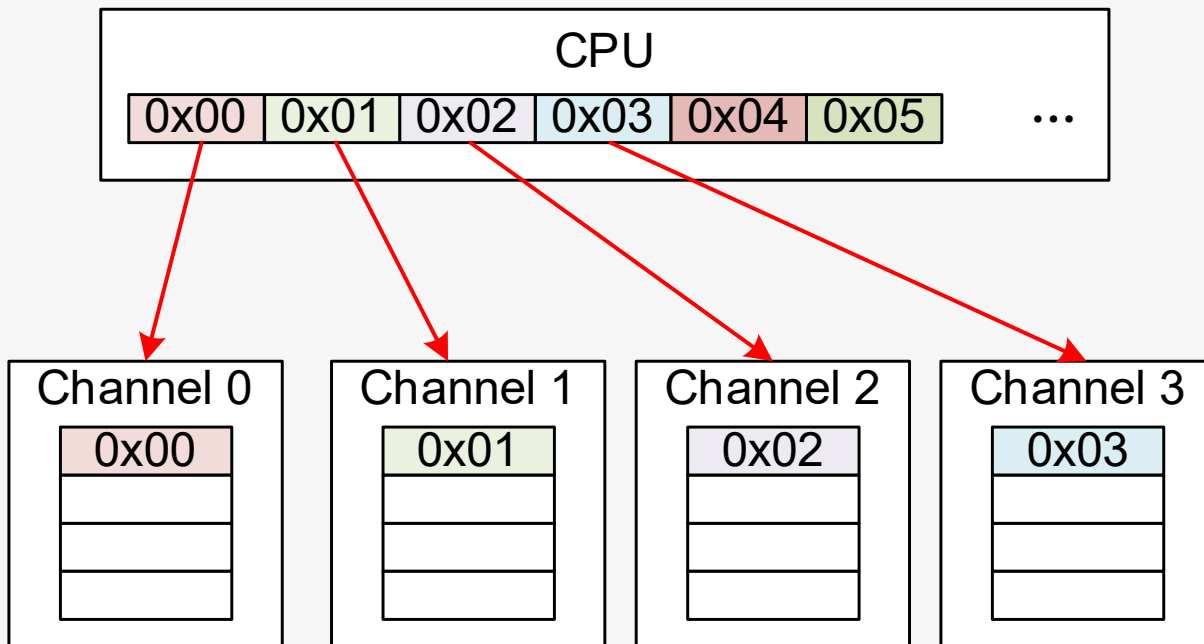
- And Access in Parallel/pipeline



Cache Line   (64B)

| CPU | | | | | |
|---|---|---|---|---|---|
| 0x00 | 0x01 | 0x02 | 0x03 | 0x04 | 0x05 | ... |

Channel 0

Channel 1

Channel 2

Channel 3

# Memory Interleaving
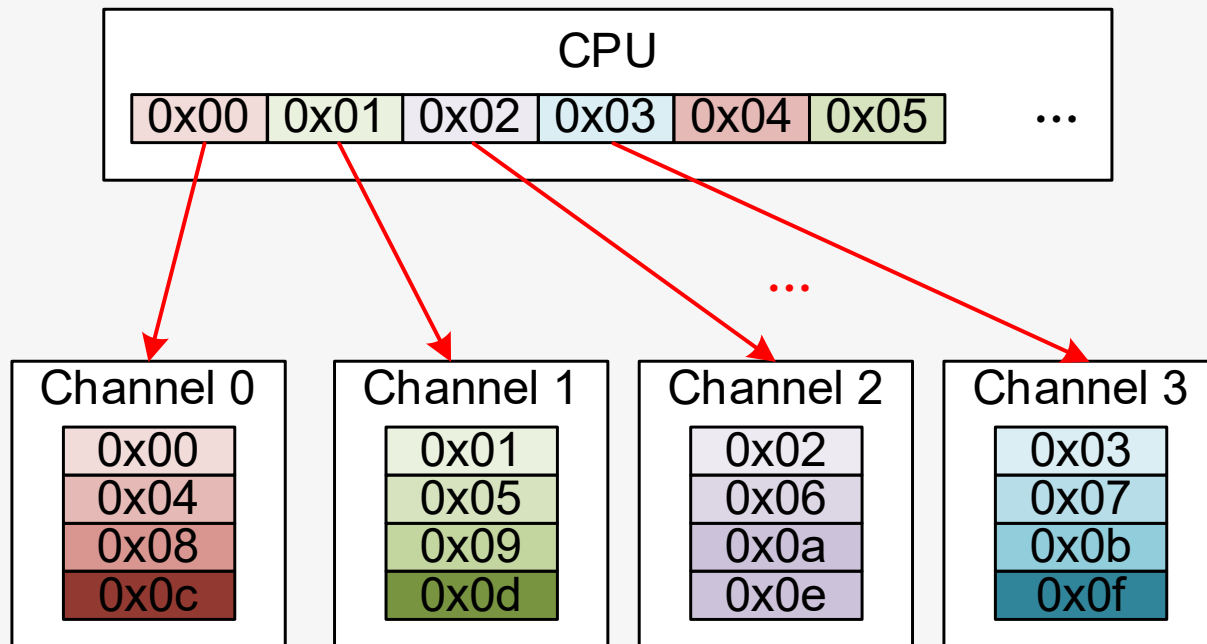
To increase CPU bandwidth:

- Adjacent data block stored in different channels

- And Access in Parallel/pipeline

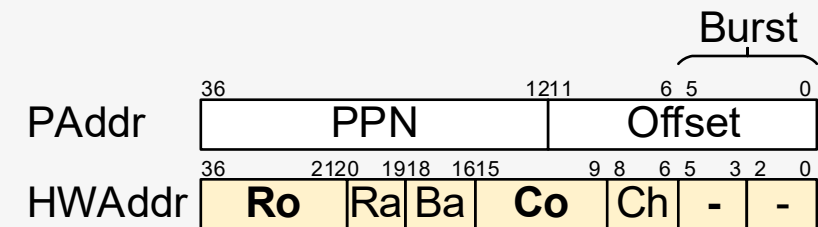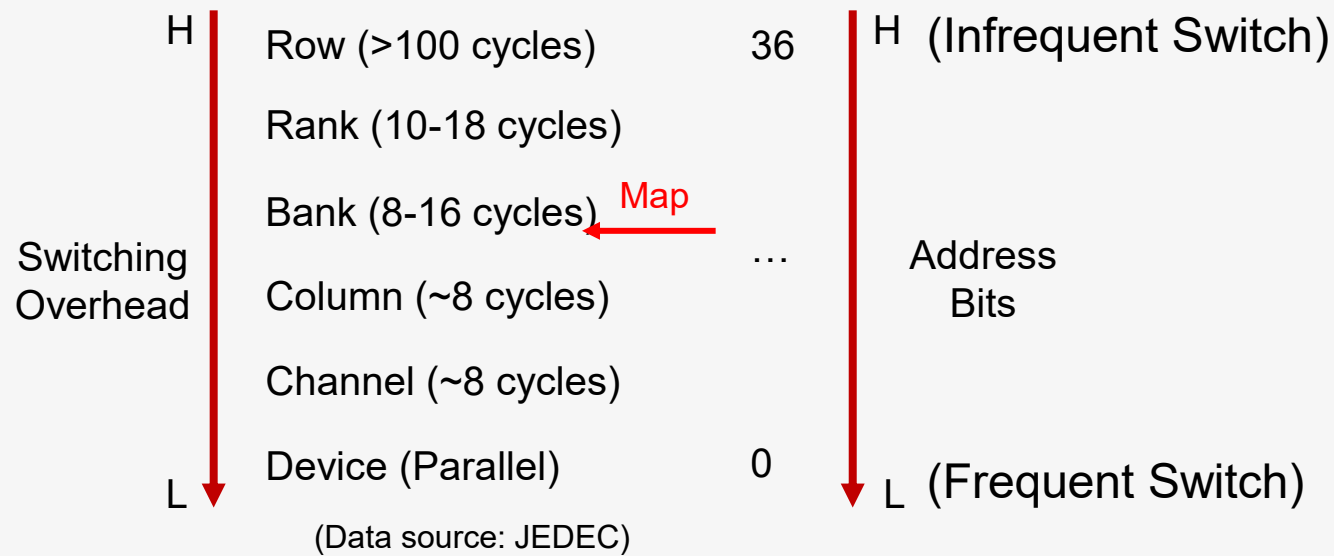# Memory Interleaving

To increase CPU bandwidth:

- Adjacent data block stored in different channels

- And Access in Parallel/pipeline
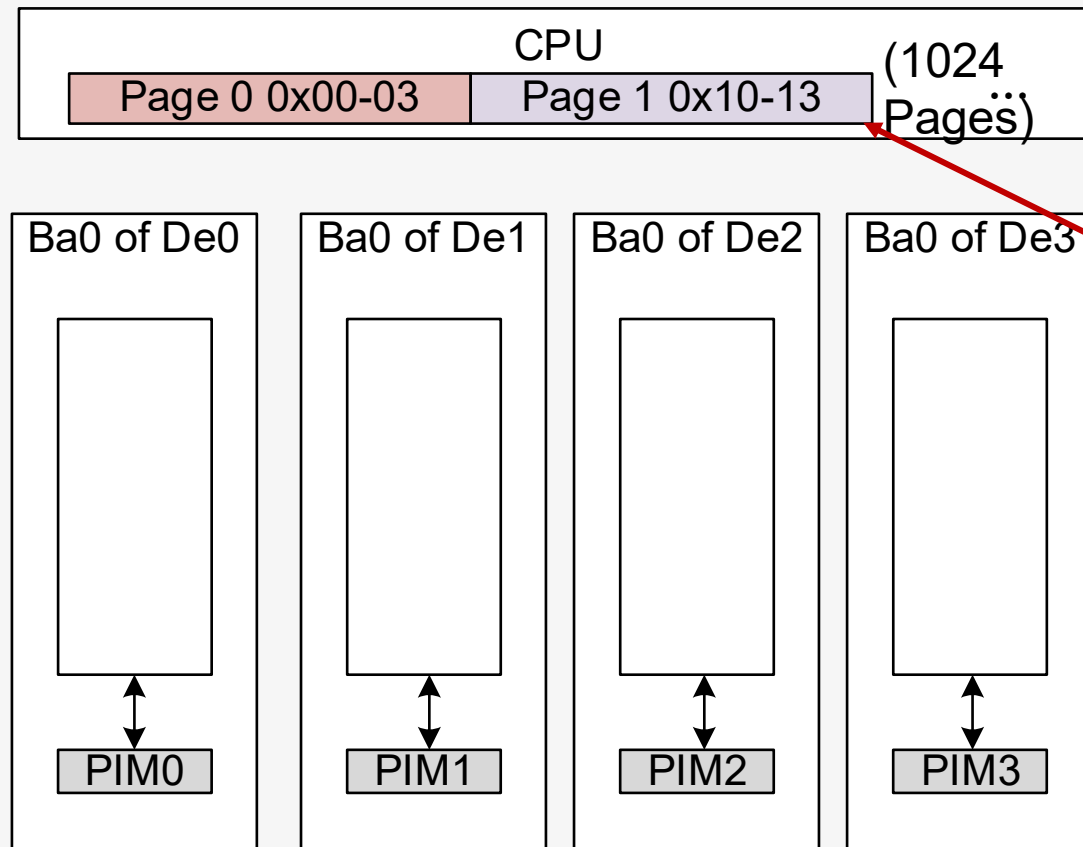
# Memory Interleaving with Address Mapping

- According to **Switch Overhead**

- Map **lower**-order bits to **low**-switch-overhead levels, **higher**-order bits to **high**-switch-overhead levels

H | Row (>100 cycles)
Rank (10-18 cycles)
Bank (8-16 cycles) — Map
Switching Overhead | Column (~8 cycles)
Channel (~8 cycles)
L | Device (Parallel)
(Data source: JEDEC)

36 | H (Infrequent Switch)
...
Address Bits
0 | L (Frequent Switch)

Burst

PAddr | 36 ... 1211 ... 6 5 ... 0
PPN | Offset

HWAddr | 36 ... 2120 1918 1615 ... 9 8 ... 6 5 ... 3 2 ... 0
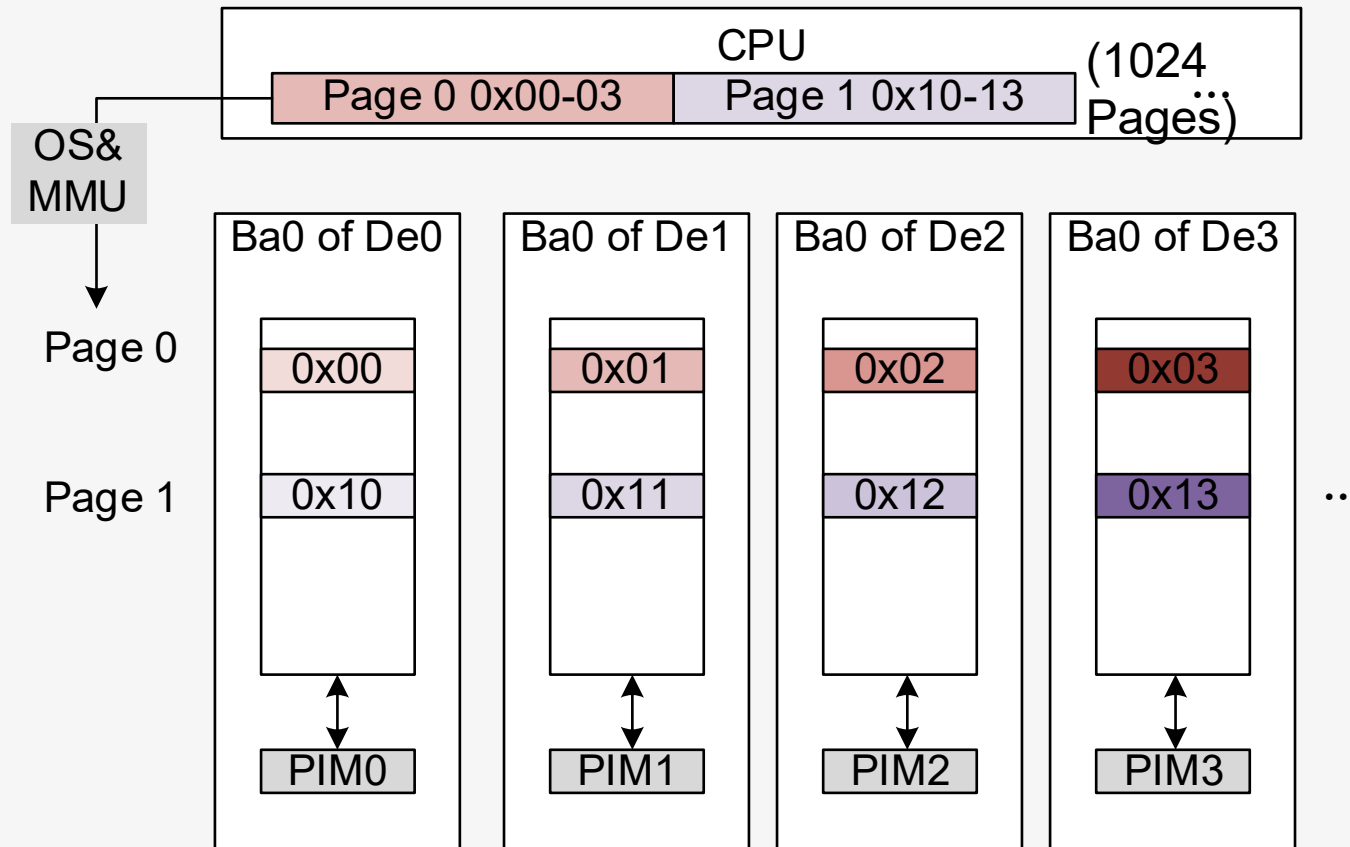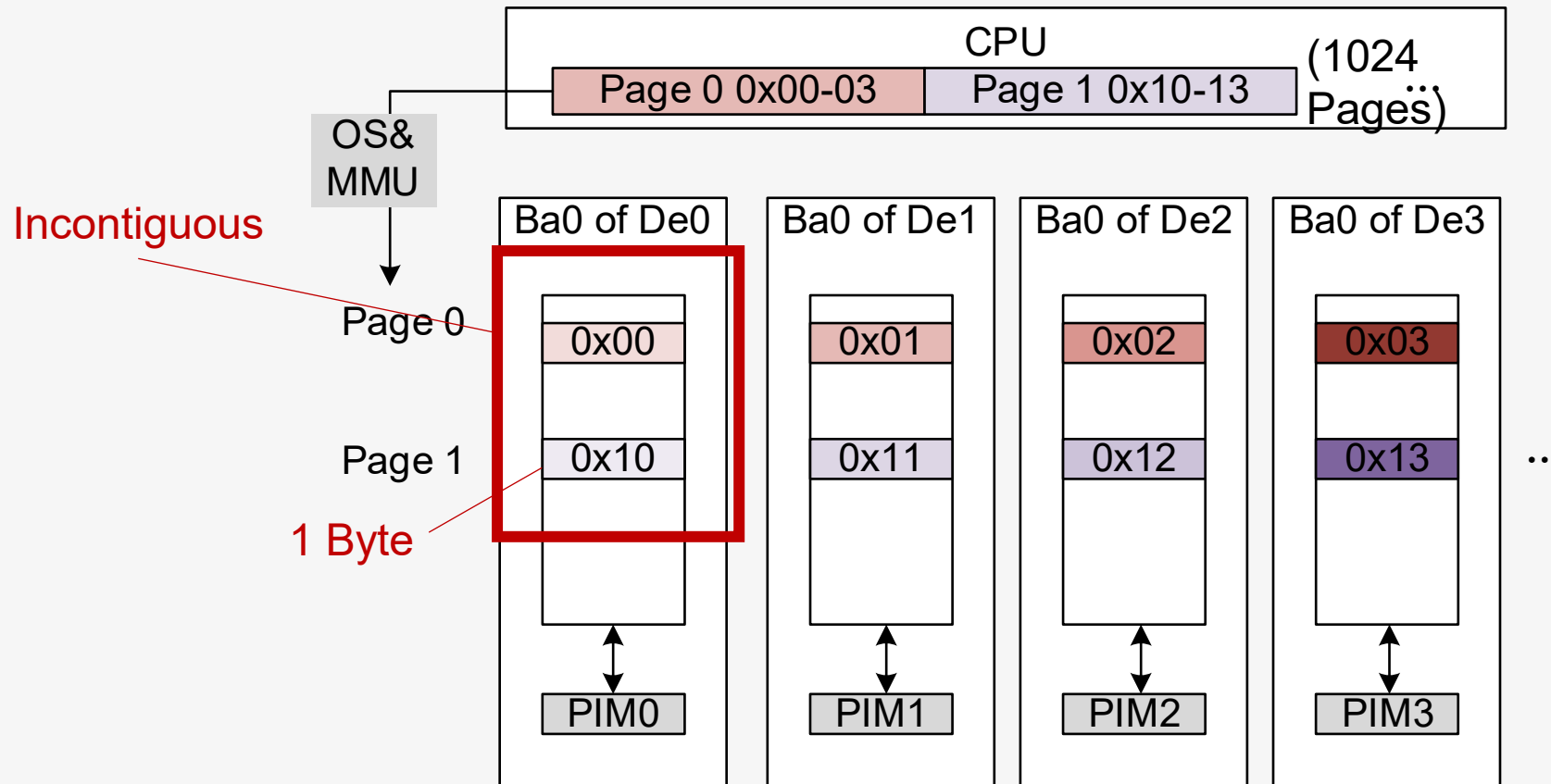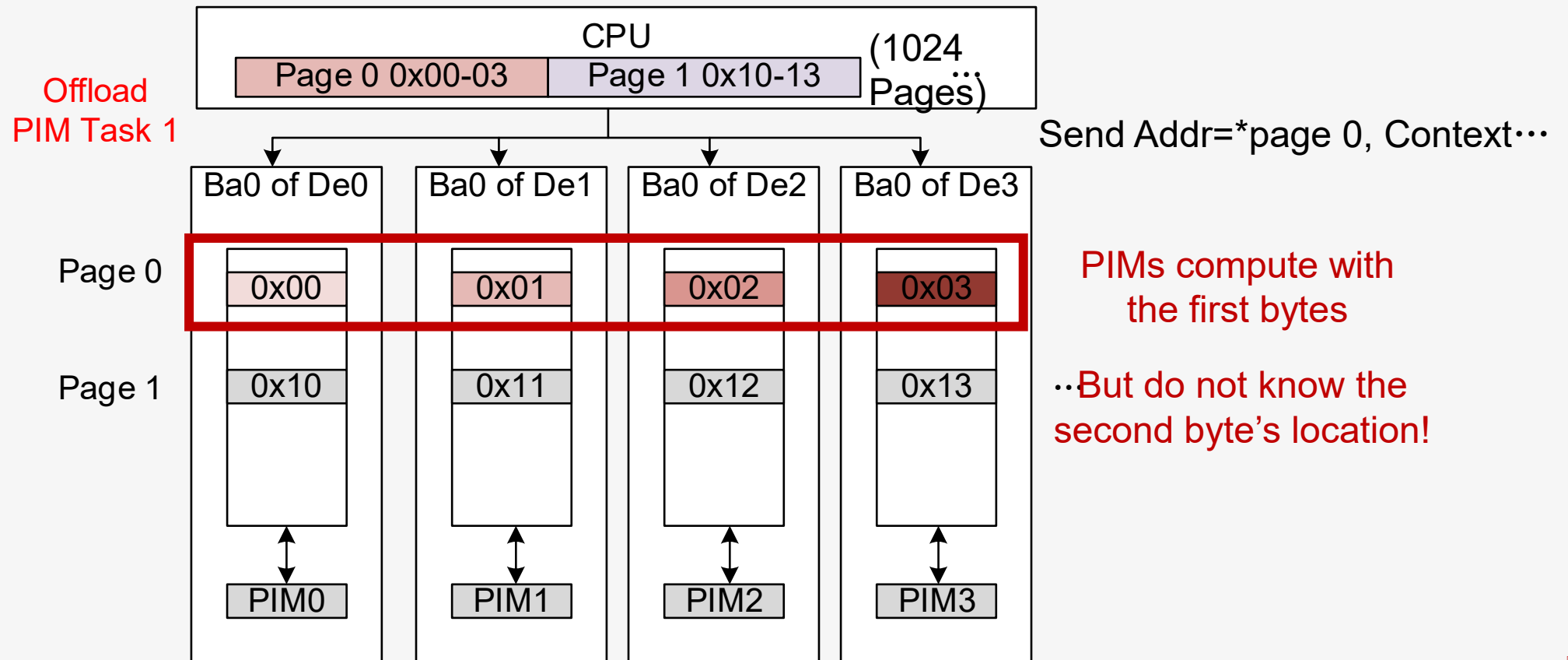**Ro** | Ra | Ba | **Co** | Ch | - | -

