

# Computer Architecture

## Lecture 4a: Programming a Real-world PIM Architecture

Prof. Onur Mutlu

ETH Zürich

Fall 2023

6 October 2023

# Sub-Agenda: In-Memory Computation

---

- Major Trends Affecting Main Memory
- The Need for Intelligent Memory Controllers
  - Bottom Up: Push from Circuits and Devices
  - Top Down: Pull from Systems and Applications
- Processing in Memory: Two Directions
  - **Processing using Memory**
  - Processing near Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

# PIM Review and Open Problems

---

## A Modern Primer on Processing in Memory

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

*SAFARI Research Group*

<sup>a</sup>*ETH Zürich*

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

<sup>c</sup>*University of Illinois at Urbana-Champaign*

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

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun,

**"A Modern Primer on Processing in Memory"**

*Invited Book Chapter in Emerging Computing: From Devices to Systems - Looking Beyond Moore and Von Neumann, Springer, to be published in 2021.*

# Two PIM Approaches

## 5.2. Two Approaches: Processing Using Memory (PUM) vs. Processing Near Memory (PNM)

Many recent works take advantage of the memory technology innovations that we discuss in Section 5.1 to enable and implement PIM. We find that these works generally take one of two approaches, which are categorized in Table 1: (1) *processing using memory* or (2) *processing near memory*. We briefly describe each approach here. Sections 6 and 7 will provide example approaches and more detail for both.

Table 1: Summary of enabling technologies for the two approaches to PIM used by recent works. Adapted from [341] and extended.

Approach	Example Enabling Technologies
Processing Using Memory	SRAM DRAM Phase-change memory (PCM) Magnetic RAM (MRAM) Resistive RAM (RRAM)/memristors
Processing Near Memory	Logic layers in 3D-stacked memory Silicon interposers Logic in memory controllers Logic in memory chips (e.g., near bank) Logic in memory modules Logic near caches Logic near/in storage devices

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna,  
and Rachata Ausavarungnirun,

[\*\*"A Modern Primer on Processing in Memory"\*\*](#)

*Invited Book Chapter in [Emerging Computing: From Devices to Systems - Looking Beyond Moore and Von Neumann](#),*

Springer, to be published in 2021.

[\[Tutorial Video on "Memory-Centric Computing Systems" \(1 hour 51 minutes\)\]](#)

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Major Trends Affecting Main Memory</b>	<b>4</b>
<b>3</b>	<b>The Need for Intelligent Memory Controllers to Enhance Memory Scaling</b>	<b>6</b>
<b>4</b>	<b>Perils of Processor-Centric Design</b>	<b>9</b>
<b>5</b>	<b>Processing-in-Memory (PIM): Technology Enablers and Two Approaches</b>	<b>11</b>
5.1	New Technology Enablers: 3D-Stacked Memory and Non-Volatile Memory . . . . .	12
5.2	Two Approaches: Processing Using Memory (PUM) vs. Processing Near Memory (PNM) . . . . .	13
<b>6</b>	<b>Processing Using Memory (PUM)</b>	<b>14</b>
6.1	RowClone . . . . .	14
6.2	Ambit . . . . .	15
6.3	SIMDRAM . . . . .	17
6.4	Gather-Scatter DRAM . . . . .	18
6.5	In-DRAM Security Primitives . . . . .	18
<b>7</b>	<b>Processing Near Memory (PNM)</b>	<b>20</b>
7.1	Tesseract: Coarse-Grained Application-Level PNM Acceleration of Graph Processing . . . . .	20
7.2	Function-Level PNM Acceleration of Mobile Consumer Workloads . . . . .	21
7.3	Programmer-Transparent Function-Level PNM Acceleration of GPU Applications . . . . .	22
7.4	Instruction-Level PNM Acceleration with PIM-Enabled Instructions (PEI) . . . . .	23
7.5	Function-Level PNM Acceleration of Genome Analysis Workloads . . . . .	24
7.6	Application-Level PNM Acceleration of Time Series Analysis . . . . .	26
<b>8</b>	<b>Enabling the Adoption of PIM</b>	<b>26</b>
8.1	Programming Models and Code Generation for PIM . . . . .	26
8.2	PIM Runtime: Scheduling and Data Mapping . . . . .	27
8.3	Memory Coherence . . . . .	29
8.4	Virtual Memory Support . . . . .	30
8.5	Data Structures for PIM . . . . .	30
8.6	Benchmarks and Simulation Infrastructures . . . . .	31
8.7	Real PIM Hardware Systems and Prototypes . . . . .	33
8.8	Security Considerations . . . . .	36
<b>9</b>	<b>Other Resources on PIM</b>	<b>37</b>
<b>10</b>	<b>Conclusion and Future Outlook</b>	<b>37</b>

## 1. Introduction

Main memory, built using the Dynamic Random Access Memory (DRAM) technology, is a major component in nearly all computing systems, including servers, cloud platforms, mobile/embedded devices, and sensor systems. Across all of these systems, the data working set sizes of modern applications are rapidly growing, while the need for fast analysis of such data is increasing. Thus, main memory is becoming an increasingly significant bottleneck across a wide variety of computing systems and applications [1–26]. Alleviating the main memory bottleneck requires the memory capacity, energy, cost, and performance to all scale in an efficient manner across technology generations. Unfortunately, it has become increasingly difficult in recent years, especially the past decade, to scale all of these dimensions [1, 2, 27–59], and thus the main memory bottleneck has been worsening.

A major reason for the main memory bottleneck is the high energy and latency cost associated with *data movement*. In modern computers, to perform any operation on data that resides in main memory, the processor must retrieve the data from main memory. This requires the memory controller to issue commands to a DRAM module across a relatively slow and power-hungry off-chip bus (known as the *memory channel*). The DRAM module sends the requested data across the memory channel, after which the data is placed in the caches and registers. The CPU can perform computation on the data once the data is in its registers. Data movement from the DRAM to the CPU incurs long latency and consumes a significant amount of energy [7–9, 60–64]. These costs are often exacerbated by the fact that much of the data brought into the caches is *not reused* by the CPU [62, 63, 65, 66], providing little benefit in return for the high latency and energy cost.

The cost of data movement is a fundamental issue with the *processor-centric* nature of contemporary computer systems. The CPU is considered to be the master in the system, and computation is performed only in the processor (and accelerators). In contrast, data storage and communication units, including the main memory, are treated as unintelligent workers that are incapable of computation. As a result of this processor-centric design paradigm, data moves a lot in the system between the computation units and communication/storage units so that computation can be done on it. With the increasingly *data-centric* nature of contemporary and emerging applications, the processor-centric design paradigm leads to great inefficiency in performance, energy and cost. For example, most of the real estate within a single compute

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun,

[\*\*"A Modern Primer on Processing in Memory"\*\*](#)

*Invited Book Chapter in [Emerging Computing: From Devices to Systems - Looking Beyond Moore and Von Neumann](#),*

Springer, to be published in 2021.

[[Tutorial Video on "Memory-Centric Computing Systems"](#) (1 hour 51 minutes)]

# PIM Becomes Real

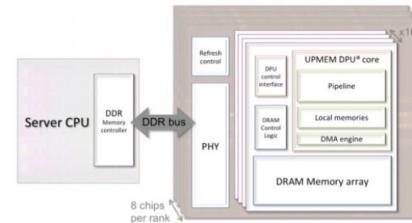
- UPMEM, founded in January 2015, announces the first real-world PIM architecture in 2016
- UPMEM's PIM-enabled DIMMs start getting commercialized in 2019
- In early 2021, Samsung announces FIMDRAM at ISSCC conference
- Samsung's LP-DDR5 and DIMM-based PIM announced a few months later
- In early 2022, SK Hynix announces AiM at ISSCC conference

The screenshot shows the header of the eeNews AUTOMOTIVE website. The logo 'eeNews' is in orange with 'AUTOMOTIVE' in blue below it. The navigation bar has tabs for 'ABOUT', 'FEATURED ARTICLES', 'LEARNING CENTER', and 'NEWS'. Below the navigation bar, there are links for 'About eeNews Europe Automotive' and 'Contact'.

Startup plans to embed processors in DRAM

October 13, 2016 // By Peter Clarke

[Email](#) [print](#) [Share](#) [in Share](#) [reddit](#)

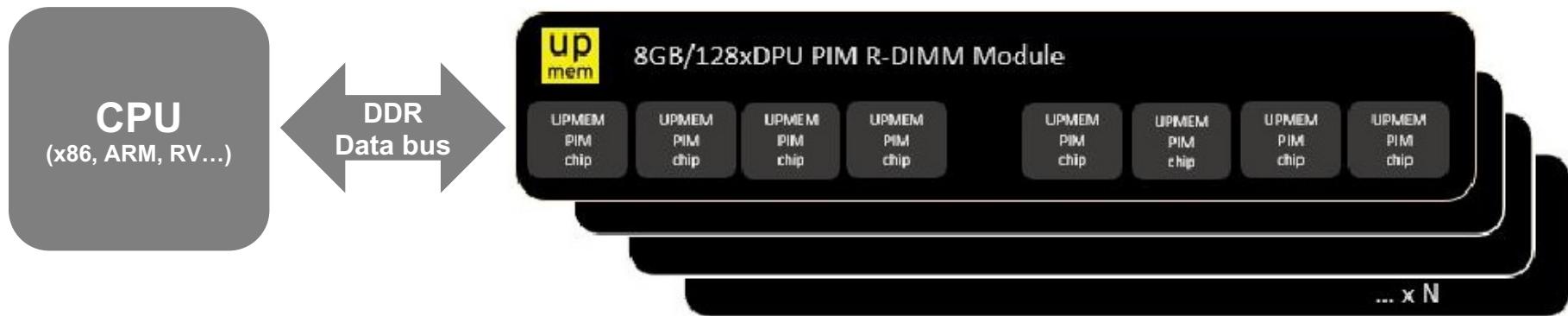


Fabless chip company Upmem SAS (Grenoble, France), founded in January 2015, is developing a microprocessor for use in data-intensive applications in the datacenter that will sit embedded in DRAM to be close to the data.

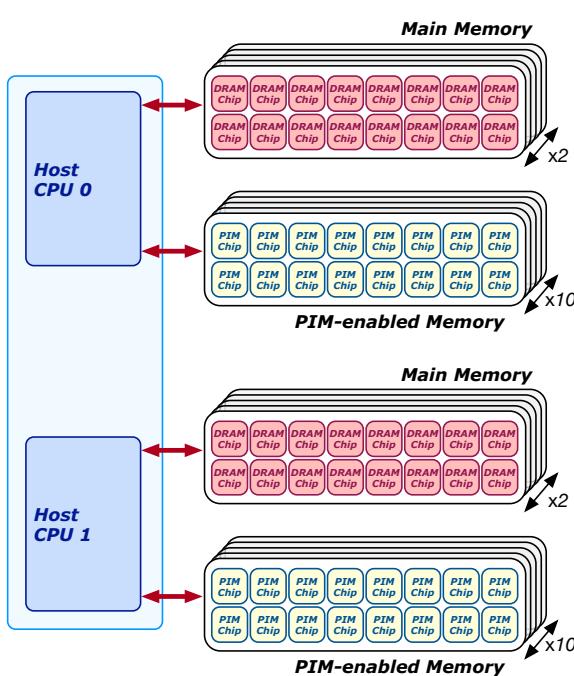
Placing hundreds or thousands of processing elements in DRAM able to perform work for a controlling server CPU could have a revolutionary impact on how data

# 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



# 2,560-DPU Processing-in-Memory System



Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture

JUAN GÓMEZ-LUNA, ETH Zürich, Switzerland

IZZAT EL HAJI, American University of Beirut, Lebanon

IVAN FERNANDEZ, ETH Zürich, Switzerland and University of Málaga, Spain

CHRISTINA GIANNOULA, ETH Zürich, Switzerland and NTUA, Greece

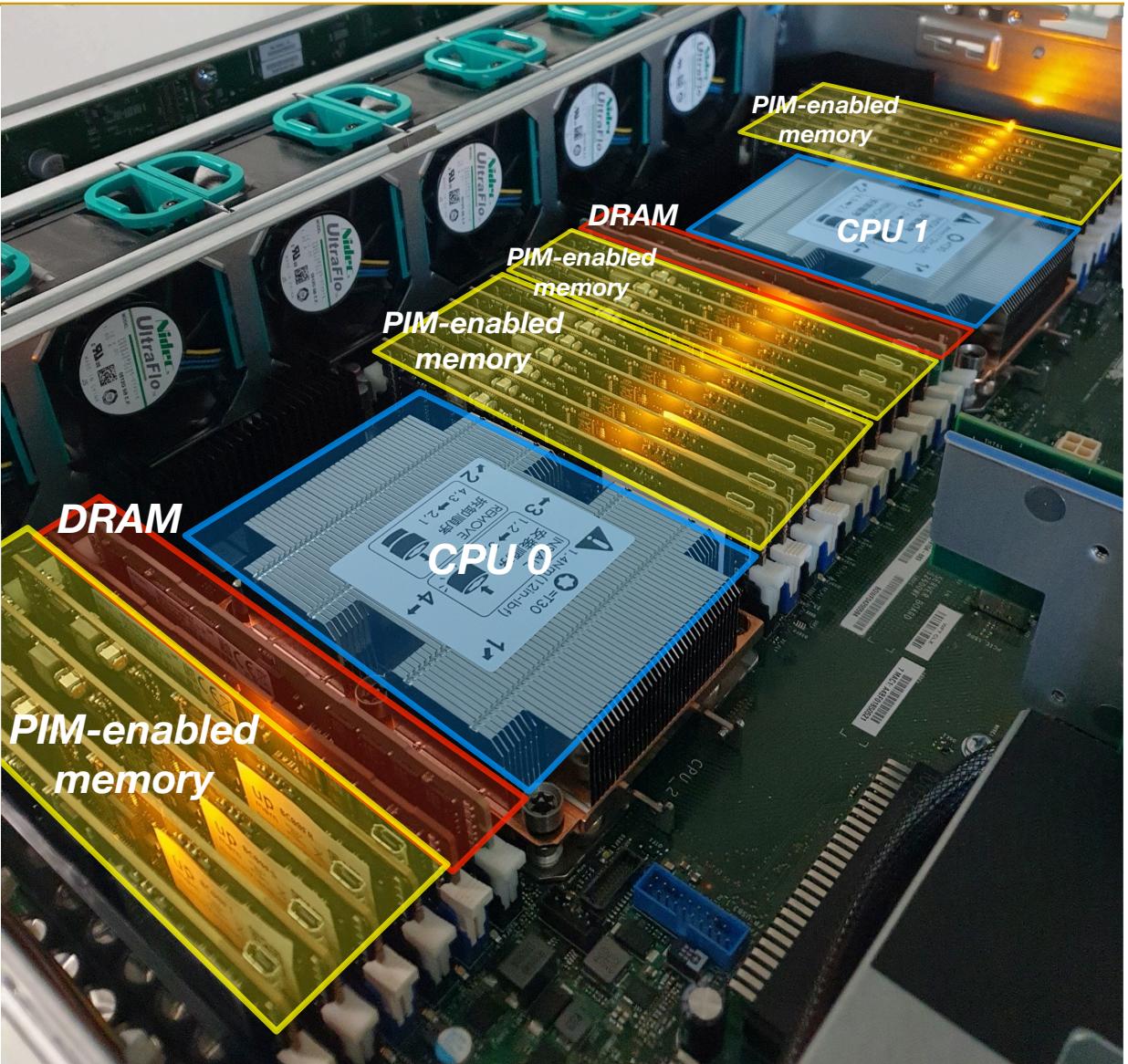
GERALDO F. OLIVEIRA, ETH Zürich, Switzerland

ONUR MUTLU, ETH Zürich, Switzerland

Many modern workloads, such as neural networks, databases, and graph processing, are fundamentally memory-bound. For such workloads, the data movement between main memory and CPU cores imposes a significant overhead in terms of both latency and energy. A major reason is that this communication happens through a narrow bus with high latency and limited bandwidth, and the low data reuse in memory-bound workloads is insufficient to amortize the cost of main memory access. Fundamentally addressing this data movement bottleneck requires a paradigm where the memory system assumes an active role in computing by integrating processing capabilities. This paradigm is known as *processing-in-memory* (PIM).

Recent research explores different forms of PIM architectures, motivated by the emergence of new 3D-stacked memory technologies that integrate memory with a logic layer where processing elements can be easily placed. Past works evaluate these architectures in simulation or, at best, with simplified hardware prototypes. In contrast, the UPMEM company has designed and manufactured the first publicly-available real-world 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.

This paper provides the first comprehensive analysis of the first publicly-available real-world PIM architecture. We make two key contributions. First, we conduct an experimental characterization of the UPMEM-based PIM system using microbenchmarks to assess various architecture limits such as compute throughput and memory bandwidth, yielding new insights. Second, we present PrIM (Processing-In-Memory benchmarks), a benchmark suite of 16 workloads from different application domains (e.g., dense/sparse linear algebra, databases, data analytics, graph processing, neural networks, bioinformatics, image processing), which we identify as memory-bound. We evaluate the performance and scaling characteristics of PrIM benchmarks on the UPMEM PIM architecture, and compare their performance and energy consumption to their state-of-the-art CPU and GPU counterparts. Our extensive evaluation conducted on two real UPMEM-based PIM systems with 640 and 2,556 DPUs provides new insights about suitability of different workloads to the PIM system, programming recommendations for software designers, and suggestions and hints for hardware and architecture designers of future PIM systems.



# Samsung Function-in-Memory DRAM (2021)



## Samsung Develops Industry's First High Bandwidth Memory with AI Processing Power

Korea on February 17, 2021

Audio



Share



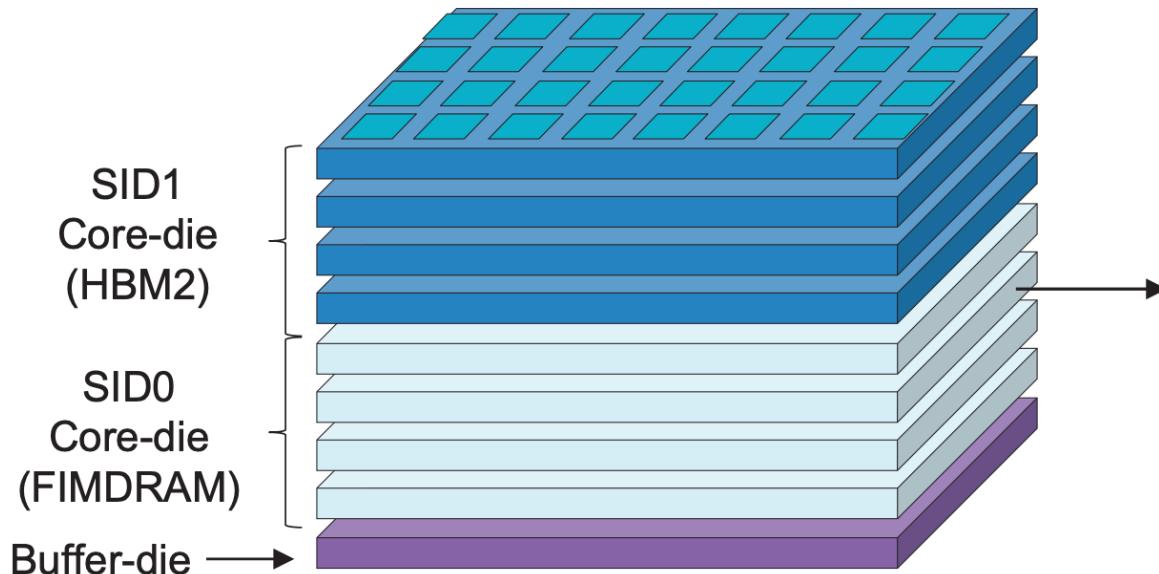
*The new architecture will deliver over twice the system performance and reduce energy consumption by more than 70%*

Samsung Electronics, the world leader in advanced memory technology, today announced that it has developed the industry's first High Bandwidth Memory (HBM) integrated with artificial intelligence (AI) processing power – the HBM-PIM. The new processing-in-memory (PIM) architecture brings powerful AI computing capabilities inside high-performance memory, to accelerate large-scale processing in data centers, high performance computing (HPC) systems and AI-enabled mobile applications.

Kwangil Park, senior vice president of Memory Product Planning at Samsung Electronics stated, "Our groundbreaking HBM-PIM is the industry's first programmable PIM solution tailored for diverse AI-driven workloads such as HPC, training and inference. We plan to build upon this breakthrough by further collaborating with AI solution providers for even more advanced PIM-powered applications."

# Samsung Function-in-Memory DRAM (2021)

## ■ FIMDRAM based on HBM2



[3D Chip Structure of HBM with FIMDRAM]

### Chip Specification

128DQ / 8CH / 16 banks / BL4

32 PCU blocks (1 FIM block/2 banks)

1.2 TFLOPS (4H)

**FP16 ADD /  
Multiply (MUL) /  
Multiply-Accumulate (MAC) /  
Multiply-and- Add (MAD)**

ISSCC 2021 / SESSION 25 / DRAM / 25.4

25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Cheon Kwon<sup>1</sup>, Suk Han Lee<sup>1</sup>, Jaehoon Lee<sup>1</sup>, Sang-Hyuk Kwon<sup>1</sup>, Je Min Ryu<sup>1</sup>, Jong-Pil Son<sup>1</sup>, Seongil Oh<sup>1</sup>, Hak-Soo Yu<sup>1</sup>, Haesuk Lee<sup>1</sup>, Soo Young Kim<sup>1</sup>, Youngmin Cho<sup>1</sup>, Jin Guk Kim<sup>1</sup>, Jongyoon Choi<sup>1</sup>, Hyun-Sung Shin<sup>1</sup>, Jin Kim<sup>1</sup>, BengSeng Phuah<sup>1</sup>, Hyo Young Kim<sup>1</sup>, Myeong Jun Song<sup>1</sup>, Ahn Choi<sup>1</sup>, Daeho Kim<sup>1</sup>, Soo Young Kim<sup>1</sup>, Eun-Bong Kim<sup>1</sup>, David Wang<sup>2</sup>, Shinhwa Kang<sup>1</sup>, Yuhwan Ro<sup>3</sup>, Seungwoo Seo<sup>3</sup>, JoonHo Song<sup>3</sup>, Jaeyoun Youn<sup>1</sup>, Kyomin Sohn<sup>1</sup>, Nam Sung Kim<sup>1</sup>

<sup>1</sup>Samsung Electronics, Hwaseong, Korea

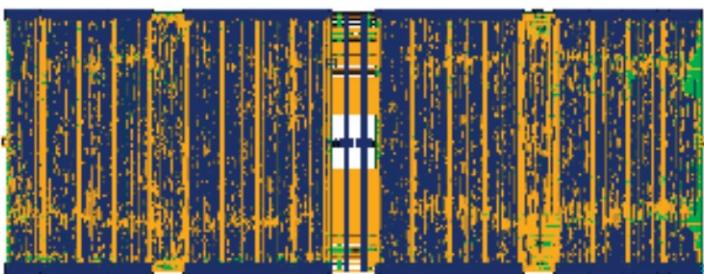
<sup>2</sup>Samsung Electronics, San Jose, CA

<sup>3</sup>Samsung Electronics, Suwon, Korea

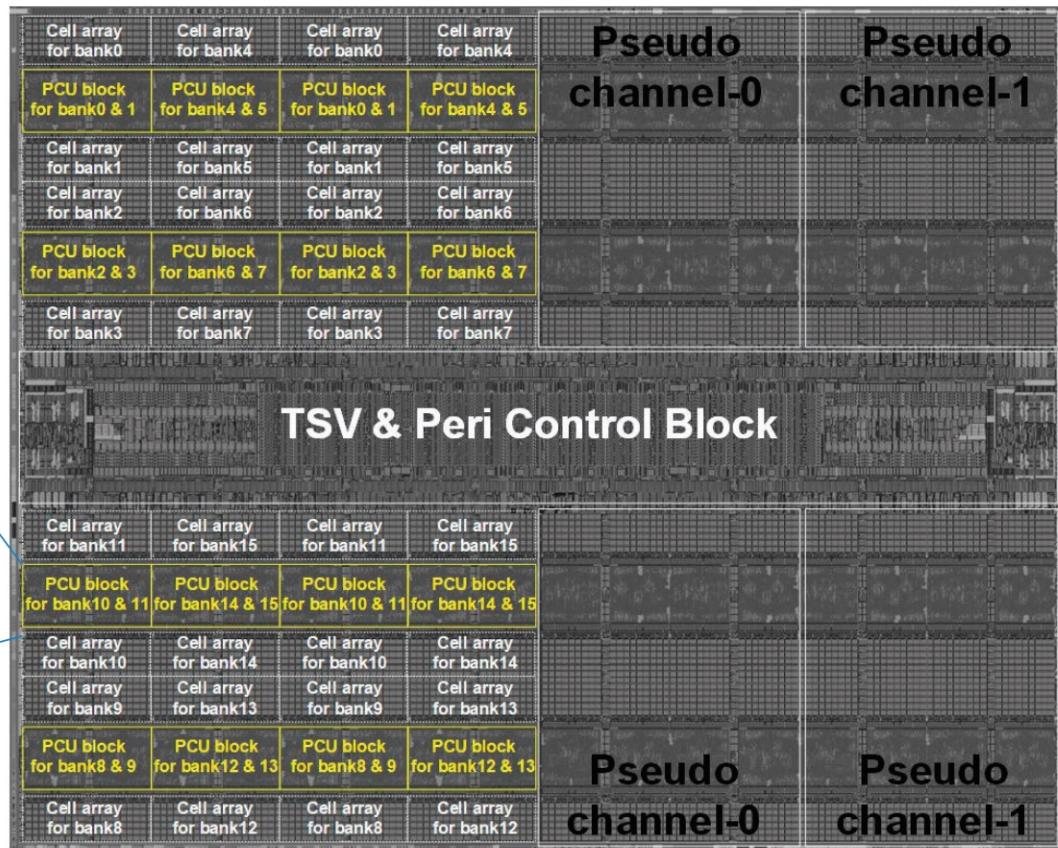
# Samsung Function-in-Memory DRAM (2021)

## Chip Implementation

- Mixed design methodology to implement FIMDRAM
  - Full-custom + Digital RTL



[Digital RTL design for PCU block]



ISSCC 2021 / SESSION 25 / DRAM / 25.4

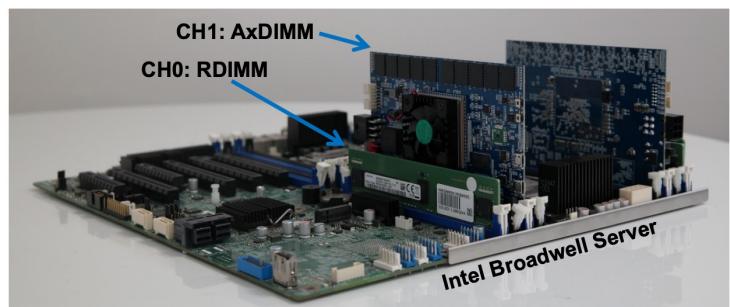
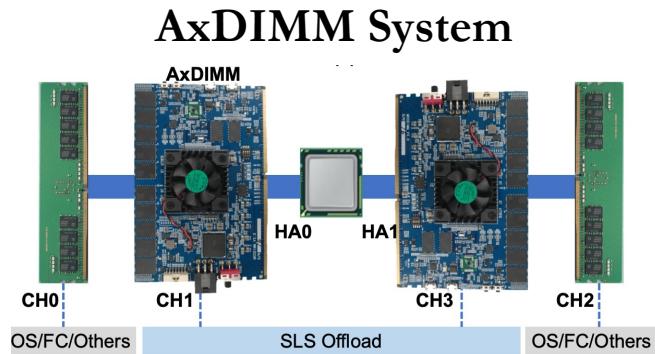
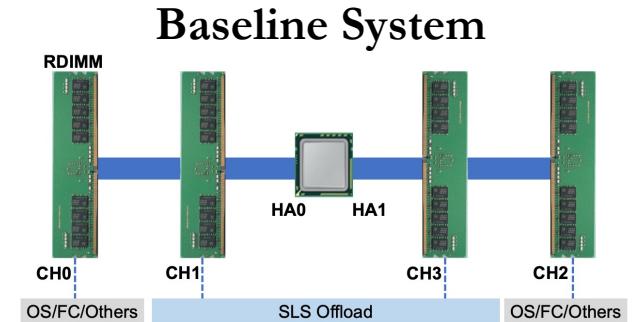
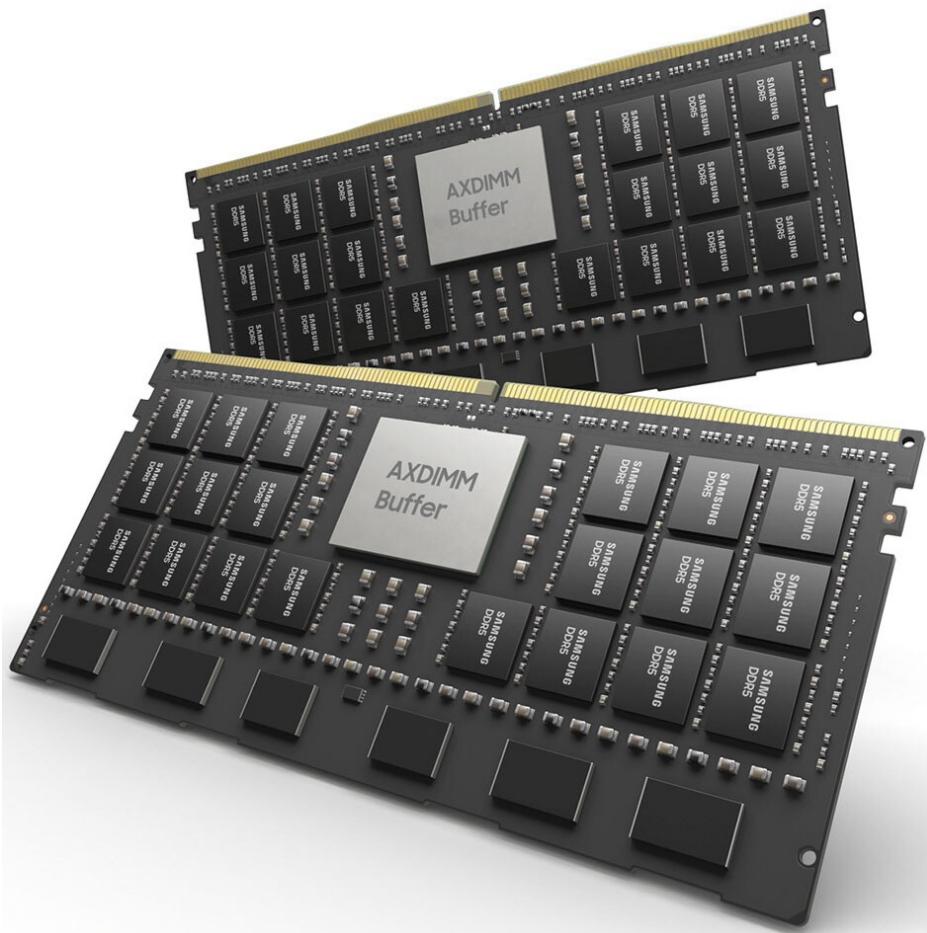
25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Cheon Kwon<sup>1</sup>, Suk Han Lee<sup>1</sup>, Jaehoon Lee<sup>1</sup>, Sang-Hyuk Kwon<sup>1</sup>, Je Min Ryu<sup>1</sup>, Jong-Pil Son<sup>1</sup>, Seongil O<sup>1</sup>, Hak-Soo Yu<sup>1</sup>, Haesuk Lee<sup>1</sup>, Soo Young Kim<sup>1</sup>, Youngmin Cho<sup>1</sup>, Jin Guk Kim<sup>1</sup>, Jongyoon Choi<sup>1</sup>, Hyun-Sung Shin<sup>1</sup>, Jin Kim<sup>1</sup>, BengSeng Phua<sup>2</sup>, HyoungMin Kim<sup>1</sup>, Myeong Jun Song<sup>1</sup>, Ahn Choi<sup>1</sup>, Daeho Kim<sup>1</sup>, SooYoung Kim<sup>1</sup>, Eun-Bong Kim<sup>1</sup>, David Wang<sup>2</sup>, Shinhwang Kang<sup>1</sup>, Yuhwan Ro<sup>1</sup>, Seungwoo Seo<sup>1</sup>, JoonHo Song<sup>1</sup>, Jaeyoun Youn<sup>1</sup>, Kyomin Sohn<sup>1</sup>, Nam Sung Kim<sup>1</sup>

<sup>1</sup>Samsung Electronics, Hwasung, Korea  
<sup>2</sup>Samsung Electronics, San Jose, CA  
<sup>3</sup>Samsung Electronics, Suwon, Korea

# Samsung AxDIMM (2021)

- DIMM-based PIM
  - DLRM recommendation system



# SK Hynix Accelerator-in-Memory (2022)

SK hynix NEWSROOM

ENG

INSIGHT SK hynix STORY PRESS CENTER MULTIMEDIA

Search



## SK hynix Develops PIM, Next-Generation AI Accelerator

February 16, 2022



Seoul, February 16, 2022

SK hynix (or “the Company”, [www.skhynix.com](http://www.skhynix.com)) announced on February 16 that it has developed PIM\*, a next-generation memory chip with computing capabilities.

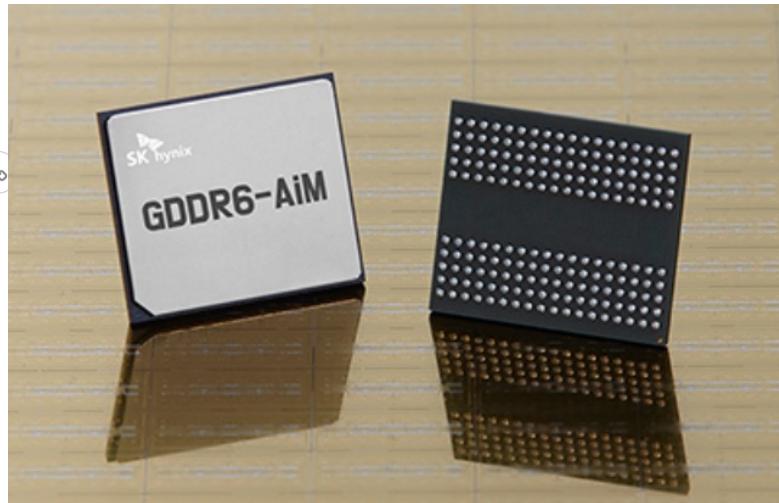
\*PIM(Processing In Memory): A next-generation technology that provides a solution for data congestion issues for AI and big data by adding computational functions to semiconductor memory

It has been generally accepted that memory chips store data and CPU or GPU, like human brain, process data. SK hynix, following its challenge to such notion and efforts to pursue innovation in the next-generation smart memory, has found a breakthrough solution with the development of the latest technology.

SK hynix plans to showcase its PIM development at the world's most prestigious semiconductor conference, 2022 ISSCC\*, in San Francisco at the end of this month. The company expects continued efforts for innovation of this technology to bring the memory-centric computing, in which semiconductor memory plays a central role, a step closer to the reality in devices such as smartphones.

\*ISSCC: The International Solid-State Circuits Conference will be held virtually from Feb. 20 to Feb. 24 this year with a theme of “Intelligent Silicon for a Sustainable World”

For the first product that adopts the PIM technology, SK hynix has developed a sample of GDDR6-AiM (Accelerator\* in memory). The GDDR6-AiM adds computational functions to GDDR6\* memory chips, which process data at 16Gbps. A combination of GDDR6-AiM with CPU or GPU instead of a typical DRAM makes certain computation speed 16 times faster. GDDR6-AiM is widely expected to be adopted for machine learning, high-performance computing, and big data computation and storage.



11.1 A 1nym 1.25V 8Gb, 16Gb/s/pin GDDR6-based Accelerator-in-Memory supporting 1TFLOPS MAC Operation and Various Activation Functions for Deep-Learning Applications

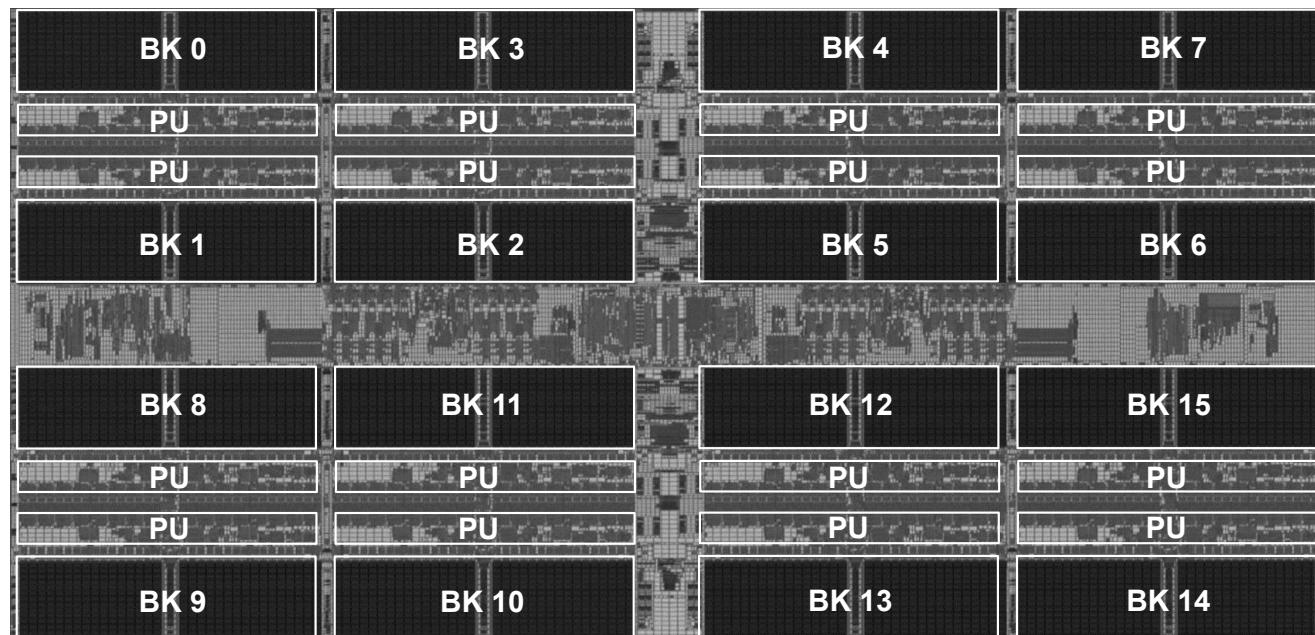
Seongju Lee, SK hynix, Icheon, Korea

In Paper 11.1, SK Hynix describes an 1nym, GDDR6-based accelerator-in-memory with a command set for deep-learning operation. The 8Gb design achieves a peak throughput of 1TFLOPS with 1GHz MAC operations and supports major activation functions to improve accuracy.

