

Computer Architecture

Lecture 3: Processing near Memory

Prof. Onur Mutlu
ETH Zürich
Fall 2023
5 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

Three Key Systems Trends

- 1. Data access is a major bottleneck**
 - ❑ Applications are increasingly data hungry
- 2. Energy consumption is a key limiter**
- 3. Data movement energy dominates compute**
 - ❑ Especially true for off-chip to on-chip movement

Observation and Opportunity

- High latency and high energy caused by data movement
 - Long, energy-hungry interconnects
 - Energy-hungry electrical interfaces
 - Movement of large amounts of data
- Opportunity: Minimize data movement by performing computation directly (near) where the data resides
 - Processing in memory (PIM)
 - In-memory computation/processing
 - Near-data processing (NDP)
 - General concept applicable to any data storage & movement unit (caches, SSDs, main memory, network, controllers)

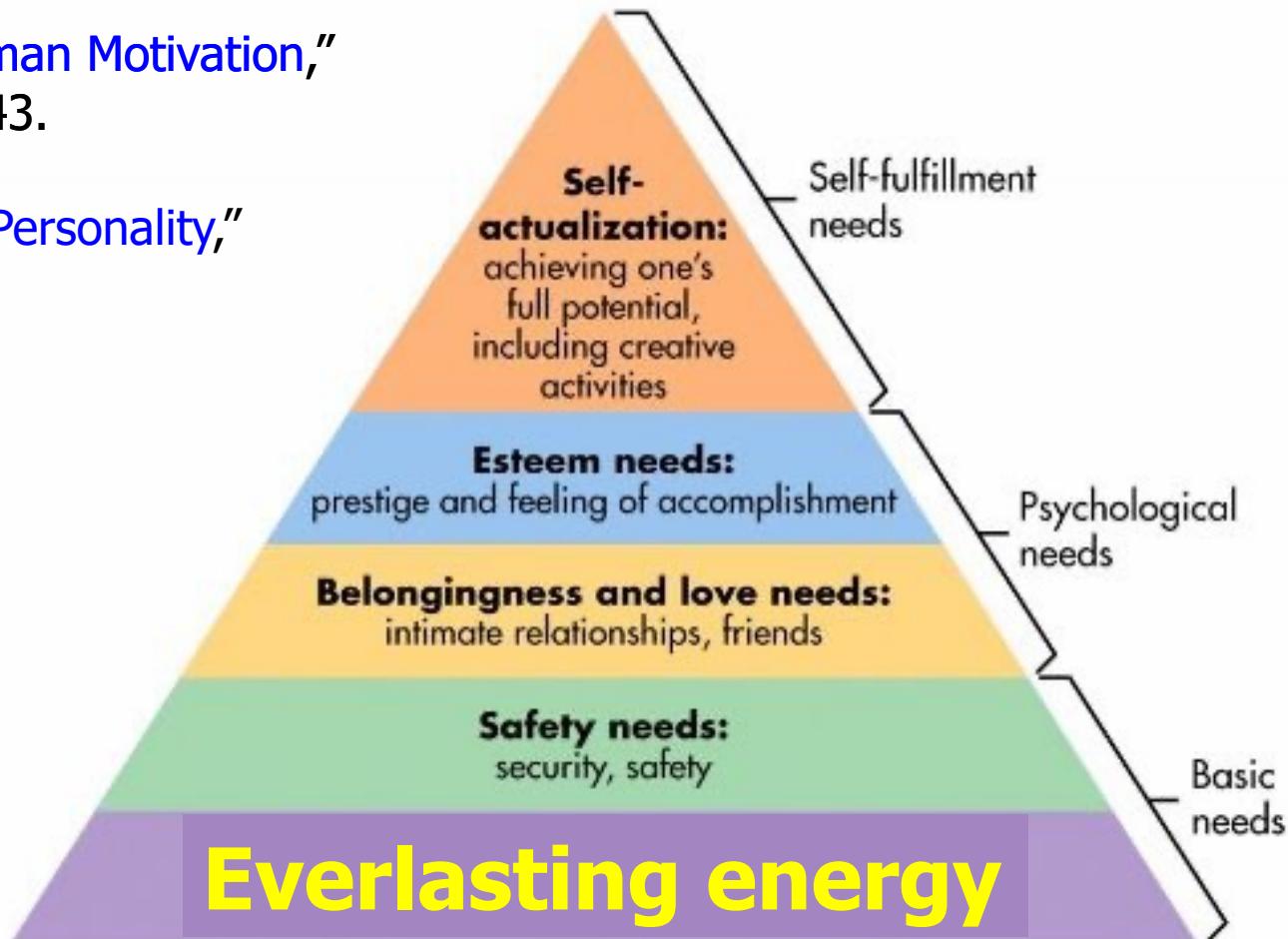
Four Key Issues in Future Platforms

- Fundamentally Secure/Reliable/Safe Architectures
- Fundamentally Energy-Efficient Architectures
 - Memory-centric (Data-centric) Architectures
- Fundamentally Low-Latency Architectures
- Architectures for AI/ML, Genomics, Medicine, Health

Maslow's (Human) Hierarchy of Needs, Revisited

Maslow, "A Theory of Human Motivation,"
Psychological Review, 1943.

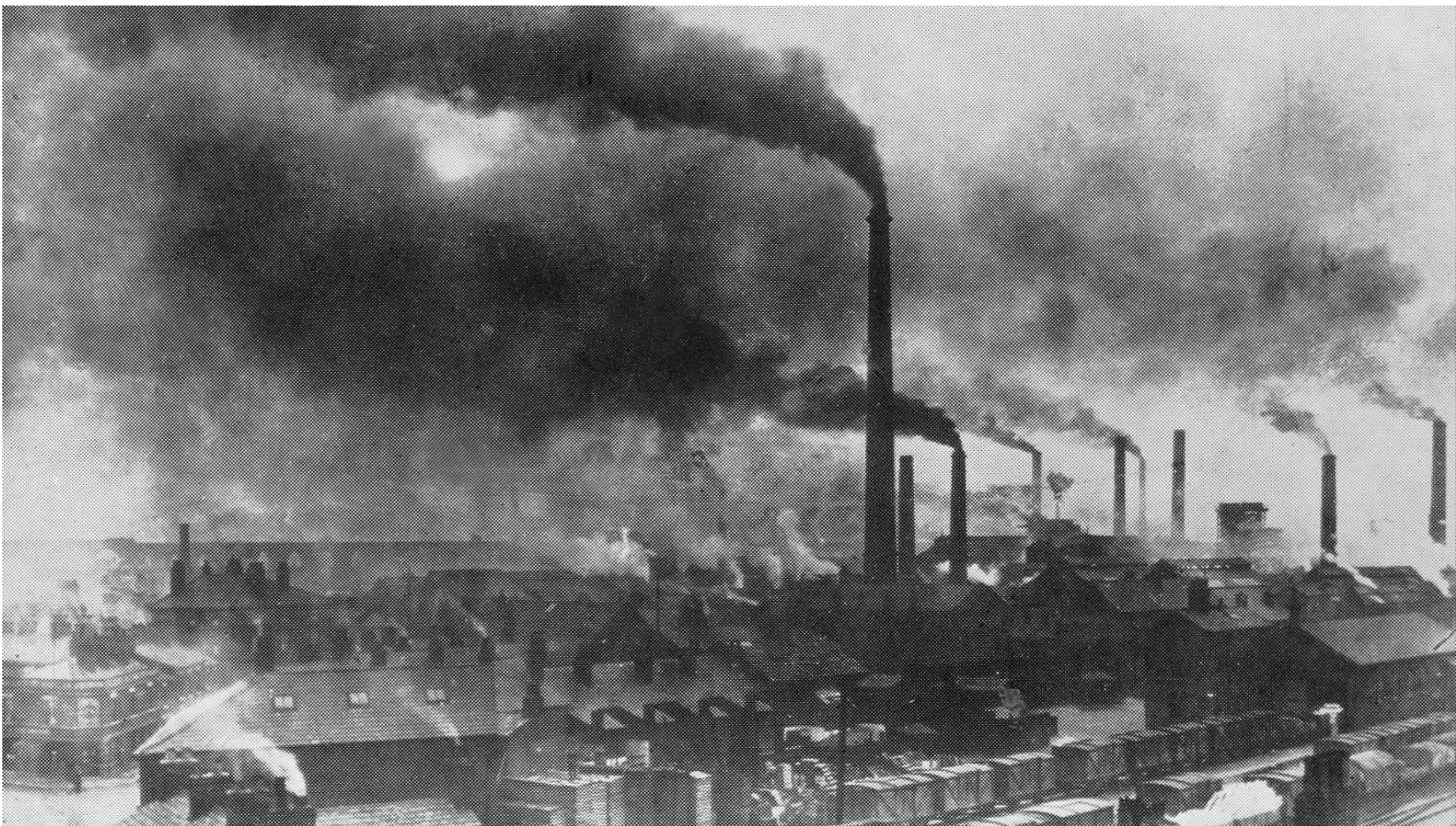
Maslow, "Motivation and Personality,"
Book, 1954-1970.



Do We Want This?



Or This?



Challenge and Opportunity for Future

High Performance,

Energy Efficient,

Sustainable

(All at the Same Time)

The Problem

Data access is the major performance and energy bottleneck

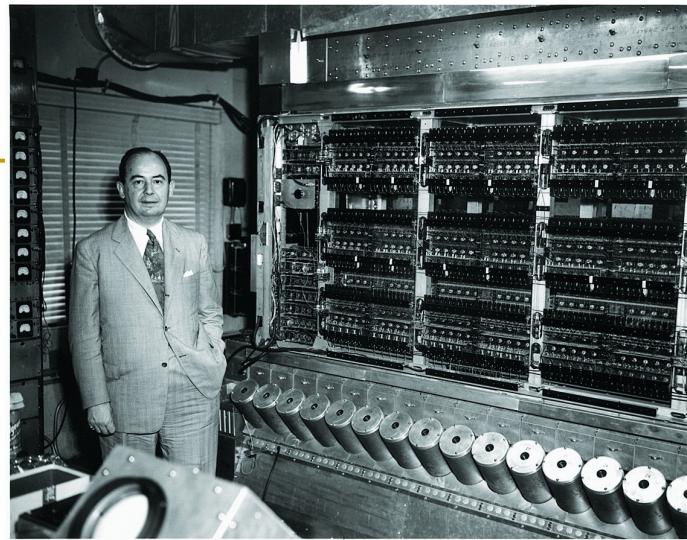
Our current
design principles
cause great energy waste
(and great performance loss)

The Problem

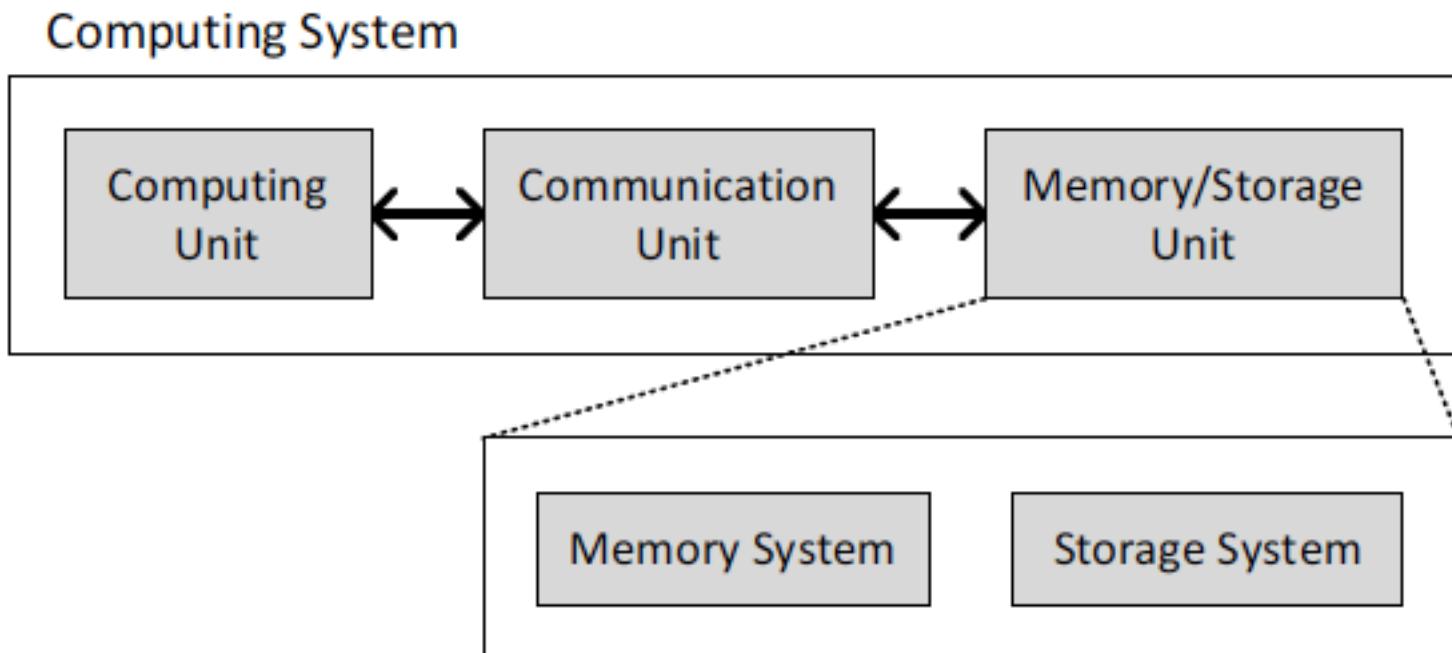
Processing of data
is performed
far away from the data

A Computing System

- Three key components
- Computation
- Communication
- Storage/memory

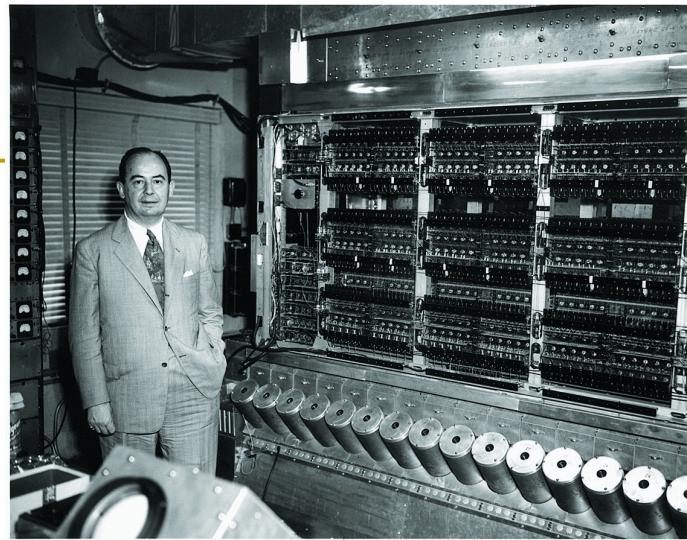


Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.