# PIM vs Memory Management: Incompatibility

- Memory interleaving, virtual memory limit the length of contiguous data block visible to PIM, limit offload granularity

Memory interleaving, virtual memory limit the length of contiguous data block visible to PIM, limit offload granularity

# PIM vs Memory Management: Incompatibility

Memory interleaving, virtual memory limit the length of contiguous data block visible to PIM, limit offload granularity
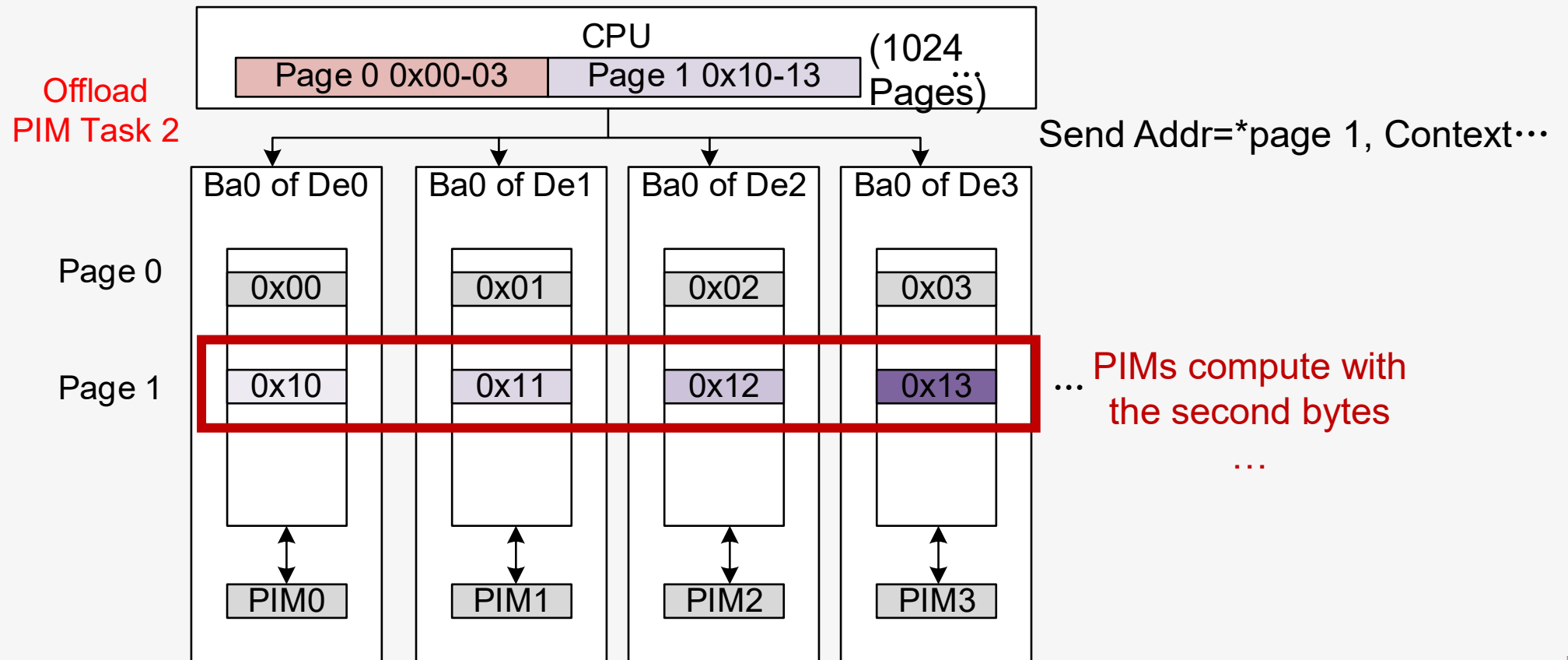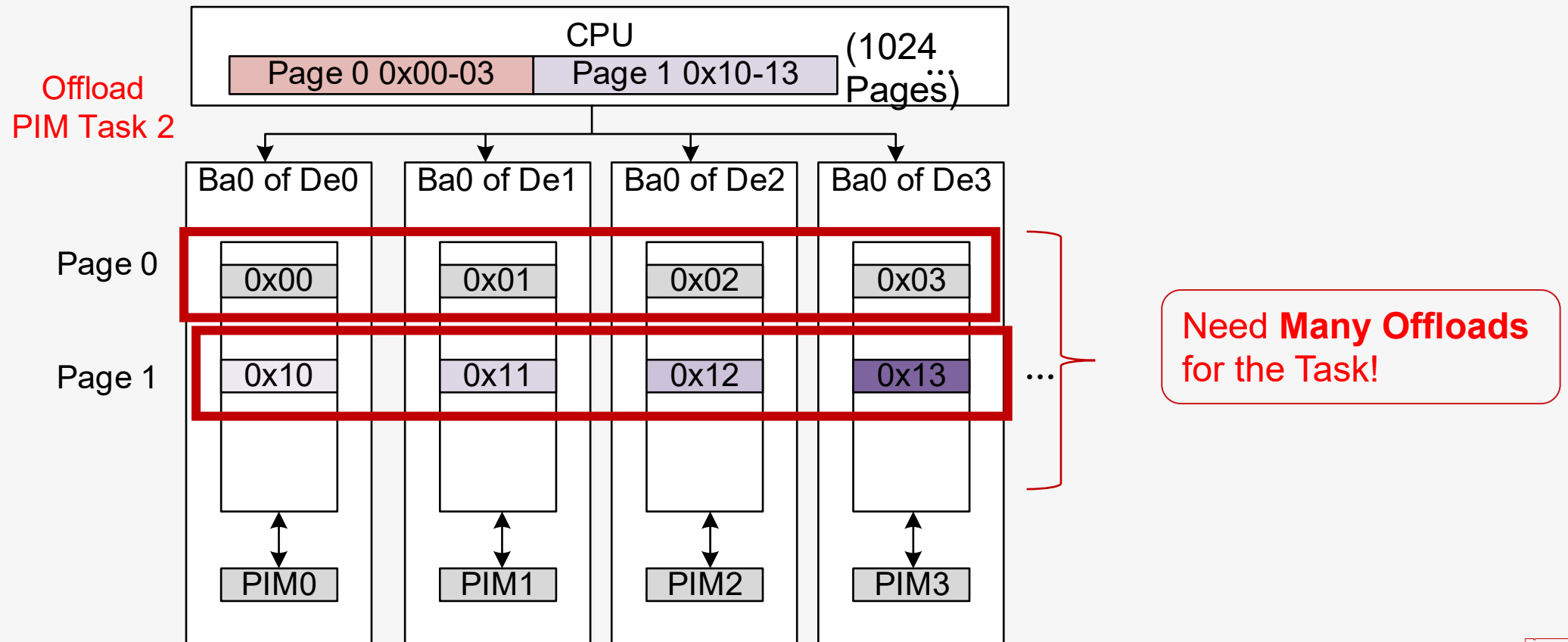
# PIM vs Memory Management: Incompatibility

- Memory interleaving, virtual memory limit the length of contiguous data block visible to PIM, limit offload granularity

# PIM vs Memory Management: Incompatibility

Memory interleaving, virtual memory limit the length of contiguous data block visible to PIM, limit offload granularity
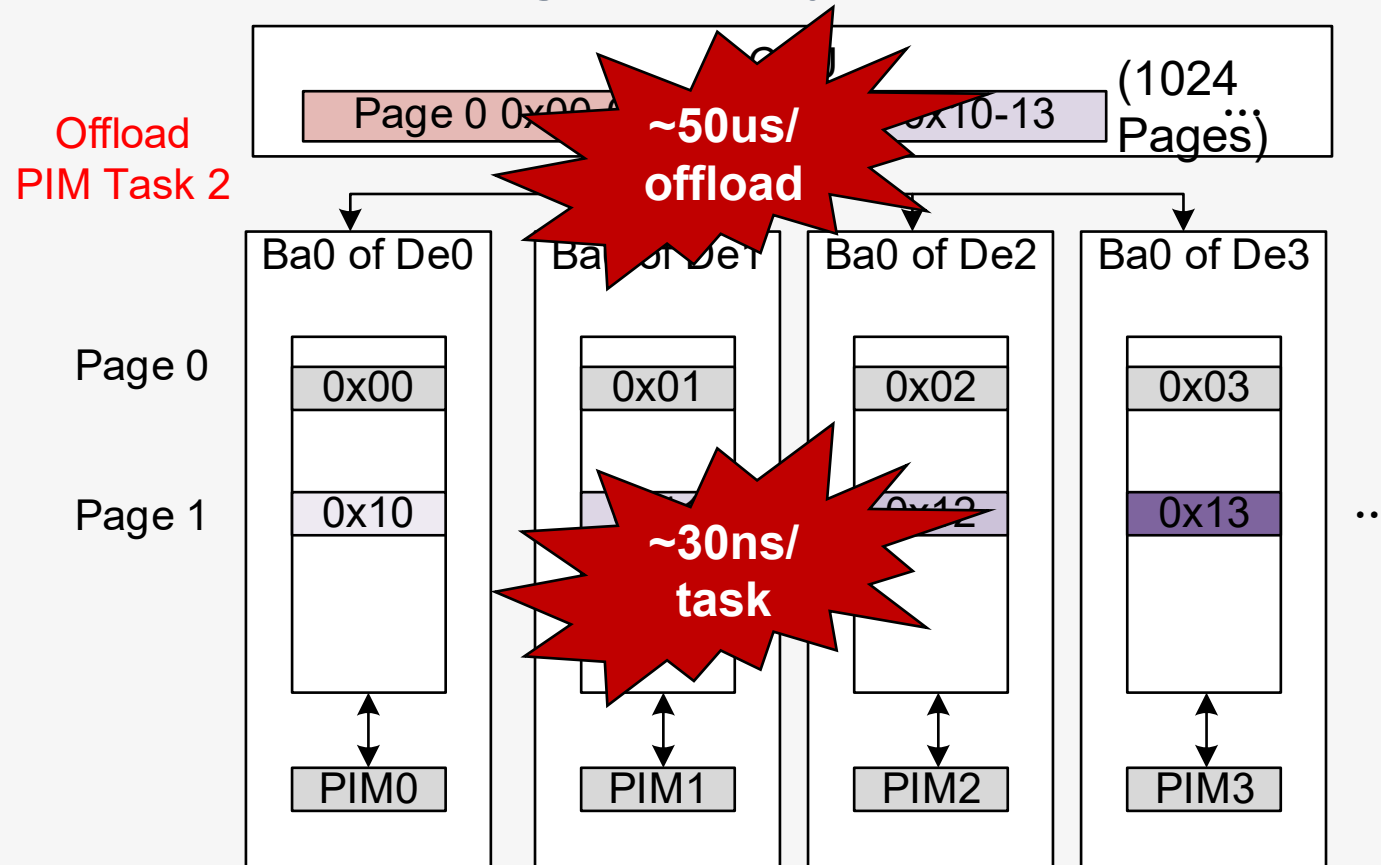
# PIM vs Memory Management: Incompatibility

- Memory interleaving, virtual memory limit the length of contiguous data block visible to PIM, limit offload granularity

Memory interleaving, virtual memory limit the length of contiguous data block visible to PIM, limit offload granularity



Offload PIM Task 2

Page 0 0x00 ... 0x10-13 (1024 Pages)

~50us/offload

Ba0 of De0 | Ba0 of De1 | Ba0 of De2 | Ba0 of De3

Page 0: 0x00 | 0x01 | 0x02 | 0x03

Page 1: 0x10 | 0x11 | 0x13

~30ns/task

PIM0 | PIM1 | PIM2 | PIM3

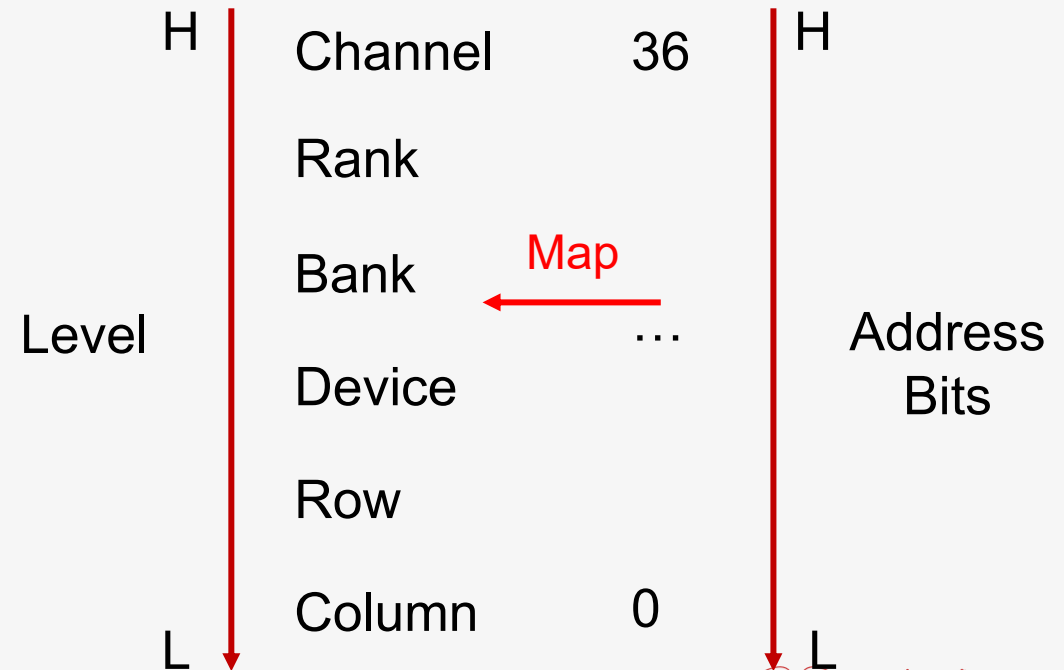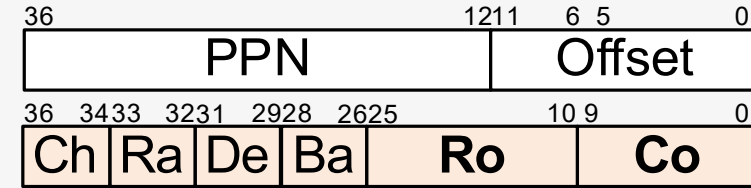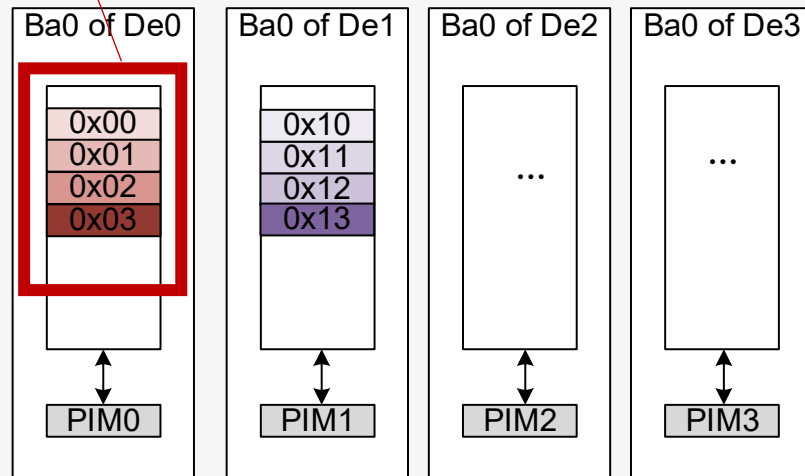**Need 1024 Offloads for the Task!**

50us offload
30ns task
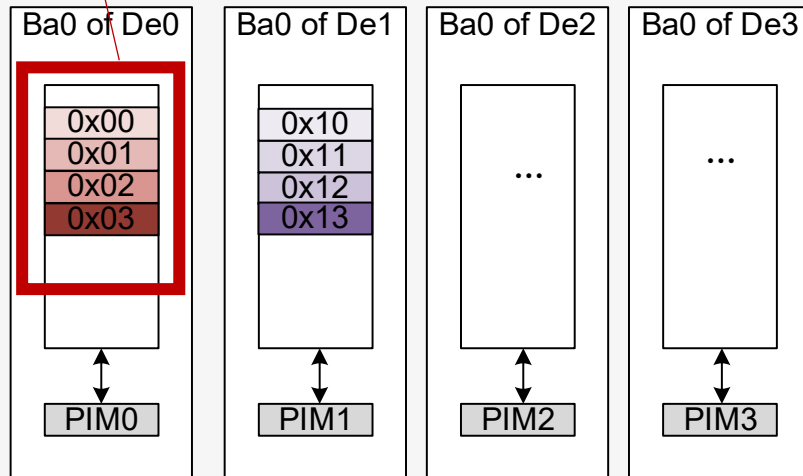**Offload >> Task**
**In-efficient Offload**

# PIM's ideal address mapping

PIM prefers contiguous data block: address mapping according to level

Only Need One Offload: Less offload overhead
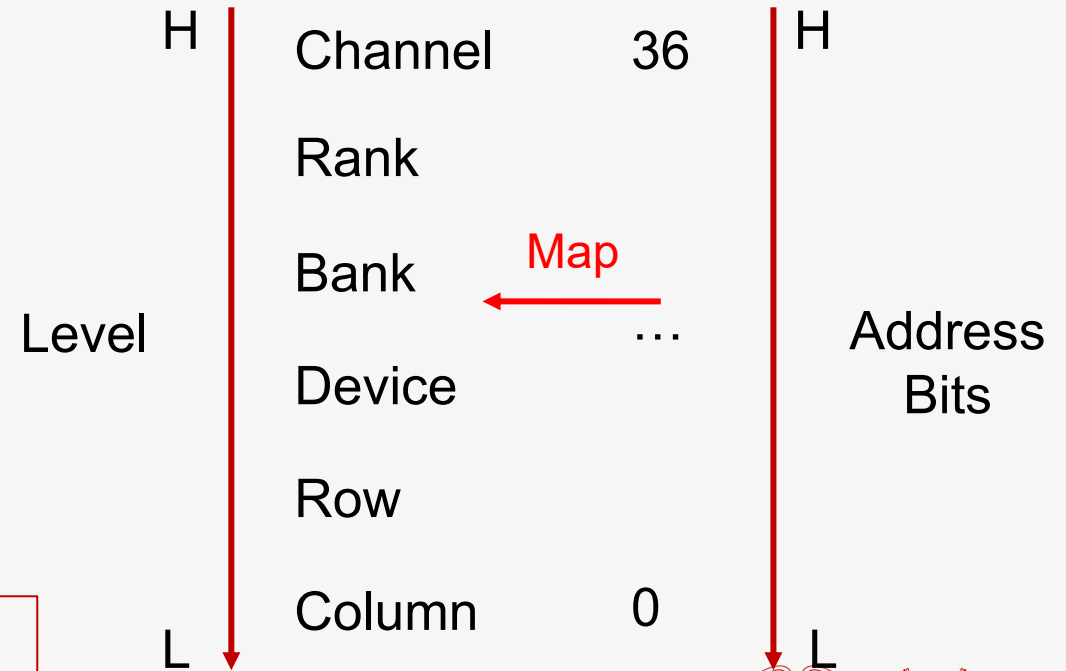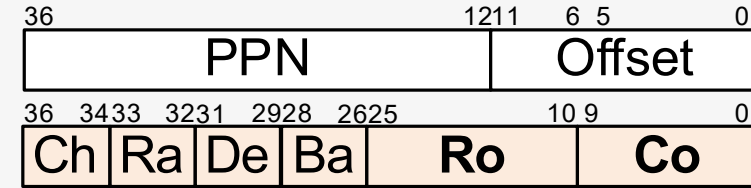
Contiguous
Data Block