# SK Hynix AiM: Chip Implementation (2022)

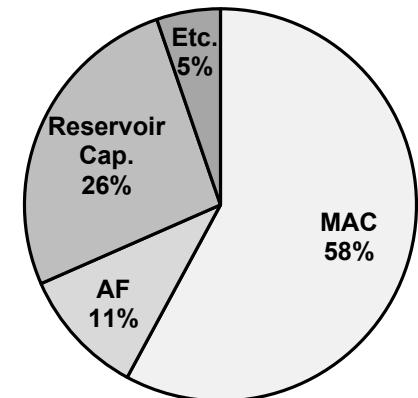
- 4 Gb AiM die with 16 processing units (PUs)

AiM Die Photograph



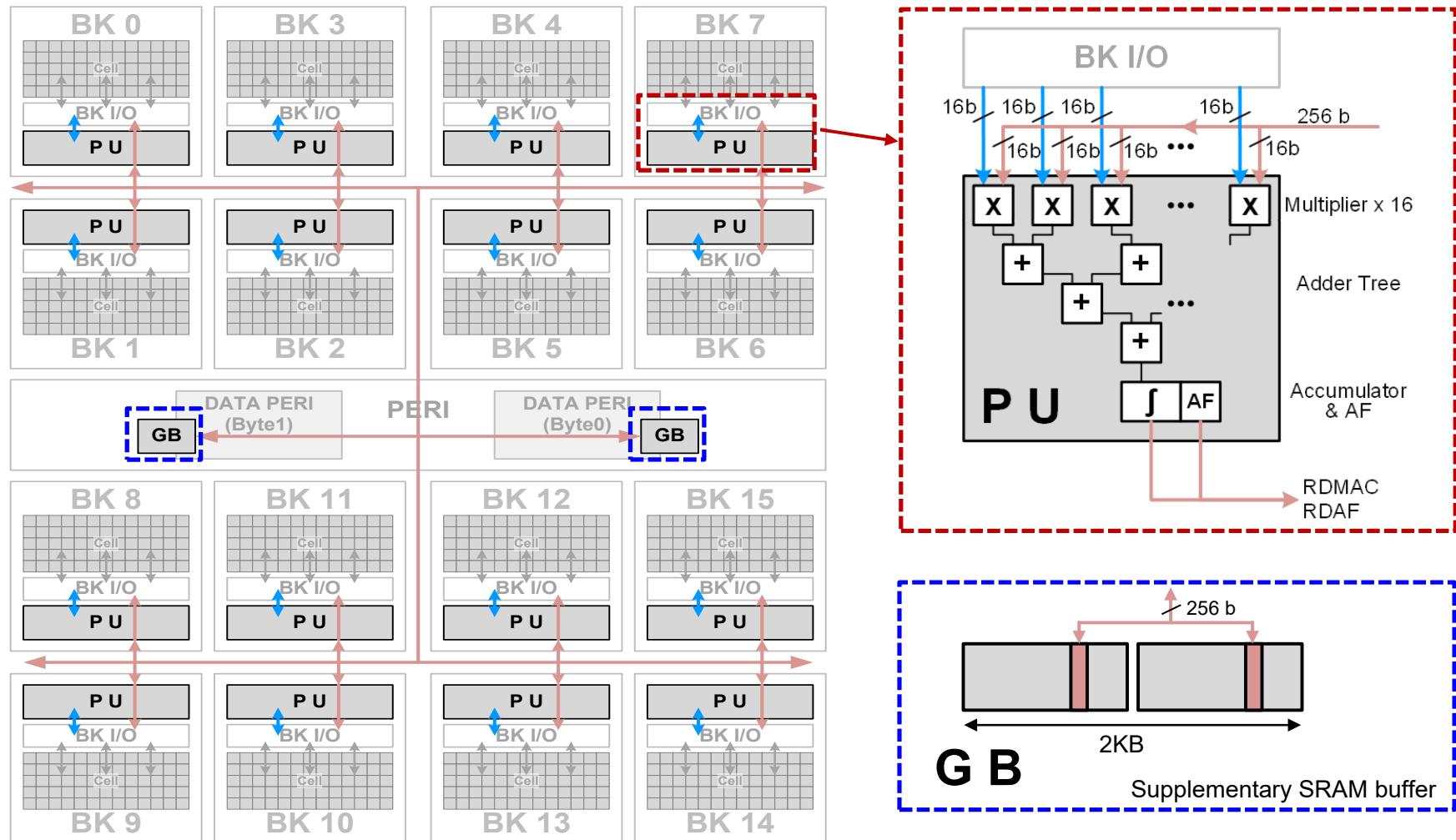
1 Process Unit (PU) Area

Total	0.19mm <sup>2</sup>
MAC	0.11mm <sup>2</sup>
Activation Function (AF)	0.02mm <sup>2</sup>
Reservoir Cap.	0.05mm <sup>2</sup>
Etc.	0.01mm <sup>2</sup>



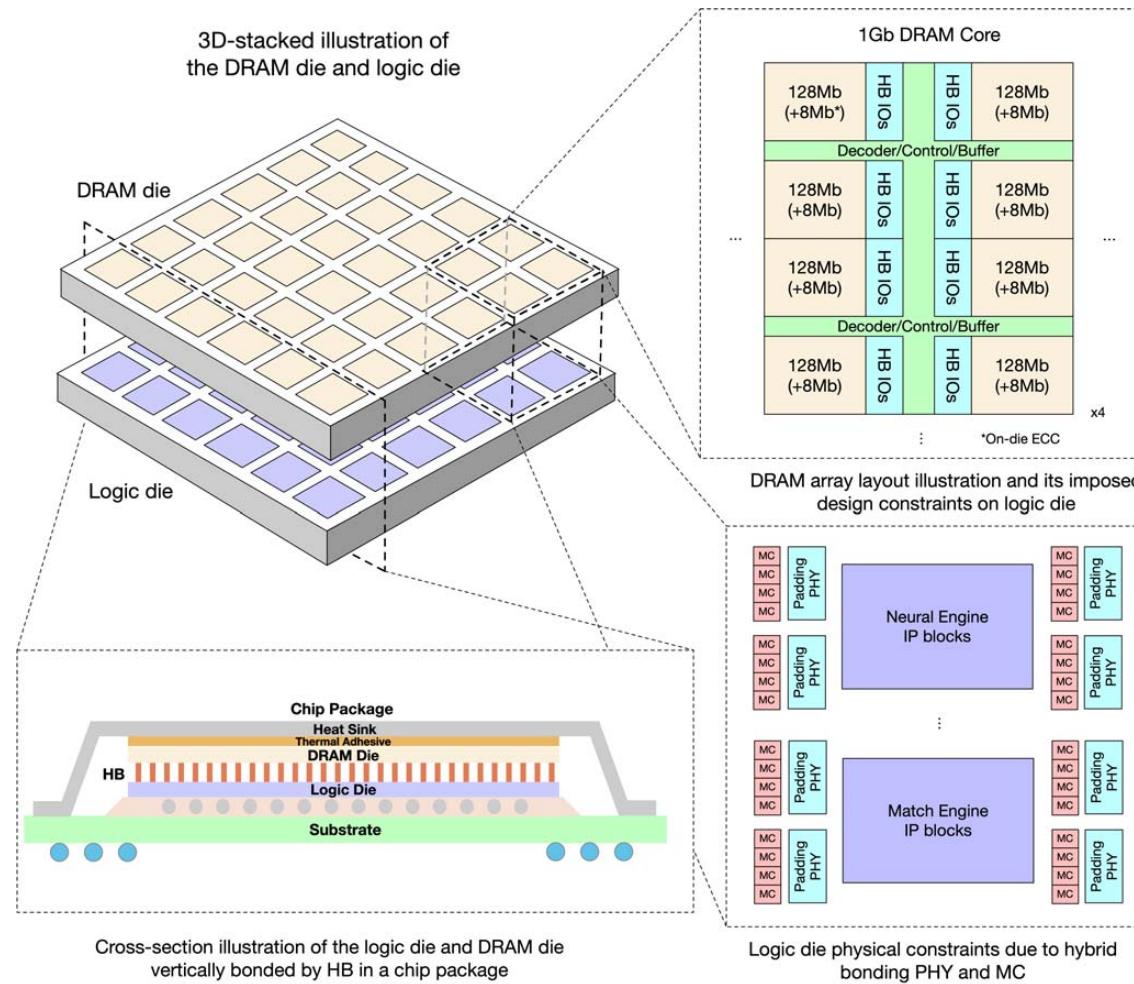
# SK Hynix AiM: System Organization (2022)

## GDDR6-based AiM architecture



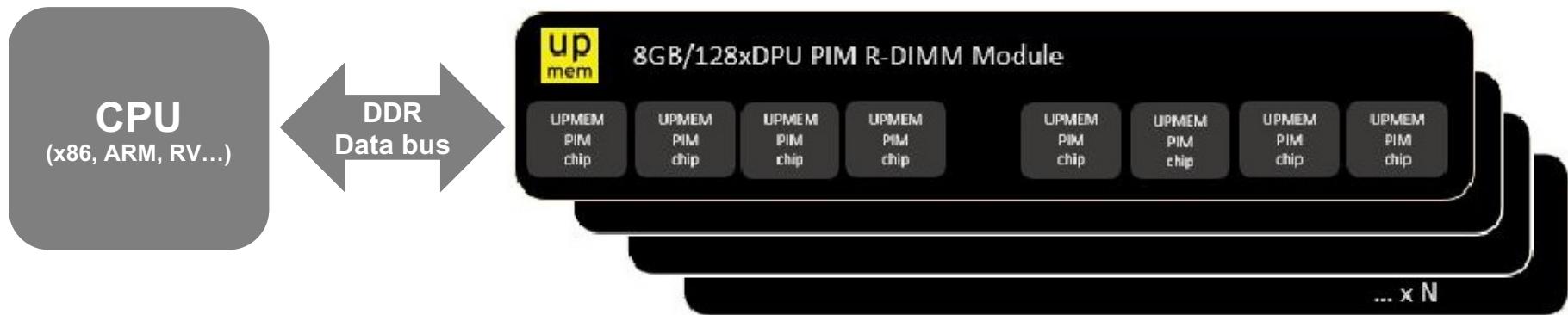
# Alibaba HB-PNM: Overall Architecture (2022)

- 3D-stacked logic die and DRAM die vertically bonded by hybrid bonding (HB)



# 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

---

## Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System

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>, AND 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

Corresponding author: Juan Gómez-Luna (e-mail: juang@ethz.ch).

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

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

# 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'. The repository has 2 stars, 1 fork, and 1 commit from Juan Gomez Luna. It contains 168 lines (132 sloc) and 5.79 KB of code. The README.md file is open, displaying the following content:

## 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.

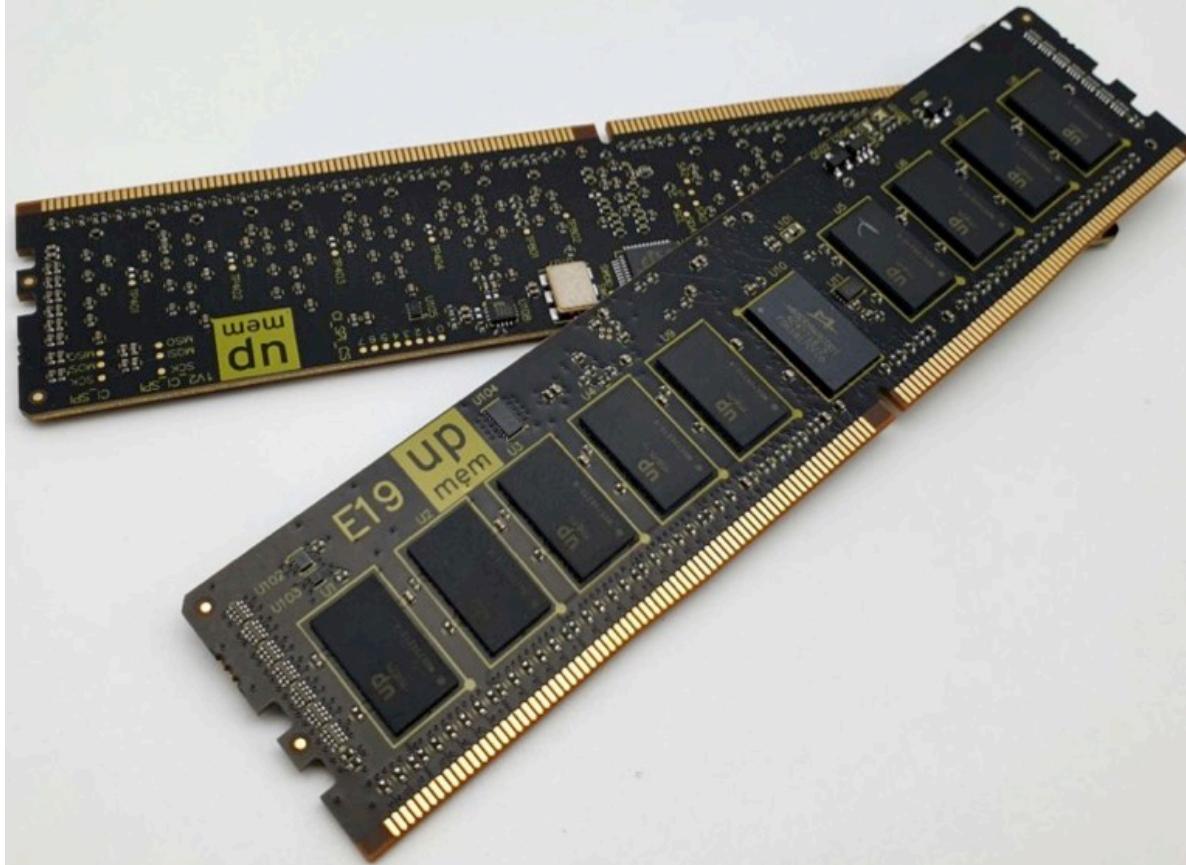
PrIM provides a common set of workloads to evaluate the UPMEM PIM architecture with and can be useful for programming, architecture and system researchers all alike to improve multiple aspects of future PIM hardware and software. The workloads have different characteristics, exhibiting heterogeneity 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 that can be used to assess various architecture limits such as compute throughput and memory bandwidth.

# UPMEM DIMMs

---

- E19: 8 chips/DIMM (1 rank). DPUs @ 267 MHz
- P21: 16 chips/DIMM (2 ranks). DPUs @ 350 MHz



# PIM's Promises

## UPMEM PIM massive benefits

- Massive speed-up
  - Massive additional compute & bandwidth
- Massive energy gains
  - Most data movement on chip
- Low cost
  - ~300\$ of additional DRAM silicon
  - Affordable programming
- Massive ROI / TCO gains

Energy efficiency when computing on or off memory chip		Server + PIM DRAM	Server + normal DRAM
DRAM to processor 64-bit operand	pJ	~150	~3000*
Operation	pJ	~20	~10*
Server consumption	W	<b>~700W</b>	~300W
speed-up		<b>~ x20</b>	x1
energy gain		<b>~ x10</b>	x1
TCO gain		<b>~ x10</b>	x1

\*Exascale Computing Trends: Adjusting to the "New Normal" for Computer Architecture; John Shalf, Computing in Science & engineering, 2013

# 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

### 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

But the PIM Graal is worth it !

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



# UPMEM Patent

(12) **United States Patent**  
Devaux et al.

(10) **Patent No.: US 10,324,870 B2**  
(45) **Date of Patent: Jun. 18, 2019**

(54) **MEMORY CIRCUIT WITH INTEGRATED PROCESSOR**

(71) Applicant: **UPMEM**, Grenoble (FR)

(72) Inventors: **Fabrice Devaux**, La Conversion (CH);  
**Jean-François Roy**, Grenoble (FR)

(73) Assignee: **UPMEM**, Grenoble (FR)

(\*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 0 days.

(21) Appl. No.: **15/551,418**

(22) PCT Filed: **Feb. 12, 2016**

(56)

## References Cited

### U.S. PATENT DOCUMENTS

5,666,485	A *	9/1997	Suresh .....	G06F 13/1605
				710/113
6,463,001	B1	10/2002	Williams .....	
7,349,277	B2 *	3/2008	Kinsley .....	G11C 11/406
				365/193
8,438,358	B1 *	5/2013	Kraipak .....	G11C 7/04
				711/167

(Continued)

### FOREIGN PATENT DOCUMENTS

EP	0780768	A1	6/1997
JP	H03109661	A	5/1991
WO	2010/141221	A1	12/2010

(57)

## ABSTRACT

A memory circuit having: a memory array including one or more memory banks; a first processor; and a processor control interface for receiving data processing commands directed to the first processor from a central processor, the processor control interface being adapted to indicate to the central processor when the first processor has finished accessing one or more of the memory banks of the memory array, these memory banks becoming accessible to the central processor.

# 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

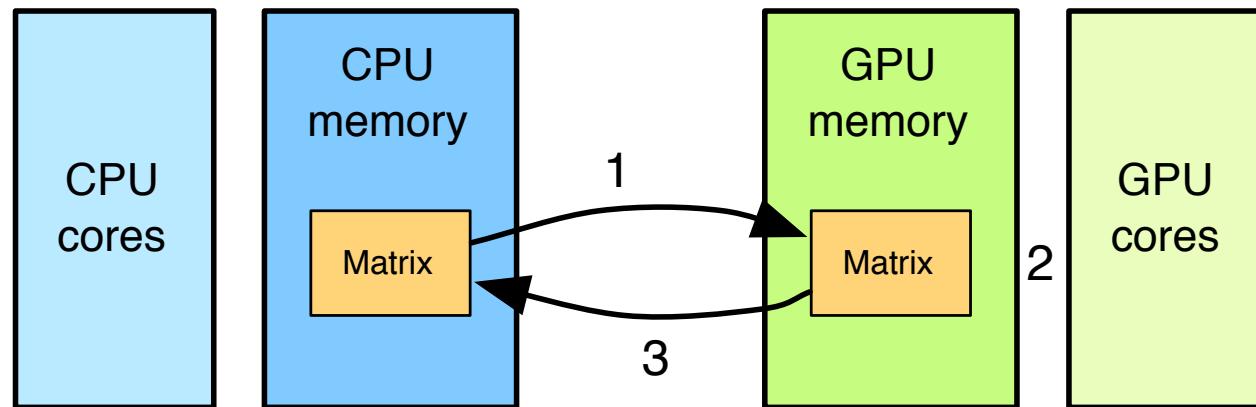
# Accelerator Model (I)

---

- 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

# GPU Computing

- Computation is offloaded to the GPU
- Three steps
  - CPU-GPU data transfer (1)
  - GPU kernel execution (2)
  - GPU-CPU data transfer (3)



# Accelerator Model (II)

- FIG. 6 is a flow diagram representing operations in a method of delegating a processing task to a DRAM processor according to an example embodiment

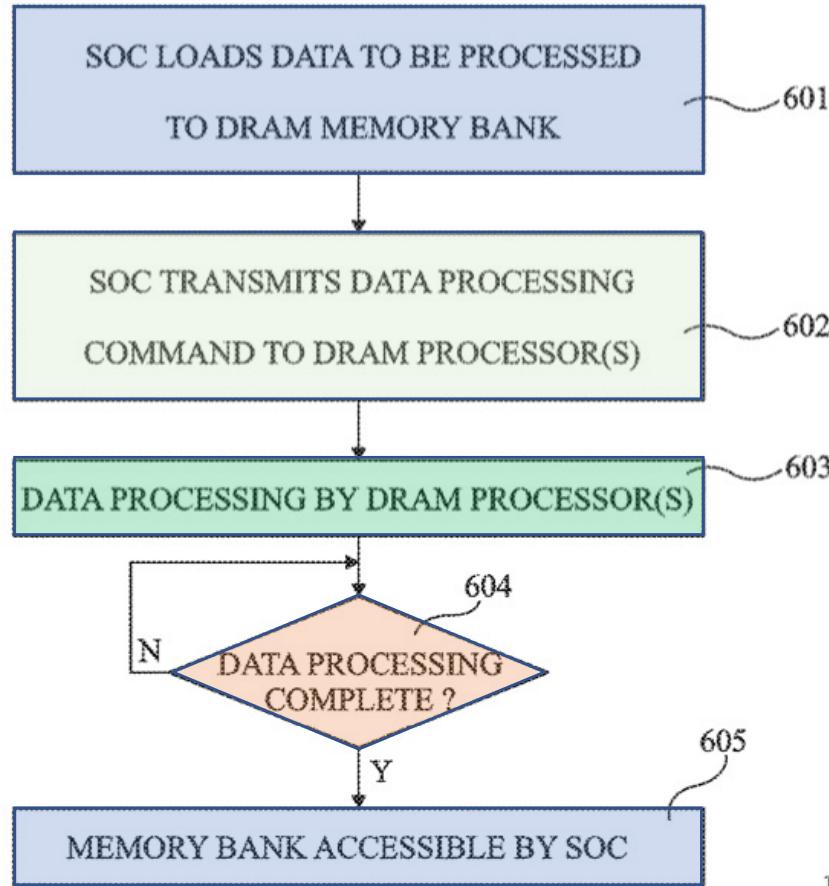


Fig 6

# System Organization (I)

- FIG. 1 schematically illustrates a computing system comprising DRAM circuits having integrated processors according to an example embodiment

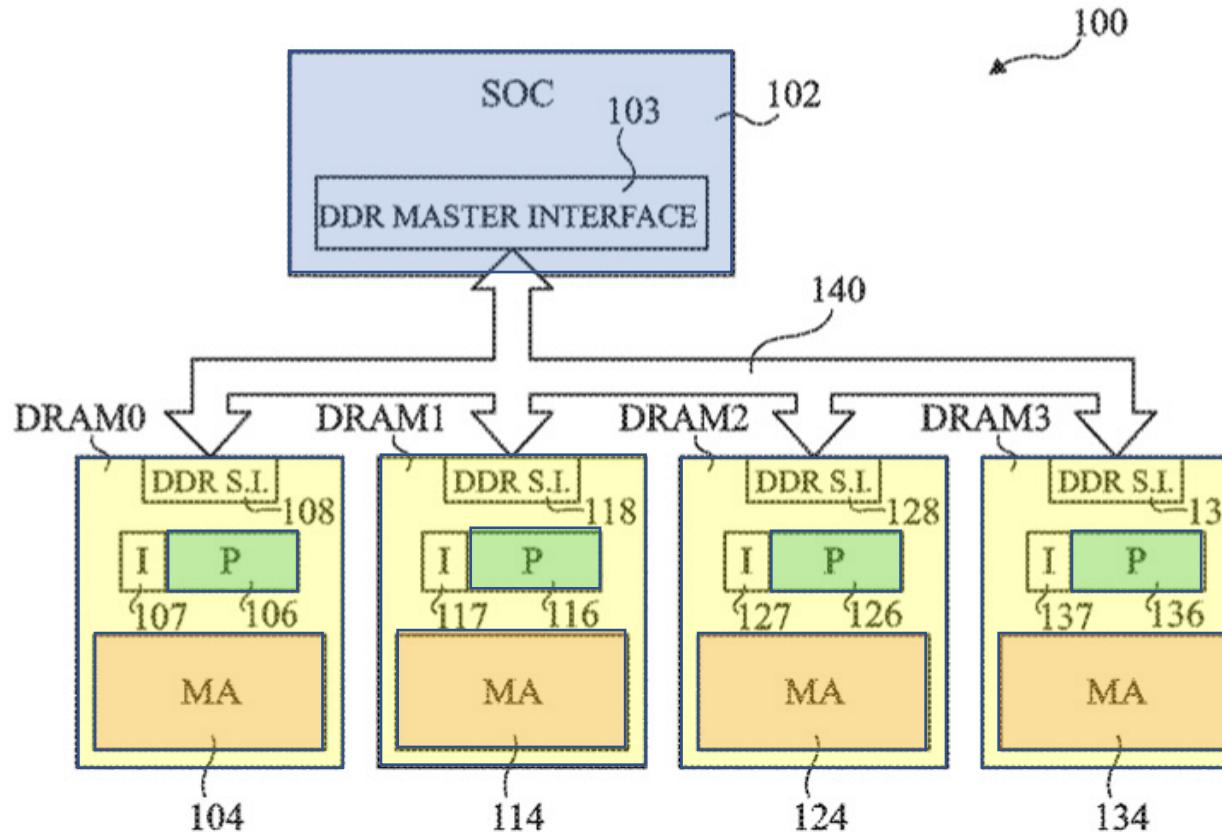
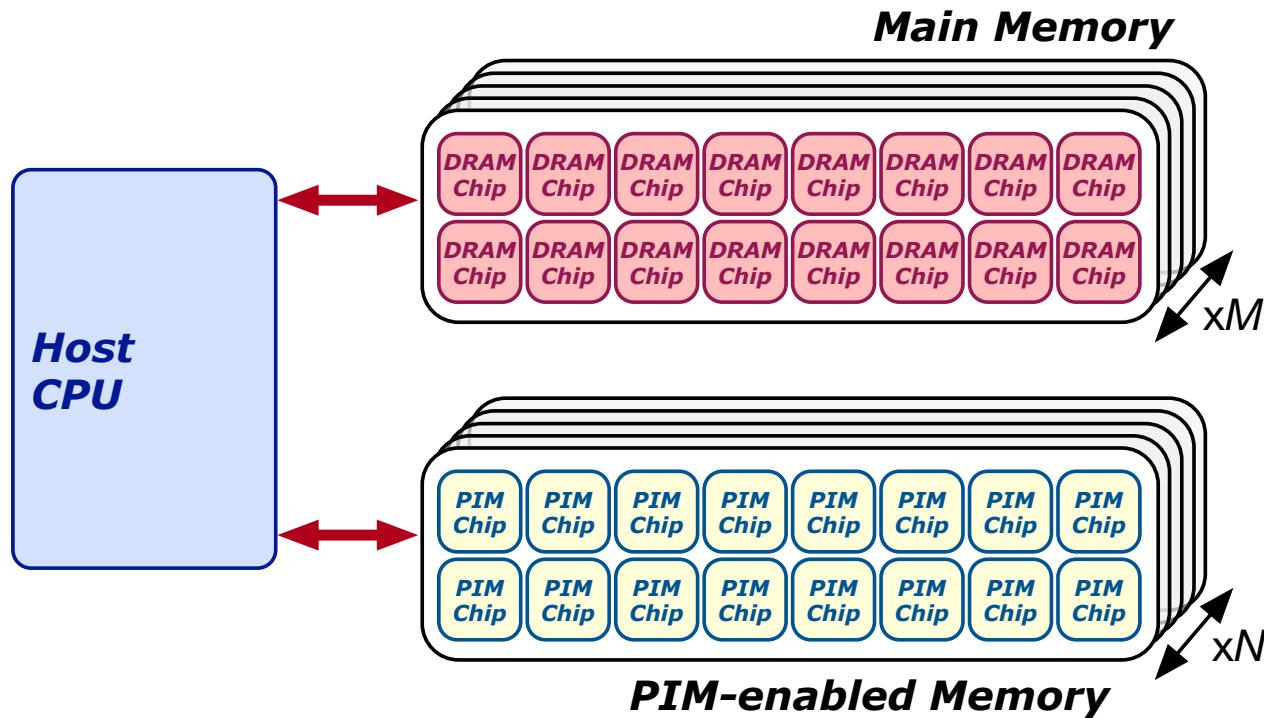


Fig 1

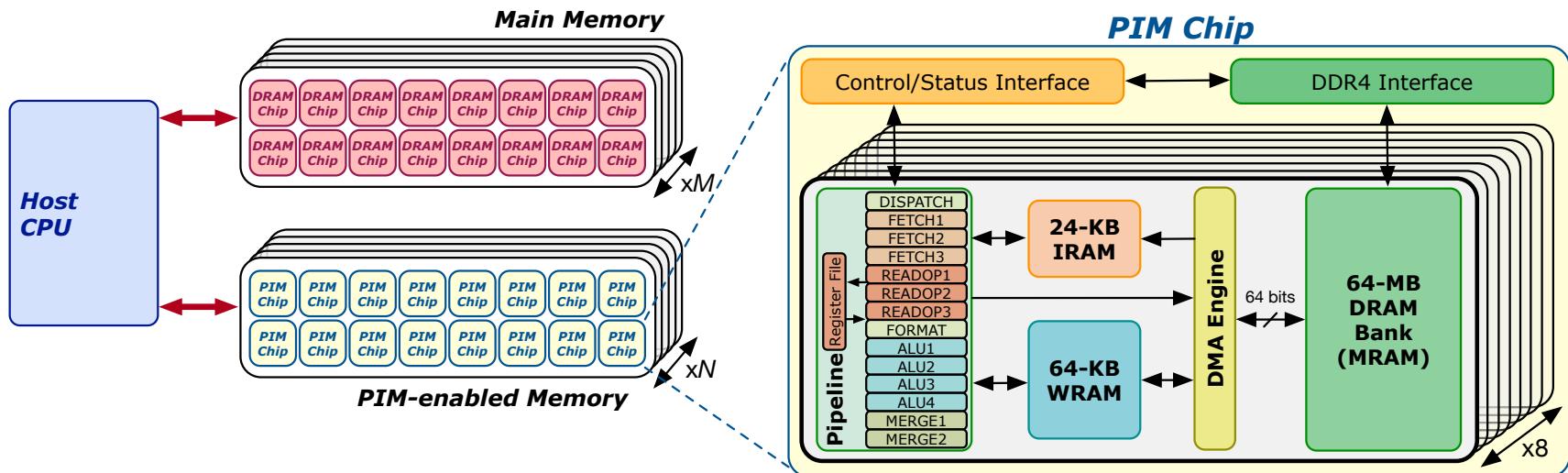
# System Organization (II)

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



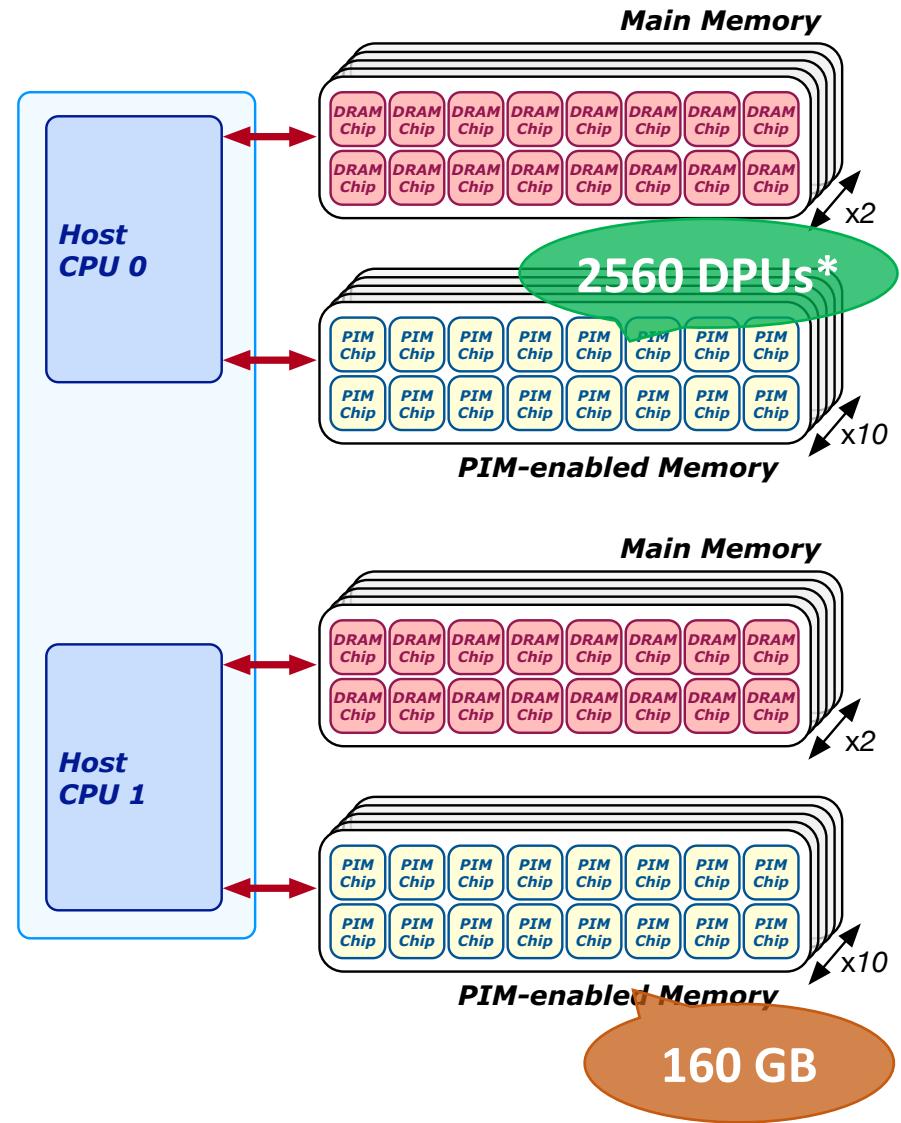
# System Organization (III)

- 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



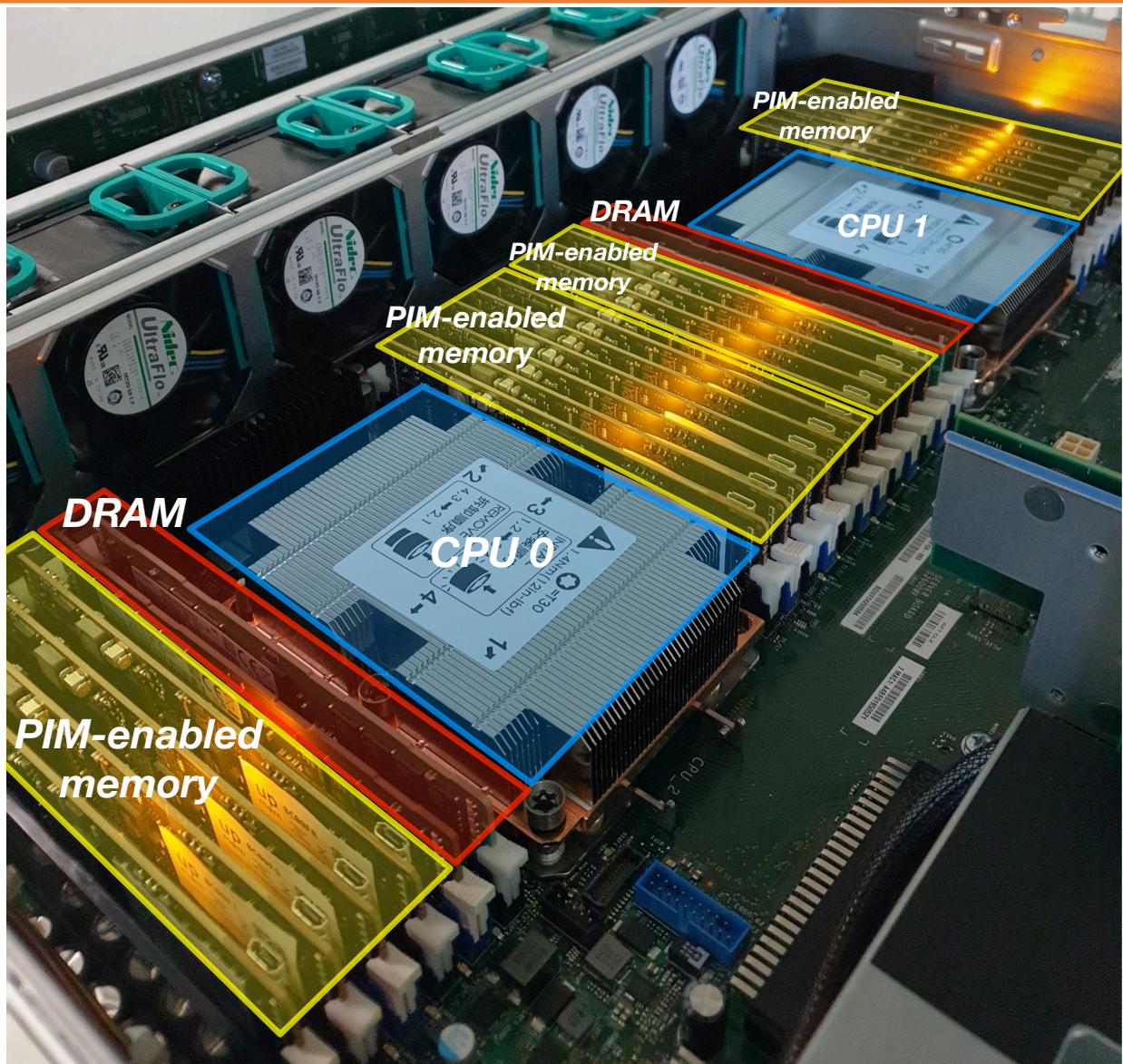
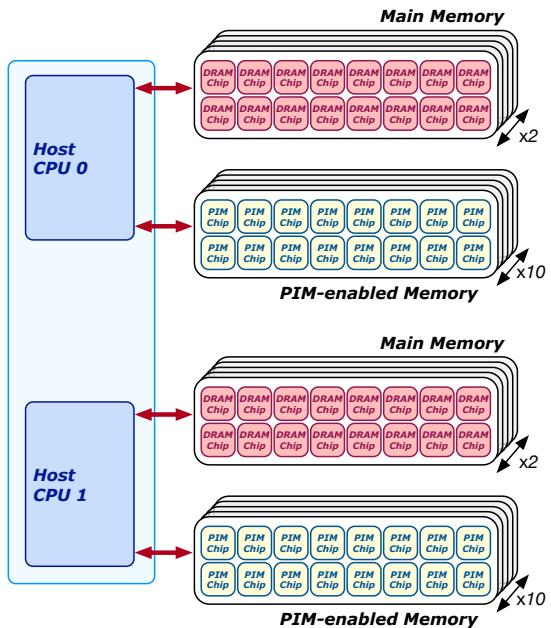
# 2,560-DPU System (I)

- 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



\* There are 4 faulty DPUs in the system that we use in our experiments. Thus, the maximum number of DPUs we can use is 2,556.