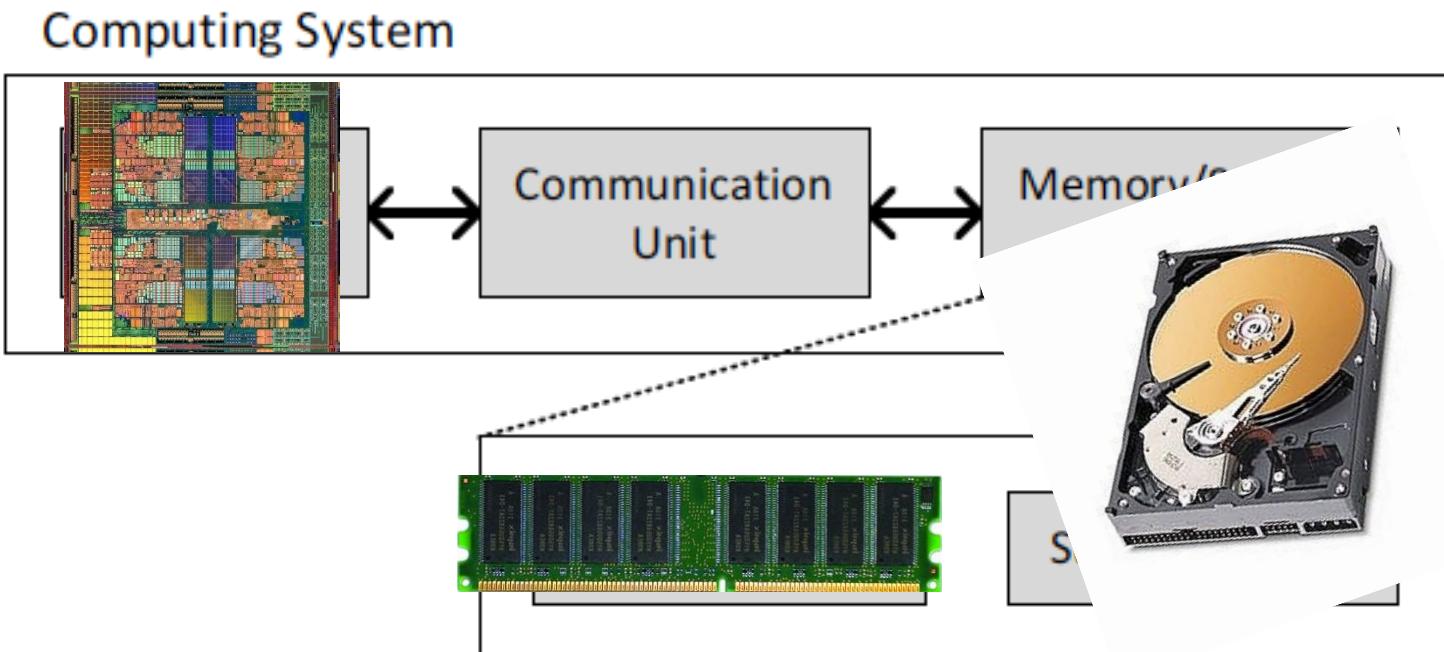


A Computing System

- Three key components
- Computation
- Communication
- Storage/memory

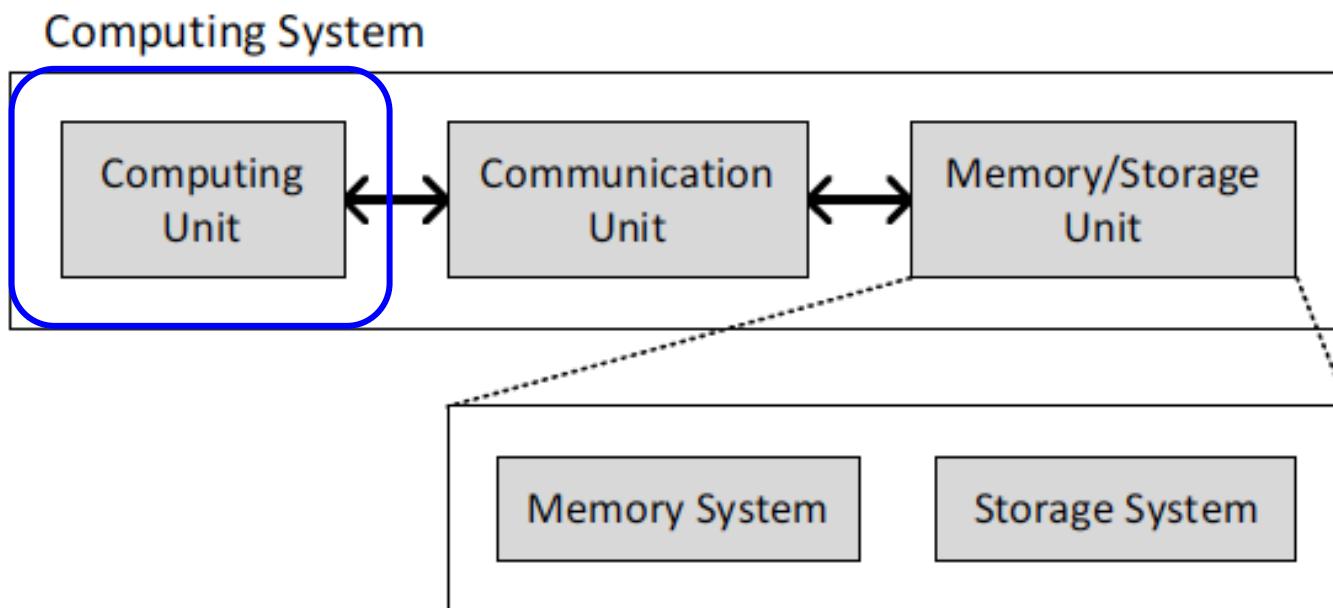


Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.



Today's Computing Systems

- Are overwhelmingly processor centric
- All data processed in the processor → at great system cost
- Processor is heavily optimized and is considered the master
- Data storage units are dumb and are largely unoptimized (except for some that are on the processor die)



It's the Memory, Stupid!

- **"It's the Memory, Stupid!"** (Richard Sites, MPR, 1996)

RICHARD SITES

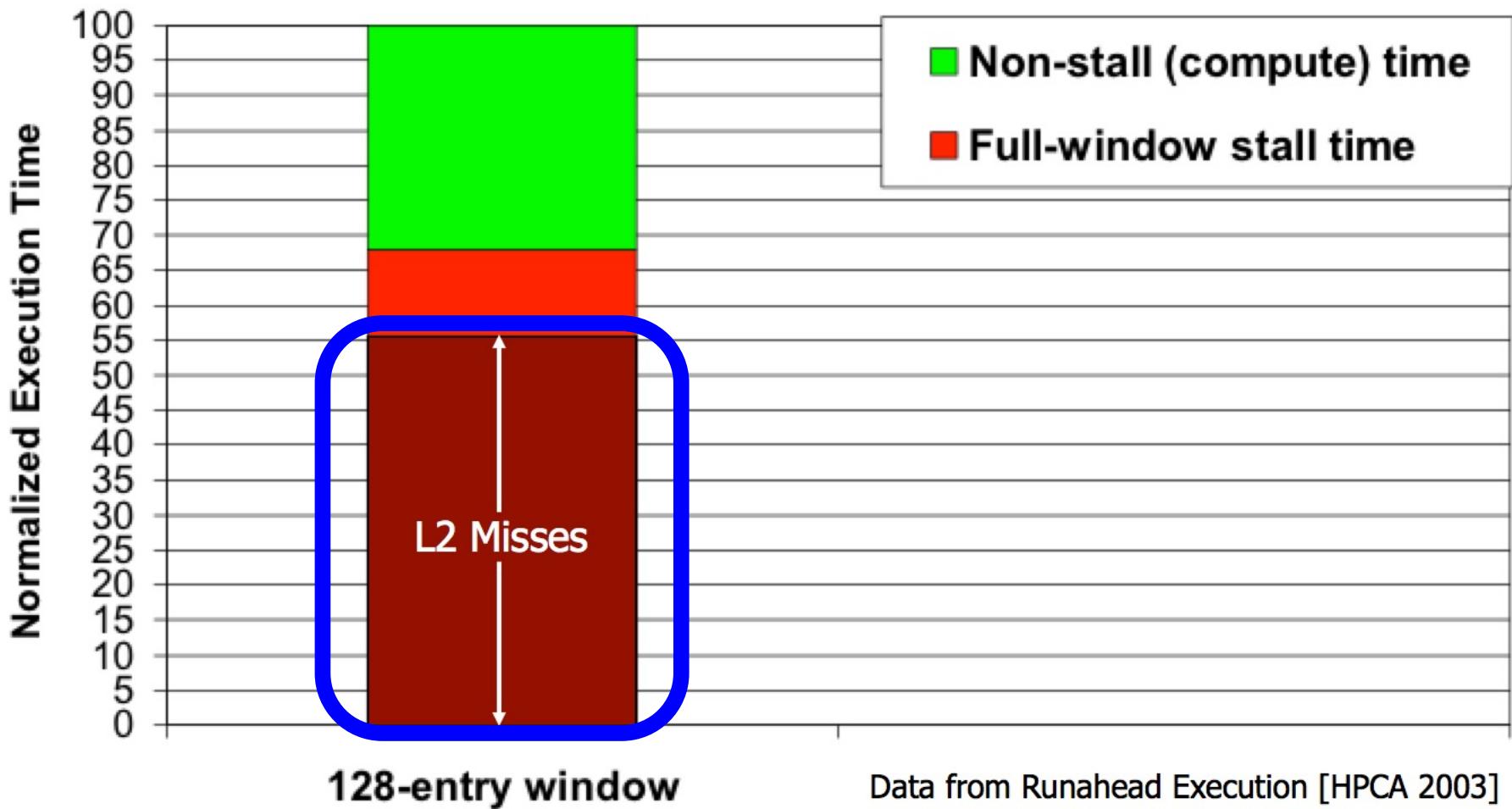
It's the Memory, Stupid!

When we started the Alpha architecture design in 1988, we estimated a 25-year lifetime and a relatively modest 32% per year compounded performance improvement of implementations over that lifetime (1,000 \times total). We guestimated about 10 \times would come from CPU clock improvement, 10 \times from multiple instruction issue, and 10 \times from multiple processors.

5 , 1996  MICROPROCESSOR REPORT

I expect that over the coming decade memory subsystem design will be the *only* important design issue for microprocessors.

The Performance Perspective



The Performance Perspective

- Onur Mutlu, Jared Stark, Chris Wilkerson, and Yale N. Patt,
"Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors"
Proceedings of the 9th International Symposium on High-Performance Computer Architecture (HPCA), Anaheim, CA, February 2003. [Slides \(pdf\)](#)
One of the 15 computer architecture papers of 2003 selected as Top Picks by IEEE Micro.

Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors

Onur Mutlu § Jared Stark † Chris Wilkerson ‡ Yale N. Patt §

§ECE Department
The University of Texas at Austin
{onur,patt}@ece.utexas.edu

†Microprocessor Research
Intel Labs
jared.w.stark@intel.com

‡Desktop Platforms Group
Intel Corporation
chris.wilkerson@intel.com

The Memory Bottleneck

- Onur Mutlu, Jared Stark, Chris Wilkerson, and Yale N. Patt,
"Runahead Execution: An Effective Alternative to Large Instruction Windows"

*IEEE Micro, Special Issue: Micro's Top Picks from Microarchitecture Conferences (**MICRO TOP PICKS**)*, Vol. 23, No. 6, pages 20-25, November/December 2003.

