

# Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez Luna, Izzat El Hajj,  
Ivan Fernandez, Christina Giannoula,  
Geraldo F. Oliveira, Onur Mutlu

<https://arxiv.org/pdf/2105.03814.pdf>

<https://github.com/CMU-SAFARI/prim-benchmarks>

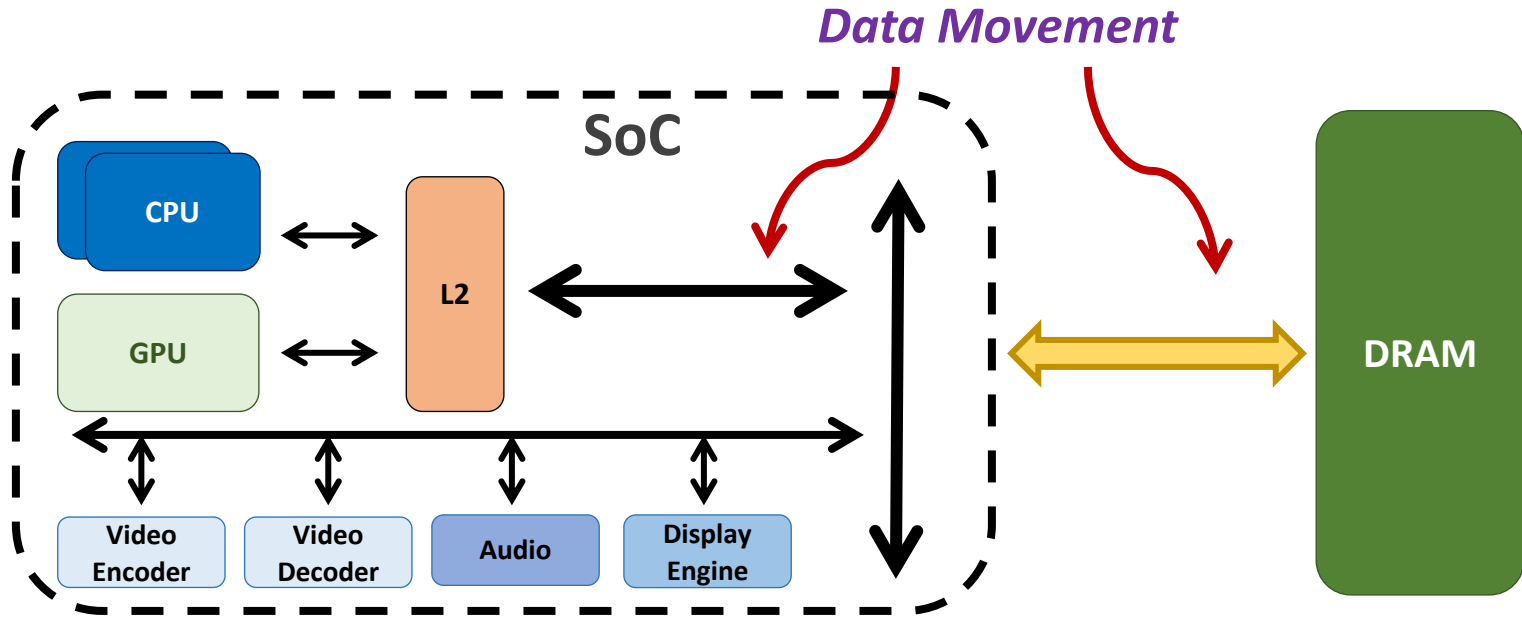
# Executive Summary

---

- **Data movement** between memory/storage units and compute units is a major contributor to execution time and energy consumption
- **Processing-in-Memory** (PIM) is a paradigm that can tackle the **data movement bottleneck**
  - Though explored for +50 years, technology challenges prevented the successful materialization
- UPMEM has designed and fabricated **the first publicly-available real-world PIM architecture**
  - DDR4 chips embedding in-order multithreaded DRAM Processing Units (DPUs)
- Our work:
  - **Introduction** to UPMEM programming model and PIM architecture
  - **Microbenchmark-based characterization** of the DPU
  - Benchmarking and **workload suitability** study
- Main contributions:
  - Comprehensive **characterization and analysis of the first commercially-available PIM architecture**
  - **PrIM** (Processing-In-Memory) benchmarks:
    - 16 workloads that are memory-bound in conventional processor-centric systems
    - Strong and weak scaling characteristics
  - Comparison to **state-of-the-art CPU and GPU**
- Takeaways:
  - Workload characteristics for **PIM suitability**
  - **Programming** recommendations
  - Suggestions and hints for **hardware and architecture designers** of future PIM systems
  - **PrIM**: (a) programming samples, (b) evaluation and comparison of current and future PIM systems

# Data Movement in Computing Systems

- Data movement dominates performance and is a major system energy bottleneck
- Total system energy: data movement accounts for
  - 62% in consumer applications\*,
  - 40% in scientific applications★,
  - 35% in mobile applications☆



\* Boroumand et al., "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks," ASPLOS 2018

★ Kestor et al., "Quantifying the Energy Cost of Data Movement in Scientific Applications," IISWC 2013

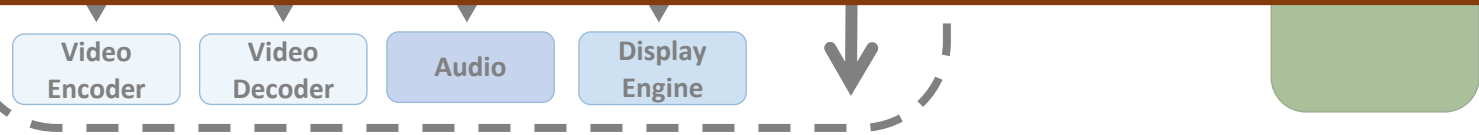
☆ Pandiyan and Wu, "Quantifying the energy cost of data movement for emerging smart phone workloads on mobile platforms," IISWC 2014

# Data Movement in Computing Systems

- Data movement dominates performance and is a major system energy bottleneck
- Total system energy: data movement accounts for
  - 62% in consumer applications\*,

Compute systems should be more data-centric

Processing-In-Memory proposes  
computing where it makes sense  
(where data resides)



\* Boroumand et al., "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks," ASPLOS 2018

\* Kestor et al., "Quantifying the Energy Cost of Data Movement in Scientific Applications," IISWC 2013

\* Pandiyan and Wu, "Quantifying the energy cost of data movement for emerging smart phone workloads on mobile platforms," IISWC 2014



# UPMEM Processing-in-DRAM Engine (2019)

- **Processing in DRAM Engine**
- Includes **standard DIMM modules**, with a **large number of DPU processors** combined with DRAM chips.
- Replaces **standard DIMMs**
  - DDR4 R-DIMM modules
    - 8GB+128 DPUs (16 PIM chips)
    - Standard 2x-nm DRAM process
  - **Large amounts of** compute & memory bandwidth



# Understanding a Modern PIM Architecture

---

## Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich   <sup>2</sup>American University of Beirut   <sup>3</sup>University of Malaga   <sup>4</sup>National Technical University of Athens

<https://arxiv.org/pdf/2105.03814.pdf>

<https://github.com/CMU-SAFARI/prim-benchmarks>

# Observations, Recommendations, Takeaways

## GENERAL PROGRAMMING RECOMMENDATIONS

1. Execute on the *DRAM Processing Units (DPUs)* **portions of parallel code** that are as long as possible.
2. Split the workload into **independent data blocks**, which the DPUs operate on independently.
3. Use **as many working DPUs** in the system as possible.
4. Launch at least **11 tasklets (i.e., software threads)** per DPU.

## PROGRAMMING RECOMMENDATION 1

For data movement between the DPU's MRAM bank and the WRAM, **use large DMA transfer sizes when all the accessed data is going to be used.**

## KEY OBSERVATION 7

**Larger CPU-DPU and DPU-CPU transfers** between the host main memory and the DRAM Processing Unit's Main memory (MRAM) banks **result in higher sustained bandwidth.**

## KEY TAKEAWAY 1

**The UPMEM PIM architecture is fundamentally compute bound.** As a result, **the most suitable work-loads are memory-bound.**

# PrIM Repository

- All microbenchmarks, benchmarks, and scripts
- <https://github.com/CMU-SAFARI/prim-benchmarks>

The screenshot shows the GitHub repository page for `CMU-SAFARI/prim-benchmarks`. At the top, the repository name is displayed with icons for Unwatch (2), Star (2), and Fork (1). Below this is a navigation bar with links for Code, Issues, Pull requests, Actions, Projects, Wiki, Security, Insights, and Settings. The main content area shows the `prim-benchmarks / README.md` file selected. Below the file name, it indicates the latest commit by Juan Gomez Luna, titled "PrIM -- first commit", dated 9 days ago. The repository has 1 contributor. The README content is displayed below, starting with the title "PrIM (Processing-In-Memory Benchmarks)".

**PrIM (Processing-In-Memory Benchmarks)**

PrIM is the first benchmark suite for a real-world processing-in-memory (PIM) architecture. PrIM is developed to evaluate, analyze, and characterize the first publicly-available real-world processing-in-memory (PIM) architecture, the [UPMEM](#) PIM architecture. The UPMEM PIM architecture combines traditional DRAM memory arrays with general-purpose in-order cores, called DRAM Processing Units (DPUs), integrated in the same chip.

PrIM provides a common set of workloads to evaluate the UPMEM PIM architecture with and can be useful for programming, architecture and system researchers all alike to improve multiple aspects of future PIM hardware and software. The workloads have different characteristics, exhibiting heterogeneity in their memory access patterns, operations and data types, and communication patterns. This repository also contains baseline CPU and GPU implementations of PrIM benchmarks for comparison purposes.

PrIm also includes a set of microbenchmarks can be used to assess various architecture limits such as compute throughput and memory bandwidth.

# Outline

---

- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PRIM Benchmarks
  - Roofline Model
  - Benchmark Diversity
- Evaluation
  - Strong and Weak Scaling
  - Comparison to CPU and GPU
- Key Takeaways

# Outline

---

- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PRIM Benchmarks
  - Roofline Model
  - Benchmark Diversity
- Evaluation
  - Strong and Weak Scaling
  - Comparison to CPU and GPU
- Key Takeaways

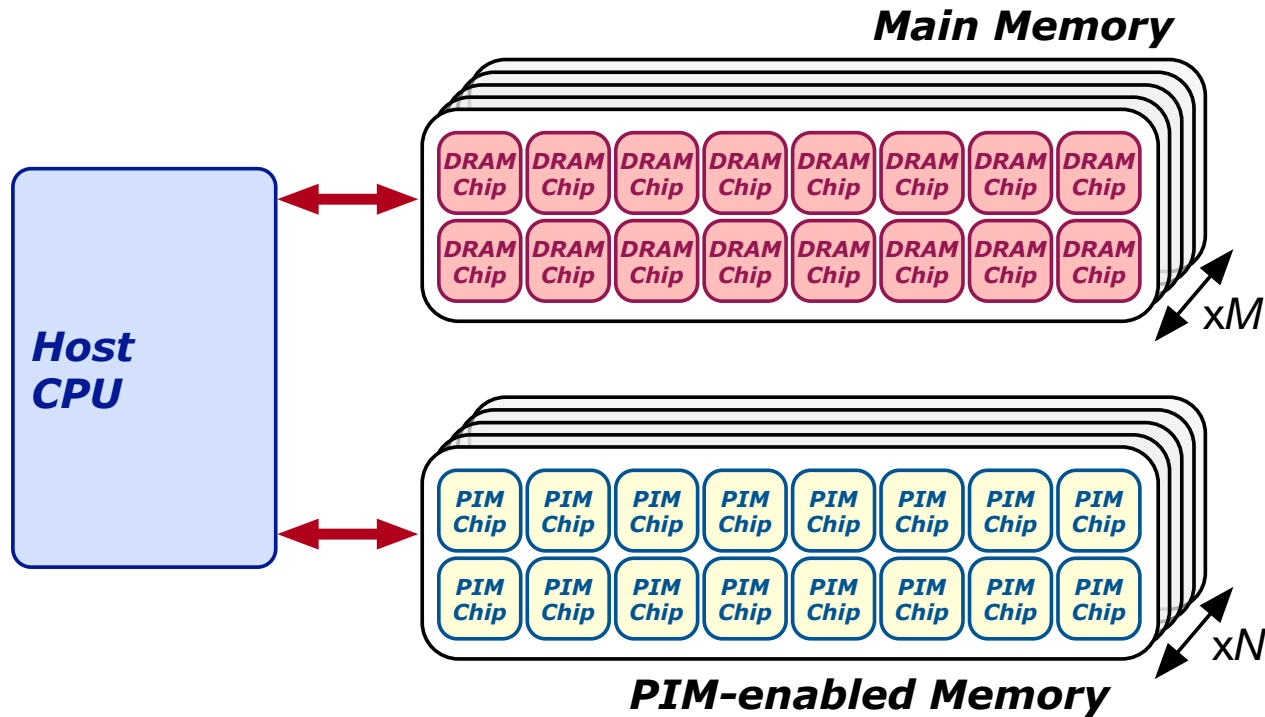
# Accelerator Model

---

- UPMEM DIMMs coexist with conventional DIMMs
- Integration of UPMEM DIMMs in a system follows an **accelerator model**
- UPMEM DIMMs can be seen as a **loosely coupled accelerator**
  - Explicit data movement between the main processor (host CPU) and the accelerator (UPMEM)
  - Explicit kernel launch onto the UPMEM processors
- This resembles GPU computing

# System Organization (I)

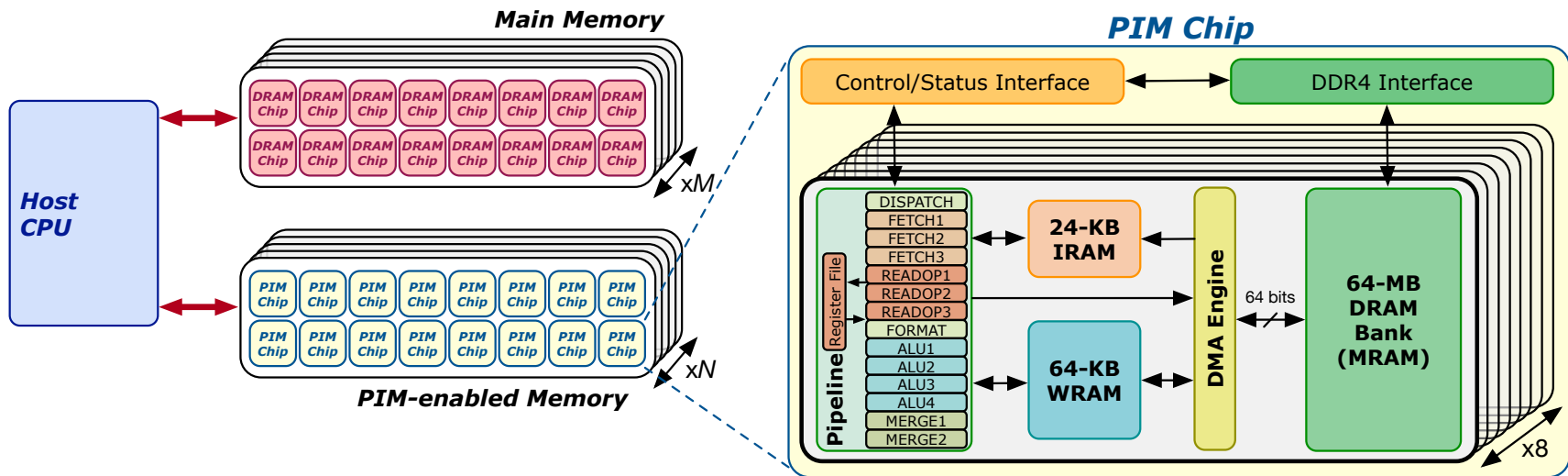
- In a UPMEM-based PIM system UPMEM DIMMs coexist with regular DDR4 DIMMs





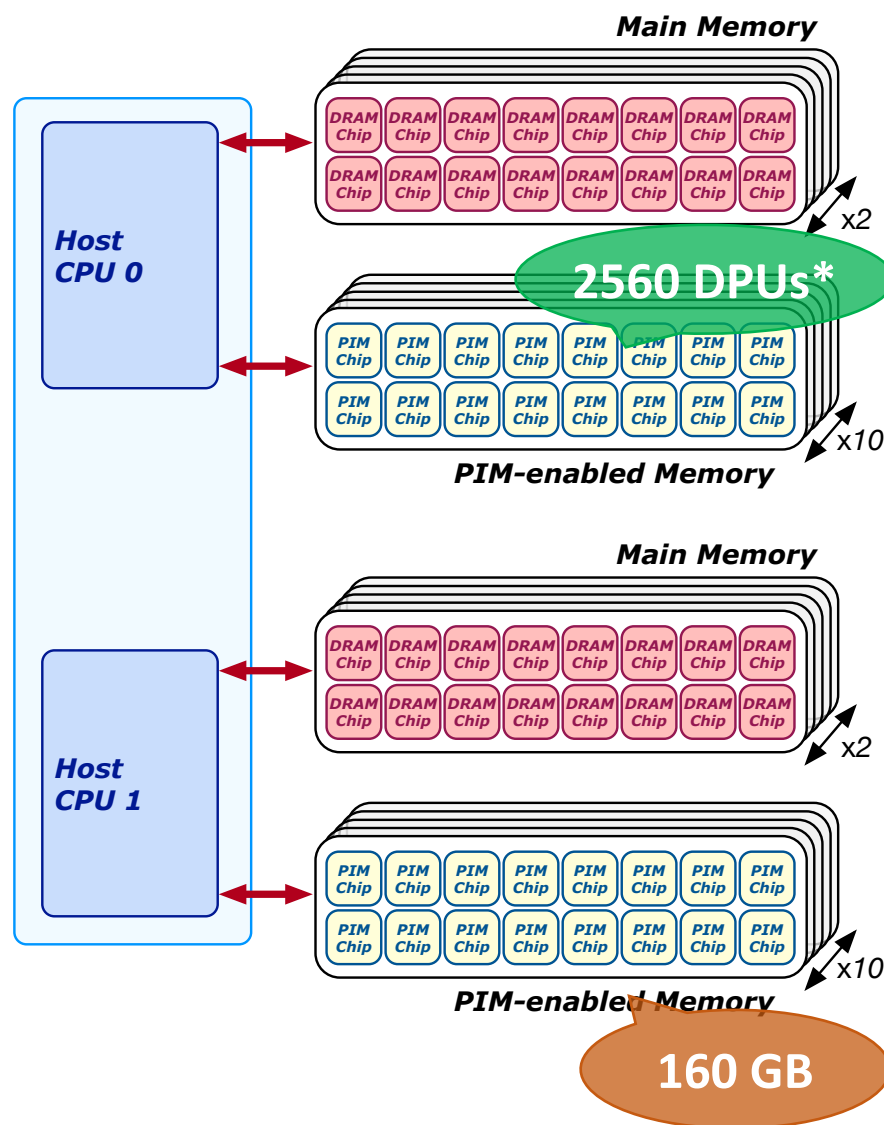
# System Organization (II)

- A UPMEM DIMM contains 8 or 16 chips
  - Thus, 1 or 2 ranks of 8 chips each
- Inside each PIM chip there are:
  - 8 64MB banks per chip: Main RAM (MRAM) banks
  - 8 DRAM Processing Units (DPUs) in each chip, 64 DPUs per rank

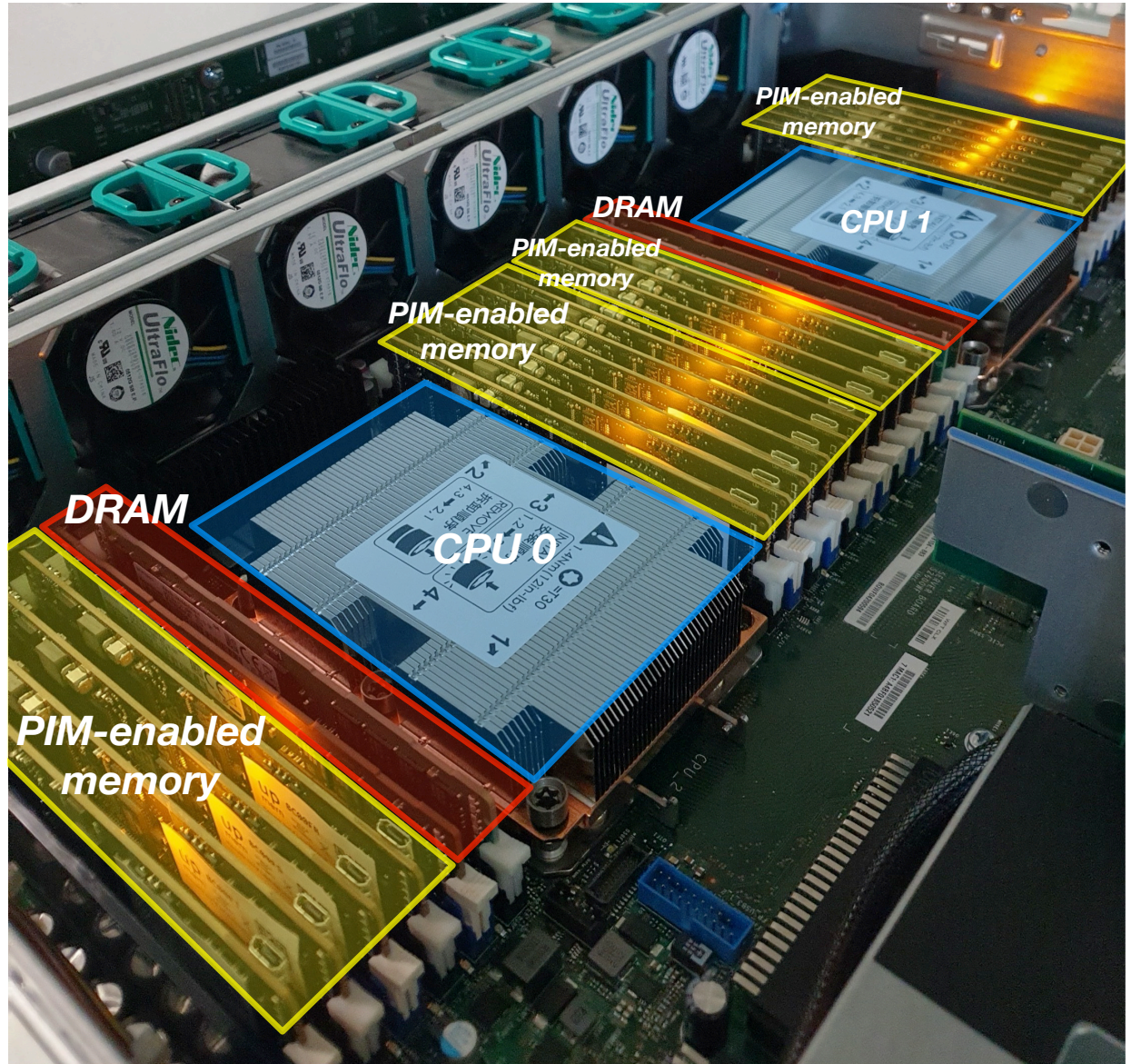
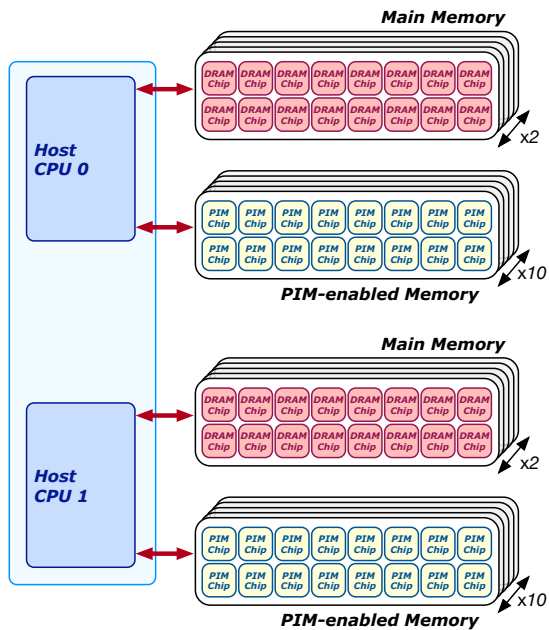


# 2,560-DPU System (I)

- UPMEM-based PIM system with 20 UPMEM DIMMs of 16 chips each (40 ranks)
  - P21 DIMMs
  - Dual x86 socket
    - UPMEM DIMMs coexist with regular DDR4 DIMMs
  - 2 memory controllers/socket (3 channels each)
  - 2 conventional DDR4 DIMMs on one channel of one controller

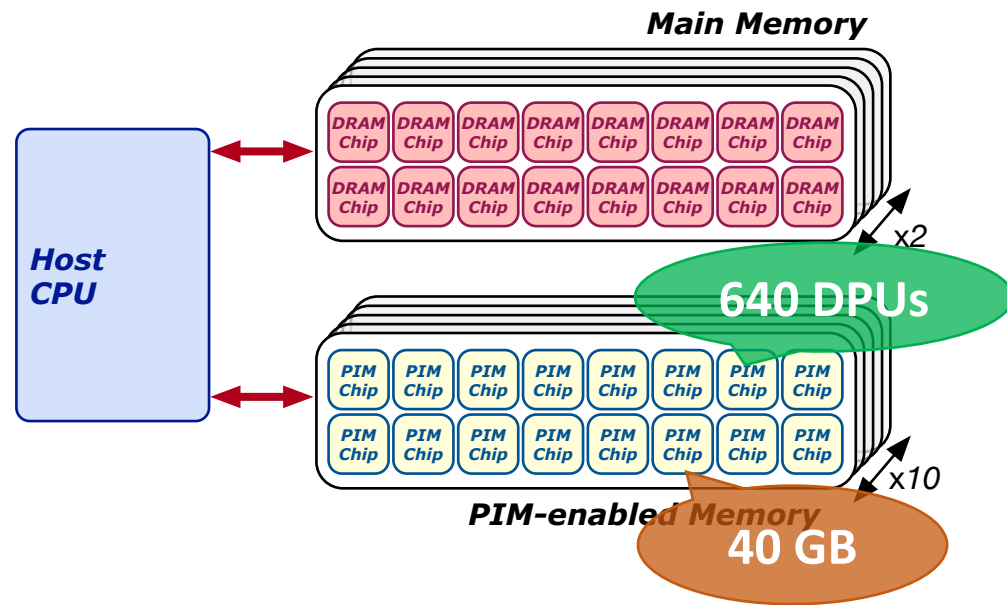


# 2,560-DPU System (II)



# 640-DPU System

- UPMEM-based PIM system with 10 UPMEM DIMMs of 8 chips each (10 ranks)
  - E19 DIMMs
  - x86 socket
    - 2 memory controllers (3 channels each)
    - 2 conventional DDR4 DIMMs on one channel of one controller





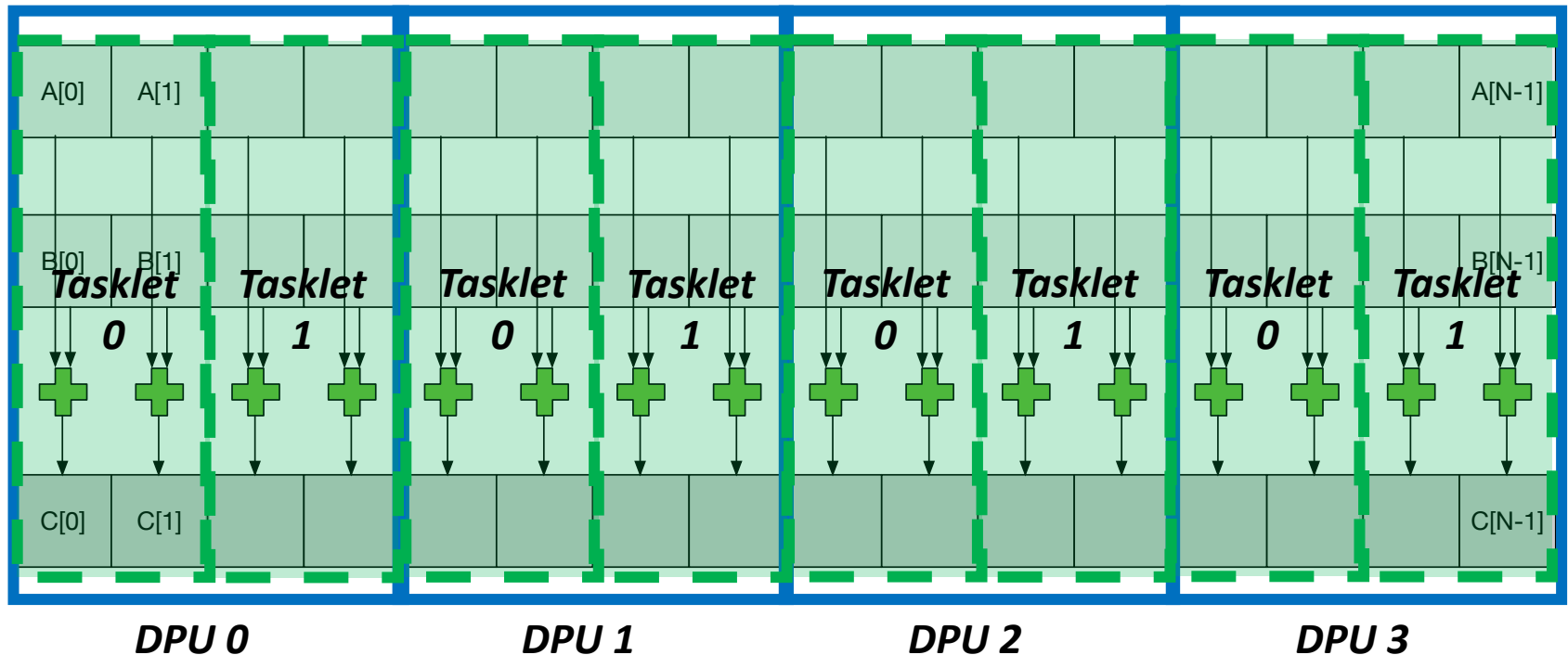
# Outline

---

- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PRIM Benchmarks
  - Roofline Model
  - Benchmark Diversity
- Evaluation
  - Strong and Weak Scaling
  - Comparison to CPU and GPU
- Key Takeaways

# Vector Addition (VA)

- Our first programming example
- We partition the input arrays across:
  - DPUs
  - Tasklets, i.e., software threads running on a DPU



# General Programming Recommendations

- From UPMEM programming guide\*, presentations★, and white papers☆

## ***GENERAL PROGRAMMING RECOMMENDATIONS***

1. Execute on the *DRAM Processing Units (DPUs)* **portions of parallel code** that are as long as possible.
2. Split the workload into **independent data blocks**, which the DPUs operate on independently.
3. Use **as many working DPUs** in the system as possible.
4. Launch at least **11 tasklets (i.e., software threads)** per DPU.

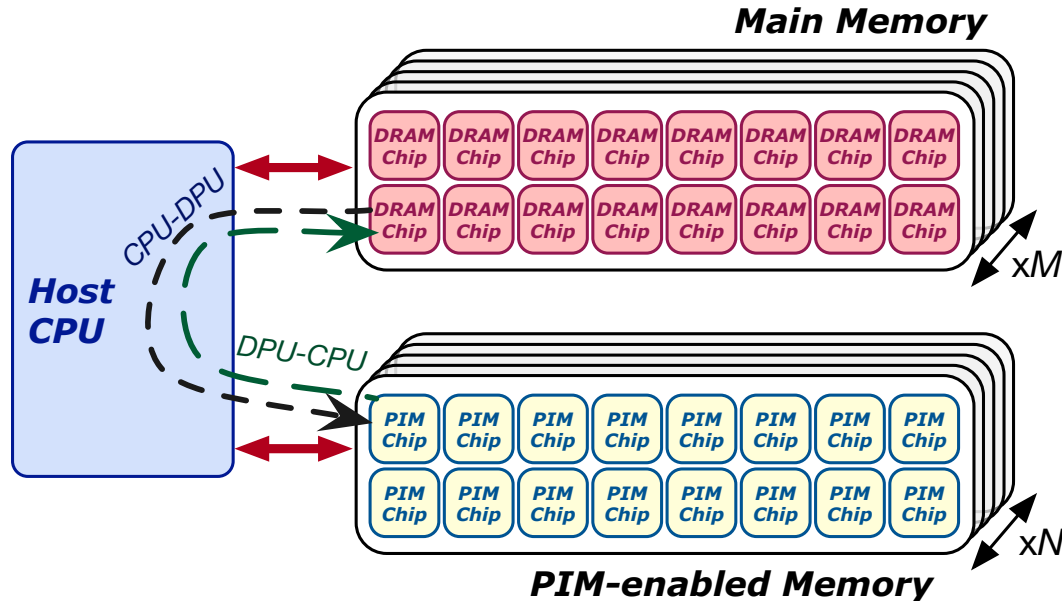
\* <https://sdk.upmem.com/2021.1.1/index.html>

★ F. Devaux, "The true Processing In Memory accelerator," HotChips 2019. doi: 10.1109/HOTCHIPS.2019.8875680

☆ UPMEM, "Introduction to UPMEM PIM. Processing-in-memory (PIM) on DRAM Accelerator," White paper

# CPU-DPU/DPU-CPU Data Transfers

- CPU-DPU and DPU-CPU transfers
  - Between host CPU's main memory and DPUs' MRAM banks

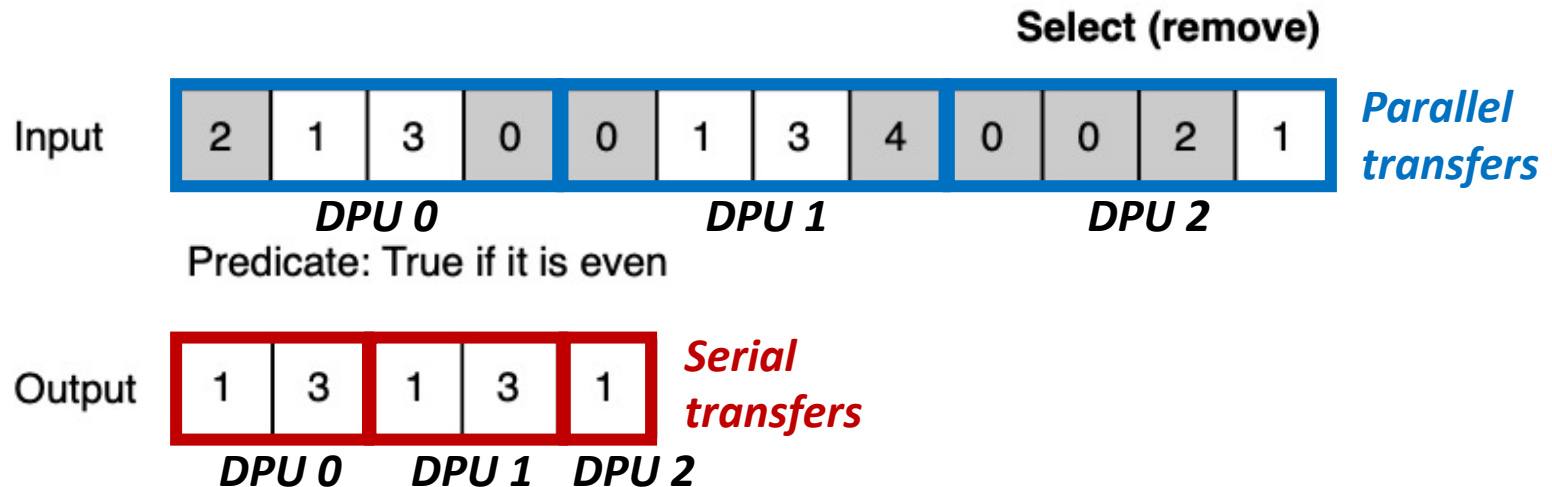


- Serial CPU-DPU/DPU-CPU transfers:
  - A single DPU (i.e., 1 MRAM bank)
- Parallel CPU-DPU/DPU-CPU transfers:
  - Multiple DPUs (i.e., many MRAM banks)
- Broadcast CPU-DPU transfers:
  - Multiple DPUs with a single buffer