# 2,560-DPU System (II)

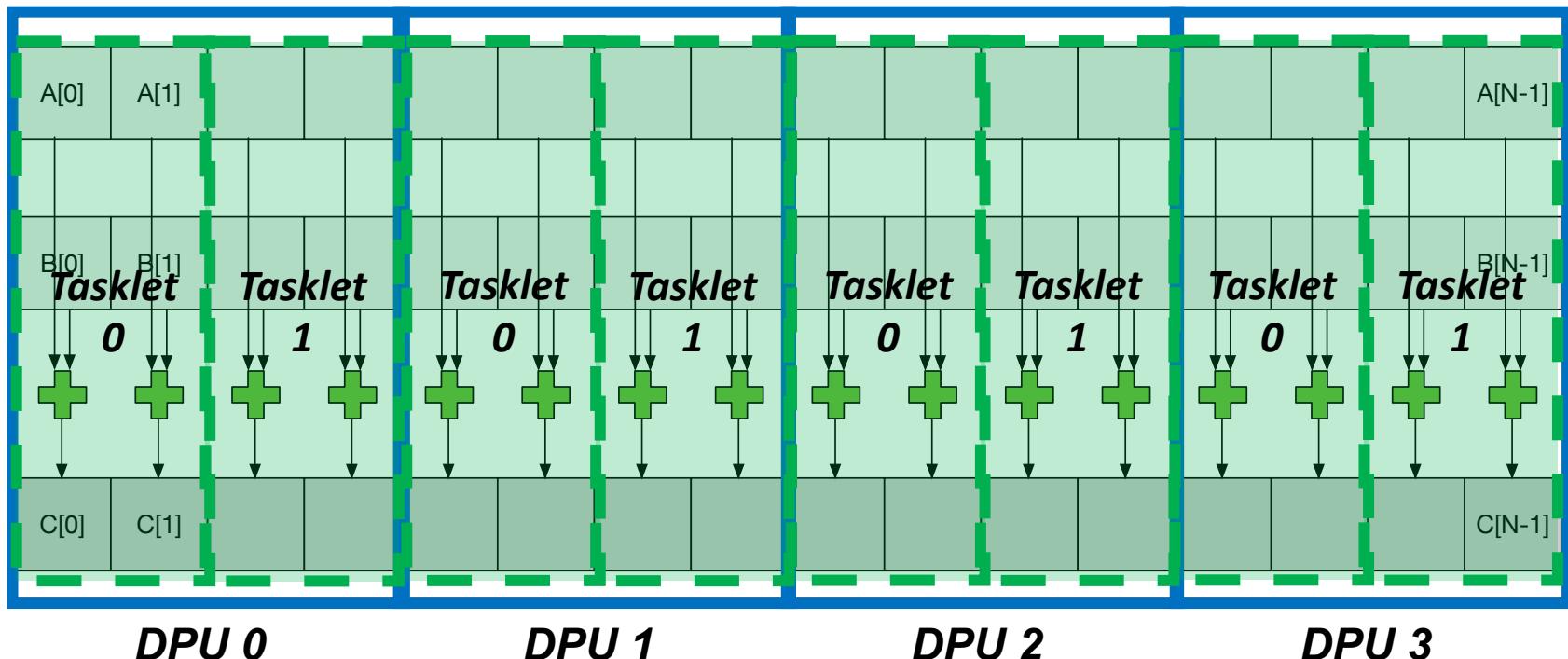


# Outline

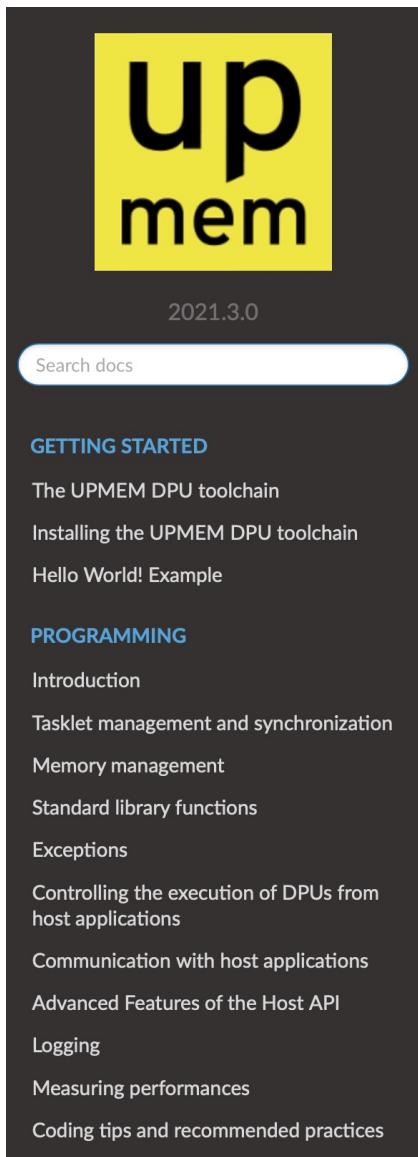
- Introduction
  - Accelerator Model
  - UPMEM-based PIM System Overview
- UPMEM PIM Programming
  - Vector Addition
  - CPU-DPU Data Transfers
  - Inter-DPU Communication
  - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
  - Arithmetic Throughput
  - WRAM and MRAM Bandwidth
- PrIM Benchmarks
  - 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



# UPMEM SDK Documentation



[Home](#) » User Manual

## User Manual

### Getting started

- [The UPMEM DPU toolchain](#)
  - [Notes before starting](#)
  - [The toolchain purpose](#)
  - [dpu-upmem-dpure-clang](#)
    - [Limitations](#)
  - [The DPU Runtime Library](#)
  - [The Host Library](#)
  - [dpu-lldb](#)
- [Installing the UPMEM DPU toolchain](#)
  - [Dependencies](#)
    - [Python](#)
  - [Installation packages](#)
    - [Installation from tar.gz binary archive](#)
  - [Functional simulator](#)
- [Hello World! Example](#)
  - [Purpose](#)
  - [Writing and building the program](#)
  - [Running and testing hello world](#)
  - [Creating a host application to drive the program](#)

# 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

# DPU Allocation

- `dpu_alloc()` allocates a number of DPUs
  - Creates a `dpu_set`

```
1  struct dpu_set_t dpu_set, dpu;
2  uint32_t nr_of_dpus;
3
4  // Allocate DPUs
5  DPU_ASSERT(dpu_alloc(NR_DPUS, NULL, &dpu_set));
6
7  DPU_ASSERT(dpu_get_nr_dpus(dpu_set, &nr_of_dpus));
8  printf("Allocated %d DPU(s)\n", nr_of_dpus);
9
```

Can we allocate different DPU sets over the course of a program?

Yes, we can. We show an example next

We deallocate a DPU set with `dpu_free()`

# DPU Allocation: Needleman-Wunsch (NW)

- In NW we change the number of DPUs in the DPU set as computation progresses

```
1 // Top-left computation on DPUs
2 ▼ for (unsigned int blk = 1; blk <= (max_cols-1)/BL; blk++) {
3
4     // If nr_of_blocks are lower than max_dpus,
5     // set nr_of_dpus to be equal with nr_of_blocks
6     unsigned nr_of_blocks = blk;
7 ▼     if (nr_of_blocks < max_dpus) {
8         DPU_ASSERT(dpu_free(dpu_set));
9         DPU_ASSERT(dpu_alloc(nr_of_blocks, NULL, &dpu_set));
10        DPU_ASSERT(dpu_load(dpu_set, DPU_BINARY, NULL));
11        DPU_ASSERT(dpu_get_nr_dpus(dpu_set, &nr_of_dpus));
12    } else if (nr_of_dpus == max_dpus) {
13        ;
14    } else {
15        DPU_ASSERT(dpu_free(dpu_set));
16        DPU_ASSERT(dpu_alloc(max_dpus, NULL, &dpu_set));
17        DPU_ASSERT(dpu_load(dpu_set, DPU_BINARY, NULL));
18        DPU_ASSERT(dpu_get_nr_dpus(dpu_set, &nr_of_dpus));
19 ▲    }
20
21
22 ▲ }
```

# Load DPU Binary

- `dpu_load()` loads a program in all DPUs of a `dpu_set`

```
1 // Define the DPU Binary path as DPU_BINARY here
2 #ifndef DPU_BINARY
3 #define DPU_BINARY "./bin/dpu_code"
4 #endif
5
6 ...
7
8 // Load binary
9 DPU_ASSERT(dpu_load(dpu_set, DPU_BINARY, NULL));
10
```

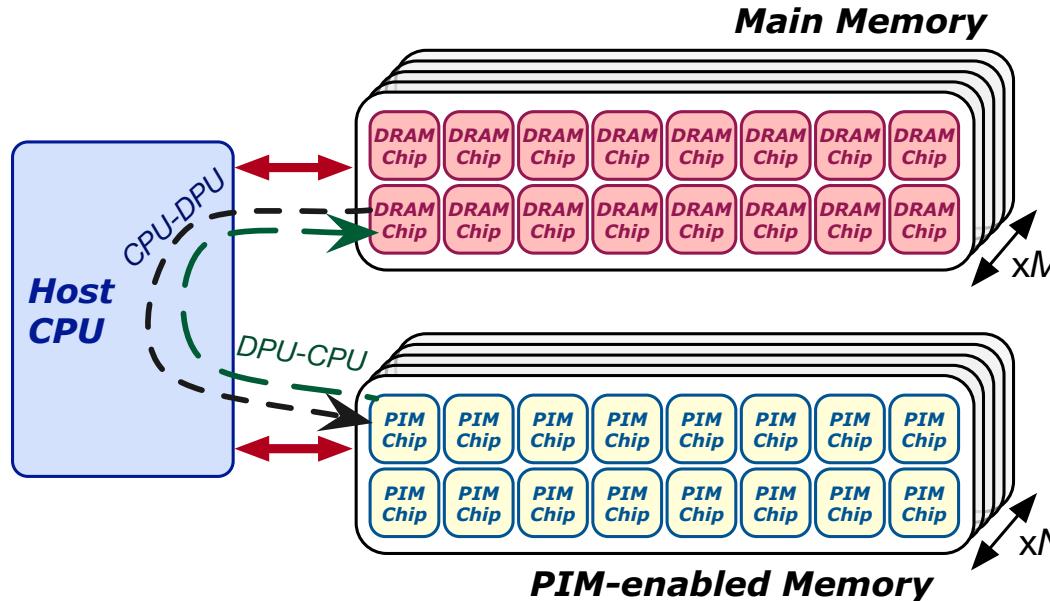
Is it possible to launch different kernels onto different DPUs?

Yes, it is possible. This enables:

- Workloads with **task-level parallelism**
- Different programs using different DPU sets

# 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

# Serial Transfers

- `dpu_copy_to()`;
- `dpu_copy_from()`;
- We transfer (part of) a buffer to/from each DPU in the `dpu_set`
- `DPU_MRAM_HEAP_POINTER_NAME`: Start of the MRAM range that can be freely accessed by applications
  - We do not allocate MRAM explicitly

```
1 ▼ DPU_FOREACH (dpu_set, dpu) {  
2     DPU_ASSERT(dpu_copy_to(dpu, DPU_MRAM_HEAP_POINTER_NAME  
3     DPU_ASSERT(dpu_copy_to(dpu, DPU_MRAM_HEAP_POINTER_NAME  
4     i++;  
5 ▲ }  
6  
Offset within MRAM    Pointer to main memory    Transfer size
```

# Parallel Transfers

- We push different buffers to/from a DPU set in one transfer
  - All buffers need to be of the same size
- First, prepare (`dpu_prepare_xfer`);  
then, push (`dpu_push_xfer`)
- Direction:
  - `DPU_XFER_TO_DPU`
  - `DPU_XFER_FROM_DPU`

```
1 ▼ DPU_FOREACH(dpu_set, dpu, i) {  
2     DPU_ASSERT(dpu_prepare_xfer(dpu, [bufferA + input_size_dpu_8bytes * i]))  
3 }  
4 DPU_ASSERT(dpu_push_xfer(dpu_set, DPU_XFER_TO_DPU, DPU_MRAM_HEAP_POINTER_NAME, [0], [input_size_dpu_8bytes * sizeof(T)], DPU_XFER_DEFAULT));  
5  
6 ▼ DPU_FOREACH(dpu_set, dpu, i) {  
7     DPU_ASSERT(dpu_prepare_xfer(dpu, [bufferB + input_size_dpu_8bytes * i]))  
8 }  
9 DPU_ASSERT(dpu_push_xfer(dpu_set, DPU_XFER_TO_DPU, DPU_MRAM_HEAP_POINTER_NAME, [input_size_dpu_8bytes * sizeof(T)], [input_size_dpu_8bytes * sizeof(T)], DPU_XFER_DEFAULT));  
10
```

Pointer to main memory      Offset within MRAM      Transfer size  
Direction

# Broadcast Transfers

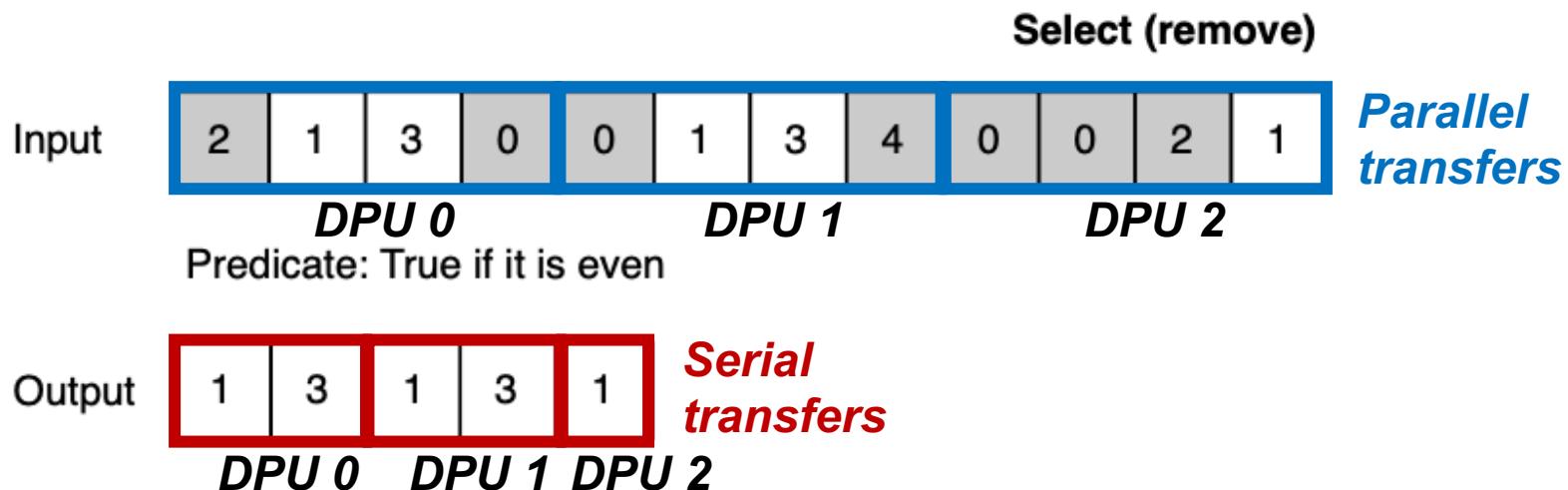
---

- `dpu_broadcast_to()`;
  - Only CPU to DPU
- We transfer the same buffer to all DPUs in the `dpu_set`

```
1 DPU_ASSERT(dpu_broadcast_to(dpu_set, DPU_MRAM_HEAP_POINTER_NAME, 0, bufferA, input_size_dpu * sizeof(T), DPU_XFER_DEFAULT));  
2                                         Pointer to main memory           Transfer size
```

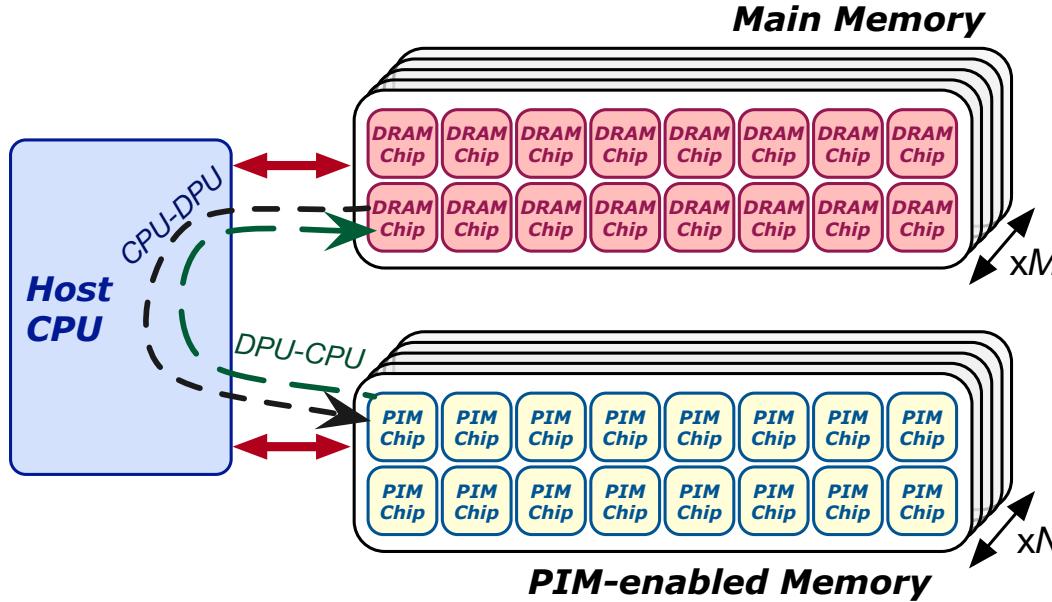
# Different Types of Transfers in a Program

- 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?

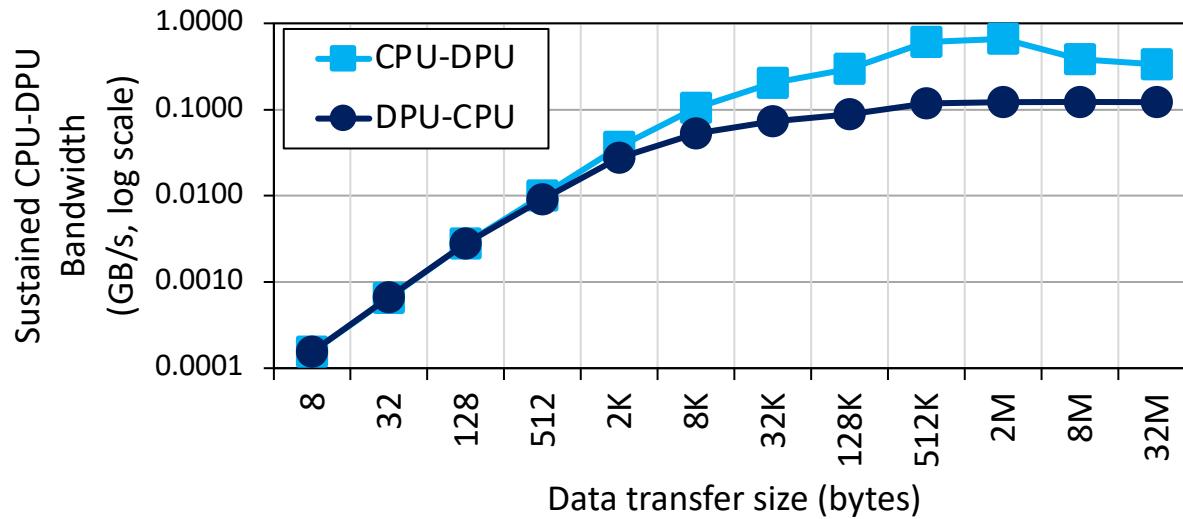
- With a microbenchmark, we obtain the **sustained bandwidth of all types of CPU-DPU and DPU-CPU transfers**
- Two experiments:
  - 1 DPU: variable CPU-DPU and DPU-CPU transfer size (**8 bytes to 32 MB**)
  - 1 rank: 32 MB CPU-DPU and DPU-CPU transfers to/from a set of **1 to 64 MRAM banks** within the same rank
- Preliminary experiments with more than one rank
  - Channel-level parallelism

**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

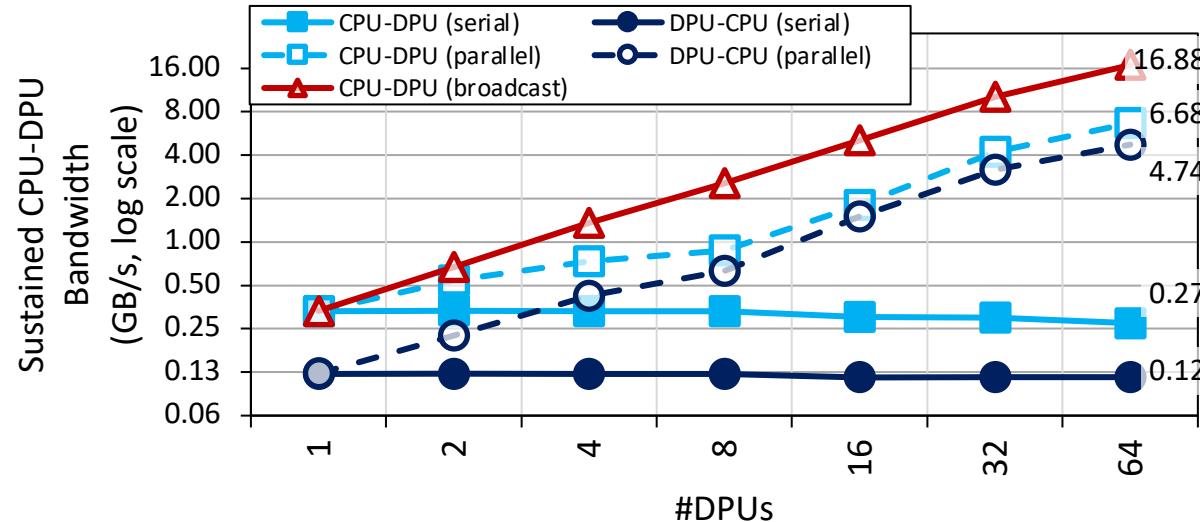


## KEY OBSERVATION 7

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 (I)

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

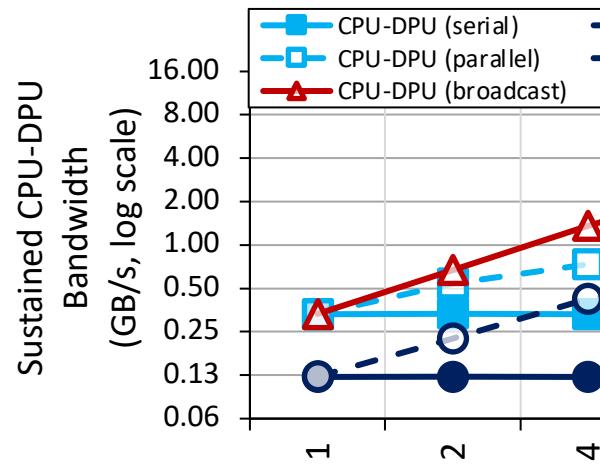


## KEY OBSERVATION 8

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 (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.**

# “Transposing” Library

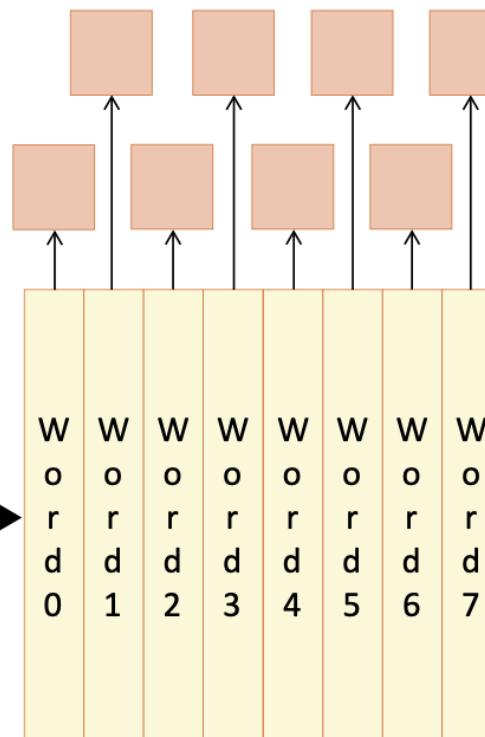
## The library feeds DPUs with correct data

Eight 64-bit “horizontal” words  
are turned into 8 vertical words,  
feeding 8 different DRAM chips

This way DPUs see full 64-bit  
words, not chunk of them

Word 0
Word 1
Word 2
Word 3
Word 4
Word 5
Word 6
Word 7

Library →



DRAM chip  
have 8-bit  
data bus

The transformation, a 8x8  
matrix transposition, is  
done by the library inside  
a 64-byte cache line, thus  
very efficiently.

# Microbenchmark: CPU-DPU

- CPU-DPU (serial/parallel/broadcast) and DPU-CPU (serial/parallel)

CMU-SAFARI / **prim-benchmarks**

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main / **prim-benchmarks / Microbenchmarks / CPU-DPU /**

Juan Gomez Luna PrIM -- first commit 3de4b49 7 days ago History

..

File	Commit Message	Time
dpu	PrIM -- first commit	7 days ago
host	PrIM -- first commit	7 days ago
support	PrIM -- first commit	7 days ago
Makefile	PrIM -- first commit	7 days ago
run.sh	PrIM -- first commit	7 days ago

# DPU Kernel Launch

- `dpu_launch()` launches a kernel on a `dpu_set`
  - `DPU_SYNCHRONOUS` suspends the application until the kernel finishes
  - `DPU_ASYNCHRONOUS` returns the control to the application
    - `dpu_sync` or `dpu_status` to check kernel completion

```
1  printf("Run program on DPU(s) \n");
2  // Run DPU kernel
3  DPU_ASSERT(dpu_launch(dpu_set, DPU_SYNCHRONOUS));
4
```

What does the asynchronous execution enable?

Some ideas:

- **Task-level parallelism:** concurrent execution of different kernels on different DPU sets
- Concurrent **heterogeneous computation** on CPU and DPUs

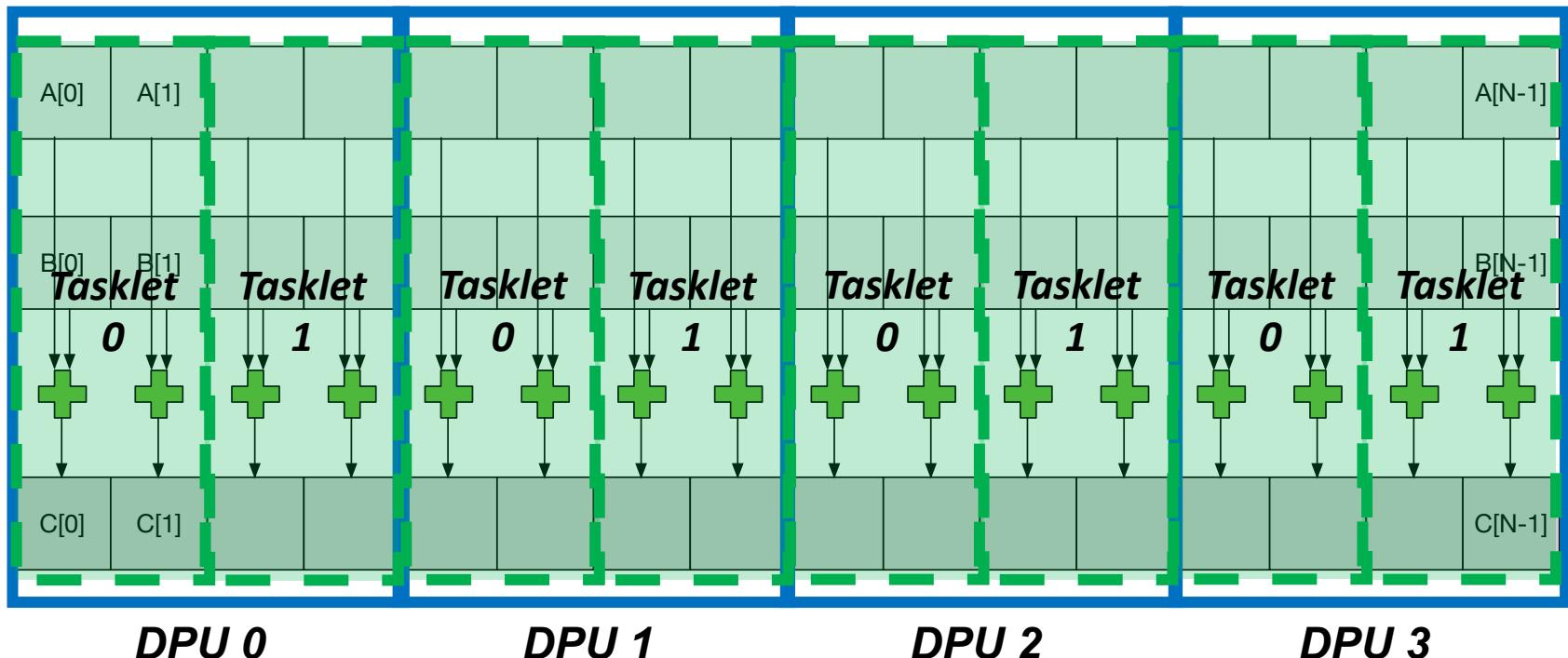
# How to Pass Parameters to the Kernel?

- We can use serial and parallel transfers
  - We pass them directly to the scratchpad memory of the DPU
    - Working RAM (WRAM): We introduce it in the next slides
  - This is useful for input parameters and some results

```
1 // In DPU WRAM (dpu/task.c)
2 __host dpu_arguments_t DPU_INPUT_ARGUMENTS;
3 __host dpu_results_t DPU_RESULTS[NR_TASKLETS];
4
5
6
7 // Host code (host/app.c)
8 #ifdef SERIAL
9     DPU_FOREACH (dpu_set, dpu) {
10         DPU_ASSERT(dpu_copy_to(dpu, "DPU_INPUT_ARGUMENTS", 0, (const void *)&input_arguments[i], sizeof(input_arguments[0])));
11         i++;
12     }
13 #else
14     DPU_FOREACH(dpu_set, dpu, i) {
15         DPU_ASSERT(dpu_prepare_xfer(dpu, &input_arguments[i]));
16     }
17     DPU_ASSERT(dpu_push_xfer(dpu_set, DPU_XFER_TO_DPU, "DPU_INPUT_ARGUMENTS", 0, sizeof(input_arguments[0]), DPU_XFER_DEFAULT));
18 #endif
```

# Recall: Vector Addition (VA)

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



# Programming a DPU Kernel (I)

- Vector addition

```
1 // Vector addition kernel          Tasklet ID
2 int main_kernel1() {               Tasklet ID
3     unsigned int tasklet_id = me();      Size of vector tile processed by a DPU
4     uint32_t input_size_dpu_bytes = DPU_INPUT_ARGUMENTS.size; // Input size per DPU in bytes
5     uint32_t input_size_dpu_bytes_transfer = DPU_INPUT_ARGUMENTS.transfer_size; // Transfer input size per DPU in bytes
6
7     // Address of the current processing block in MRAM
8     uint32_t base_tasklet = tasklet_id << BLOCK_SIZE_LOG2;           MRAM addresses of arrays A and B
9     uint32_t mram_base_addr_A = (uint32_t)DPU_MRAM_HEAP_POINTER;
10    uint32_t mram_base_addr_B = (uint32_t)(DPU_MRAM_HEAP_POINTER + input_size_dpu_bytes_transfer);
11
12    // Initialize a local cache to store the MRAM block
13    T *cache_A = (T *) mem_alloc(BLOCK_SIZE);                         WRAM allocation
14    T *cache_B = (T *) mem_alloc(BLOCK_SIZE);
15
16    for(unsigned int byte_index = base_tasklet; byte_index < input_size_dpu_bytes; byte_index += BLOCK_SIZE * NR_TASKLETS){
17        // Bound checking
18        uint32_t l_size_bytes = (byte_index + BLOCK_SIZE >= input_size_dpu_bytes) ? (input_size_dpu_bytes - byte_index) : BLOCK_SIZE;
19
20        // Load cache with current MRAM block
21        mram_read((__mram_ptr void const*)(mram_base_addr_A + byte_index), cache_A, l_size_bytes);           MRAM-WRAM DMA
22        mram_read((__mram_ptr void const*)(mram_base_addr_B + byte_index), cache_B, l_size_bytes);           transfers
23
24        // Compute vector addition
25        vector_addition(cache_B, cache_A, l_size_bytes >> DIV);           Vector addition (see next slide)
26
27        // Write cache to current MRAM block
28        mram_write(cache_B, (__mram_ptr void*)(mram_base_addr_B + byte_index), l_size_bytes);           WRAM-MRAM DMA transfer
29
30    }
31
32    return 0;
}
```

# Programming a DPU Kernel (II)

---

- Vector addition

```
1 // vector_addition: Computes the vector addition of a cached block
2 ▼ static void vector_addition(T *bufferB, T *bufferA, unsigned int l_size) {
3
4 ▼     for (unsigned int i = 0; i < l_size; i++){
5         bufferB[i] += bufferA[i];
6 ▲     }
7
8 ▲ }
```

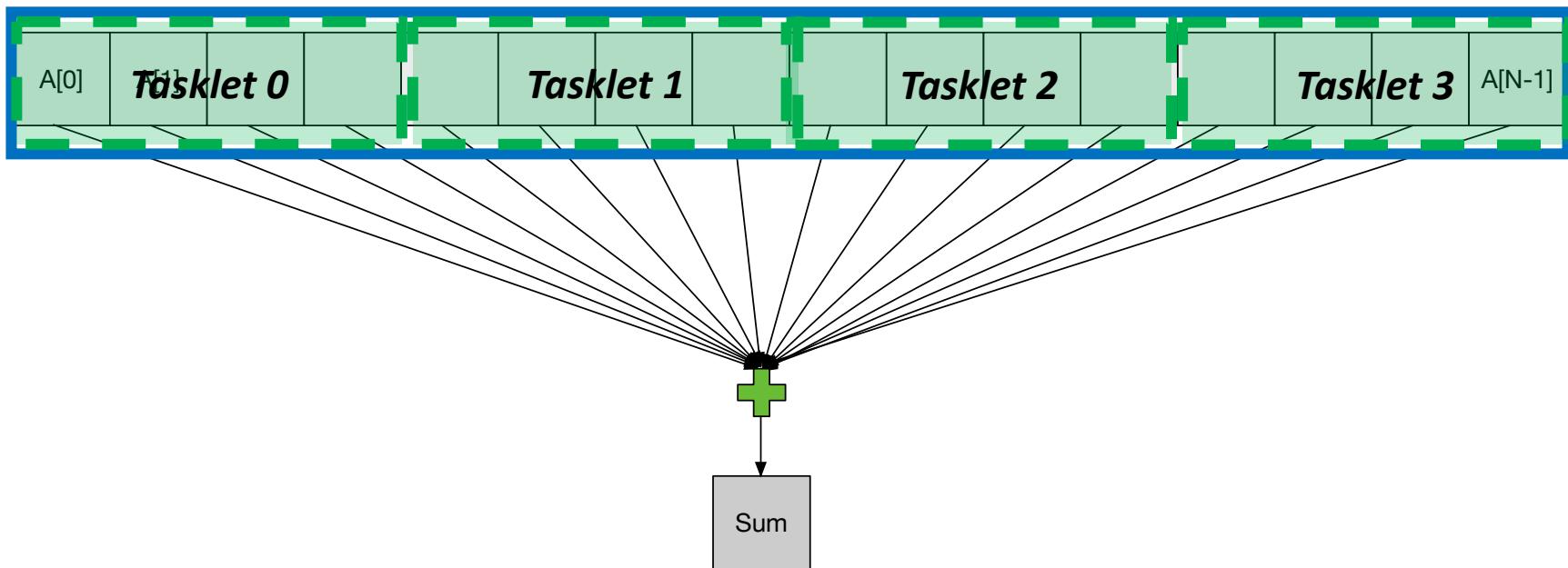
# Programming a DPU Kernel (III)

---

- A **tasklet** is the software abstraction of a hardware thread
- Each tasklet can have its **own memory space in WRAM**
  - Tasklets can also share data in WRAM by sharing pointers
- Tasklets within the same DPU can **synchronize**
  - Mutual exclusion
    - `mutex_lock(); mutex_unlock();`
  - Handshakes
    - `handshake_wait_for(); handshake_notify();`
  - Barriers
    - `barrier_wait();`
  - Semaphores
    - `sem_give(); sem_take();`

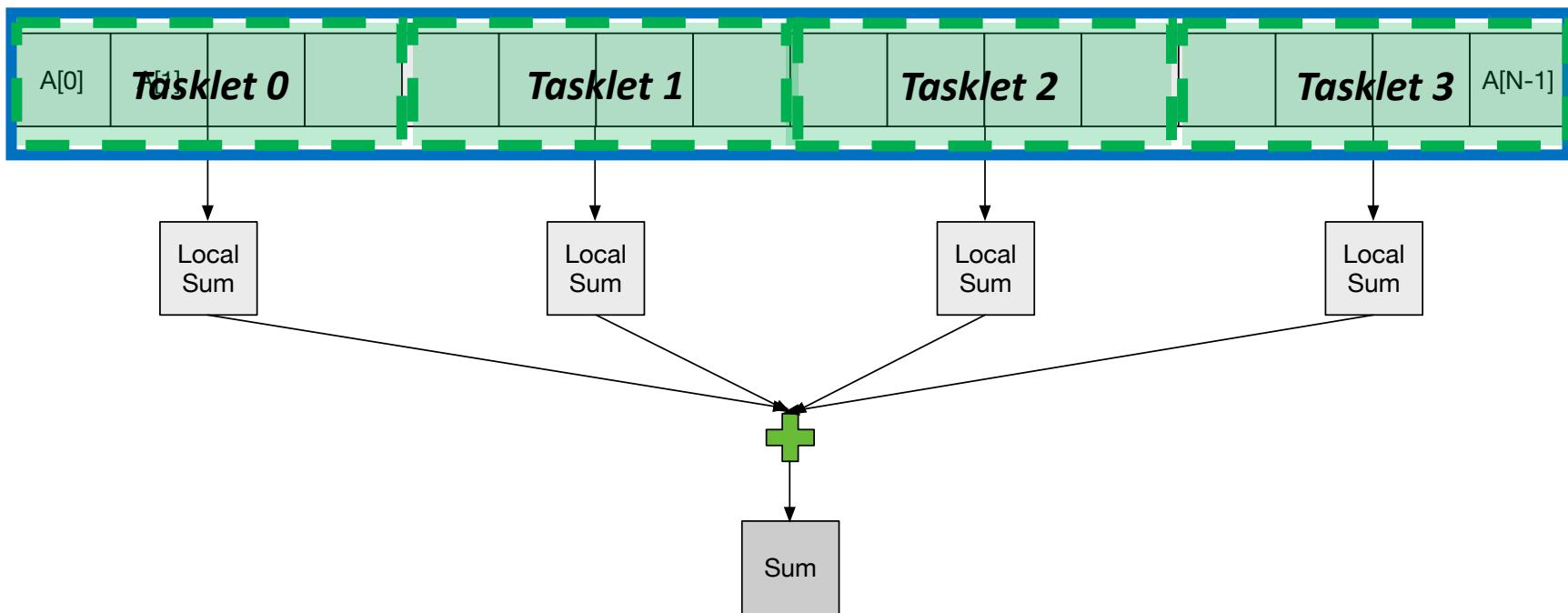
# Parallel Reduction (I)

- Tasklets in a DPU can work together on a parallel reduction



# Parallel Reduction (II)

- Each tasklet computes a local sum



# Parallel Reduction (III)

- Each tasklet computes a local sum

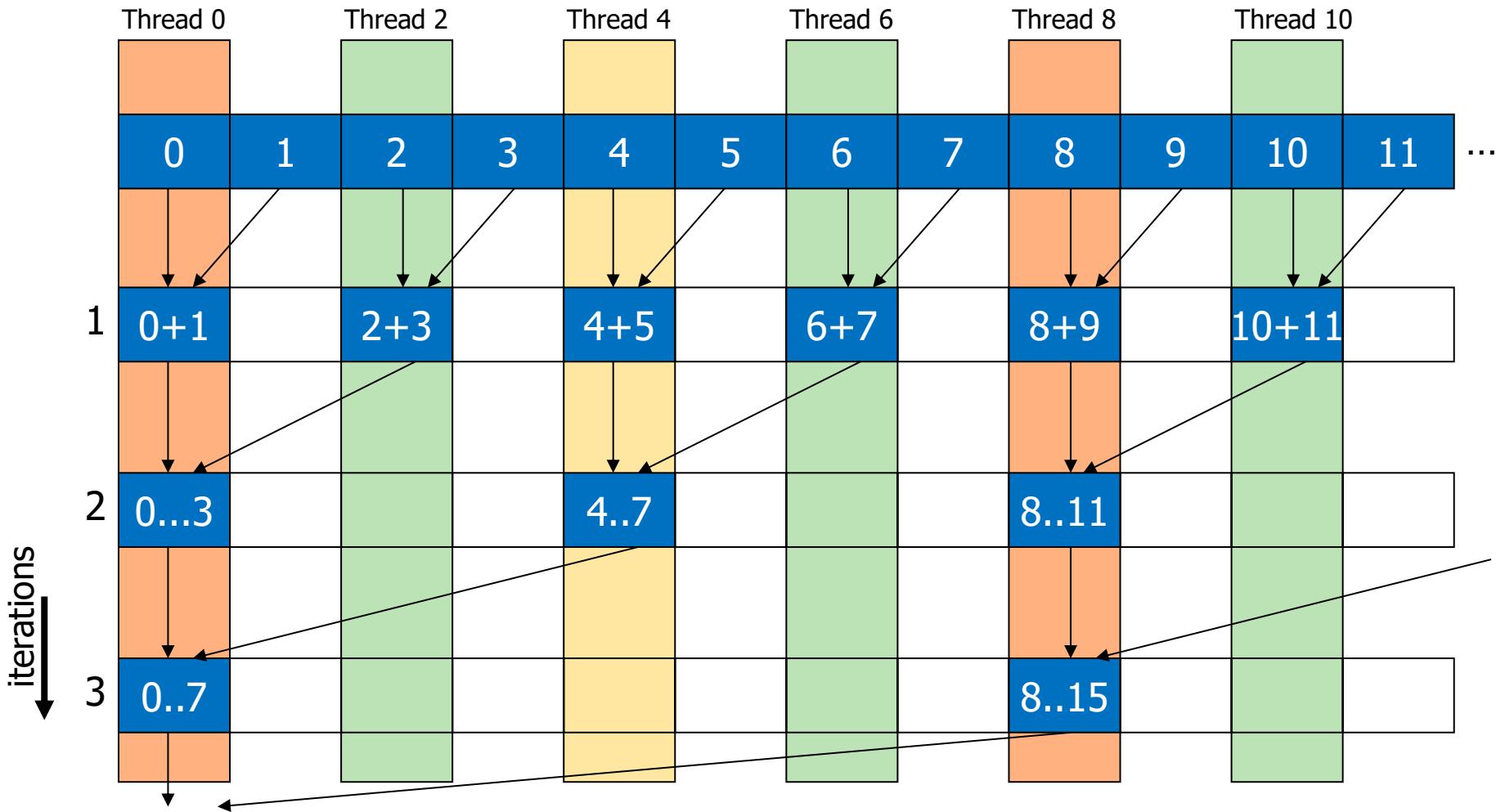
```
1 ▼ for(unsigned int byte_index = base_tasklet; byte_index < input_size_dpu_bytes; byte_index += BLOCK_SIZE * NR_TASKLETS){  
2  
3     // Bound checking  
4     uint32_t l_size_bytes = (byte_index + BLOCK_SIZE >= input_size_dpu_bytes) ? (input_size_dpu_bytes - byte_index) : BLOCK_SIZE;  
5  
6     // Load cache with current MRAM block  
7     mram_read((__mram_ptr void const*)(mram_base_addr_A + byte_index), cache_A, l_size_bytes);  
8  
9     // Reduction in each tasklet  
10    l_count += reduction(cache_A, l_size_bytes >> DIV);    Accumulate in a local sum  
11  
12 ▲ }  
13 // Copy local count to shared array in WRAM  
14 message[tasklet_id] = l_count;  Copy local sum into WRAM
```

# Final Reduction

- A single tasklet can perform the final reduction

```
1 ▼ for(unsigned int byte_index = base_tasklet; byte_index < input_size_dpu_bytes; byte_index += BLOCK_SIZE * NR_TASKLETS){  
2  
3     // Bound checking  
4     uint32_t l_size_bytes = (byte_index + BLOCK_SIZE >= input_size_dpu_bytes) ? (input_size_dpu_bytes - byte_index) : BLOCK_SIZE;  
5  
6     // Load cache with current MRAM block  
7     mram_read((__mram_ptr void const*)(mram_base_addr_A + byte_index), cache_A, l_size_bytes);  
8  
9     // Reduction in each tasklet  
10    l_count += reduction(cache_A, l_size_bytes >> DIV); Accumulate in a local sum  
11  
12 ▲ }  
13 // Copy local count to shared array in WRAM  
14 message[tasklet_id] = l_count; Copy local sum into WRAM  
  
1 // Single-thread reduction  
2 // Barrier  
3 barrier_wait(&my_barrier); Barrier synchronization  
4  
5 ▼ if(tasklet_id == 0){  
6     #pragma unroll  
7     ▼ for (unsigned int each_tasklet = 1; each_tasklet < NR_TASKLETS; each_tasklet++){  
8         message[0] += message[each_tasklet]; Sequential accumulation  
9     }  
10  
11     // Total count in this DPU  
12     result->t_count = message[0];  
13 ▲ }
```

# Vector Reduction: Naïve Mapping



# Using Barriers: Tree-Based Reduction

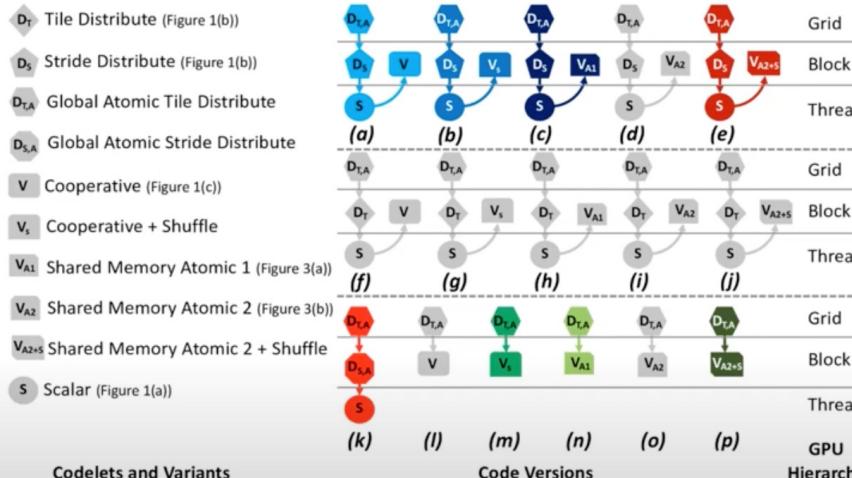
- Multiple tasklets can perform a tree-based reduction
  - After every iteration tasklets synchronize with a barrier
  - Half of the tasklets retire at the end of an iteration

```
1 // Barrier
2 barrier_wait(&my_barrier);
3
4 #pragma unroll
5 ▼ for (unsigned int offset = 1; offset < NR_TASKLETS; offset <= 1){
6
7 ▼   if((tasklet_id & (2*offset - 1)) == 0){
8     message[tasklet_id] += message[tasklet_id + offset]; } "offset" tasklets working
9 ▲ }
10
11 // Barrier
12 barrier_wait(&my_barrier); } Barrier synchronization
13 ▲ }
```

A handshake-based tree-based reduction is also possible.  
We can compare single-tasklet, barrier-based,  
and handshake-based versions\*

# Parallel Reduction on GPU

## Search Space of Parallel Reduction



Over 85 different versions possible!