PIM prefers contiguous data block: address mapping according to level



Only Need One Offload: Less offload overhead

Contiguous Data Block

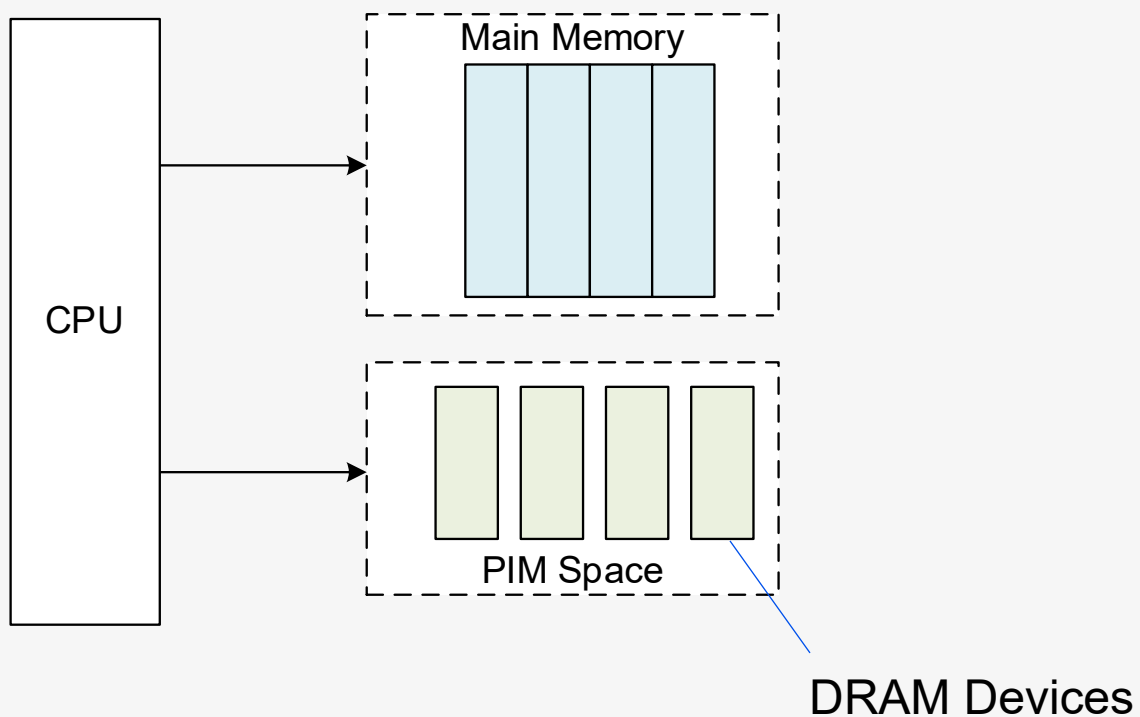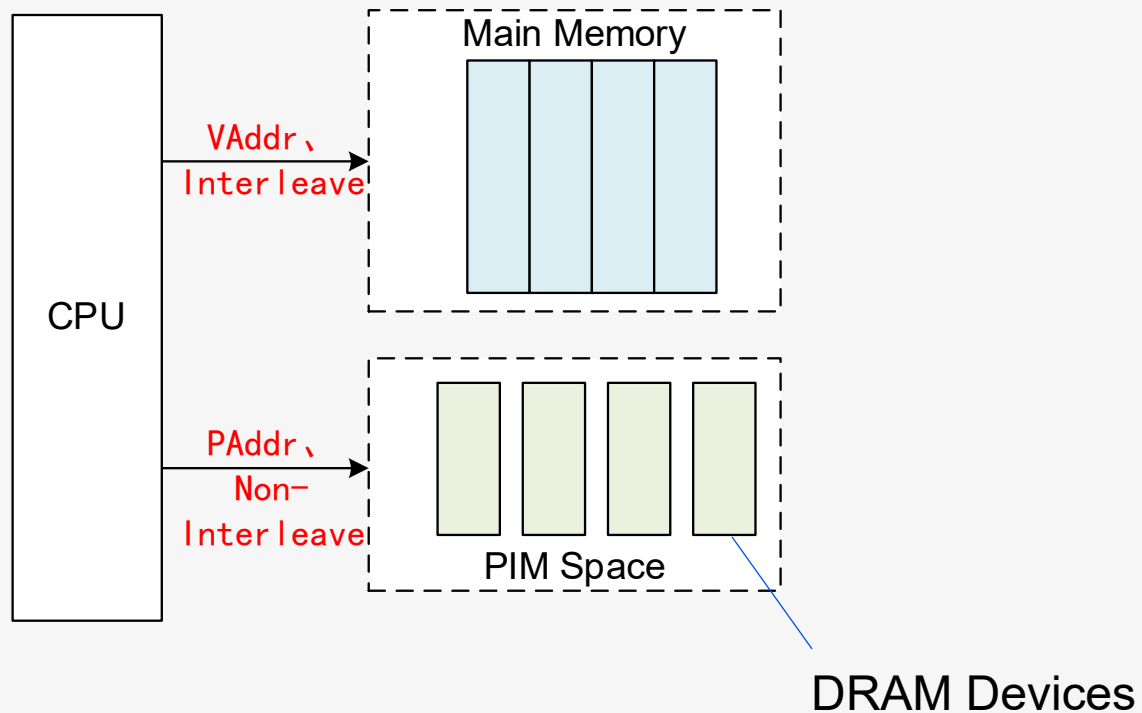CPU & PIM require different Data Layouts!

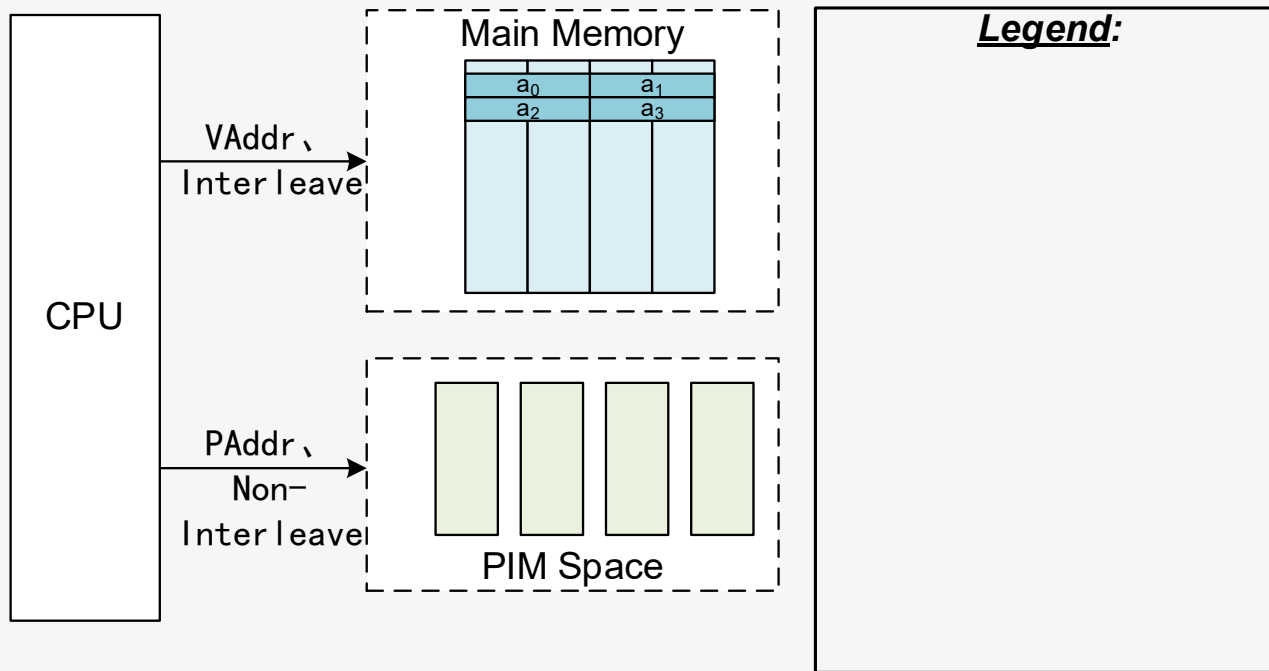# Measures of current PIM systems

## Measure 1: Isolated Memory Space

# Measures of current PIM systems

**Measure 1: Isolated Memory Space**

## Measure 1: Isolated Memory Space



Example: MLP

# Measures of current PIM systems

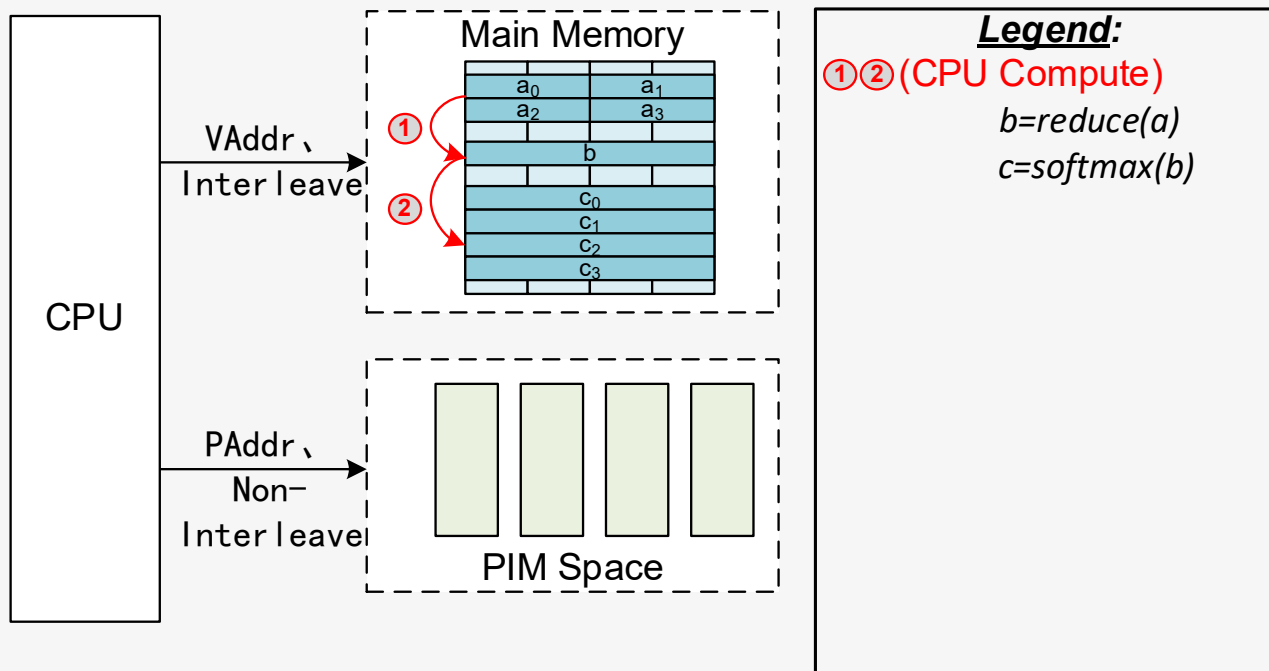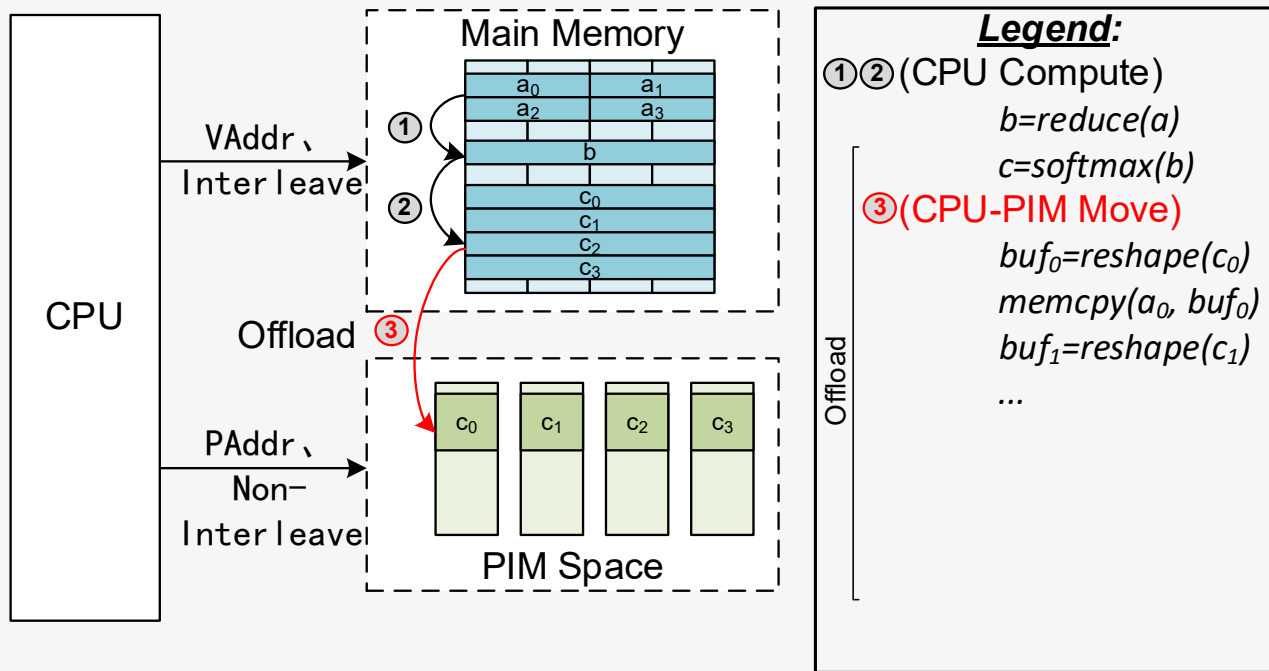## Measure 1: Isolated Memory Space



Example: MLP

# Measures of current PIM systems

⊛Measure 1: Isolated Memory Space **(Resulting in Data Transfer)**



Example: MLP

# Measures of current PIM systems

Measure 1: Isolated Memory Space **(Resulting in Data Transfer)**
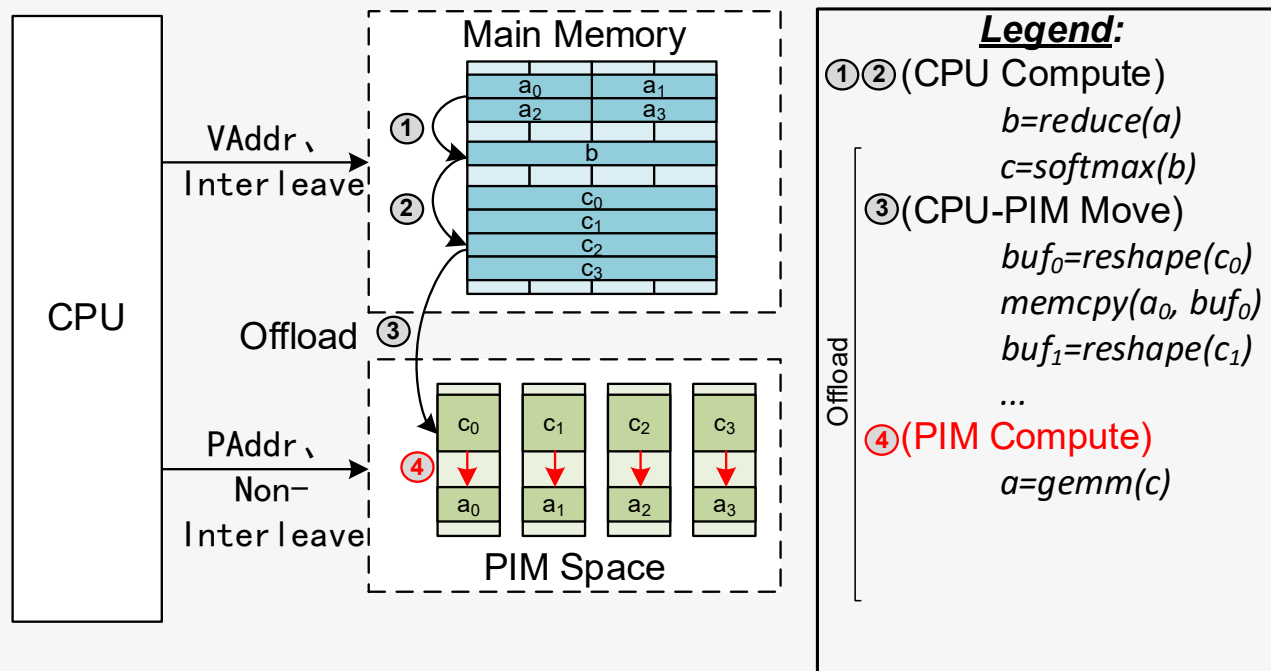


Example: MLP

◉Measure 1: Isolated Memory Space **(Resulting in Data Transfer)**



Example: MLP

# Measures of current PIM systems

Measure 1: Isolated Memory Space **(Resulting in Data Transfer)**



**Legend:**
①② (CPU Compute)
$\quad b=reduce(a)$
$\quad c=softmax(b)$
③ (CPU-PIM Move)
$\quad buf_0=reshape(c_0)$
$\quad memcpy(a_0, buf_0)$
$\quad buf_1=reshape(c_1)$
$\quad ...$
④ (PIM Compute)
$\quad a=gemm(c)$
⑤ (PIM-CPU Move)
$\quad buf_0=reshape(c_0)$
$\quad memcpy(a_0, buf_0)$
$\quad ...$

Additional Memory copy

## Measure 1: Isolated Memory Space (Resulting in Data Transfer)



**Legend:**

①②(CPU Compute)
$b=reduce(a)$
$c=softmax(b)$

③(CPU-PIM Move)
$buf_0=reshape(c_0)$
$memcpy(a_0, buf_0)$
$buf_1=reshape(c_1)$
...

④(PIM Compute)
$a=gemm(c)$

⑤(PIM-CPU Move)
$buf_0=reshape(c_0)$
$memcpy(a_0, buf_0)$
...

Data Re-layout
(Because of different data layout)

Memcpy

Additional Memory copy and data re-layout

Re-layout with Software

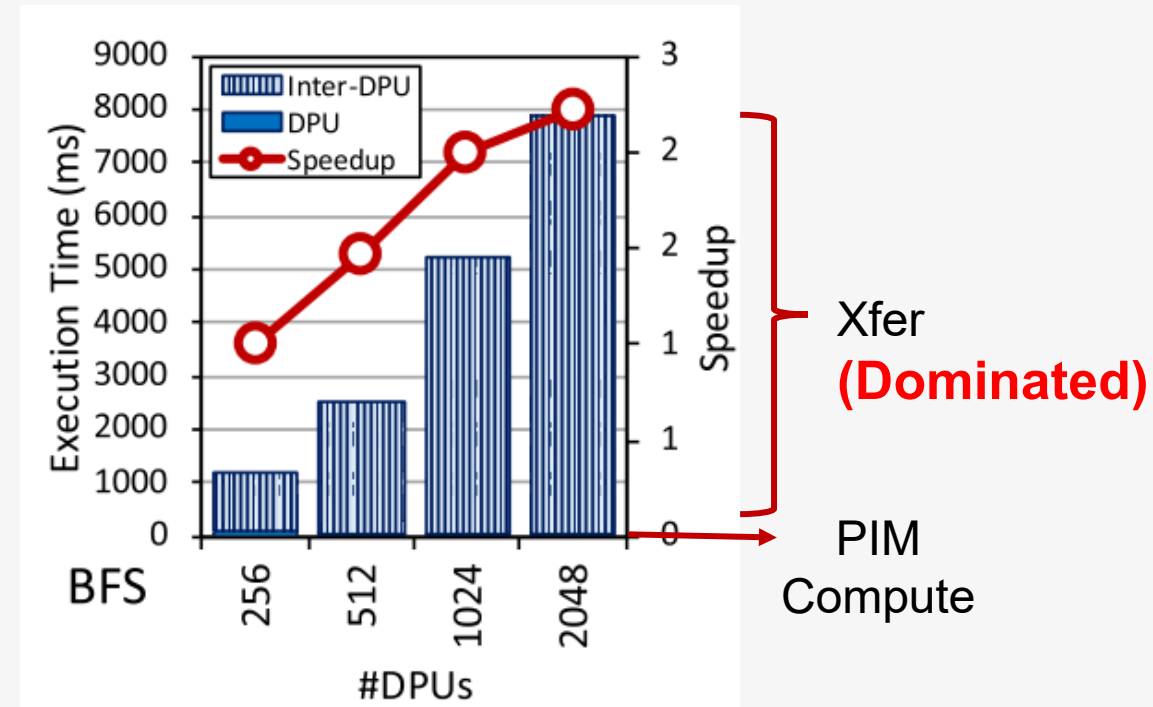Measure 1: Isolated Memory Space (Resulting in Data Transfer)



Additional Memory copy and data re-layout

Gómez-Luna et.al. 2022

# Measures of current PIM systems

◉ Measure 2: Switch off channel & rank interleaving

- Reduce re-layout overhead, but damage CPU bandwidth

- The data re-layout overhead still accounts for 70% of data transfer time
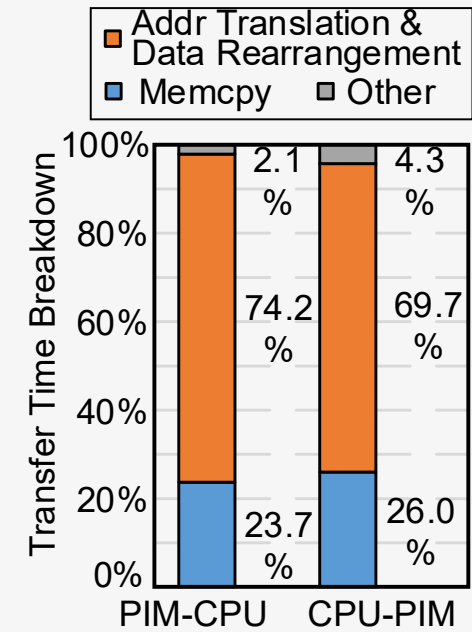
Save Computations