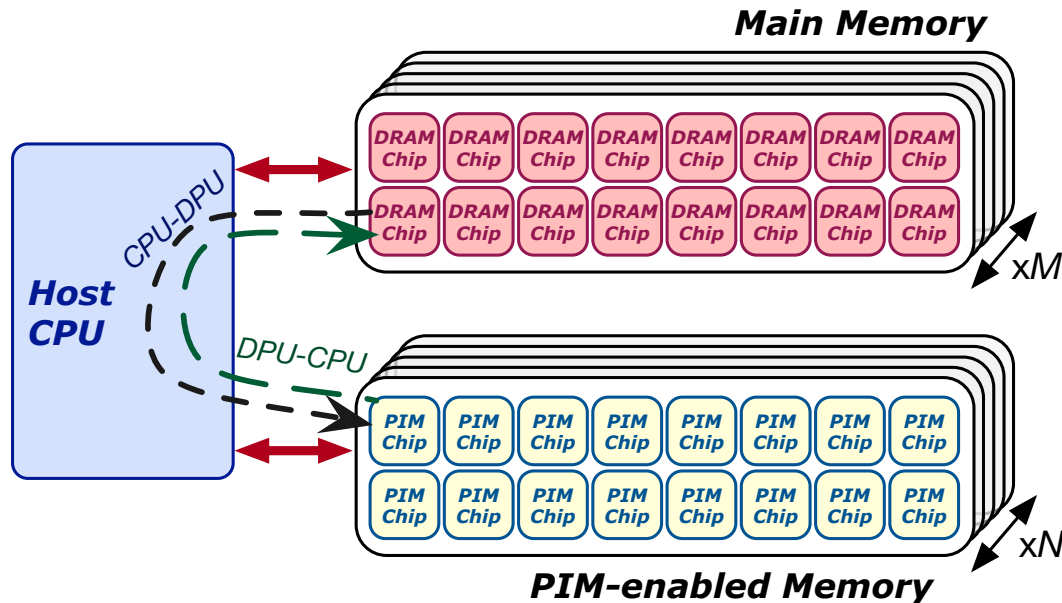
# Different Types of Transfers in a Program

- An example benchmark that uses both parallel and serial transfers
- Select (SEL)
  - Remove even values



# Inter-DPU Communication

- There is **no direct communication channel** between DPUs



- Inter-DPU communication** takes place via the **host CPU** using **CPU-DPU** and **DPU-CPU** transfers
- Example communication patterns:
  - Merging of partial results to obtain the final result
    - Only **DPU-CPU** transfers
  - Redistribution of intermediate results for further computation
    - DPU-CPU** transfers and **CPU-DPU** transfers

# How Fast are these Data Transfers?

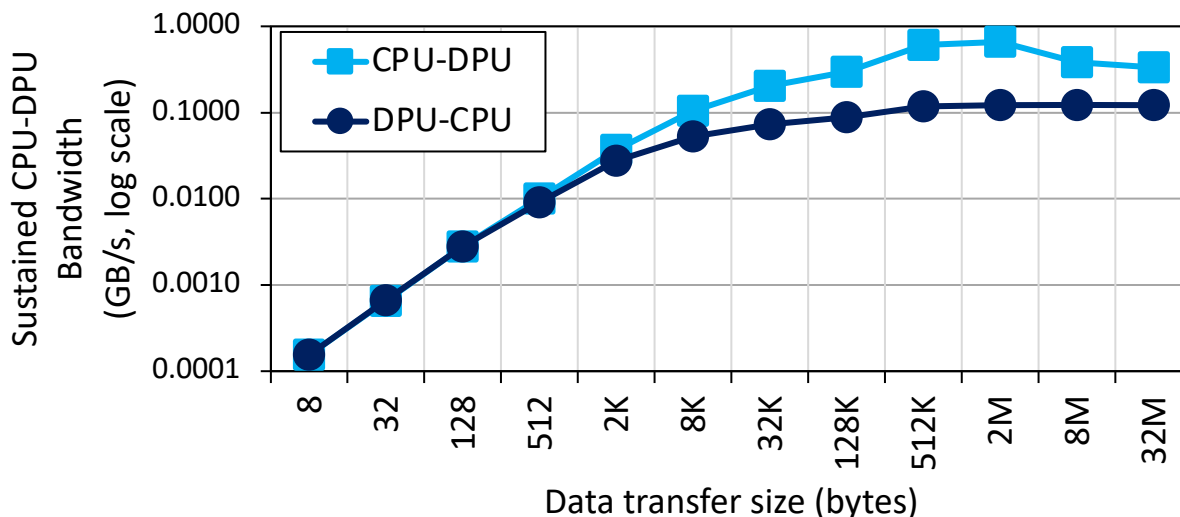
- With a microbenchmark, we obtain the **sustained bandwidth of all types of CPU-DPU and DPU-CPU transfers**
- Two experiments:
  - **1 DPU**: variable CPU-DPU and DPU-CPU transfer size (**8 bytes to 32 MB**)
  - **1 rank**: 32 MB CPU-DPU and DPU-CPU transfers to/from a set of **1 to 64 MRAM banks** within the same rank
- We do not experiment with more than one rank
  - Preliminary experiments show that the UPMEM SDK\* only parallelizes transfers within the same rank

**DDR4 bandwidth** bounds the maximum transfer bandwidth

The cost of the **transfers can be amortized**,  
if enough computation is run on the DPUs

# CPU-DPU/DPU-CPU Transfers: 1 DPU

- Data transfer size varies between 8 bytes and 32 MB

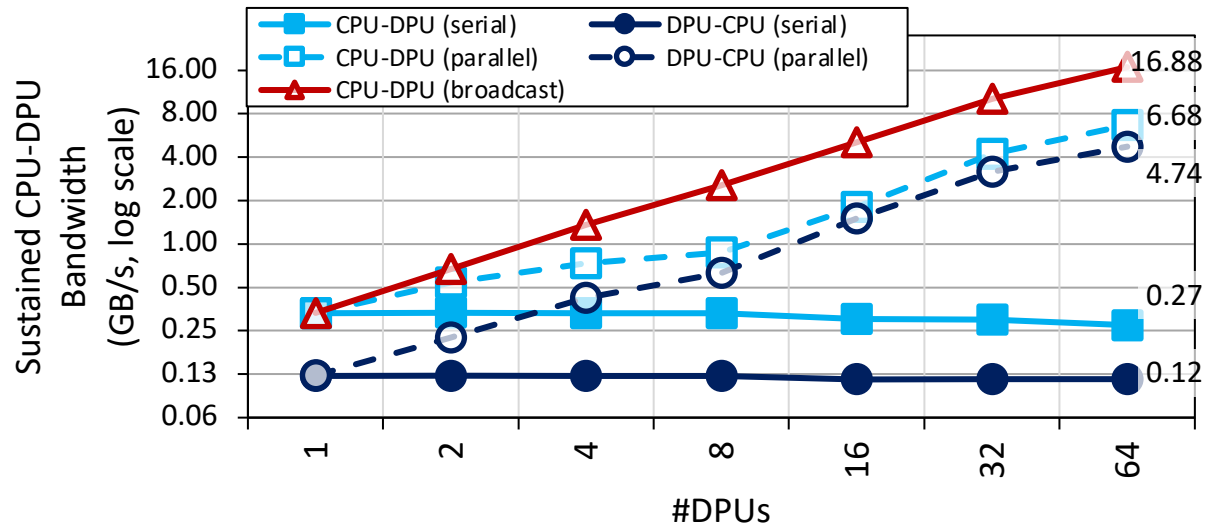


## KEY OBSERVATION 7

**Larger CPU-DPU and DPU-CPU transfers** between the host main memory and the DRAM Processing Unit's Main memory (MRAM) banks **result in higher sustained bandwidth.**

# CPU-DPU/DPU-CPU Transfers: 1 Rank

- CPU-DPU (serial/parallel/**broadcast**) and DPU-CPU (serial/parallel)
- The number of DPUs varies between 1 and 64



## KEY OBSERVATION 8

The **sustained bandwidth of parallel CPU-DPU and DPU-CPU transfers** between the host main memory and the DRAM Processing Unit's Main memory (MRAM) banks **increases with the number of DRAM Processing Units inside a rank.**

# CPU-DPU/DPU-CPU Transfers: 1 Rank

- CPU-DPU (serial/parallel/**broadcast**) and DPU-CPU (serial/parallel)
- The number of DPUs varies between 1 and 64



## Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>American University of Beirut <sup>3</sup>University of Malaga <sup>4</sup>National Technical University of Athens

### KEY OBSERVATION 8

The sustained bandwidth of parallel CPU-DPU and DPU-CPU transfers between the host main memory and the DPA-M Processing

Unit of

DL <https://arxiv.org/pdf/2105.03814.pdf>  
<https://github.com/CMU-SAFARI/prim-benchmarks>

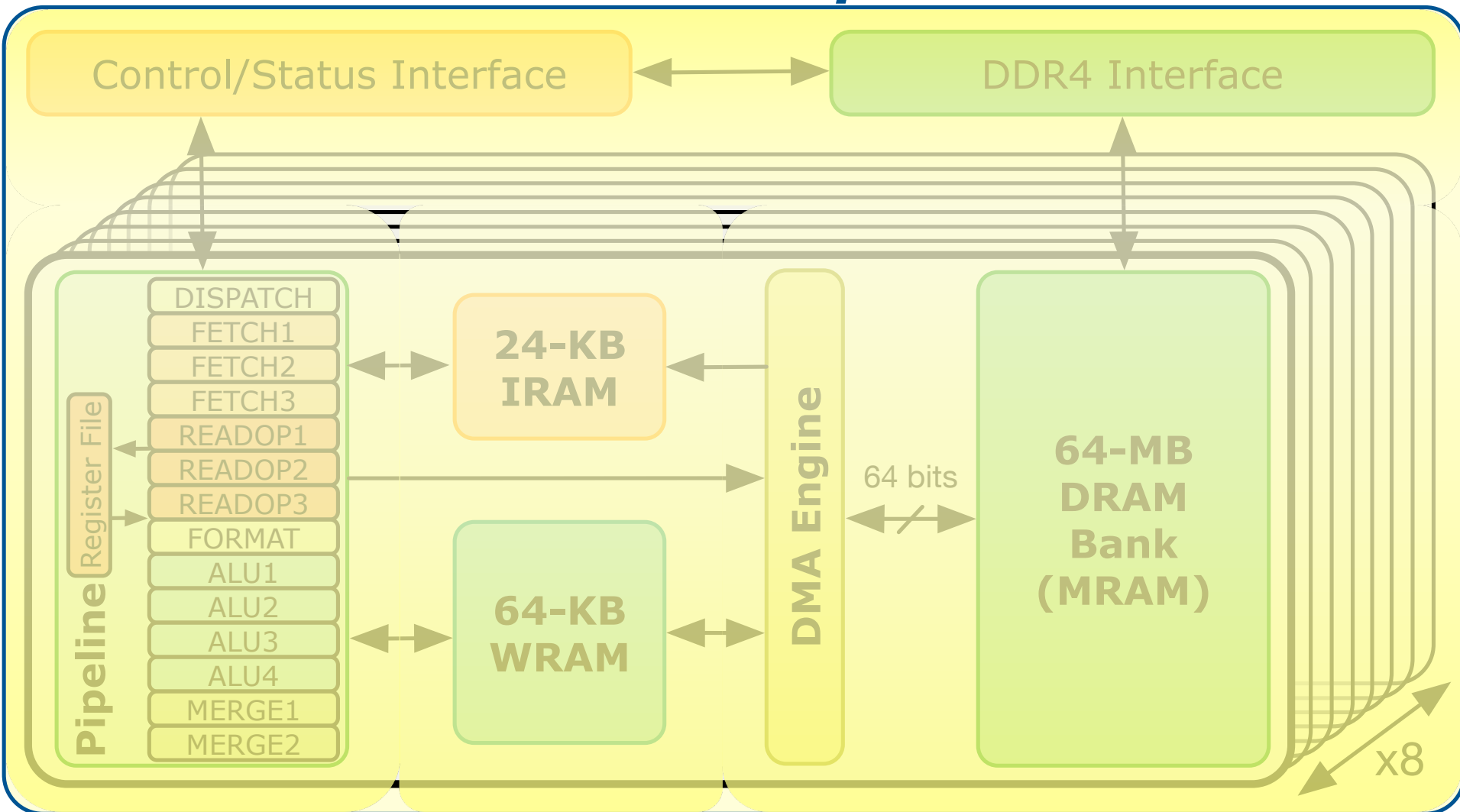
# Outline

---

- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PRIM Benchmarks
  - Roofline Model
  - Benchmark Diversity
- Evaluation
  - Strong and Weak Scaling
  - Comparison to CPU and GPU
- Key Takeaways

# DRAM Processing Unit

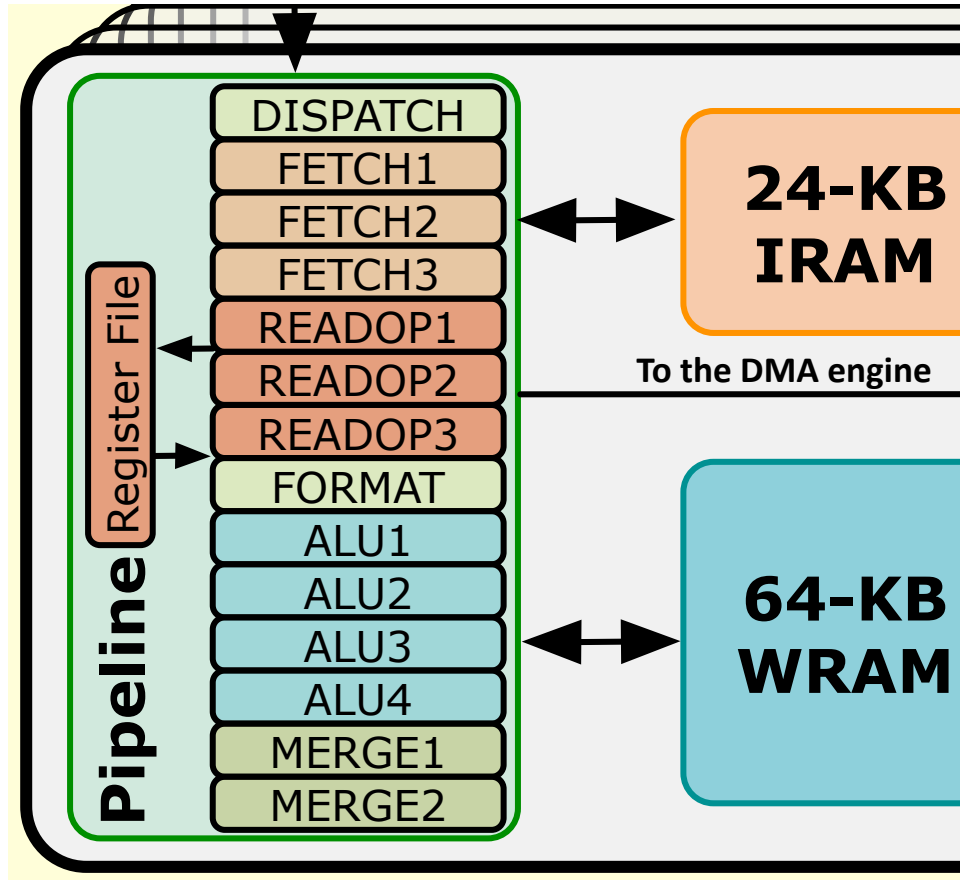
## *PIM Chip*





# DPU Pipeline

- In-order pipeline
  - Up to 350 MHz
- Fine-grain multithreaded
  - 24 hardware threads
- 14 pipeline stages
  - **DISPATCH**: Thread selection
  - **FETCH**: Instruction fetch
  - **READOP**: Register file
  - **FORMAT**: Operand formatting
  - **ALU**: Operation and WRAM
  - **MERGE**: Result formatting



# Arithmetic Throughput: Microbenchmark

---

- Goal
  - Measure the maximum arithmetic throughput for different datatypes and operations
- Microbenchmark
  - We stream over an array in WRAM and perform read-modify-write operations
  - Experiments on one DPU
  - We vary the number of tasklets from 1 to 24
  - Arithmetic operations: add, subtract, multiply, divide
  - Datatypes: int32, int64, float, double
- We measure cycles with an accurate cycle counter that the SDK provides
  - We include WRAM accesses (including address calculation) and arithmetic operation

# Microbenchmark for INT32 ADD Throughput

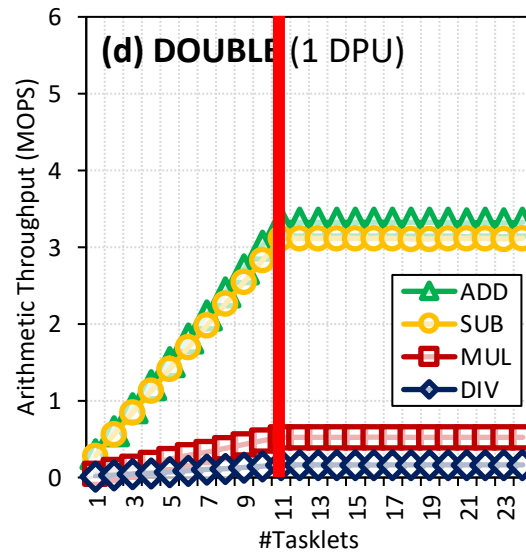
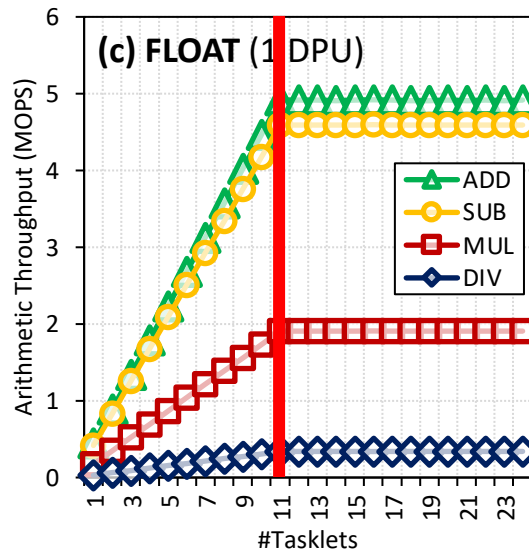
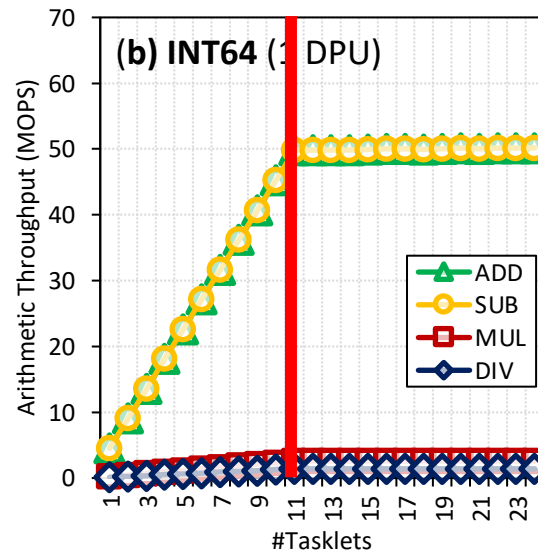
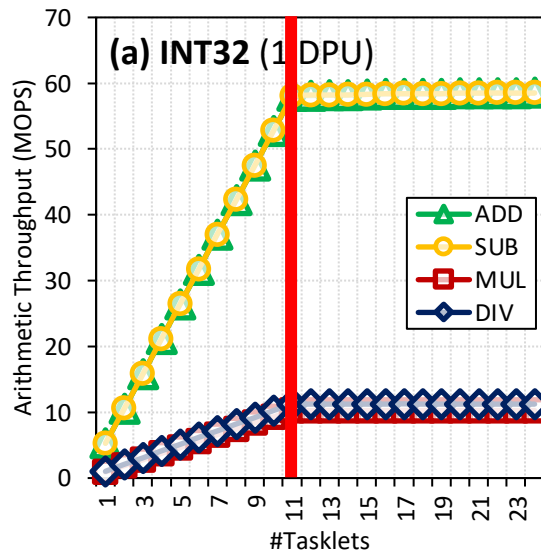
C-based code

```
1  #define SIZE 256
2  int* bufferA = mem_alloc(SIZE * sizeof(int));
3  for(int i = 0; i < SIZE; i++){
4      int temp = bufferA[i];
5      temp += scalar;
6      bufferA[i] = temp;
7  }
```

Compiled code  
(UPMEM DPU ISA)

```
1  move r2, 0
2  .LBB0_1:                                // Loop header
3  lsl_add r3, r0, r2, 2                    // Address calculation
4  lw r4, r3, 0                            // Load from WRAM
5  add r4, r4, r1                           // Add
6  sw r3, 0, r4                            // Store to WRAM
7  add r2, r2, 1                           // Index update
8  jneq r2, 256, .LBB0_1                   // Conditional jump
```

# Arithmetic Throughput: 11 Tasklets

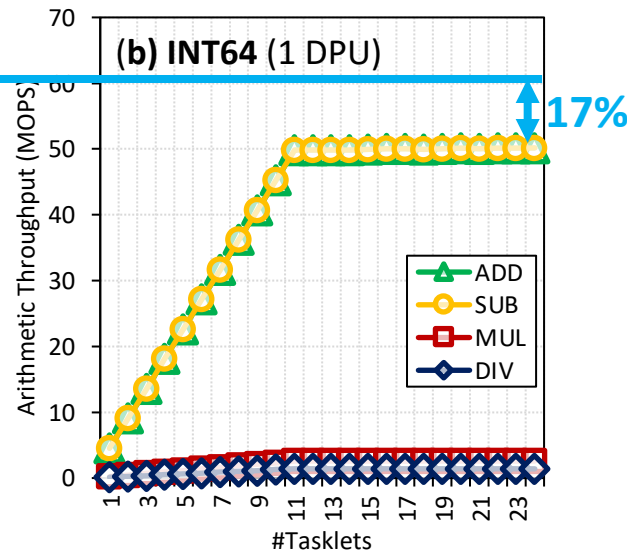
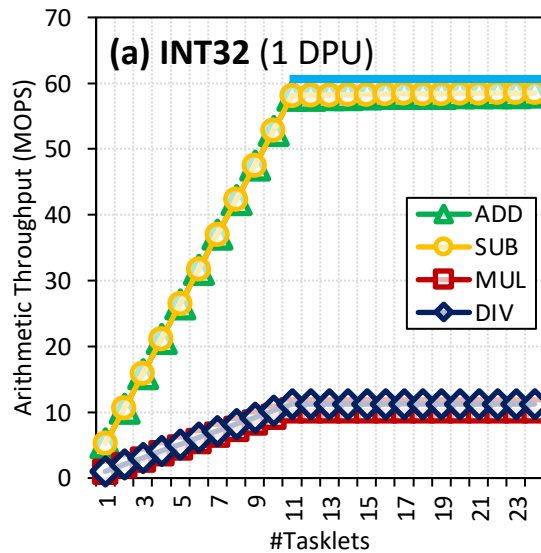


## KEY OBSERVATION 1

**The arithmetic throughput of a DRAM Processing Unit saturates at 11 or more tasklets.**

This observation is consistent for different datatypes (INT32, INT64, UINT32, UINT64, FLOAT, DOUBLE) and operations (ADD, SUB, MUL, DIV).

# Arithmetic Throughput: ADD/SUB



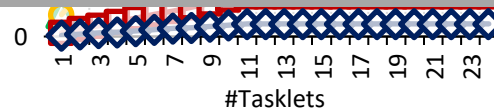
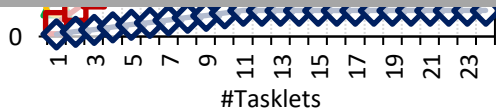
INT32 ADD/SUB are  
17% faster than  
INT64 ADD/SUB

Can we explain the peak throughput?

Peak throughput at 11 tasklets.

One instruction retires every cycle when the pipeline is full

$$\text{Arithmetic Throughput (in OPS)} = \frac{\text{frequency}_{\text{DPU}}}{\#instructions}$$



# Arithmetic Throughput: #Instructions

- Compiler explorer: <https://dpu.dev>

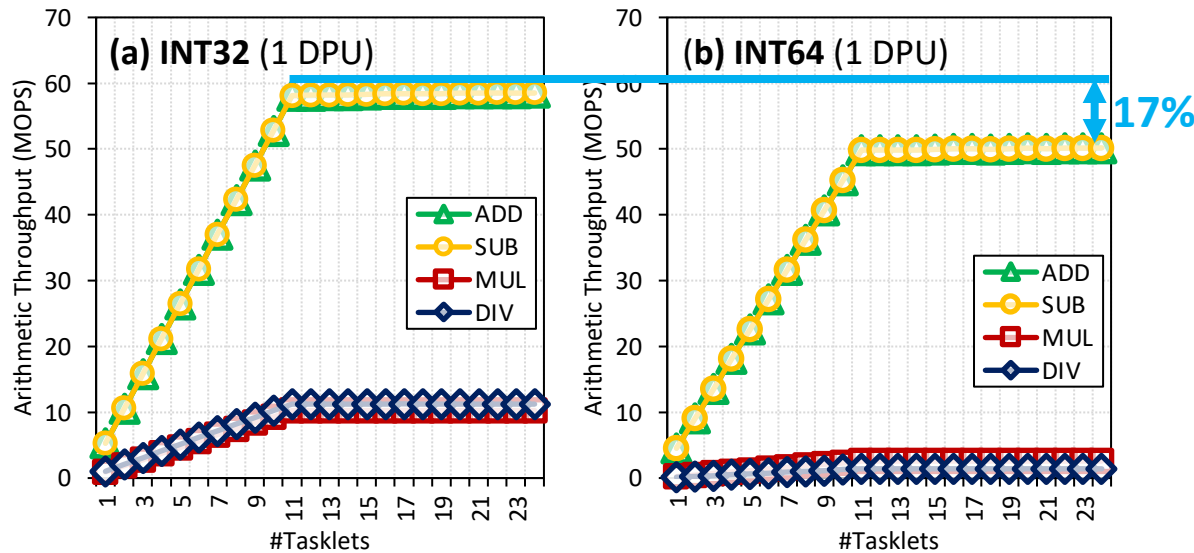
```
1  #define BLOCK_SIZE 1024
2
3  typedef int T;
4  void Benchmark__32bits(T *cache_A, T scalar) {
5      for (int i = 0; i < BLOCK_SIZE / sizeof(T); i++){
6          // WRAM READ
7          T temp = cache_A[i];
8
9          temp += scalar; // ADD
10
11         // WRAM WRITE
12         cache_A[i] = temp;
13     }
14 }
15
16 typedef long T_long;
17 void Benchmark__64bits(T_long *cache_A, T_long scalar) {
18     for (int i = 0; i < BLOCK_SIZE / sizeof(T_long); i++){
19         // WRAM READ
20         T_long temp = cache_A[i];
21
22         temp += scalar; // ADD
23     }
24
25
26
27
```

A ▾ ☐ 11010 ☐ ./a.out ☒ .LX0: ☒ .text ☒ // ☐ \

```
1 Benchmark__32bits:
2     move r2, 0
3 .LBB0_1:
4     lsl_add r3, r0, r2, 2
5     lw r4, r3, 0
6     add r4, r4, r1
7     sw r3, 0, r4
8     add r2, r2, 1
9     jneq r2, 256, .LBB0_1
10    jump r23
11 Benchmark__64bits:
12     move r1, 0
13 .LBB1_1:
14     lsl_add r4, r0, r1, 3
15     ld d6, r4, 0
16     add r7, r7, r3
17     addc r6, r6, r2
18     sd r4, 0, d6
19     add r1, r1, 1
20     jneq r1, 128, .LBB1_1
21    jump r23
```

6 instructions in the 32-bit ADD/SUB microbenchmark  
7 instructions in the 64-bit ADD/SUB microbenchmark

# Arithmetic Throughput: ADD/SUB



INT32 ADD/SUB are  
17% faster than  
INT64 ADD/SUB

Can we explain the peak throughput?

Peak throughput at 11 tasklets.

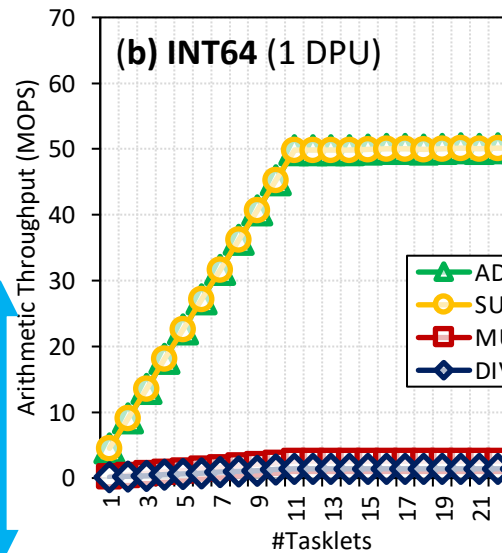
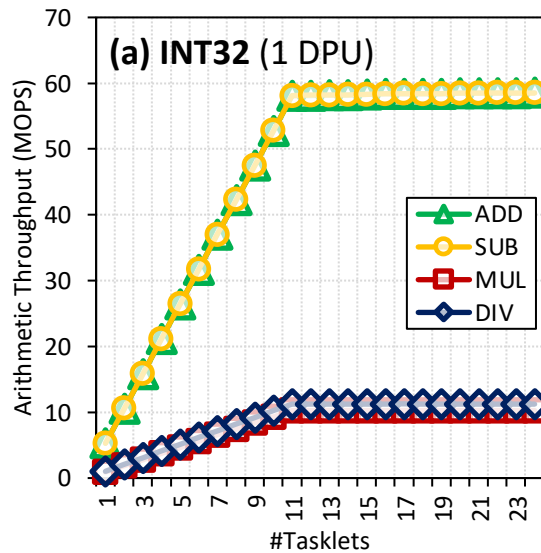
One instruction retires every cycle when the pipeline is full

$$\text{Arithmetic Throughput (in OPS)} = \frac{\text{frequency}_{DPU}}{\#instructions}$$

64-bit ADD/SUB: 7 instructions → 50.00 MOPS

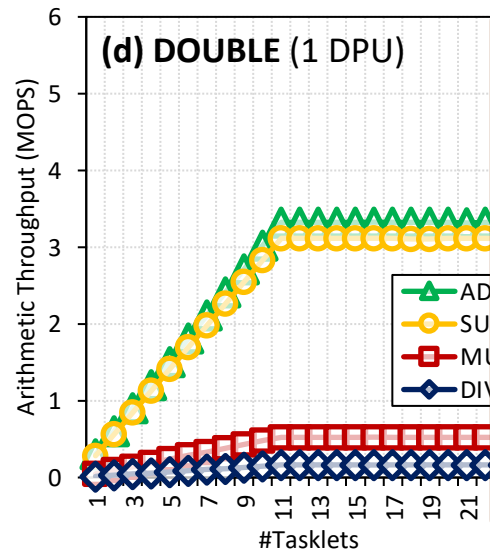
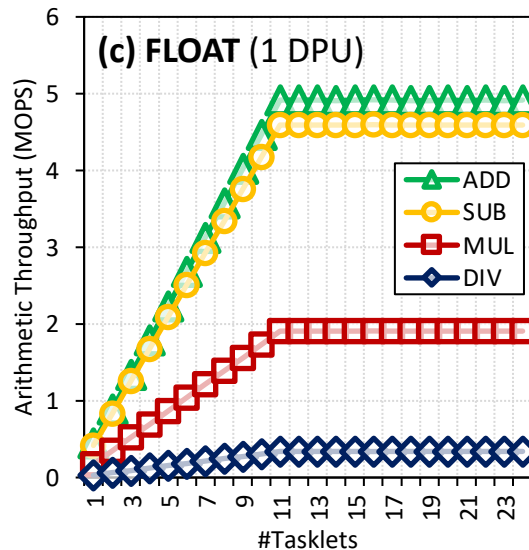
at  $\text{frequency}_{DPU} = 350 \text{ MHz}$

# Arithmetic Throughput: MUL/DIV



Huge throughput difference between ADD/SUB and MUL/DIV

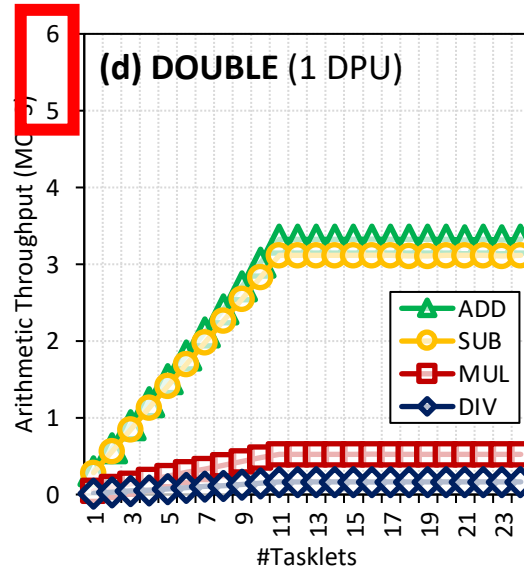
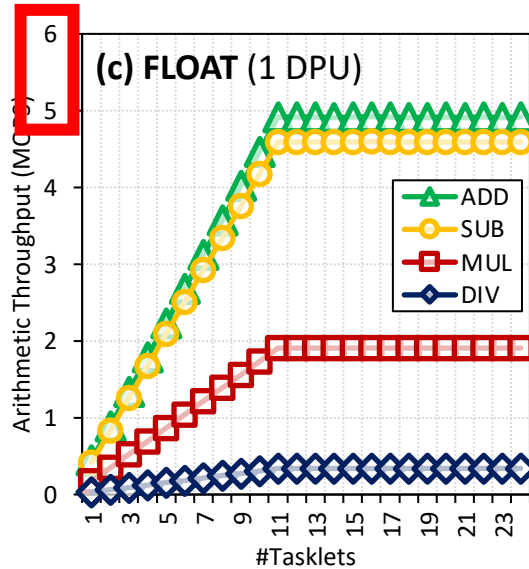
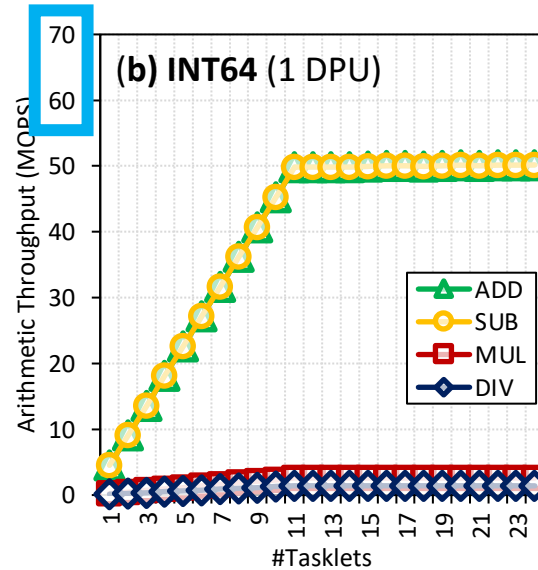
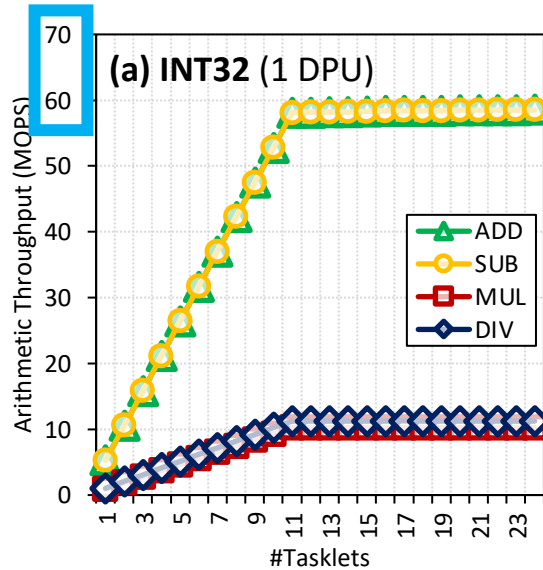
DPU's do *not* have a 32-bit multiplier



MUL/DIV implementation is based on an instruction that performs *bit shifting and addition* in 1 cycle (MUL/DIV take a maximum of 32 instructions)



# Arithmetic Throughput: Native Support

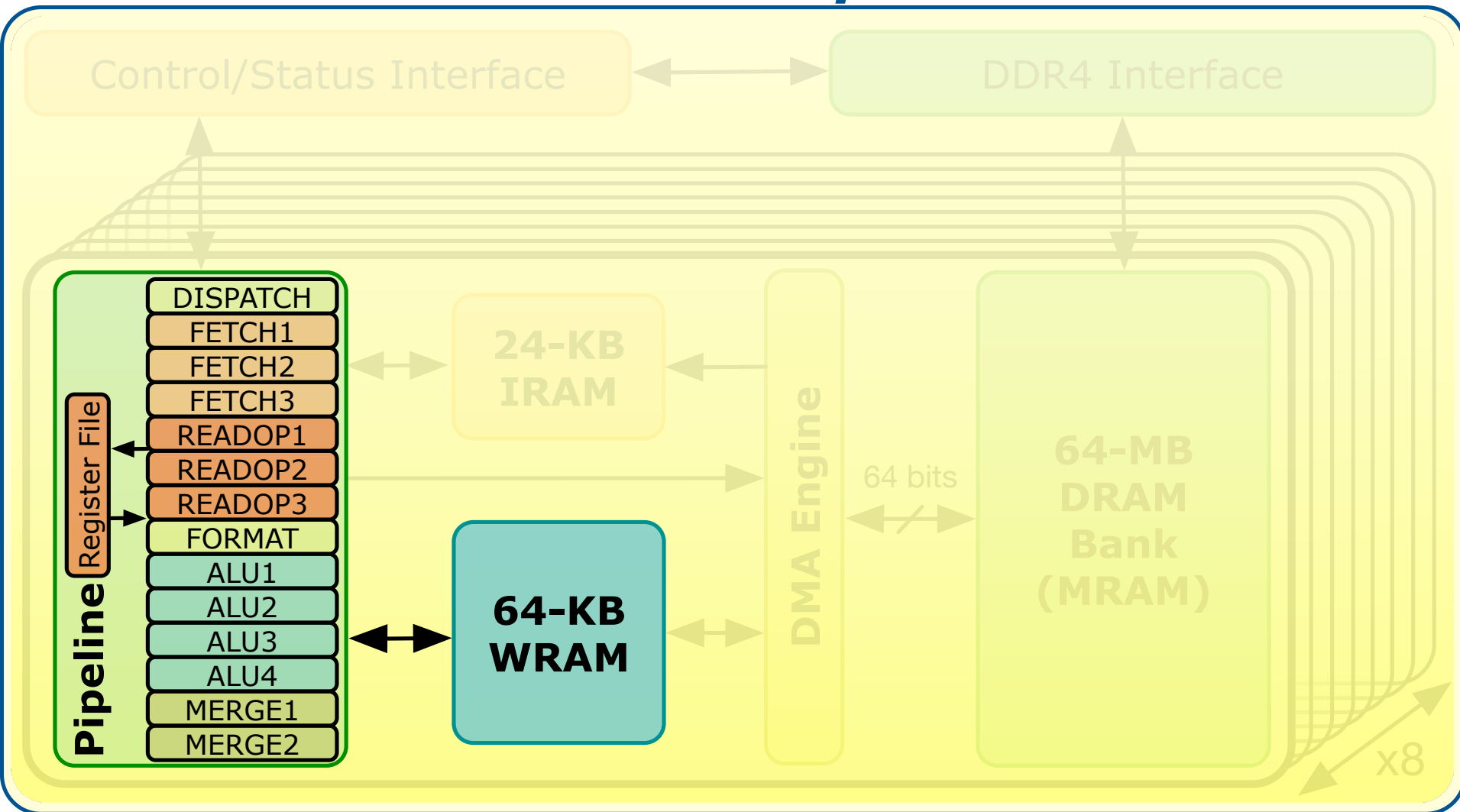


## KEY OBSERVATION 2

- DPUs provide **native hardware support for 32- and 64-bit integer addition and subtraction**, leading to high throughput for these operations.
- DPUs do ***not* natively support 32- and 64-bit multiplication and division, and floating point operations**. These operations are **emulated by the UPMEM runtime library**, leading to much lower throughput.