HetSys Course: Lecture 6: Parallel Patterns: Reduction (Spring 2022)

201 views • Premiered Apr 19, 2022

15 DISLIKE SHARE CLIP SAVE ...



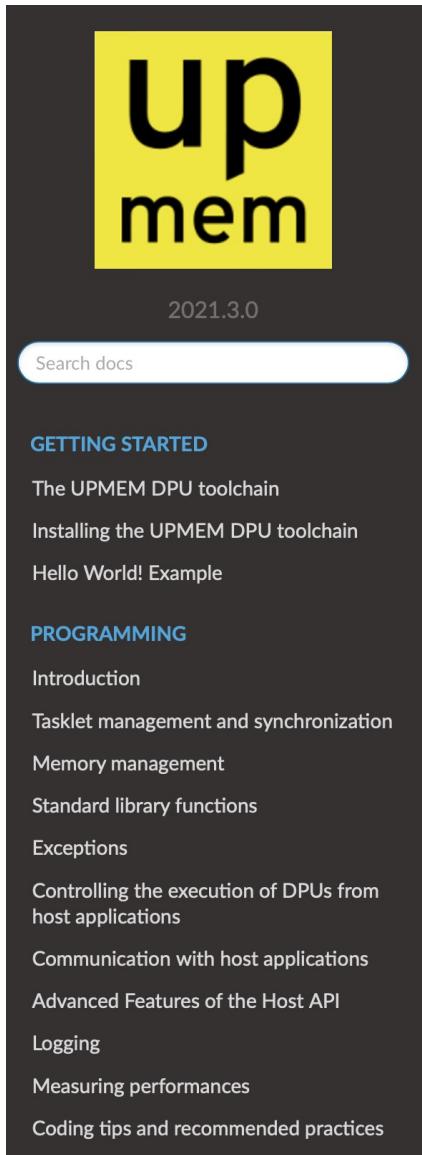
Onur Mutlu Lectures  
24.1K subscribers

SUBSCRIBED



Project & Seminar, ETH Zürich, Spring 2022  
Hands-on Acceleration on Heterogeneous Computing Systems (  
[https://safari.ethz.ch/projects\\_and\\_s...](https://safari.ethz.ch/projects_and_s...))

# UPMEM SDK Documentation



[Home](#) » User Manual

## User Manual

### Getting started

- [The UPMEM DPU toolchain](#)
  - [Notes before starting](#)
  - [The toolchain purpose](#)
  - [dpu-upmem-dpure-clang](#)
    - [Limitations](#)
  - [The DPU Runtime Library](#)
  - [The Host Library](#)
  - [dpu-lldb](#)
- [Installing the UPMEM DPU toolchain](#)
  - [Dependencies](#)
    - [Python](#)
  - [Installation packages](#)
    - [Installation from tar.gz binary archive](#)
  - [Functional simulator](#)
- [Hello World! Example](#)
  - [Purpose](#)
  - [Writing and building the program](#)
  - [Running and testing hello world](#)
  - [Creating a host application to drive the program](#)

# 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 (I)

- FIG. 4 schematically illustrates part of the computing system of FIG. 1 in more detail according to an example embodiment

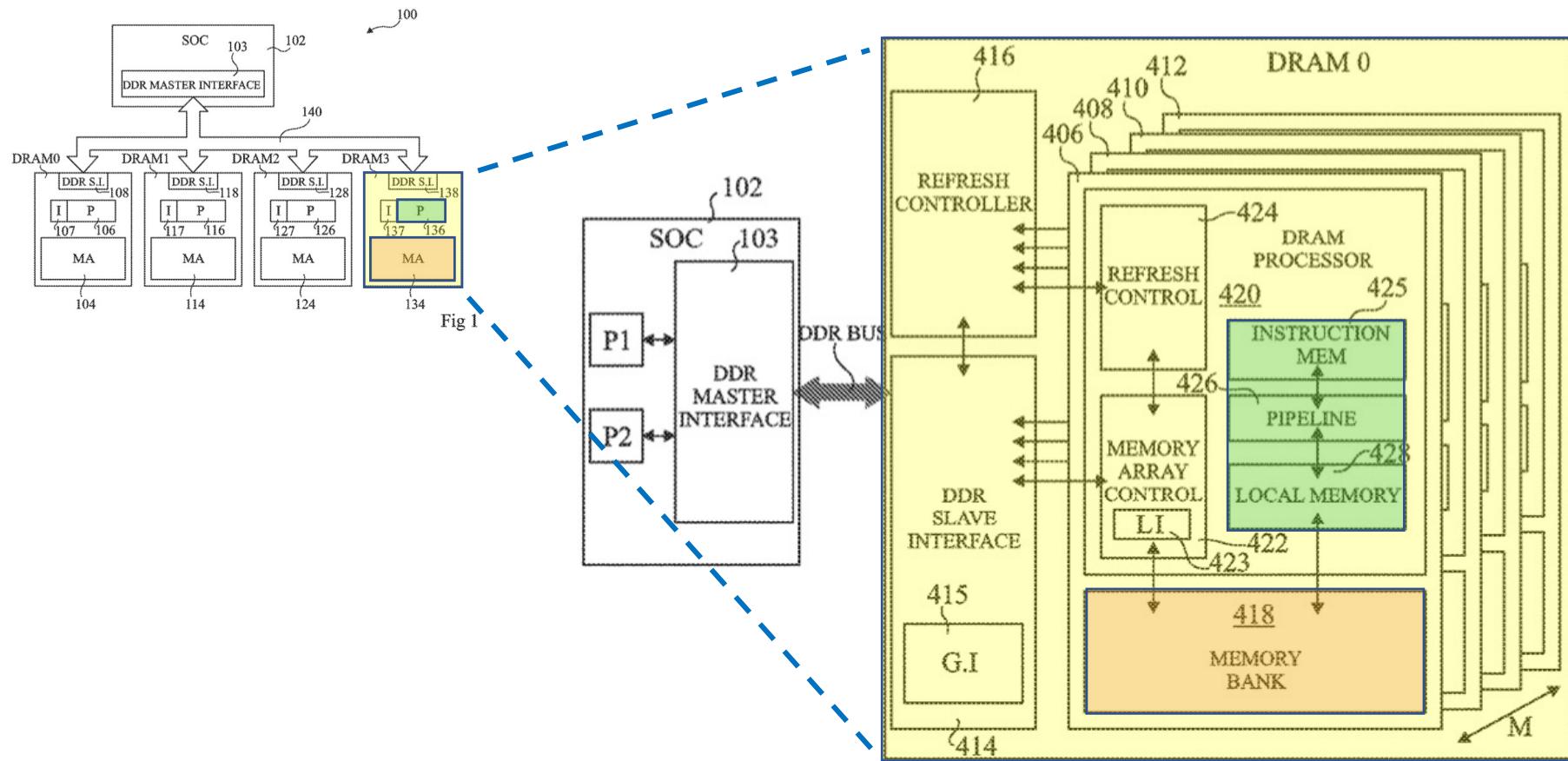
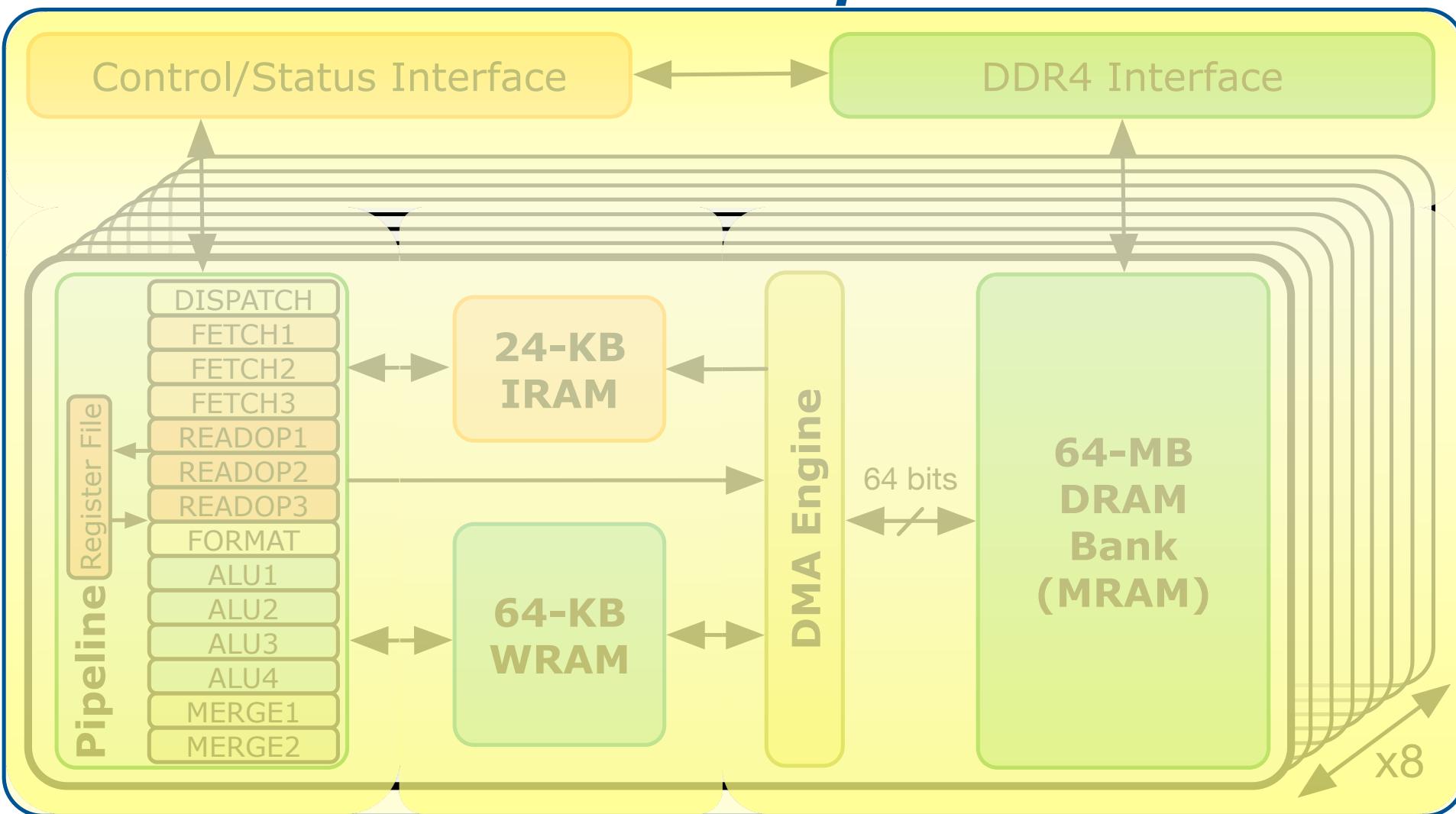


Fig 4

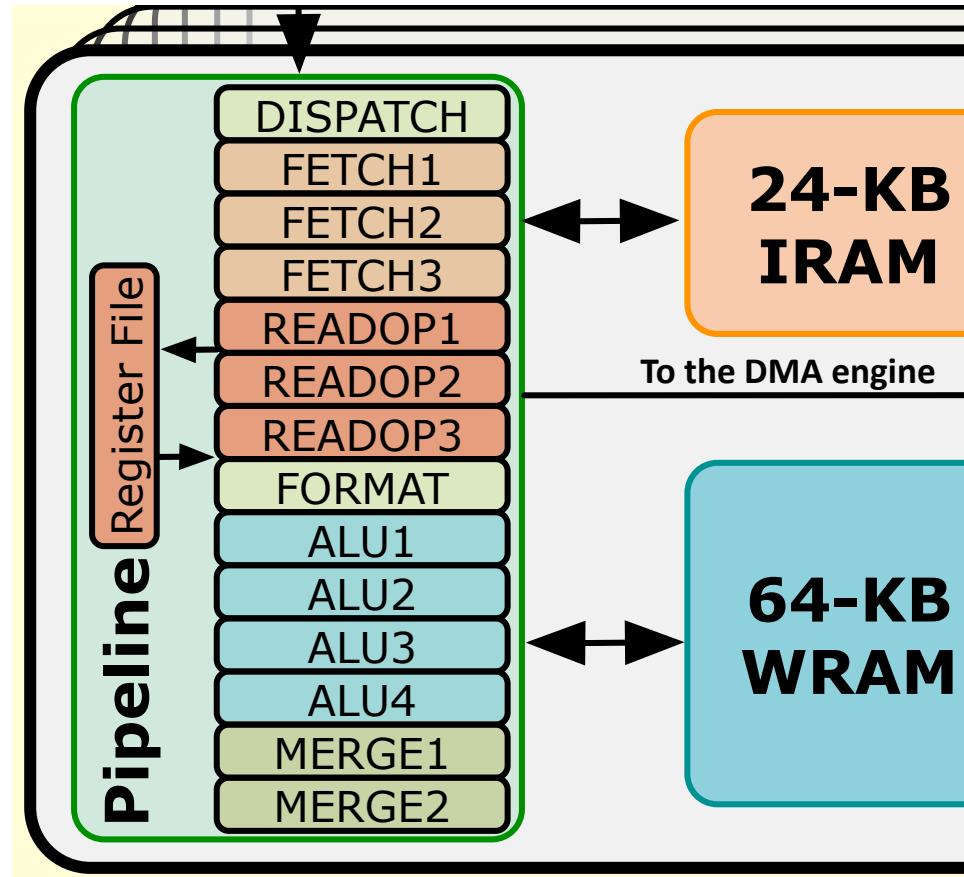
# DRAM Processing Unit (II)

## PIM Chip



# DPU Pipeline

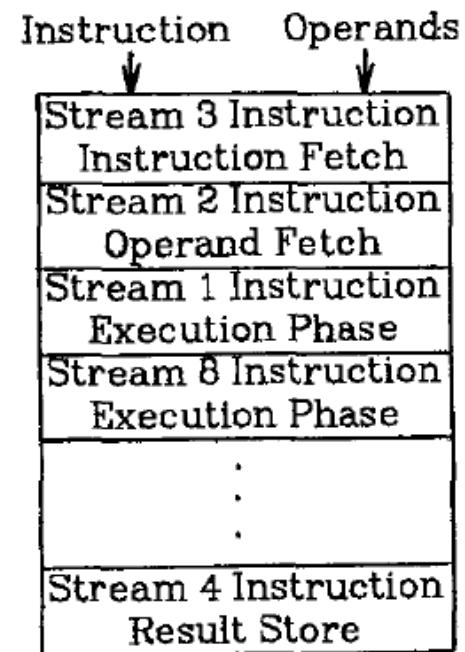
- In-order pipeline
  - Up to 425 MHz
- Fine-grain multithreaded
  - 24 hardware threads
- 14 pipeline stages
  - DISPATCH: Thread selection
  - FETCH: Instruction fetch
  - READOP: Register file
  - FORMAT: Operand formatting
  - ALU: Operation and WRAM
  - MERGE: Result formatting



# Fine-Grained Multithreading

# Fine-Grained Multithreading

- Idea: Hardware has multiple thread contexts (PC+registers).  
Each cycle, fetch engine fetches from a different thread.
  - By the time the fetched branch/instruction resolves, no instruction is fetched from the same thread
  - Branch/instruction resolution latency overlapped with execution of other threads' instructions
- + No logic needed for handling control and data dependences within a thread
- Single thread performance suffers
- Extra logic for keeping thread contexts
- Does not overlap latency if not enough threads to cover the whole pipeline



# Fine-Grained Multithreading (II)

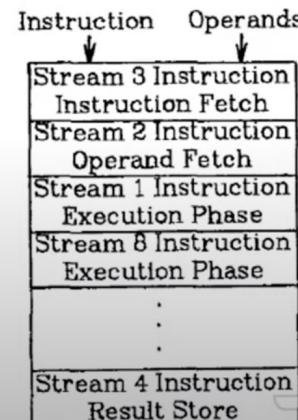
---

- Idea: Switch to another thread every cycle such that no two instructions from a thread are in the pipeline concurrently
- Tolerates the control and data dependence latencies by overlapping the latency with useful work from other threads
- Improves pipeline utilization by taking advantage of multiple threads
- Thornton, “Parallel Operation in the Control Data 6600,” AFIPS 1964.
- Smith, “A pipelined, shared resource MIMD computer,” ICPP 1978.

# Lecture on Fine-Grained Multithreading

## Fine-Grained Multithreading

- Idea: Hardware has multiple thread contexts (PC+registers).  
Each cycle, fetch engine fetches from a different thread.
  - By the time the fetched branch/instruction resolves, no instruction is fetched from the same thread
  - Branch/instruction resolution latency overlapped with execution of other threads' instructions
- + No logic needed for handling control and data dependences within a thread
  - Single thread performance suffers
  - Extra logic for keeping thread contexts
  - Does not overlap latency if not enough threads to cover the whole pipeline



◀ ▶ ⏪ ⏩ 🔍 1:38:38 / 1:57:49

63 CC HD 🔍 ⏴ ⏵ ⏴ ⏵

Onur Mutlu - Digital Design & Comp Arch - Lecture 14: Pipelined Processor Design (Spring 2021)

1,193 views • Streamed live on Apr 22, 2021

42

0

SHARE

SAVE

...



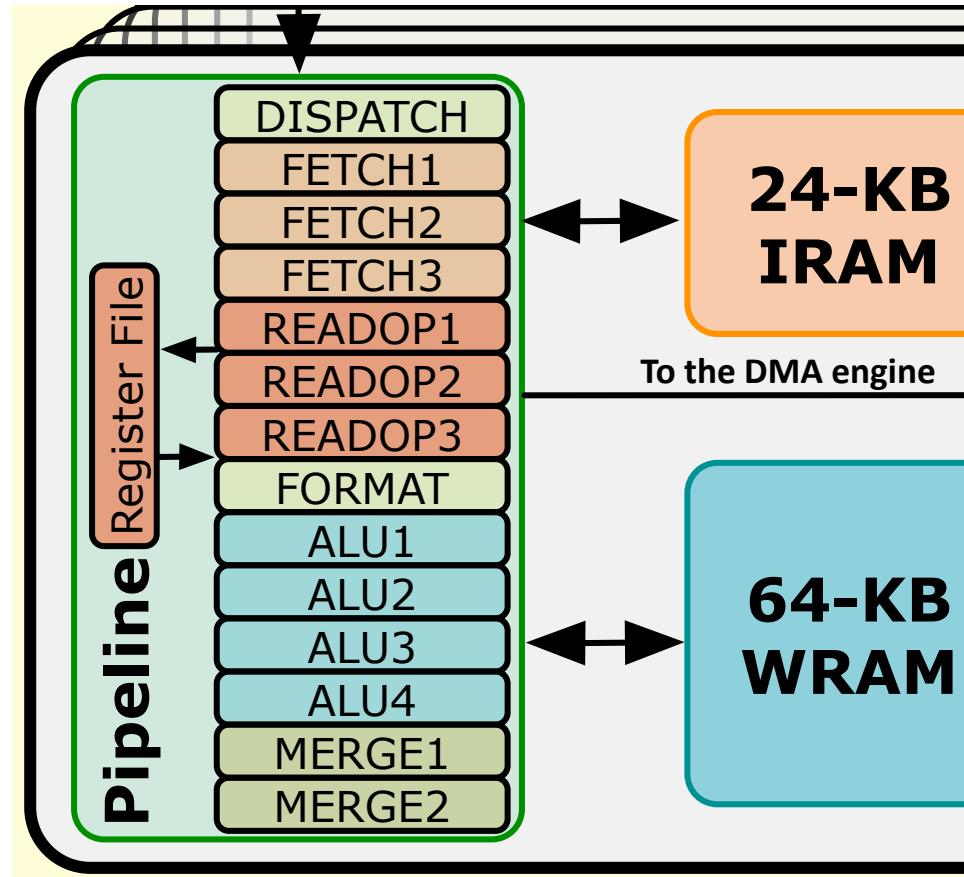
Onur Mutlu Lectures  
16.2K subscribers

ANALYTICS

EDIT VIDEO

# DPU Pipeline

- In-order pipeline
  - Up to 425 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



# DPU Instruction Set Architecture

- Specific 32-bit ISA
  - Aiming at scalar, in-order, and multithreaded implementation
  - Allowing compilation of 64-bit C code
  - LLVM/Clang compiler

The screenshot shows a documentation page for the UPMEM DPU SDK. The header includes the logo 'U' for Instruction Set Architecture, the title 'Instruction Set Architecture — UPMEM DPU SDK 2021.2.0 Documentation', and a navigation bar with a menu icon and the text 'UPMEM development tools documentation'. Below the header, there is a breadcrumb trail ('» Instruction Set Architecture') and a link to 'View page source'. The main content area features a section titled 'Instruction Set Architecture' which describes the architecture concepts required for software developers. It also mentions that the section serves as a reference manual for assembly code. Below this, there are sections for 'Resources overview' and 'Thread registers'. The 'Thread registers' section discusses the 24 hardware threads and their private resources, including 24 general purpose 32-bit registers (r0 through r23), a 16-bit wide program counter (PC), and two persistent flags (ZF). At the bottom of the page, there is a footer with a 'Display a menu' button and a link to the full URL: [https://sdk.upmem.com/2021.2.0/201\\_IS.html#](https://sdk.upmem.com/2021.2.0/201_IS.html#).

# 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

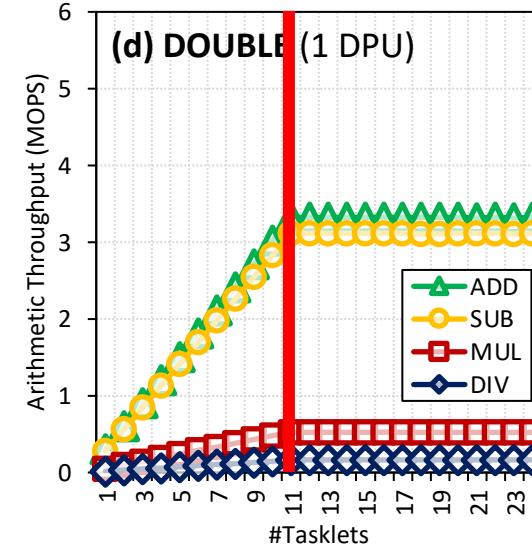
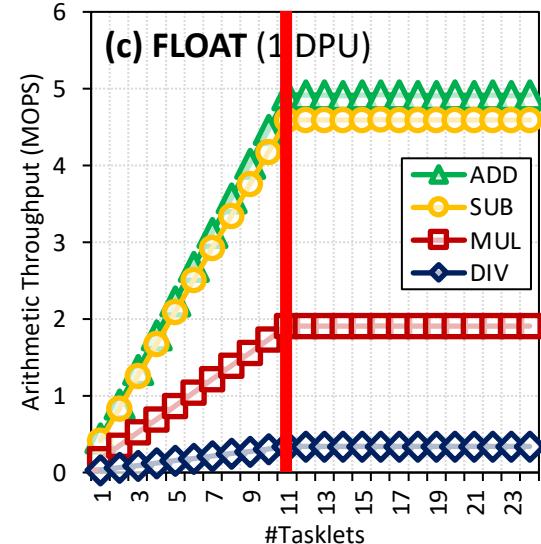
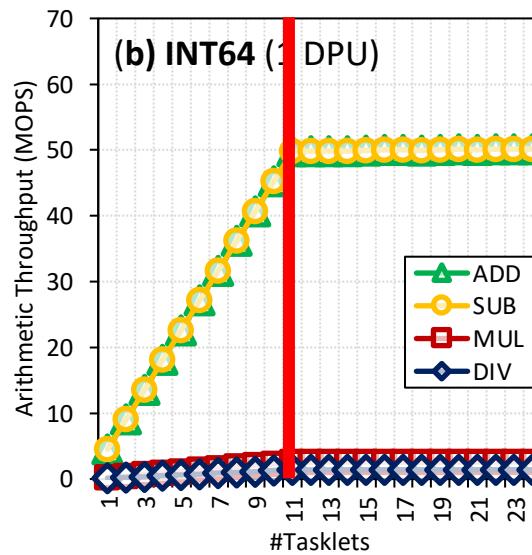
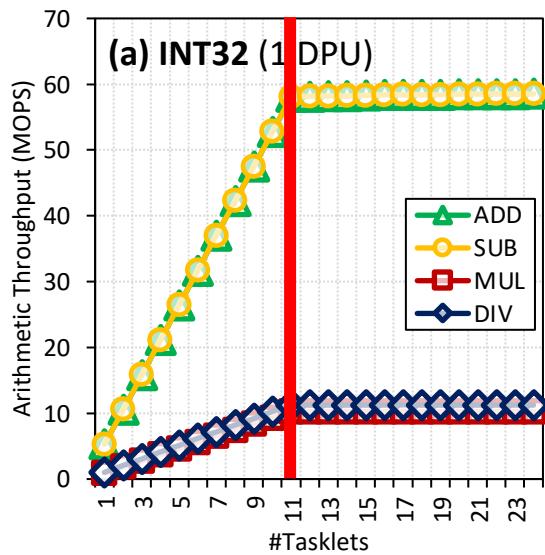
C-based code

```
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

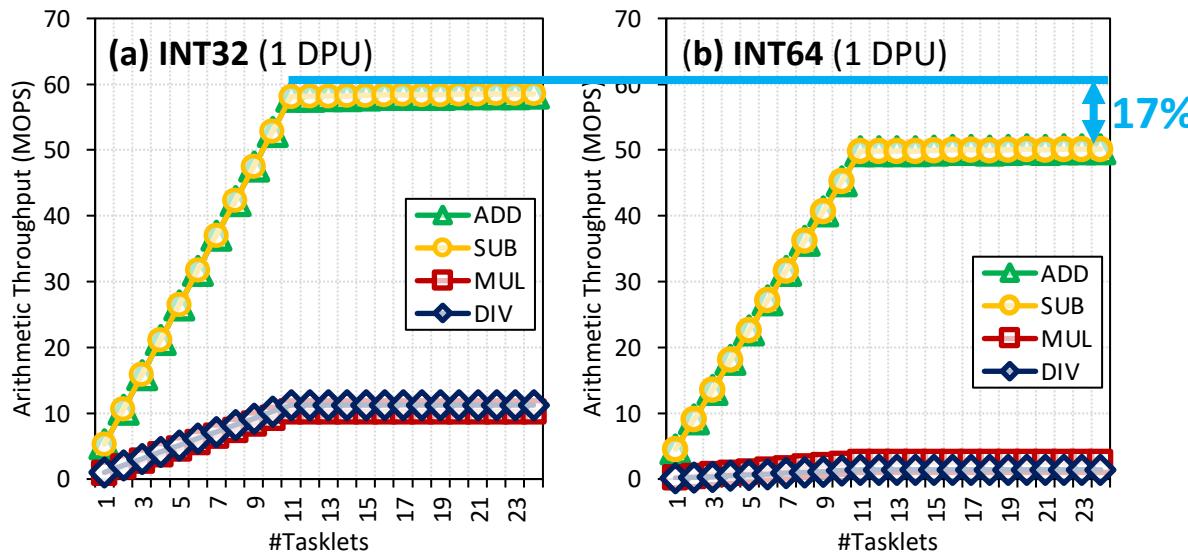


## KEY OBSERVATION 1

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

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



# 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}
```