Ch on, Ba on ⟶ Ch off, Ba on

$PAddr| = Ch \ll m \mid Ba \ll n$       $PAddr| = Ba \ll n$

| C0B0 | C0B1 | C1B0 | C1B1 |
|------|------|------|------|
| 0x00 | 0x01 | 0x02 | 0x03 |
| 0x04 | 0x05 | 0x06 | 0x07 |
| 0x08 | 0x09 | 0x0A | 0x0B |
| 0x0C | 0x0D | 0x0E | 0x0F |

| C0B0 | C0B1 | C1B0 | C1B1 |
|------|------|------|------|
| 0x00 | 0x01 | 0x08 | 0x09 |
| 0x02 | 0x03 | 0x0A | 0x0B |
| 0x04 | 0x05 | 0x0C | 0x0D |
| 0x06 | 0x07 | 0x0E | 0x0F |

Re-layout with Software

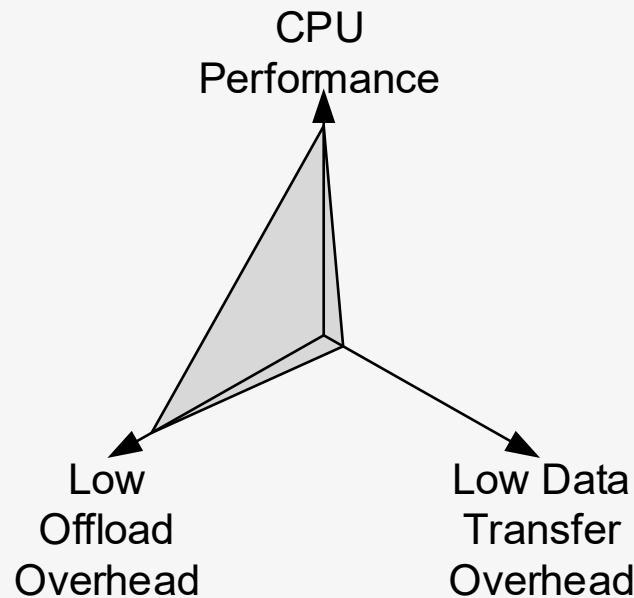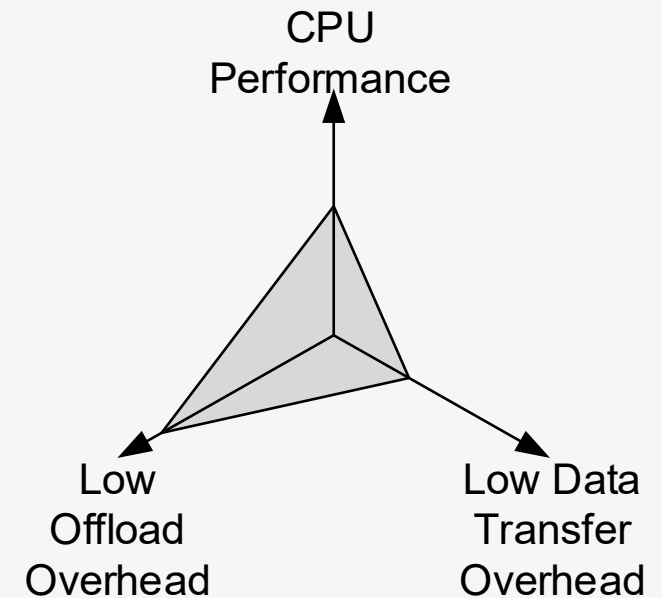| C0B0 | C0B1 | C1B0 | C1B1 |
|------|------|------|------|
| 0x00 | 0x04 | 0x08 | 0x0C |
| 0x01 | 0x05 | 0x09 | 0x0D |
| 0x02 | 0x06 | 0x0A | 0x0E |
| 0x03 | 0x07 | 0x0B | 0x0F |

UPMEM Data Transfer Breakdown

Memory spaces with single data layout.



(Naïve)
Uniform Layout
Fine-grained Offload

Isolated Memory Space
Software Re-layout

MetaPNM, PIM-HBM

Isolated Memory Space
Switch off Interleaving

AxDIMM, UPMEM

# Summary: An Impossible Triangle

◉ Memory spaces with ~~single data layout.~~



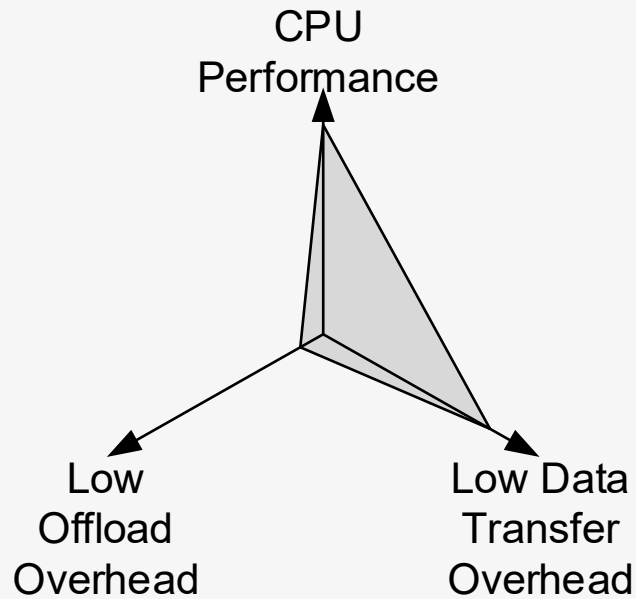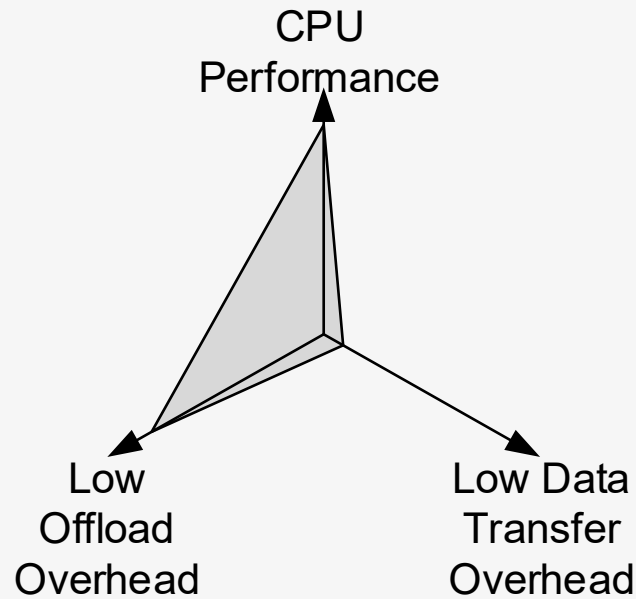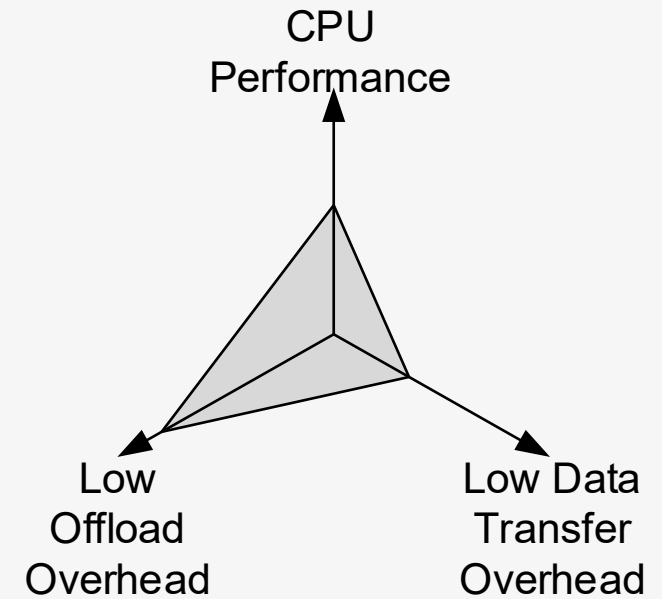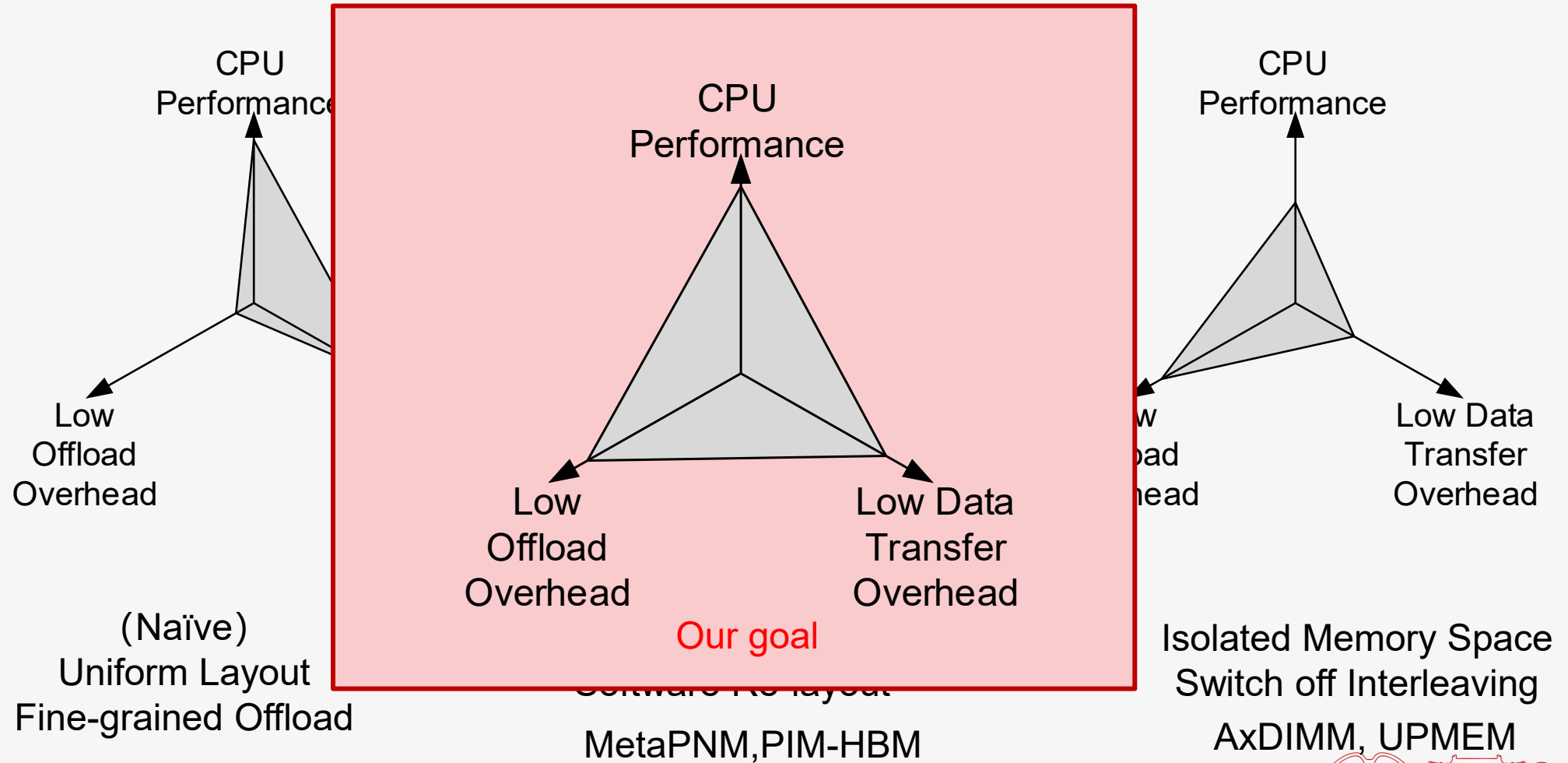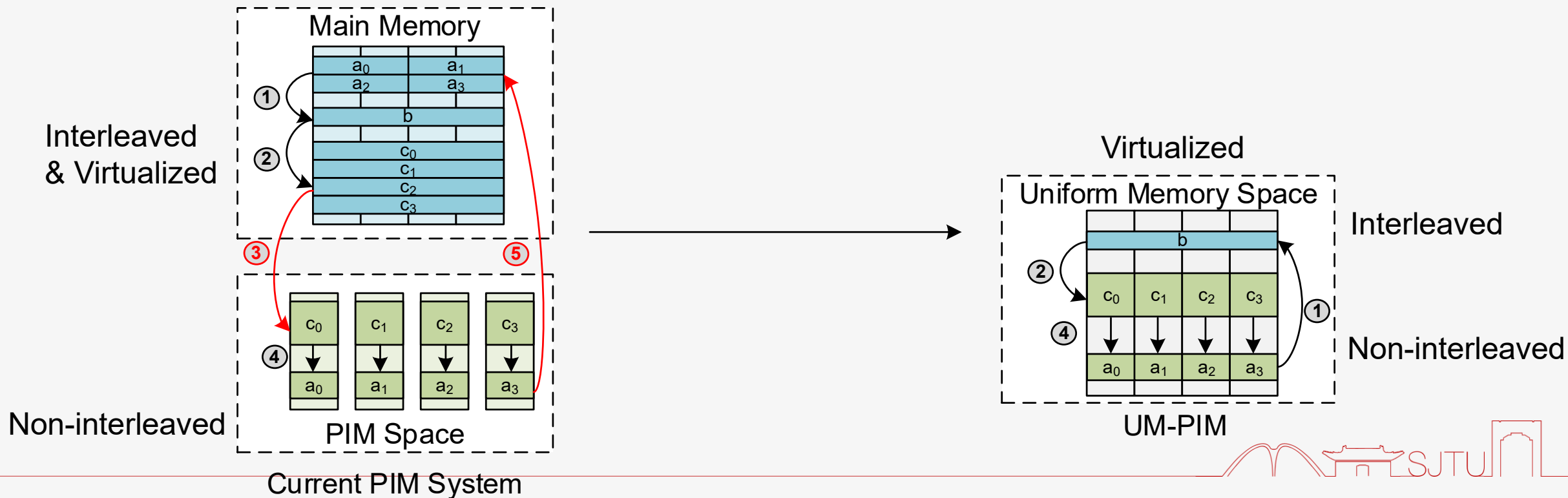|   |   |   |
|---|---|---|
| CPU Performance | CPU Performance | CPU Performance |
| Low Offload Overhead — Low Data Transfer Overhead | Low Offload Overhead — Low Data Transfer Overhead | Low Offload Overhead — Low Data Transfer Overhead |
| (Naïve) Uniform Layout Fine-grained Offload | Isolated Memory Space Software Re-layout MetaPNM,PIM-HBM | Isolated Memory Space Switch off Interleaving AxDIMM, UPMEM |

# Summary: An Impossible Triangle

◉ Memory spaces with ~~single data layout.~~



CPU
Performance

Low
Offload
Overhead

Low Data
Transfer
Overhead

(Naïve)
Uniform Layout
Fine-grained Offload

CPU
Performance

Low
Offload
Overhead

Low Data
Transfer
Overhead

Our goal

Software Re-layout

MetaPNM,PIM-HBM

CPU
Performance

Low
Offload
Overhead

Low Data
Transfer
Overhead

Isolated Memory Space
Switch off Interleaving

AxDIMM, UPMEM

Key insight: A **Uniform Memory Space** with CPU & PIM Pages (Different layout) co-existing in the space
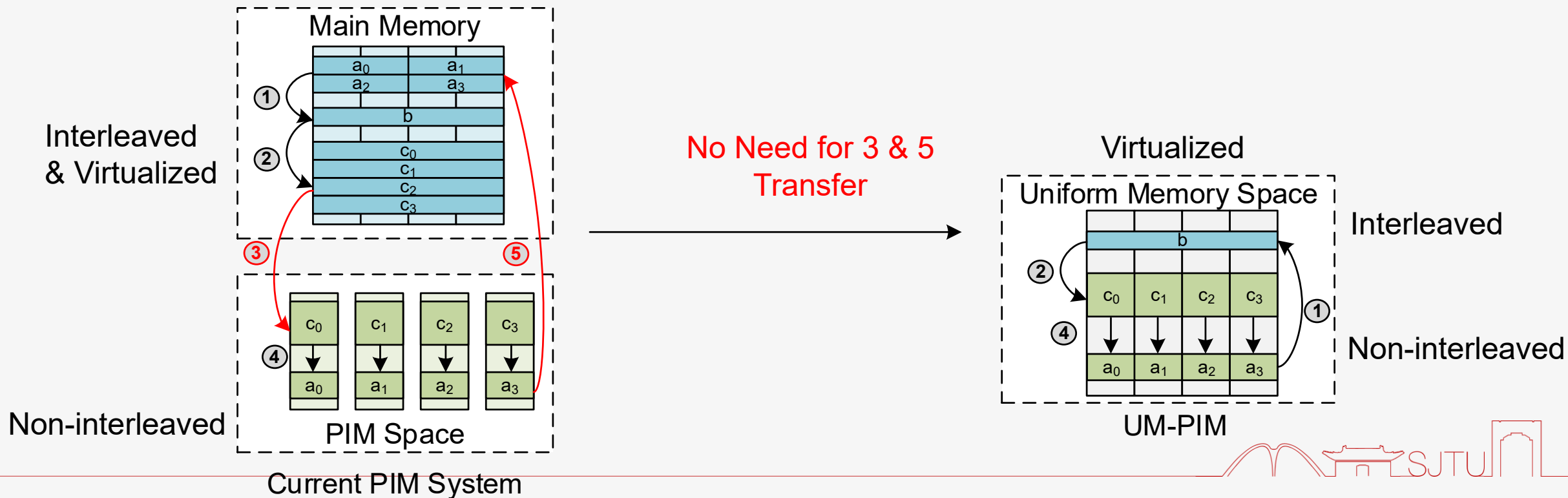
- No Data Transfer Overhead

- Satisfy CPU and PIM's data layout

- Key insight: A **Uniform Memory Space** with CPU & PIM Pages (Different layout) co-existing in the space
  - No Data Transfer Overhead
  - Satisfy CPU and PIM's data layout
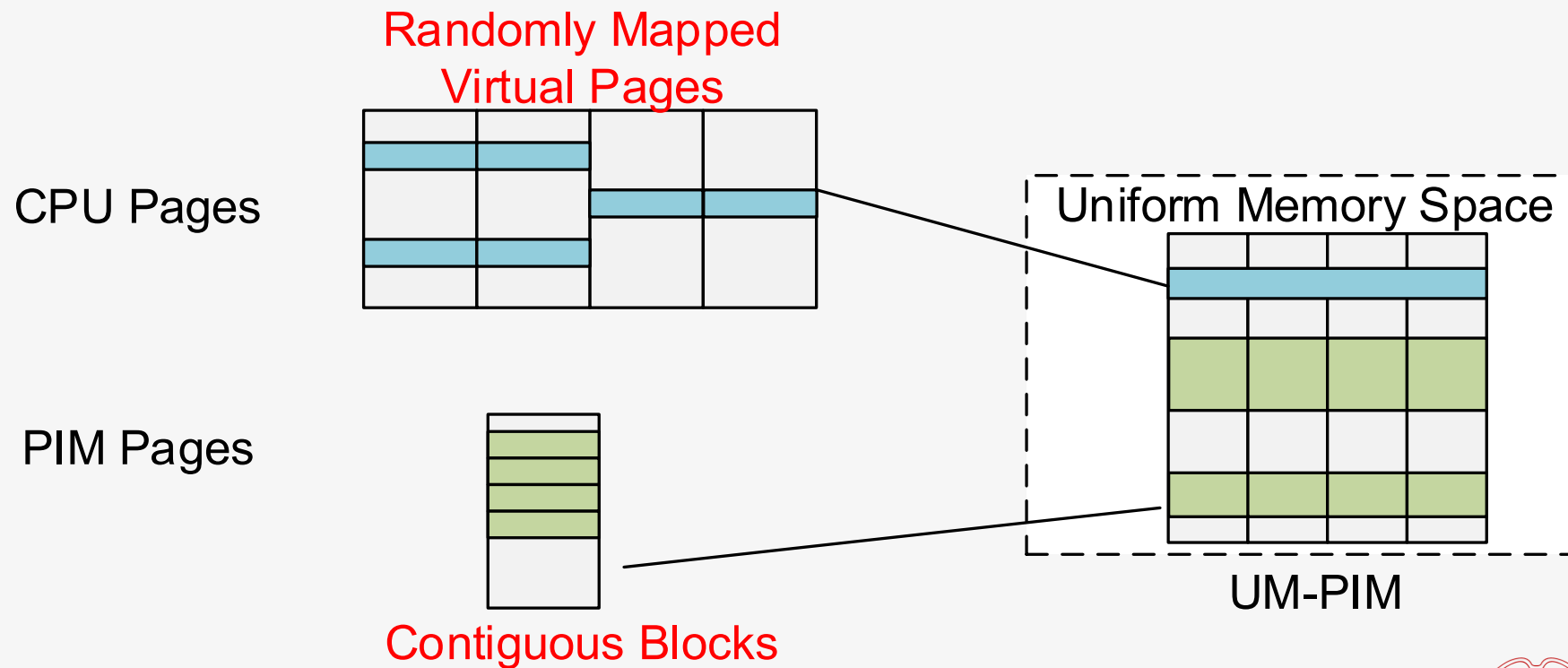
# UM-PIM Motivation