# DPU: WRAM Bandwidth

## *PIM Chip*



# WRAM Bandwidth: Microbenchmark

---

- Goal
  - Measure the **WRAM bandwidth** for the STREAM benchmark
- Microbenchmark
  - We implement the four versions of STREAM: **COPY, ADD, SCALE, and TRIAD**
  - The operations performed in ADD, SCALE, and TRIAD are **addition, multiplication, and addition+multiplication**, respectively
  - We vary the number of tasklets from 1 to 16
  - We show results for 1 DPU
- We do not include accesses to MRAM

# STREAM Benchmark in WRAM

// COPY

```
for(int i = 0; i < SIZE; i++){  
    bufferB[i] = bufferA[i];  
}
```

8 bytes read, 8 bytes written,  
no arithmetic operations

// ADD

```
for(int i = 0; i < SIZE; i++){  
    bufferC[i] = bufferA[i] + bufferB[i];  
}
```

16 bytes read, 8 bytes written,  
ADD

// SCALE

```
for(int i = 0; i < SIZE; i++){  
    bufferB[i] = scalar * bufferA[i];  
}
```

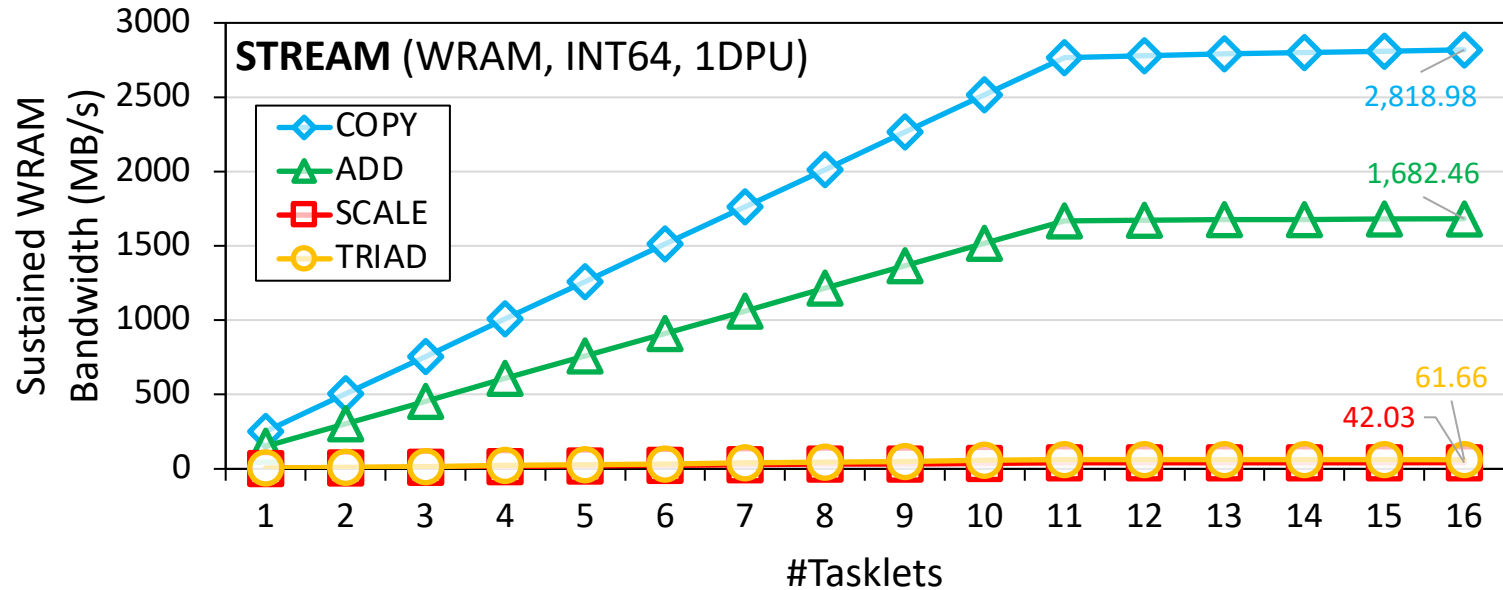
8 bytes read, 8 bytes written,  
MUL

// TRIAD

```
for(int i = 0; i < SIZE; i++){  
    bufferC[i] = bufferA[i] + scalar * bufferB[i];  
}
```

16 bytes read, 8 bytes written,  
MUL, ADD

# WRAM Bandwidth: STREAM

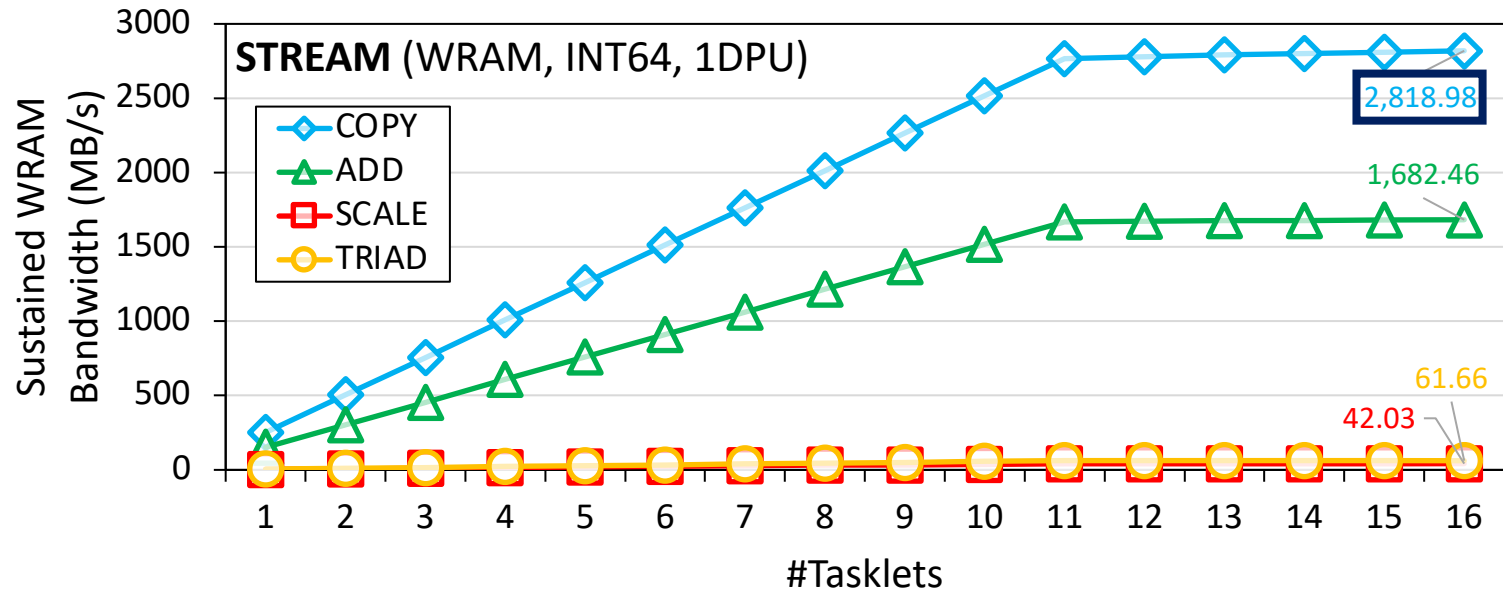


How can we estimate the bandwidth?

Assuming that the pipeline is full, and *Bytes* is the number of bytes read and written:

$$WRAM\ Bandwidth\ \left(in\ \frac{B}{S}\right) = \frac{Bytes \times frequency_{DPU}}{\#instructions}$$

# WRAM Bandwidth: COPY



**COPY** executes **2 instructions** (WRAM load and store).  
With 11 tasklets, **11 × 16 bytes** in **22 cycles**:

$$WRAM \text{ Bandwidth } \left( in \frac{B}{S} \right) = 2,800 \frac{MB}{s} \text{ at } 350 \text{ MHz}$$

# WRAM Bandwidth: Access Patterns

- All 8-byte WRAM loads and stores take one cycle when the DPU pipeline is full

## KEY OBSERVATION 3

### Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

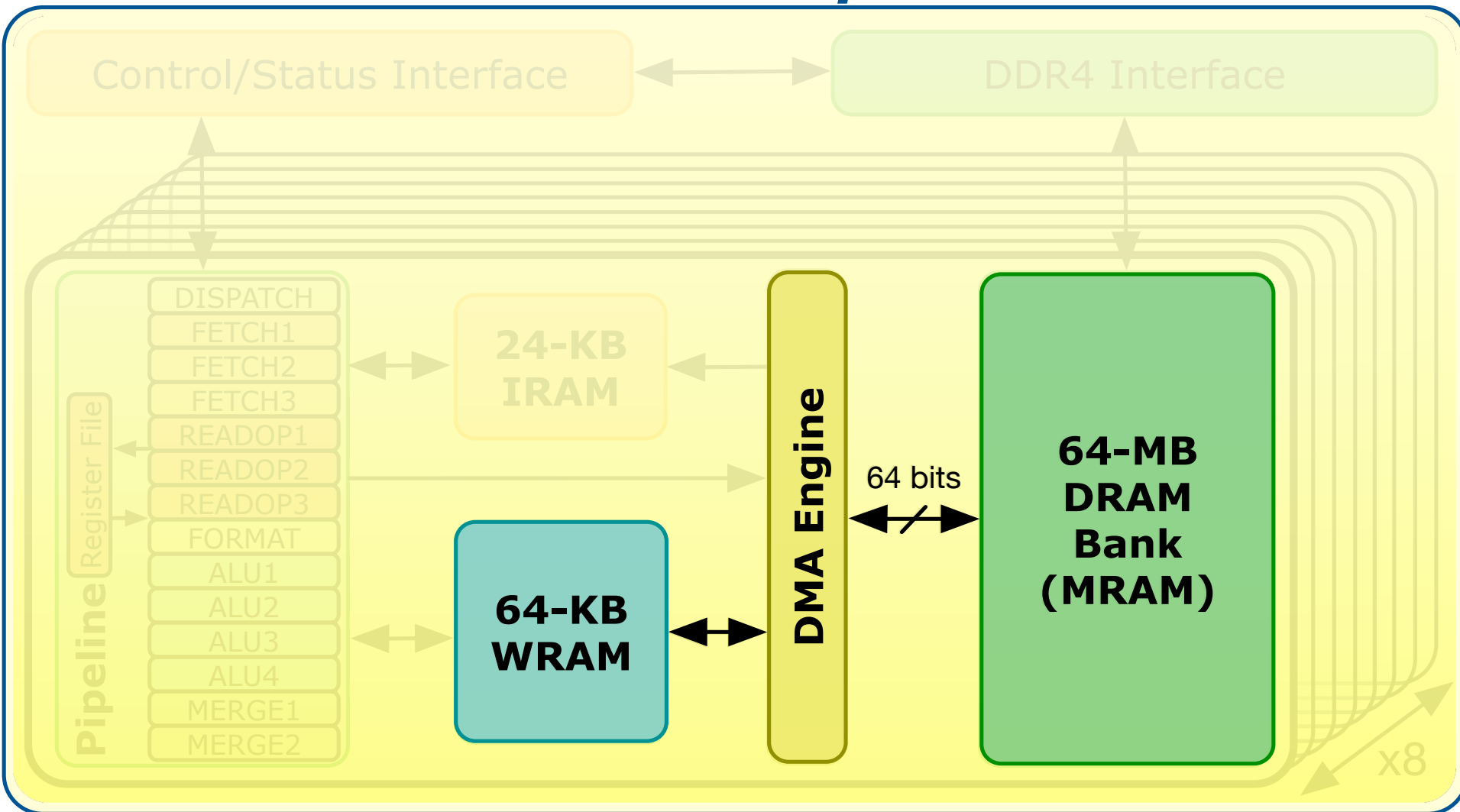
Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich   <sup>2</sup>American University of Beirut   <sup>3</sup>University of Malaga   <sup>4</sup>National Technical University of Athens

- Microbenchmark: `c[a[i]] = b[a[i]];`
  - Unit-stride: `a[i] = a[i-1] + 1;`
  - Si <https://arxiv.org/pdf/2105.03814.pdf>
  - R <https://github.com/CMU-SAFARI/prim-benchmarks>

# DPU: MRAM Latency and Bandwidth

## *PIM Chip*



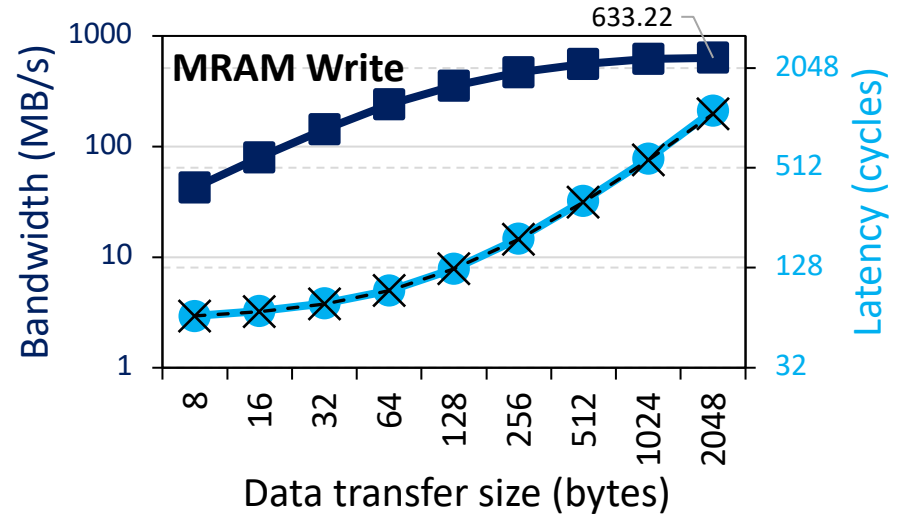
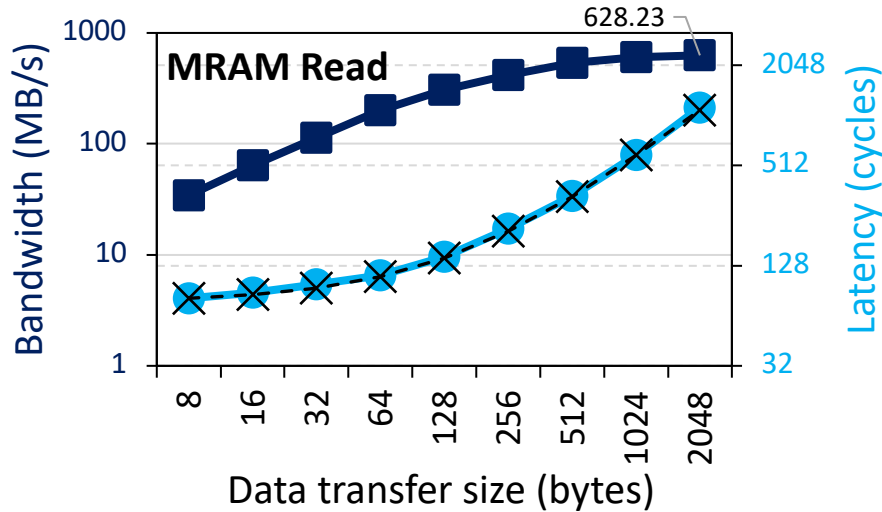


# MRAM Bandwidth

---

- Goal
  - Measure **MRAM bandwidth** for different access patterns
- Microbenchmarks
  - Latency of a single DMA transfer for different transfer sizes
    - `mram_read();` // MRAM-WRAM DMA transfer
    - `mram_write();` // WRAM-MRAM DMA transfer
  - **STREAM** benchmark
    - COPY, COPY-DMA
    - ADD, SCALE, TRIAD
  - **Strided** access pattern
    - Coarse-grain strided access
    - Fine-grain strided access
  - **Random** access pattern (GUPS)
- We do include accesses to MRAM

# MRAM Read and Write Latency (I)



$$MRAM \text{ Bandwidth } \left( \text{in } \frac{B}{S} \right) = \frac{\text{size} \times \text{frequency}_{DPU}}{MRAM \text{ Latency}}$$

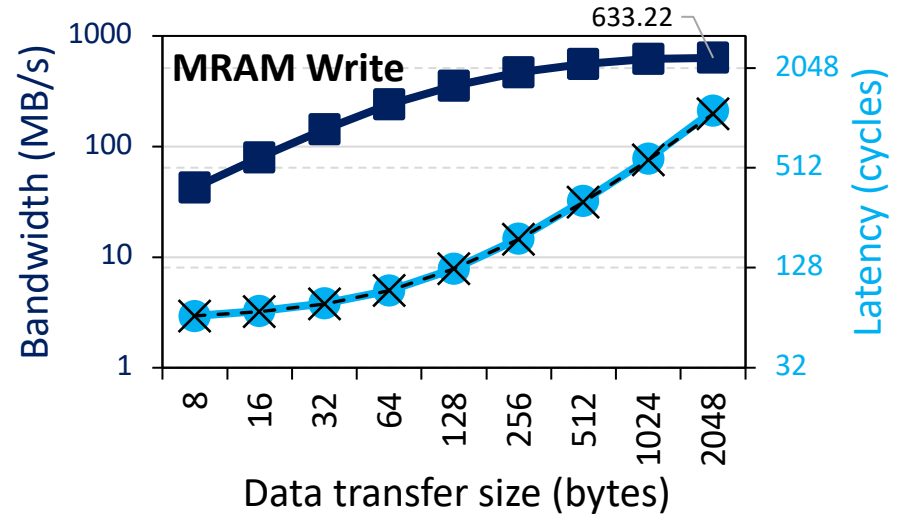
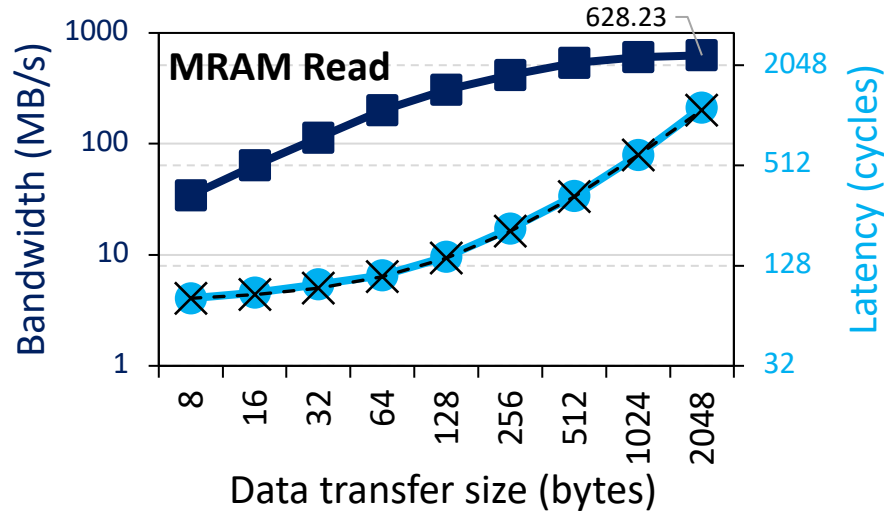
We can model the MRAM latency with a linear expression

$$MRAM \text{ Latency (in cycles)} = \alpha + \beta \times \text{size}$$

In our measurements,  $\beta$  equals 0.5 cycles/byte.

Theoretical maximum MRAM bandwidth = 700 MB/s at 350 MHz

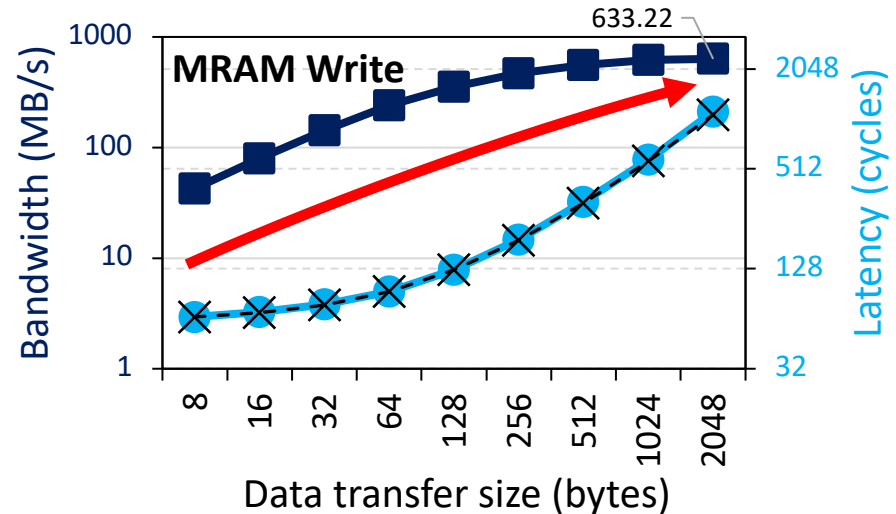
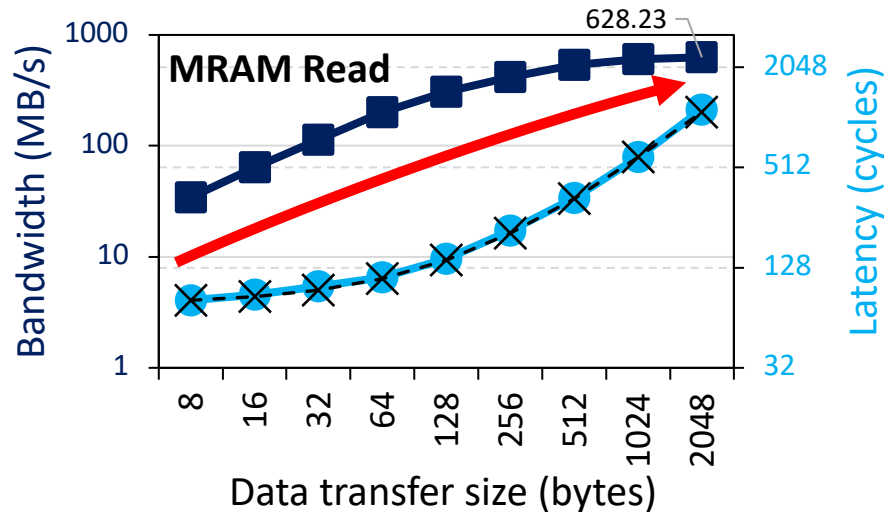
# MRAM Read and Write Latency (II)



## KEY OBSERVATION 4

- The DPU's Main memory (MRAM) bank access latency increases **linearly** with the transfer size.
- The maximum theoretical MRAM **bandwidth** is 2 bytes per cycle.

# MRAM Read and Write Latency (III)



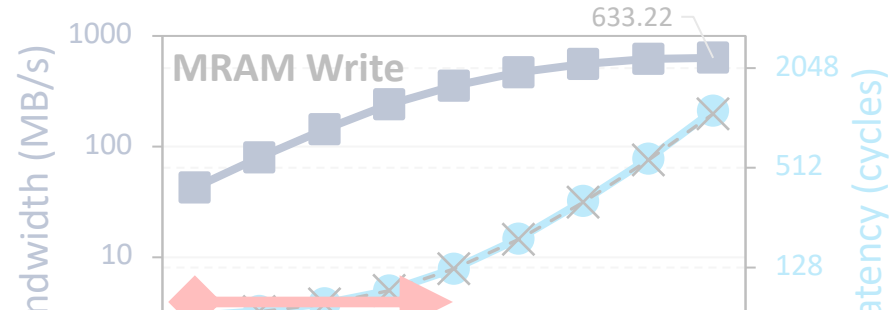
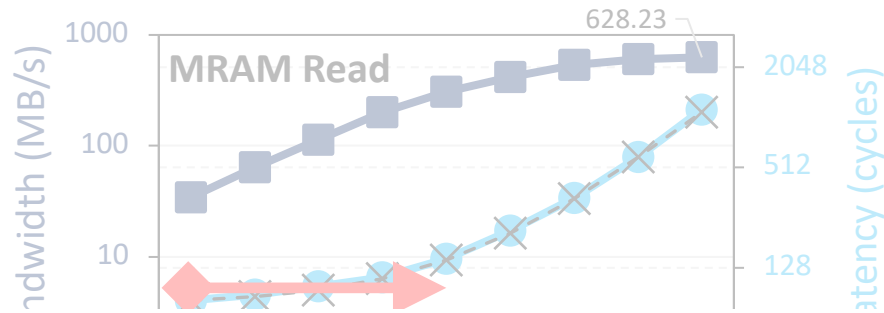
Read and write accesses to MRAM are symmetric

The sustained MRAM bandwidth increases with data transfer size

## **PROGRAMMING RECOMMENDATION 1**

For data movement between the DPU's MRAM bank and the WRAM, **use large DMA transfer sizes when all the accessed data is going to be used.**

# MRAM Read and Write Latency (IV)



## Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>American University of Beirut <sup>3</sup>University of Malaga <sup>4</sup>National Technical University of Athens  
desired data is in WRAM before issuing a new MRAM access).

### PROGRAMMING RECOMMENDATION 3

Choose the data transfer size between the MRAM bank and the WRAM based on the <https://arxiv.org/pdf/2105.03814.pdf> maintained  
MRAM dictated <https://github.com/CMU-SAFARI/prim-benchmarks> which is

# MRAM Bandwidth

---

- Goal
  - Measure **MRAM bandwidth** for different access patterns
- Microbenchmarks
  - **Latency of a single DMA transfer** for different transfer sizes
    - `mram_read();` // MRAM-WRAM DMA transfer
    - `mram write();` // WRAM-MRAM DMA transfer
  - **STREAM** benchmark
    - COPY, COPY-DMA
    - ADD, SCALE, TRIAD
  - **Strided** access pattern
    - Coarse-grain strided access
    - Fine-grain strided access
  - **Random** access pattern (GUPS)
- We do include accesses to MRAM

# STREAM Benchmark in MRAM

---

```
// COPY
// Load current MRAM block to WRAM
mram_read((__mram_ptr void const*)mram_address_A, bufferA,
          SIZE * sizeof(uint64_t));

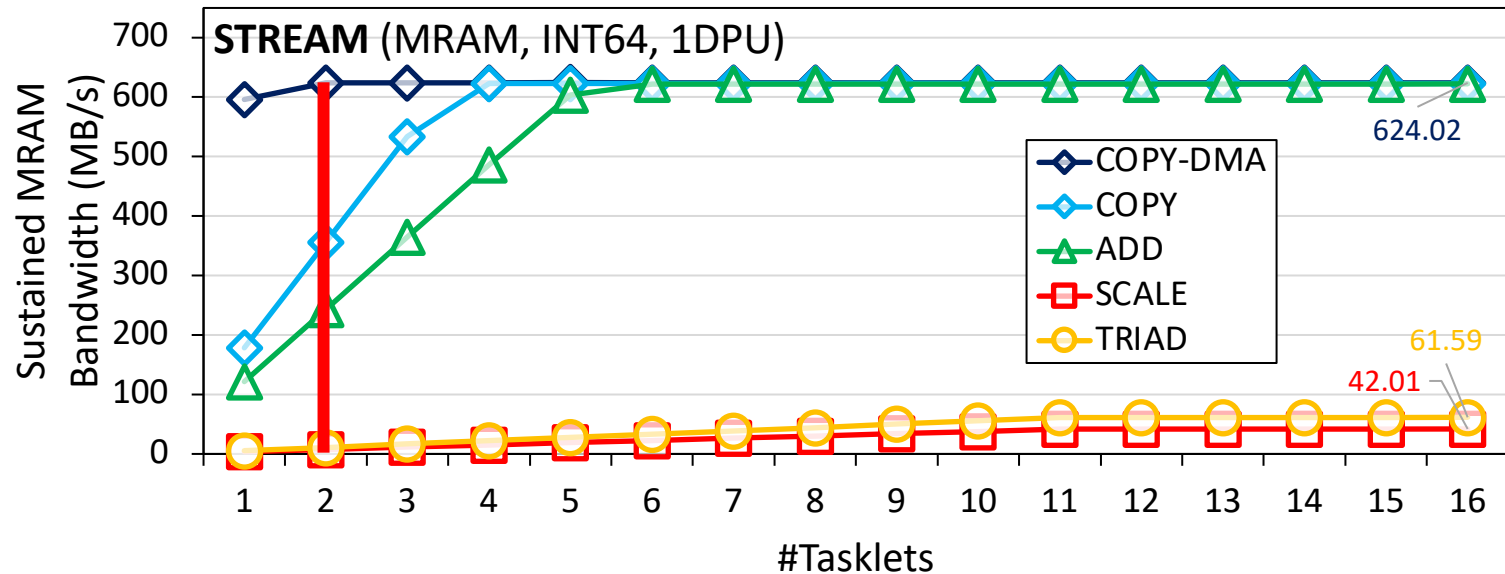
for(int i = 0; i < SIZE; i++){
    bufferB[i] = bufferA[i];
}

// Write WRAM block to MRAM
mram_write(bufferB, (__mram_ptr void*)mram_address_B,
           SIZE * sizeof(uint64_t));

// COPY-DMA
// Load current MRAM block to WRAM
mram_read((__mram_ptr void const*)mram_address_A, bufferA,
          SIZE * sizeof(uint64_t));

// Write WRAM block to MRAM
mram_write(bufferB, (__mram_ptr void*)mram_address_B,
           SIZE * sizeof(uint64_t));
```

# STREAM Benchmark: COPY-DMA



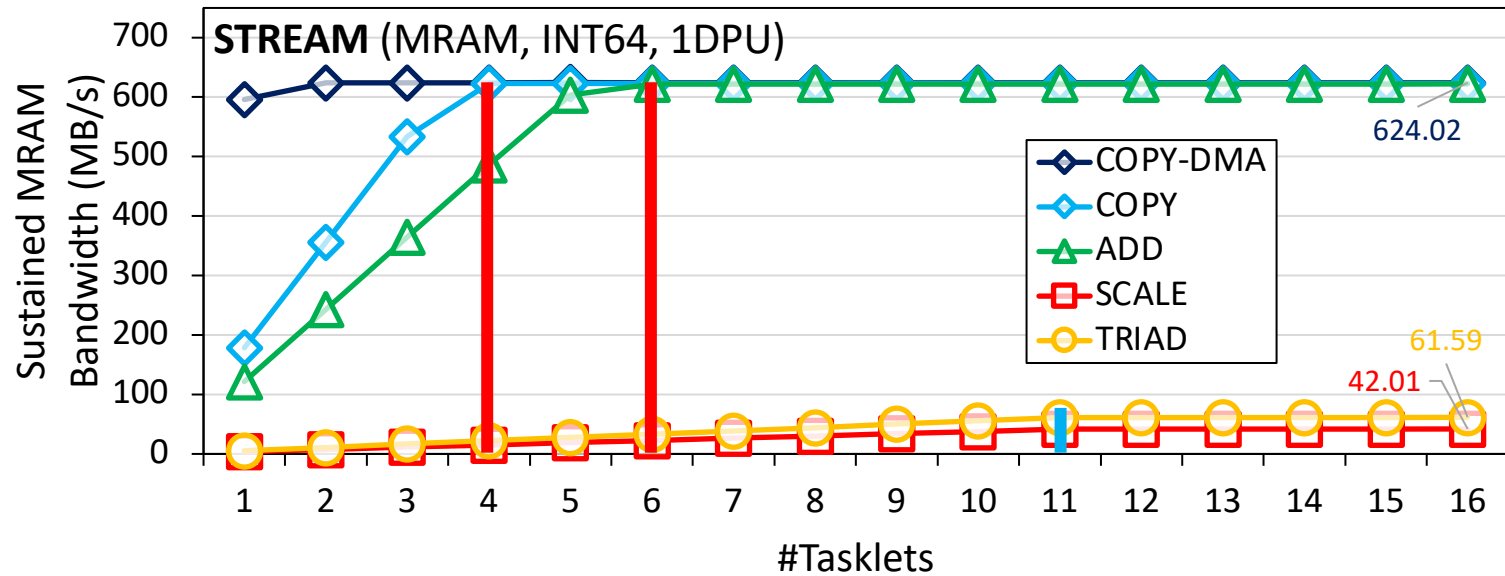
The sustained bandwidth of **COPY-DMA** is close to the theoretical maximum (700 MB/s): **~1.6 TB/s for 2,556 DPUs**

**COPY-DMA** saturates with **two tasklets**, even though the DMA engine can perform only one transfer at a time

Using **two or more tasklets** guarantees that there is always a DMA request enqueued to keep the DMA engine busy



# STREAM Benchmark: Bandwidth Saturation (I)



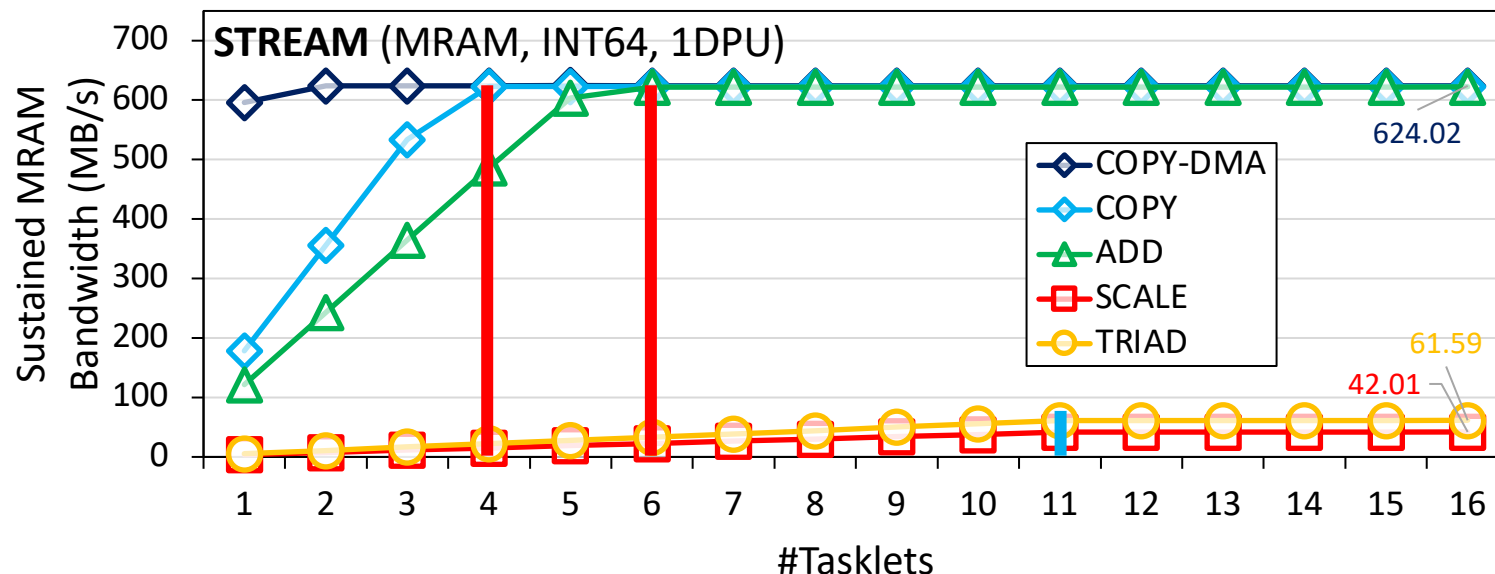
**COPY** and **ADD** saturate at 4 and 6 tasklets, respectively