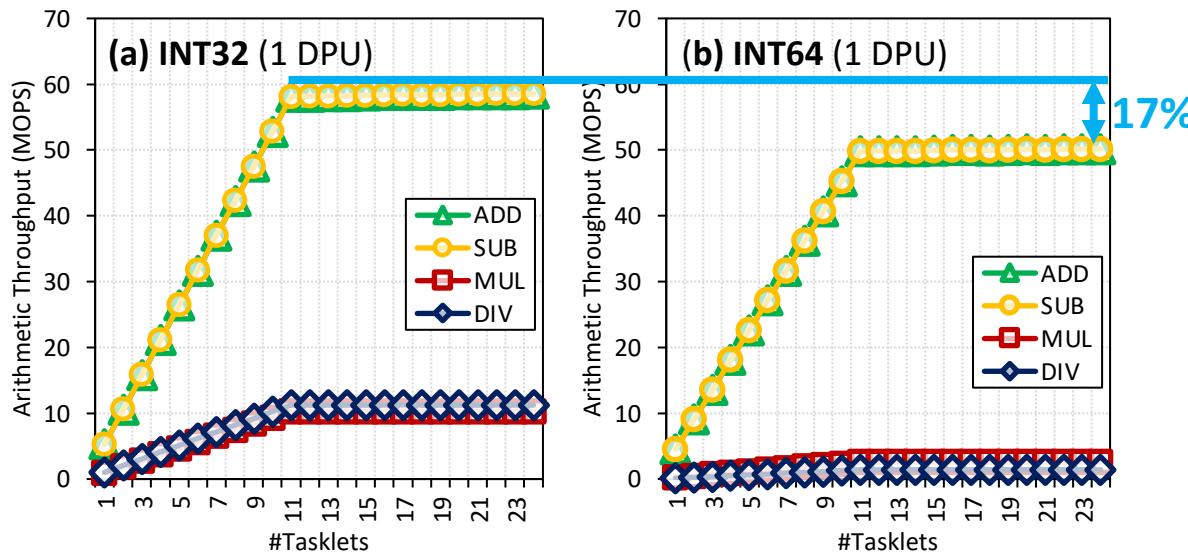


```
A □ 11010 □ ./a.out ✓ .LX0: ✓ .text ✓ // □ \
1 Benchmark_32bits:
2     move r2, 0
3 .LBB0_1:
4     lsl_add r3, r0, r2, 2
5     lw r4, r3, 0
6     add r4, r4, r1
7     sw r3, 0, r4
8     add r2, r2, 1
9     jneq r2, 256, .LBB0_1
10    jump r23
11 Benchmark_64bits:
12     move r1, 0
13 .LBB1_1:
14     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?

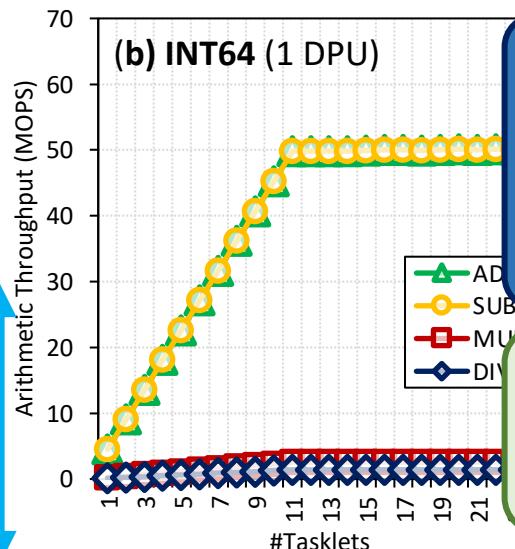
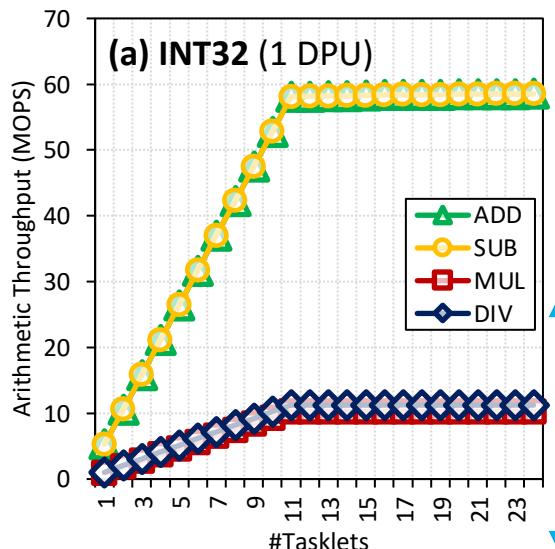
Peak throughput at 11 tasklets.

One instruction retires every cycle when the pipeline is full

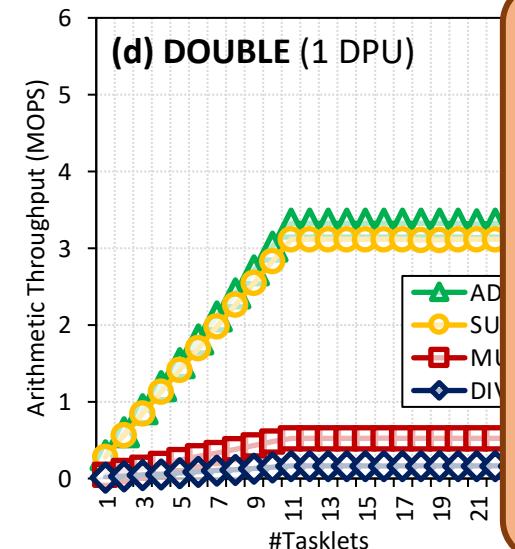
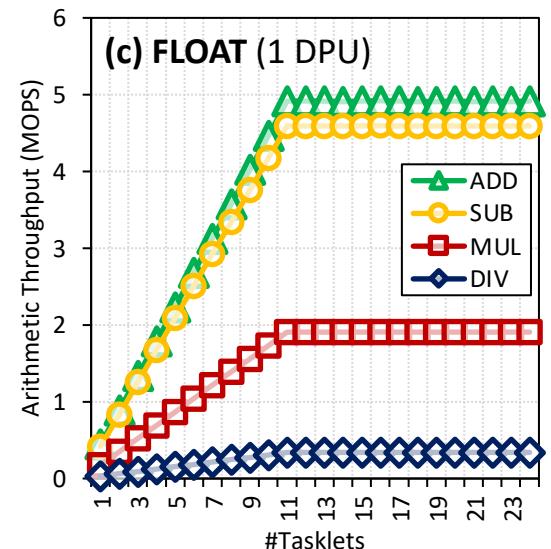
$$\text{Arithmetic Throughput (in OPS)} = \frac{\text{frequency}_{DPU}}{\#\text{instructions}}$$

64-bit ADD/SUB: 7 instructions  $\rightarrow$  50.00 MOPS  
at  $\text{frequency}_{DPU} = 350$  MHz

# Arithmetic Throughput: MUL/DIV

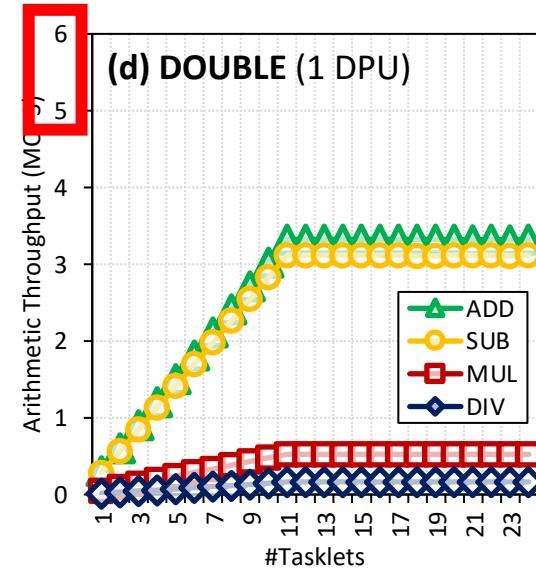
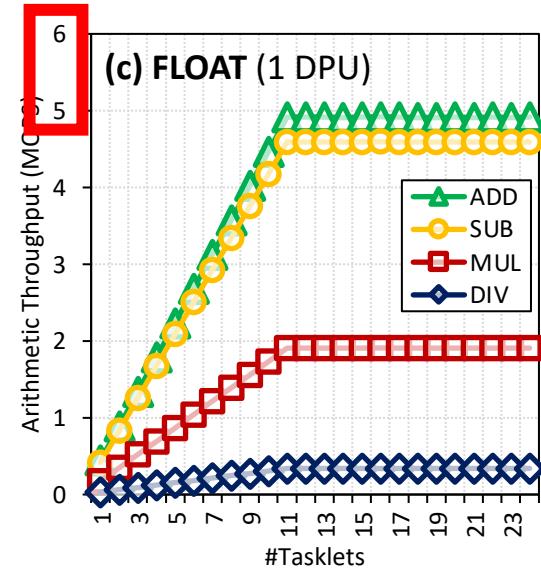
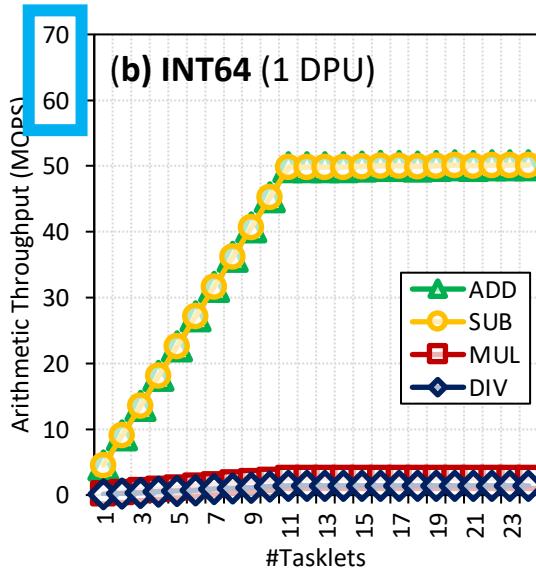
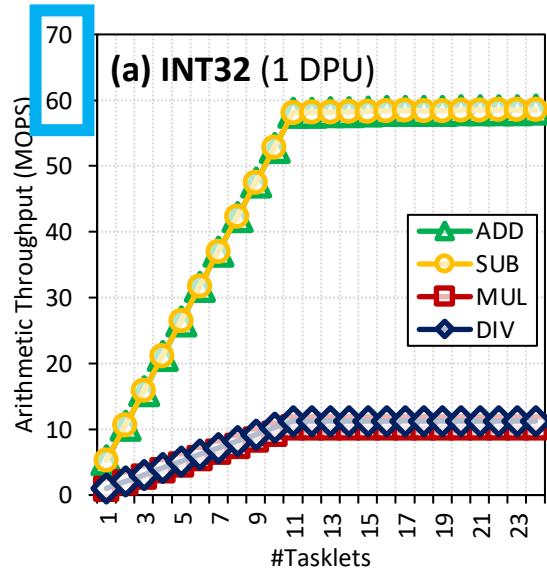


Huge throughput difference between ADD/SUB and MUL/DIV



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.

# Microbenchmark: Arithmetic Throughput

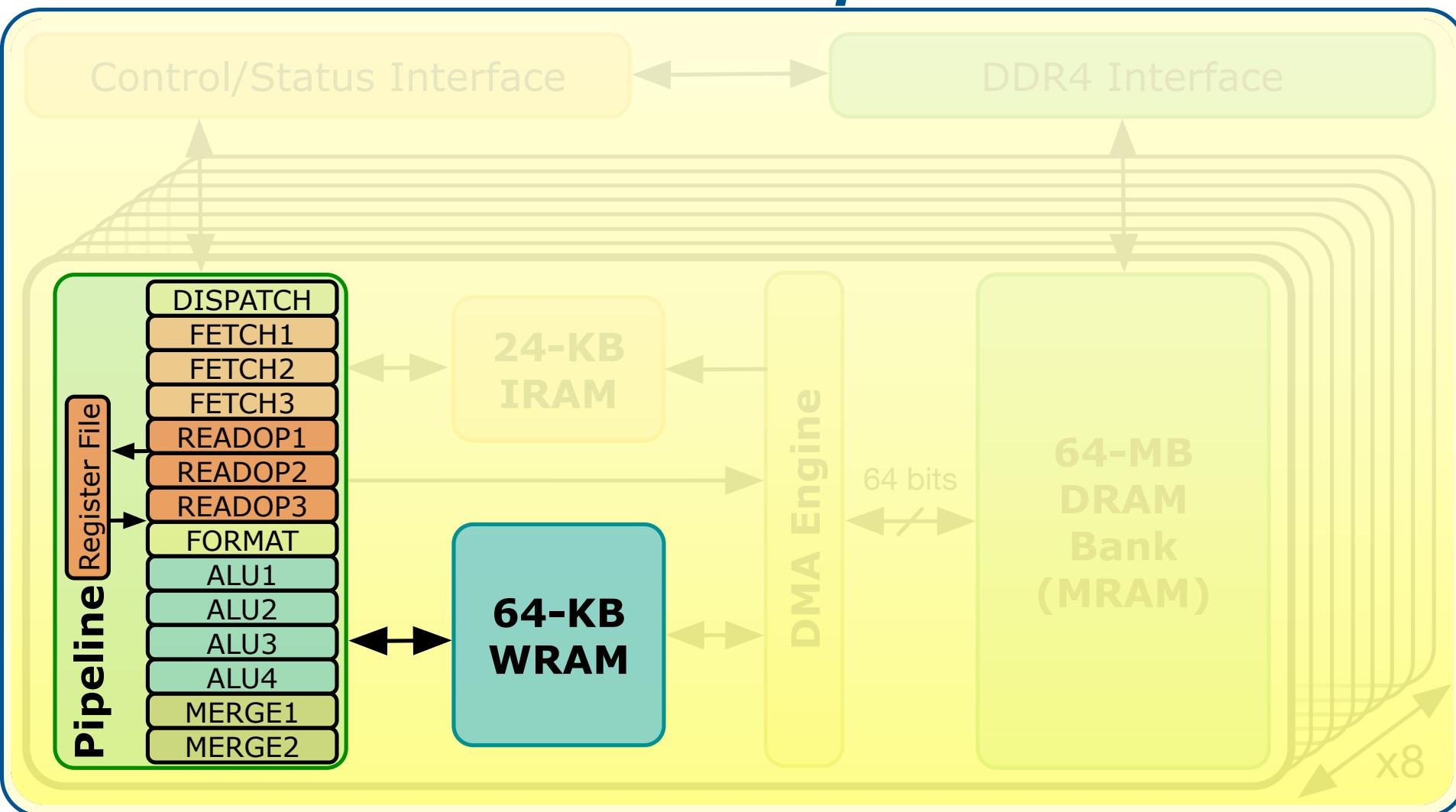
- Arithmetic throughput for different operations and datatypes

The screenshot shows a GitHub repository page for `CMU-SAFARI / prim-benchmarks`. The repository has 2 stars and 1 fork. The main tab is selected, showing the commit history for the `main` branch. The commits are as follows:

Commit	Author	Date	Message
3de4b49	Juan Gomez Luna	9 days ago	PrIM -- first commit
..			
dpu		9 days ago	PrIM -- first commit
host		9 days ago	PrIM -- first commit
support		9 days ago	PrIM -- first commit
Makefile		9 days ago	PrIM -- first commit
run.sh		9 days ago	PrIM -- first commit

# 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];
}

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

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

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

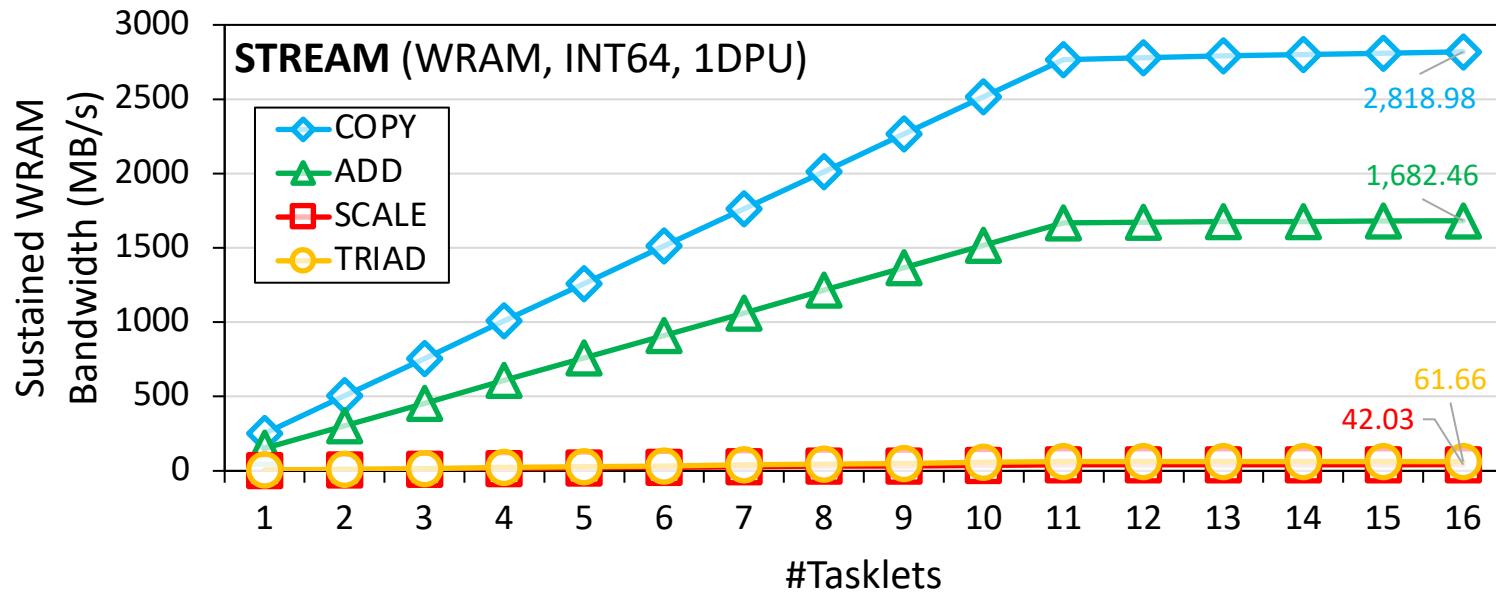
8 bytes read, 8 bytes written,  
no arithmetic operations

16 bytes read, 8 bytes written,  
ADD

8 bytes read, 8 bytes written,  
MUL

16 bytes read, 8 bytes written,  
MUL, ADD

# WRAM Bandwidth: STREAM

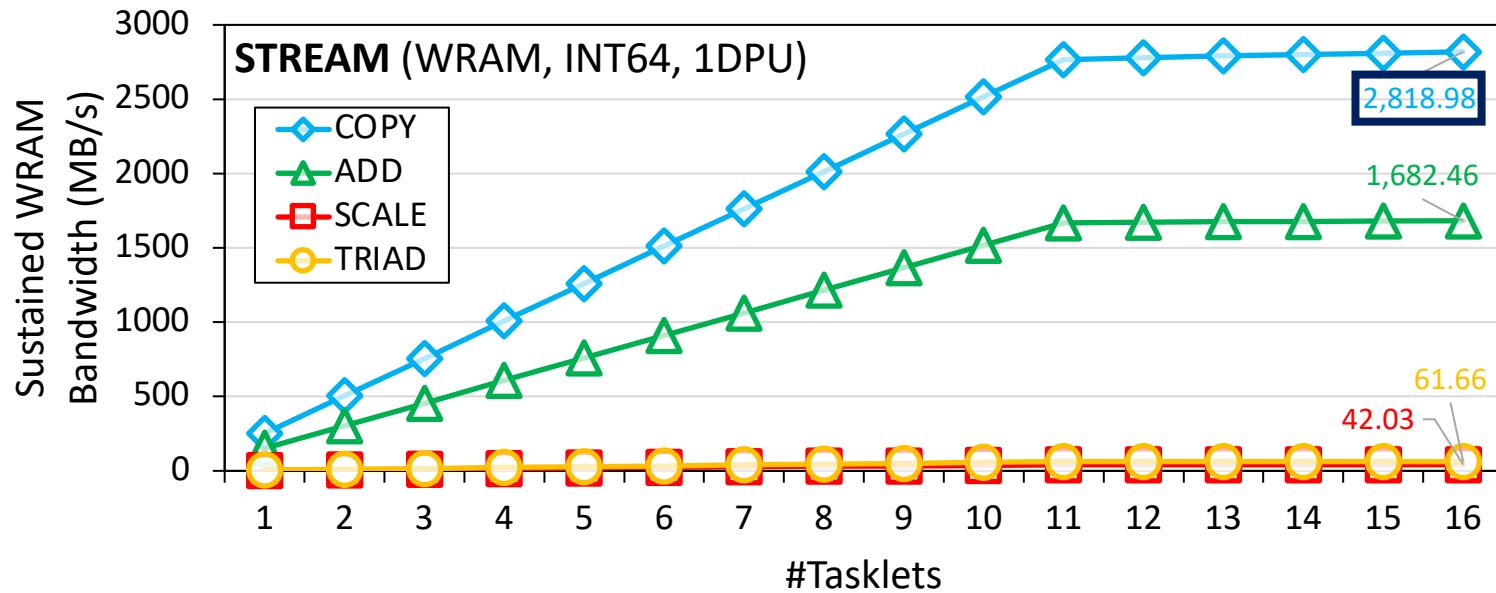


How can we estimate the bandwidth?

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

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

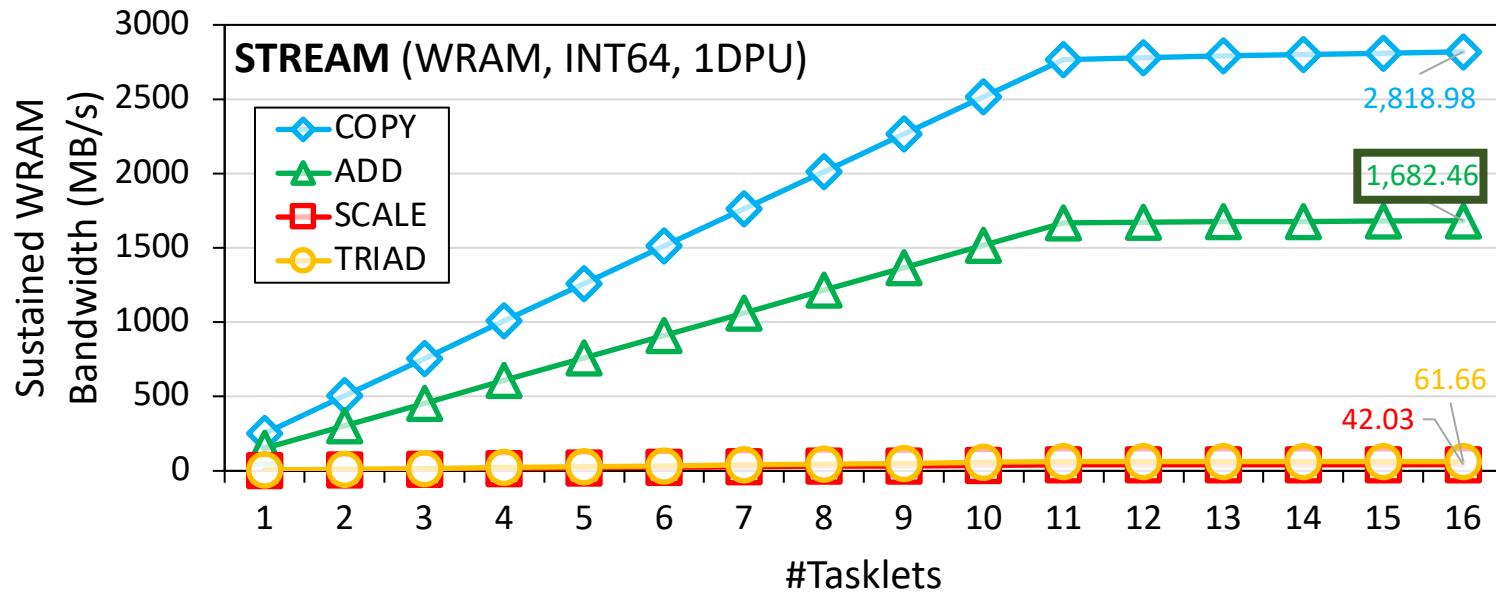
# WRAM Bandwidth: COPY



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

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

# WRAM Bandwidth: ADD



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

**ADD** executes **5 instructions** (2 ld, add, addc, sd).

With 11 tasklets, **11 × 24 bytes in 55 cycles**:

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

# WRAM Bandwidth: Access Patterns

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

## KEY OBSERVATION 3

The sustained bandwidth provided by the DPU's internal Working memory (WRAM) is **independent of the memory access pattern** (either streaming, strided, or random access pattern).

**All 8-byte WRAM loads and stores take one cycle**, when the DPU's pipeline is full (i.e., with 11 or more tasklets).

- Microbenchmark: `c[a[i]]=b[a[i]];`
  - Unit-stride: `a[i]=a[i-1]+1;`
  - Strided: `a[i]=a[i-1]+stride;`
  - Random: `a[i]=rand();`

# Microbenchmark: STREAM and WRAM

- STREAM benchmark and WRAM access patterns

Screenshot of a GitHub repository page for "CMU-SAFARI / prim-benchmarks".

The repository has 2 stars, 1 fork, and 1 issue.

Navigation tabs: Code (selected), Issues, Pull requests, Actions, Projects, Wiki, Security, Insights, Settings.

Breadcrumbs: CMU-SAFARI / prim-benchmarks / Microbenchmarks / STREAM /

File navigation: Go to file, Add file, ...

Commit history for the STREAM folder:

Author	Message	Date	Actions
Juan Gomez Luna	PrIM -- first commit	3de4b49 9 days ago	History
...			
dpu	PrIM -- first commit	9 days ago	
host	PrIM -- first commit	9 days ago	
support	PrIM -- first commit	9 days ago	
Makefile	PrIM -- first commit	9 days ago	
run.sh	PrIM -- first commit	9 days ago	

Breadcrumbs for the WRAM folder:

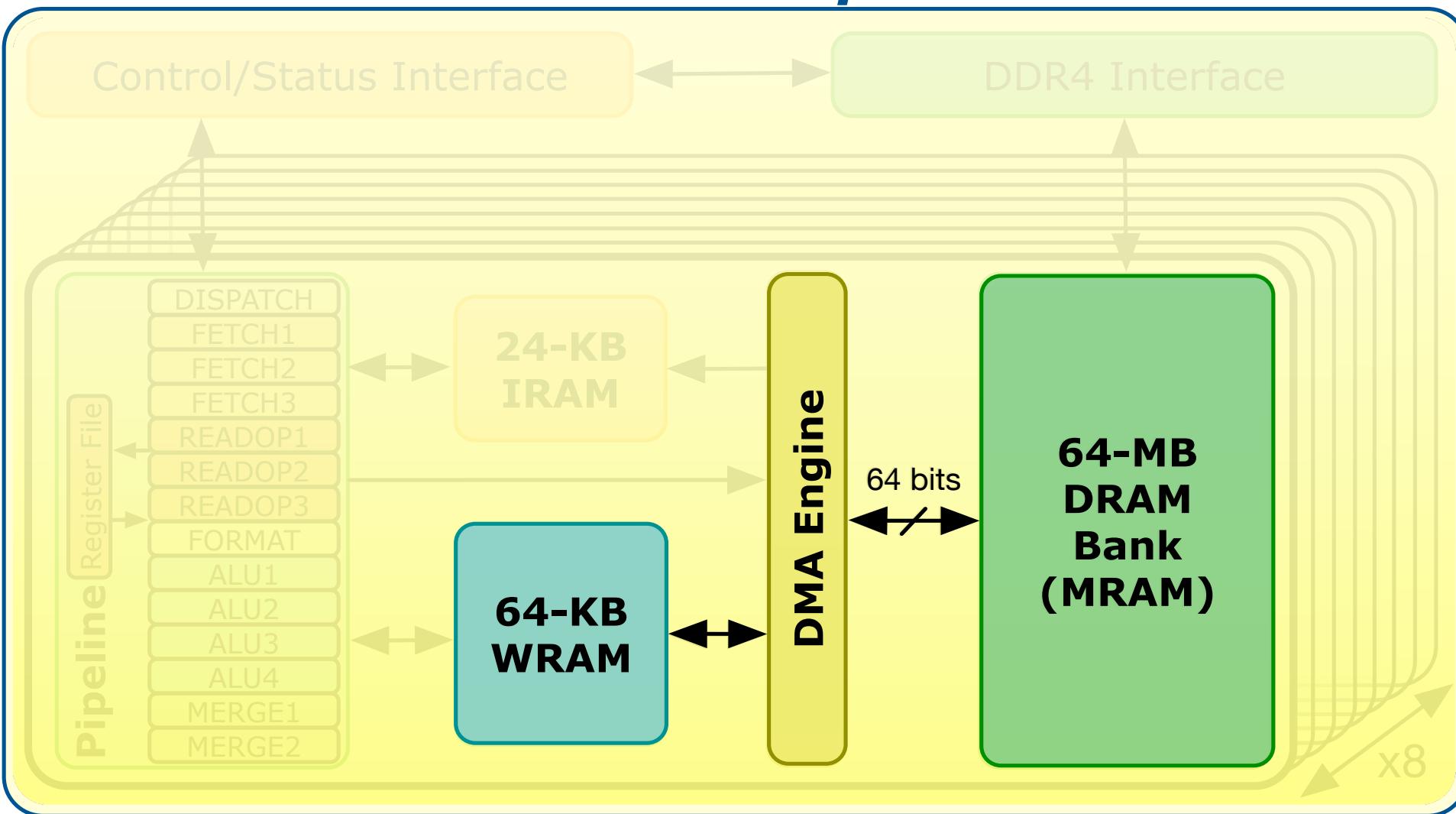
File navigation: Go to file, Add file, ...

Commit history for the WRAM folder:

Author	Message	Date	Actions
Juan Gomez Luna	PrIM -- first commit	3de4b49 9 days ago	History
...			

# DPU: MRAM Latency and Bandwidth

## PIM Chip

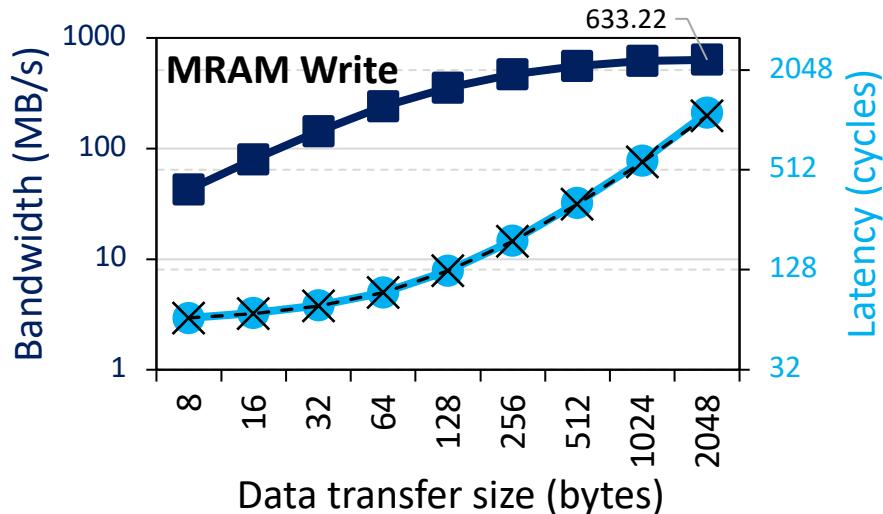
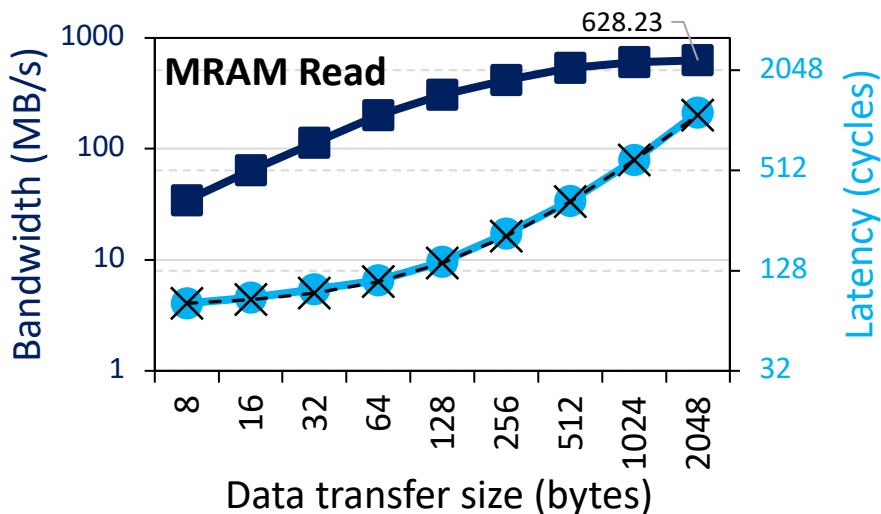


# MRAM Bandwidth

---

- 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)



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

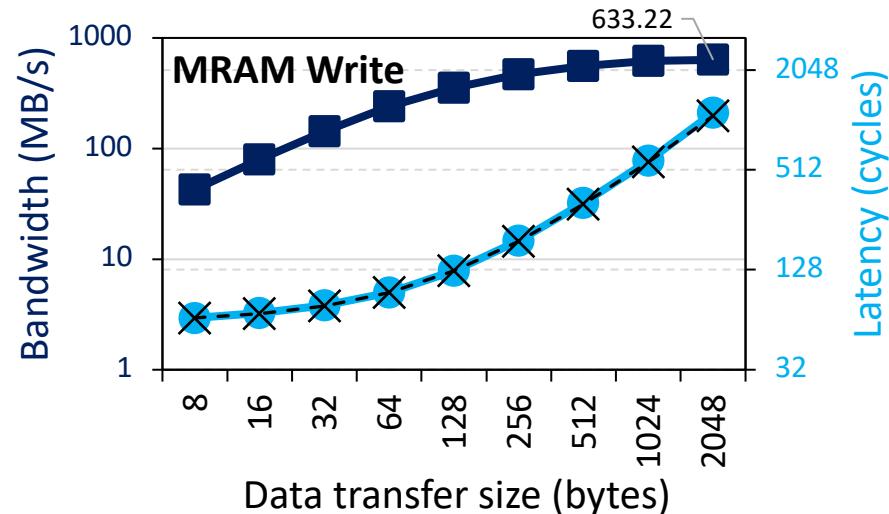
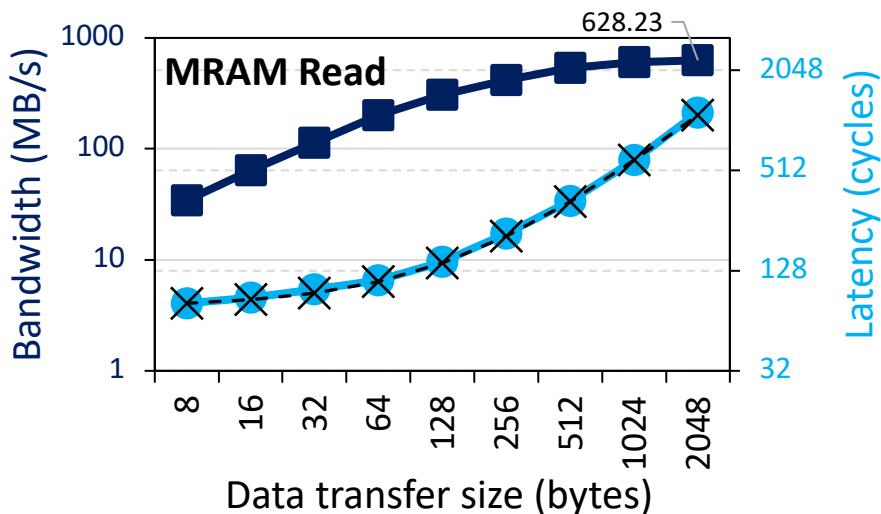
We can model the MRAM latency with a linear expression

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

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

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

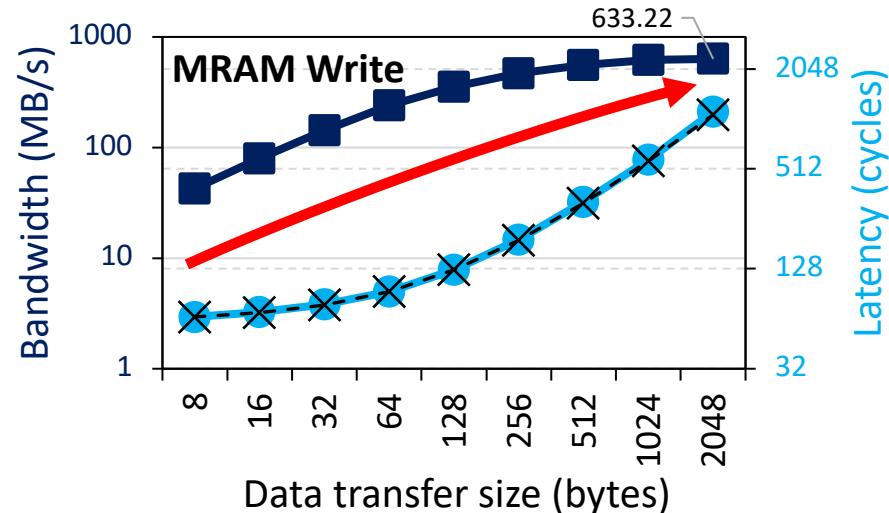
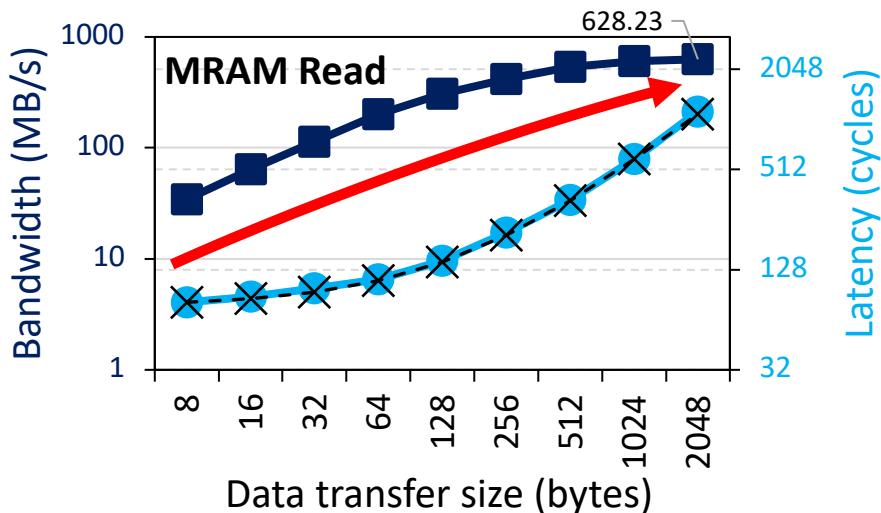
# MRAM Read and Write Latency (II)



## 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)



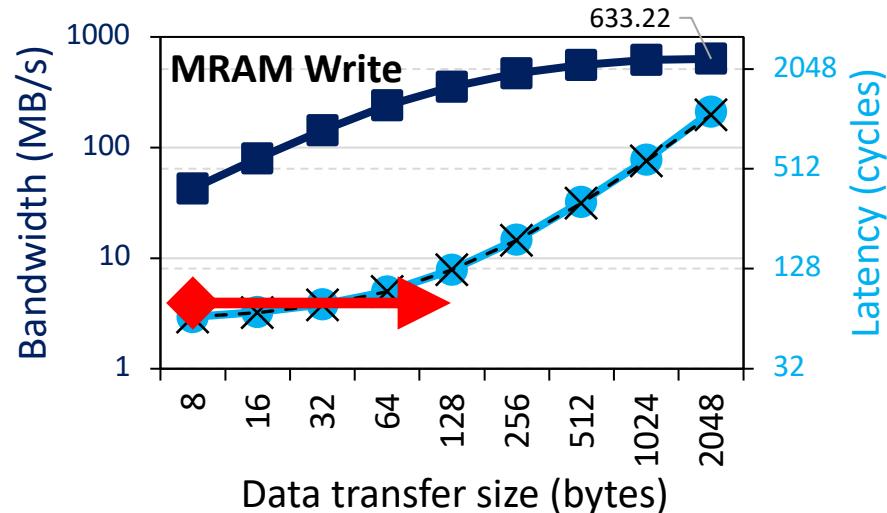
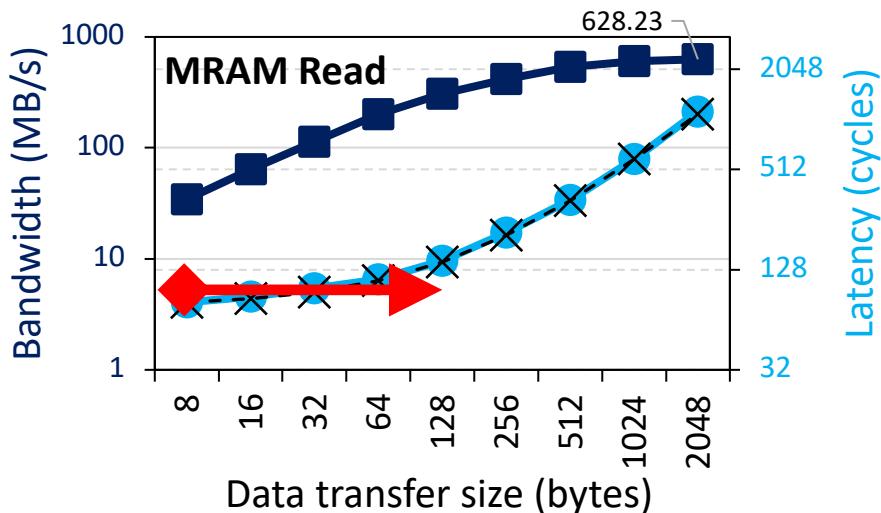
Read and write accesses to MRAM are symmetric

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)



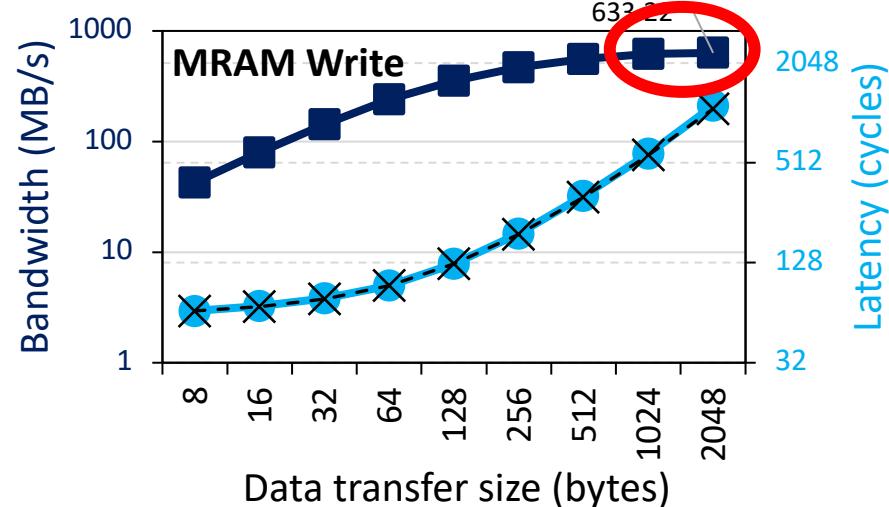
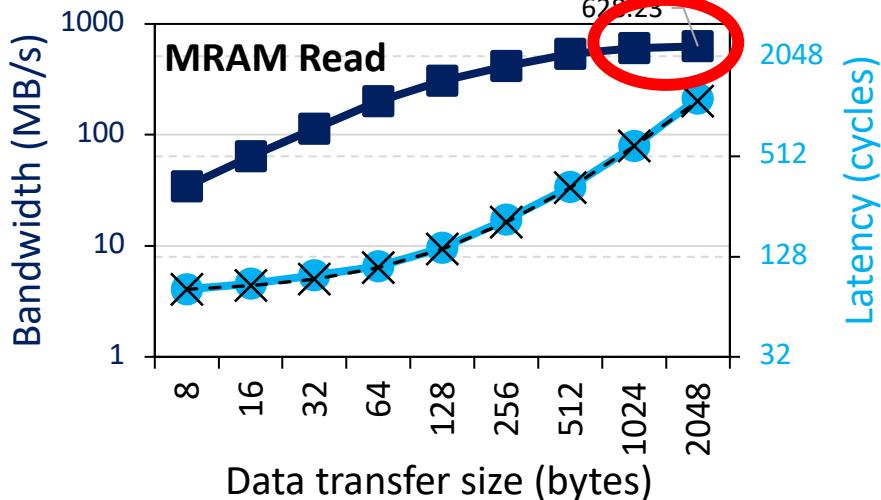
MRAM latency changes slowly between 8 and 128 bytes

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).

# 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));

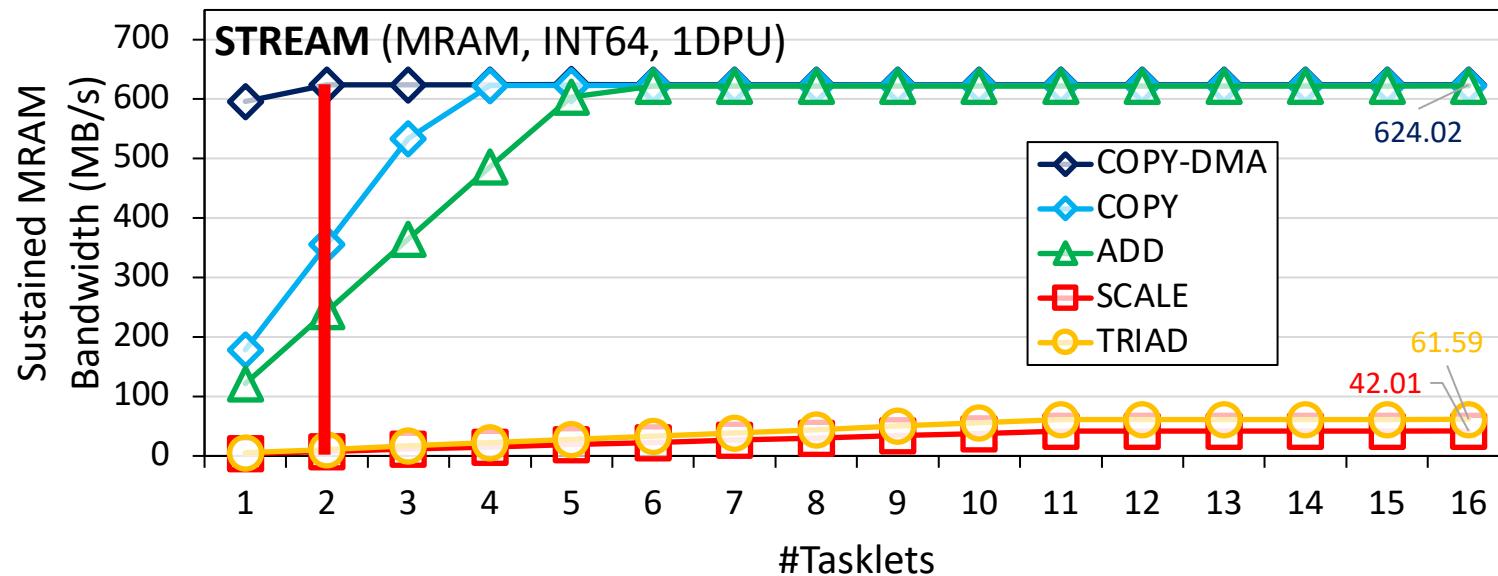
for(int i = 0; i < SIZE; i++){
    bufferB[i] = bufferA[i];
}

// 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(bufferA, (__mram_ptr void*)mram_address_B,
           SIZE * sizeof(uint64_t));
```