RUNAHEAD EXECUTION: AN EFFECTIVE ALTERNATIVE TO LARGE INSTRUCTION WINDOWS

An Informal Interview on Memory

- Madeleine Gray and Onur Mutlu,

"It's the memory, stupid': A conversation with Onur Mutlu"

HiPEAC info 55, HiPEAC Newsletter, October 2018.

[Shorter Version in Newsletter]

[Longer Online Version with References]

'It's the memory, stupid': A conversation with Onur Mutlu

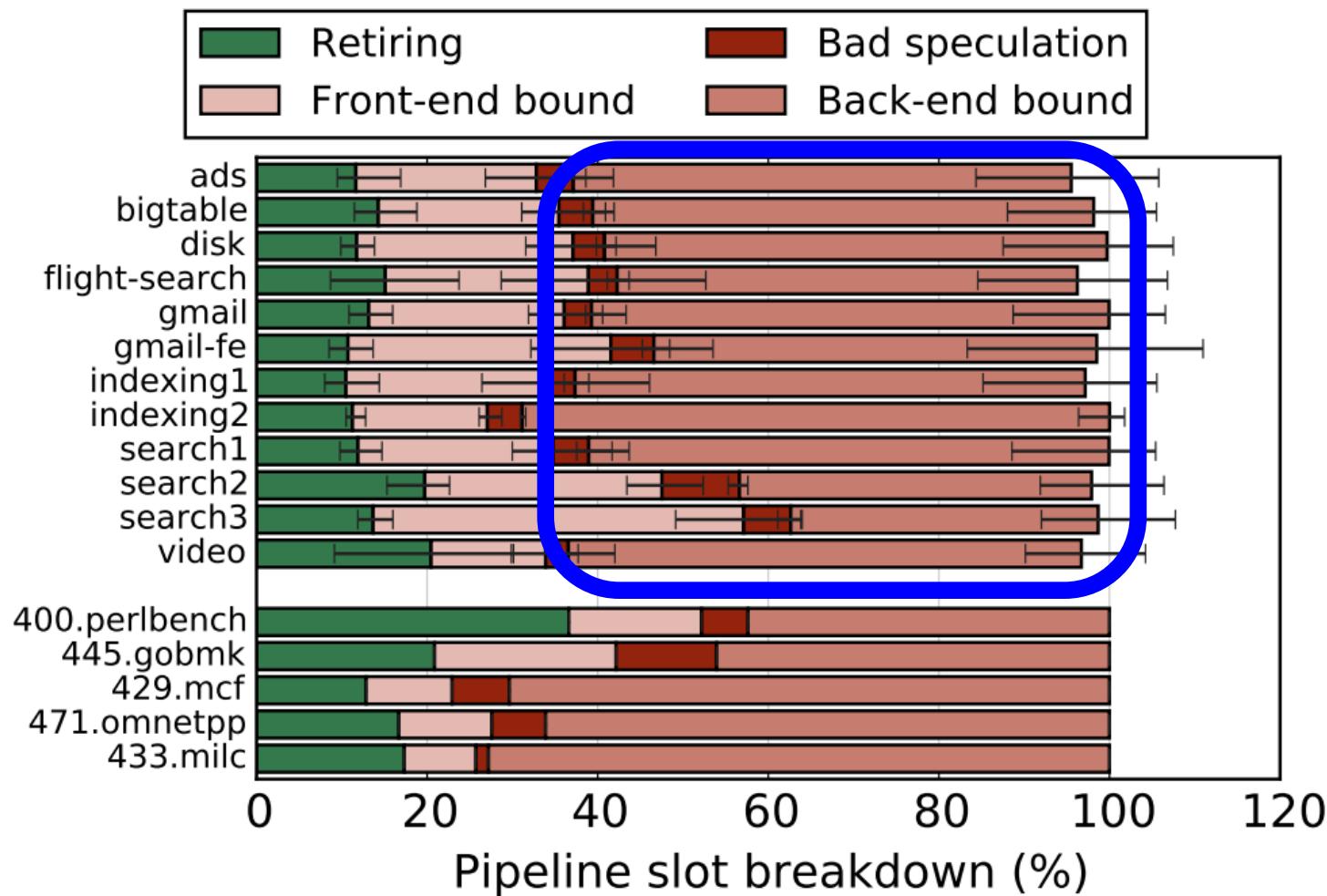
'We're beyond computation; we know how to do computation really well, we can optimize it, we can build all sorts of accelerators ... but the memory – how to feed the data, how to get the data into the accelerators – is a huge problem.'

This was how ETH Zürich and Carnegie Mellon Professor Onur Mutlu opened his course on memory systems and memory-centric computing systems at HiPEAC's summer school, ACACES18. A prolific publisher – he recently bagged the top spot on the International Symposium on Computer Architecture (ISCA) hall of fame – Onur is passionate about computation and communication that are efficient and secure by design. In advance of our Computing Systems Week focusing on data centres, storage, and networking, which takes place next week in Heraklion, HiPEAC picked his brains on all things data-based.



The Performance Perspective (Today)

- All of Google's Data Center Workloads (2015):



The Performance Perspective (Today)

- All of Google's Data Center Workloads (2015):

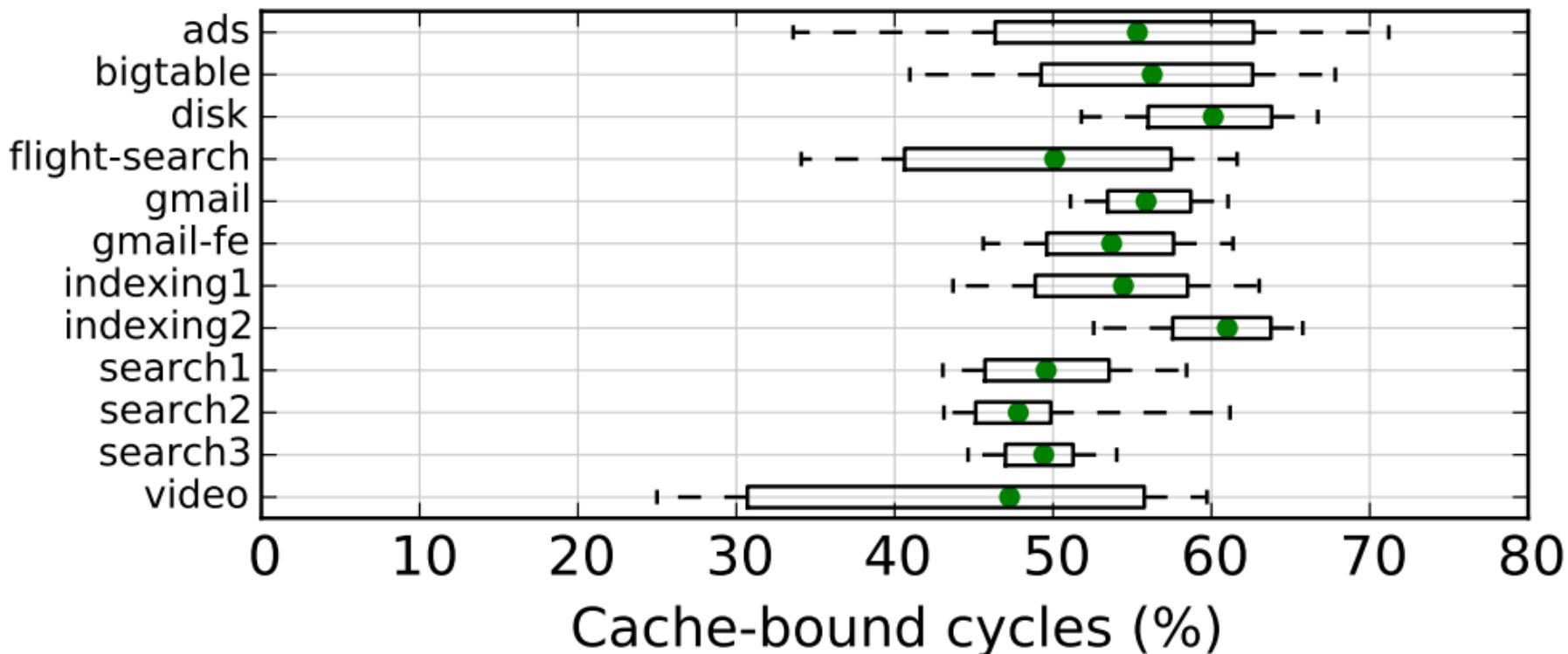
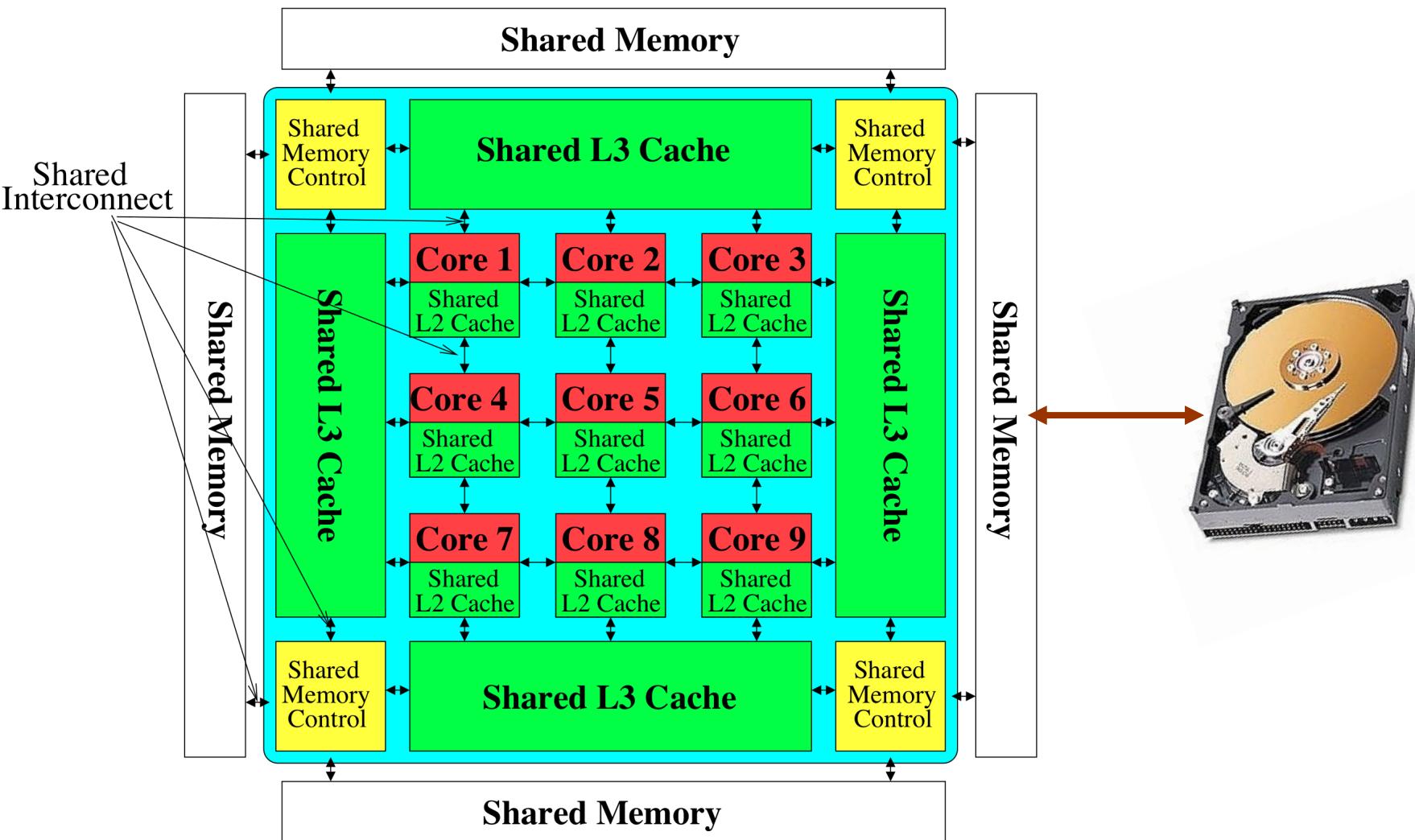


Figure 11: Half of cycles are spent stalled on caches.