◉ Key insight: A **Uniform Memory Space**

◉ Challenges:

- **Co-existence (PIM's contiguous data block, CPU's virtual memory)**



Randomly Mapped Virtual Pages

CPU Pages

PIM Pages

Contiguous Blocks
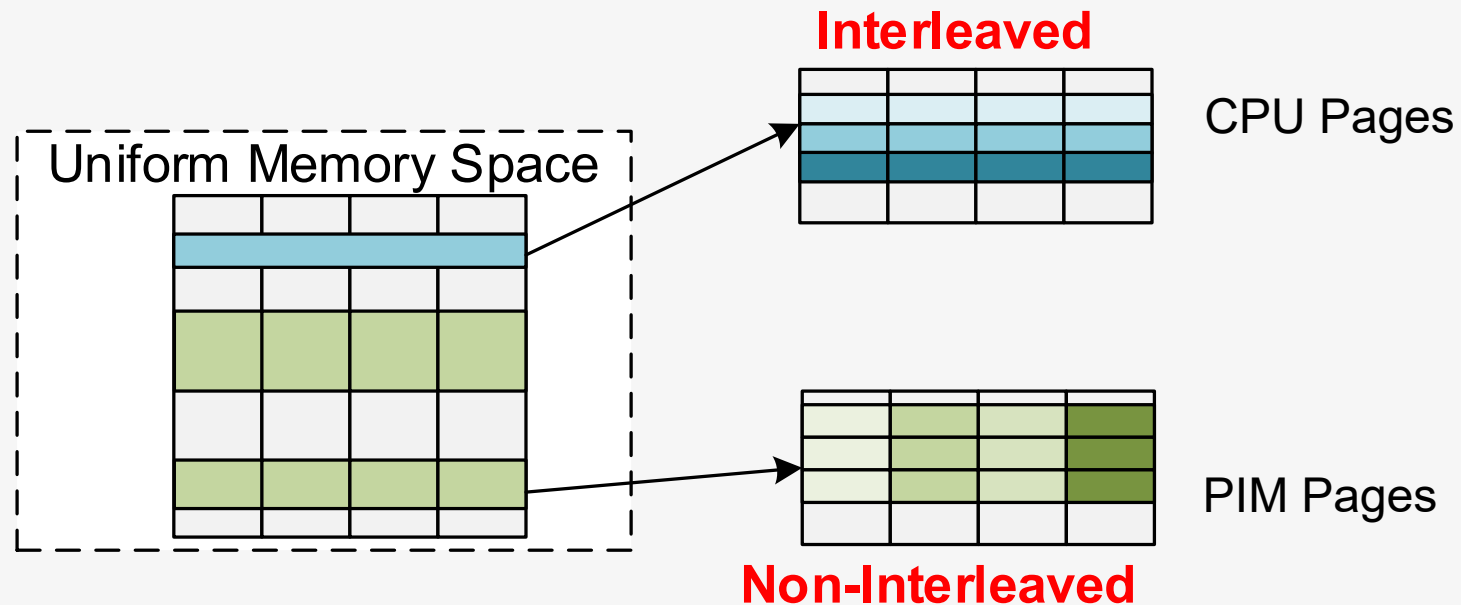
Uniform Memory Space

UM-PIM

# UM-PIM Motivation

◉ Key insight: A **Uniform Memory Space**

◉ Challenges:

- Co-existence (PIM's contiguous data block, CPU's virtual memory)
- **Two different data layout**

# UM-PIM Motivation

- Key insight: A **Uniform Memory Space**

- Challenges:

  - Co-existence (PIM's contiguous data block, CPU's virtual memory)

  - Two different data layout
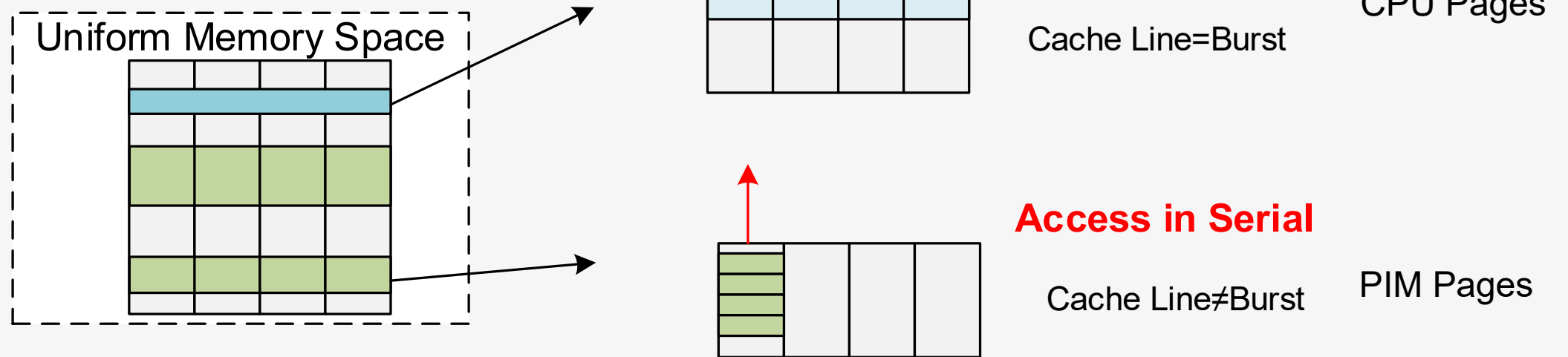
  - **Accelerate CPU accessing PIM Page**
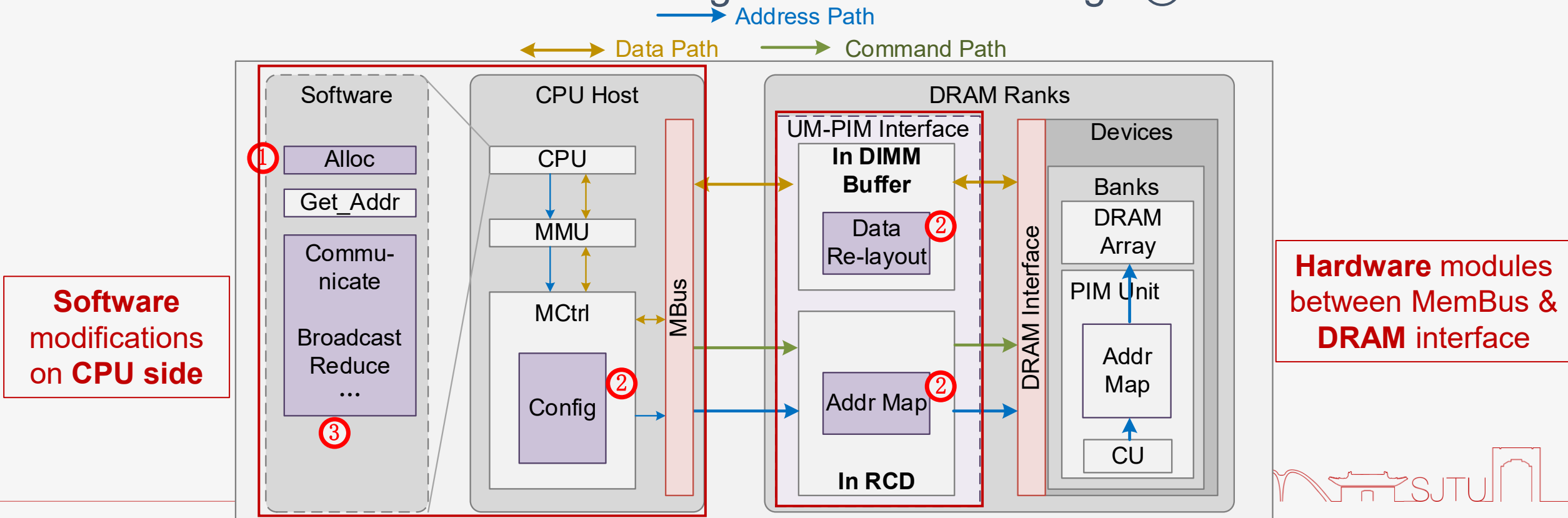
Outline

Background: Process-in-memory and Memory Interleaving

UM-PIM: DRAM-based PIM with Uniform & Shared Memory Space

Evaluation

# UM-PIM Overview

- A. Memory Management: Chunk-based management of PIM Pages ①

- B. Data Layout: Address Mapping and Data Re-layout ②

- C. Accelerate CPU Access PIM Page: HW-SW Co-design ③

# A. Memory Management – CPU Malloc ①

Challenge 1: Co-existence of virtualized
CPU Page & contiguous PIM Page

◉ CPU Pages: keep current management strategy

◉ PIM Pages: Chunk-based Management

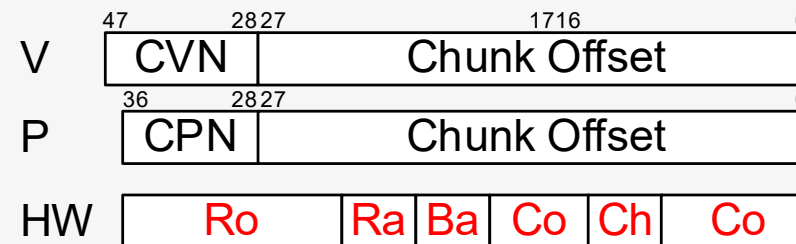Challenge 1: Co-existence of virtualized CPU Page & contiguous PIM Page
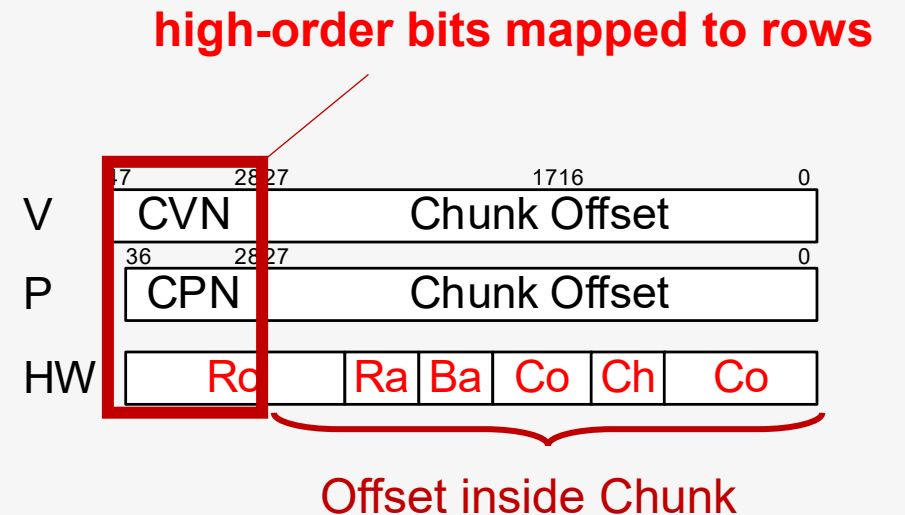
## PIM Pages: Chunk-based Management

- Let high-order bits mapped to rows and Allocate a **Huge Page** (Chunk)

A **huge page** 0x2..00000~0x3..00000

| | 47 28 | 27 | 17 16 | 0 |
|---|---|---|---|---|
| V | CVN | | Chunk Offset | |

| | 36 28 | 27 | | 0 |
|---|---|---|---|---|
| P | CPN | | Chunk Offset | |

| | | | | | | |
|---|---|---|---|---|---|---|
| HW | Ro | Ra | Ba | Co | Ch | Co |

Switch on interleaving

**PIM Pages: Chunk-based Management**

Challenge 1: Co-existence of virtualized CPU Page & contiguous PIM Page

- Let high-order bits mapped to rows and Allocate a **Huge Page** (Chunk)

A **huge page** 0x2..00000~0x3..00000

**high-order bits mapped to rows**

|   | 47      28 | 27          1716          0 |
|---|---|---|
| V | CVN | Chunk Offset |
| P | 36      28 | 27                         0 |
|   | CPN | Chunk Offset |
| HW | Ro | Ra \| Ba \| Co \| Ch \| Co |

Offset inside Chunk

# A. Memory Management – CPU Malloc ①

Challenge 1: Co-existence of virtualized CPU Page & contiguous PIM Page
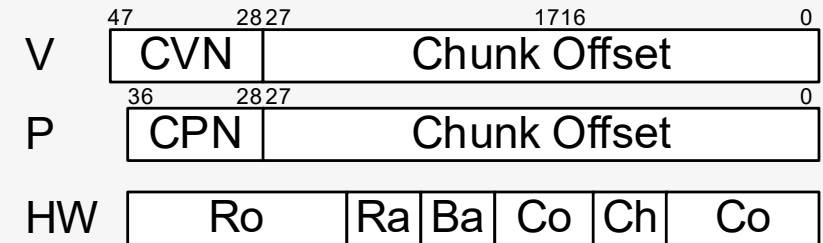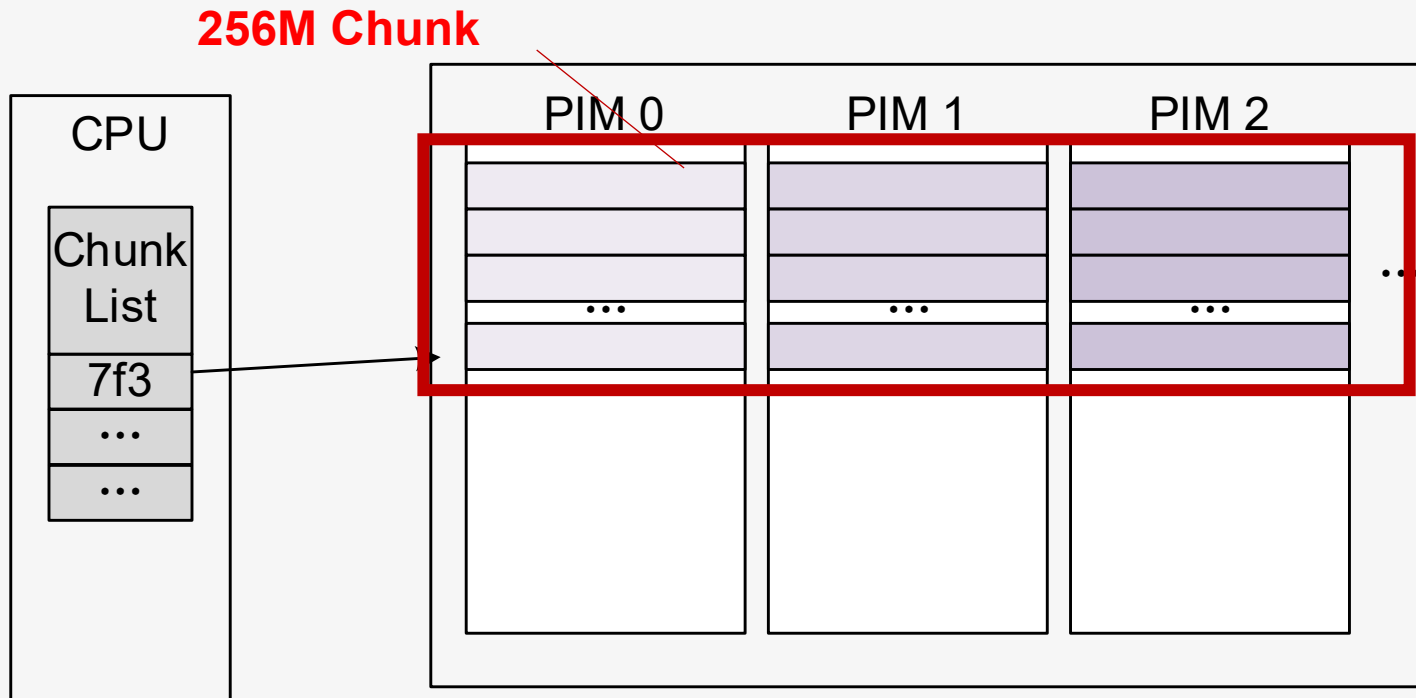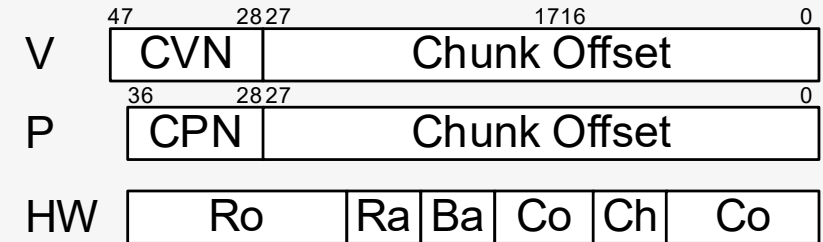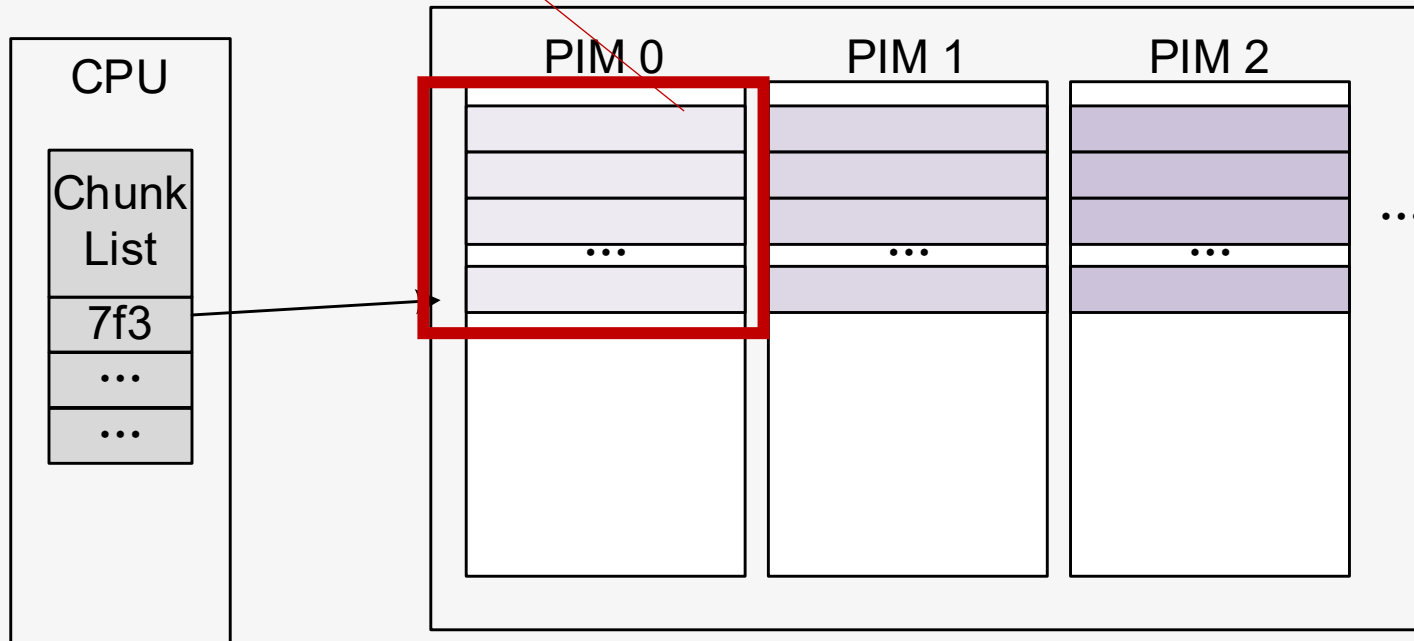
⊕ PIM Pages: Chunk-based Management

- Let high-order bits mapped to rows and Allocate a Huge Page (Chunk):
  - The Chunk Evenly distributed on every PIMs

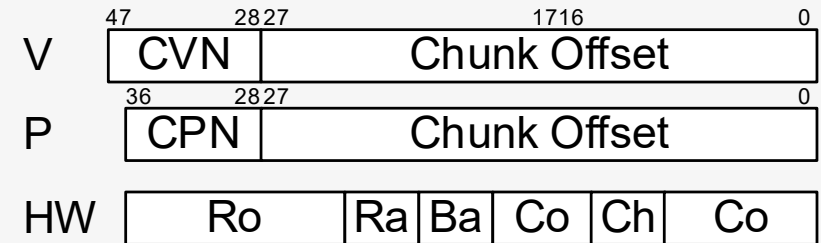A huge page 0x2..00000~0x3..00000



**256M Chunk**

Challenge 1: Co-existence of virtualized CPU Page & contiguous PIM Page

## PIM Pages: Chunk-based Management

- Let high-order bits mapped to rows and Allocate a Huge Page (Chunk):
  - The Chunk Evenly distributed on every PIMs