# STREAM Benchmark: COPY-DMA

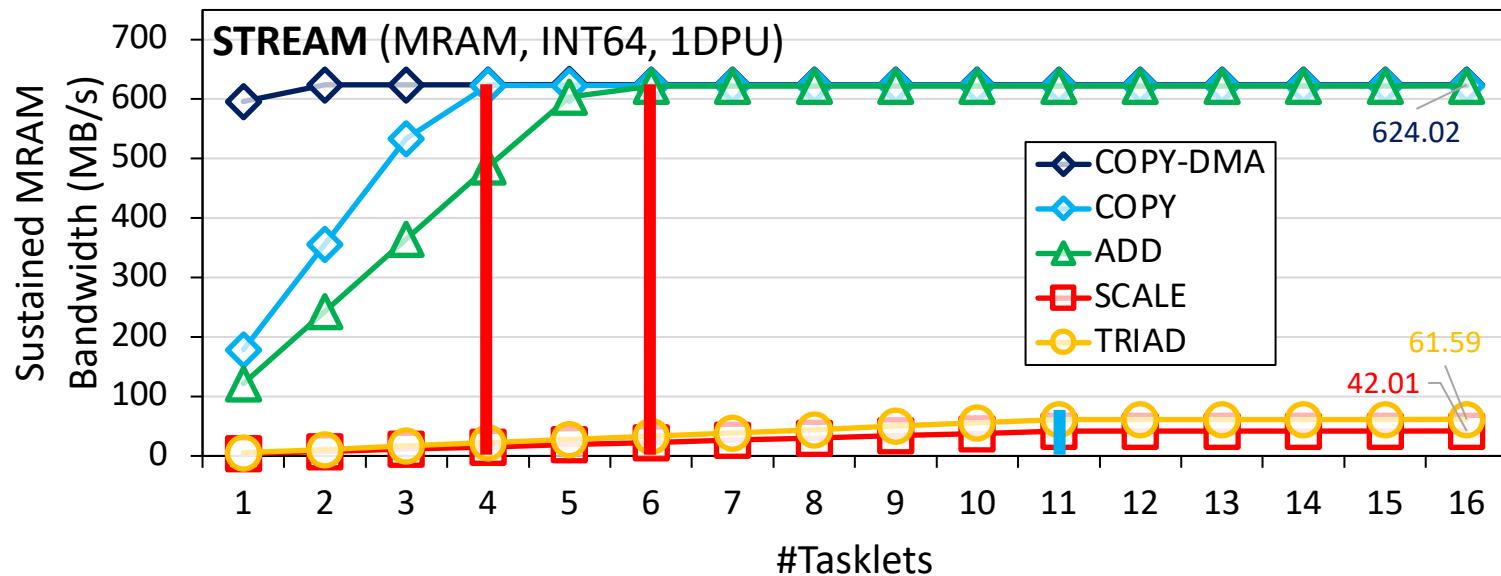


The sustained bandwidth of **COPY-DMA** is close to the theoretical maximum (700 MB/s): ~1.6 TB/s for 2,556 DPUs

**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)



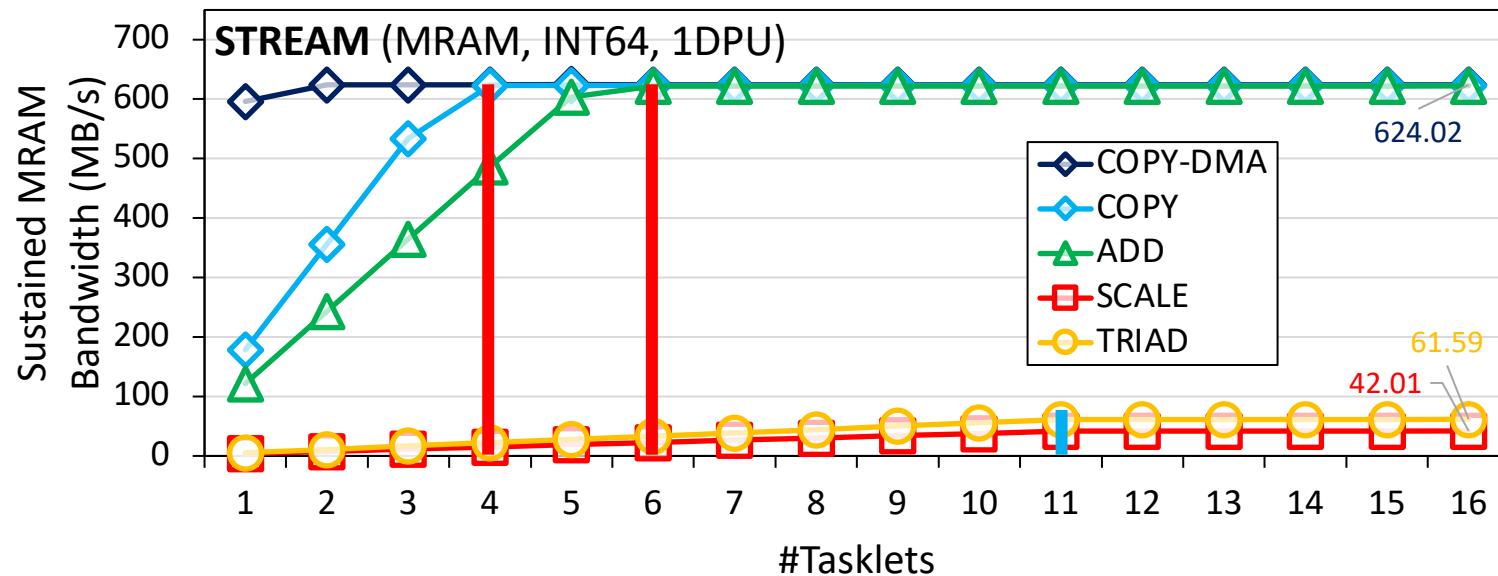
**COPY** and **ADD** saturate at 4 and 6 tasklets, respectively