Perils of Processor-Centric Design

- Grossly-imbalanced systems
 - Processing done only in **one place**
 - Everything else just stores and moves data: **data moves a lot**
 - Energy inefficient
 - Low performance
 - Complex
- Overly complex and bloated processor (and accelerators)
 - To tolerate data access from memory
 - Complex hierarchies and mechanisms
 - Energy inefficient
 - Low performance
 - Complex

Perils of Processor-Centric Design

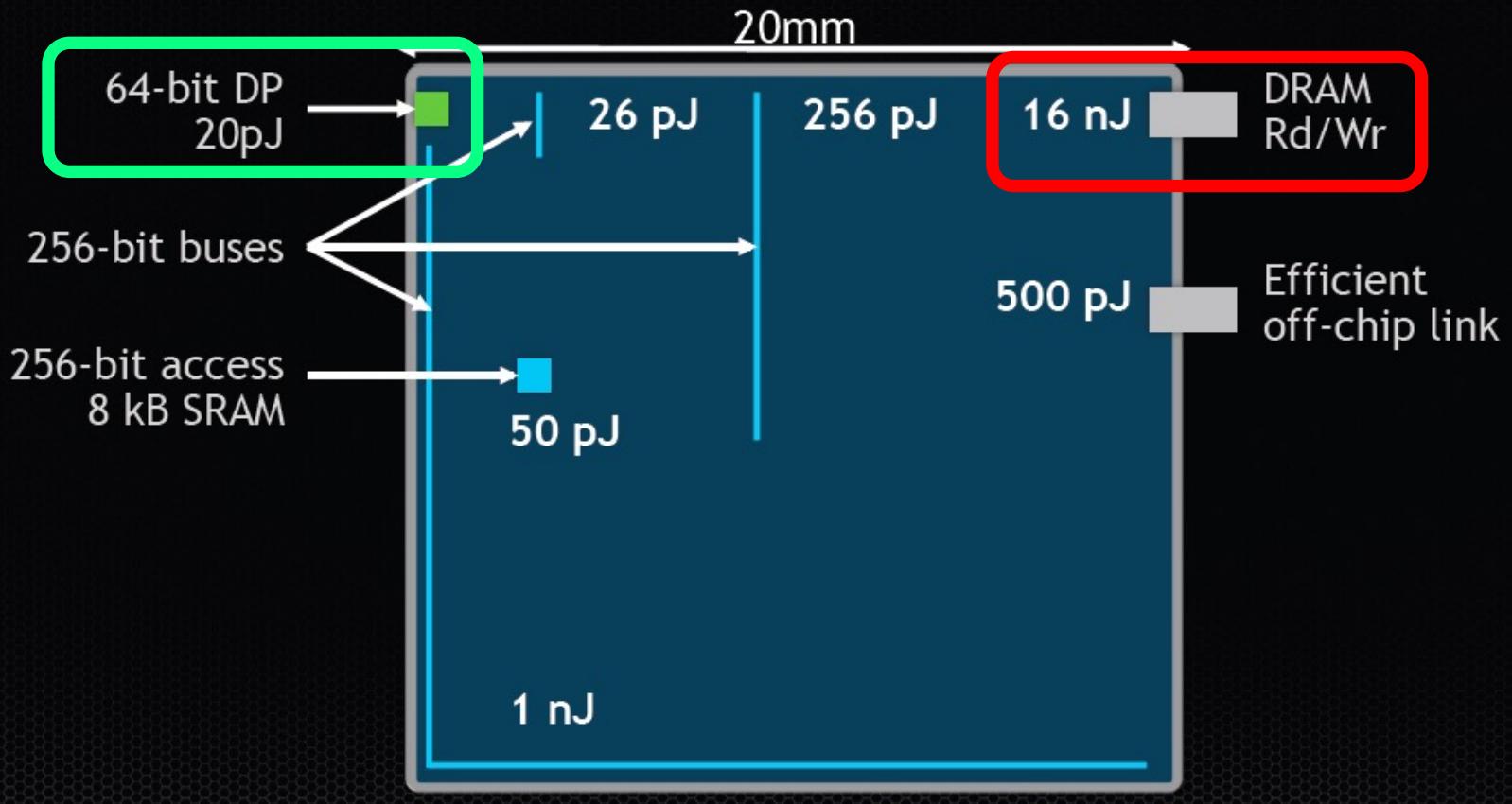


Most of the system is dedicated to storing and moving data

The Energy Perspective

Communication Dominates Arithmetic

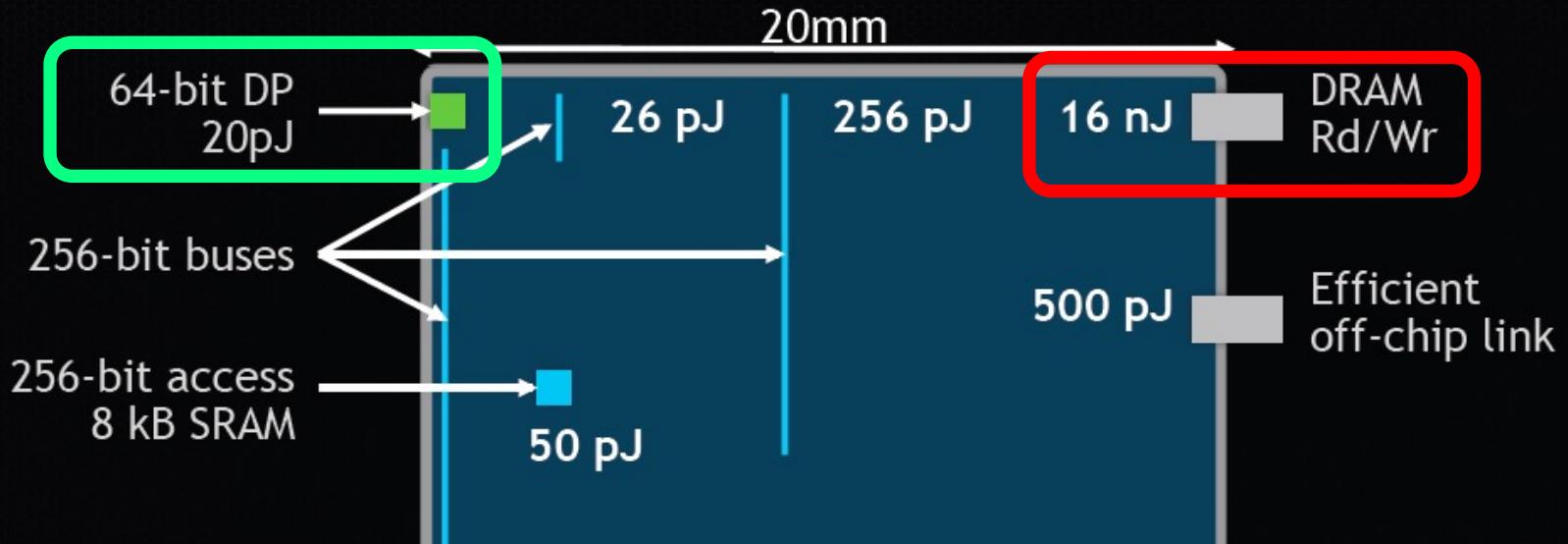
Dally, HiPEAC 2015



Data Movement vs. Computation Energy

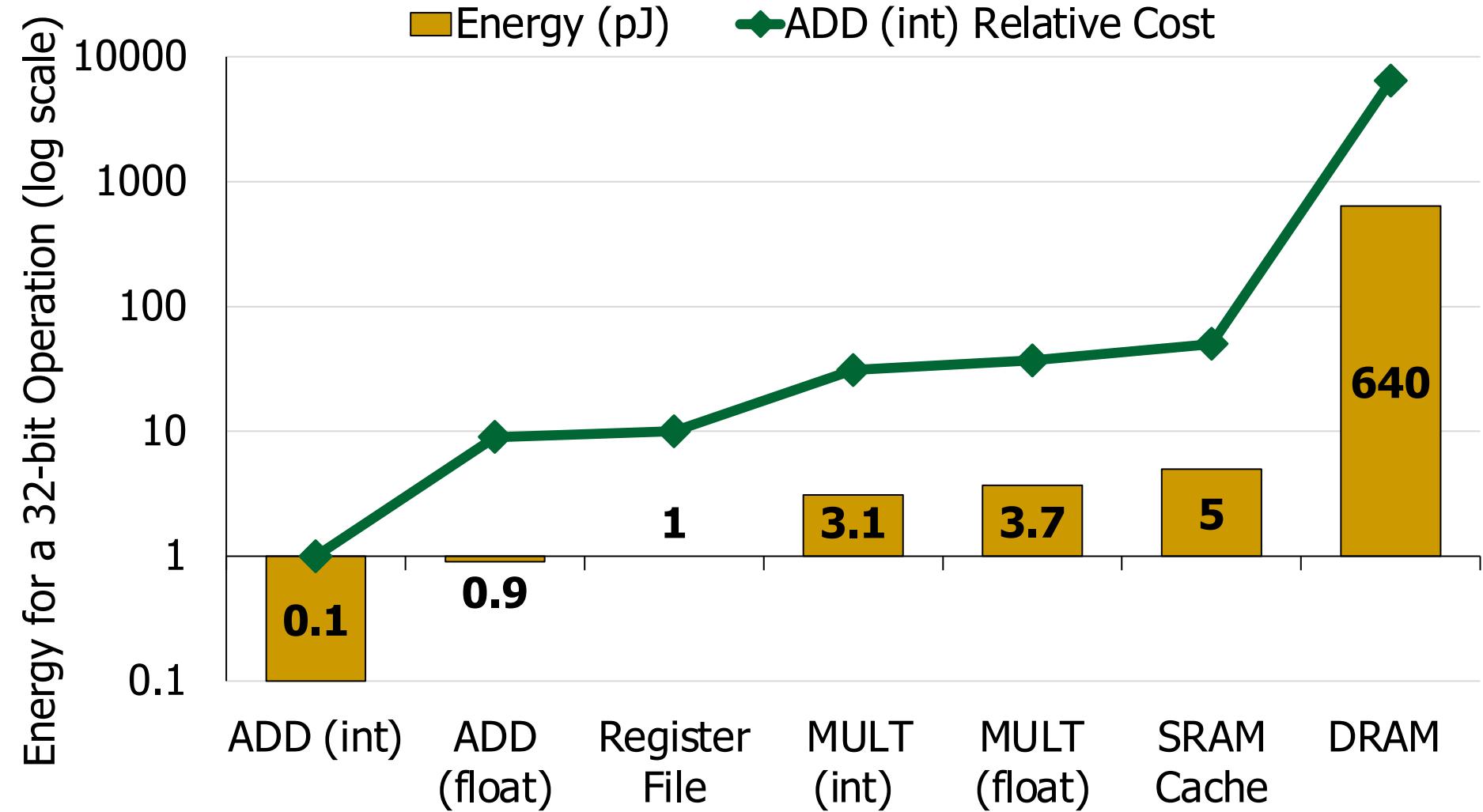
Communication Dominates Arithmetic

Dally, HiPEAC 2015

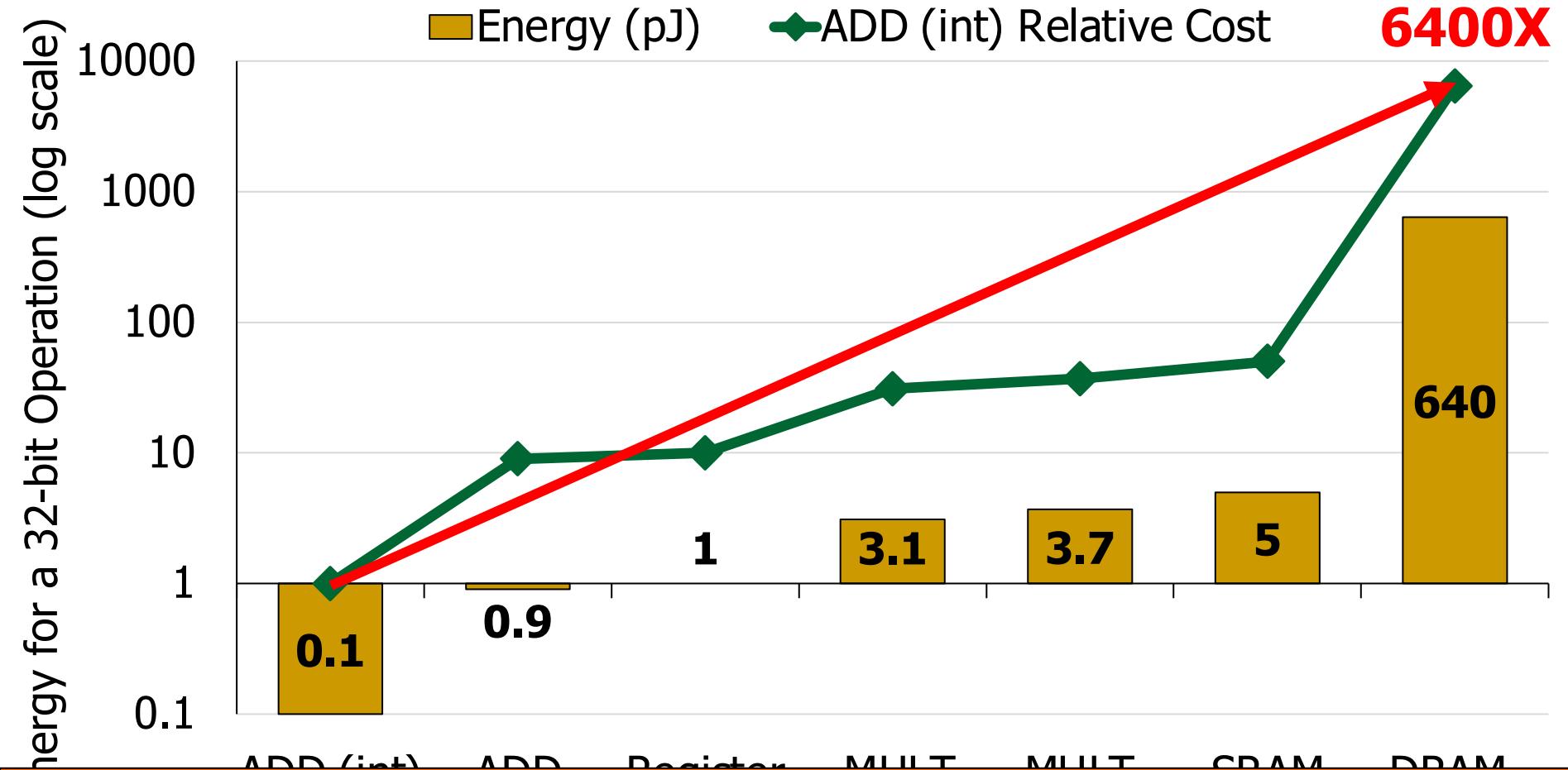


A memory access consumes \sim 100-1000X
the energy of a complex addition

Data Movement vs. Computation Energy



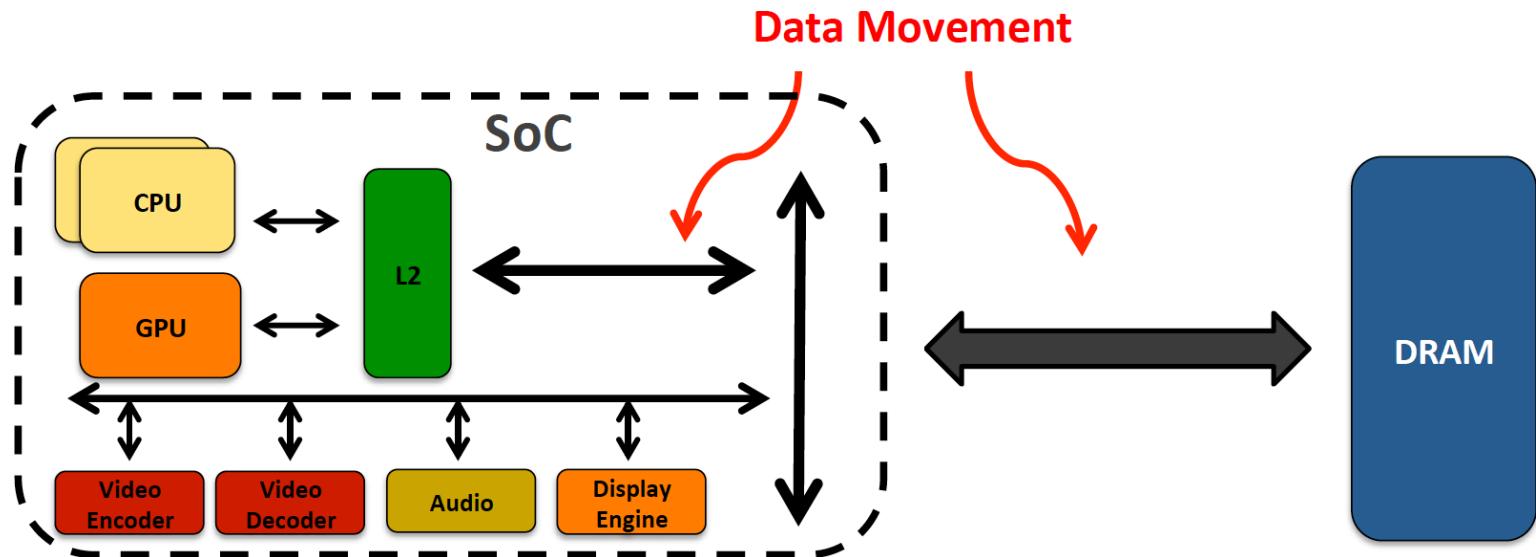
Data Movement vs. Computation Energy



A memory access consumes 6400X
the energy of a simple integer addition

Data Movement vs. Computation Energy

- Data movement is a major system energy bottleneck