**SCALE** and **TRIAD** saturate at 11 tasklets

The latency of MRAM accesses becomes longer than the pipeline latency after 4 and 6 tasklets for **COPY** and **ADD**, respectively

The pipeline latency of **SCALE** and **TRIAD** is longer than the MRAM latency for any number of tasklets (both use costly MUL)

# STREAM Benchmark: Bandwidth Saturation (II)



## KEY OBSERVATION 5

- **When the access latency to an MRAM bank for a streaming benchmark (COPY-DMA, COPY, ADD) is larger than the pipeline latency** (i.e., execution latency of arithmetic operations and WRAM accesses), the performance of the DPU saturates at a number of tasklets smaller than 11. **This is a memory-bound workload.**
- **When the pipeline latency for a streaming benchmark (SCALE, TRIAD) is larger than the MRAM access latency**, the performance of a DPU saturates at 11 tasklets. **This is a compute-bound workload.**

# MRAM Bandwidth

---

- Goal
  - Measure **MRAM bandwidth** for different access patterns
- Microbenchmarks
  - **Latency of a single DMA transfer** for different transfer sizes
    - `mram_read();` // MRAM-WRAM DMA transfer
    - `mram_write();` // WRAM-MRAM DMA transfer
  - **STREAM** benchmark
    - COPY, COPY-DMA
    - ADD, SCALE, TRIAD
  - **Strided** access pattern
    - Coarse-grain strided access
    - Fine-grain strided access
  - **Random** access pattern (GUPS)
- We do include accesses to MRAM

# Strided and Random Access to MRAM

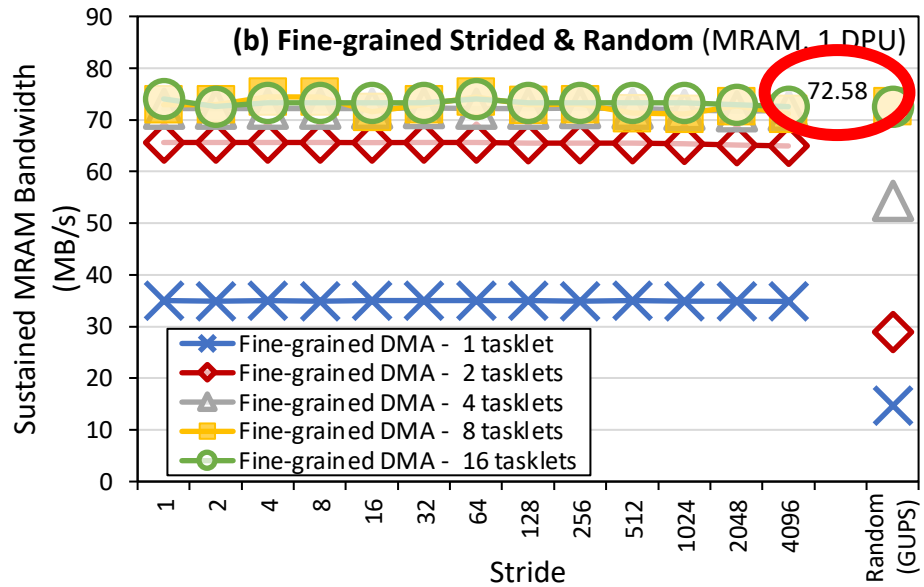
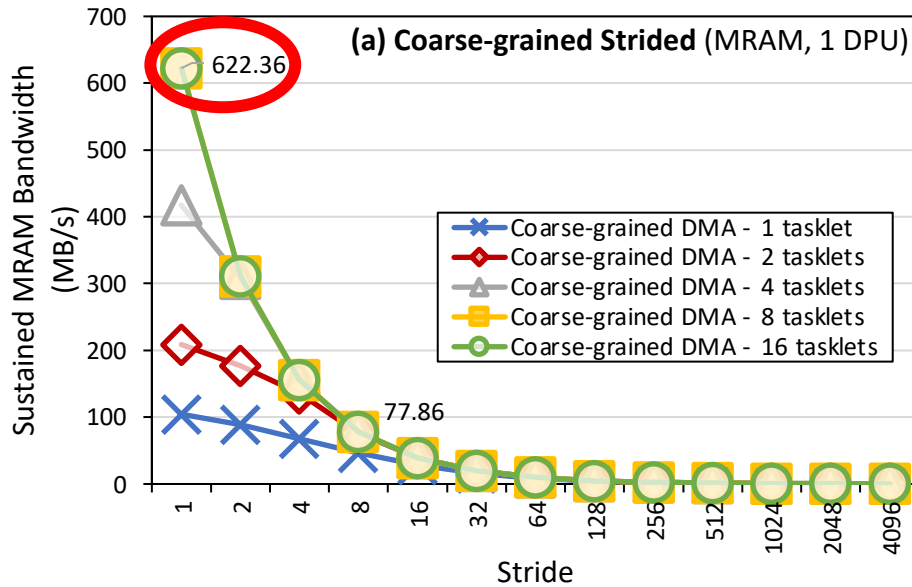
```
// COARSE-GRAINED STRIDED ACCESS
// Load current MRAM block to WRAM
mram_read((__mram_ptr void const*)mram_address_A, bufferA,
          SIZE * sizeof(uint64_t));
mram_read((__mram_ptr void const*)mram_address_B, bufferB,
          SIZE * sizeof(uint64_t));

for(int i = 0; i < SIZE; i += stride){
    bufferB[i] = bufferA[i];
}
// Write WRAM block to MRAM
mram_write(bufferB, (__mram_ptr void*)mram_address_B,
          SIZE * sizeof(uint64_t));

// FINE-GRAINED STRIDED & RANDOM ACCESS
for(int i = 0; i < SIZE; i += stride){
    int index = i * sizeof(uint64_t);
    // Load current MRAM element to WRAM
    mram_read((__mram_ptr void const*)(mram_address_A + index), bufferA,
              sizeof(uint64_t));

    // Write WRAM element to MRAM
    mram_write(bufferA, (__mram_ptr void*)(mram_address_B + index),
              sizeof(uint64_t));
}
```

# Strided and Random Accesses (I)

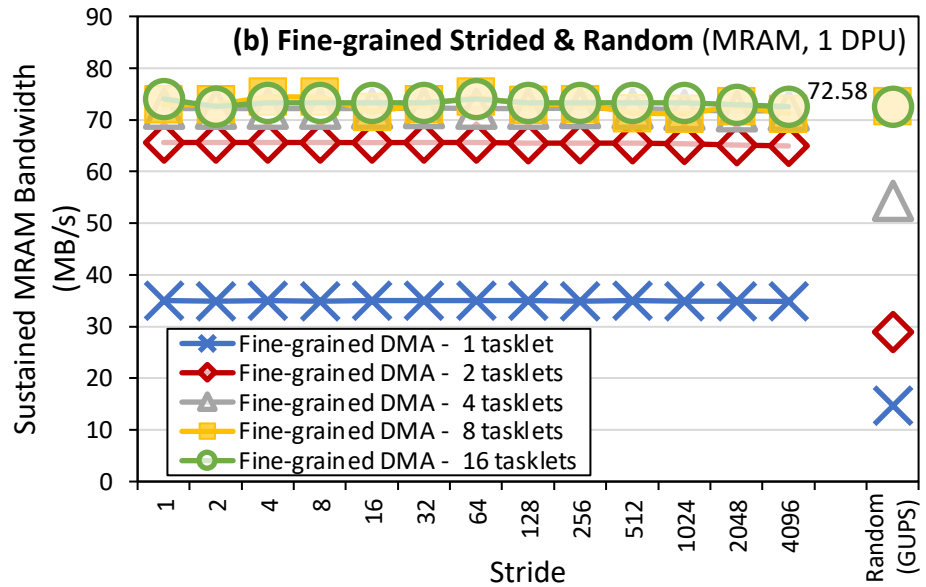
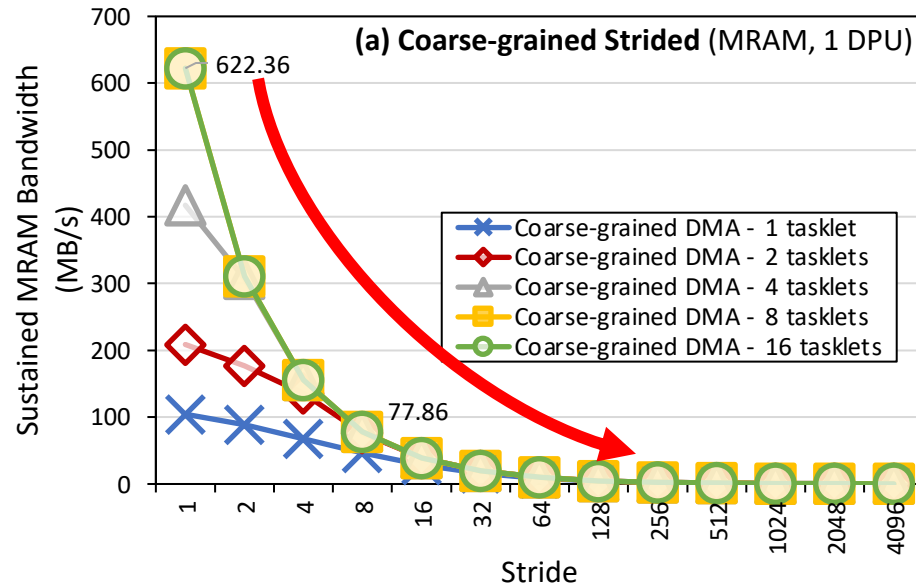


Large difference in maximum sustained bandwidth between coarse-grained and fine-grained DMA

Coarse-grained DMA uses 1,024-byte transfers, while fine-grained DMA uses 8-byte transfers

Random access achieves very similar maximum sustained bandwidth to fine-grained strided approach

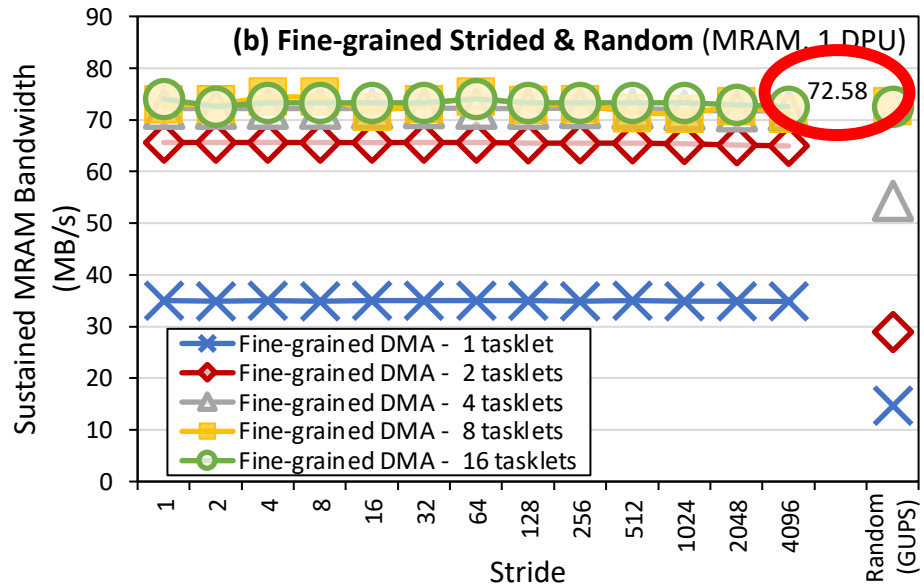
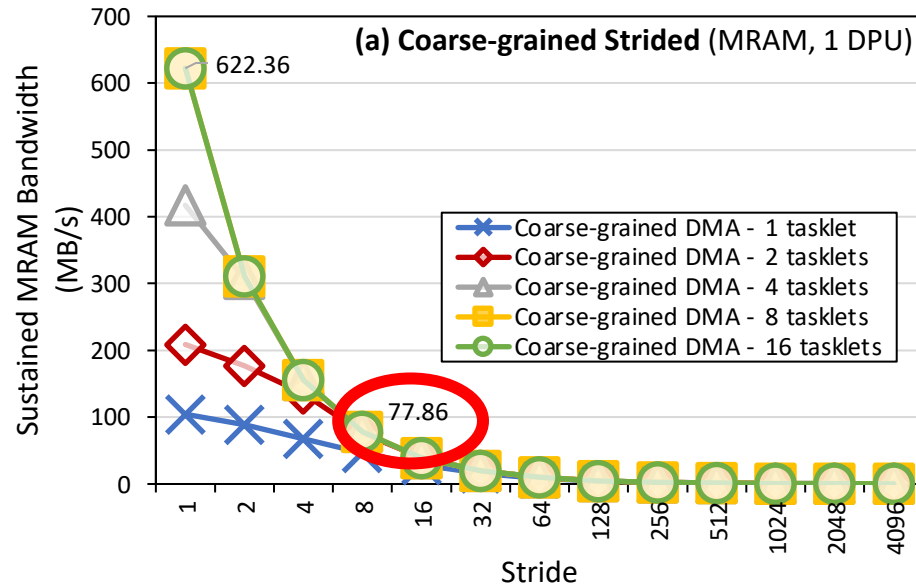
# Strided and Random Accesses (II)



The sustained MRAM bandwidth of coarse-grained DMA decreases as the stride increases

The effective utilization of the transferred data decreases as the stride becomes larger (e.g., a stride 4 means that only one fourth of the transferred data is used)

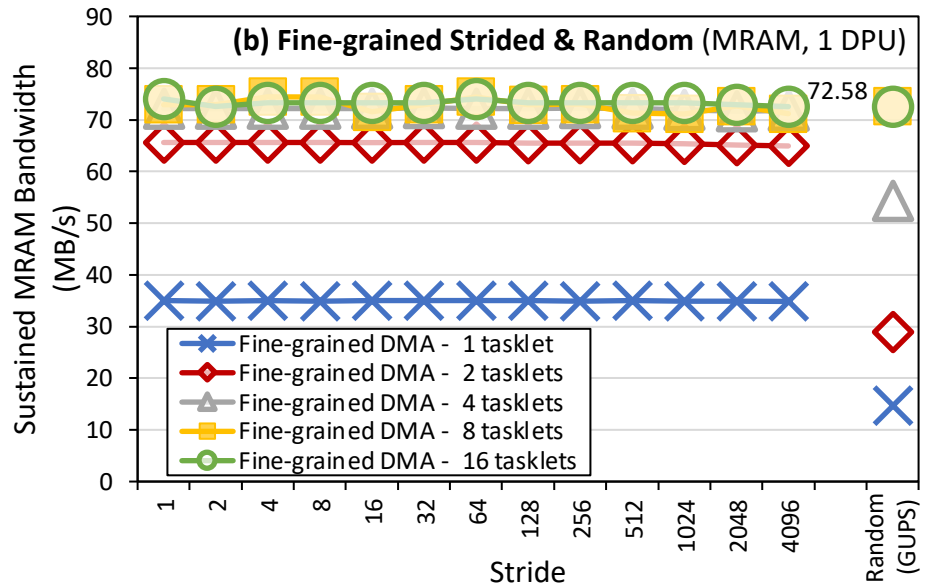
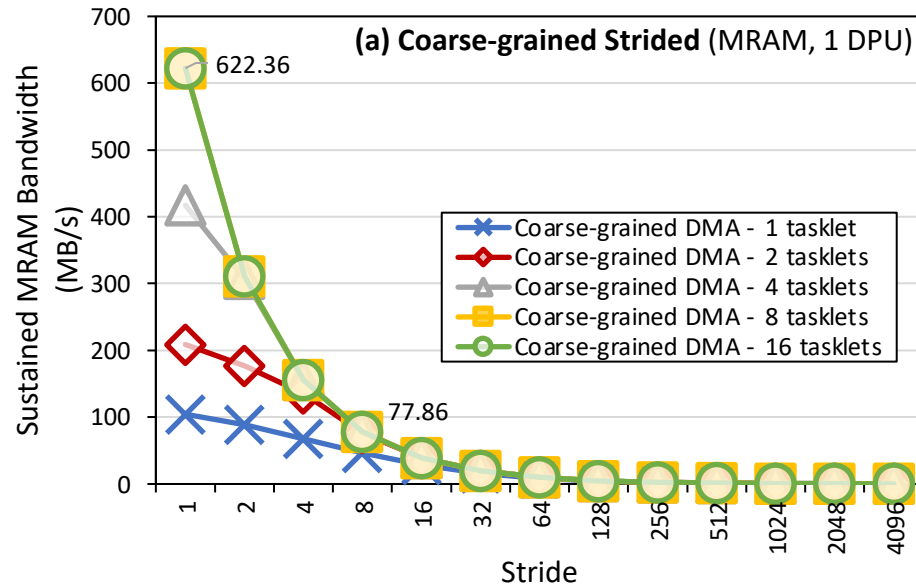
# Strided and Random Accesses (III)



For a stride of 16 or larger, the fine-grained DMA approach achieves higher bandwidth

With stride 16, only one sixteenth of the maximum sustained bandwidth (622.36 MB/s) of coarse-grained DMA is effectively used, which is lower than the bandwidth of fine-grained DMA (72.58 MB/s)

# Strided and Random Accesses (IV)



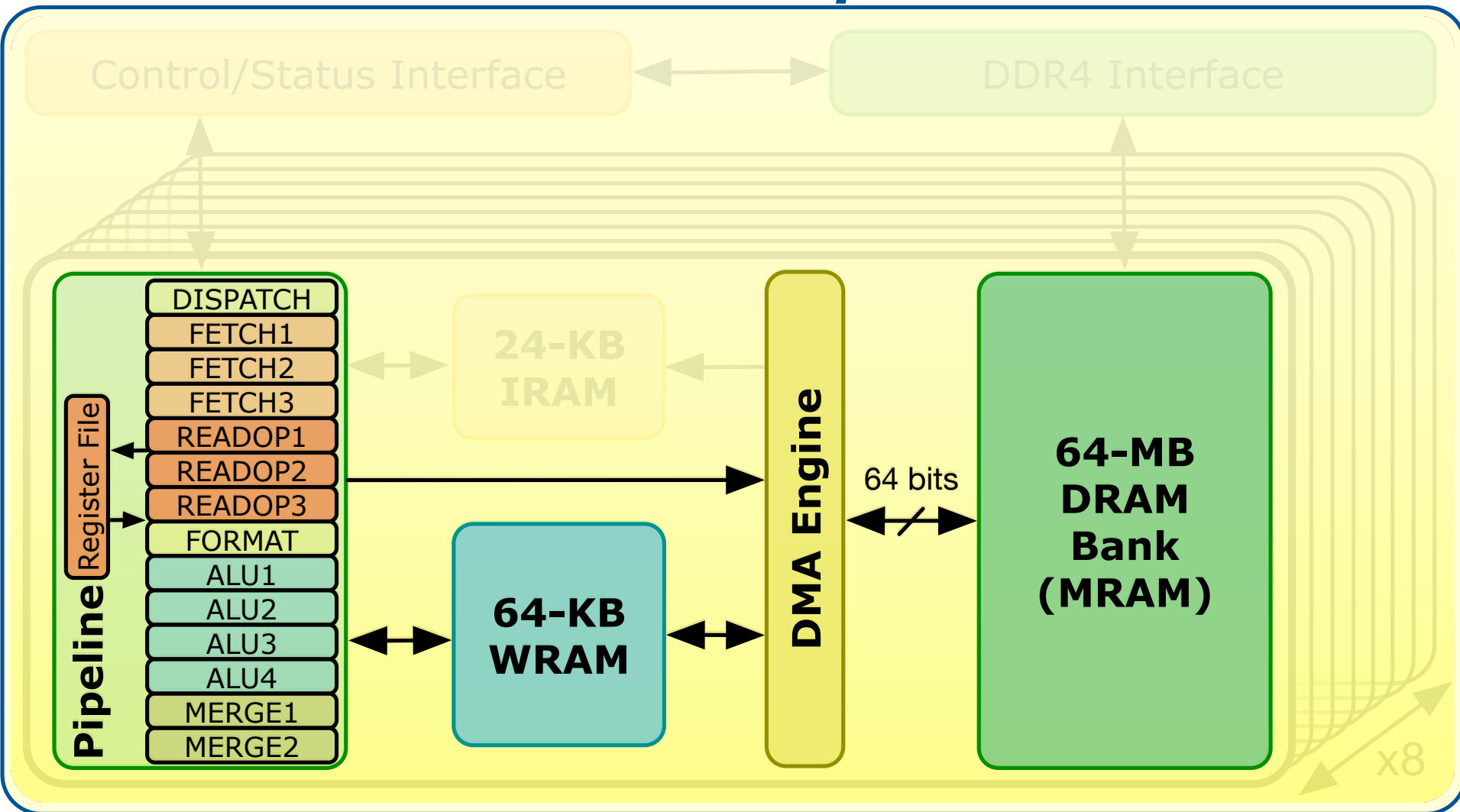
## PROGRAMMING RECOMMENDATION 4

- For strided access patterns with a **stride smaller than 16 8-byte elements**, **fetch a large contiguous chunk** (e.g., 1,024 bytes) from a DPU's MRAM bank.
- For strided access patterns with **larger strides and random access patterns**, **fetch only the data elements that are needed** from an MRAM bank.



# DPU: Arithmetic Throughput vs. Operational Intensity

## *PIM Chip*



# Arithmetic Throughput vs. Operational Intensity (I)

---

- Goal
  - Characterize **memory-bound regions** and **compute-bound regions** for different datatypes and operations
- Microbenchmark
  - We **load one chunk of an MRAM array into WRAM**
  - Perform a **variable number of operations on the data**
  - **Write back** to MRAM
- The experiment is inspired by the **Roofline model**\*
- We define **operational intensity** (OI) as the number of arithmetic operations performed per byte accessed from MRAM (OP/B)
- The pipeline latency changes with the operational intensity, but the MRAM access latency is fixed

# Arithmetic Throughput vs. Operational Intensity (II)

```
int repetitions = input_repeat >= 1.0 ? (int)input_repeat : 1;
int stride      = input_repeat >= 1.0 ? 1 : (int)(1 / input_repeat);

// Load current MRAM block to WRAM
mram_read((__mram_ptr void const*)mram_address_A, bufferA, SIZE * sizeof(T));

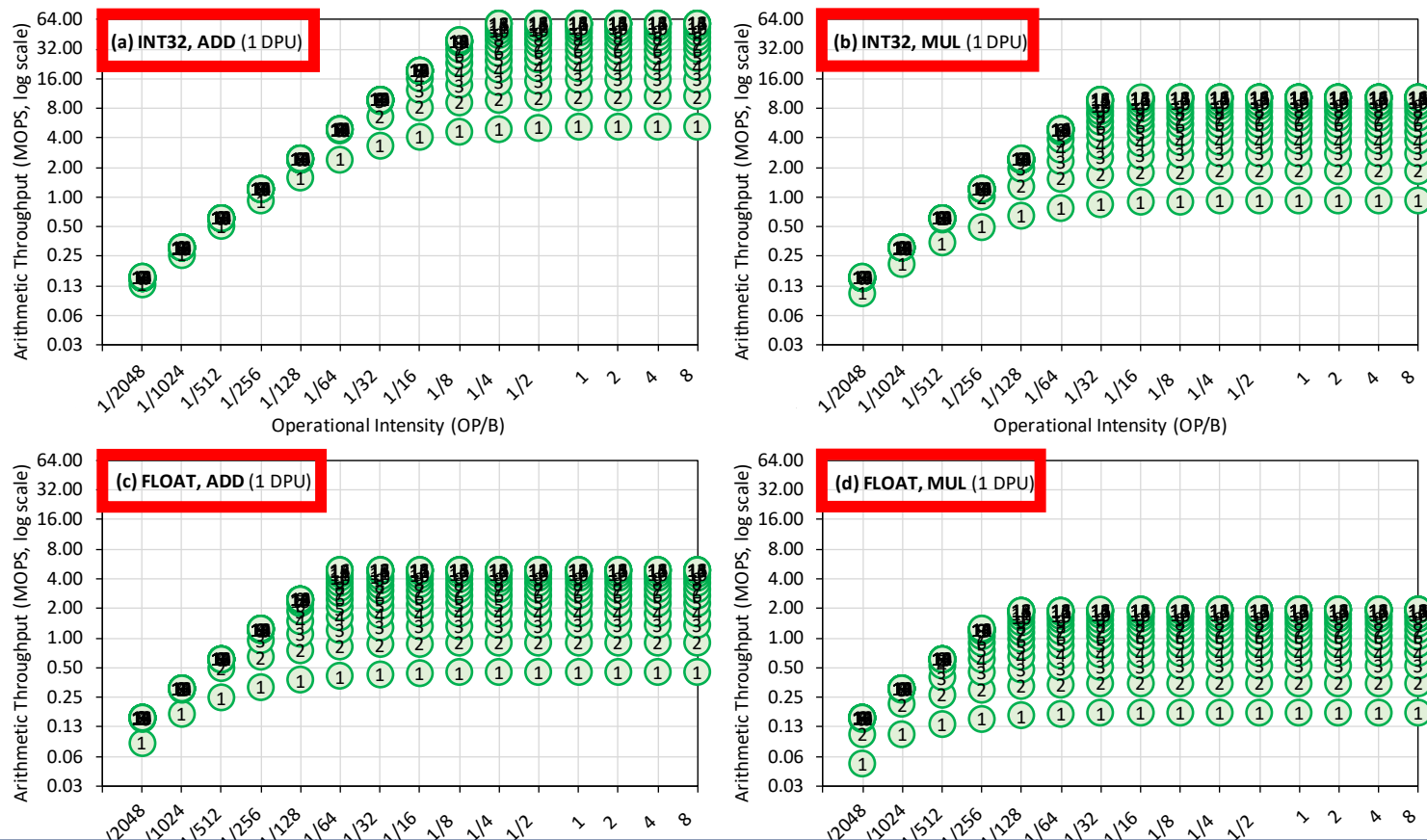
// Update
for(int r = 0; r < repetitions; r++){
    for(int i = 0; i < SIZE; i+=stride){
#ifdef ADD
        bufferA[i] += scalar; // ADD
#elif SUB
        bufferA[i] -= scalar; // SUB
#elif MUL
        bufferA[i] *= scalar; // MUL
#elif DIV
        bufferA[i] /= scalar; // DIV
#endif
    }
}

// Write WRAM block to MRAM
mram_write(bufferA, (__mram_ptr void*)mram_address_B, SIZE * sizeof(T));
```

input\_repeat greater or equal to 1 indicates the (integer) number of repetitions per input element

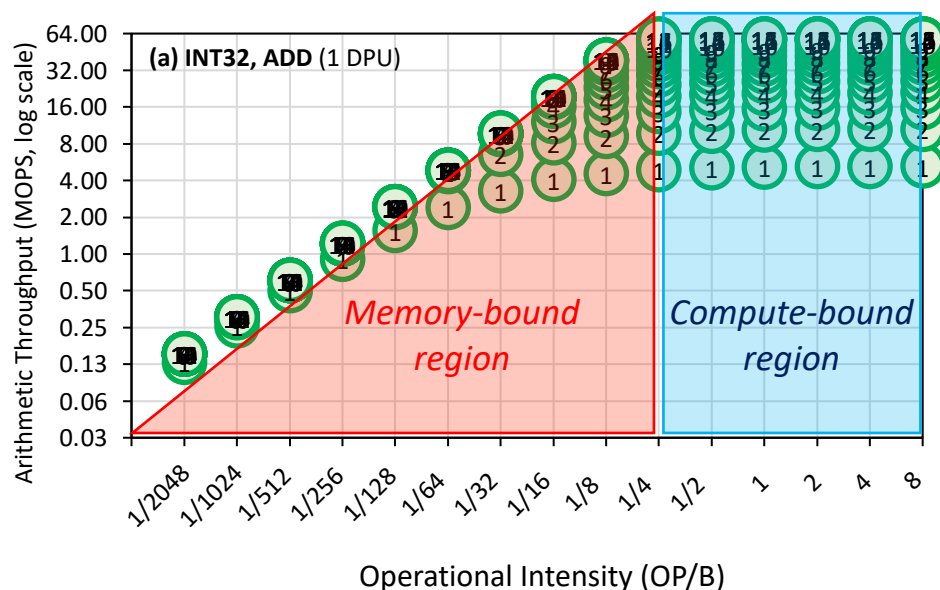
input\_repeat smaller than 1 indicates the fraction of elements that are updated

# Arithmetic Throughput vs. Operational Intensity (III)



We show results of arithmetic throughput vs. operational intensity for (a) 32-bit integer ADD, (b) 32-bit integer MUL, (c) 32-bit floating-point ADD, and (d) 32-bit floating-point MUL (results for other datatypes and operations show similar trends)

# Arithmetic Throughput vs. Operational Intensity (IV)



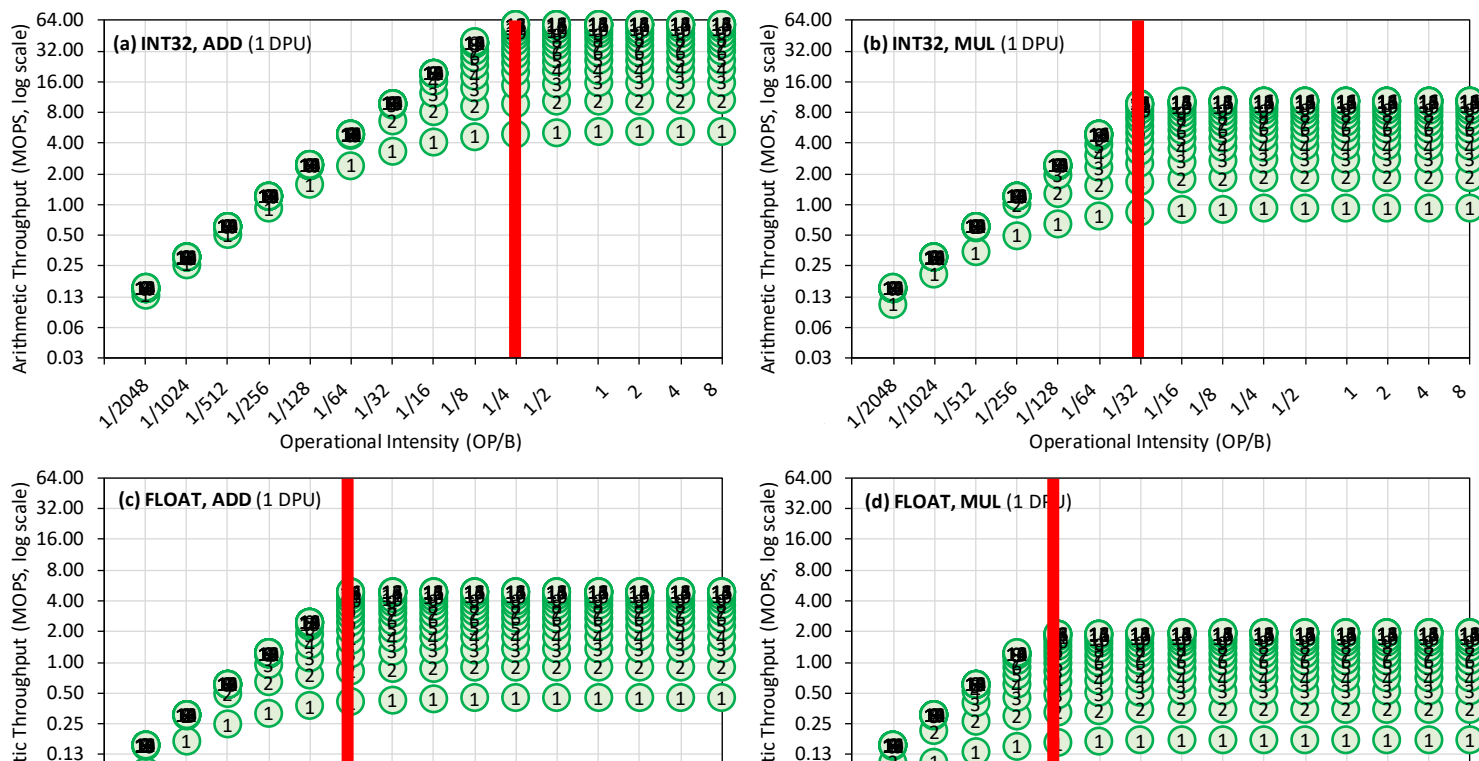
In the **memory-bound region**, the arithmetic throughput increases with the operational intensity

In the **compute-bound region**, the arithmetic throughput is flat at its maximum

The **throughput saturation point** is the operational intensity where the transition between the memory-bound region and the compute-bound region happens

The throughput saturation point is as low as  $\frac{1}{4}$  OP/B, i.e., **1 integer addition per every 32-bit element** fetched

# Arithmetic Throughput vs. Operational Intensity (V)



## KEY OBSERVATION 6

The arithmetic throughput of a DRAM Processing Unit (DPU) saturates at low or very low operational intensity (e.g., 1 integer addition per 32-bit element). Thus, the DPU is fundamentally a compute-bound processor. We expect most real-world workloads be compute-bound in the UPMEM PIM architecture.