**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

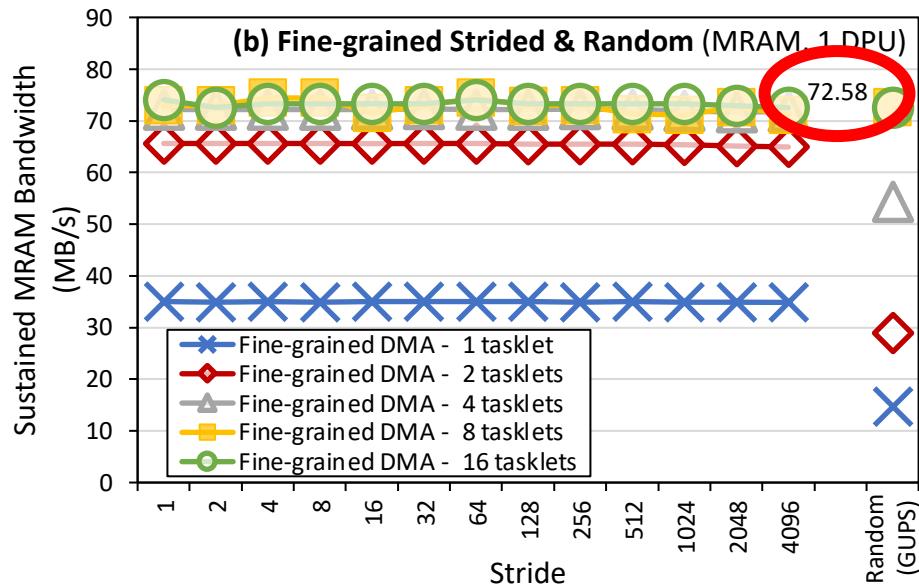
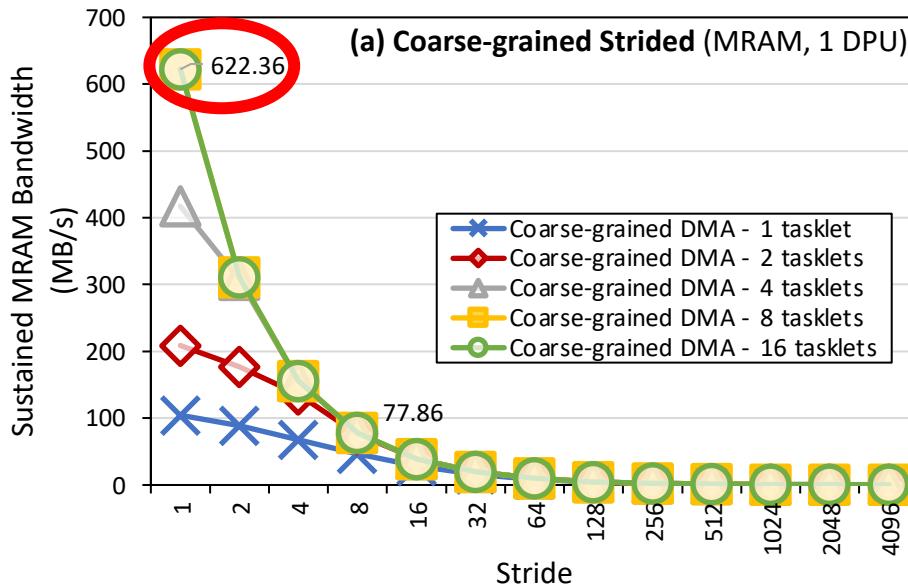
```
// COARSE-GRAINED STRIDED ACCESS
// Load current MRAM block to WRAM
mram_read((__mram_ptr void const*)mram_address_A, bufferA,
           SIZE * sizeof(uint64_t));
mram_read((__mram_ptr void const*)mram_address_B, bufferB,
           SIZE * sizeof(uint64_t));

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)

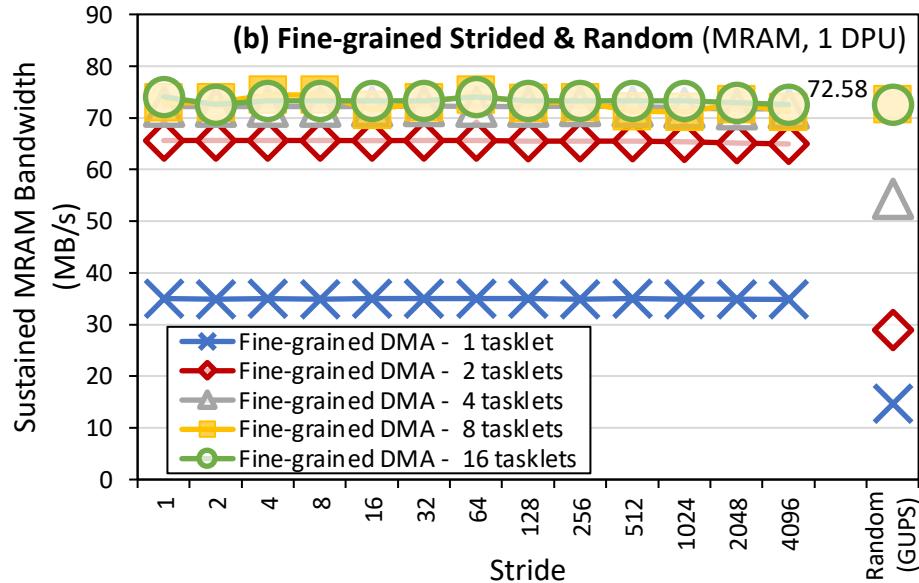
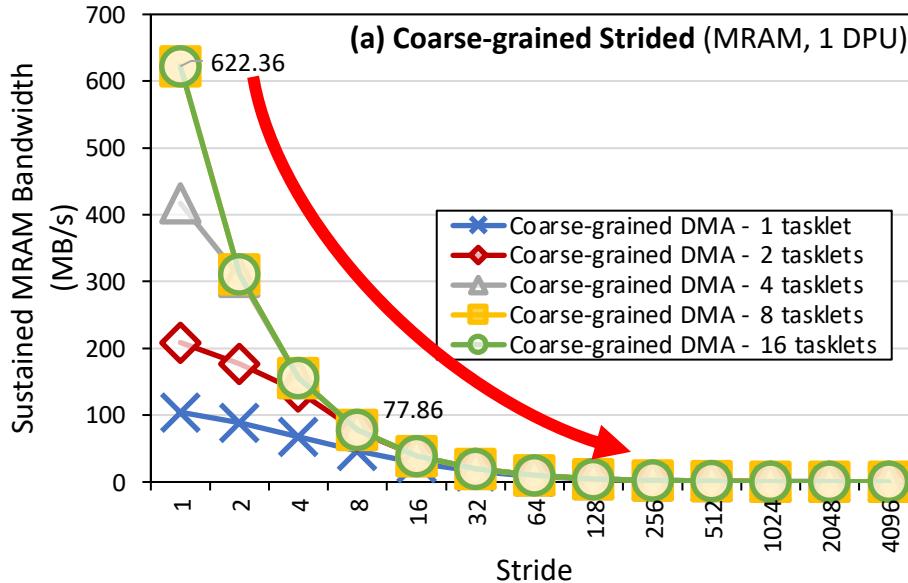


Large difference in maximum sustained bandwidth between coarse-grained and fine-grained DMA

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

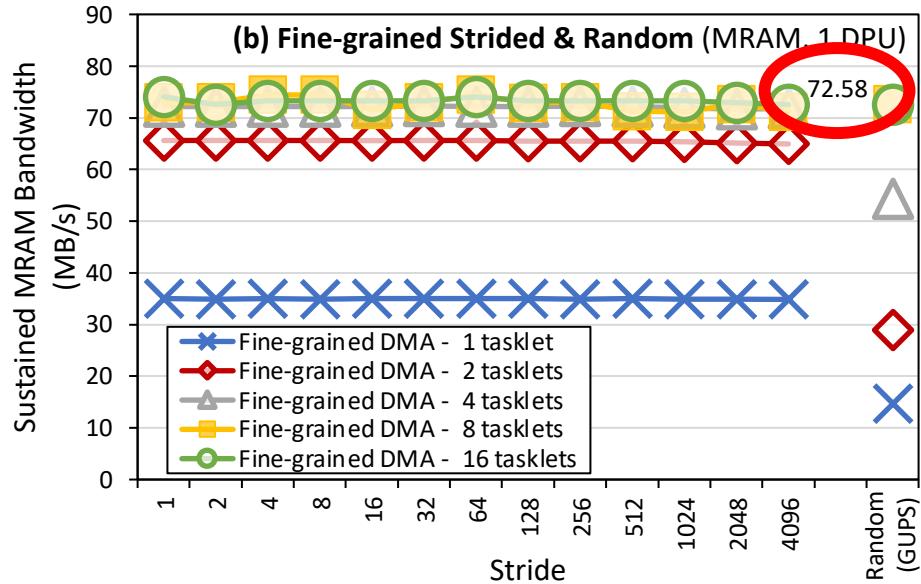
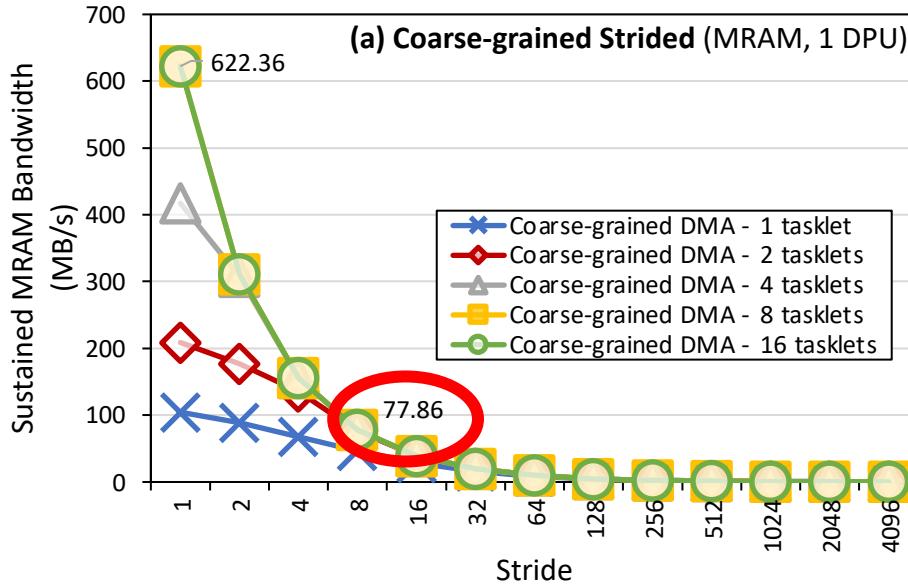
# Strided and Random Accesses (II)



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)

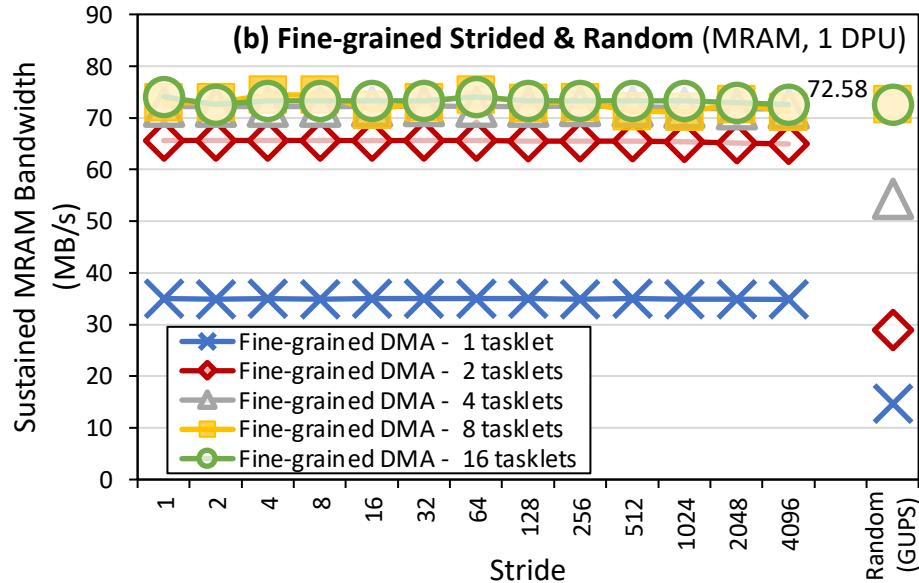
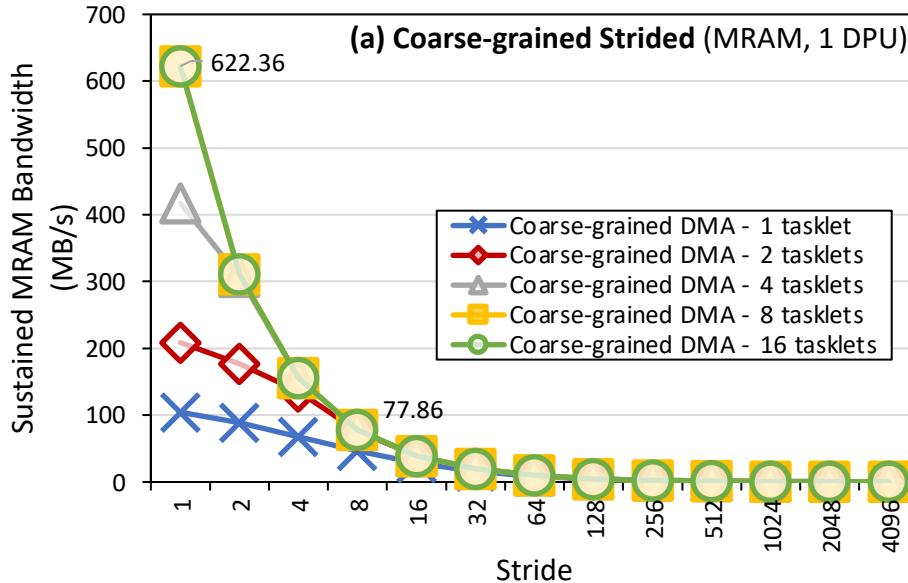
# Strided and Random Accesses (III)



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.

# Microbenchmark: Strided and Random

- Strided and random accesses to MRAM

CMU-SAFARI / prim-benchmarks

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main prim-benchmarks / Microbenchmarks / STRIDED / Go to file Add file ...

main prim-benchmarks / Microbenchmarks / Random-GUPS / Go to file Add file ...

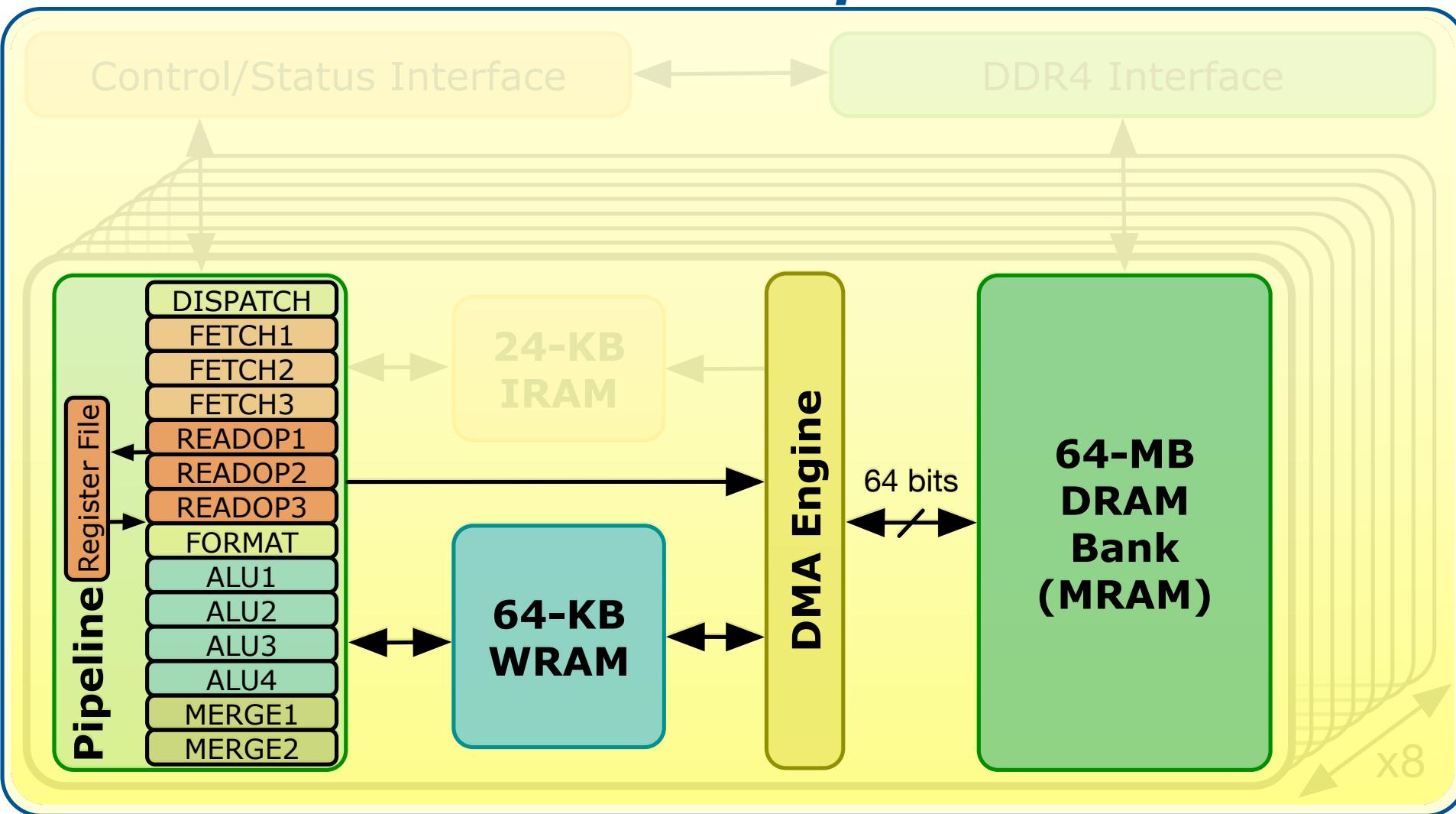
Juan Gomez Luna PrIM -- first commit 3de4b49 9 days ago History

..

dpu	PrIM -- first commit	9 days ago
host	PrIM -- first commit	9 days ago
support	PrIM -- first commit	9 days ago
Makefile	PrIM -- first commit	9 days ago
run.sh	PrIM -- first commit	9 days ago

# 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));

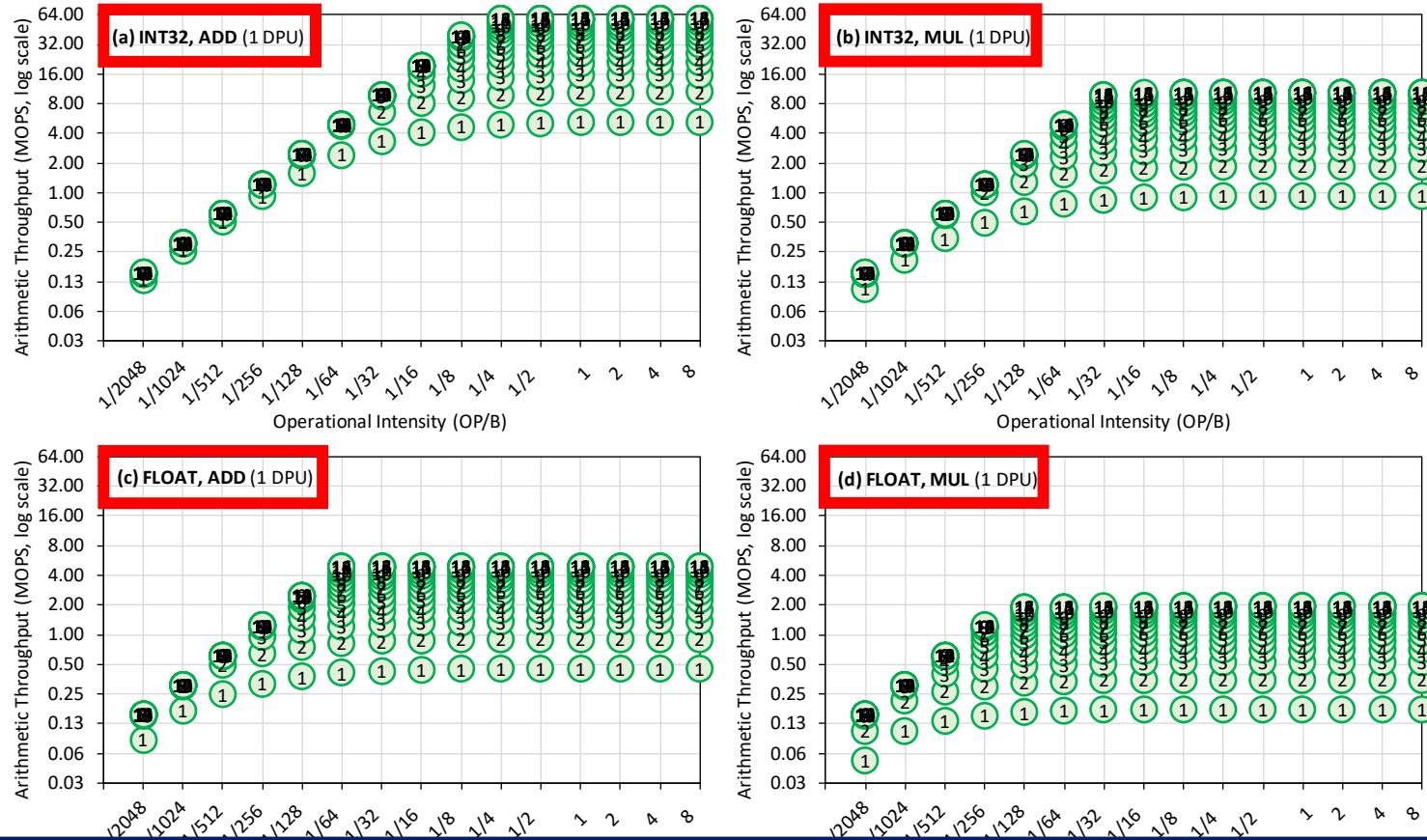
// Update
for(int r = 0; r < repetitions; r++){
    for(int i = 0; i < SIZE; i+=stride){
#ifdef ADD
        bufferA[i] += scalar; // ADD
#elif SUB
        bufferA[i] -= scalar; // SUB
#elif MUL
        bufferA[i] *= scalar; // MUL
#elif DIV
        bufferA[i] /= scalar; // DIV
#endif
    }
}

// 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

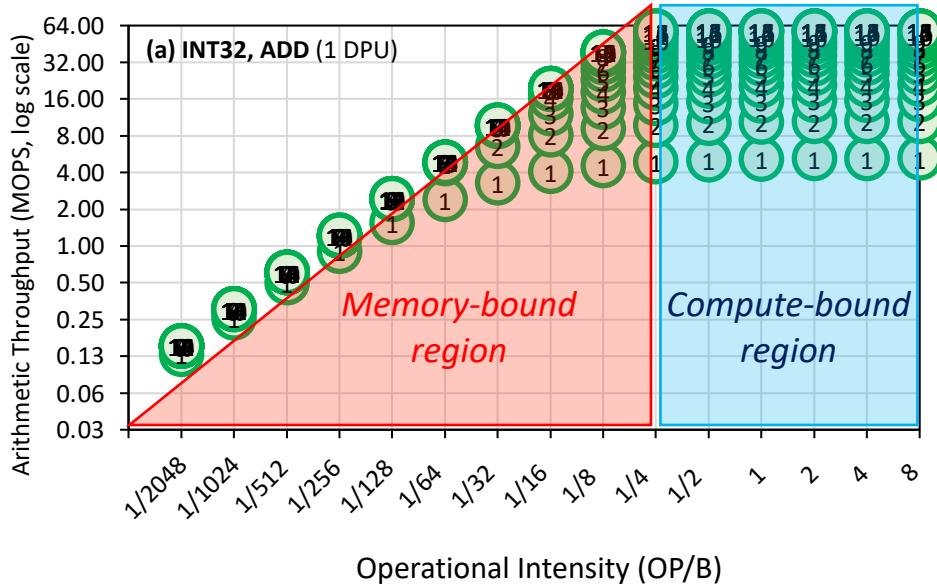
input\_repeat smaller than 1 indicates the fraction of elements that are updated

# 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)



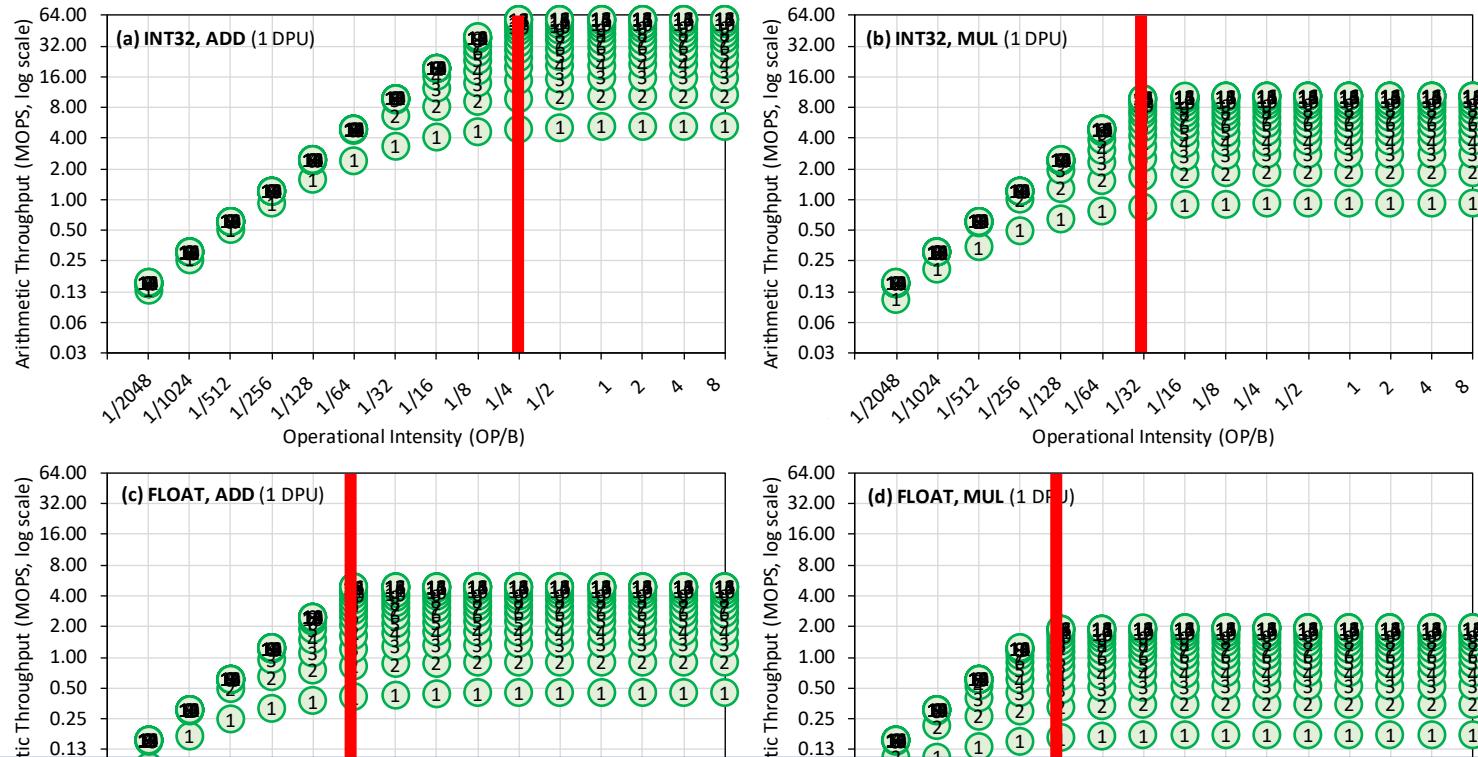
In the **memory-bound region**, the arithmetic throughput increases with the operational intensity

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.

# Microbenchmark: Arithmetic Throughput vs. Operational Intensity

- Arithmetic Throughput versus Operational Intensity

The screenshot shows a GitHub repository page for **CMU-SAFARI / prim-benchmarks**. The repository has 2 stars and 1 fork. The navigation bar includes links for Code, Issues, Pull requests, Actions, Projects, Wiki, Security, Insights, and Settings. The main content area shows the commit history for the **main** branch, which includes a single commit from **Juan Gomez Luna** titled "PrIM -- first commit" (commit hash: 3de4b49, 9 days ago). Below this commit, there are entries for files: dpu, host, support, Makefile, and run.sh, all with the same timestamp of 9 days ago.

File	Commit Message	Time
dpu	PrIM -- first commit	9 days ago
host	PrIM -- first commit	9 days ago
support	PrIM -- first commit	9 days ago
Makefile	PrIM -- first commit	9 days ago
run.sh	PrIM -- first commit	9 days ago

# 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

---

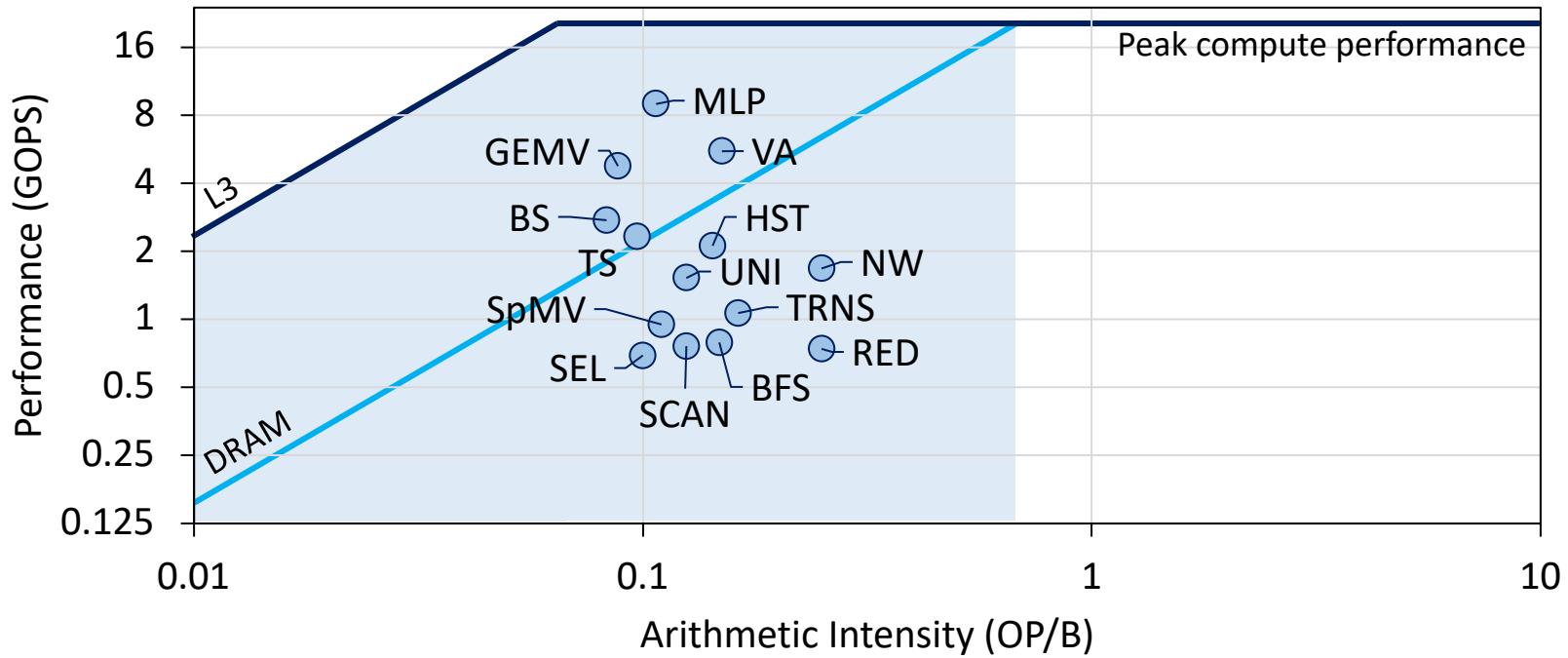
- Goal
  - A **common set of workloads** that can be used to
    - evaluate the UPMEM PIM architecture,
    - compare software improvements and compilers,
    - compare future PIM architectures and hardware
- Two key selection criteria:
  - Selected workloads from **different application domains**
  - **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-Wunck	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-ctrl-id)	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			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

# PrIM Benchmarks

- 16 benchmarks and scripts for evaluation
- <https://github.com/CMU-SAFARI/prim-benchmarks>

The screenshot shows the GitHub repository interface for 'CMU-SAFARI / prim-benchmarks'. The 'Code' tab is selected. At the top, it displays 'main' (1 branch, 0 tags), 'Go to file', 'Add file', and a green 'Code' button. Below is a table of commits:

Author	Commit Message	Date	Commits
Juan Gomez Luna	PrIM -- first commit	3de4b49 15 days ago	2 commits
	BFS	PrIM -- first commit	15 days ago
	BS	PrIM -- first commit	15 days ago
	GEMV	PrIM -- first commit	15 days ago
	HST-L	PrIM -- first commit	15 days ago
	HST-S	PrIM -- first commit	15 days ago
	MLP	PrIM -- first commit	15 days ago
	Microbenchmarks	PrIM -- first commit	15 days ago
	NW	PrIM -- first commit	15 days ago
	RED	PrIM -- first commit	15 days ago
	SCAN-RSS	PrIM -- first commit	15 days ago
	SCAN-SSA	PrIM -- first commit	15 days ago
	SEL	PrIM -- first commit	15 days ago
	SpMV	PrIM -- first commit	15 days ago
	TRNS	PrIM -- first commit	15 days ago
	TS	PrIM -- first commit	15 days ago
	UNI	PrIM -- first commit	15 days ago
	VA	PrIM -- first commit	15 days ago
	LICENSE	PrIM -- first commit	15 days ago
	README.md	PrIM -- first commit	15 days ago
	run_strong_full.py	PrIM -- first commit	15 days ago
	run_strong_rank.py	PrIM -- first commit	15 days ago
	run_weak.py	PrIM -- first commit	15 days ago

# 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

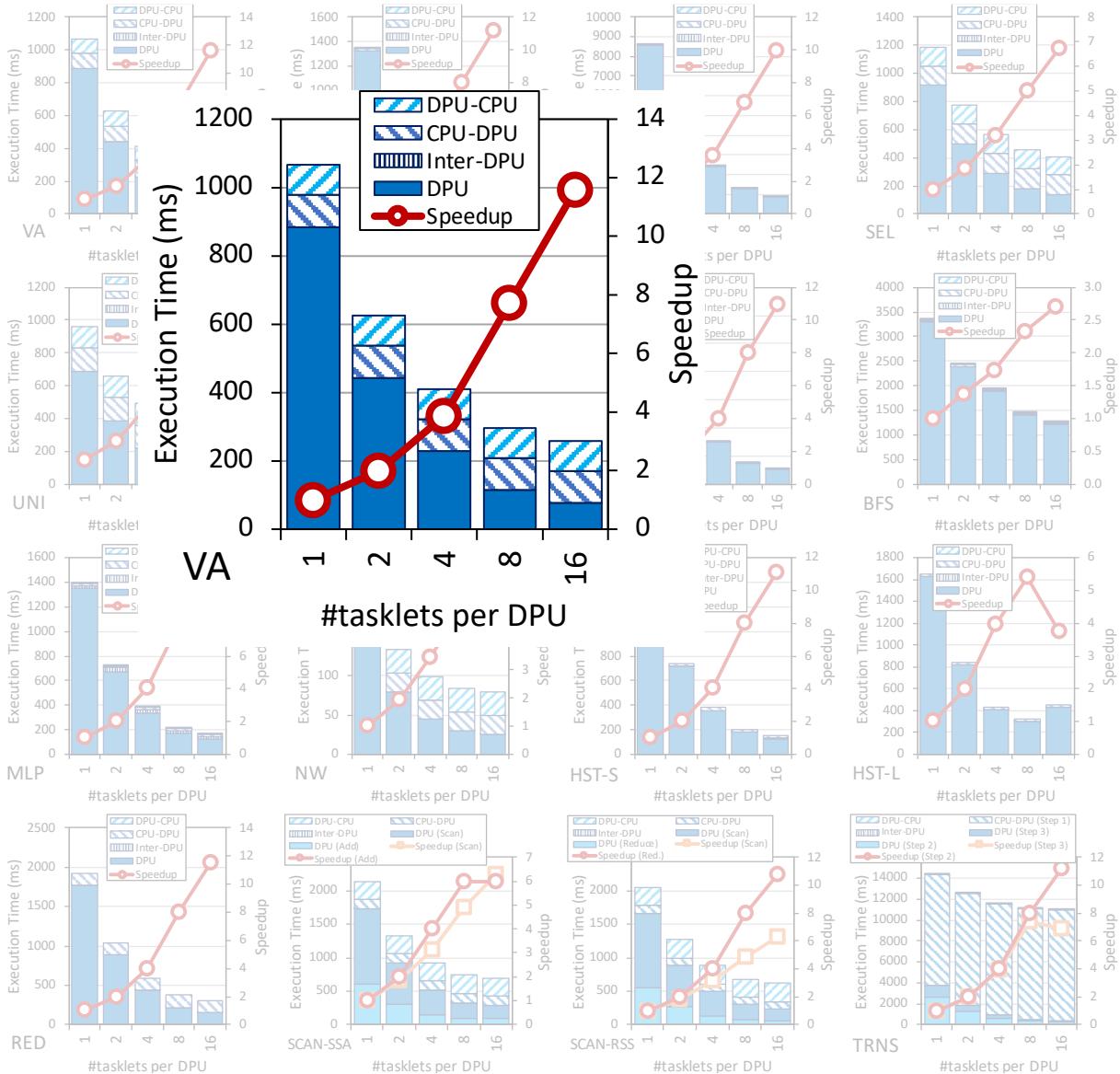
Benchmark	Strong Scaling Dataset	Weak Scaling Dataset	MRAM-WRAM Transfer Sizes
VA	1 DPU-1 rank: 2.5M elem. (10 MB)   32 ranks: 160M elem. (640 MB)	2.5M elem./DPU (10 MB)	1024 bytes
GEMV	1 DPU-1 rank: $8192 \times 1024$ elem. (32 MB)   32 ranks: $163840 \times 4096$ elem. (2.56 GB)	$1024 \times 2048$ 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 = $\frac{2560}{\#DPUs}$ /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: $1536 \times 1024$ input image [256] (6 MB)   32 ranks: $64 \times$ input image	$1536 \times 1024$ input image [256]/DPU (6 MB)	1024 bytes
HST-L	1 DPU-1 rank: $1536 \times 1024$ input image [256] (6 MB)   32 ranks: $64 \times$ input image	$1536 \times 1024$ 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: $12288 \times 16 \times 64 \times 8$ (768 MB)   32 ranks: $12288 \times 16 \times 2048 \times 8$ (24 GB)	$12288 \times 16 \times 1 \times 8$ /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>

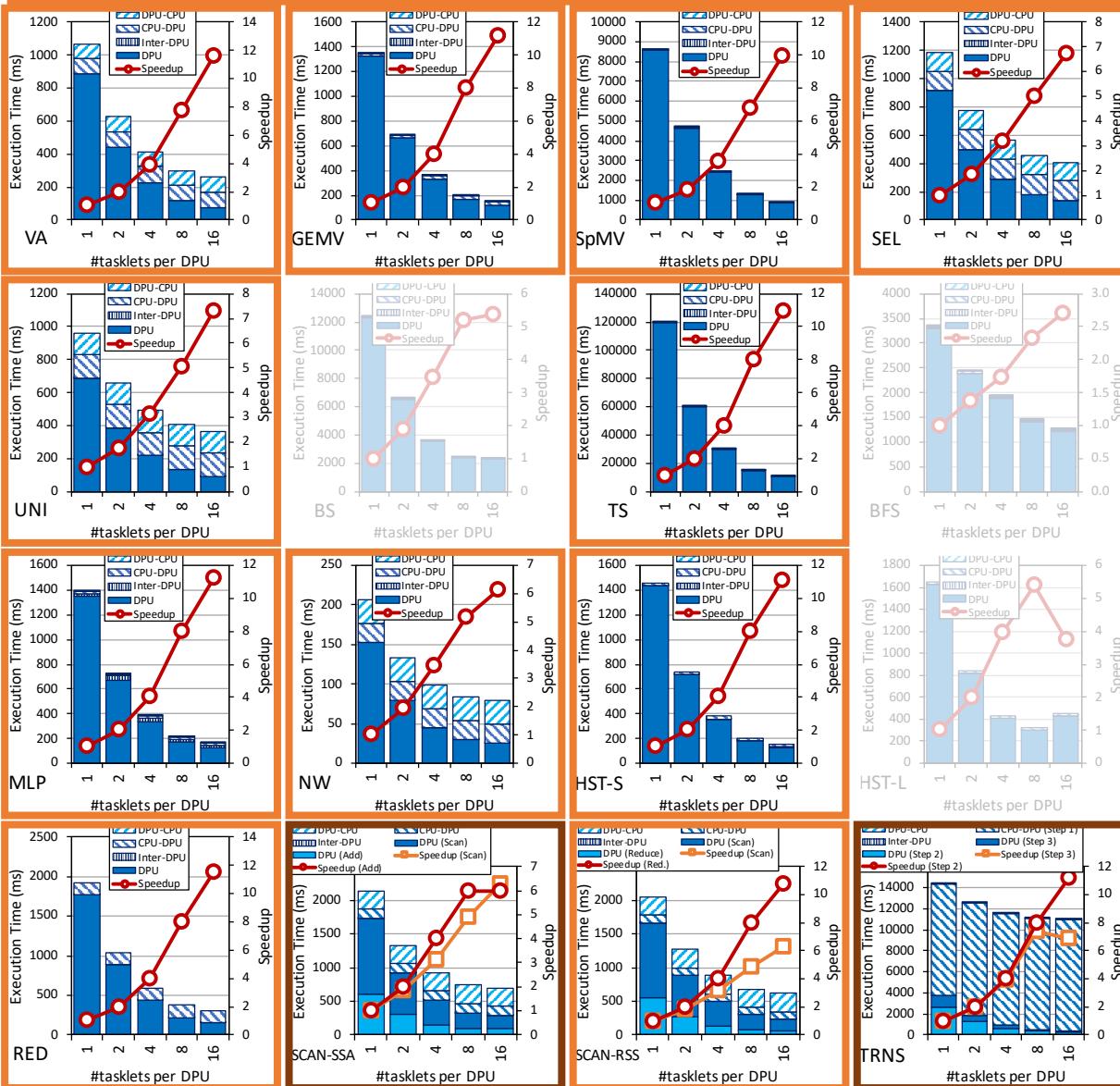
# Strong Scaling: 1 DPU (I)

- 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 kernels from 16 benchmarks), as it fully utilizes the DPU's pipeline.

# Strong Scaling: 1 DPU (III)

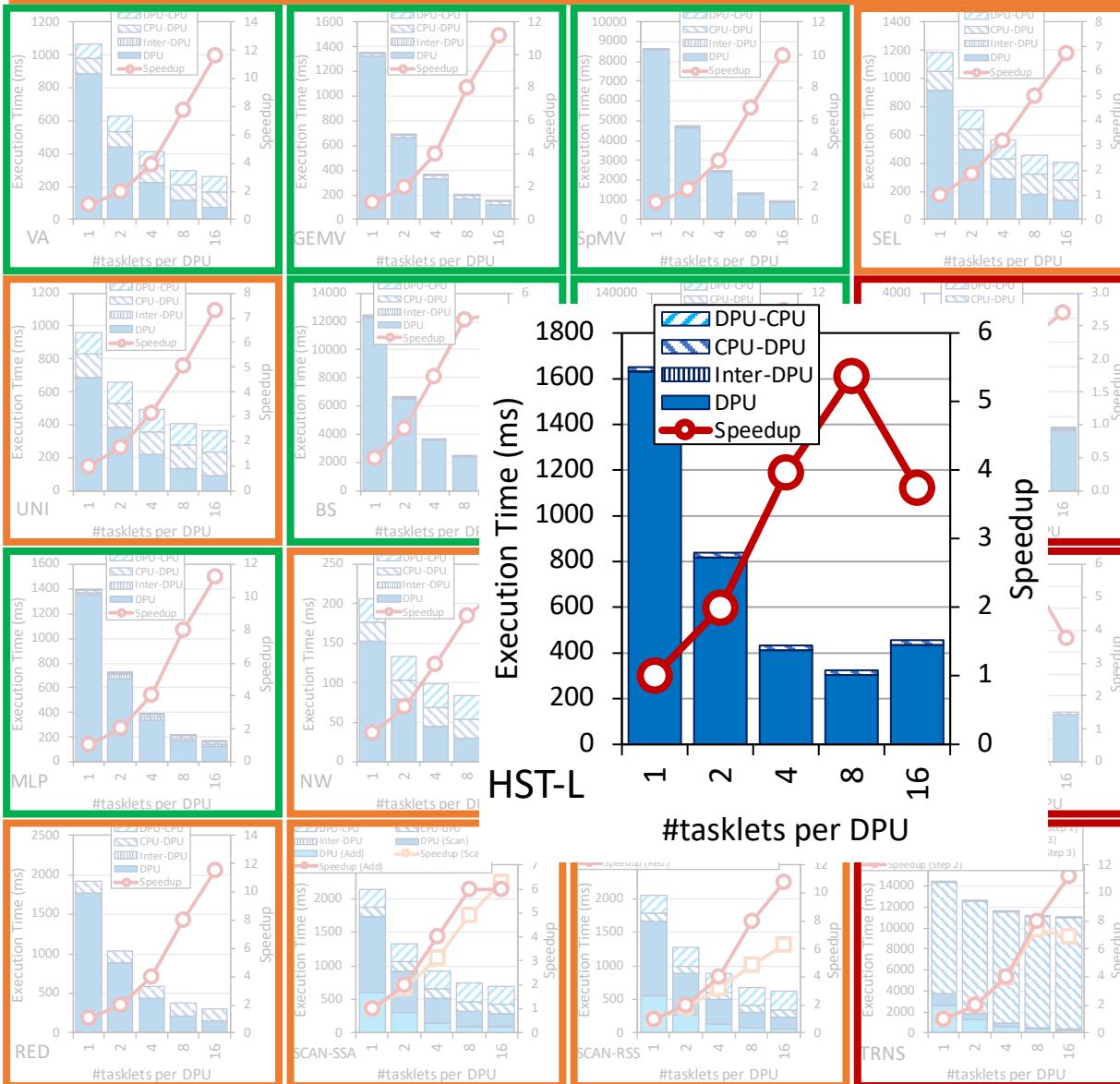


VA, GEMV, SpMV, BS, TS, MLP, HST-S do not use intra-DPU synchronization primitives

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 intra-DPU synchronization primitives

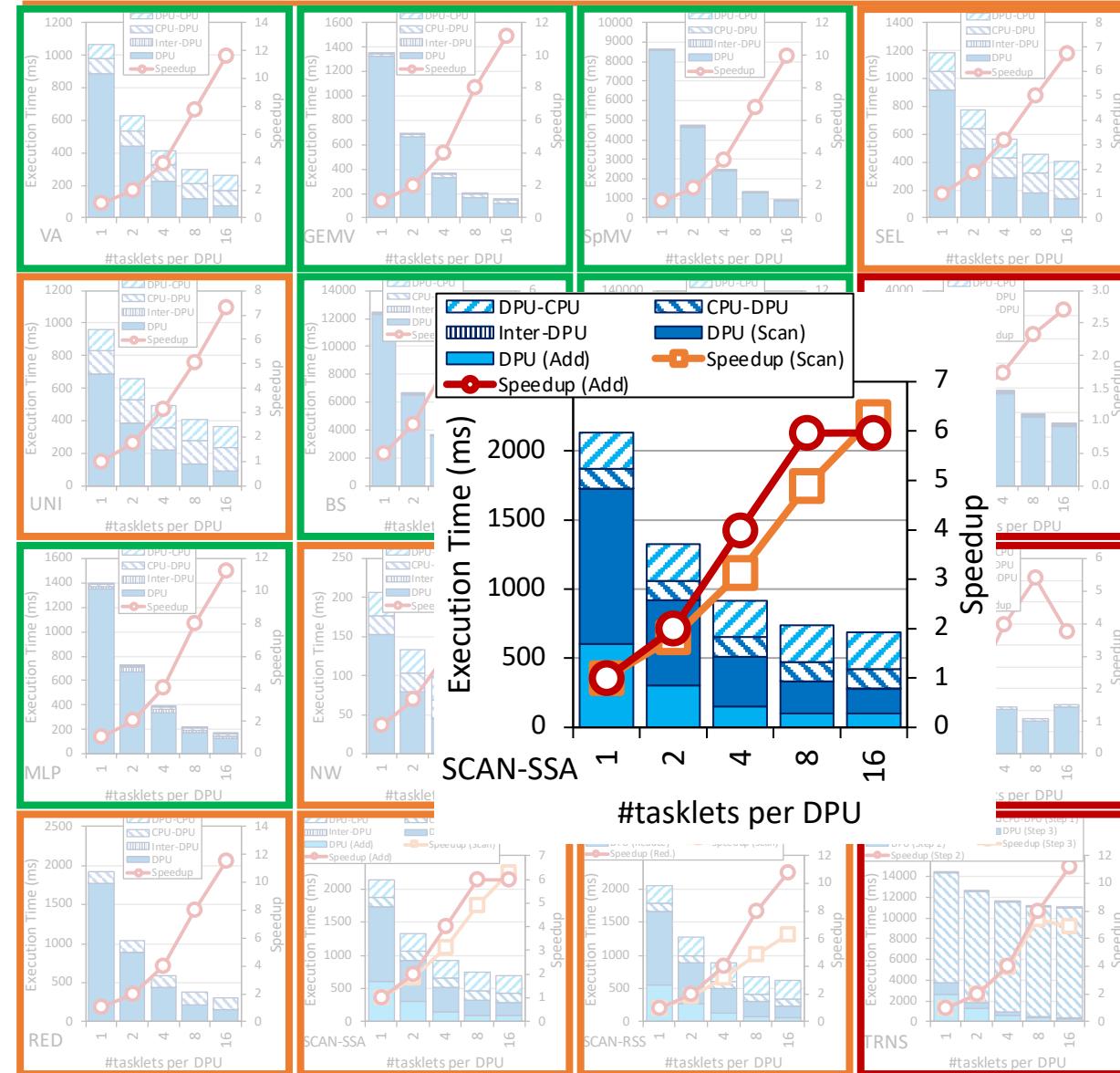