 - Comprises 41% of mobile system energy during web browsing [2]
 - Costs ~115 times as much energy as an ADD operation [1, 2]



[1]: Reducing data Movement Energy via Online Data Clustering and Encoding (MICRO'16)

[2]: Quantifying the energy cost of data movement for emerging smart phone workloads on mobile platforms (IISWC'14)

Energy Waste in Mobile Devices

- Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu,
"Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks"

Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Williamsburg, VA, USA, March 2018.

**62.7% of the total system energy
is spent on data movement**

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand¹

Rachata Ausavarungnirun¹

Aki Kuusela³

Saugata Ghose¹

Eric Shiu³

Allan Knies³

Youngsok Kim²

Rahul Thakur³

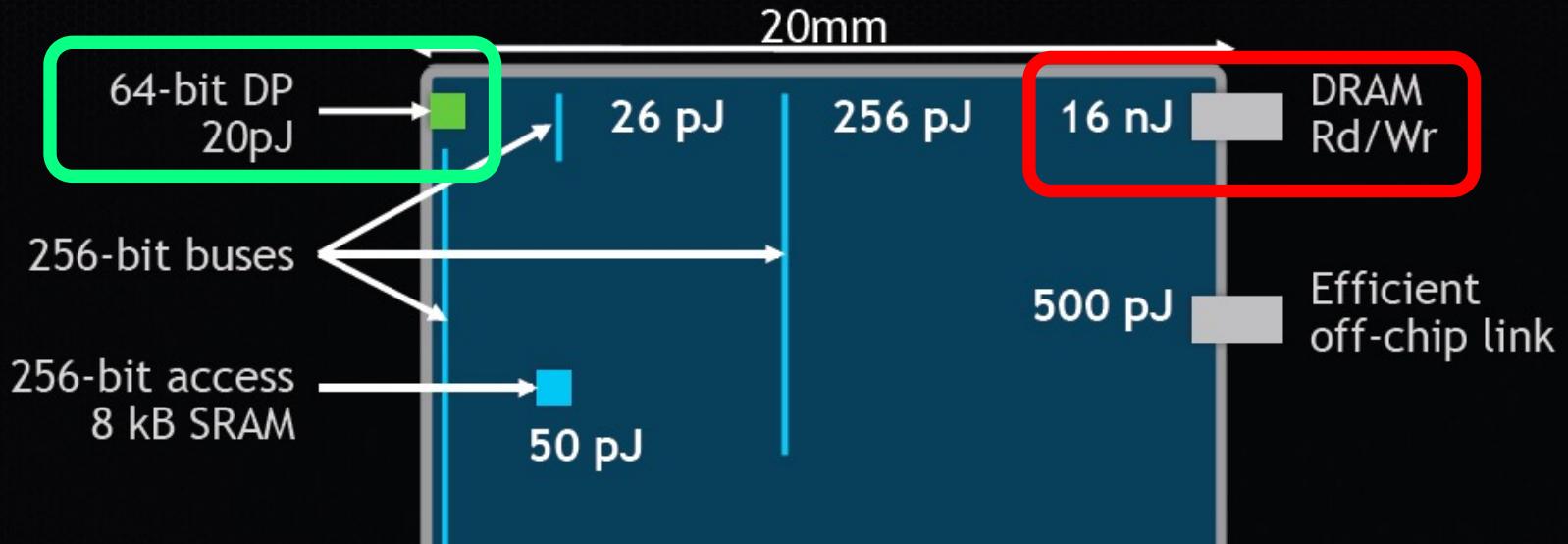
Daehyun Kim^{4,3}

Onur Mutlu^{5,1}

We Do Not Want to Move Data!

Communication Dominates Arithmetic

Dally, HiPEAC 2015



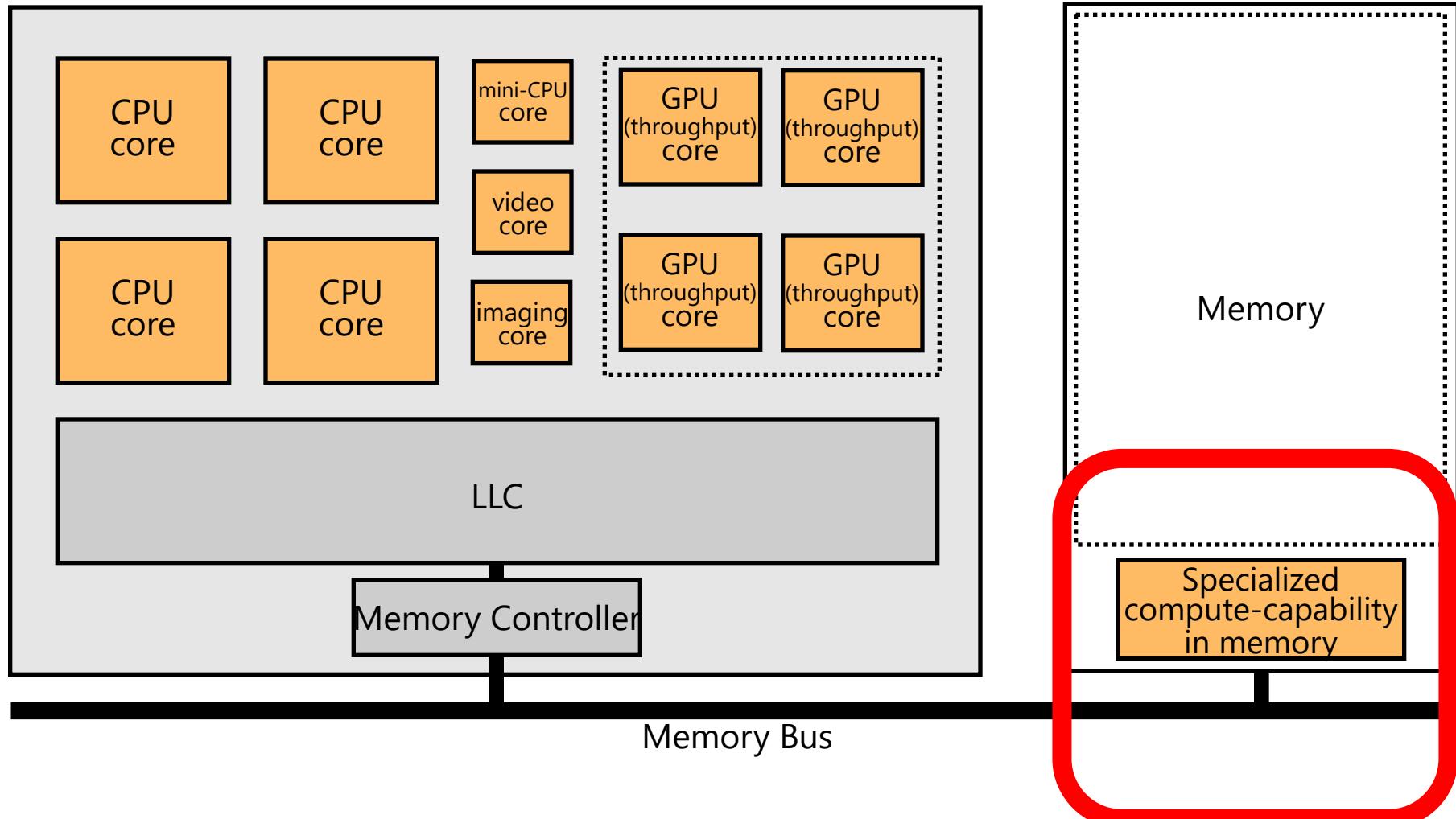
A memory access consumes \sim 100-1000X
the energy of a complex addition

We Need to Think Differently
from the Past Approaches

We Need A Paradigm Shift To ...