# Outline

---

- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PRIM Benchmarks
  - Roofline Model
  - Benchmark Diversity
- Evaluation
  - Strong and Weak Scaling
  - Comparison to CPU and GPU
- Key Takeaways

# PrIM Benchmarks

---

- Goal
  - A **common set of workloads** that can be used to
    - evaluate the UPMEM PIM architecture,
    - compare software improvements and compilers,
    - compare future PIM architectures and hardware
- Two key selection criteria:
  - Selected workloads from **different application domains**
  - **Memory-bound workloads** on processor-centric architectures
- 14 different workloads, 16 different benchmarks\*

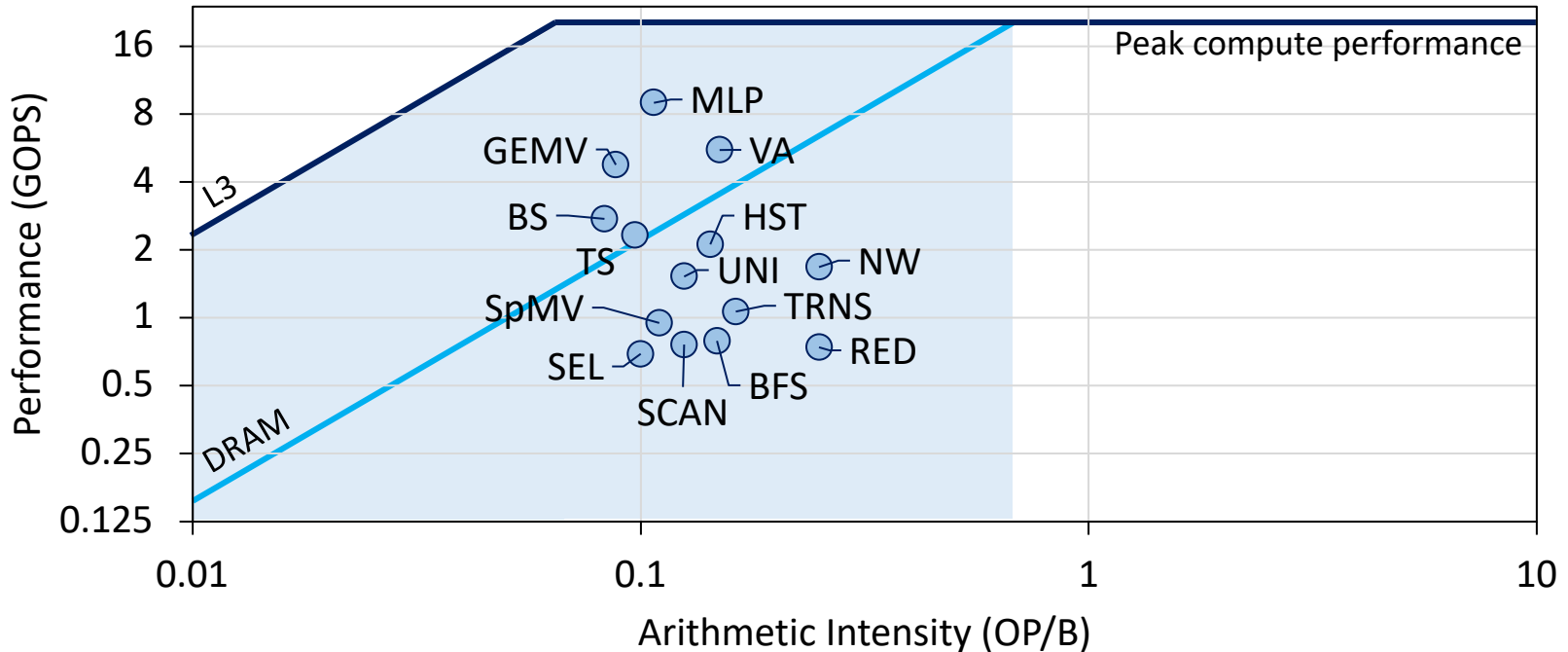


# PrIM Benchmarks: Application Domains

Domain	Benchmark	Short name
Dense linear algebra	Vector Addition	VA
	Matrix-Vector Multiply	GEMV
Sparse linear algebra	Sparse Matrix-Vector Multiply	SpMV
Databases	Select	SEL
	Unique	UNI
Data analytics	Binary Search	BS
	Time Series Analysis	TS
Graph processing	Breadth-First Search	BFS
Neural networks	Multilayer Perceptron	MLP
Bioinformatics	Needleman-Wunsch	NW
Image processing	Image histogram (short)	HST-S
	Image histogram (large)	HST-L
Parallel primitives	Reduction	RED
	Prefix sum (scan-scan-add)	SCAN-SSA
	Prefix sum (reduce-scan-scan)	SCAN-RSS
	Matrix transposition	TRNS

# Roofline Model

- Intel Advisor on an Intel Xeon E3-1225 v6 CPU



All workloads fall in the **memory-bound area of the Roofline**

# PrIM Benchmarks: Diversity

- PrIM benchmarks are diverse:
  - Memory access patterns
  - Operations and datatypes
  - Communication/synchronization

Domain	Benchmark	Short name	Memory access pattern			Computation pattern		Communication/synchronization	
			Sequential	Strided	Random	Operations	Datatype	Intra-DPU	Inter-DPU
Dense linear algebra	Vector Addition	VA	Yes			add	int32_t		
	Matrix-Vector Multiply	GEMV	Yes			add, mul	uint32_t		
Sparse linear algebra	Sparse Matrix-Vector Multiply	SpMV	Yes		Yes	add, mul	float		
Databases	Select	SEL	Yes			add, compare	int64_t	handshake, barrier	Yes
	Unique	UNI	Yes			add, compare	int64_t	handshake, barrier	Yes
Data analytics	Binary Search	BS	Yes		Yes	compare	int64_t		
	Time Series Analysis	TS	Yes			add, sub, mul, div	int32_t		
Graph processing	Breadth-First Search	BFS	Yes		Yes	bitwise logic	uint64_t	barrier, mutex	Yes
Neural networks	Multilayer Perceptron	MLP	Yes			add, mul, compare	int32_t		
Bioinformatics	Needleman-Wunsch	NW	Yes	Yes		add, sub, compare	int32_t	barrier	Yes
Image processing	Image histogram (short)	HST-S	Yes		Yes	add	uint32_t	barrier	Yes
	Image histogram (long)	HST-L	Yes		Yes	add	uint32_t	barrier, mutex	Yes
Parallel primitives	Reduction	RED	Yes	Yes		add	int64_t	barrier	Yes
	Prefix sum (scan-scan-add)	SCAN-SSA	Yes			add	int64_t	handshake, barrier	Yes
	Prefix sum (reduce-scan-scan)	SCAN-RSS	Yes			add	int64_t	handshake, barrier	Yes
	Matrix transposition	TRNS	Yes		Yes	add, sub, mul	int64_t	mutex	

# PrIM Benchmarks: Inter-DPU Communication

Domain	Benchmark	Short name	Memory access pattern			Computation pattern		Communication/synchronization	
			Sequential	Strided	Random	Operations	Datatype	Intra-DPU	Inter-DPU
Dense linear algebra	Vector Addition	VA	Yes			add	int32_t		
	Matrix-Vector Multiply	GEMV	Yes			add, mul	uint32_t		
Sparse linear algebra	Sparse Matrix-Vector Multiply	SpMV	Yes		Yes	add, mul	float		
Databases	Select	SEL	Yes			add, compare	int64_t	handshake, barrier	Yes
	Unique	UNI	Yes			add, compare	int64_t	handshake, barrier	Yes
Data analytics	Binary Search	BS	Yes		Yes	compare	int64_t		
	Time Series Analysis	TS	Yes			add, sub, mul, div	int32_t		
Graph processing	Breadth-First Search	BFS	Yes		Yes	bitwise logic	uint64_t	barrier, mutex	Yes
Neural networks	Multilayer Perceptron	MLP	Yes			add, mul, compare	int32_t		
Bioinformatics	Needleman-Wunsch	NW	Yes	Yes		add, sub, compare	int32_t	barrier	Yes
	Image histogram (short)	HST-S	Yes		Yes	add	uint32_t	barrier	Yes
Image processing	Image histogram (long)	HST-L	Yes		Yes	add	uint32_t	barrier, mutex	Yes
	Reduction	RED	Yes	Yes		add	int64_t	barrier	Yes
Parallel primitives	Prefix sum (scan-scan-add)	SCAN-SSA	Yes			add	int64_t	handshake, barrier	Yes
	Prefix sum (reduce-scan-scan)	SCAN-RSS	Yes			add	int64_t	handshake, barrier	Yes
	Matrix transposition	TRANS	Yes		Yes	add, sub, mul	int64_t	mutex	

• Inter-DPU communication

- Result merging:

• SEL, UNI, HST-S, HST-L, RED

• Only DPU-CPU transfers

- Redistribution of intermediate results:

• BFS, MLP, NW, SCAN-SSA, SCAN-RSS

• DPU-CPU and CPU-DPU transfers

# Outline

---

- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PRIM Benchmarks
  - Roofline Model
  - Benchmark Diversity
- Evaluation
  - Strong and Weak Scaling
  - Comparison to CPU and GPU
- Key Takeaways

# Evaluation Methodology

---

- We evaluate the 16 PrIM benchmarks on two UPMEM-based systems:
  - 2,556-DPU system
  - 640-DPU system
- Strong and weak scaling experiments on the 2,556-DPU system
  - 1 DPU with different numbers of tasklets
  - 1 rank (strong and weak)
  - Up to 32 ranks

*Strong scaling* refers to how the execution time of a program solving a particular problem varies with the number of processors for a fixed problem size

*Weak scaling* refers to how the execution time of a program solving a particular problem varies with the number of processors for a fixed problem size per processor

# Evaluation Methodology

---

- We evaluate the **16 PrIM benchmarks** on two **UPMEM-based systems**:
  - 2,556-DPU system
  - 640-DPU system
- **Strong and weak scaling experiments** on the 2,556-DPU system
  - **1 DPU** with different numbers of tasklets
  - **1 rank** (strong and weak)
  - Up to **32 ranks**
- Comparison of both UPMEM-based PIM systems to **state-of-the-art CPU and GPU**
  - Intel Xeon E3-1240 CPU
  - NVIDIA Titan V GPU

# Datasets

- Strong and weak scaling experiments

Benchmark	Strong Scaling Dataset	Weak Scaling Dataset	MRAM-WRAM Transfer Sizes
VA	1 DPU-1 rank: 2.5M elem. (10 MB)   32 ranks: 160M elem. (640 MB)	2.5M elem./DPU (10 MB)	1024 bytes
GEMV	1 DPU-1 rank: 8192 × 1024 elem. (32 MB)   32 ranks: 163840 × 4096 elem. (2.56 GB)	1024 × 2048 elem./DPU (8 MB)	1024 bytes
SpMV	<i>bcsstk30</i> [253] (12 MB)	<i>bcsstk30</i> [253]	64 bytes
SEL	1 DPU-1 rank: 3.8M elem. (30 MB)   32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
UNI	1 DPU-1 rank: 3.8M elem. (30 MB)   32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
BS	2M elem. (16 MB). 1 DPU-1 rank: 256K queries. (2 MB)   32 ranks: 16M queries. (128 MB)	2M elem. (16 MB). 256K queries./DPU (2 MB).	8 bytes
TS	256 elem. query. 1 DPU-1 rank: 512K elem. (2 MB)   32 ranks: 32M elem. (128 MB)	512K elem./DPU (2 MB)	256 bytes
BFS	<i>loc-gowalla</i> [254] (22 MB)	<i>rMat</i> [255] ( $\approx 100K$ vertices and 1.2M edges per DPU)	8 bytes
MLP	3 fully-connected layers. 1 DPU-1 rank: 2K neurons (32 MB)   32 ranks: $\approx 160K$ neur. (2.56 GB)	3 fully-connected layers. 1K neur./DPU (4 MB)	1024 bytes
NW	1 DPU-1 rank: 2560 bps (50 MB), large/small sub-block = $\frac{2560}{\#DPU_s}/2$   32 ranks: 64K bps (32 GB), l./s. = 32/2	512 bps/DPU (2MB), l./s. = 512/2	8, 16, 32, 40 bytes
HST-S	1 DPU-1 rank: 1536 × 1024 input image [256] (6 MB)   32 ranks: 64 × input image	1536 × 1024 input image [256]/DPU (6 MB)	1024 bytes
HST-L	1 DPU-1 rank: 1536 × 1024 input image [256] (6 MB)   32 ranks: 64 × input image	1536 × 1024 input image [256]/DPU (6 MB)	1024 bytes
RED	1 DPU-1 rank: 6.3M elem. (50 MB)   32 ranks: 400M elem. (3.1 GB)	6.3M elem./DPU (50 MB)	1024 bytes
SCAN-SSA	1 DPU-1 rank: 3.8M elem. (30 MB)   32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
SCAN-RSS	1 DPU-1 rank: 3.8M elem. (30 MB)   32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
TRNS	1 DPU-1 rank: 12288 × 16 × 64 × 8 (768 MB)   32 ranks: 12288 × 16 × 2048 × 8 (24 GB)	12288 × 16 × 1 × 8/DPU (12 MB)	128, 1024 bytes

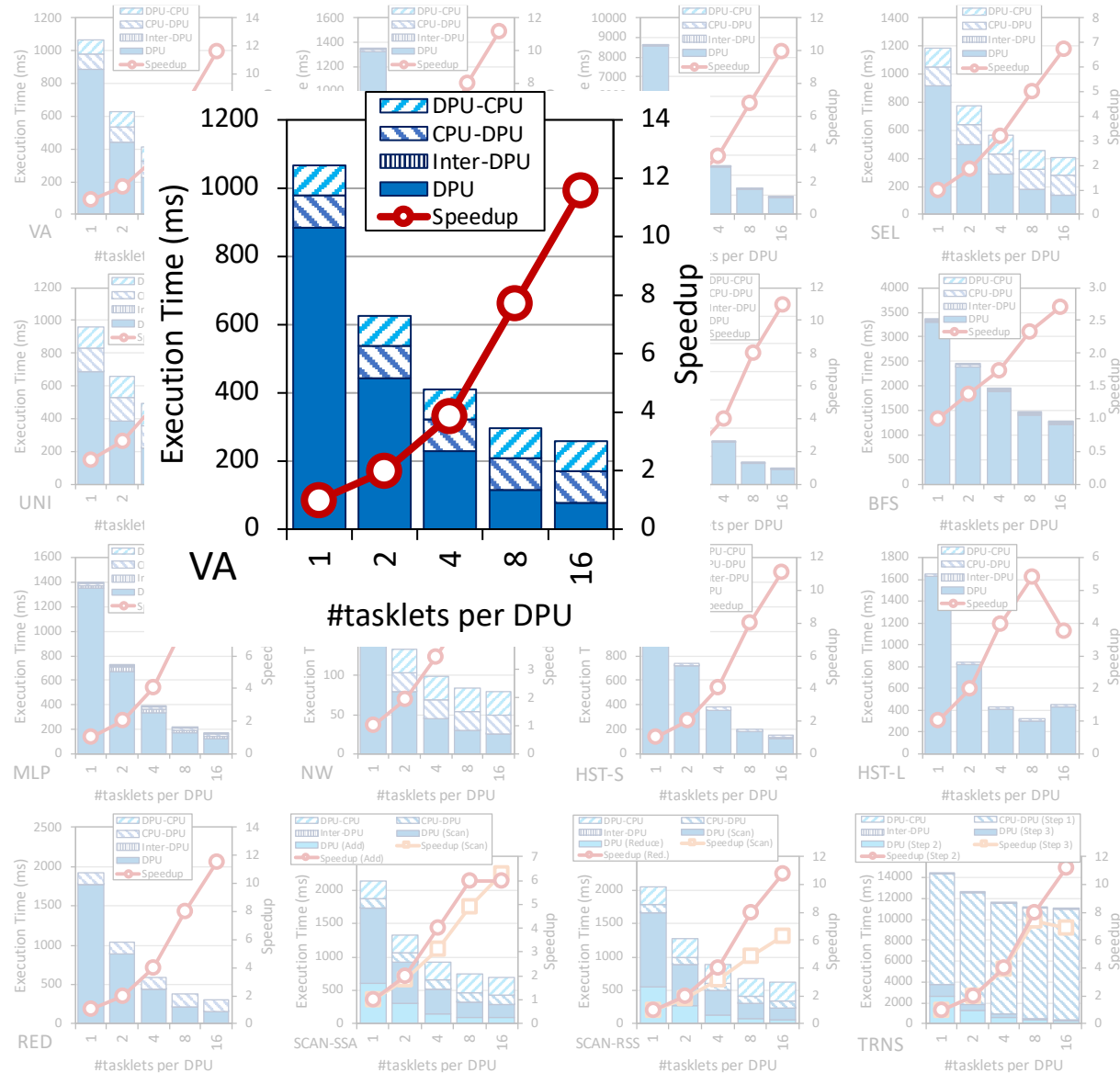
The **PrIM benchmarks** repository includes all datasets and scripts used in our evaluation  
<https://github.com/CMU-SAFARI/prim-benchmarks>



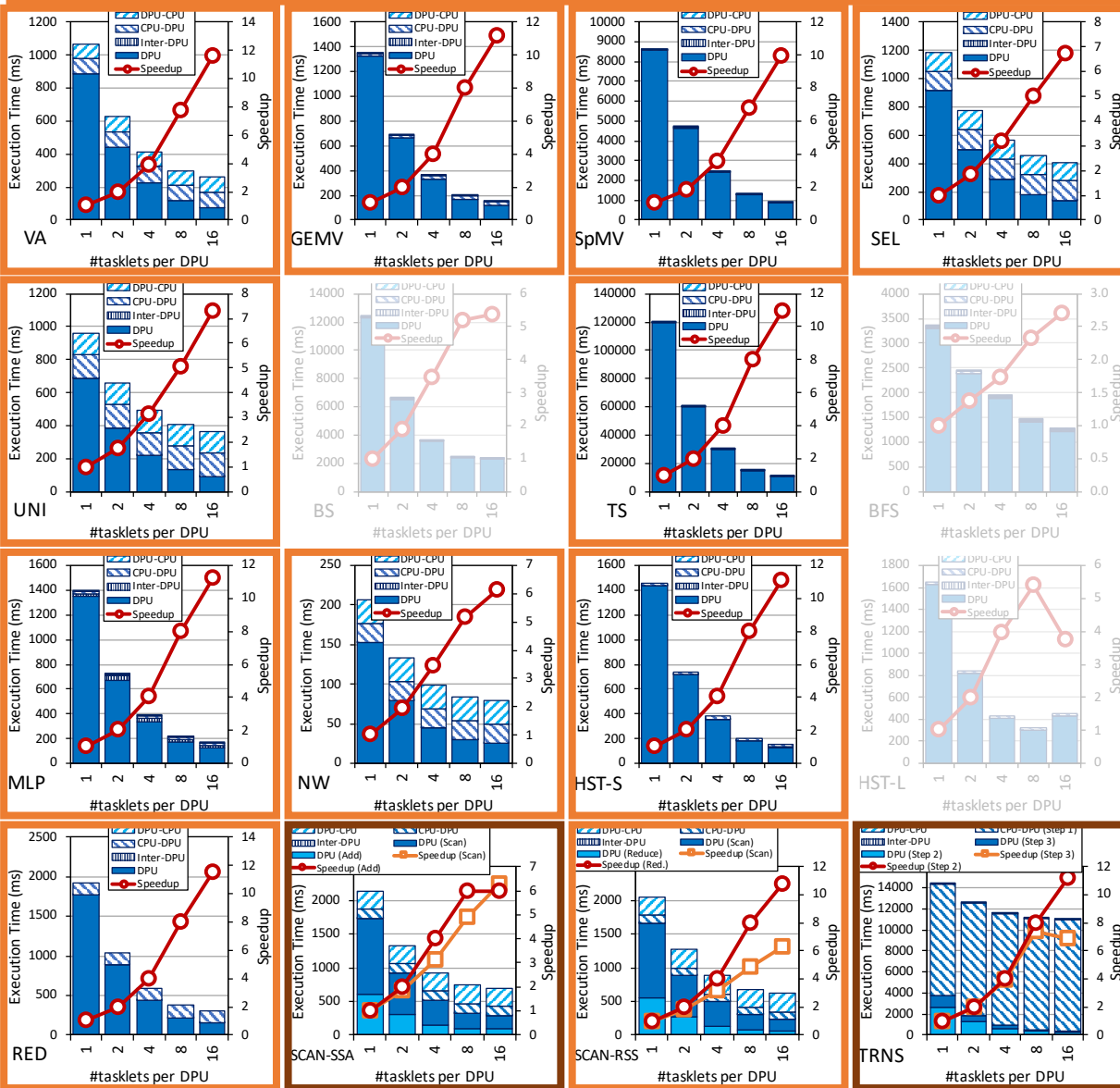
# Strong Scaling: 1 DPU (I)

- Strong scaling experiments on 1 DPU

- We set the **number of tasklets to 1, 2, 4, 8, and 16**
- We show the breakdown of execution time:
  - **DPU**: Execution time on the DPU
  - **Inter-DPU**: Time for inter-DPU communication via the host CPU
  - **CPU-DPU**: Time for CPU to DPU transfer of input data
  - **DPU-CPU**: Time for DPU to CPU transfer of final results
- **Speedup over 1 tasklet**



# Strong Scaling: 1 DPU (II)



VA, GEMV, SpMV, SEL, UNI, TS, MLP, NW, HST-S, RED, SCAN-SSA (Scan kernel), SCAN-RSS (both kernels), and TRNS (Step 2 kernel), the best performing number of tasklets is 16

Speedups 1.5-2.0x as we double the number of tasklets from 1 to 8. Speedups 1.2-1.5x from 8 to 16, since the pipeline throughput saturates at 11 tasklets

## KEY OBSERVATION 10

A number of tasklets greater than 11 is a good choice for most real-world workloads we tested (16 kernels out of 19 benchmarks from 16 kernels from 16 benchmarks), as it fully utilizes the DPU's pipeline.

# Strong Scaling: 1 DPU (III)

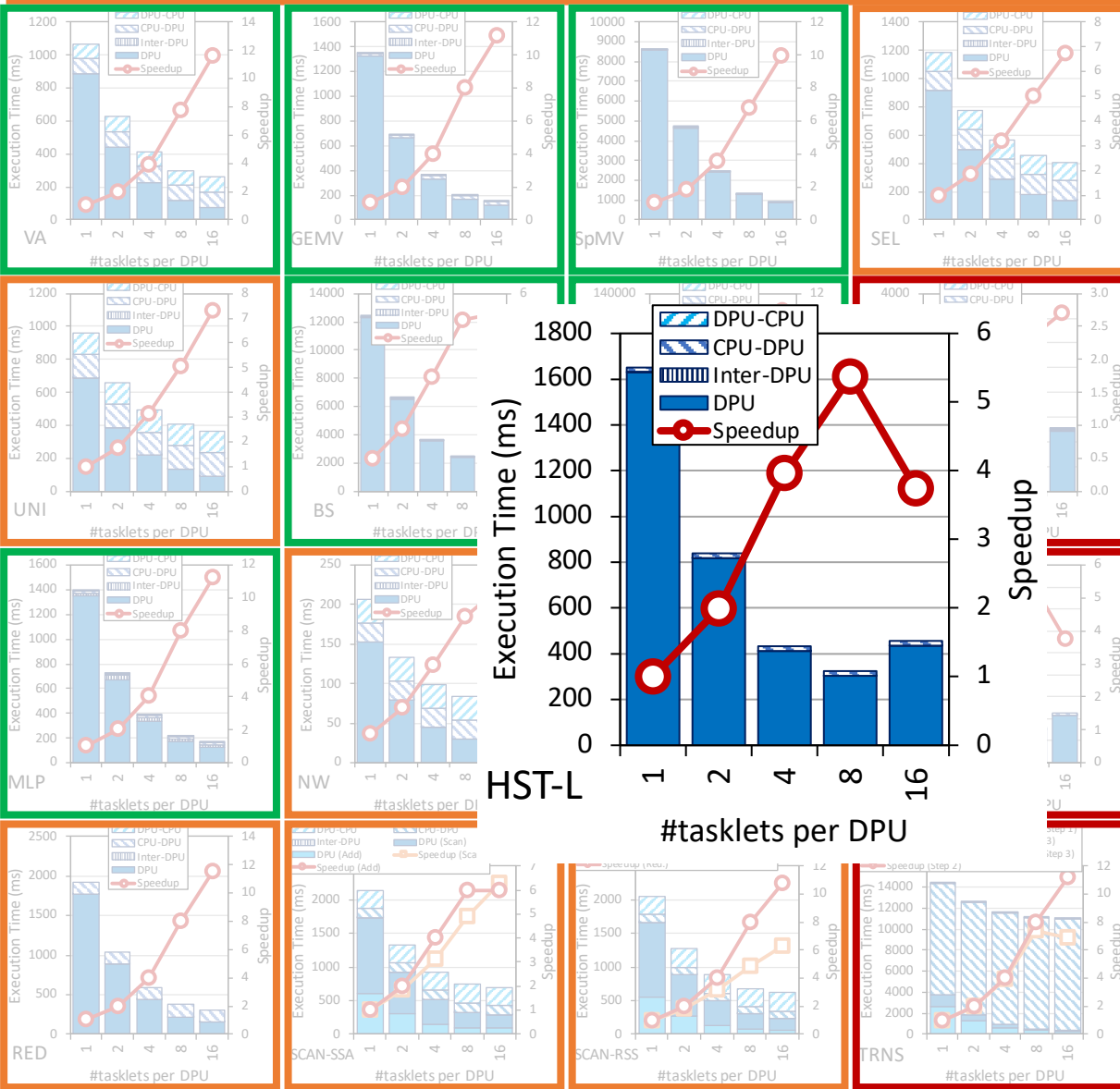


VA, GEMV, SpMV, BS, TS, MLP, HST-S do not use intra-DPU synchronization primitives

In SEL, UNI, NW, RED, SCAN-SSA (Scan kernel), SCAN-RSS (both kernels), synchronization is lightweight

BFS, HST-L, TRNS (Step 3) use mutexes, which cause contention when accessing shared data structures

# Strong Scaling: 1 DPU (IV)



VA, GEMV, SpMV, BS, TS, MLP, HST-S do not use synchronization primitives

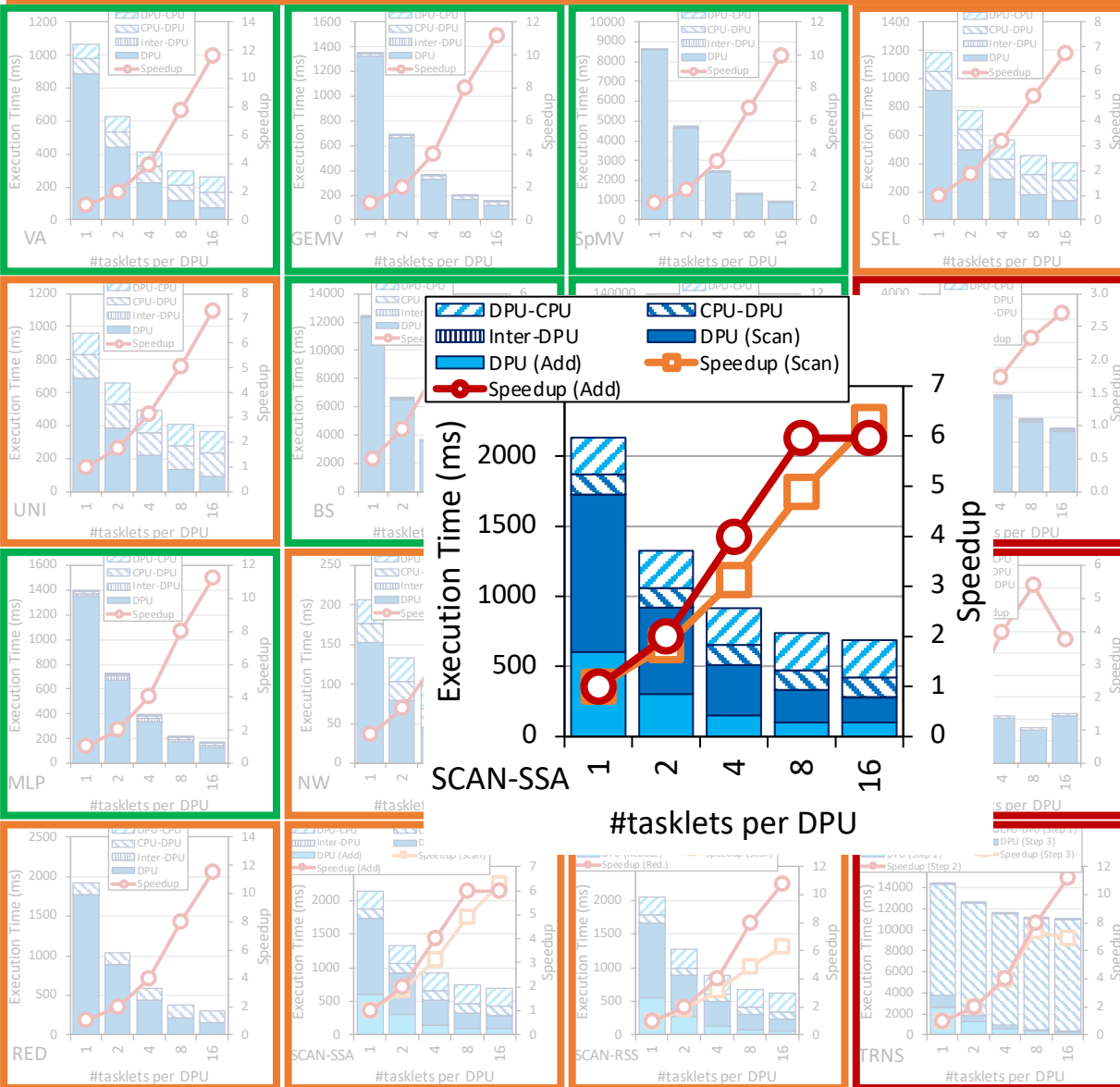
In SEL, UNI, NW, RED, SCAN-SSA (Scan kernel), SCAN-RSS (both kernels), synchronization is lightweight

BFS, HST-L, TRNS (Step 3) use mutexes, which cause contention when accessing shared data structures

## KEY OBSERVATION 11

Intensive use of intra-DPU synchronization across tasklets (e.g., mutexes, barriers, handshakes) may limit scalability, sometimes causing the best performing number of tasklets to be lower than 11.

# Strong Scaling: 1 DPU (V)

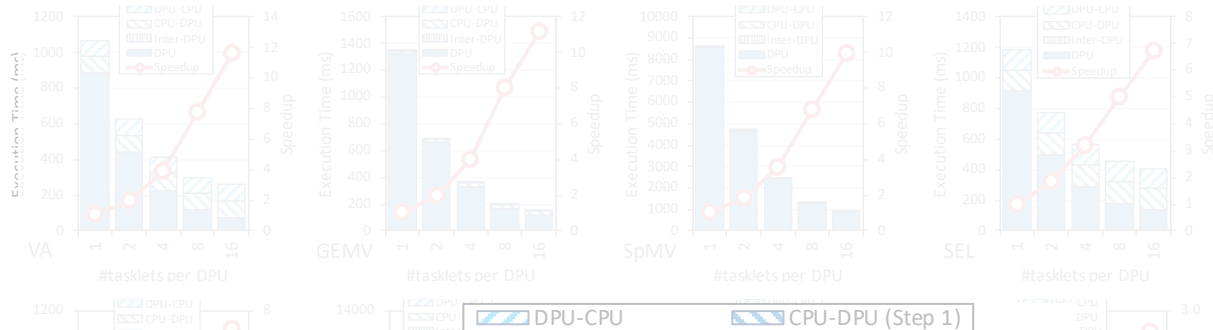


SCAN-SSA (Add kernel) is **not compute-intensive**. Thus, performance saturates with less than 11 tasklets (recall STREAM ADD). BS shows similar behavior

## KEY OBSERVATION 12

**Most real-world workloads are in the compute-bound region of the DPU (all kernels except SCAN-SSA (Add kernel) and BS), i.e., the pipeline latency dominates the MRAM access latency.**

# Strong Scaling: 1 DPU (VI)



The amount of time spent on CPU-DPU and DPU-CPU transfers is low compared to the time spent on DPU execution

## Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich

<sup>2</sup>American University of Beirut

<sup>3</sup>University of Malaga

<sup>4</sup>National Technical University of Athens



Transferring large data chunks from/to the host CPU is preferred for input data and output results due to bandwidth constraints.

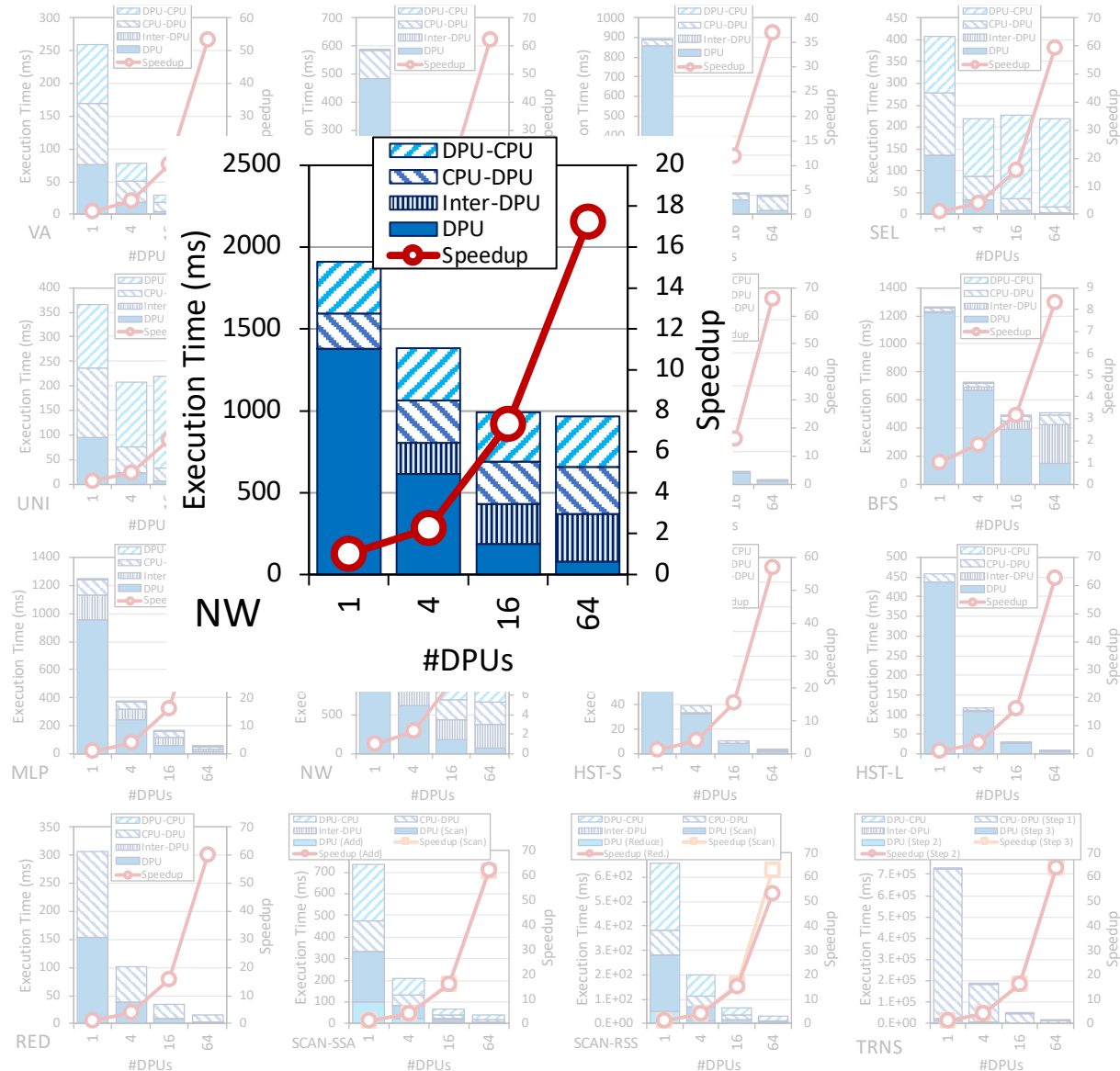
<https://arxiv.org/pdf/2105.03814.pdf>

<https://github.com/CMU-SAFARI/prim-benchmarks>

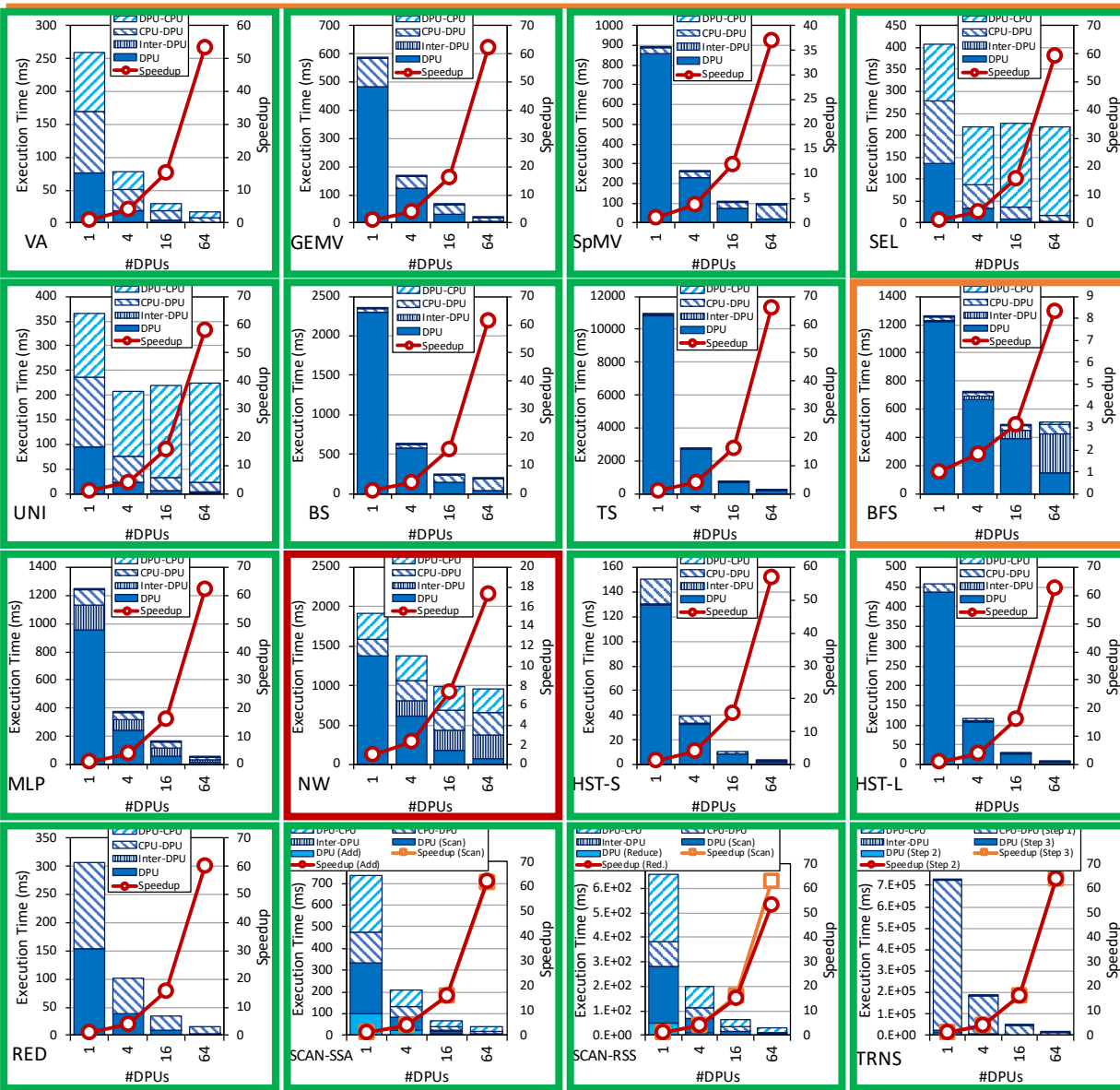


# Strong Scaling: 1 Rank (I)

- Strong scaling experiments on 1 rank
  - We set the number of tasklets to the best performing one
  - The number of DPUs is 1, 4, 16, 64
  - We show the breakdown of execution time:
    - DPU: Execution time on the DPU
    - Inter-DPU: Time for inter-DPU communication via the host CPU
    - CPU-DPU: Time for CPU to DPU transfer of input data
    - DPU-CPU: Time for DPU to CPU transfer of final results
  - Speedup over 1 DPU