In SEL, UNI, NW, RED, SCAN-SSA (Scan kernel), SCAN-RSS (both kernels), synchronization is lightweight

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

## 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)

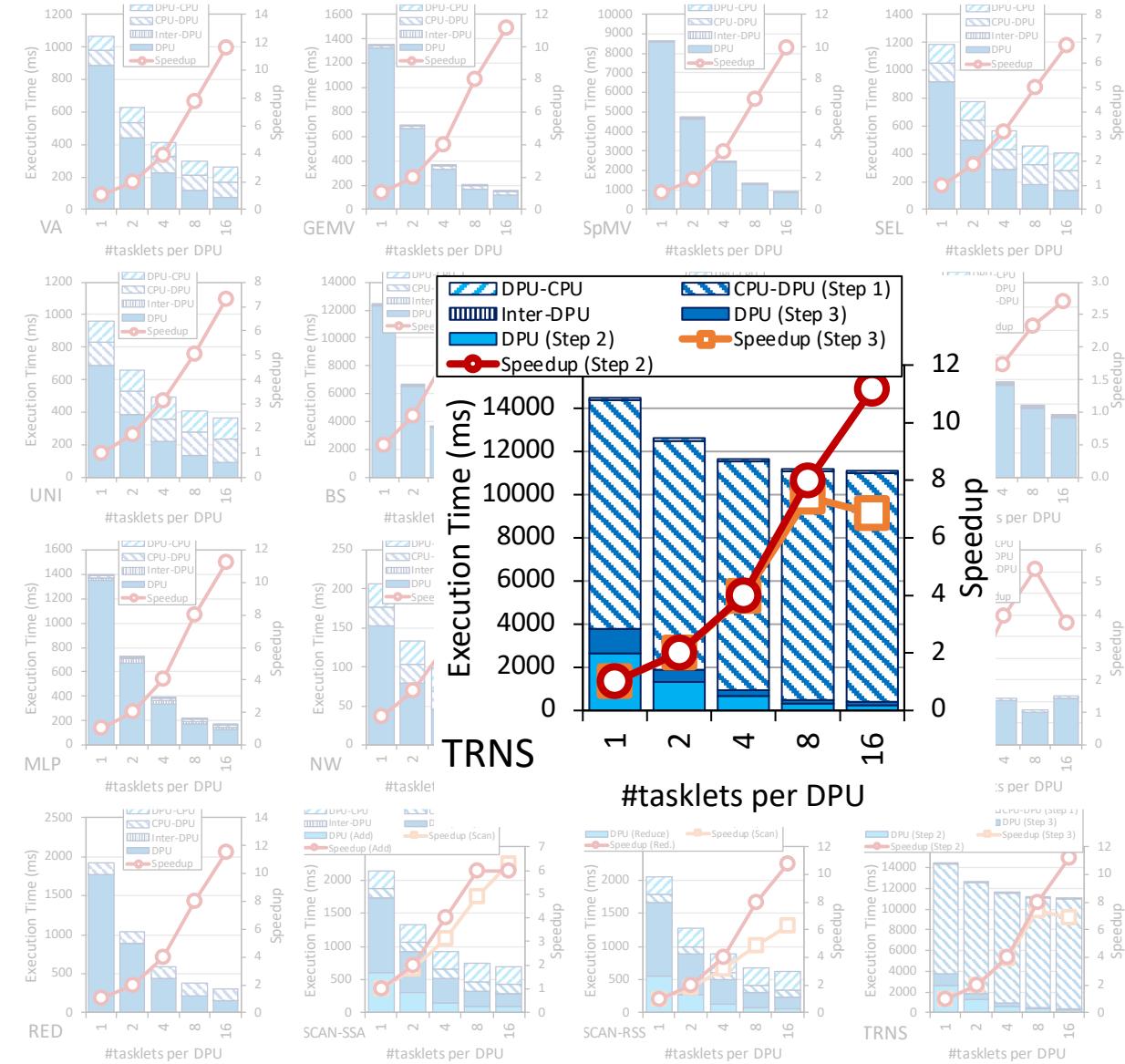


SCAN-SSA (Add kernel) is **not compute-intensive**. Thus, performance saturates with less than 11 tasklets (recall STREAM ADD).  
BS shows similar behavior

## 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

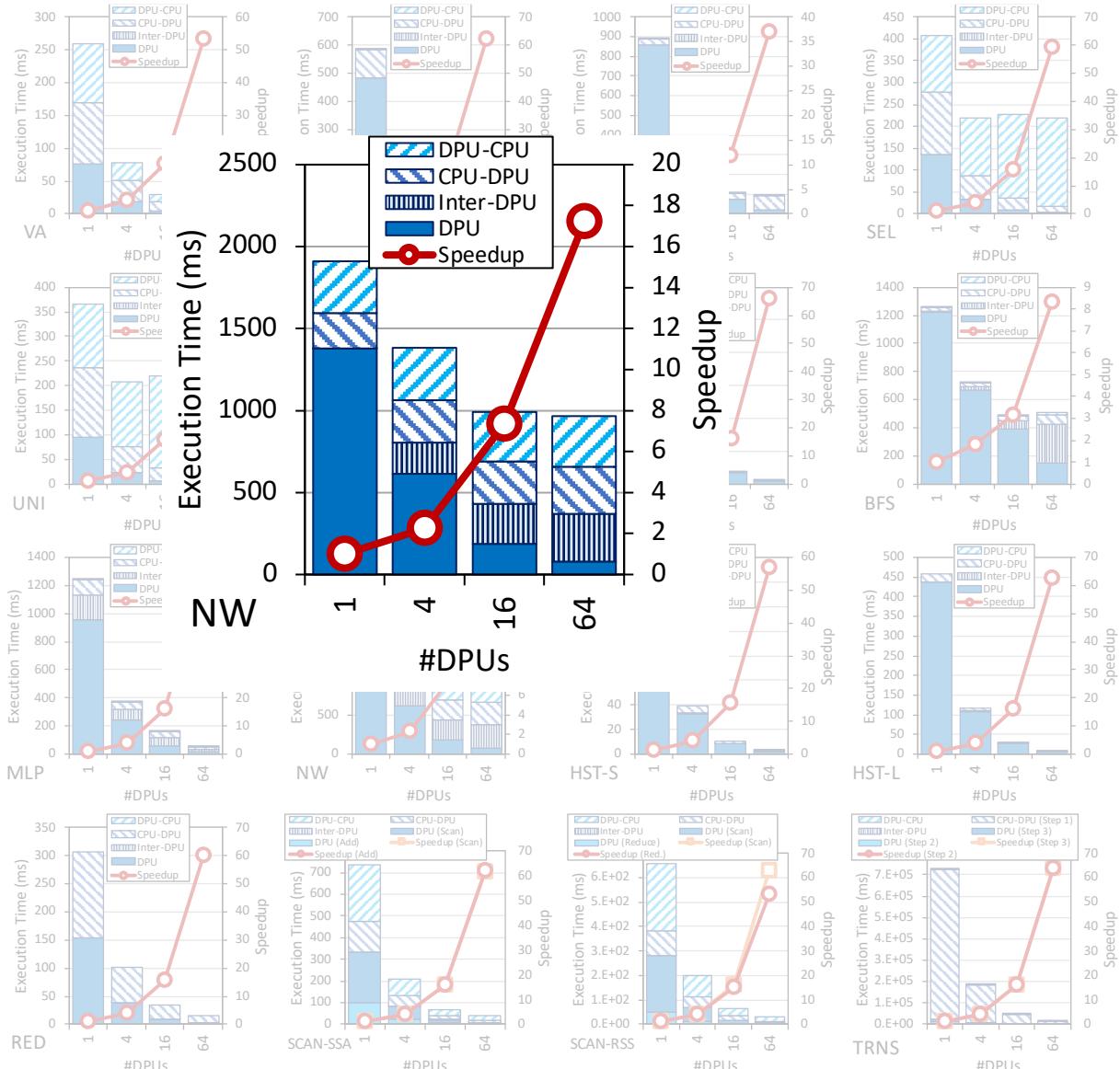
TRNS performs step 1 of the matrix transposition via the CPU-DPU transfer.  
Using small transfers (8 elements) does not exploit full CPU-DPU bandwidth

## KEY OBSERVATION 13

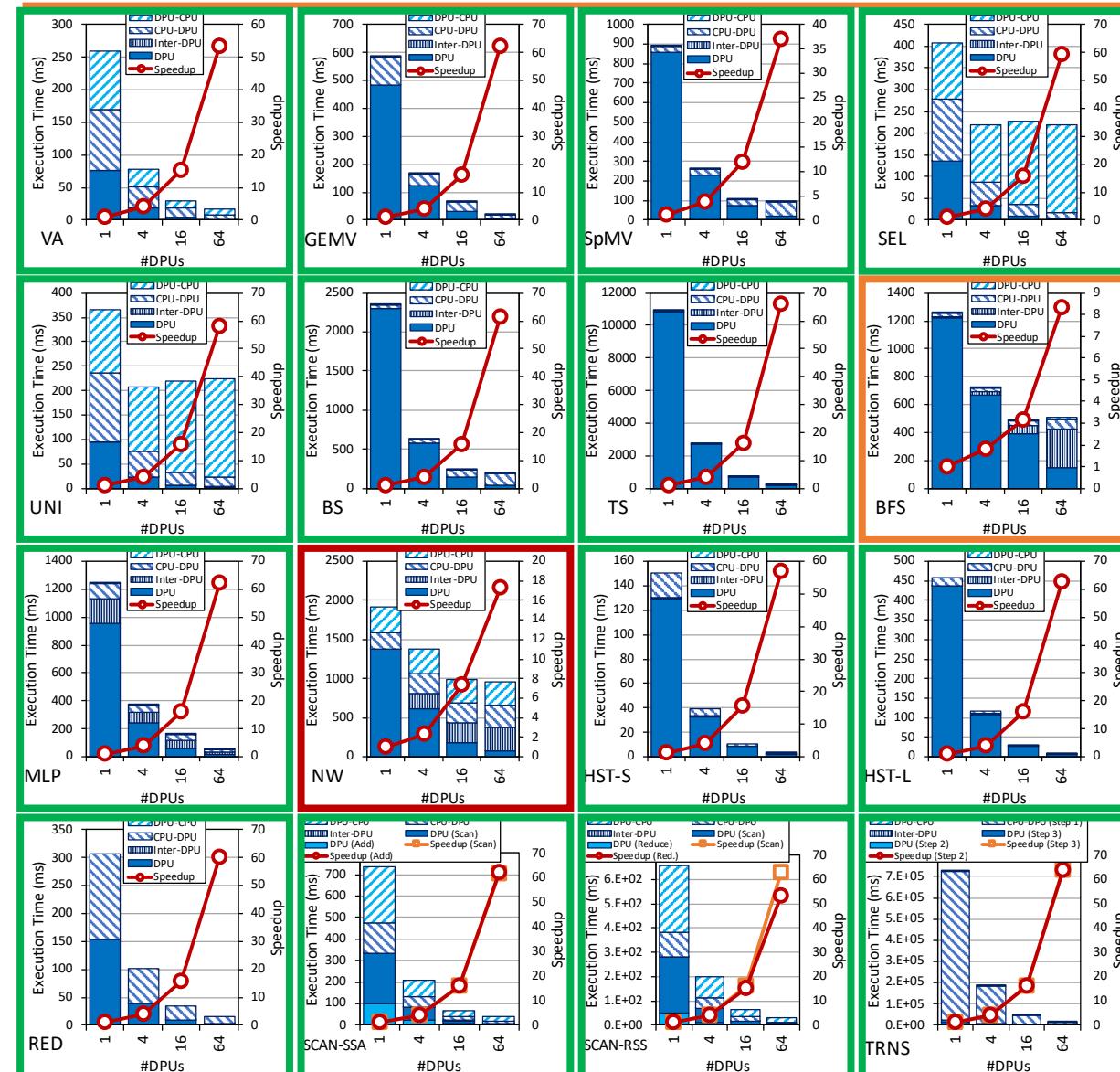
Transferring large data chunks from/to the host CPU is preferred for input data and output results due to higher sustained CPU-DPU/DPU-CPU bandwidths.

# 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

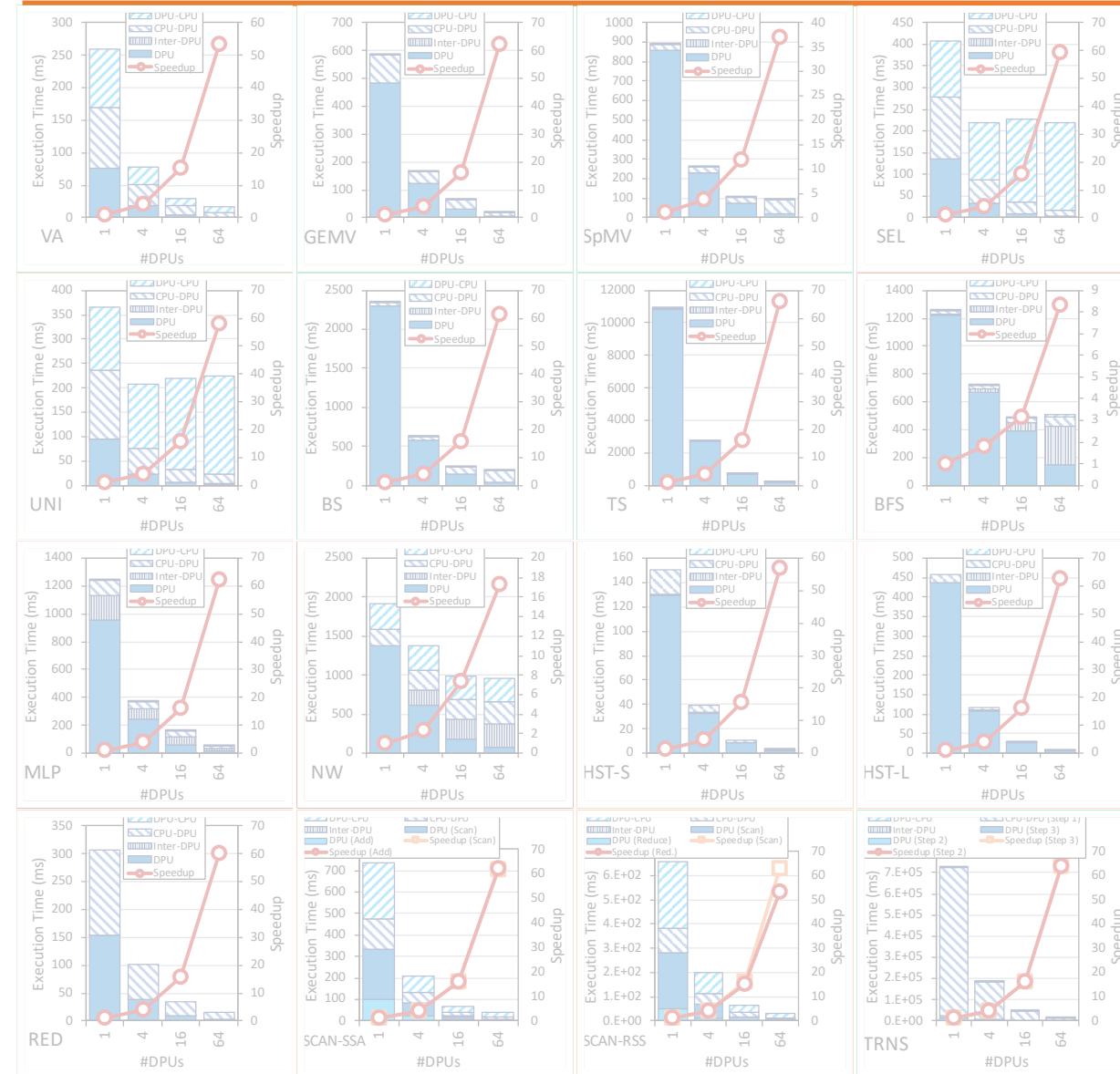
Scaling is sublinear for BFS and NW

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)

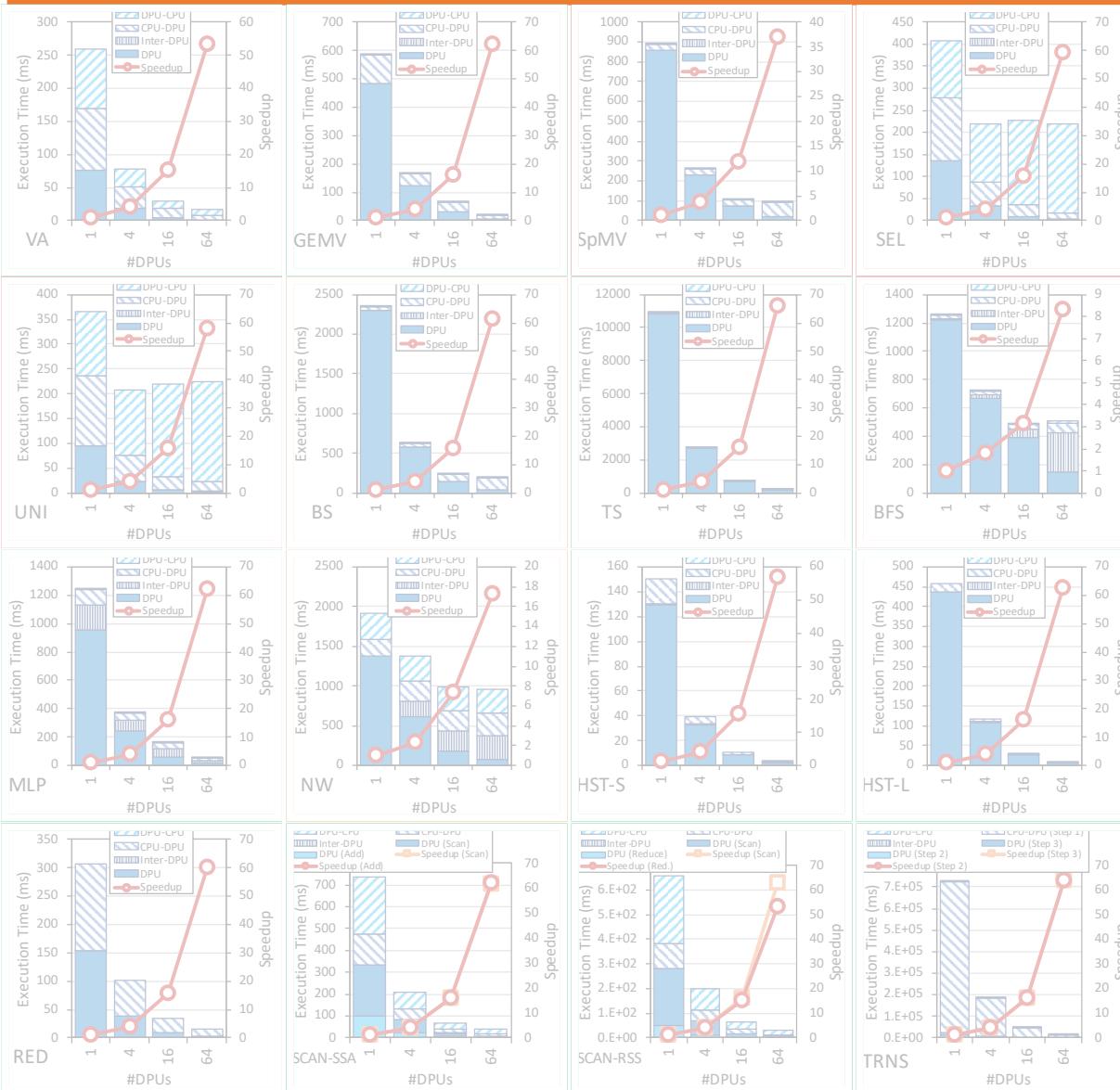


VA, GEMV, SpMV, BS, TS, TRNS do not need inter-DPU synchronization

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

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

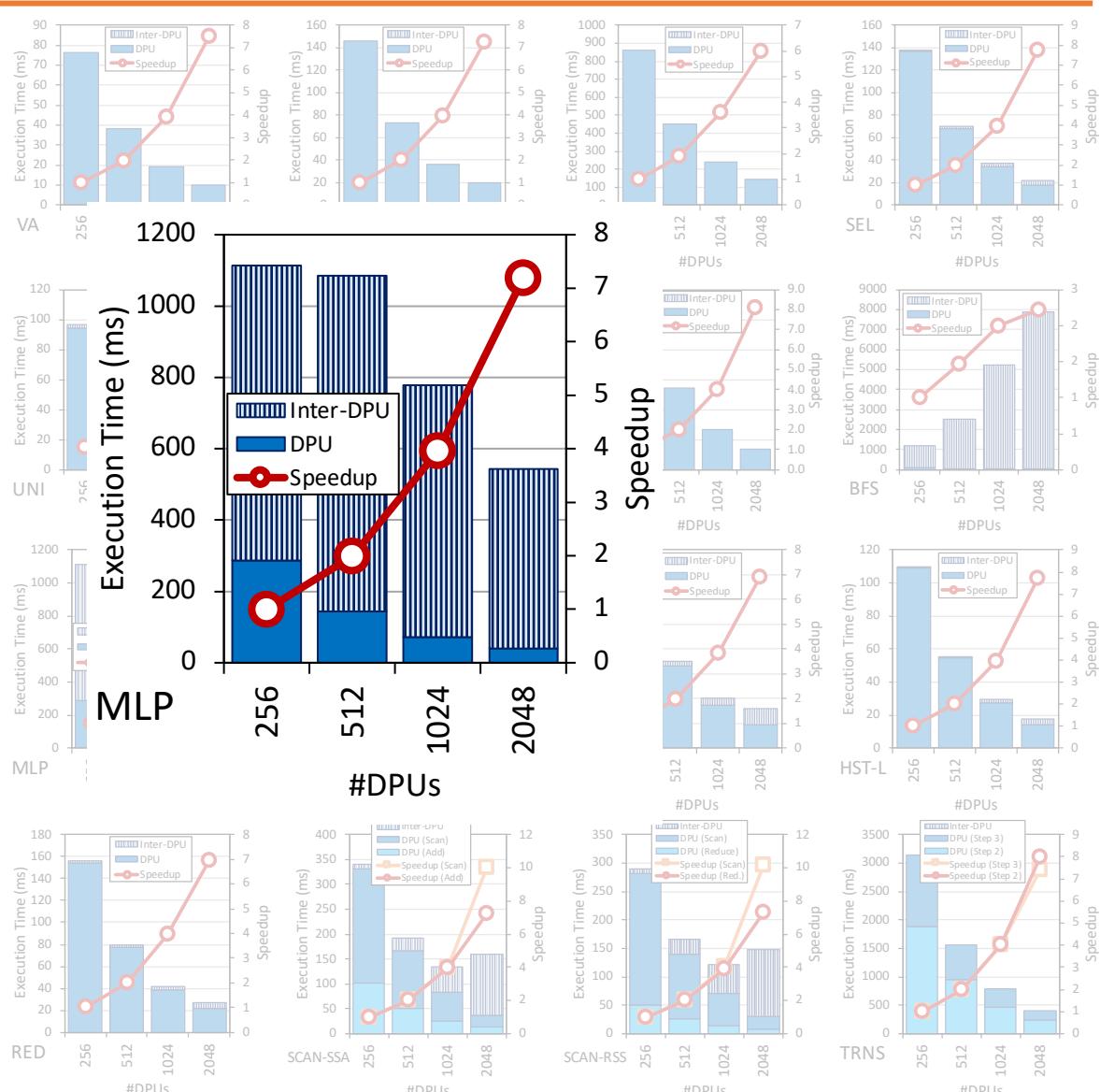
## PROGRAMMING RECOMMENDATION 5

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.

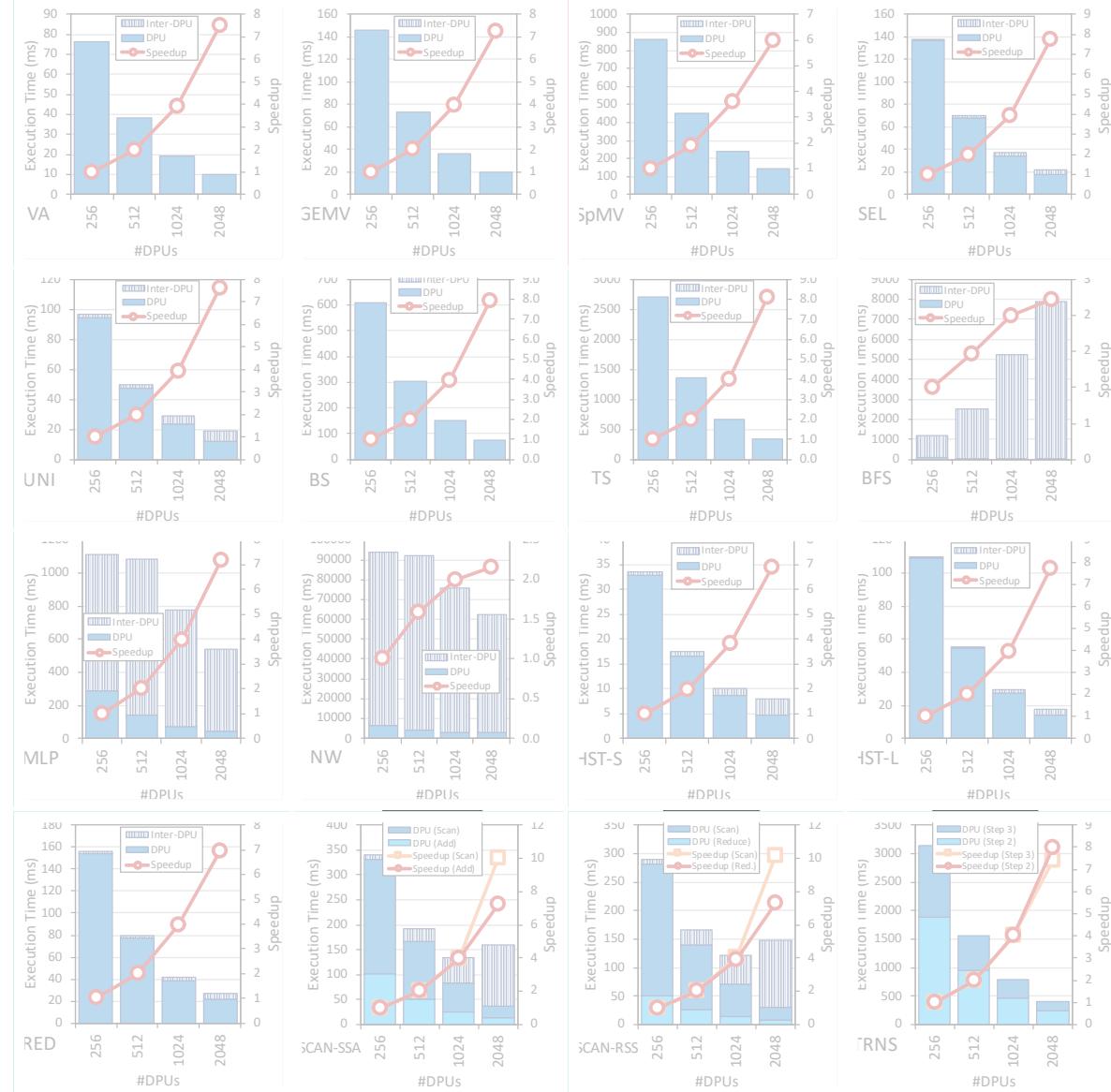
# Strong Scaling: 32 Ranks (I)

- 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)



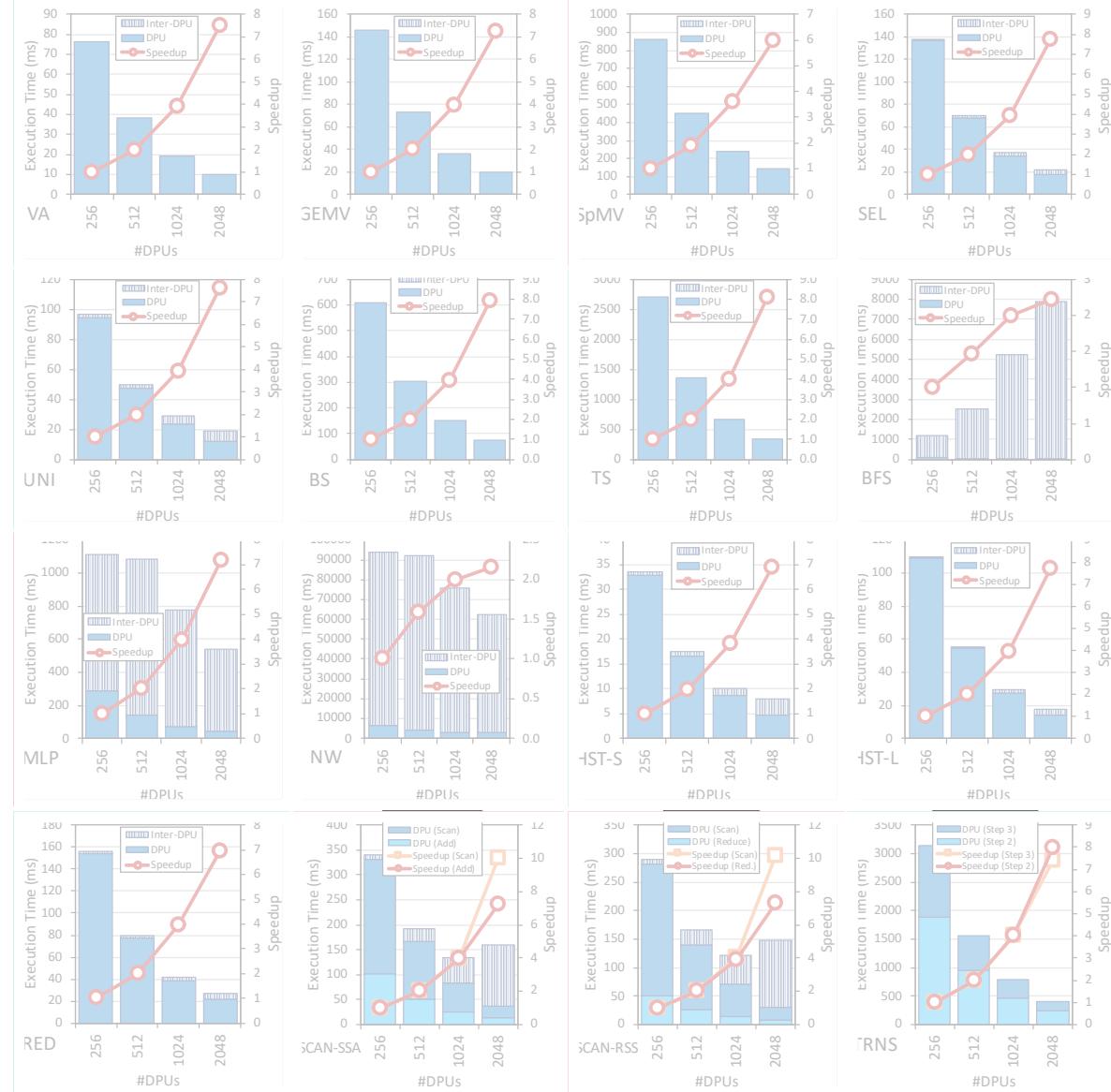
VA, GEMV, SEL, UNI, BS, TS, MLP, HST-S, HSTS-L, RED, SCAN-SSA (both kernel), SCAN-RSS (both kernels), and TRNS (both kernels) **scale linearly with the number of DPUs**

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

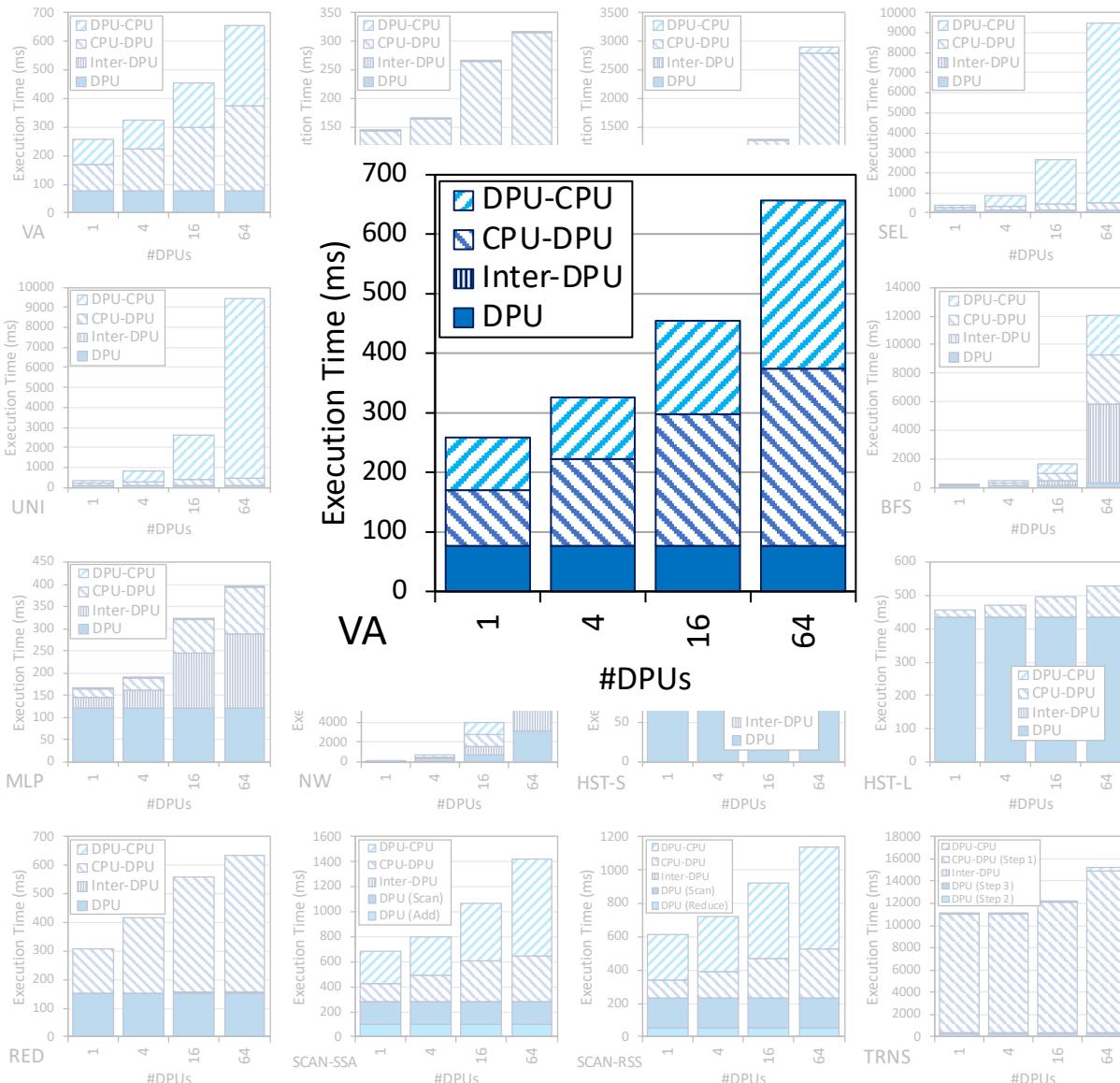
**The overhead of merging partial results from DPUs in the host CPU is tolerable** across all PrIM benchmarks that need it.

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

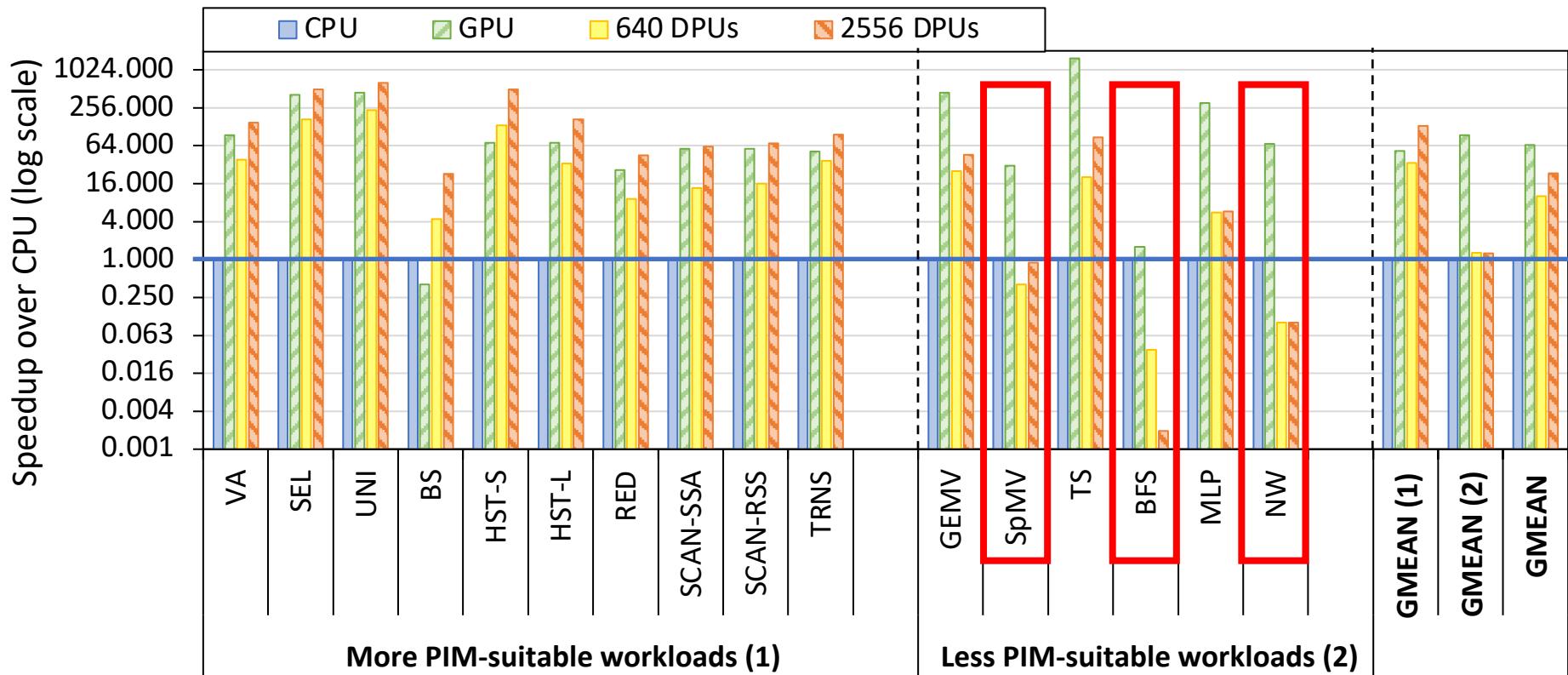
**Sustained bandwidth of parallel CPU-DPU/DPU-CPU transfers inside a rank of DPUs increases sublinearly with the number of DPUs.**

# 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

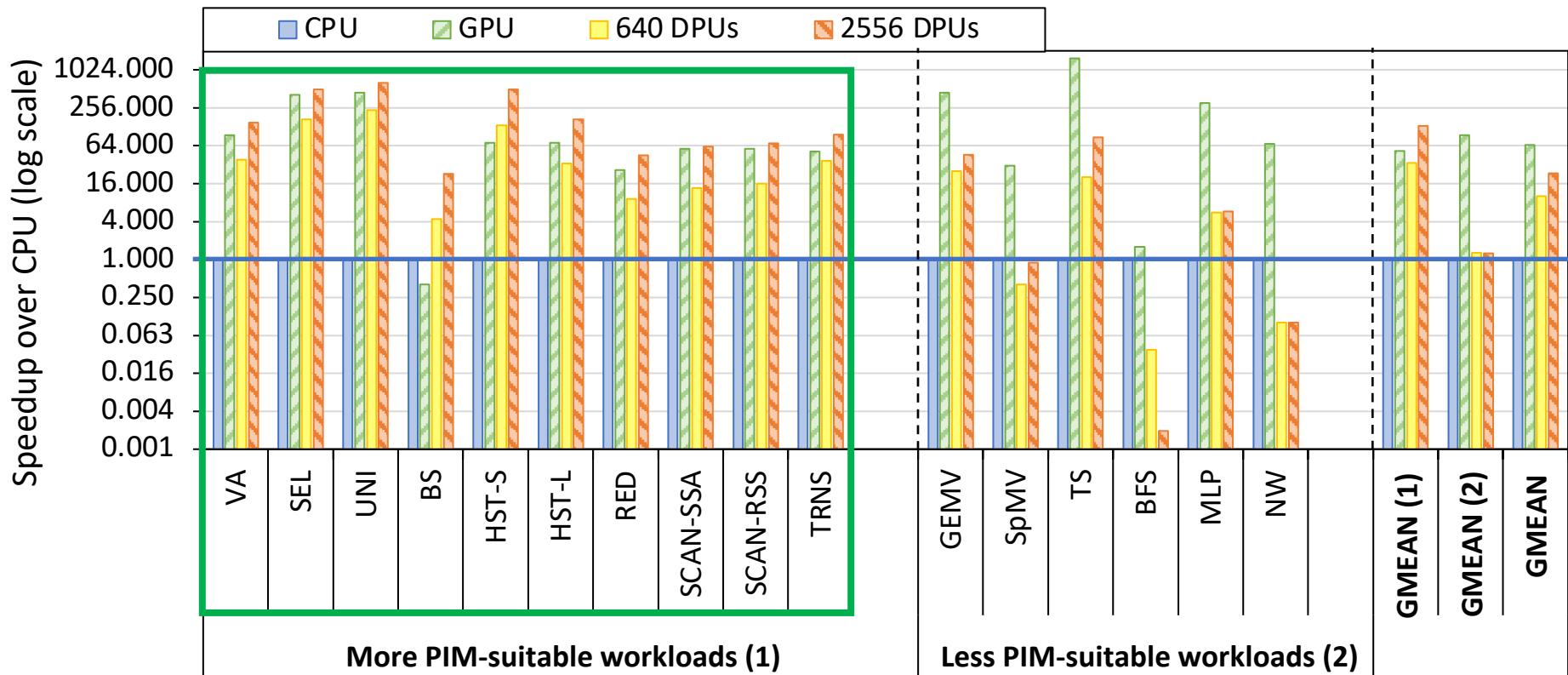
# CPU/GPU: Performance Comparison (I)



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

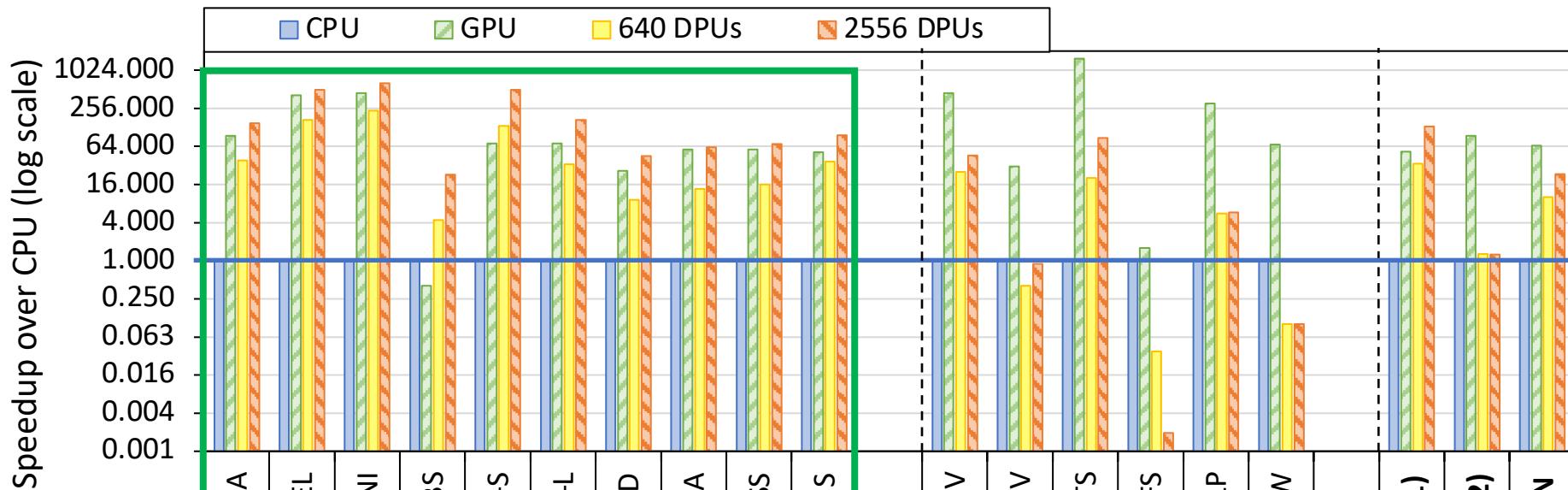
# CPU/GPU: Performance Comparison (II)



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)



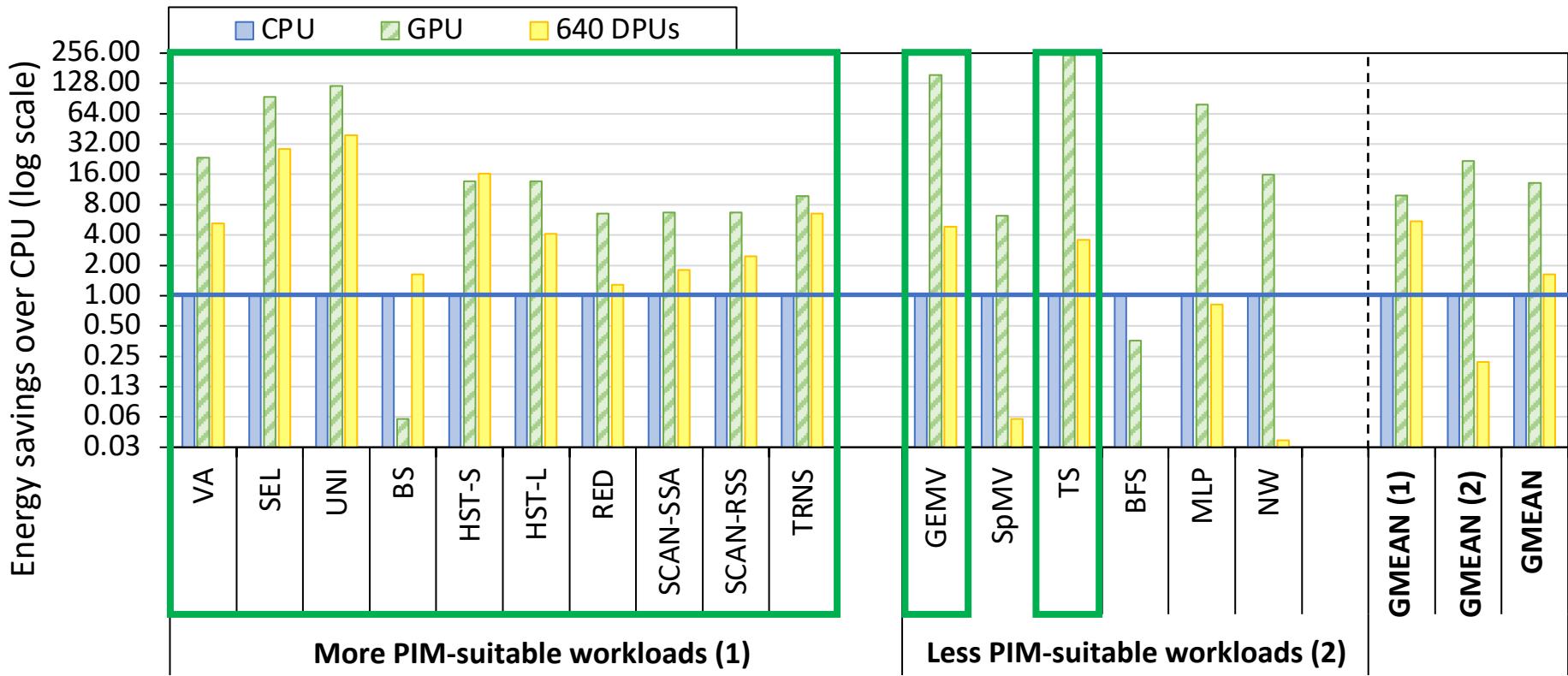
## KEY OBSERVATION 19

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.

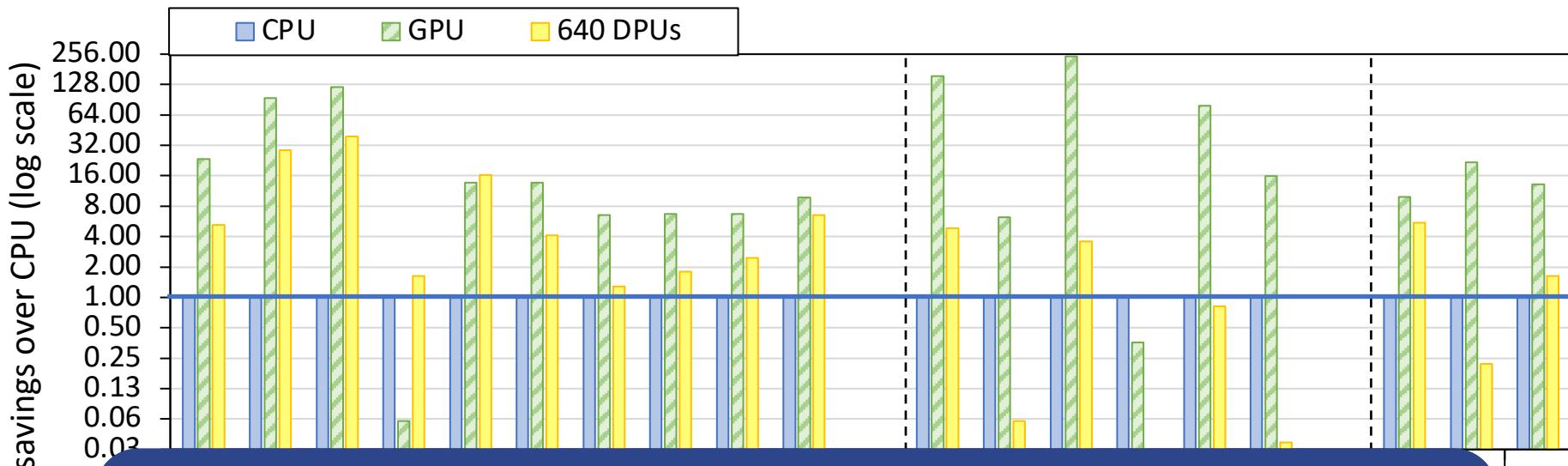
# CPU/GPU: Energy Comparison (I)



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)



## KEY OBSERVATION 20

The UPMEM-based PIM system provides large energy savings over a state-of-the-art CPU due to higher performance (thus, lower static energy) and less data movement between memory and processors.

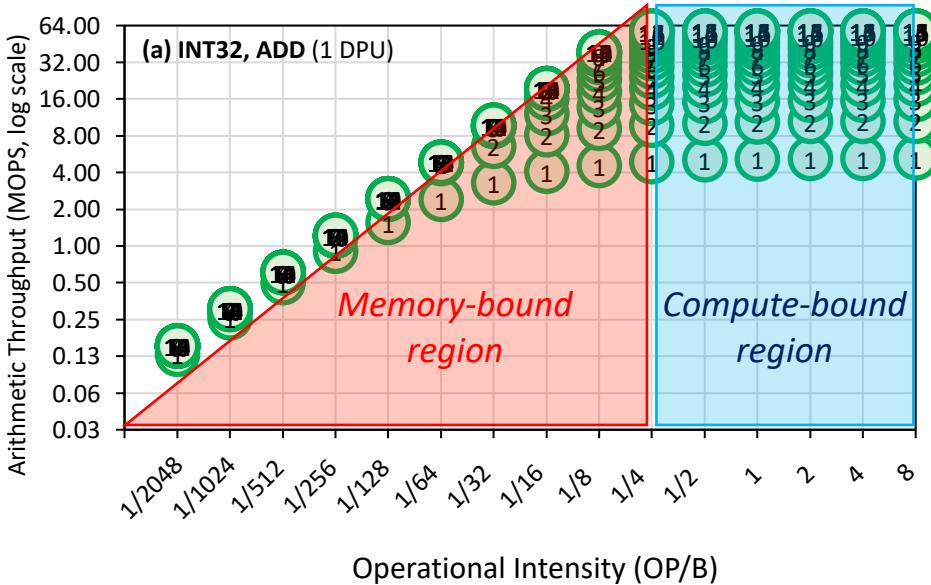
The UPMEM-based PIM system provides energy savings over a state-of-the-art CPU/GPU on workloads where it outperforms the CPU/GPU. This is because the source of both performance improvement and energy savings is the same: **the significant reduction in data movement between the memory and the processor cores**, which the UPMEM-based PIM system can provide for PIM-suitable workloads.

GMEAN

# 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

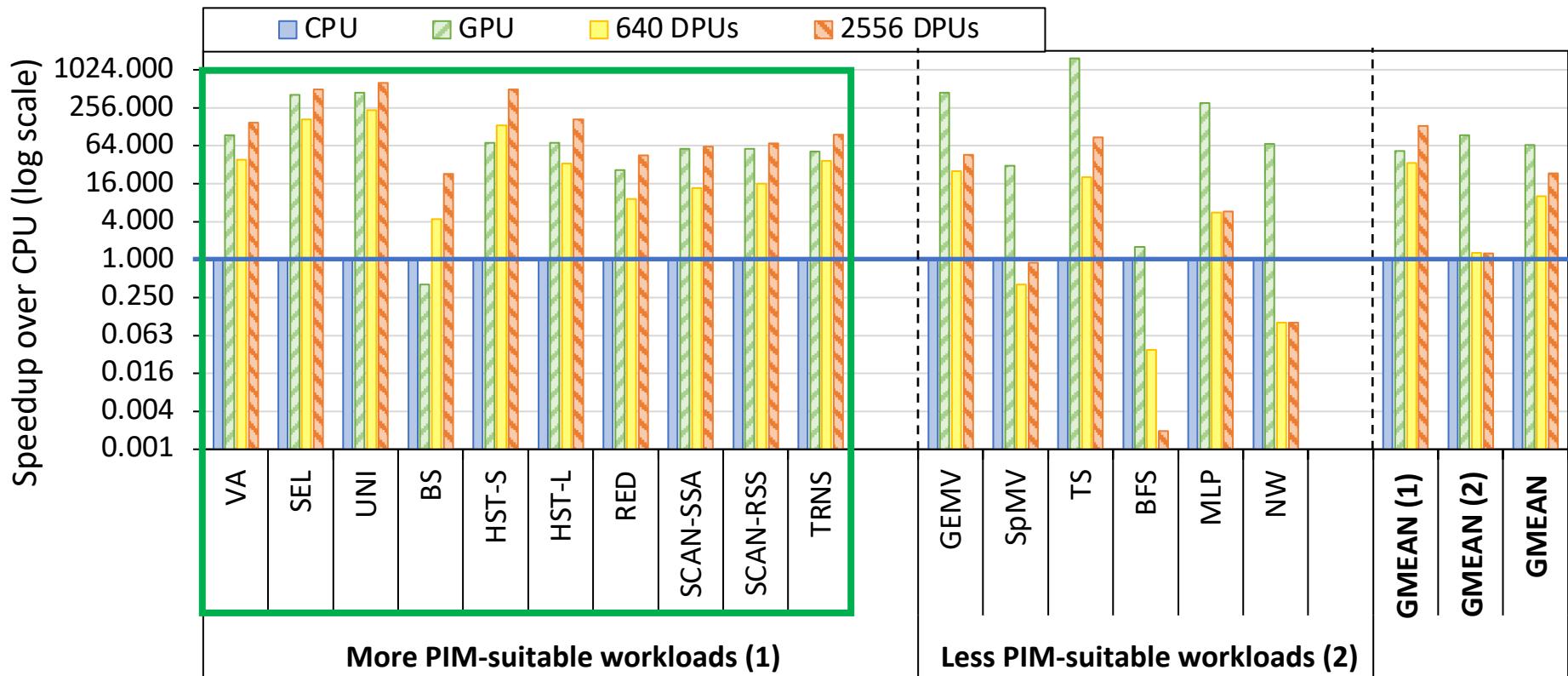


The throughput saturation point is as low as  $\frac{1}{4}$  OP/B,  
i.e., 1 integer addition per every 32-bit element fetched

## KEY TAKEAWAY 1

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

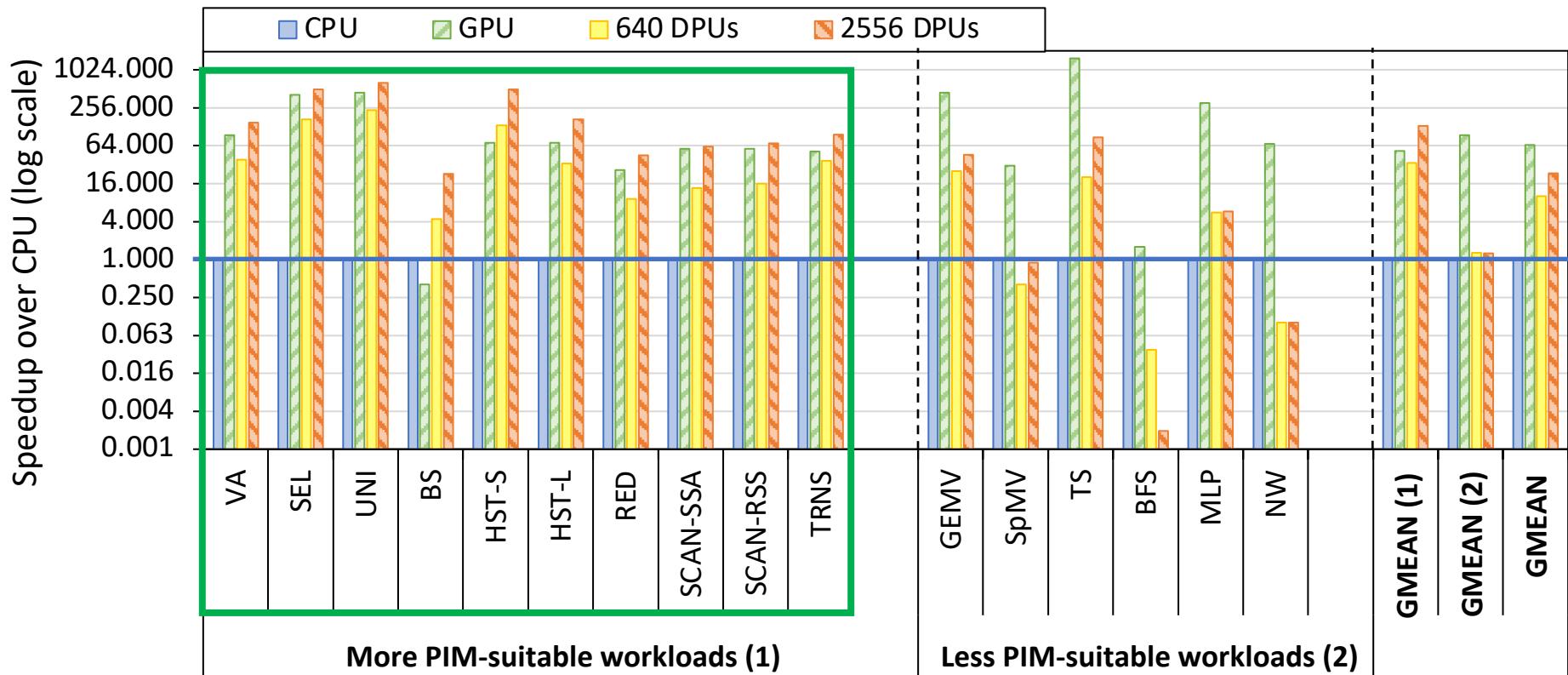
# Key Takeaway 2



## 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** (by  $23.2\times$  on 2,556 DPs for 16 PrIM benchmarks) **and energy efficiency on most of PrIM benchmarks**.
- UPMEM-based PIM systems **outperform state-of-the-art GPUs on a majority of PrIM benchmarks** (by  $2.54\times$  on 2,556 DPs for 10 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.

# Understanding a Modern PIM Architecture

---

## Benchmarking a New Paradigm: Experimental Analysis and Characterization of a Real Processing-in-Memory System

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>, AND 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

Corresponding author: Juan Gómez-Luna (e-mail: juang@ethz.ch).

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

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

# Short arXiv Version

---

## Benchmarking Memory-Centric Computing Systems: Analysis of Real Processing-in-Memory Hardware

Juan Gómez-Luna

*ETH Zürich*

Izzat El Hajj

*American University  
of Beirut*

Ivan Fernandez

*University  
of Malaga*

Christina Giannoula

*National Technical  
University of Athens*

Geraldo F. Oliveira

*ETH Zürich*

Onur Mutlu

*ETH Zürich*

<https://arxiv.org/pdf/2110.01709.pdf>

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

# Long arXiv Version

---

## Benchmarking a New Paradigm: An Experimental Analysis of a Real Processing-in-Memory Architecture

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>

# 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'. The page includes the repository name at the top, a header with 'Code' being the active tab, and various repository statistics like 2 stars, 1 fork, and 1 issue. Below the header is a navigation bar with links for 'Issues', 'Pull requests', 'Actions', 'Projects', 'Wiki', 'Security', 'Insights', and 'Settings'. The main content area displays the 'README.md' file. It shows a commit from 'Juan Gomez Luna' with the message 'PrIM -- first commit' and a timestamp of 'Latest commit 3de4b49 9 days ago'. There is also a link to 'History'. Below the commit is information about the file: '168 lines (132 sloc) | 5.79 KB' with options to 'Raw', 'Blame', 'Copy', 'Edit', and 'Delete'. The content of the README.md file is titled 'PrIM (Processing-In-Memory Benchmarks)' and describes the repository as the first benchmark suite for a real-world processing-in-memory (PIM) architecture. It highlights the UPMEM PIM architecture and its combination of DRAM memory arrays and general-purpose cores. The text also mentions that PrIM provides a common set of workloads for researchers and can be used for programming, architecture, and system research. A note at the bottom states that PrIM includes microbenchmarks for assessing architecture limits.

CMU-SAFARI / prim-benchmarks

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 Copy Edit Delete

## 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.

# Computer Architecture

## Lecture 4a: Programming a Real-world PIM Architecture

Prof. Onur Mutlu

ETH Zürich

Fall 2023

6 October 2023

We Covered Until This Point  
in the Lecture