A huge page 0x2..00000~0x3..00000
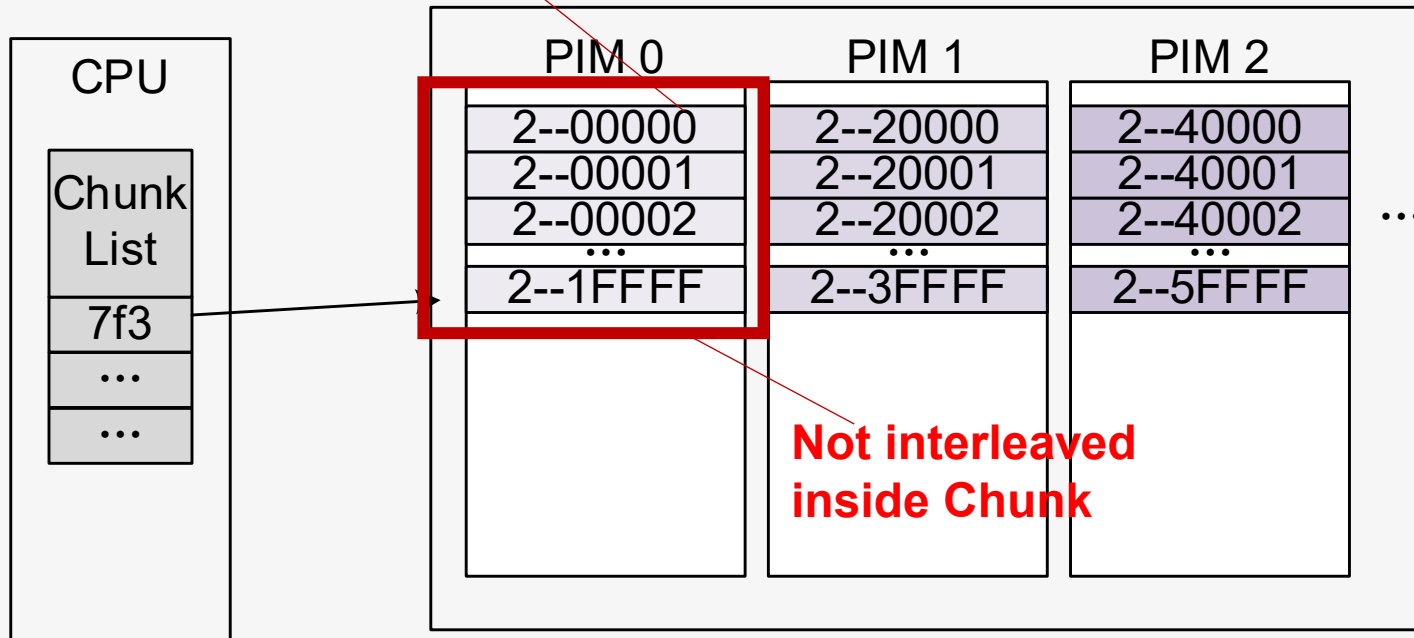
**256M Chunk = 128 k PIM Page on each PIM**

Challenge 1: Co-existence of virtualized CPU Page & contiguous PIM Page

## PIM Pages: Chunk-based Management

- Let high-order bits mapped to rows and Allocate a Huge Page (Chunk):
  - The Chunk Evenly distributed on every PIMs

A huge page 0x2..00000~0x3..00000
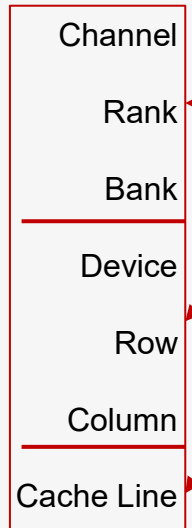
256M Chunk = 128 k PIM Page on each PIM



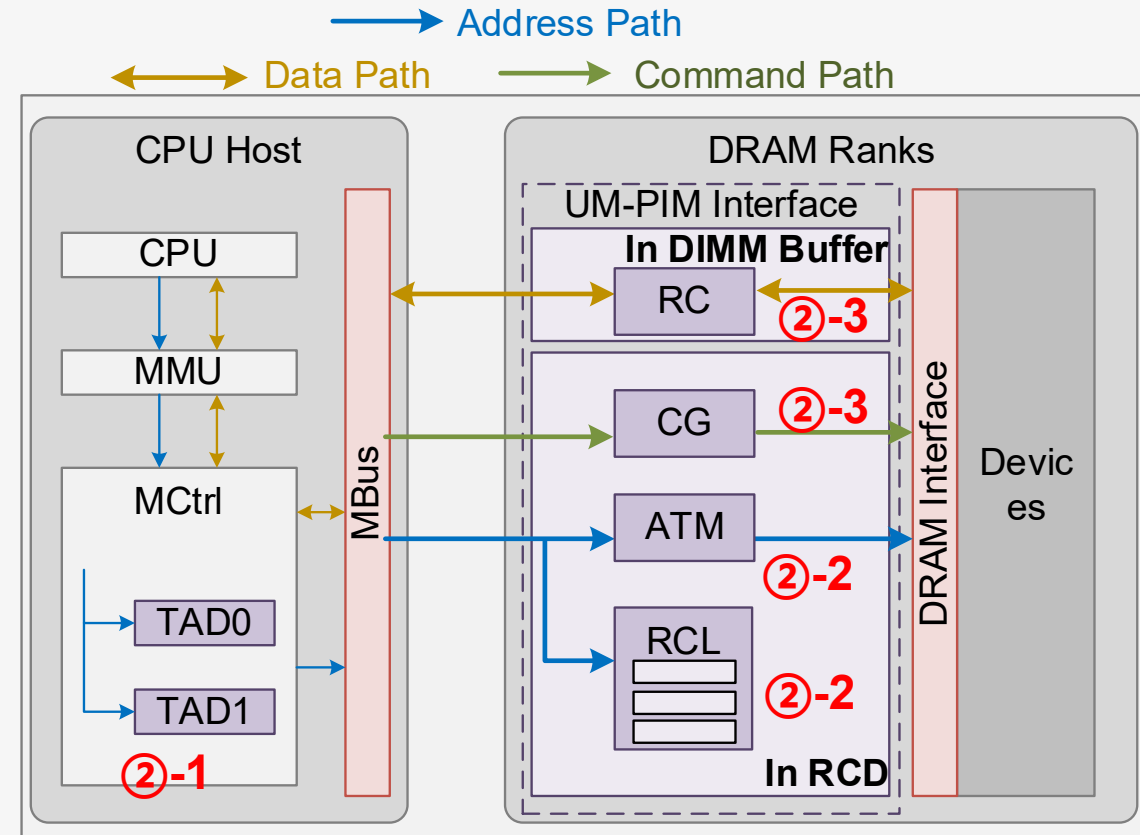**Not interleaved inside Chunk**

# B Data Layout

֍ Integrate Hardware Modules between DRAM Interface & Memory Bus, on DRAM side

֍ **3 Levels**:

- Cache Line: Dynamic Address Mapping
  - ②-1 Above Bank: Config Mem Ctrl
  - ②-2 Below Bank: DRAM-side HW (ATM &RCL)
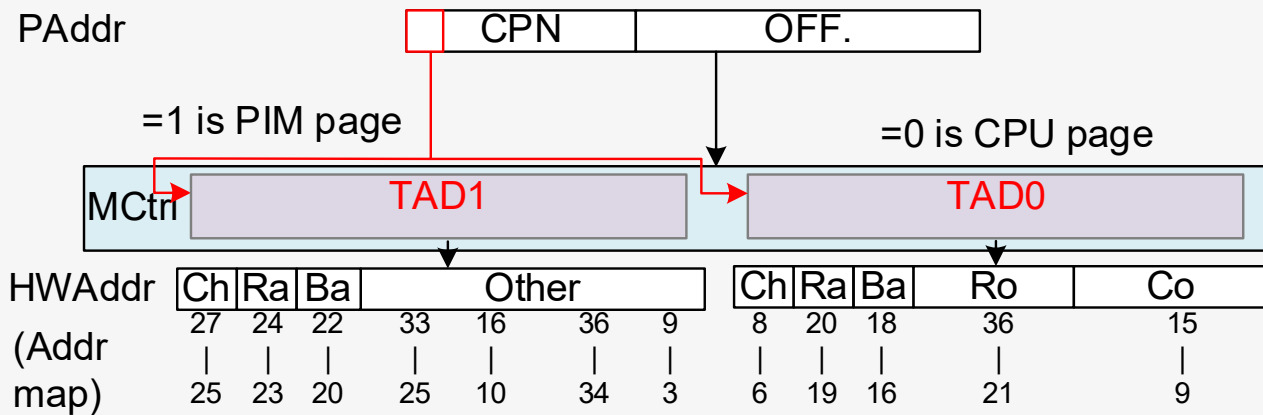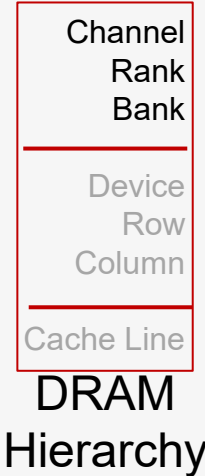- ②-3 Inside Cache Line: Data re-layout
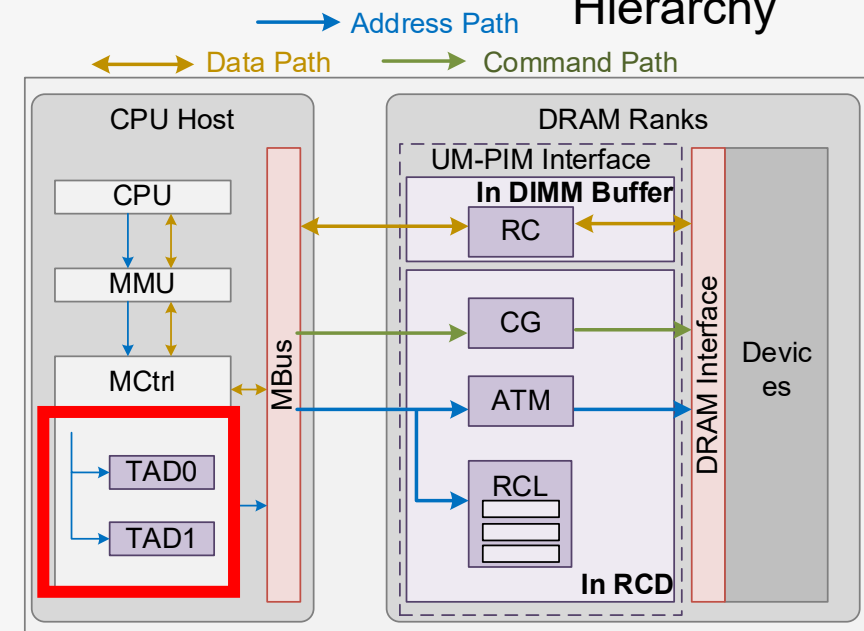  - DRAM-side HW (RC&CG)



DRAM Hierarchy

# B Data Layout — Address Mapping ②-1

- **(Above Bank)** Add a Address Alias (TAD 1) in MemCtrl
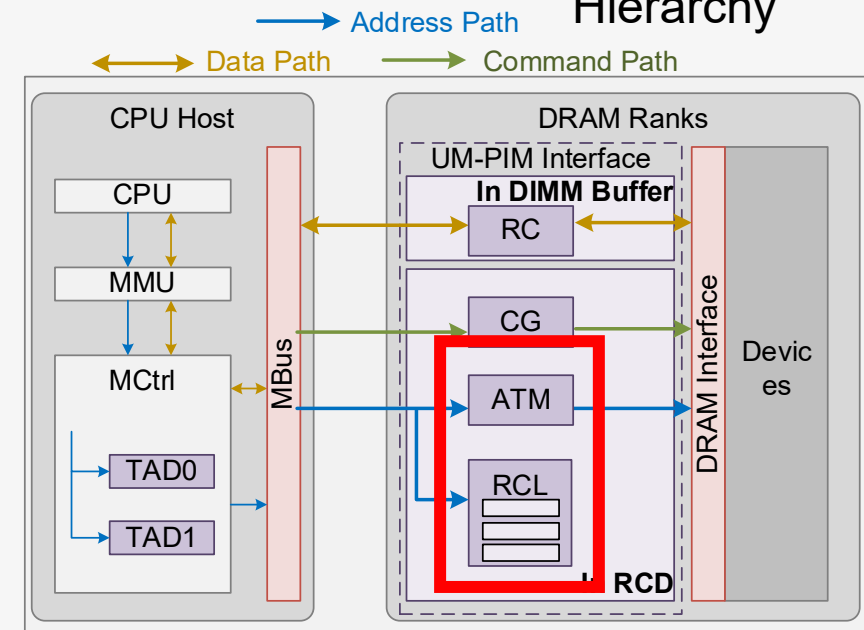  - TAD1 have different address mapping from the origin TAD0
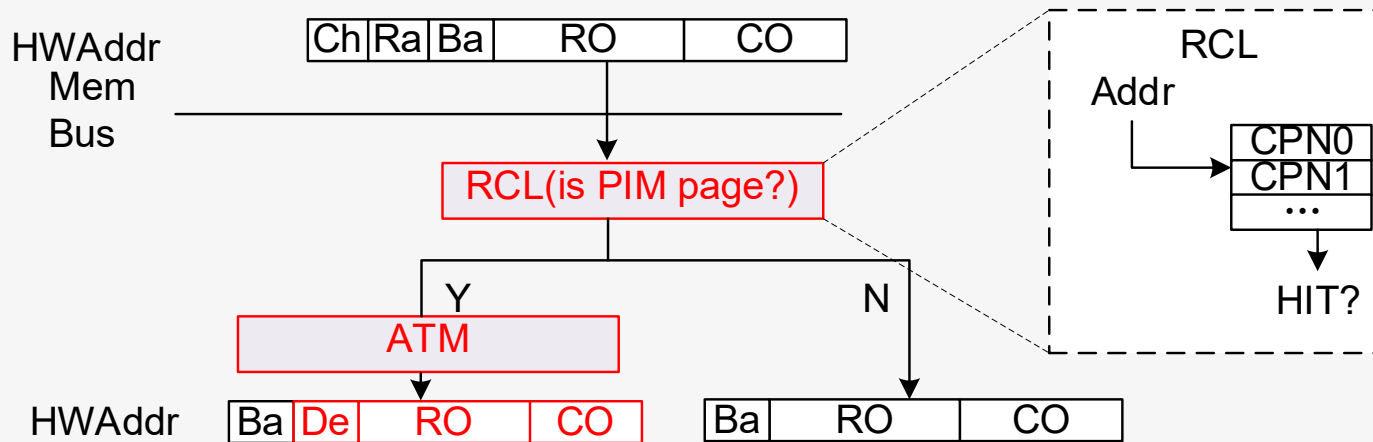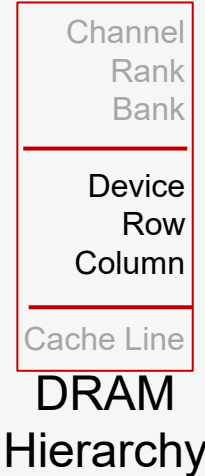


TAD0 & TAD1 have different address mapping

**(Below Bank)** Hardware module on DRAM side

- RCL: Check whether the Address belongs to PIM Page
- ATM: If is PIM address, map the address bits to Device Column and Row

**(Below Bank)** Hardware module on DRAM side

- RCL: Check whether the Address belongs to PIM Page

- ATM: If is PIM address, map the address bits to Device Column and Row

Address: 0xB0020001



DRAM Hierarchy

**(Below Bank)** Hardware module on DRAM side

- RCL: Check whether the Address belongs to PIM Page

- ATM: If is PIM address, map the address bits to Device Column and Row



Address: 0xB0020001

CPN=0xB

Channel
Rank
Bank

Device
Row
Column

Cache Line

DRAM
Hierarchy

# B Data Layout — Address Mapping ②-2

**(Below Bank)** Hardware module on DRAM side

- RCL: Check whether the Address belongs to PIM Page

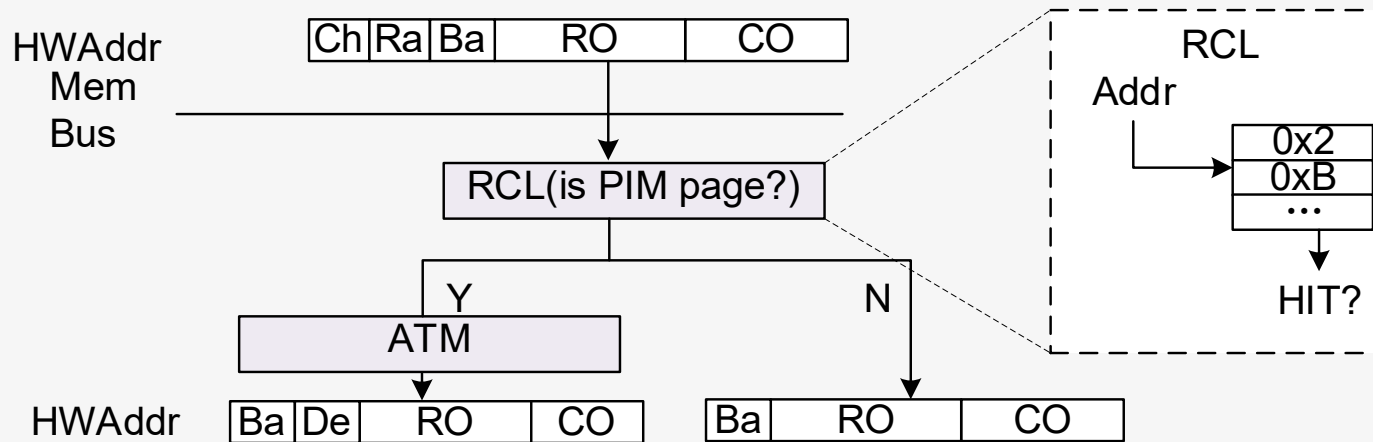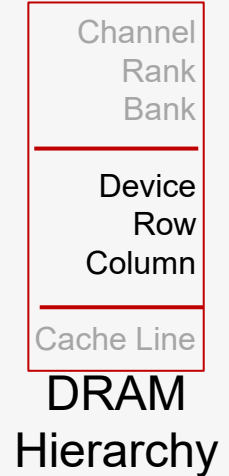- ATM: If is PIM address, map the address bits to Device Column and Row



Address: 0xB0020001

CPN=0xB

Channel
Rank
Bank

Device
Row
Column

Cache Line

DRAM Hierarchy

HWAddr
Mem
Bus

| Ch | Ra | Ba | RO | CO |

RCL(is PIM page?)

Y

ATM

| Ba | De | RO | CO |

N

| Ba | RO | CO |

RCL
Addr

0x2
0xB
...

HIT?

HIT=1, is **PIM Page**

**(Below Bank)** Hardware module on DRAM side

- RCL: Check whether the Address belongs to PIM Page
- ATM: If is PIM address, map the address bits to Device Column and Row
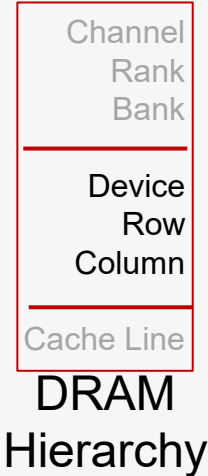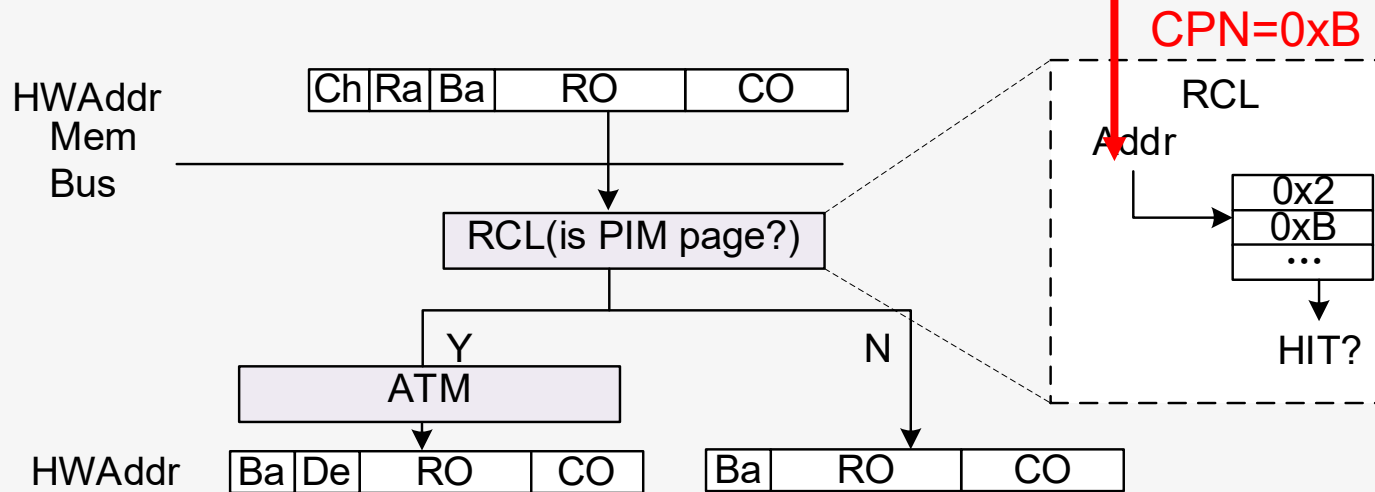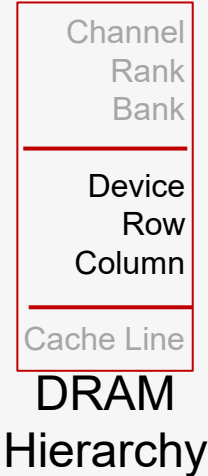


Address: 0xB0020001

CPN=0xB

HWAddr
Mem
Bus

| Ch | Ra | Ba | RO | CO |

RCL(is PIM page?)