# Strong Scaling: 1 Rank (II)



VA, GEMV, SpMV, SEL, UNI, BS, TS, MLP, HST-S, HSTS-L, RED, SCAN-SSA (both kernel), SCAN-RSS (both kernels), and TRNS (both kernels) scale linearly with the number of DPUs

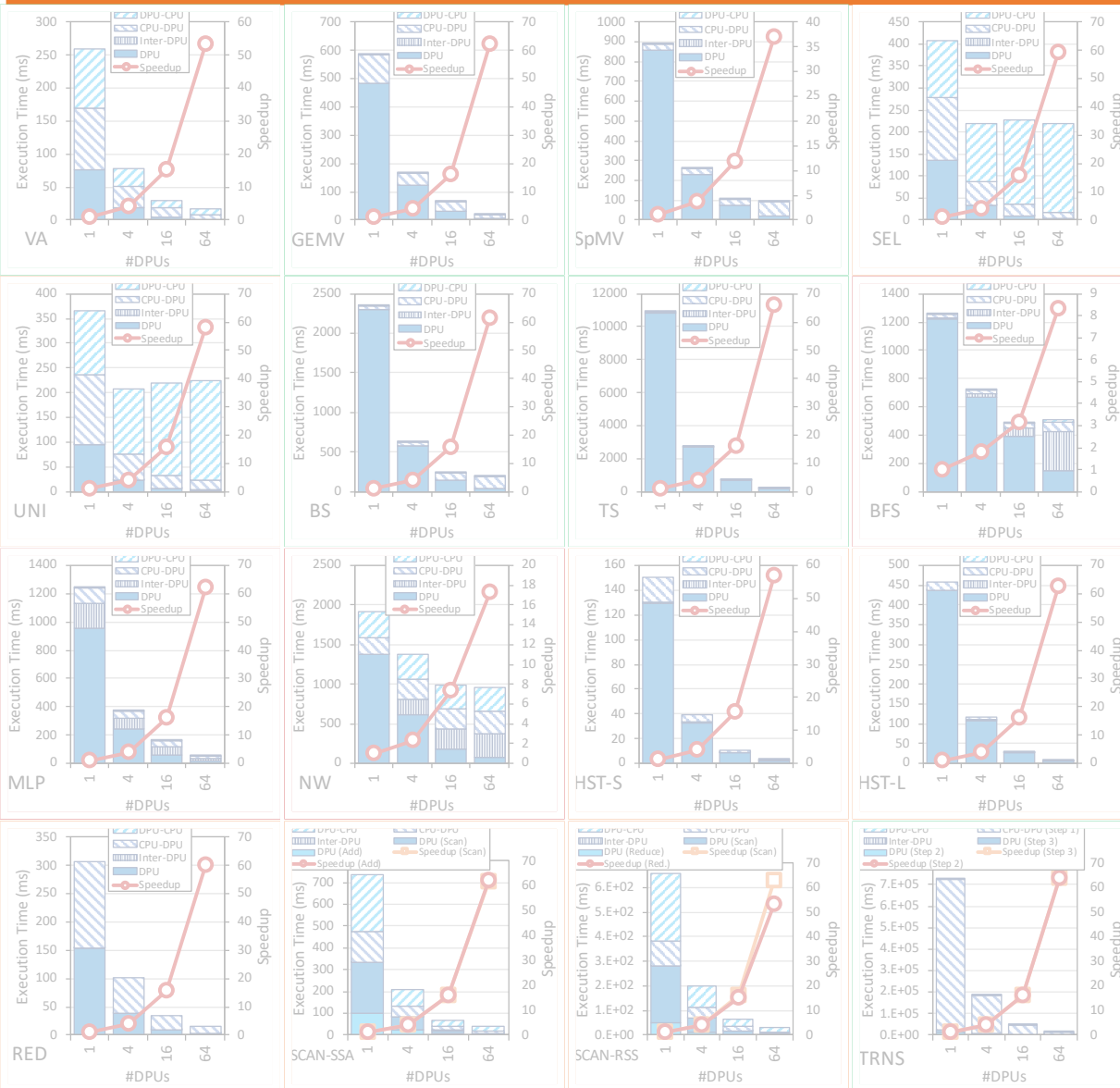
Scaling is sublinear for BFS and NW

BFS suffers **load imbalance** due to irregular graph topology

NW computes a diagonal of a 2D matrix in each iteration.  
More DPUs does not mean more parallelization in shorter diagonals.



# Strong Scaling: 1 Rank (III)

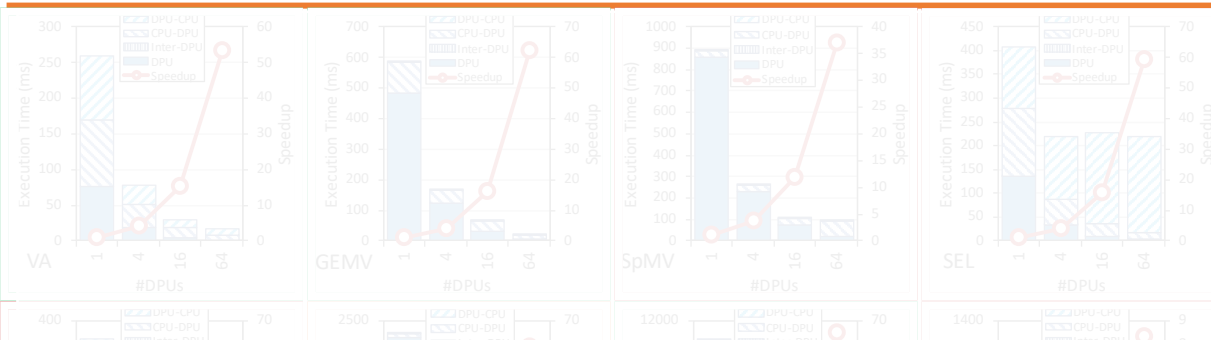


VA, GEMV, SpMV, BS, TS, TRNS **do not need inter-DPU synchronization**

SEL, UNI, HST-S, HST-L, RED, SCAN-SSA, SCAN-RSS **need inter-DPU synchronization but 64 DPUs still obtain the best performance**

BFS, MLP, NW require **heavy inter-DPU synchronization**, involving DPU-CPU and CPU-DPU transfers

# Strong Scaling: 1 Rank (IV)



VA, GEMV, TS, MLP, HST-S, HST-L, RED, SCAN-SSA, SCAN-RSS, TRNS use parallel transfers. CPU-DPU and DPU-CPU transfer times decrease as we increase the number of DPUs

## Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich

<sup>2</sup>American University of Beirut

<sup>3</sup>University of Malaga

<sup>4</sup>National Technical University of Athens



size per DPU is not fixed

### PROGRAMMING RECOMMENDATION 5

Parallel CPU-DPU/DPU-CPU transfers inside a rank of DPUs

<https://arxiv.org/pdf/2105.03814.pdf>

<https://github.com/CMU-SAFARI/prim-benchmarks>

for real-  
en all  
of the same

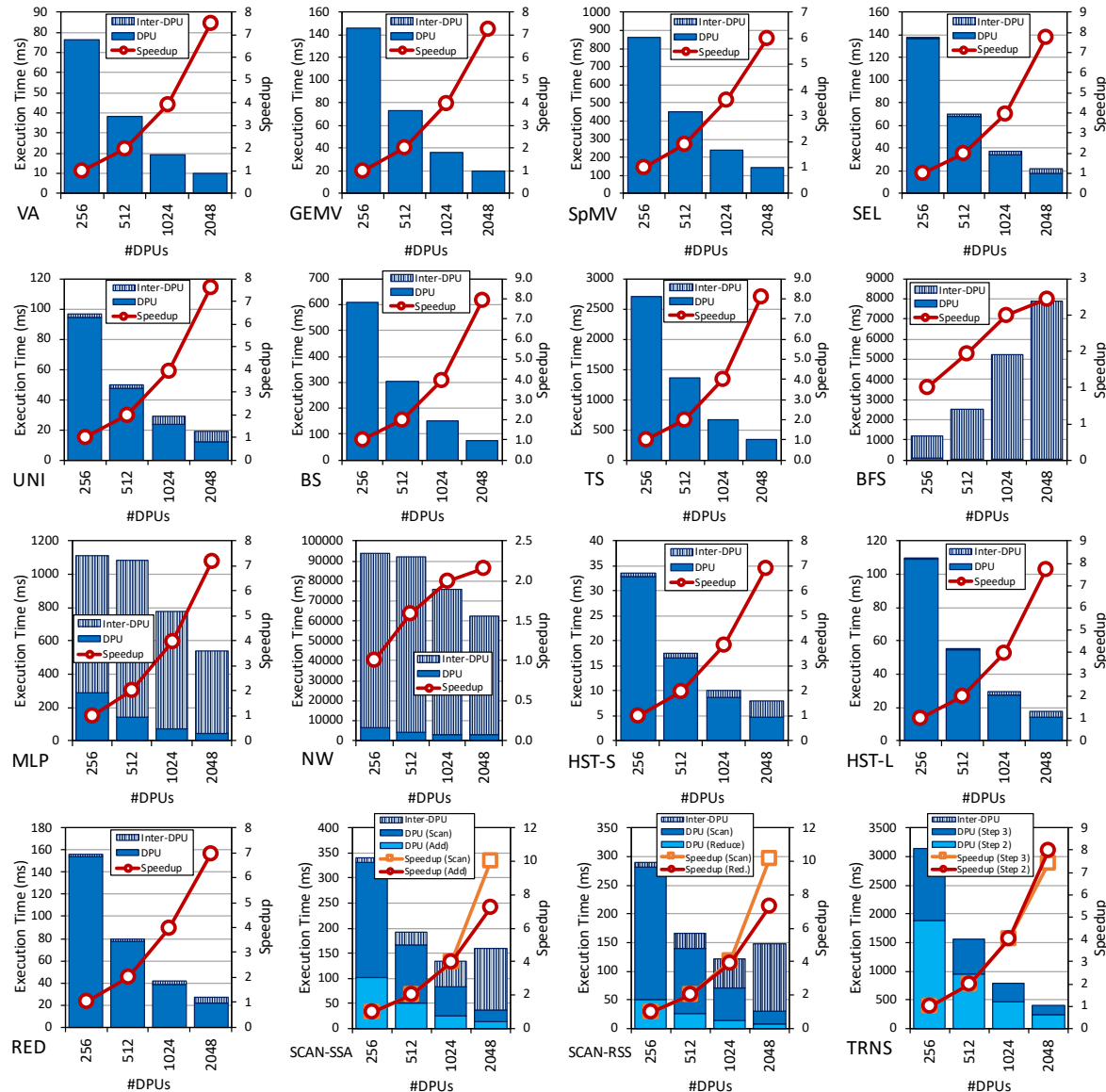
# Strong Scaling: 32 Ranks

- Strong scaling experiments on 32 rank

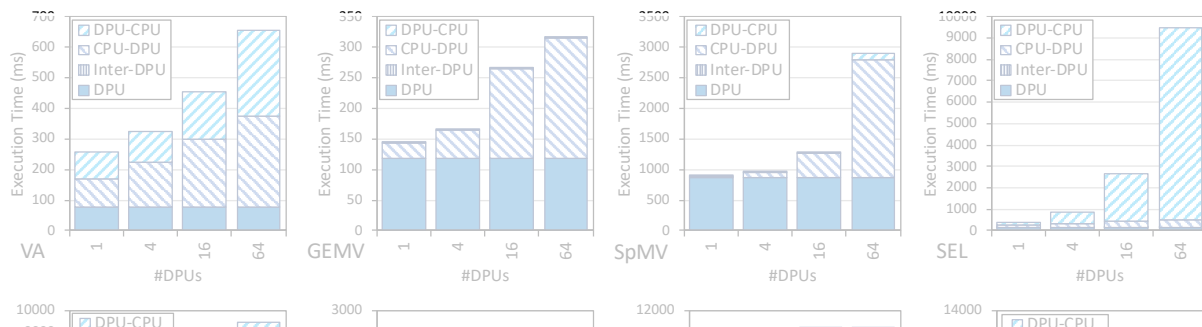
- We set the number of tasklets to the best performing one
- The number of DPUs is 256, 512, 1024, 2048
- We show the breakdown of execution time:

- DPU: Execution time on the DPU
- Inter-DPU: Time for inter-DPU communication via the host CPU
- We do not show CPU-DPU/DPU-CPU transfer times

- Speedup over 256 DPUs



# Weak Scaling: 1 Rank



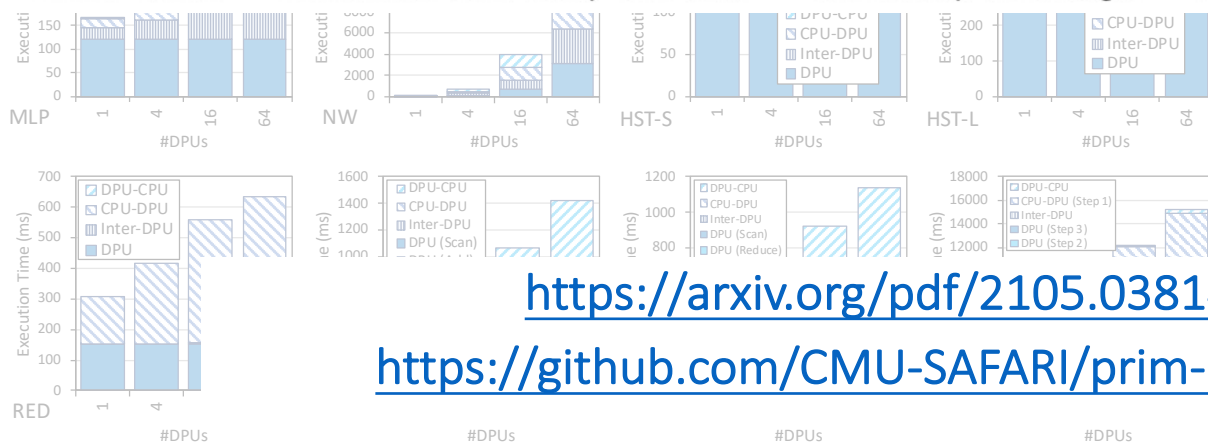
## KEY OBSERVATION 17

Equally-sized problems assigned to different DPUs and little/no inter-DPU

## Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>American University of Beirut <sup>3</sup>University of Malaga <sup>4</sup>National Technical University of Athens



## KEY OBSERVATION 18

Sustained bandwidth of parallel CPU-DPU/DPU-CPU transfers inside a rank of bilinearly DPUs.

<https://arxiv.org/pdf/2105.03814.pdf>

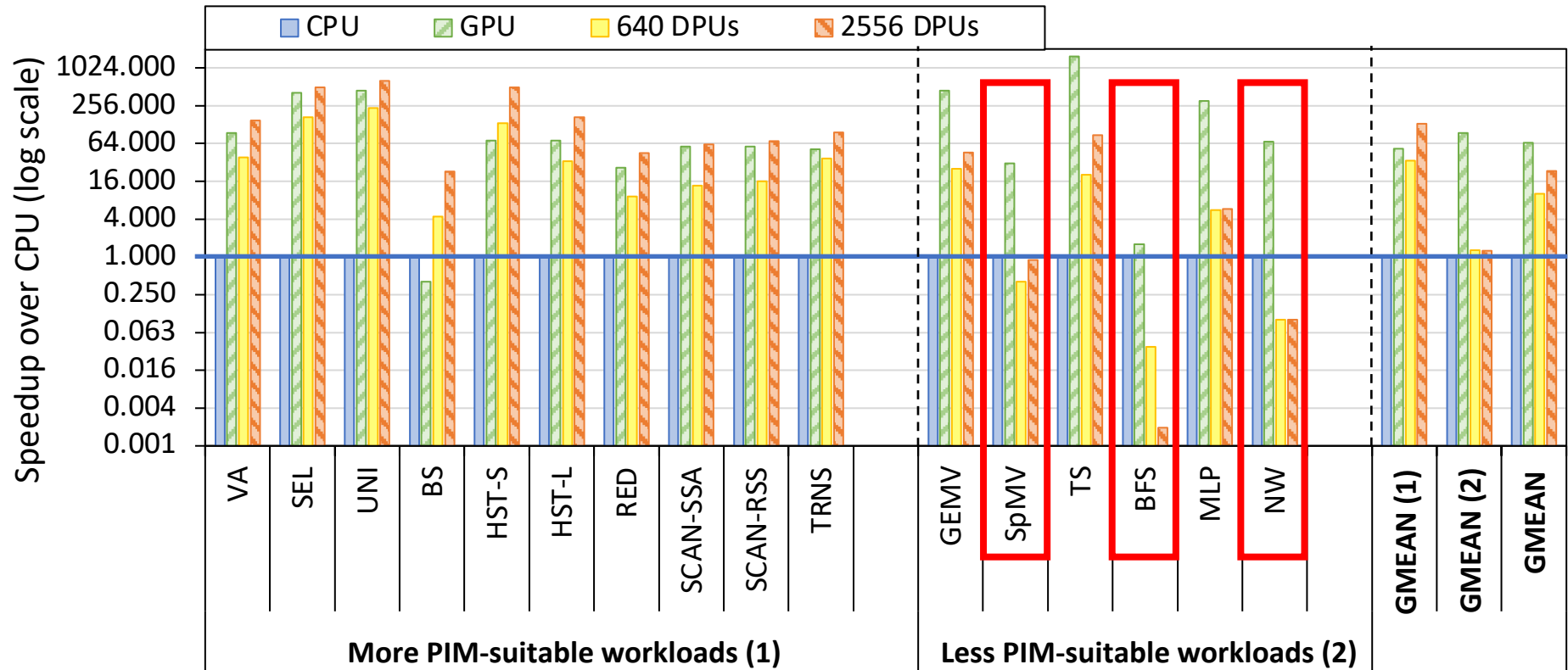
<https://github.com/CMU-SAFARI/prim-benchmarks>

# CPU/GPU: Evaluation Methodology

---

- Comparison of both UPMEM-based PIM systems **to state-of-the-art CPU and GPU**
  - Intel Xeon E3-1240 CPU
  - NVIDIA Titan V GPU
- We use **state-of-the-art CPU and GPU counterparts** of PrIM benchmarks
  - <https://github.com/CMU-SAFARI/prim-benchmarks>
- We use the **largest dataset that we can fit in the GPU memory**
- We show overall execution time, including DPU kernel time and inter DPU communication

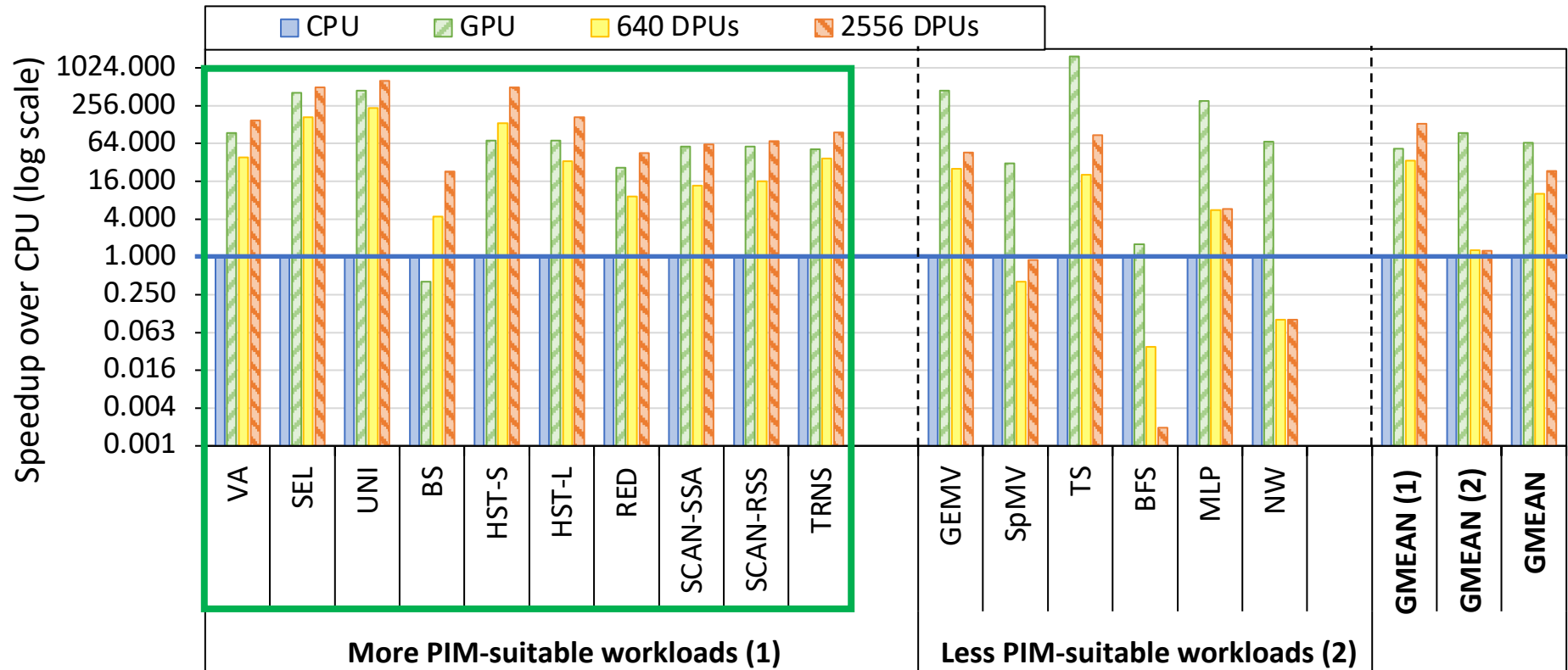
# CPU/GPU: Performance Comparison (I)



The 2,556-DPU and the 640-DPU systems outperform the CPU for all benchmarks except SpMV, BFS, and NW

The 2,556-DPU and the 640-DPU are, respectively, 93.0x and 27.9x faster than the CPU for 13 of the PRIM benchmarks

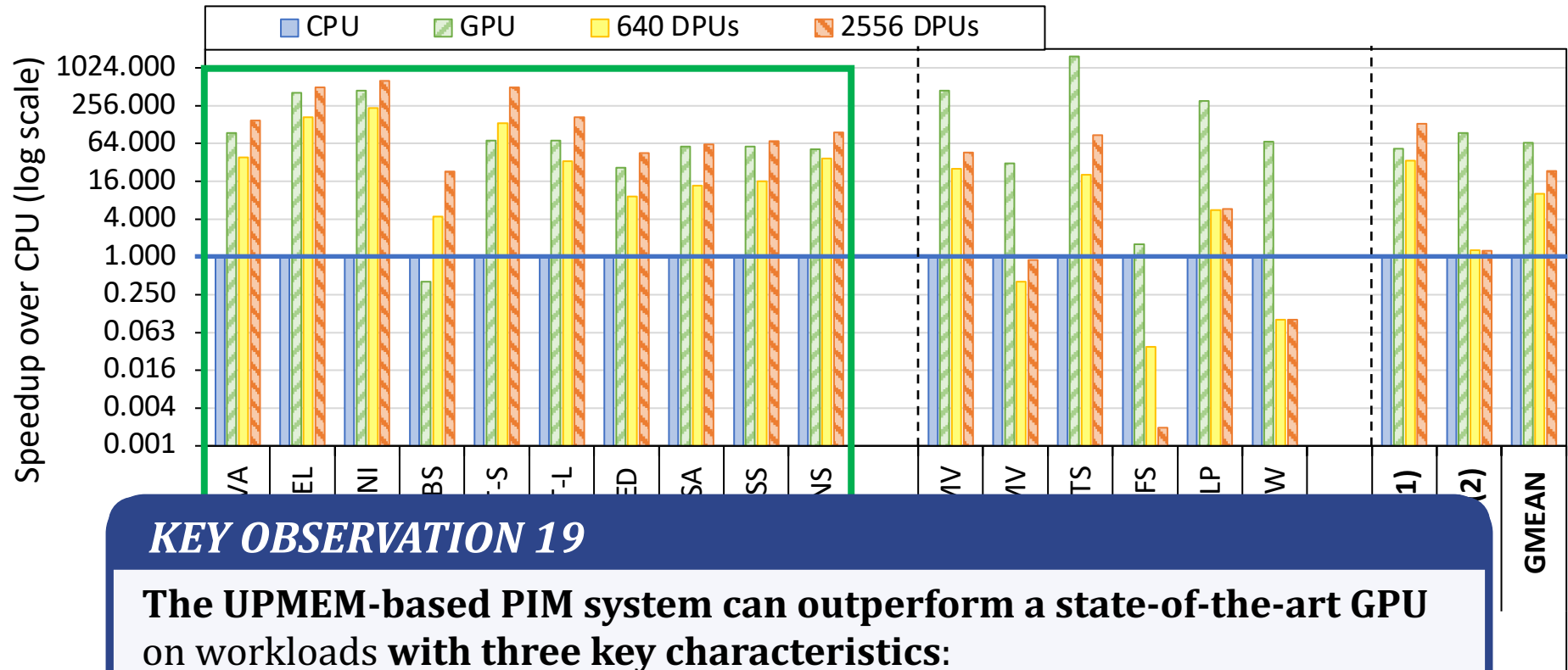
# CPU/GPU: Performance Comparison (II)



The 2,556-DPU outperforms the GPU  
for 10 PRIM benchmarks with an average of 2.54x

The performance of the 640-DPU is within 65%  
the performance of the GPU for the same 10 PRIM benchmarks

# CPU/GPU: Performance Comparison (III)



## KEY OBSERVATION 19

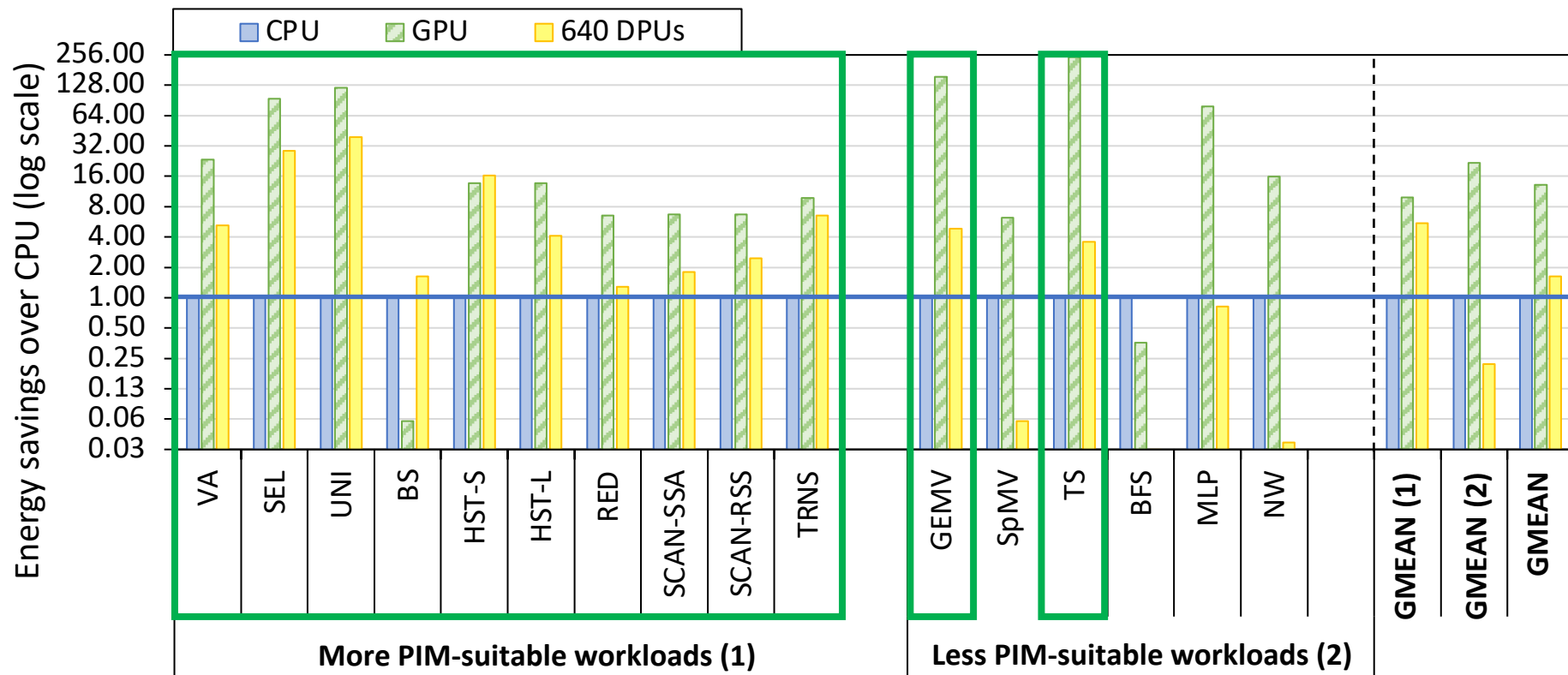
The UPMEM-based PIM system can outperform a state-of-the-art GPU on workloads **with three key characteristics**:

1. Streaming memory accesses
2. No or little inter-DPU synchronization
3. No or little use of integer multiplication, integer division, or floating point operations

These three key characteristics make a **workload potentially suitable to the UPMEM PIM architecture**.



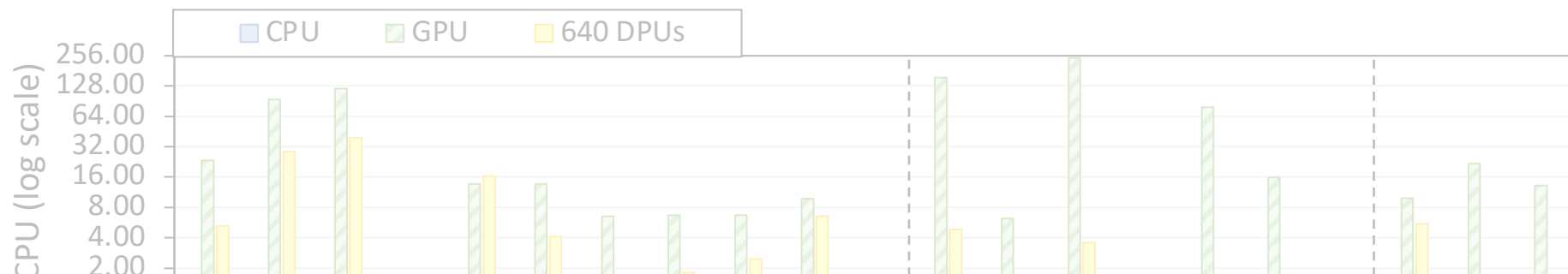
# CPU/GPU: Energy Comparison (I)



The 640-DPU system consumes **on average 1.64x less energy** than the CPU for all 16 PrIM benchmarks

For 12 benchmarks, the 640-DPU system provides energy savings of **5.23x over the CPU**

# CPU/GPU: Energy Comparison (II)



## Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez-Luna<sup>1</sup> Izzat El Hajj<sup>2</sup> Ivan Fernandez<sup>1,3</sup> Christina Giannoula<sup>1,4</sup>  
Geraldo F. Oliveira<sup>1</sup> Onur Mutlu<sup>1</sup>

<sup>1</sup>ETH Zürich <sup>2</sup>American University of Beirut <sup>3</sup>University of Malaga <sup>4</sup>National Technical University of Athens

and less data movement between memory and processors.

The UPMEM-based PIM system provides energy savings over a state-of-the-art CPU/GPU on workloads where it outperforms the CPU/GPU.

This is because the source of both performance improvement and energy

sa  
th  
sy

<https://arxiv.org/pdf/2105.03814.pdf>

<https://github.com/CMU-SAFARI/prim-benchmarks>

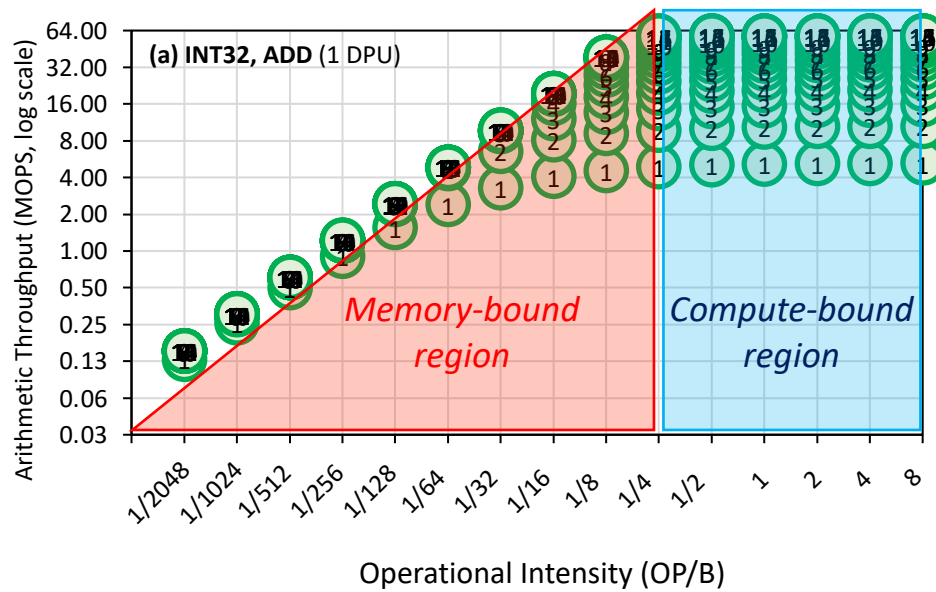
een

# Outline

---

- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PRIM Benchmarks
  - Roofline Model
  - Benchmark Diversity
- Evaluation
  - Strong and Weak Scaling
  - Comparison to CPU and GPU
- Key Takeaways

# Key Takeaway 1

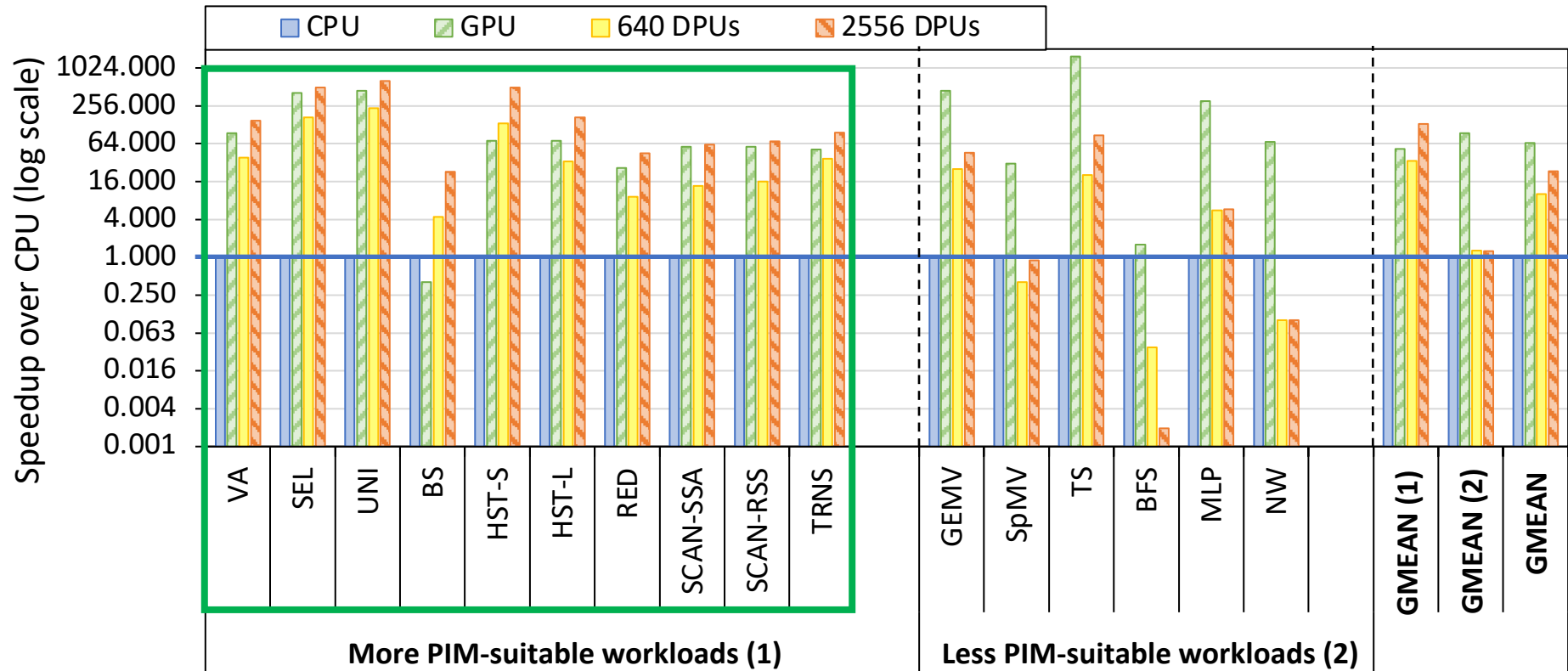


The throughput saturation point is as low as  $\frac{1}{4}$  OP/B, i.e., 1 integer addition per every 32-bit element fetched

## KEY TAKEAWAY 1

The UPMEM PIM architecture is fundamentally compute bound. As a result, the most suitable workloads are memory-bound.