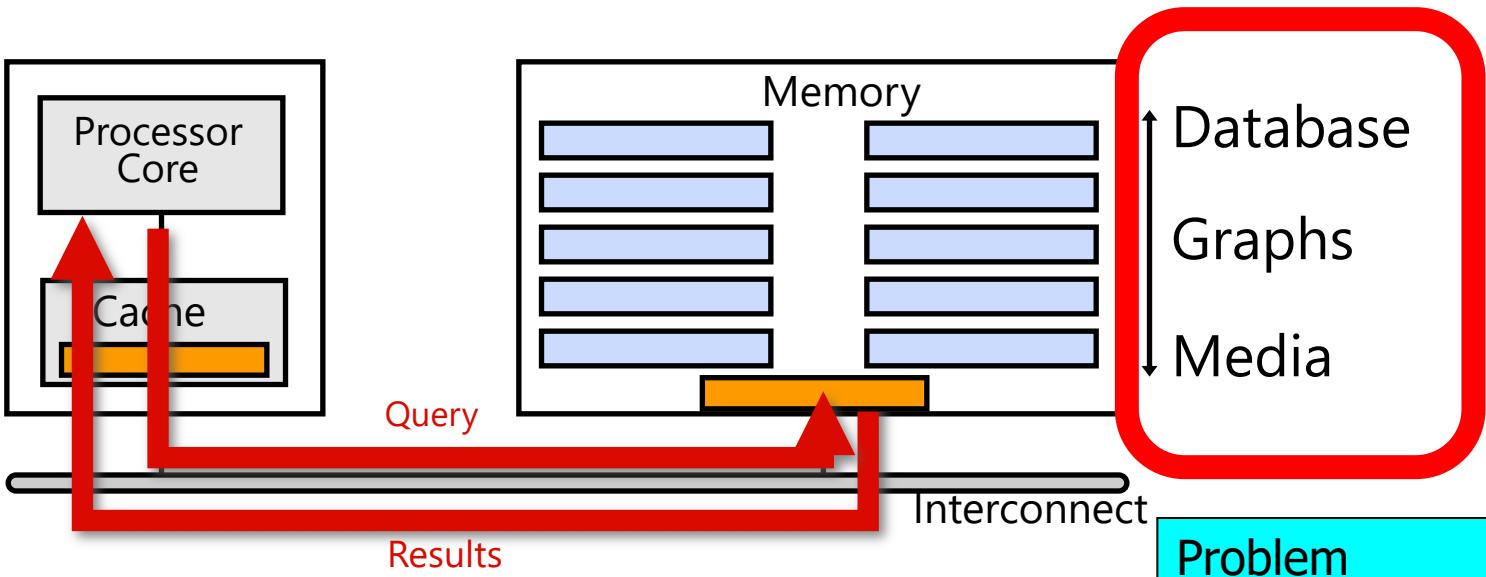
- Enable computation with **minimal data movement**
- Compute where it makes sense (**where data resides**)
- Make computing architectures more **data-centric**

Memory as an Accelerator

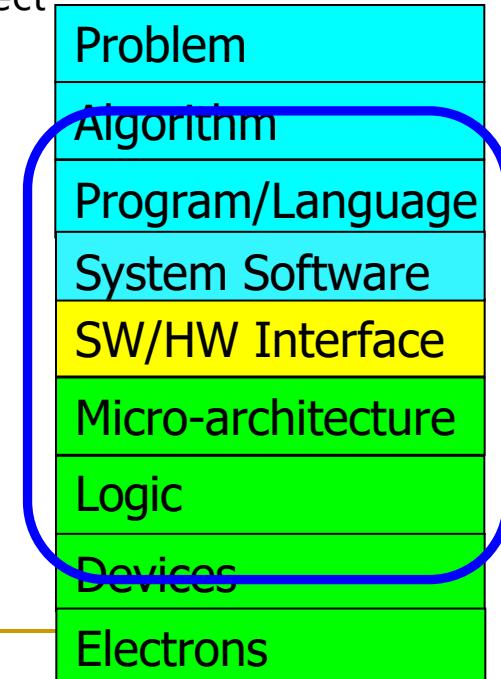


Memory similar to a “conventional” accelerator

Goal: Processing Inside Memory



- Many questions ... How do we design the:
 - compute-capable memory & controllers?
 - processor chip and in-memory units?
 - software and hardware interfaces?
 - system software, compilers, languages?
 - algorithms and theoretical foundations?



PIM Review and Open Problems

A Modern Primer on Processing in Memory

Onur Mutlu^{a,b}, Saugata Ghose^{b,c}, Juan Gómez-Luna^a, Rachata Ausavarungnirun^d

SAFARI Research Group

^a*ETH Zürich*

^b*Carnegie Mellon University*

^c*University of Illinois at Urbana-Champaign*

^d*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.

A Modern Primer on Processing in Memory

Onur Mutlu^{a,b}, Saugata Ghose^{b,c}, Juan Gómez-Luna^a, Rachata Ausavarungnirun^d

SAFARI Research Group

^aETH Zürich

^bCarnegie Mellon University

^cUniversity of Illinois at Urbana-Champaign

^dKing Mongkut's University of Technology North Bangkok

Abstract

Modern computing systems are overwhelmingly designed to move data to computation. This design choice goes directly against at least three key trends in computing that cause performance, scalability and energy bottlenecks: (1) data access is a key bottleneck as many important applications are increasingly data-intensive, and memory bandwidth and energy do not scale well, (2) energy consumption is a key limiter in almost all computing platforms, especially server and mobile systems, (3) data movement, especially off-chip to on-chip, is very expensive in terms of bandwidth, energy and latency, much more so than computation. These trends are especially severely-felt in the data-intensive server and energy-constrained mobile systems of today.

At the same time, conventional memory technology is facing many technology scaling challenges in terms of reliability, energy, and performance. As a result, memory system architects are open to organizing memory in different ways and making it more intelligent, at the expense of higher cost. The emergence of 3D-stacked memory plus logic, the adoption of error correcting codes inside the latest DRAM chips, proliferation of different main memory standards and chips, specialized for different purposes (e.g., graphics, low-power, high bandwidth, low latency), and the necessity of designing new solutions to serious reliability and security issues, such as the RowHammer phenomenon, are an evidence of this trend.

This chapter discusses recent research that aims to practically enable computation close to data, an approach we call *processing-in-memory* (PIM). PIM places computation mechanisms in or near where the data is stored (i.e., inside the memory chips, in the logic layer of 3D-stacked memory, or in the memory controllers), so that data movement between the computation units and memory is reduced or eliminated. While the general idea of PIM is not new, we discuss motivating trends in applications as well as memory circuits/technology that greatly exacerbate the need for enabling it in modern computing systems. We examine at least two promising new approaches to designing PIM systems to accelerate important data-intensive applications: (1) *processing using memory* by exploiting analog operational properties of DRAM chips to perform massively-parallel operations in memory, with low-cost changes, (2) *processing near memory* by exploiting 3D-stacked memory technology design to provide high memory bandwidth and low memory latency to in-memory logic. In both approaches, we describe and tackle relevant cross-layer research, design, and adoption challenges in devices, architecture, systems, and programming models. Our focus is on the development of in-memory processing designs that can be adopted in real computing platforms at low cost. We conclude by discussing work on solving key challenges to the practical adoption of PIM.

Keywords: memory systems, data movement, main memory, processing-in-memory, near-data processing, computation-in-memory, processing using memory, processing near memory, 3D-stacked memory, non-volatile memory, energy efficiency, high-performance computing, computer architecture, computing paradigm, emerging technologies, memory scaling, technology scaling, dependable systems, robust systems, hardware security, system security, latency, low-latency computing

Contents

1	Introduction	2
2	Major Trends Affecting Main Memory	4
3	The Need for Intelligent Memory Controllers to Enhance Memory Scaling	6
4	Perils of Processor-Centric Design	9
5	Processing-in-Memory (PIM): Technology Enablers and Two Approaches	12
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
6	Processing Using Memory (PUM)	14
6.1	RowClone	14
6.2	Ambit	15
6.3	Gather-Scatter DRAM	17
6.4	In-DRAM Security Primitives	17
7	Processing Near Memory (PNM)	18
7.1	Tesseract: Coarse-Grained Application-Level PNM Acceleration of Graph Processing	19
7.2	Function-Level PNM Acceleration of Mobile Consumer Workloads	20
7.3	Programmer-Transparent Function-Level PNM Acceleration of GPU Applications	21
7.4	Instruction-Level PNM Acceleration with PIM-Enabled Instructions (PEI) . .	21
7.5	Function-Level PNM Acceleration of Genome Analysis Workloads	22
7.6	Application-Level PNM Acceleration of Time Series Analysis	23
8	Enabling the Adoption of PIM	24
8.1	Programming Models and Code Generation for PIM	24
8.2	PIM Runtime: Scheduling and Data Mapping	25
8.3	Memory Coherence	27
8.4	Virtual Memory Support	27
8.5	Data Structures for PIM	28
8.6	Benchmarks and Simulation Infrastructures	29
8.7	Real PIM Hardware Systems and Prototypes	30
8.8	Security Considerations	30
9	Conclusion and Future Outlook	31

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, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16]. 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, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49], 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, 50, 51, 52, 53, 54]. These costs are often exacerbated by the fact that much of the data brought into the caches is *not reused* by the CPU [52, 53, 55, 56], 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 appli-

Processing Data Where It Makes Sense

Processing in/near Memory: An Old Idea

- Kautz, "Cellular Logic-in-Memory Arrays", IEEE TC 1969.

IEEE TRANSACTIONS ON COMPUTERS, VOL. C-18, NO. 8, AUGUST 1969

Cellular Logic-in-Memory Arrays

WILLIAM H. KAUTZ, MEMBER, IEEE

Abstract—As a direct consequence of large-scale integration, many advantages in the design, fabrication, testing, and use of digital circuitry can be achieved if the circuits can be arranged in a two-dimensional iterative, or cellular, array of identical elementary networks, or cells. When a small amount of storage is included in each cell, the same array may be regarded either as a logically enhanced memory array, or as a logic array whose elementary gates and connections can be "programmed" to realize a desired logical behavior.

In this paper the specific engineering features of such cellular logic-in-memory (CLIM) arrays are discussed, and one such special-purpose array, a cellular sorting array, is described in detail to illustrate how these features may be achieved in a particular design. It is shown how the cellular sorting array can be employed as a single-address, multiword memory that keeps in order all words stored within it. It can also be used as a content-addressed memory, a pushdown memory, a buffer memory, and (with a lower logical efficiency) a programmable array for the realization of arbitrary switching functions. A second version of a sorting array, operating on a different sorting principle, is also described.

Index Terms—Cellular logic, large-scale integration, logic arrays logic in memory, push-down memory, sorting, switching functions.

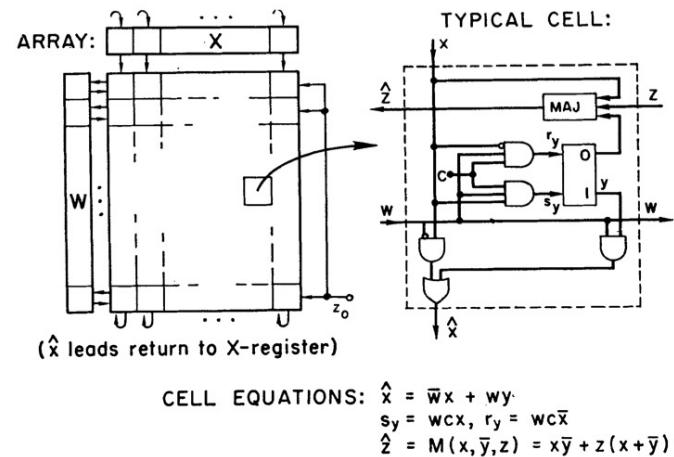


Fig. 1. Cellular sorting array I.

Processing in/near Memory: An Old Idea

- Stone, "A Logic-in-Memory Computer," IEEE TC 1970.

A Logic-in-Memory Computer

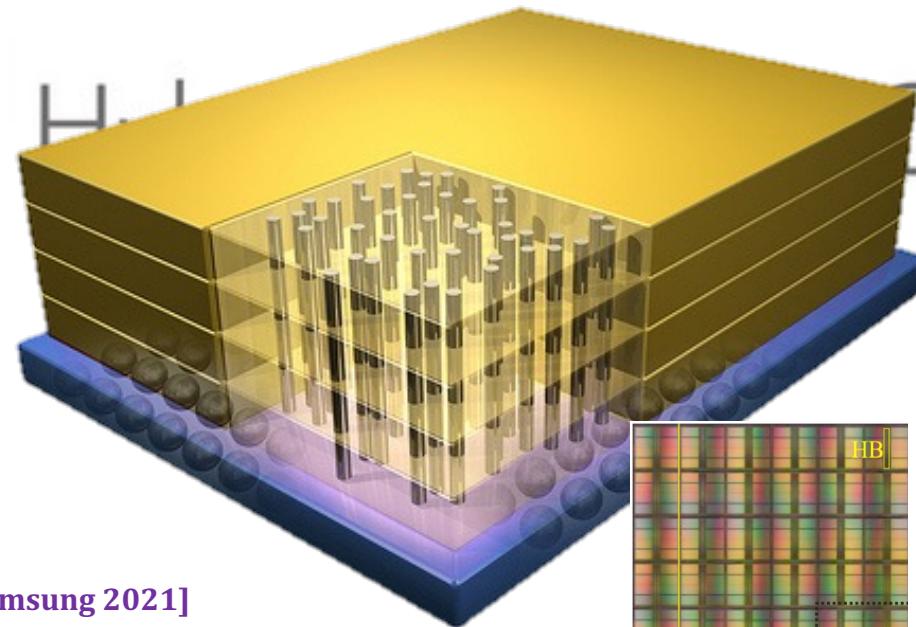
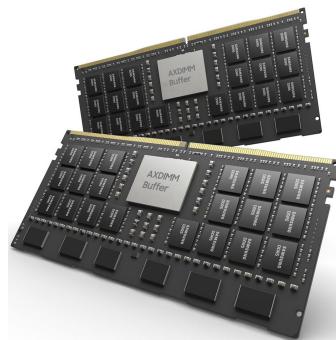
HAROLD S. STONE

Abstract—If, as presently projected, the cost of microelectronic arrays in the future will tend to reflect the number of pins on the array rather than the number of gates, the logic-in-memory array is an extremely attractive computer component. Such an array is essentially a microelectronic memory with some combinational logic associated with each storage element.

Why In-Memory Computation Today?