Y          N

ATM

HWAddr | Ba | De | RO | CO |     | Ba | RO | CO |

Convert with ATM

RCL
Addr
0x2
0xB
...
HIT?

HIT=1, is **PIM Page**

Channel
Rank
Bank

Device
Row
Column
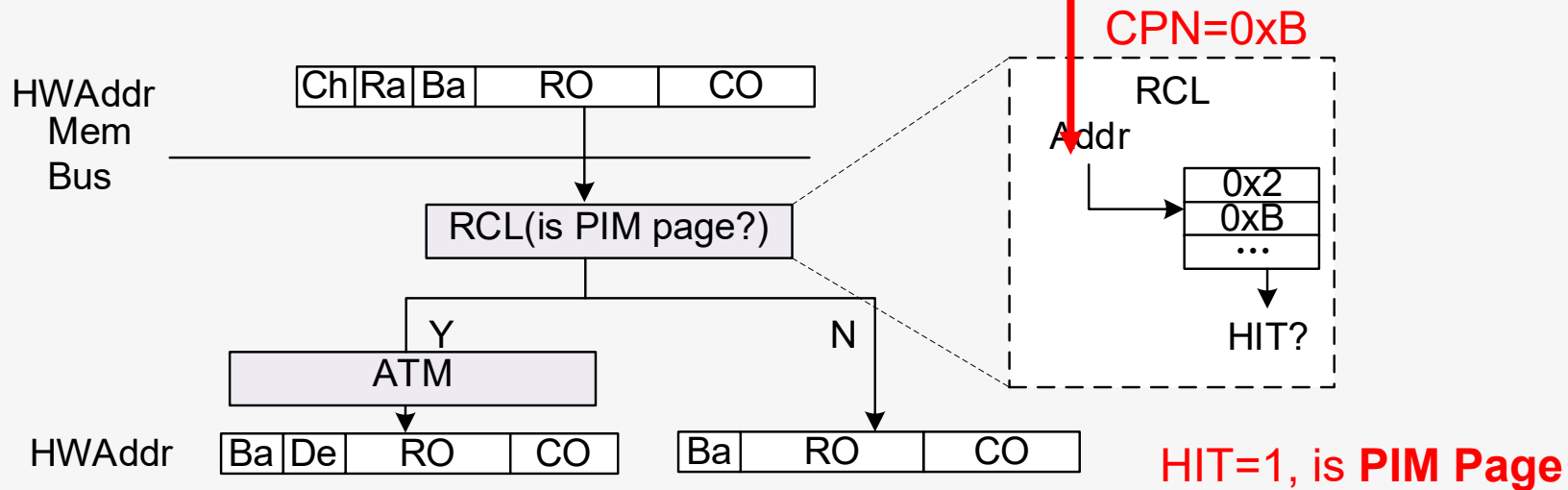
Cache Line

DRAM
Hierarchy

⊛ **(Below Bank)** Hardware module on DRAM side

- RCL: Check whether the Address belongs to PIM Page

- ATM: If is PIM address, map the address bits to Device Column and Row
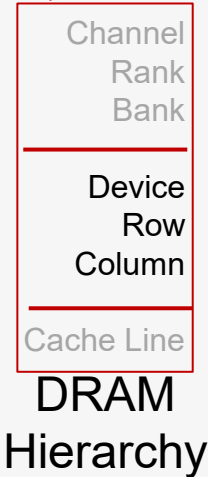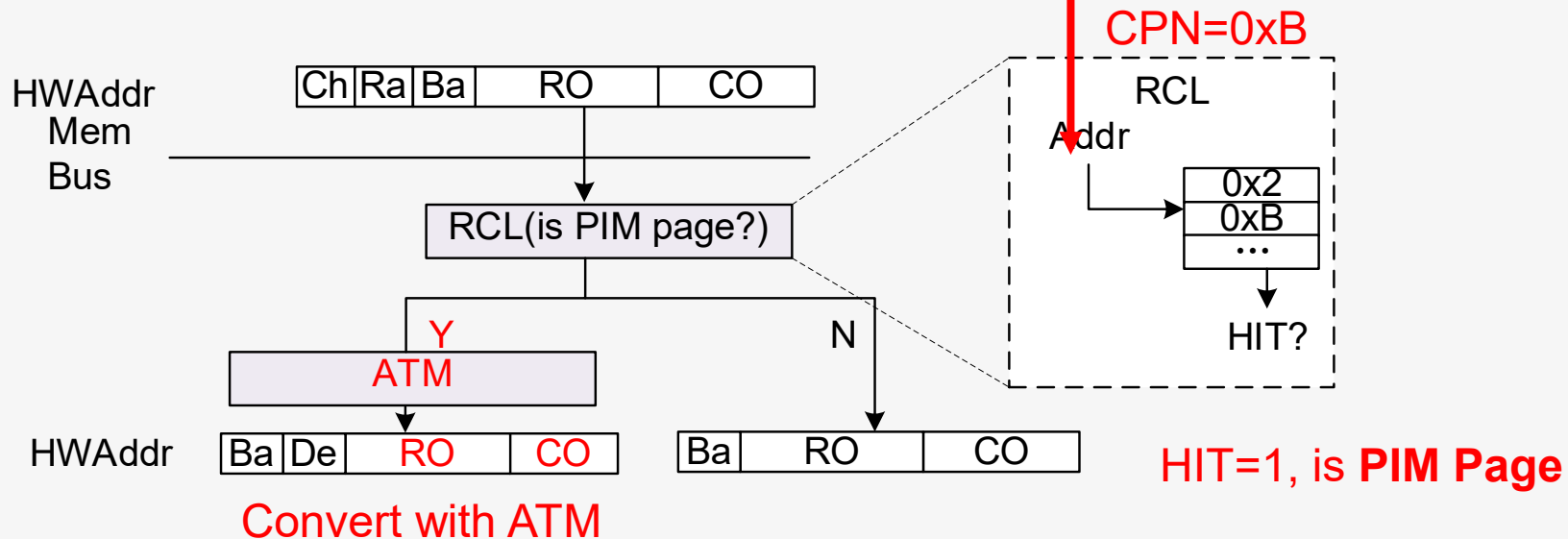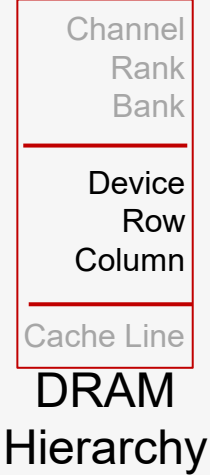


Address: 0xC0020001

CPN=0xC

HWAddr
Mem
Bus

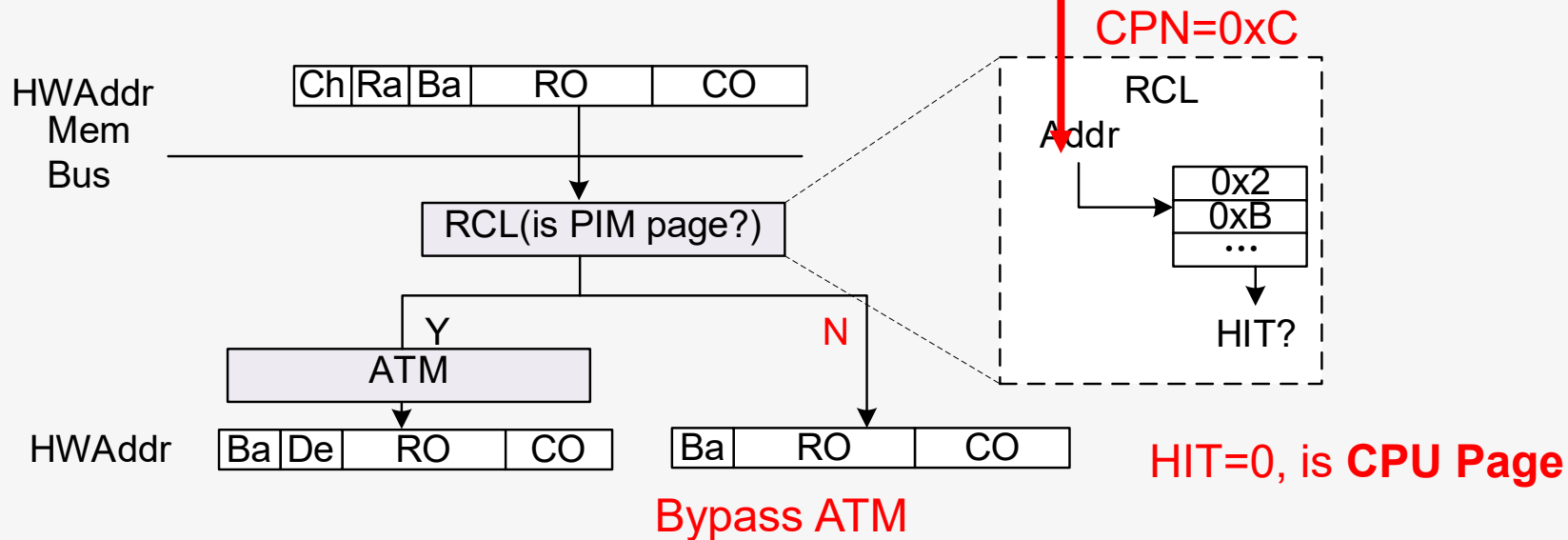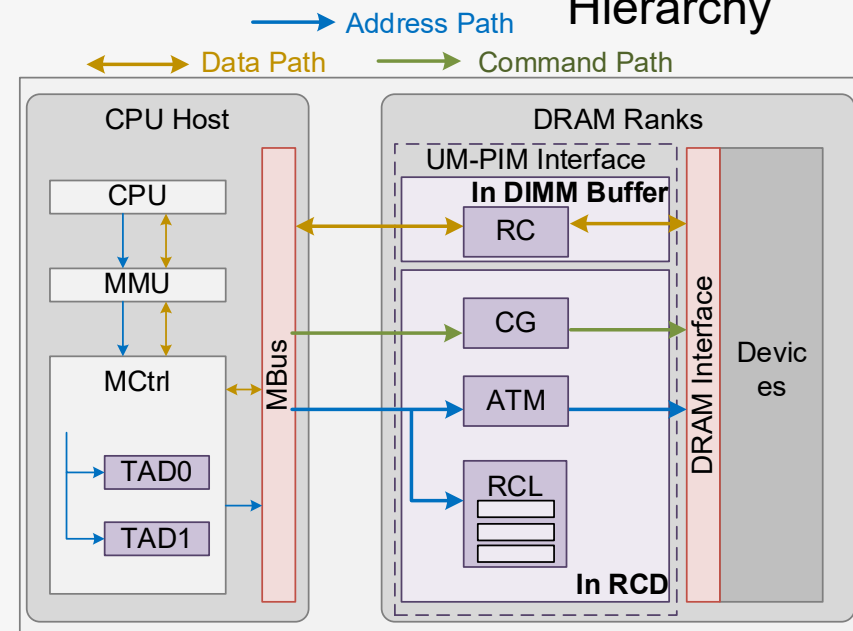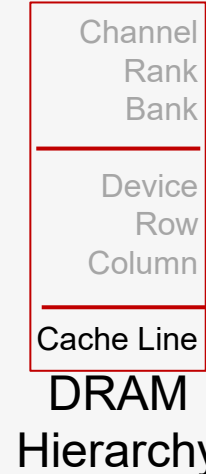| Ch | Ra | Ba | RO | CO |

RCL(is PIM page?)

Y          N

ATM

HWAddr | Ba | De | RO | CO |     | Ba | RO | CO |

Bypass ATM

RCL
Addr
0x2
0xB
...
HIT?

HIT=0, is **CPU Page**

Channel
Rank
Bank
Device
Row
Column
Cache Line

DRAM
Hierarchy

◉ **(Inside Cache Line)** Cache Line ≠ Burst problem of PIM Page

**PIM Friendly**
**Data Localization**
A Cache Line of PIM Page

**CPU Friendly**
**High Bandwidth**
**Decided by Hardware**
A DRAM Burst

| D0 | D1 | D2 | D3 |
|------|------|------|------|
| 0x00 | 0x04 | 0x08 | 0x0C |
| 0x01 | 0x05 | 0x09 | 0x0D |
| 0x02 | 0x06 | 0x0A | 0x0E |
| 0x03 | 0x07 | 0x0B | 0x0F |
| | | | |

...

D：Device

Channel
Rank
Bank

Device
Row
Column

Cache Line

DRAM Hierarchy

Address Path
Data Path    Command Path

CPU Host        DRAM Ranks
                UM-PIM Interface
CPU             In DIMM Buffer
                RC
MMU
                CG
MCtrl           ATM
TAD0
TAD1            RCL

MBus    DRAM Interface    Devices

In RCD
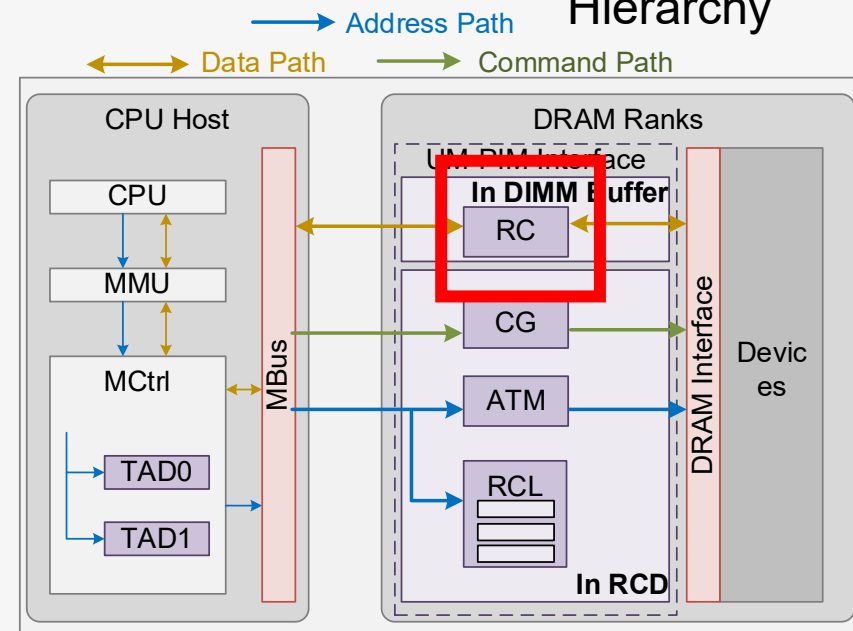
- **(Inside Cache Line)** Cache Line ≠ Burst problem of PIM Page

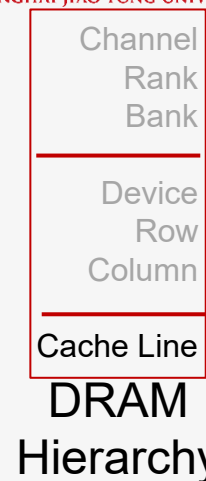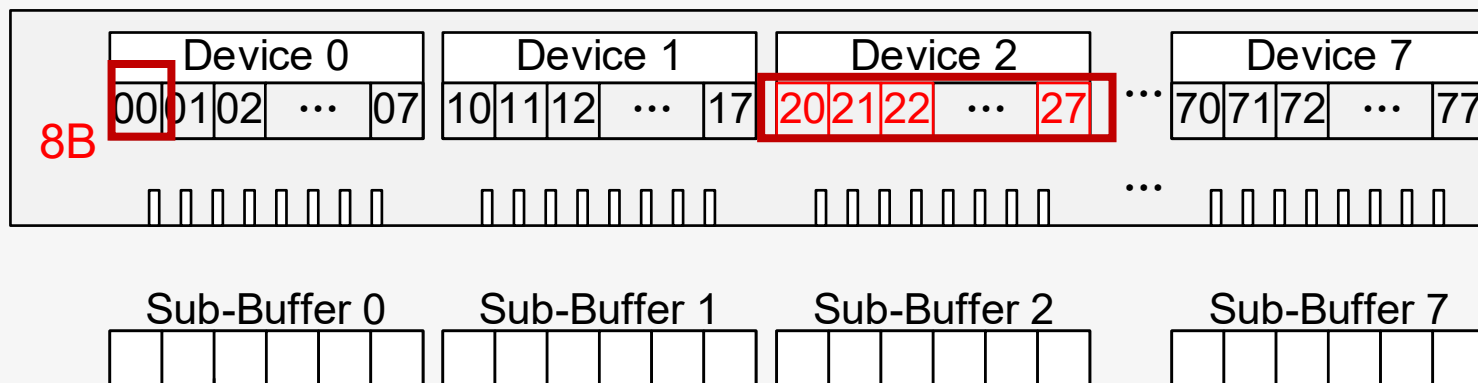- A Hardware module **RC** for data re-layout inside Cache Line



Re-layout cache (RC)

# B Data Layout — Re-layout ②-3

- **(Inside Cache Line)**

- Read a 64B Cache Line of PIM Page in Device 2

Channel
Rank
Bank

Device
Row
Column

Cache Line

DRAM
Hierarchy

One 64B Cache Line

| Device 0 | Device 1 | Device 2 | Device 7 |
|---|---|---|---|
| 00 01 02 ··· 07 | 10 11 12 ··· 17 | 20 21 22 ··· 27 | ··· 70 71 72 ··· 77 |

8B

···

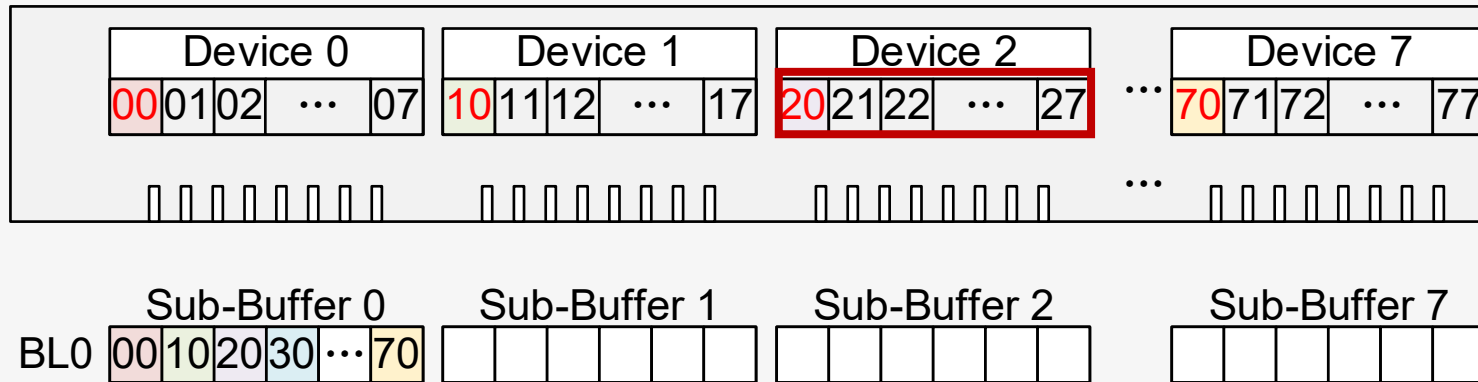| Sub-Buffer 0 | Sub-Buffer 1 | Sub-Buffer 2 | Sub-Buffer 7 |
|---|---|---|---|

- **(Inside Cache Line)**

- Read a 64B Cache Line of PIM Page in Device 2

  - Cycle 0-7: Read 8 Burst form every device, cache the Bursts in RC

Channel
Rank
Bank

Device
Row
Column

Cache Line

DRAM
Hierarchy

**(Inside Cache Line)**

Read a 64B Cache Line of PIM Page in Device 2

- Cycle 0-7: Read 8 Burst form every device, cache the Bursts in RC
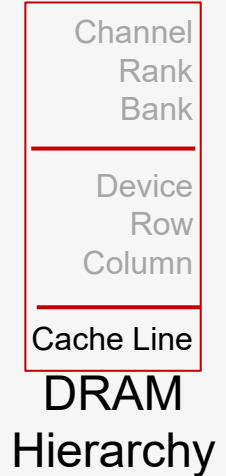
**(Inside Cache Line)**