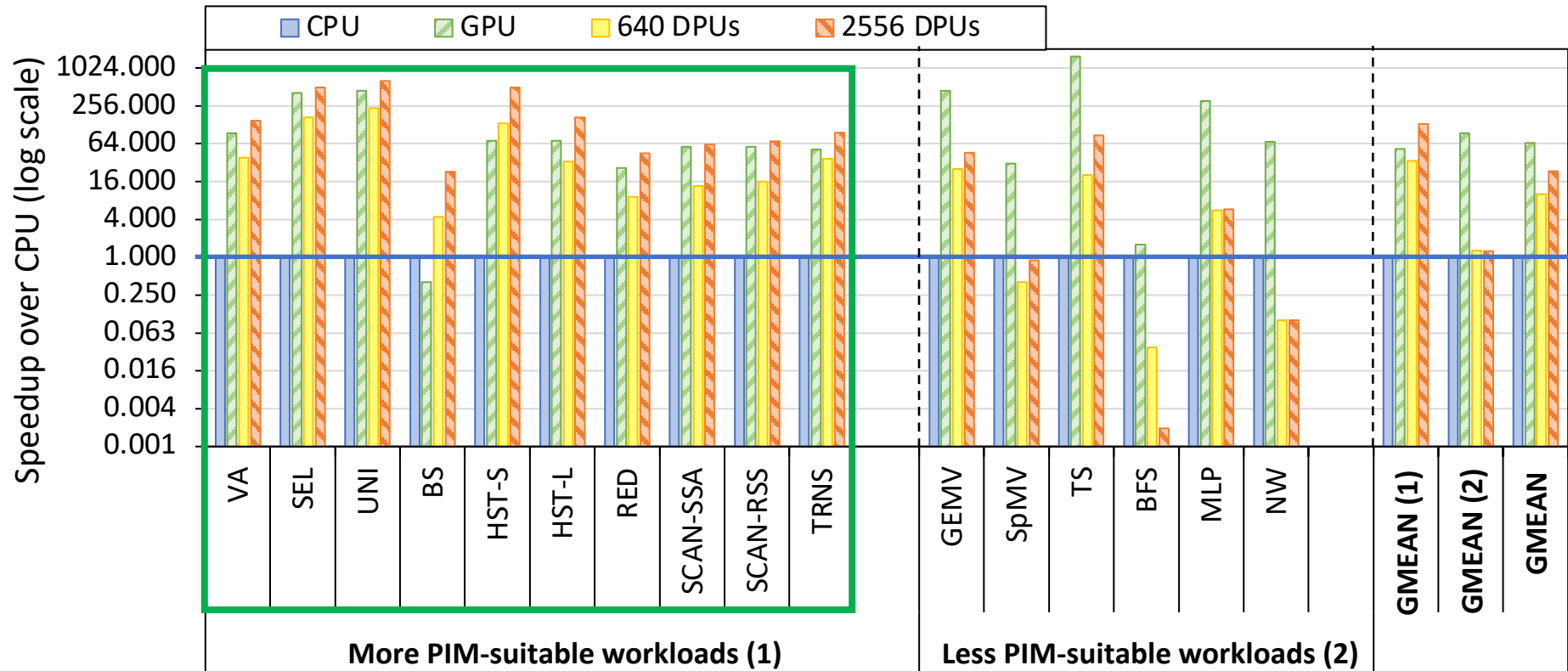
# Key Takeaway 2



## KEY TAKEAWAY 2

The most well-suited workloads for the UPMEM PIM architecture use no arithmetic operations or use only simple operations (e.g., bitwise operations and integer addition/subtraction).

# Key Takeaway 3



## KEY TAKEAWAY 3

The most well-suited workloads for the UPMEM PIM architecture require little or no communication across DPUs (inter-DPU communication).

# Key Takeaway 4

---

## ***KEY TAKEAWAY 4***

- UPMEM-based PIM systems **outperform state-of-the-art CPUs in terms of performance and energy efficiency on most of PrIM benchmarks.**
- UPMEM-based PIM systems **outperform state-of-the-art GPUs on a majority of PrIM benchmarks**, and the outlook is even more positive for future PIM systems.
- UPMEM-based PIM systems are **more energy-efficient than state-of-the-art CPUs and GPUs on workloads that they provide performance improvements** over the CPUs and the GPUs.

# Executive Summary

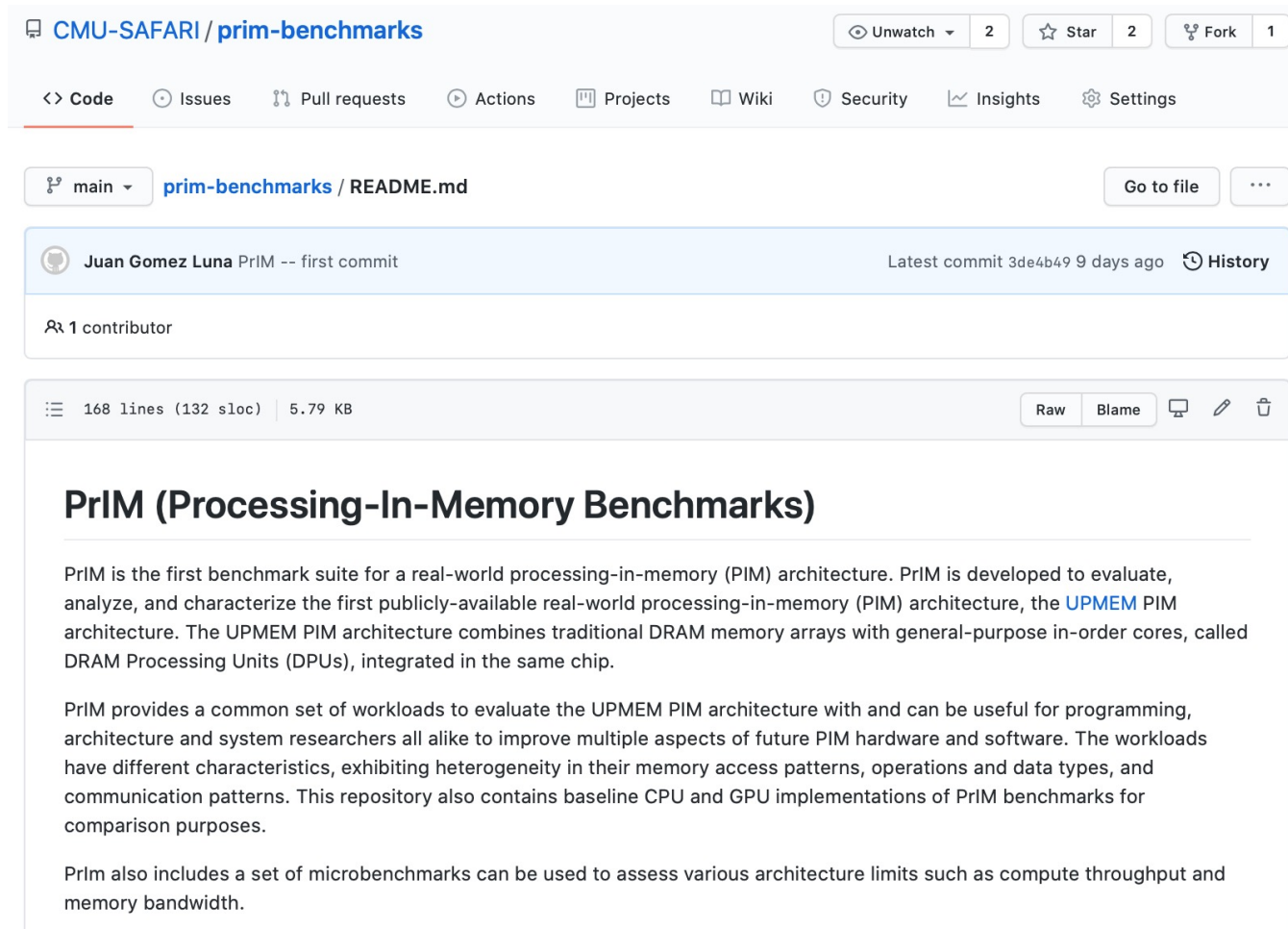
---

- **Data movement** between memory/storage units and compute units is a major contributor to execution time and energy consumption
- **Processing-in-Memory** (PIM) is a paradigm that can tackle the **data movement bottleneck**
  - Though explored for +50 years, technology challenges prevented the successful materialization
- UPMEM has designed and fabricated **the first publicly-available real-world PIM architecture**
  - **DDR4 chips embedding in-order multithreaded DRAM Processing Units (DPUs)**
- Our work:
  - **Introduction** to UPMEM programming model and PIM architecture
  - **Microbenchmark-based characterization** of the DPU
  - Benchmarking and **workload suitability** study
- Main contributions:
  - Comprehensive **characterization and analysis of the first commercially-available PIM architecture**
  - **PrIM (Processing-In-Memory) benchmarks**:
    - 16 workloads that are memory-bound in conventional processor-centric systems
    - Strong and weak scaling characteristics
  - Comparison to **state-of-the-art CPU and GPU**
- Takeaways:
  - Workload characteristics for **PIM suitability**
  - **Programming** recommendations
  - Suggestions and hints for **hardware and architecture designers** of future PIM systems
  - **PrIM**: (a) programming samples, (b) evaluation and comparison of current and future PIM systems



# PrIM Repository

- All microbenchmarks, benchmarks, and scripts
- <https://github.com/CMU-SAFARI/prim-benchmarks>



The screenshot shows the GitHub repository page for `CMU-SAFARI/prim-benchmarks`. At the top, there are navigation links for Code, Issues, Pull requests, Actions, Projects, Wiki, Security, Insights, and Settings. Below these, the repository name is displayed along with statistics: 2 Unwatch, 2 Stars, and 1 Fork. The main content area shows the `prim-benchmarks / README.md` file. The commit history section indicates the first commit by Juan Gomez Luna. The file statistics show 168 lines (132 sloc) and 5.79 KB. The README content describes the PrIM benchmark suite, its purpose, and its components.

CMU-SAFARI / `prim-benchmarks` Unwatch 2 Star 2 Fork 1

<> Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main `prim-benchmarks / README.md` Go to file ...

Juan Gomez Luna PrIM -- first commit Latest commit 3de4b49 9 days ago History

1 contributor

168 lines (132 sloc) 5.79 KB Raw Blame

## PrIM (Processing-In-Memory Benchmarks)

PrIM is the first benchmark suite for a real-world processing-in-memory (PIM) architecture. PrIM is developed to evaluate, analyze, and characterize the first publicly-available real-world processing-in-memory (PIM) architecture, the [UPMEM](#) PIM architecture. The UPMEM PIM architecture combines traditional DRAM memory arrays with general-purpose in-order cores, called DRAM Processing Units (DPUs), integrated in the same chip.

PrIM provides a common set of workloads to evaluate the UPMEM PIM architecture with and can be useful for programming, architecture and system researchers all alike to improve multiple aspects of future PIM hardware and software. The workloads have different characteristics, exhibiting heterogeneity in their memory access patterns, operations and data types, and communication patterns. This repository also contains baseline CPU and GPU implementations of PrIM benchmarks for comparison purposes.

PrIM also includes a set of microbenchmarks can be used to assess various architecture limits such as compute throughput and memory bandwidth.

# Understanding a Modern Processing-in-Memory Architecture: Benchmarking and Experimental Characterization

Juan Gómez Luna, Izzat El Hajj,  
Ivan Fernandez, Christina Giannoula,  
Geraldo F. Oliveira, Onur Mutlu

[el1goluj@gmail.com](mailto:el1goluj@gmail.com)

<https://arxiv.org/pdf/2105.03814.pdf>

<https://github.com/CMU-SAFARI/prim-benchmarks>

# Technology Challenges

## The Hurdles on the road to the Graal

- DRAM process highly constrained
  - 3x slower transistors than same node digital process
  - Logic 10 times less dense vs. ASIC process
  - Routing density dramatically lower
    - 3 metals only for routing (vs. 10+), pitch x4 larger
- Strong design choices mandatory

But the PIM Graal is worth it !

### Take away

#### DRAM vs. ASIC

- Far less performing
- Wafers 2x cheaper vs. ASIC

#### Leapfrogging Moore's law

- **Total** Energy efficiency x10
- Massive, scalable parallelism
- Very low cost

Copyright UPMEM® 2019

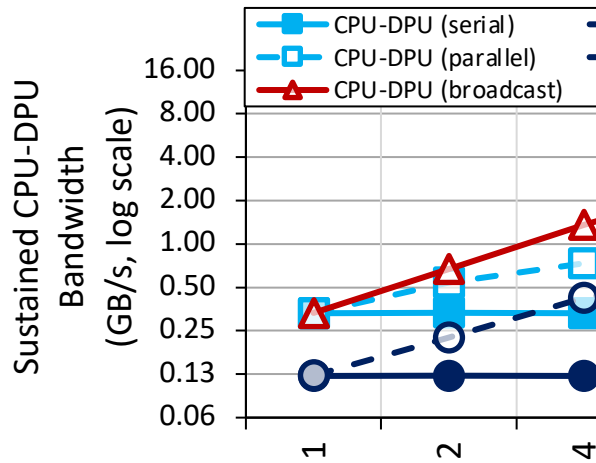
Authorized licensed use limited to: ETH BIBLIOTHEK ZURICH. Downloaded on September 04, 2020 at 13:55:41 UTC from IEEE Xplore. Restrictions apply.

HOT CHIPS 31

up  
mem

# CPU-DPU/DPU-CPU Transfers: 1 Rank (II)

- CPU-DPU (serial/parallel/**broadcast**) and DPU-CPU (serial/parallel)
- The number of DPUs varies between 1 and 64

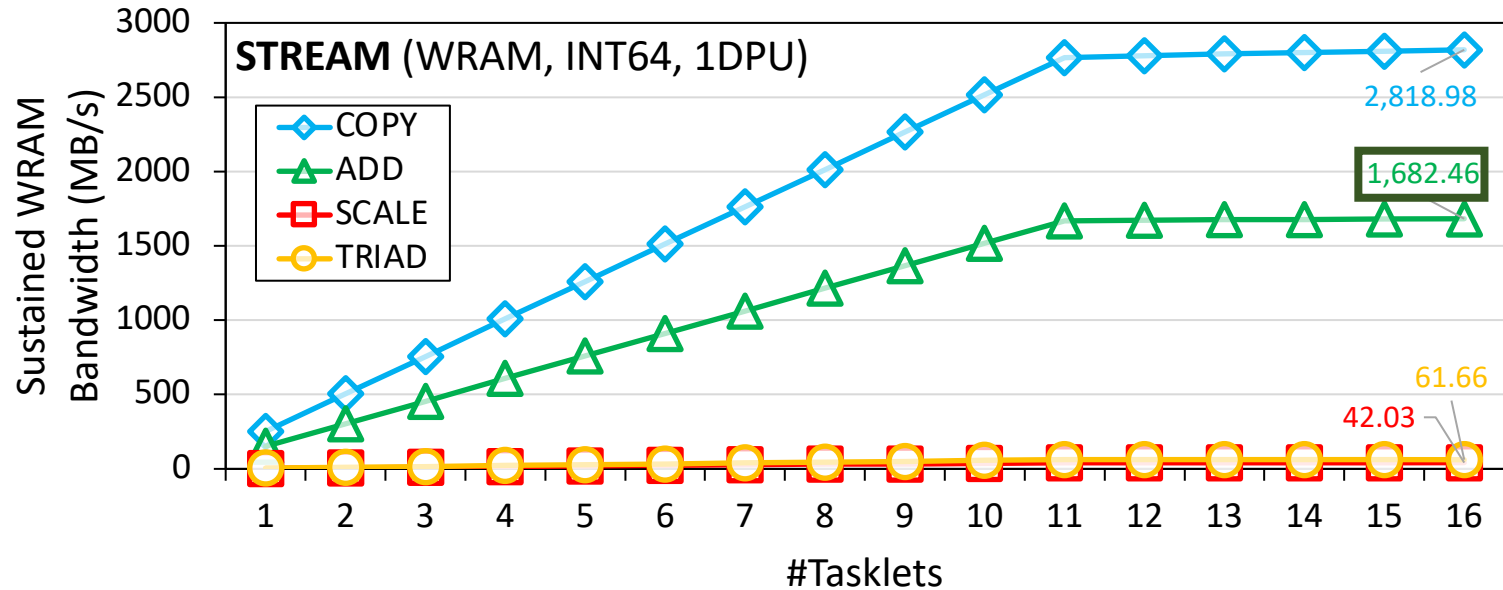


## KEY OBSERVATION 9

**The sustained bandwidth of parallel CPU-DPU transfers is higher than the sustained bandwidth of parallel DPU-CPU transfers due to different implementations of CPU-DPU and DPU-CPU transfers in the UPMEM runtime library.**

**The sustained bandwidth of broadcast CPU-DPU transfers (i.e., the same buffer is copied to multiple MRAM banks) is higher than that of parallel CPU-DPU transfers (i.e., different buffers are copied to different MRAM banks) due to higher temporal locality in the CPU cache hierarchy.**

# WRAM Bandwidth: ADD

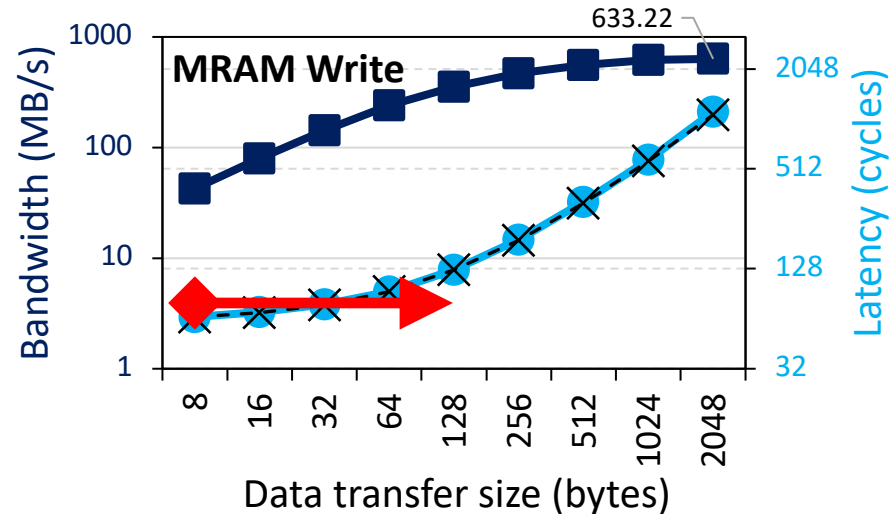
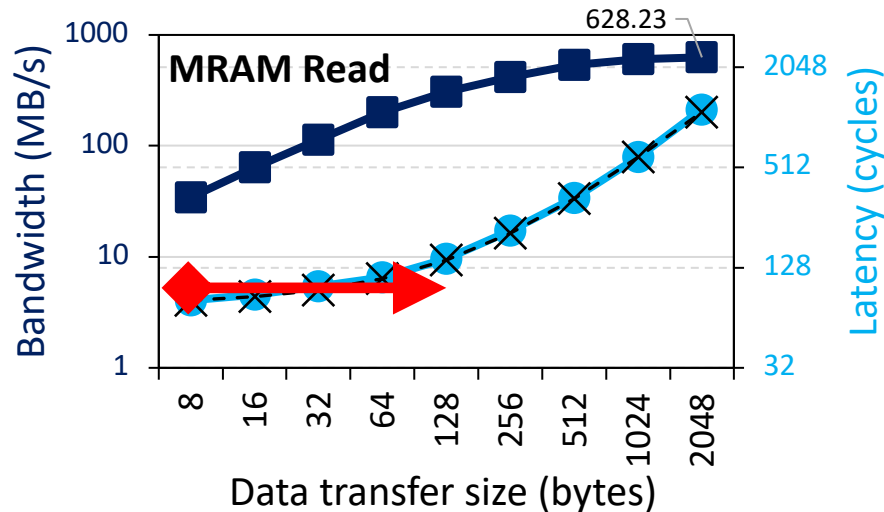


$$WRAM \text{ Bandwidth } \left( in \frac{B}{S} \right) = \frac{Bytes \times frequency_{DPU}}{\#instructions}$$

**ADD** executes 5 instructions (2 ld, add, addc, sd).  
With 11 tasklets,  $11 \times 24$  bytes in 55 cycles:

$$WRAM \text{ Bandwidth } \left( in \frac{B}{S} \right) = 1,680 \frac{MB}{s} \text{ at } 350 \text{ MHz}$$

# MRAM Read and Write Latency (IV)



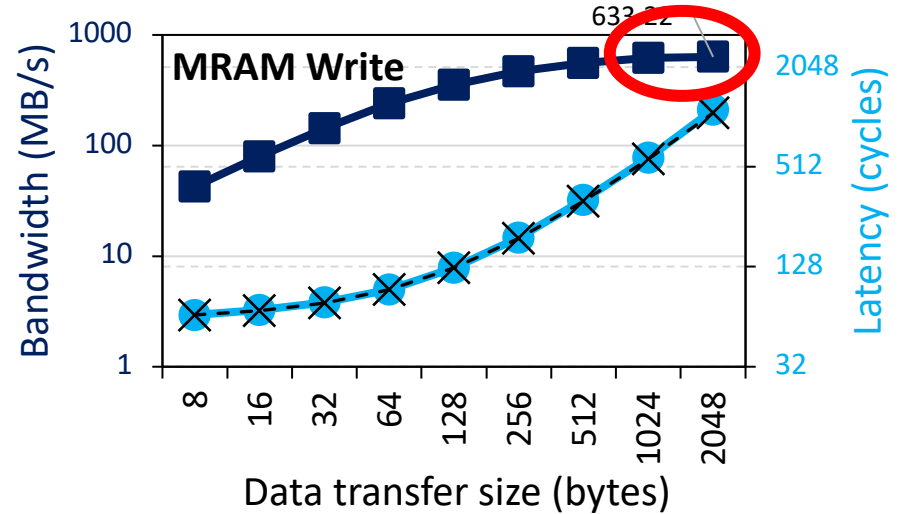
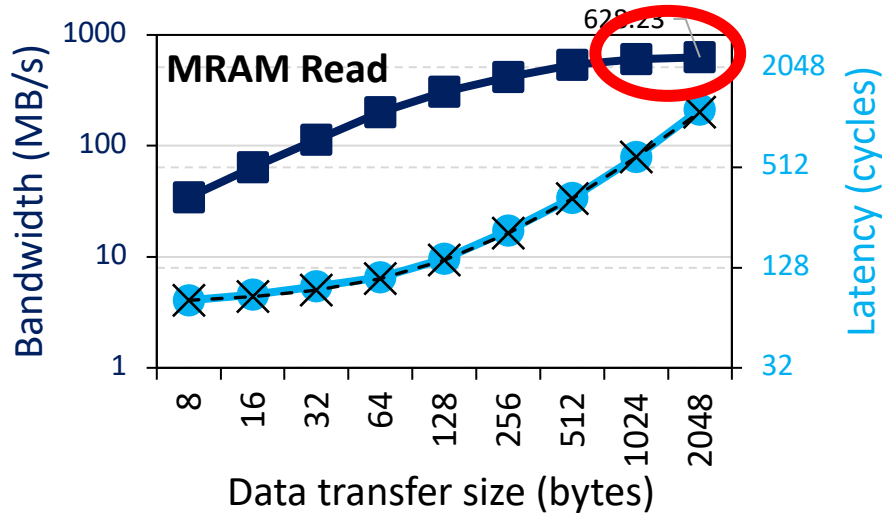
MRAM latency changes slowly between 8 and 128 bytes

For small transfers, the fixed cost ( $\alpha$ ) dominates the variable cost ( $\beta \times \text{size}$ )

## PROGRAMMING RECOMMENDATION 2

For small transfers between the MRAM bank and the WRAM, **fetch more bytes than necessary within a 128-byte limit**. Doing so increases the likelihood of finding data in WRAM for later accesses (i.e., the program can check whether the desired data is in WRAM before issuing a new MRAM access).

# MRAM Read and Write Latency (V)



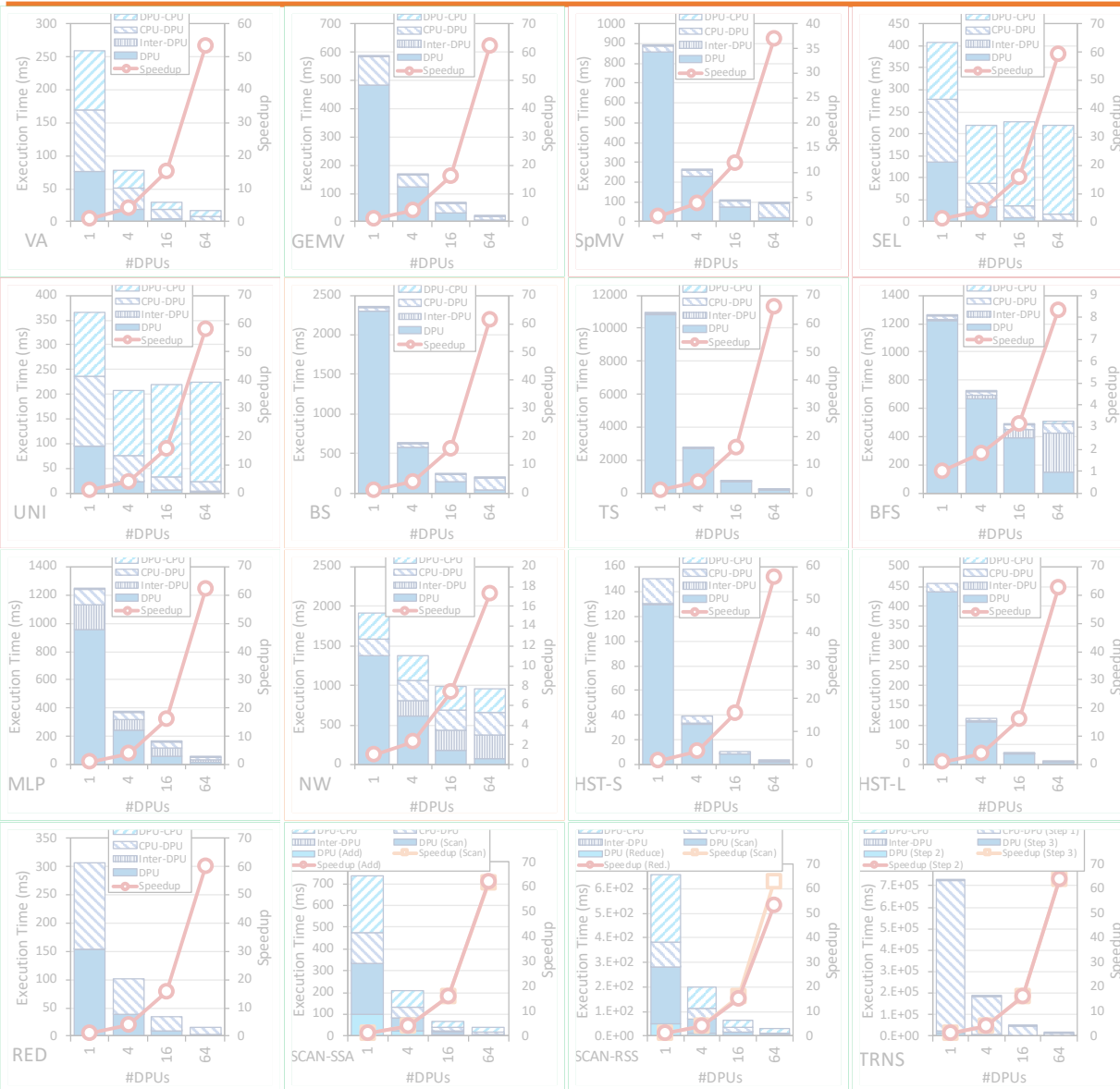
2,048-byte transfers are only 4% faster than 1,024-byte transfers

Larger transfers require more WRAM, which may limit the number of tasklets

## ***PROGRAMMING RECOMMENDATION 3***

**Choose the data transfer size between the MRAM bank and the WRAM based on the program's WRAM usage**, as it imposes a tradeoff between the sustained MRAM bandwidth and the number of tasklets that can run in the DPU (which is dictated by the limited WRAM capacity).

# Strong Scaling: 1 Rank (IV)



VA, GEMV, TS, MLP, HST-S, HST-L, RED, SCAN-SSA, SCAN-RSS, TRNS **use parallel transfers.**

CPU-DPU and DPU-CPU transfer times decrease as we increase the number of DPUs

BS, NW **use parallel transfers but do not reduce transfer times:**

- BS transfers a complete array to all DPUs.
- NW does not use all DPUs in all iterations

SpMV, SEL, UNI, BFS **cannot use parallel transfers**, as the transfer size per DPU is not fixed

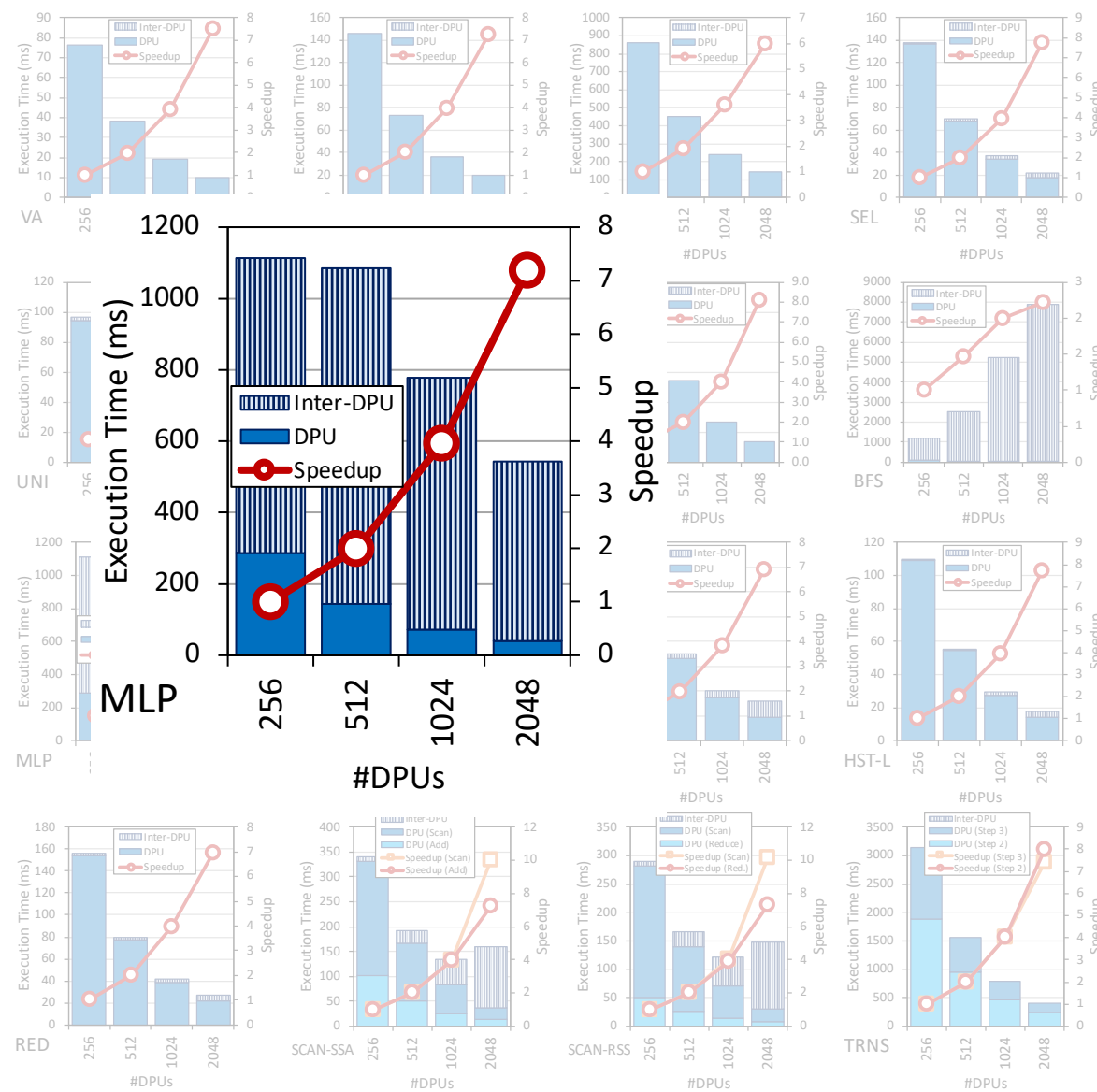
## PROGRAMMING RECOMMENDATION 5

**Parallel CPU-DPU/DPU-CPU transfers inside a rank of DPUs are recommended for real-world workloads when all transferred buffers are of the same size.**

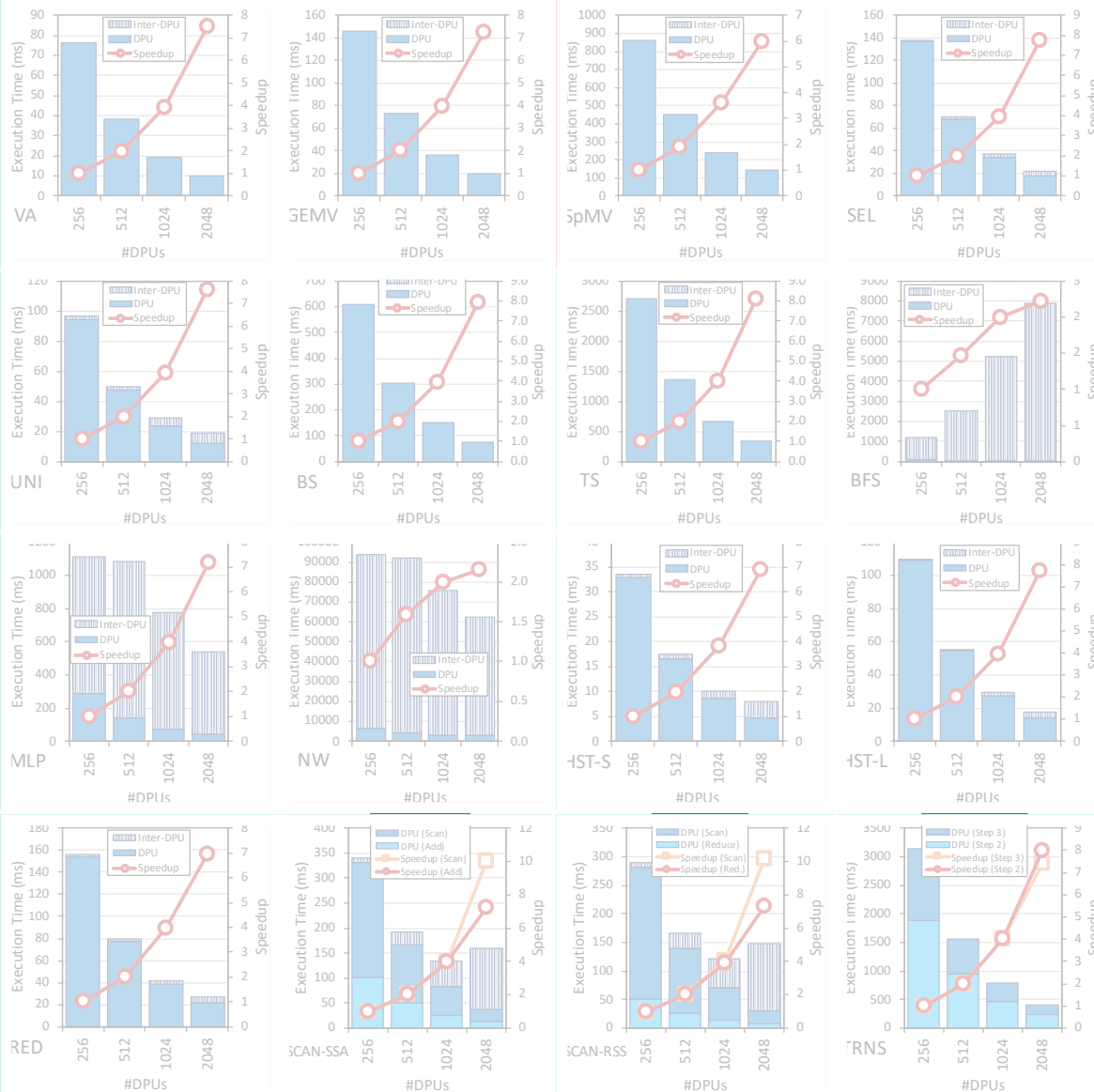


# Strong Scaling: 32 Ranks (I)

- Strong scaling experiments on 32 rank
  - We set the number of tasklets to the best performing one
  - The number of DPUs is 256, 512, 1024, 2048
  - We show the breakdown of execution time:
    - DPU: Execution time on the DPU
    - Inter-DPU: Time for inter-DPU communication via the host CPU
    - We do not show CPU-DPU/DPU-CPU transfer times
  - Speedup over 256 DPUs