- **Huge problems with Memory Technology**
 - Memory technology scaling is not going well (e.g., RowHammer)
 - Many scaling issues demand intelligence in memory
- **Huge demand from Applications & Systems**
 - Data access bottleneck
 - Energy & power bottlenecks
 - Data movement energy dominates computation energy
 - Need all at the same time: performance, energy, sustainability
 - We can improve all metrics by minimizing data movement
- **Designs are squeezed in the middle**

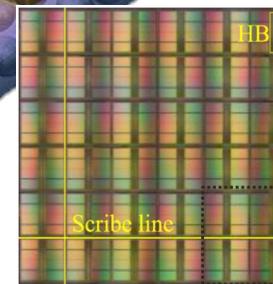
Processing-in-Memory Landscape Today



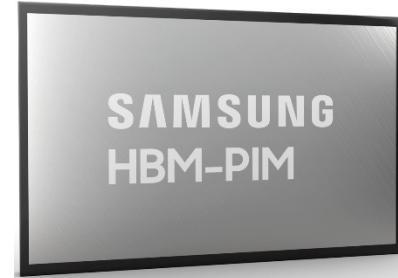
[Samsung 2021]



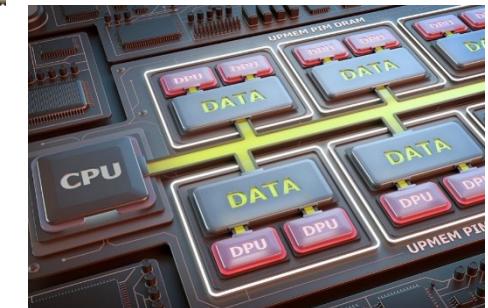
[Alibaba 2022]



[SK Hynix 2022]



[Samsung 2021]



[UPMEM 2019]

Why In-Memory Computation Today?

- **Huge problems with Memory Technology**
 - Memory technology scaling is not going well (e.g., RowHammer)
 - Many scaling issues demand intelligence in memory
- **Huge demand from Applications & Systems**
 - Data access bottleneck
 - Energy & power bottlenecks
 - Data movement energy dominates computation energy
 - Need all at the same time: performance, energy, sustainability
 - We can improve all metrics by minimizing data movement
- **Designs are squeezed in the middle**

Main Memory Needs Intelligent Controllers

Memory Scaling Issues Are Real

- Onur Mutlu,

"Memory Scaling: A Systems Architecture Perspective"

Proceedings of the 5th International Memory

Workshop (IMW), Monterey, CA, May 2013. Slides

(pptx) (pdf)

EETimes Reprint

Memory Scaling: A Systems Architecture Perspective

Onur Mutlu

Carnegie Mellon University

onur@cmu.edu

<http://users.ece.cmu.edu/~omutlu/>

Application Scaling Issues Are Real

- Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,

"A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing"

Proceedings of the 42nd International Symposium on Computer Architecture (ISCA), Portland, OR, June 2015.

[Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)]

Top Picks Honorable Mention by IEEE Micro.

Selected to the ISCA-50 25-Year Retrospective Issue covering 1996-2020 in 2023 (Retrospective (pdf) Full Issue).

A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn Sungpack Hong[§] Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi

junwhan@snu.ac.kr, sungpack.hong@oracle.com, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr

Seoul National University

[§]Oracle Labs

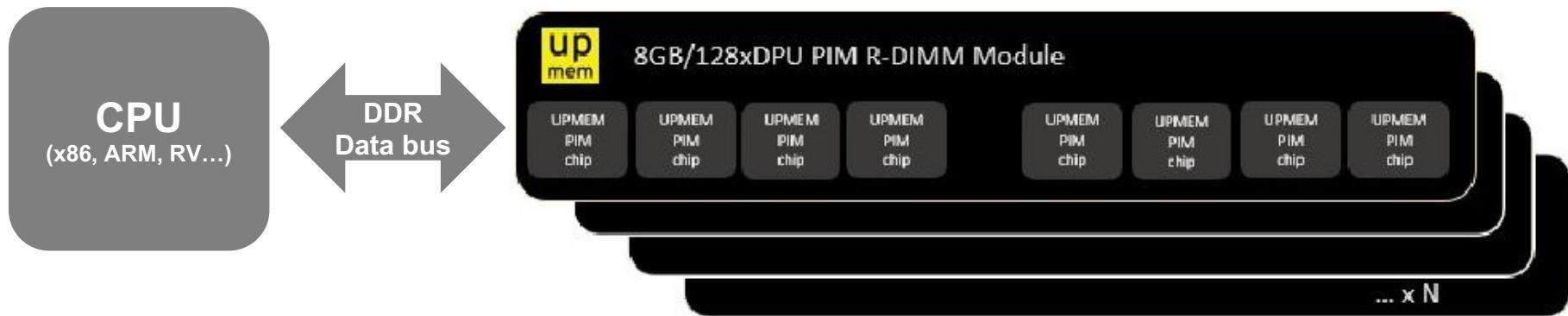
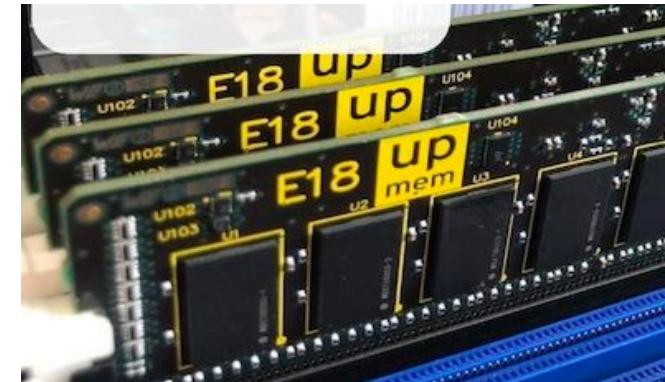
[†]Carnegie Mellon University

Eliminating the Adoption Barriers

Processing-in-Memory in the Real World

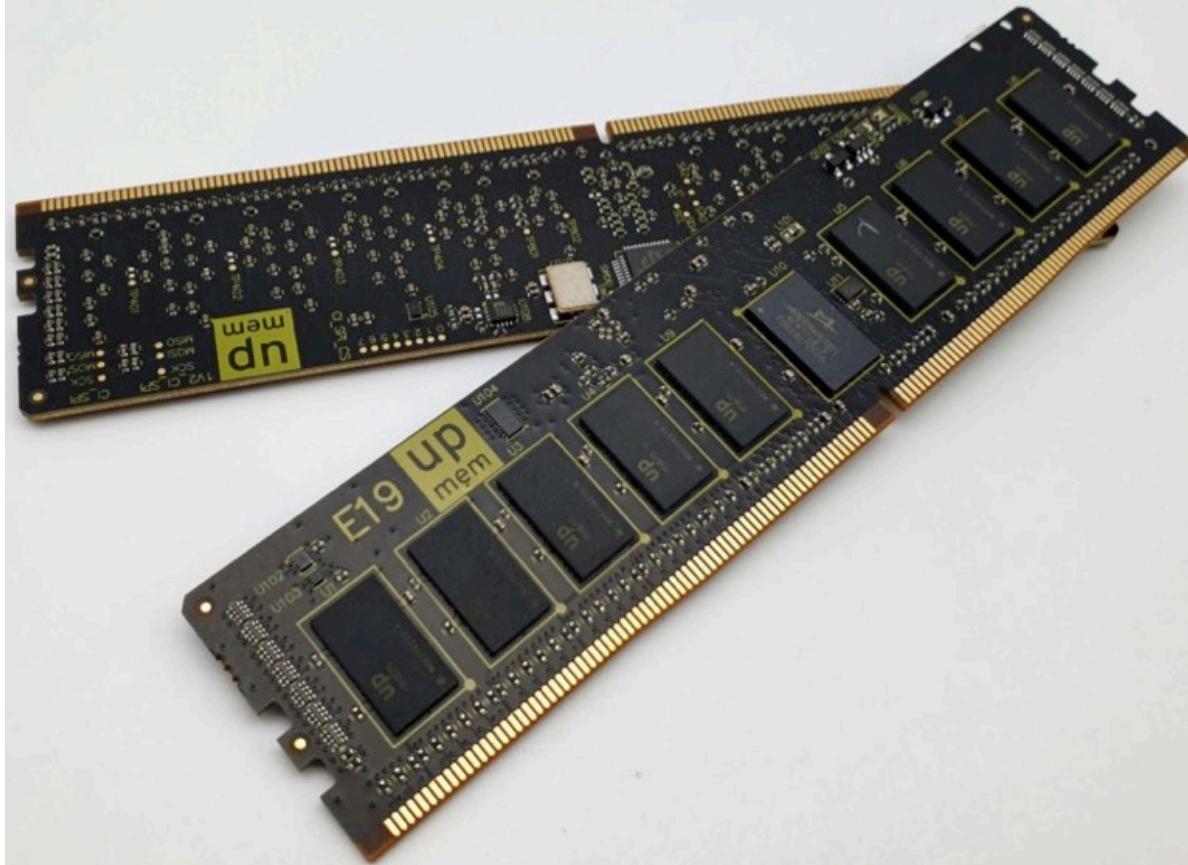
UPMEM Processing-in-DRAM Engine (2019)

- 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

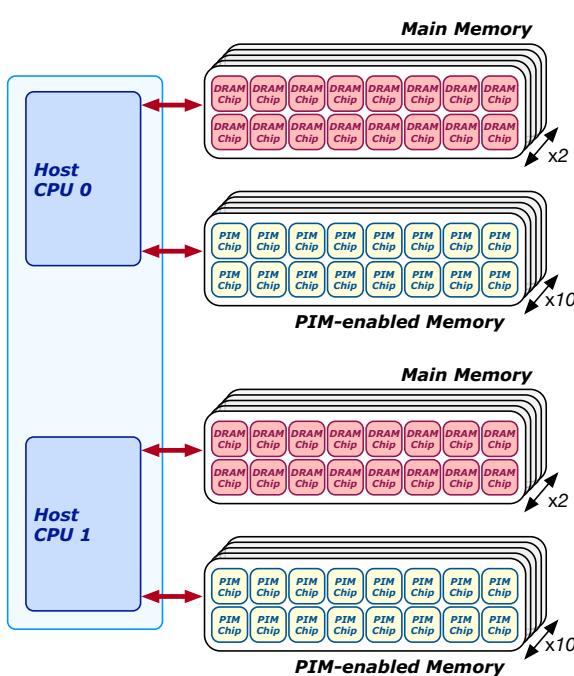


UPMEM Memory Modules

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



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

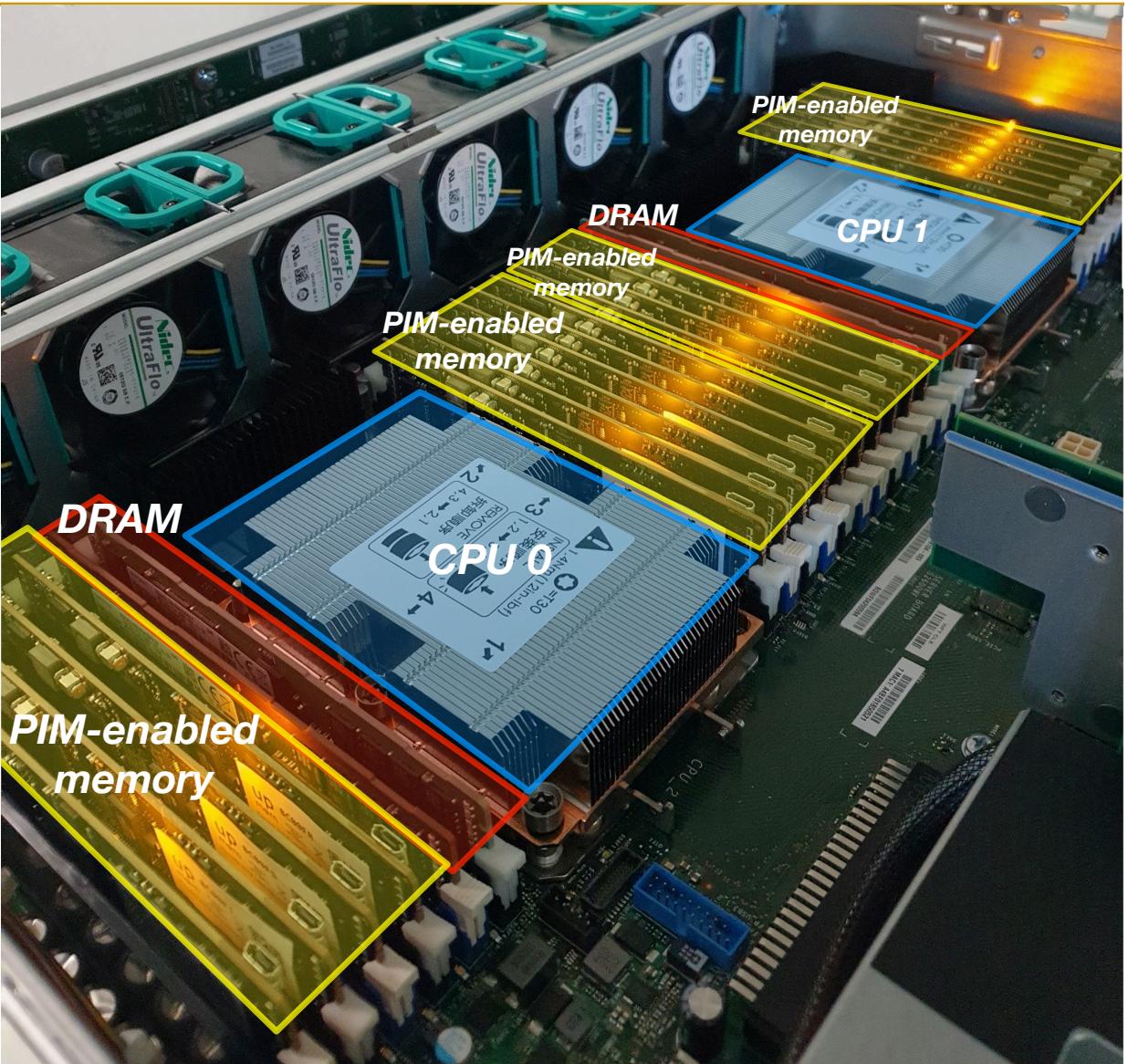
GERALDO F. OLIVEIRA, ETH Zürich, Switzerland

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.