Read a 64B Cache Line of PIM Page in Device 2

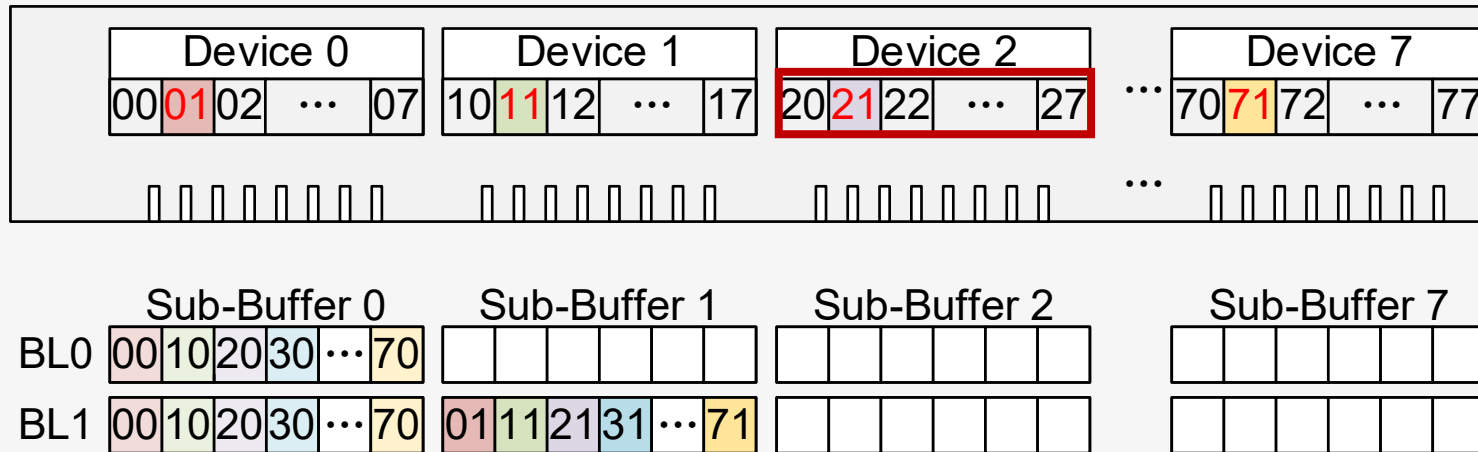- Cycle 0-7: Read 8 Burst form every device, cache the Bursts in RC

Channel
Rank
Bank

Device
Row
Column

Cache Line

DRAM
Hierarchy

**(Inside Cache Line)**

Read a 64B Cache Line of PIM Page in Device 2

- Cycle 0-7: Read 8 Burst form every device, cache the Bursts in RC

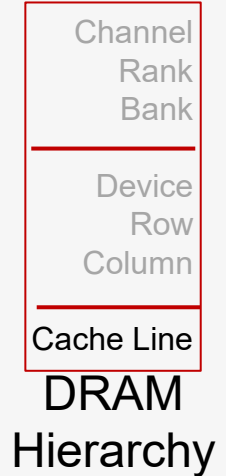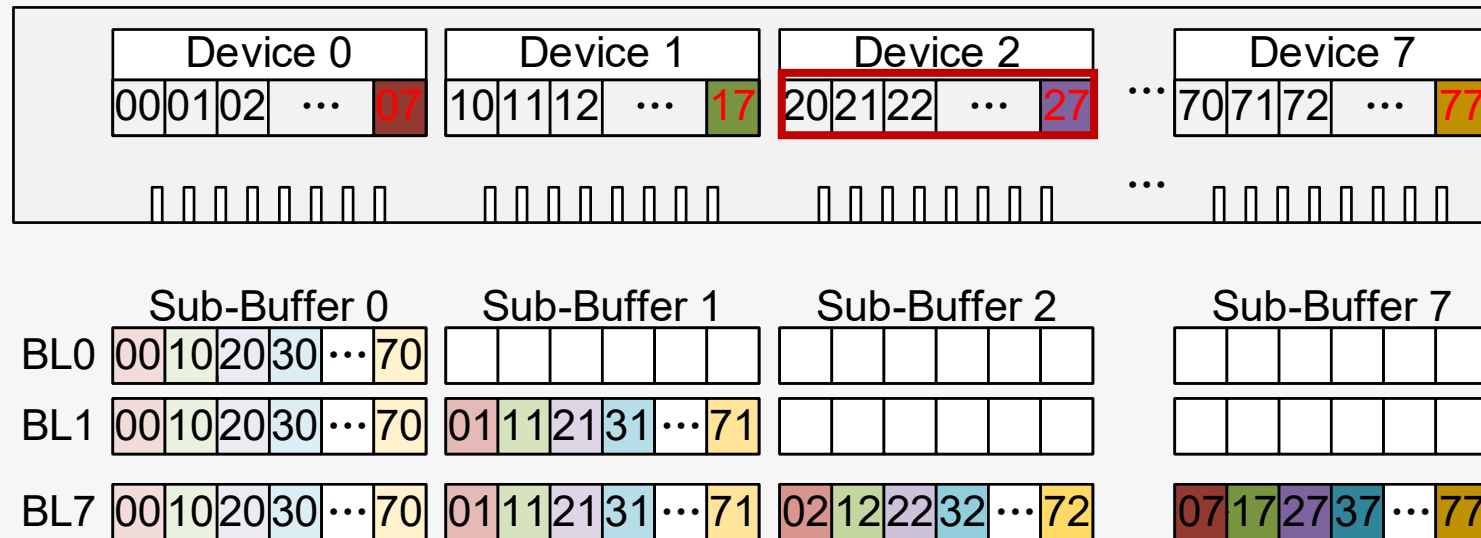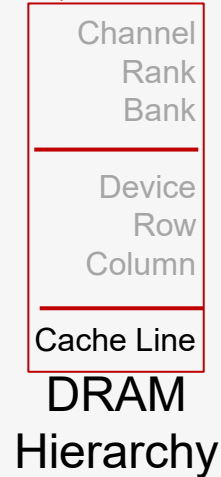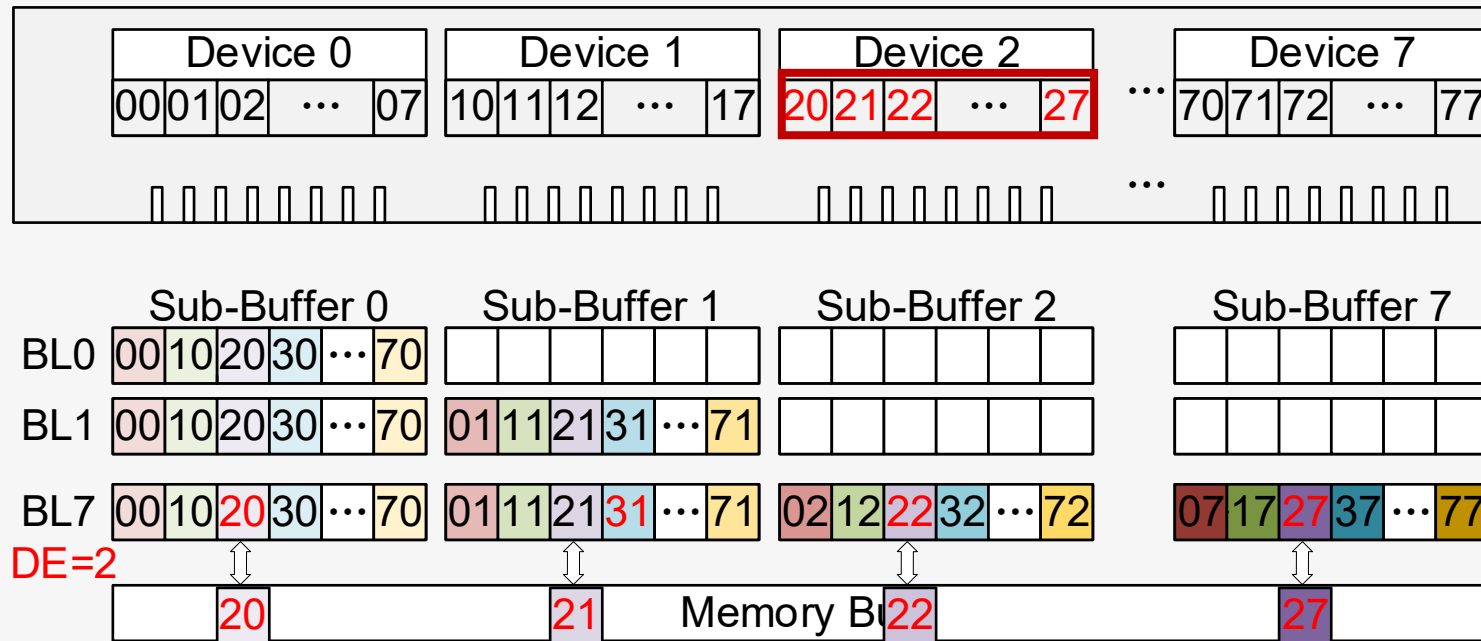- Cycle 8: Get the data from Device 2

Channel
Rank
Bank

Device
Row
Column

Cache Line

DRAM
Hierarchy

# C Accelerate CPU Access PIM Page ③

⌬ Require 8 bursts to read one cache line of PIM Page

Challenge 3: Improve CPU bandwidth when accessing PIM page

When CPU iterate over results of PIM units:

```
for de in range (#device):
    for l in range (len):
        access(de, l)
```

When CPU iterate over results of PIM units:

```
for de in range (#device):
    for l in range (len):
        access(de, l)
```

When CPU iterate over results of PIM units:

- Let iteration on devices be the inner loop

  - $dw$: device address for writing

All-Gather
(d)

PIM 0
PIM 1
PIM 2
PIM 3

```
for l in range(len/64):
 for br in range(#bank):
  for dr in range(8):
   for bw in range(#bank):
    for dw in range(8):
     memcpy(dest[bw, dw, (br*8+dr)*s+l],
            src[br, dr, l], 64)
```

Time    RC Hit Rate    Normalized Read BDW

Order of Nested Loop

**These are good iteration orders**

Outline

Background: Process-in-memory and Memory Interleaving

UM-PIM: DRAM-based PIM with Uniform & Shared Memory Space

Evaluation

# Evaluation Methodology

- Simulator & Configuration
  - CPU: GEM5 + Ramulator2;        PIM units: UPMEM DPU @ 500MHz
  - DDR4-2400, 8(Channel)x 8(Ranks)

- Benchmarks:
  - CPU Workloads: from SPEC CPU 2006
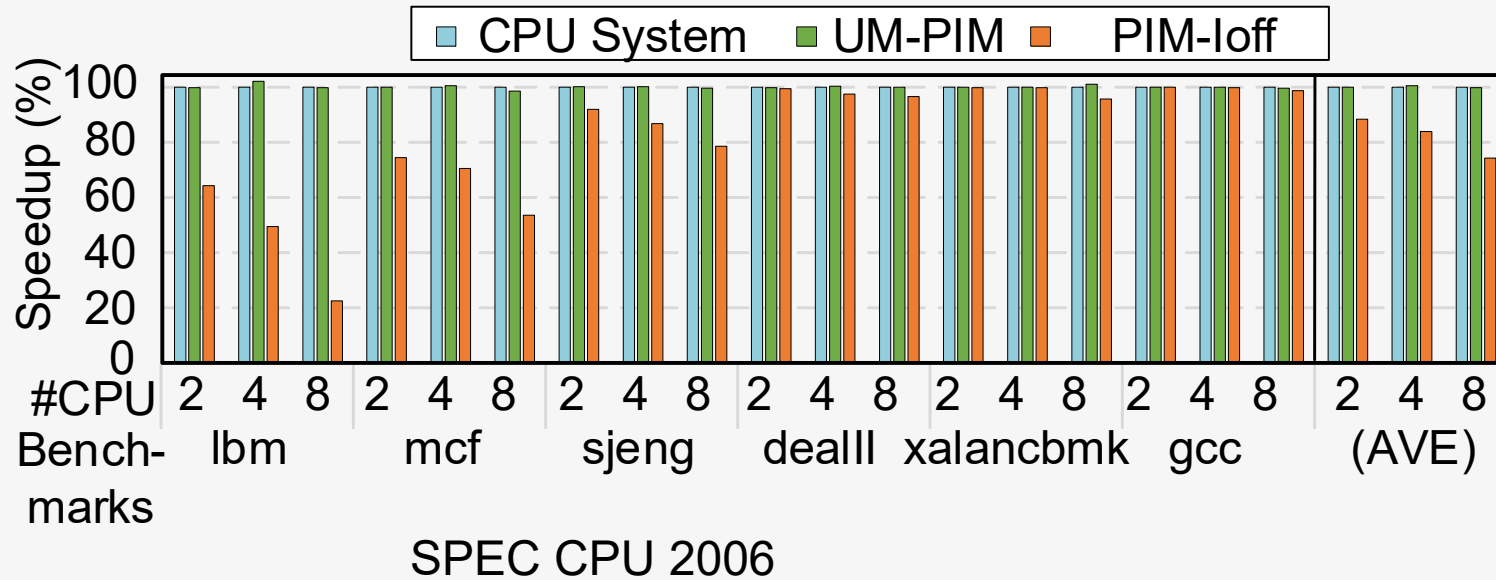    - lbm, mcf, sjeng, dealII, xalancbmk, gcc
  - PIM Workloads (requires CPU & PIM): from PRIM
    - BFS, PR, MLP, UNI, TC, SCAN-RSS/SSA, SEL, HST, NW, WFA, RL

- Baseline:
  - PIM-Ion: PIM System with Memory interleaving switched on
  - PIM-Ioff: PIM System with Rank & Channel interleaving switched off, (e.g. UPMEM)

Because of CPU Page's memory interleaving is switched on
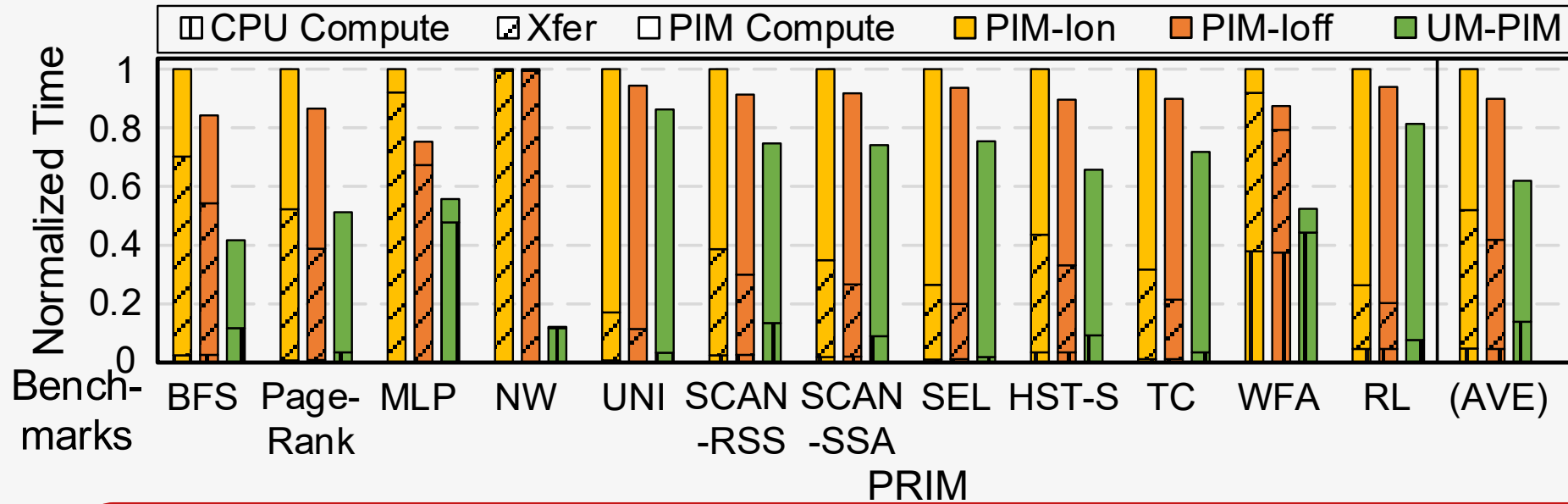


SPEC CPU 2006

PIM-Ioff: 22% performance degradation because of memory interleaving switched off

UM-PIM: **<0.1% performance loss** because CPU page is interleaved.

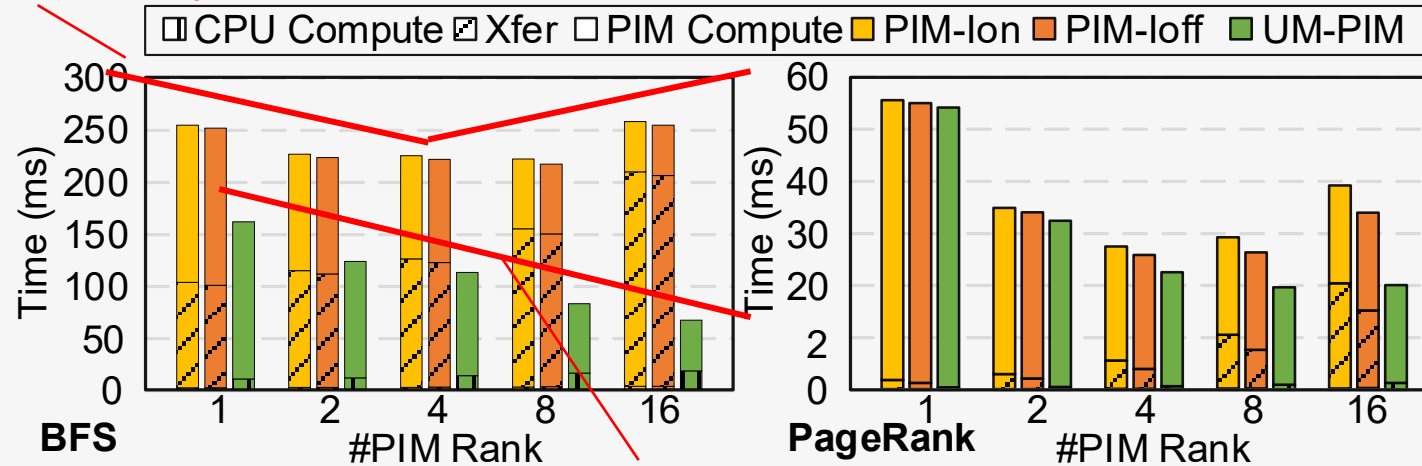●UM-PIM do not need data transfer between isolated memory spaces



Compared to PIM-Ioff (Rank & Channel interleaving switched off, UPMEM)
UM-PIM has **4.93× reduction on CPU time** (including Computing and Data Transfer), resulting in **1.96× speed up** on the whole workloads.

NW have intensive inter-PIM communication, therefore, it has the best speedup.

# Results on PIM Workloads

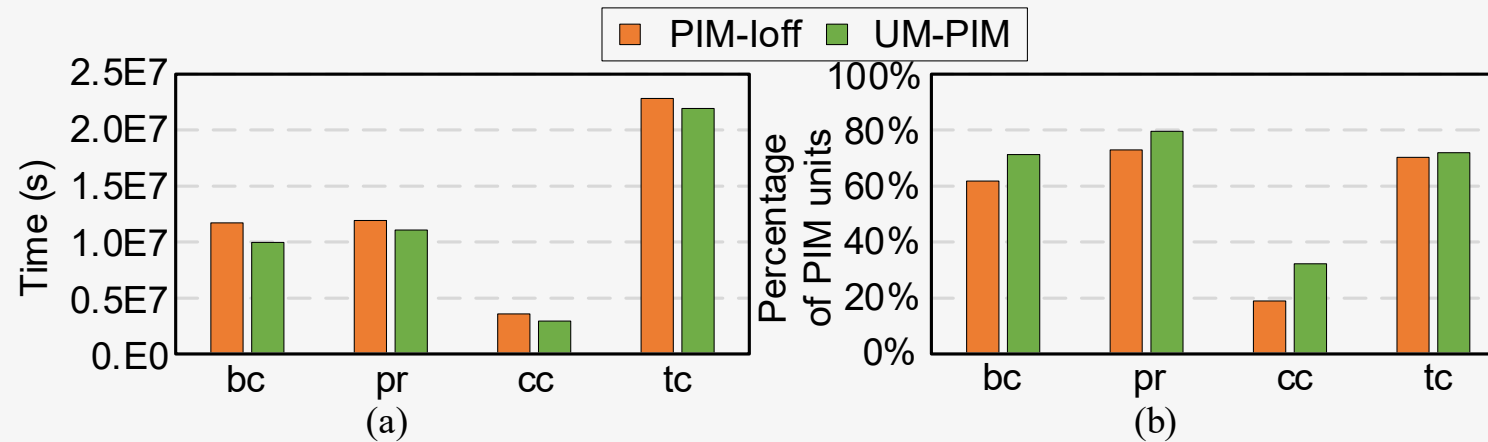⦿ UM-PIM do not need data transfer between isolated memory spaces



For workloads with intensive inter-PIM-units data transfer:
UM-PIM can **still benefit from computing using more PIM units**
In contrast, time on current PIM systems increases when using more PIM units

⊛PIM compilers decide program segment offload to PIM units according to offload overhead & compute speedup. UM-PIM reduce offload overhead because of eliminating data transfer.



UM-PIM can offload **8% more program segments** to PIM units, resulting in 1.13x speedup

# Discussion

- Extend to different Device number of DRAM:

  - Change number of SB in RC module. E.g. for x4 DRAM, RC circuit have 4 SBs.

- Extend to other DRAM type:

  - HBM and LPDDR do not have Device level in DRAM hierarchy. ATM and RC is not needed. Other APIs and memory management are still effective.

# More in the paper

- Details of CPU and PIM access PIM pages

- Extended DRAM Instruction support for UM-PIM

- Inter-PIM-units Communication APIs
    - Four inter-PIM units communications APIs like NCCL

# Conclusion

Key advantages of UM-PIM:

- **Uniform & shared memory space:** CPU and PIM pages co-exist in a uniform memory space, and data transfer is eliminated compared to isolated space design

- **Fast & Transparent data layout:** Dynamic address mapping and data re-layout for two kinds of pages are processed fast and transparently by DRAM-side hardware module

- **Compatibility with current CPU systems:** no modifications on CPU-side hardware and CPU page's memory management.