# Strong Scaling: 32 Ranks (II)



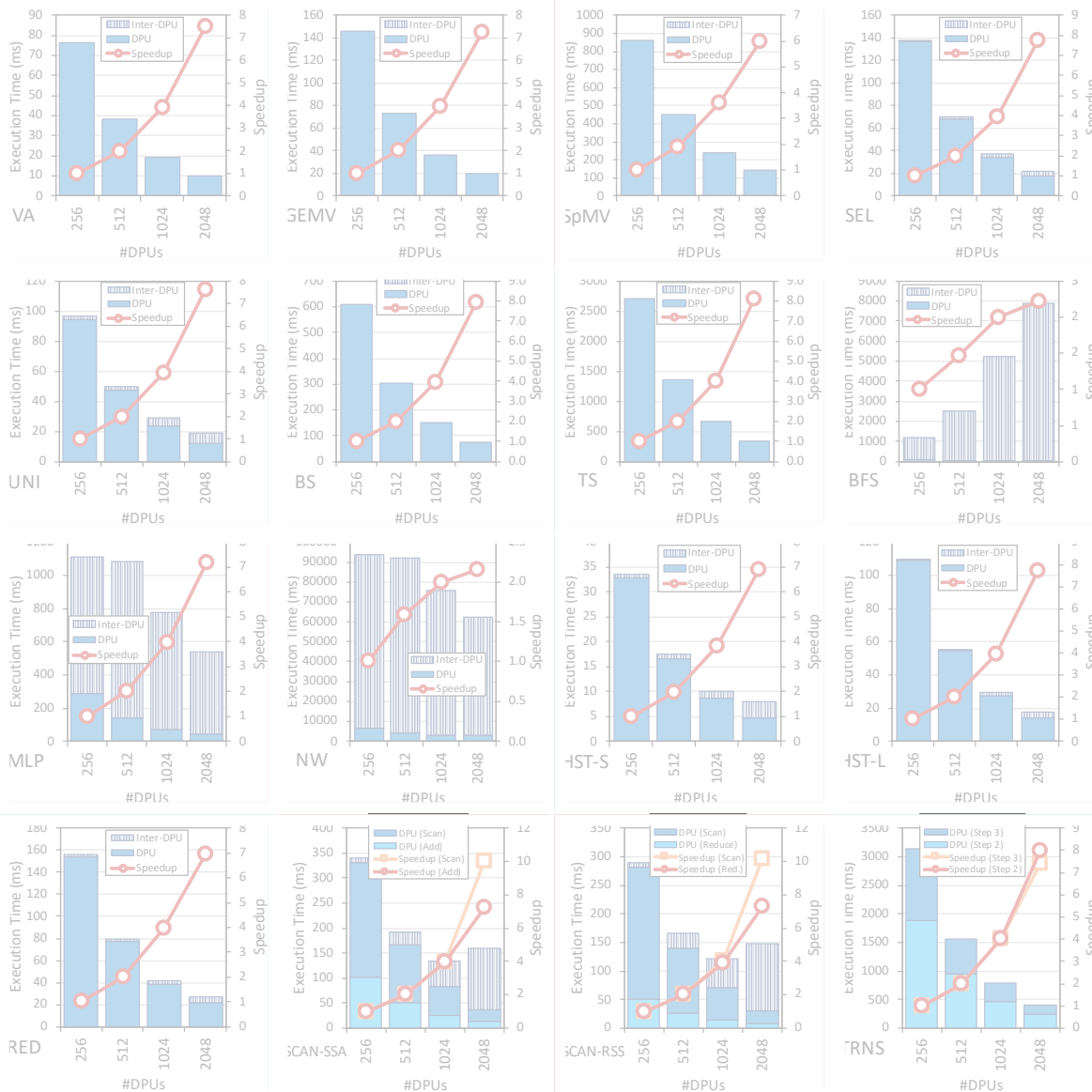
VA, GEMV, SEL, UNI, BS, TS, MLP, HST-S, HSTS-L, RED, SCAN-SSA (both kernel), SCAN-RSS (both kernels), and TRNS (both kernels) scale linearly with the number of DPUs

SpMV, BFS, NW do not scale linearly due to load imbalance

## KEY OBSERVATION 14

**Load balancing across DPUs ensures linear reduction of the execution time spent on the DPUs for a given problem size, when all available DPUs are used (as observed in strong scaling experiments).**

# Strong Scaling: 32 Ranks (III)



SEL, UNI, HST-S, HST-L, RED only need to merge final results

## KEY OBSERVATION 15

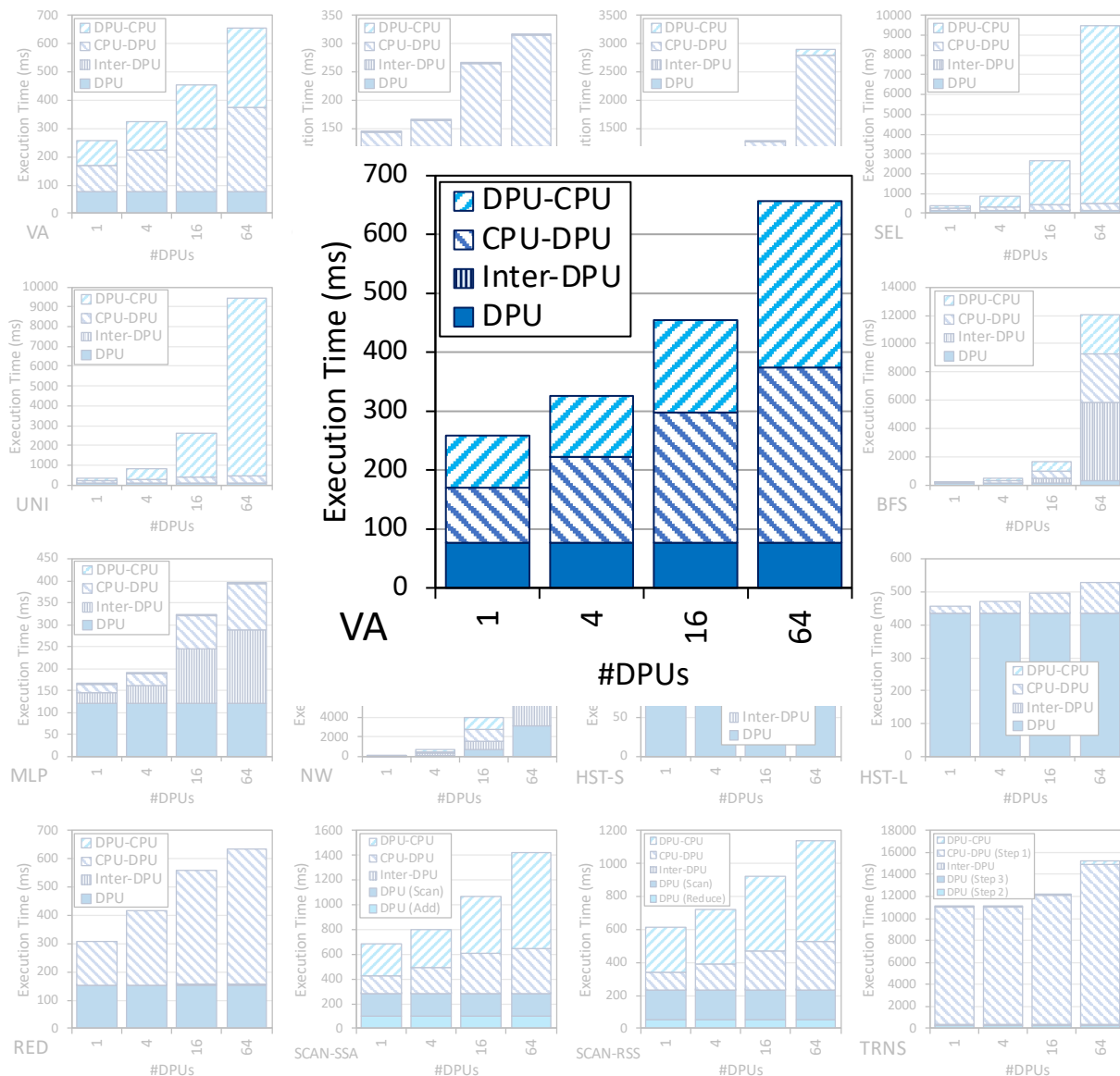
The overhead of merging partial results from DPUs in the host CPU is tolerable across all PRIM benchmarks that need it.

BFS, MLP, NW, SCAN-SSA, SCAN-RSS have more complex communication

## KEY OBSERVATION 16

Complex synchronization across DPUs (i.e., inter-DPU synchronization involving two-way communication with the host CPU) imposes significant overhead, which limits scalability to more DPUs.

# Weak Scaling: 1 Rank



## KEY OBSERVATION 17

**Equally-sized problems assigned to different DPUs and little/no inter-DPU synchronization lead to linear weak scaling of the execution time spent on the DPUs (i.e., constant execution time when we increase the number of DPUs and the dataset size accordingly).**

## KEY OBSERVATION 18

**Sustained bandwidth of parallel CPU-DPU/DPU-CPU transfers inside a rank of DPUs increases sublinearly with the number of DPUs.**

# Resources

---

- UPMEM SDK documentation
  - [https://sdk.upmem.com/master/00\\_ToolchainAtAGlance.html](https://sdk.upmem.com/master/00_ToolchainAtAGlance.html)
- Fabrice Devaux's presentation at HotChips 2019
  - <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8875680>
- Onur's lectures and talks

# Characterization of UPMEM PIM

---

- Microbenchmarks
  - Pipeline throughput
  - STREAM benchmark: WRAM, MRAM
  - Strided accesses and GUPS
  - Throughput vs. Operational intensity
  - CPU-DPU data transfers
- Real-world benchmarks
  - Dense linear algebra
  - Sparse linear algebra
  - Databases
  - Graph processing
  - Bioinformatics
  - Etc.

# Banner Colors

---

This is a question or an observation

This is an answer from, e.g., UPMEM documentation or our own research

This is an idea or a discussion starter, an opportunity for brainstorming

# DPU Sharing? Security Implications?

---

- DPUs cannot be shared across multiple CPU processes
  - There are so many DPUs in the system that there is no need for sharing
- According to UPMEM, this assumption makes things simpler

Is it possible to perform RowHammer bit flips?  
Can we attack the previous or the next application  
that runs on a DPU?

RowHammer patents and Giray's paper?



# More Questions and Ideas?

How do we handle memory coherence,  
memory oversubscription, etc.?

They are programmer's responsibility

A software library to handle  
**memory management transparently** to programmers

ASPLOS 2010

## **An Asymmetric Distributed Shared Memory Model for Heterogeneous Parallel Systems**

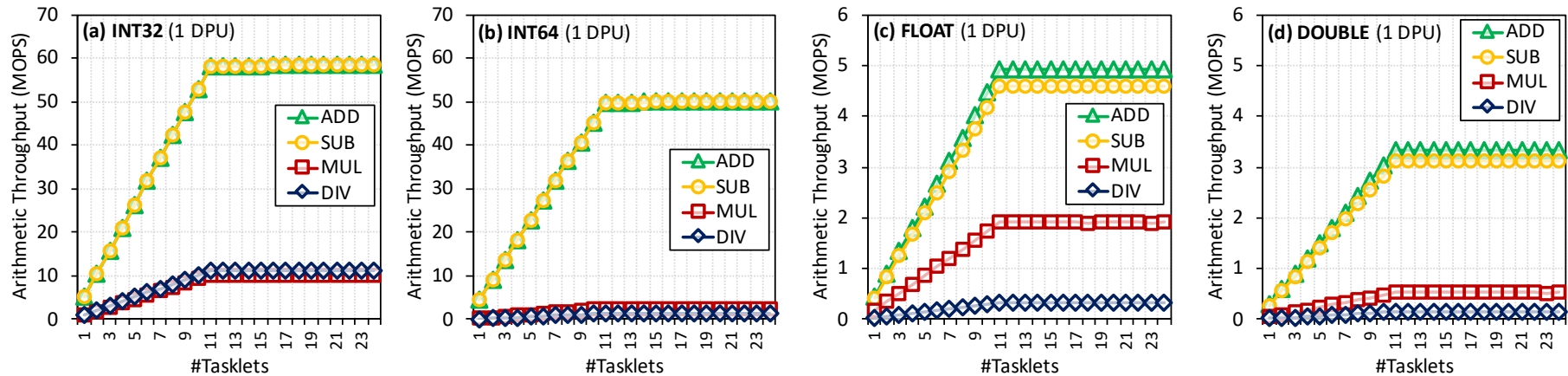
Isaac Gelado    Javier Cabezas  
                    Nacho Navarro

Universitat Politecnica de Catalunya  
{igelado, jcabezas, nacho}@ac.upc.edu

John E. Stone    Sanjay Patel  
                    Wen-mei W. Hwu

University of Illinois  
{jestone, sjp, hwu}@illinois.edu

# Arithmetic Throughput (II)

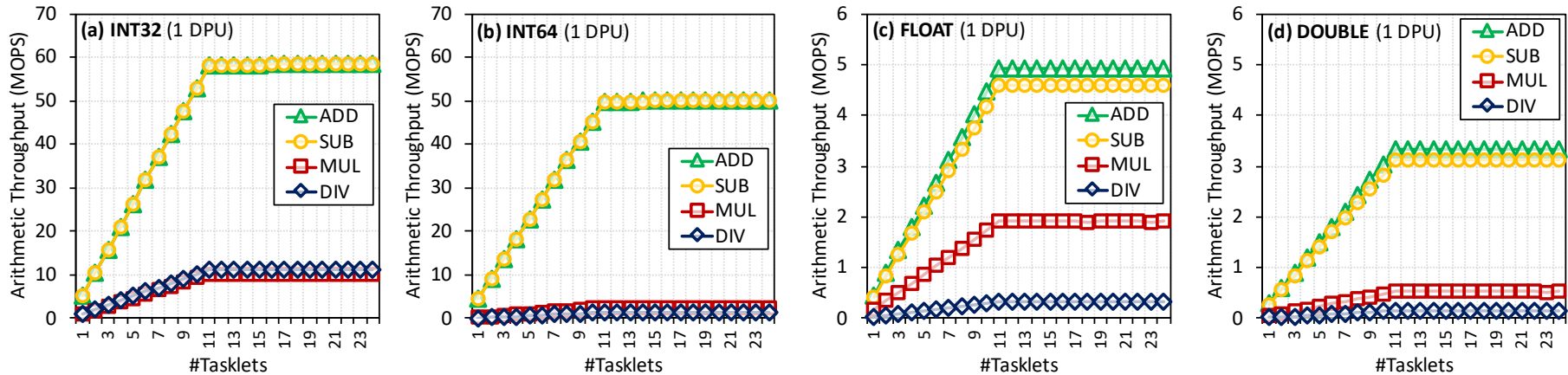


Huge throughput difference between add/sub and mul/div

DPU's do not have a 32-bit multiplier.  
mul/div implementation is based on bit shifting and addition:  
maximum of 32 cycles (instructions) to complete

There is an 8-bit multiplier in the pipeline.  
Would it be possible to use it for more efficient implementation?

# Arithmetic Throughput (III)

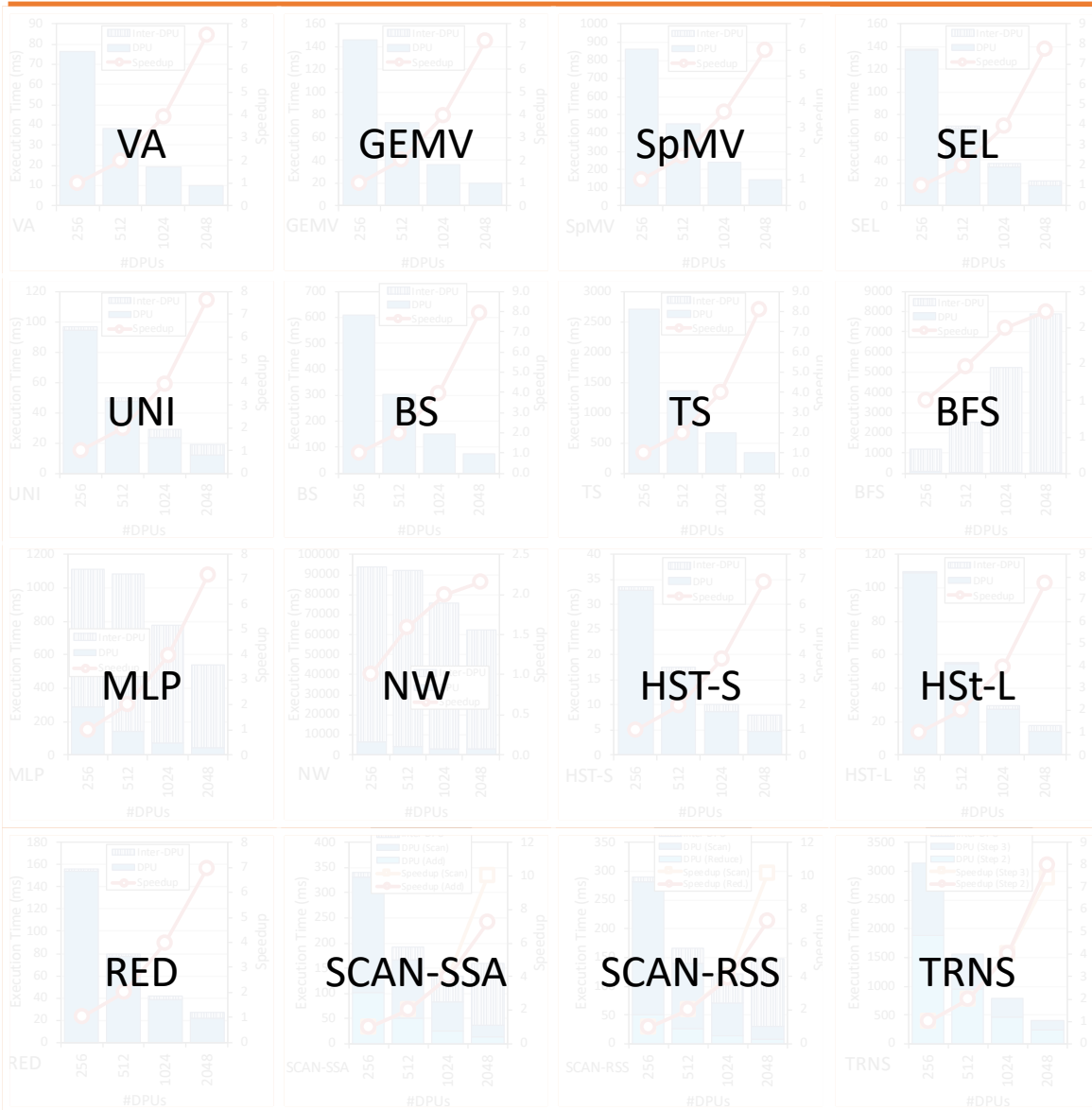


Huge throughput difference between  
int32/int64 and float/double

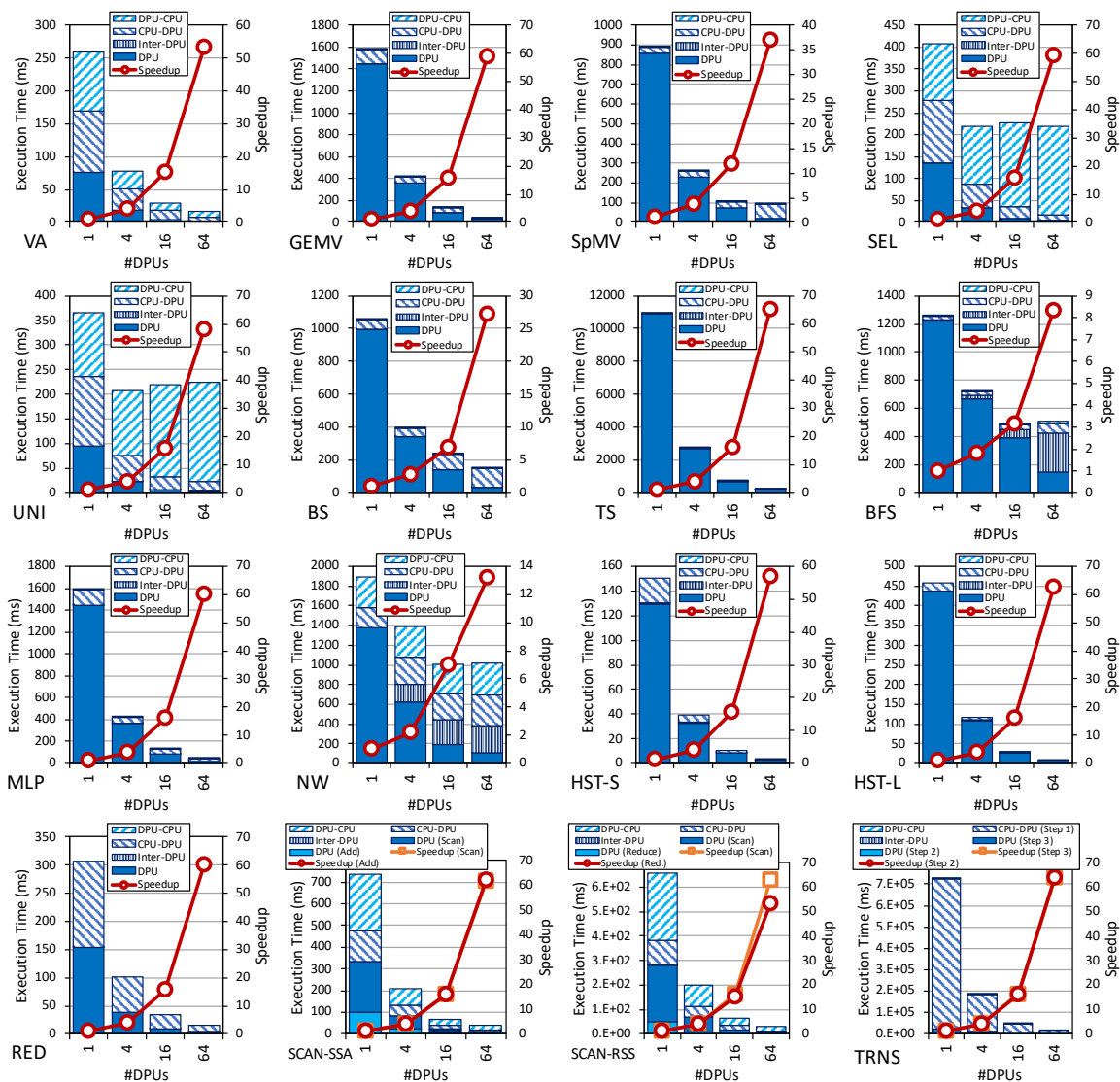
DPU's do not have floating point units.  
Software emulation for floating point computations

More efficient algorithms based on **other formats**?  
E.g., posit, TF32?

# Strong Scaling: 32 Ranks



# Strong Scaling: 1 Rank



# DSLs, High-level Programming

---

- Tangram

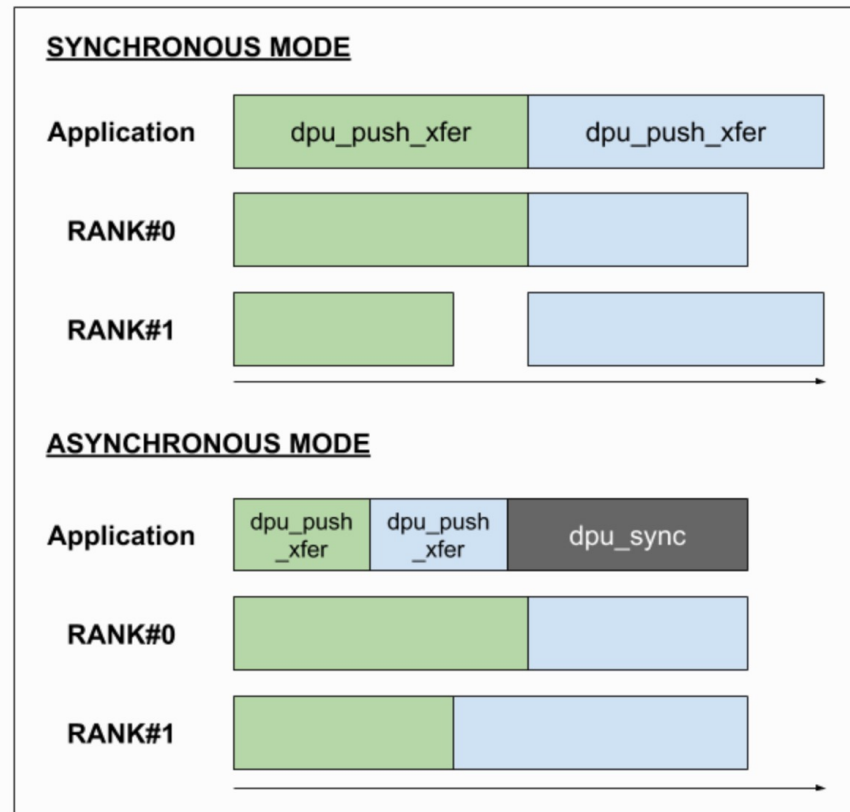
# Recap

---

It is possible:  
More complex benchmarks with **task-level parallelism**

# Backup: CPU-DPU Data Transfers

- Parallel asynchronous mode
  - Two transfers to a set of two ranks



[https://sdk.upmem.com/master/032\\_DPURuntimeService\\_HostCommunication.html#dpu-rank-transfer-interface-label](https://sdk.upmem.com/master/032_DPURuntimeService_HostCommunication.html#dpu-rank-transfer-interface-label)



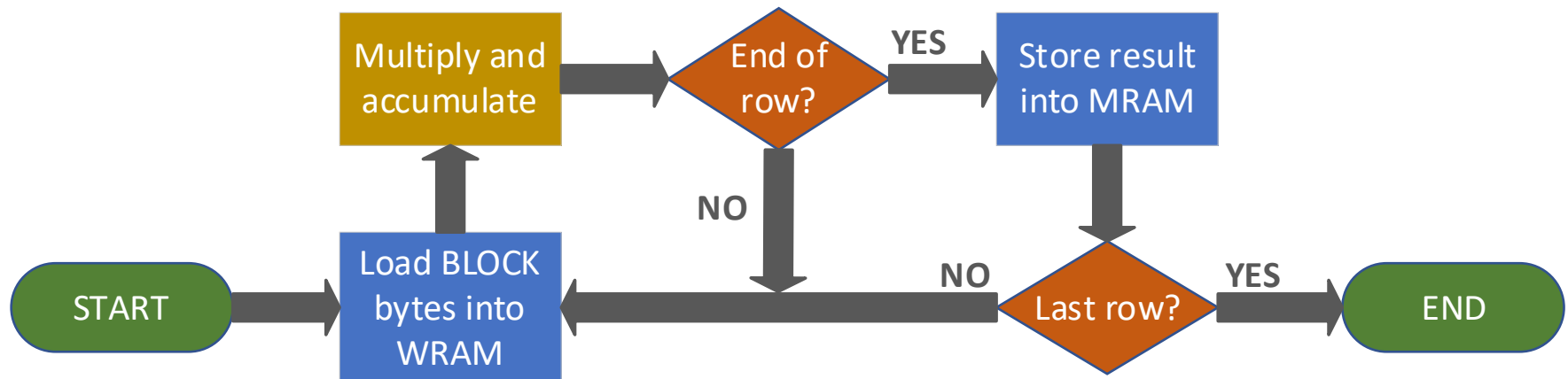
# GEMV: Parallelization Approach

- GEMV (general matrix-vector multiplication)

$$\begin{bmatrix} 1 & 0 & 2 & 0 \\ 0 & 3 & 0 & 4 \\ 0 & 0 & 5 & 0 \\ 6 & 0 & 0 & 7 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 5 \\ 1 \\ 8 \end{bmatrix} = \begin{bmatrix} 4 \\ 47 \\ 5 \\ 68 \end{bmatrix}$$

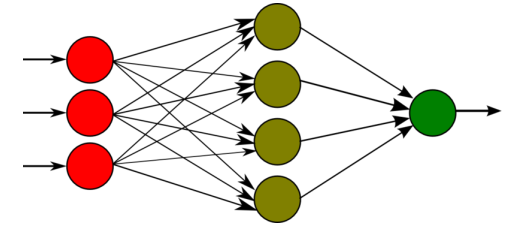
- Workload distribution

- $\text{chunk\_size} = (\text{num\_rows} / (\text{nr\_ranks} * \text{nr\_dpus}))$ , to each DPU
- $\text{chunk\_size} / \text{NR\_TASKLETS}$ , to each tasklet

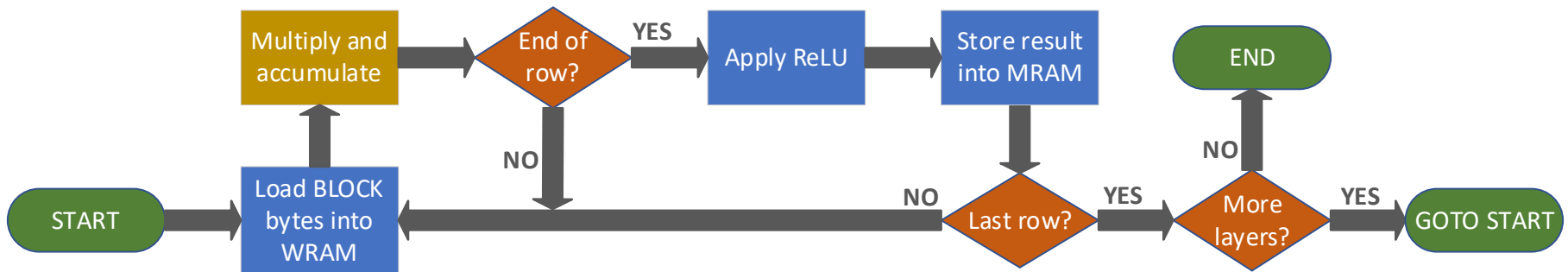


# MLP: Parallelization Approach

- MLP (multi-layer perceptron), based on GEMV

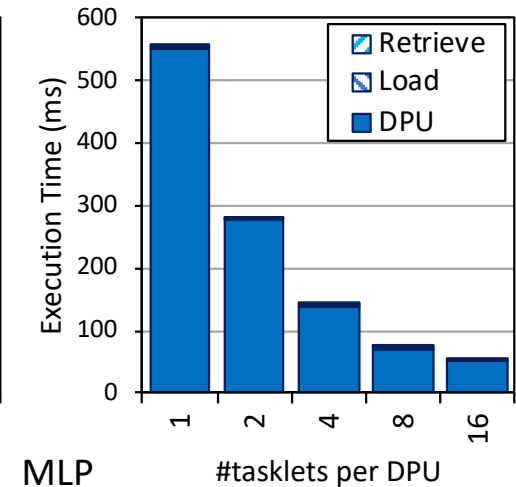
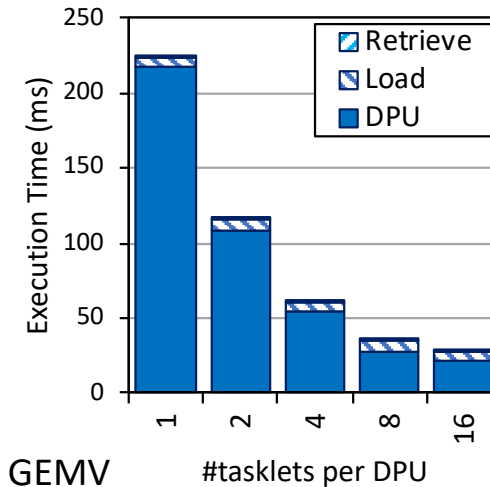


- Workload distribution
  - $\text{chunk\_size} = (\text{num\_rows} / (\text{nr\_ranks} * \text{nr\_dpus}))$ , to each DPU
  - $\text{chunk\_size} / \text{NR\_TASKLETS}$ , to each tasklet

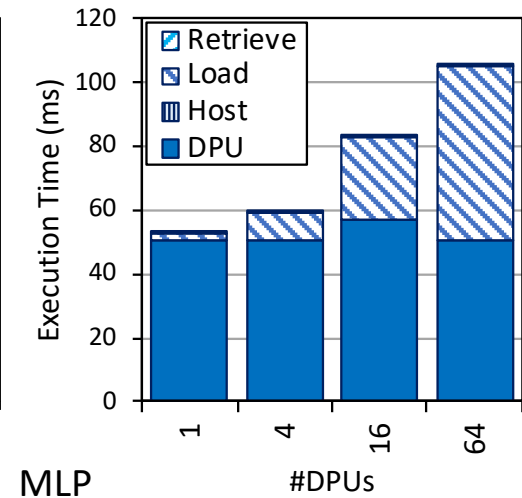
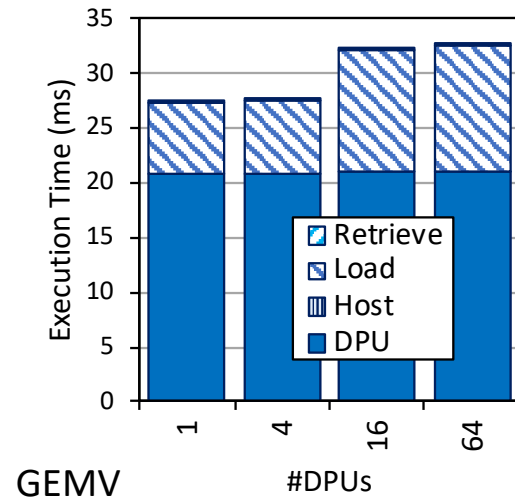


# Performance Scaling Results

- Strong scaling



- Weak scaling



# PIM Review and Open Problems

## Processing Data Where It Makes Sense: Enabling In-Memory Computation

Onur Mutlu<sup>a,b</sup>, Saugata Ghose<sup>b</sup>, Juan Gómez-Luna<sup>a</sup>, Rachata Ausavarungnirun<sup>b,c</sup>

<sup>a</sup>*ETH Zürich*

<sup>b</sup>*Carnegie Mellon University*

<sup>c</sup>*King Mongkut's University of Technology North Bangkok*

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun,  
**"Processing Data Where It Makes Sense: Enabling In-Memory  
Computation"**

*Invited paper in Microprocessors and Microsystems (**MICPRO**), June 2019.  
[[arXiv version](#)]*

# PIM Review and Open Problems (II)

## **A Workload and Programming Ease Driven Perspective of Processing-in-Memory**

Saugata Ghose<sup>†</sup>   Amirali Boroumand<sup>†</sup>   Jeremie S. Kim<sup>†§</sup>   Juan Gómez-Luna<sup>§</sup>   Onur Mutlu<sup>§†</sup>

<sup>†</sup>*Carnegie Mellon University*

<sup>§</sup>*ETH Zürich*

Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gomez-Luna, and Onur Mutlu,

**"Processing-in-Memory: A Workload-Driven Perspective"**

*Invited Article in [IBM Journal of Research & Development](#), Special Issue on Hardware for Artificial Intelligence, to appear in November 2019.*

[\[Preliminary arXiv version\]](#)