Understanding a Modern PIM Architecture (I)

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

JUAN GÓMEZ-LUNA¹, IZZAT EL HAJJ², IVAN FERNANDEZ^{1,3}, CHRISTINA GIANNOULA^{1,4},
GERALDO F. OLIVEIRA¹, AND ONUR MUTLU¹

¹ETH Zürich

²American University of Beirut

³University of Malaga

⁴National 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>

Understanding a Modern PIM Architecture (II)

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

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

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

ETH Zürich SAFARI

2:26 / 2:57:10

SAFARI Live Seminar: Understanding a Modern Processing-in-Memory Architecture

2,579 views • Streamed live on Jul 12, 2021

93 likes 0 dislikes SHARE SAVE ...



Onur Mutlu Lectures
18.7K subscribers

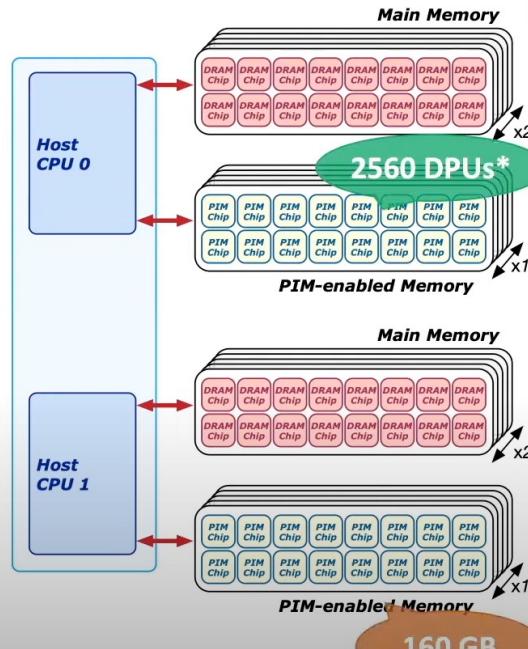
SUBSCRIBED



Longer Lectures on UPMEM PIM

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



13:12 / 31:45

* There are 4 faulty DPUs in the system that we use in our experiments. Thus, the maximum number of DPUs we can use is 1,225.



PIM Course: Lecture 3: Real-world PIM: UPMEM PIM - Fall 2022



Onur Mutlu Lectures

31.7K subscribers



Subscribed

18



Share



Clip

Save



564 views 4 months ago Livestream - P&S Data-Centric Architectures: Fundamentally Improving Performance and Energy (Fall 2022)

Projects & Seminars, ETH Zürich, Fall 2022

Data-Centric Architectures: Fundamentally Improving Performance and Energy

(https://safari.ethz.ch/projects_and_s...) Show more

https://youtu.be/p_sLhKeo6ys

<https://youtu.be/7c6x5GJG6dw>

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



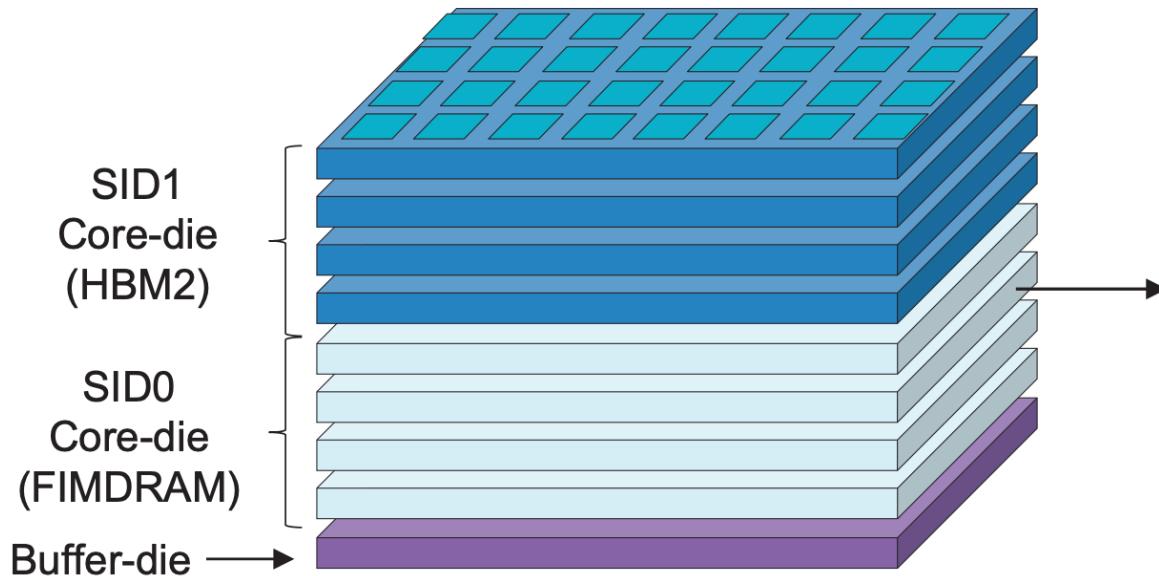
The new architecture will deliver over twice the system performance and reduce energy consumption by more than 70%

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 Kwon¹, Suk Han Lee¹, Jaehoon Lee¹, Sang-Hyuk Kwon¹, Je Min Ryu¹, Jong-Pil Son¹, Seongil Oh¹, Hak-Soo Yu¹, Haesuk Lee¹, Soo Young Kim¹, Younghmin Cho¹, Jin Guk Kim¹, Jongyoon Choi¹, Hyun-Sung Shin¹, Jin Kim¹, BengSeng Phua², Hyo Young Kim¹, Myeong Jun Song¹, Ahn Choi¹, Daeho Kim¹, Soo Young Kim¹, Eun-Bong Kim¹, David Wang², Shinhaeng Kang¹, Yuhwan Ro³, Seungwoo Seo³, JoonHo Song³, Jaeyoun Youn¹, Kyomin Sohn¹, Nam Sung Kim¹

¹Samsung Electronics, Hwaseong, Korea

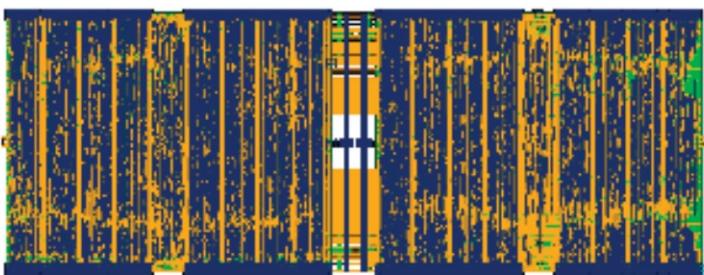
²Samsung Electronics, San Jose, CA

³Samsung Electronics, Suwon, Korea

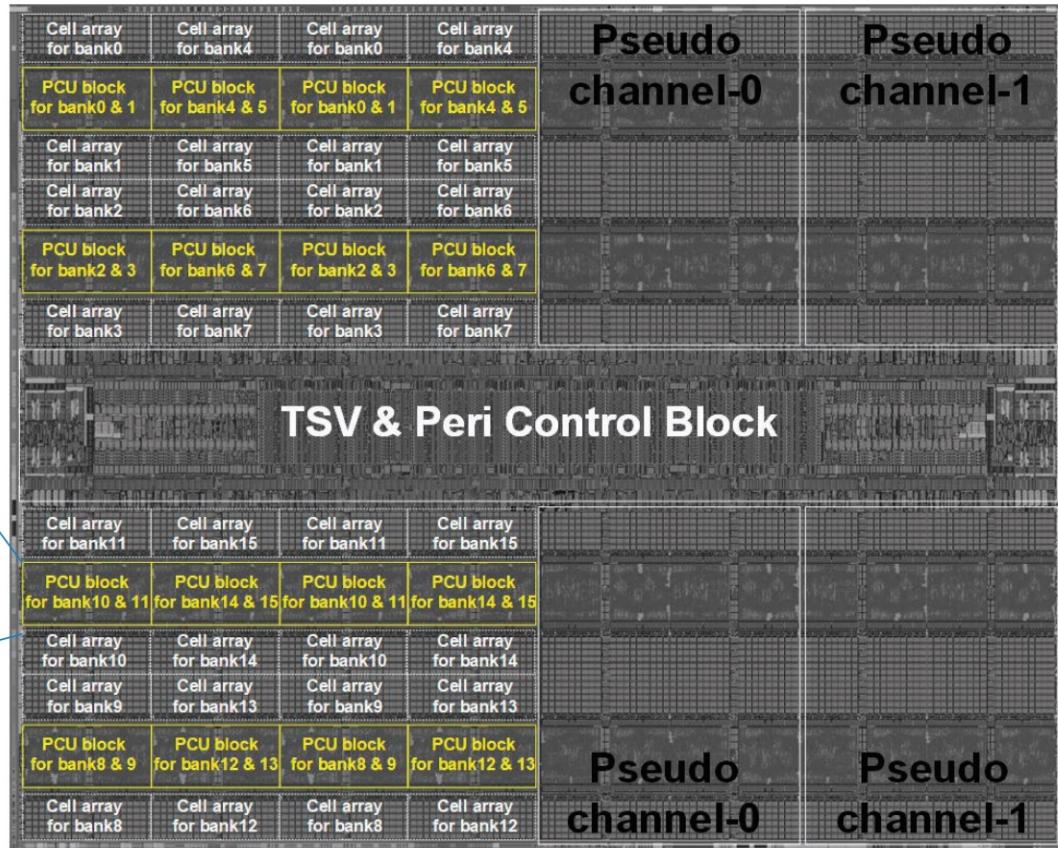
Samsung Function-in-Memory DRAM (2021)

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

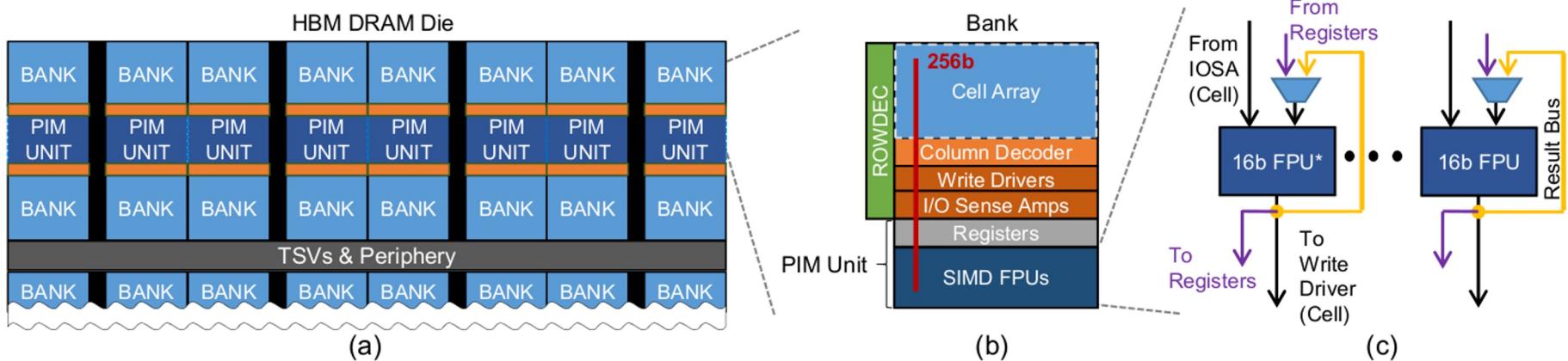
25.4 A 20nm 6GB Function-In-Memory DRAM, Based on HBM2 with a 1.2TFLOPS Programmable Computing Unit Using Bank-Level Parallelism, for Machine Learning Applications

Young-Cheon Kwon¹, Suk Han Lee¹, Jaehoon Lee¹, Sang-Hyuk Kwon¹, Je Min Ryu¹, Jong-Pil Son¹, Seongil O¹, Hak-Soo Yu¹, Haesuk Lee¹, Soo Young Kim¹, Youngmin Cho¹, Jin Guk Kim¹, Jongyoon Choi¹, Hyun-Sung Shin¹, Jin Kim¹, BengSeung Phua², HyoungMin Kim¹, Myeong Jun Song¹, Ahn Choi¹, Daeho Kim¹, SooYoung Kim¹, Eun-Bong Kim¹, David Wang¹, Shinhwang Kang¹, Yuhwan Ro¹, Seungwoo Seo¹, JoonHo Song¹, Jaeyoun Youn¹, Kyomin Sohn¹, Nam Sung Kim¹

¹Samsung Electronics, Hwasung, Korea
²Samsung Electronics, San Jose, CA
³Samsung Electronics, Suwon, Korea

FIMDRAM: System Organization

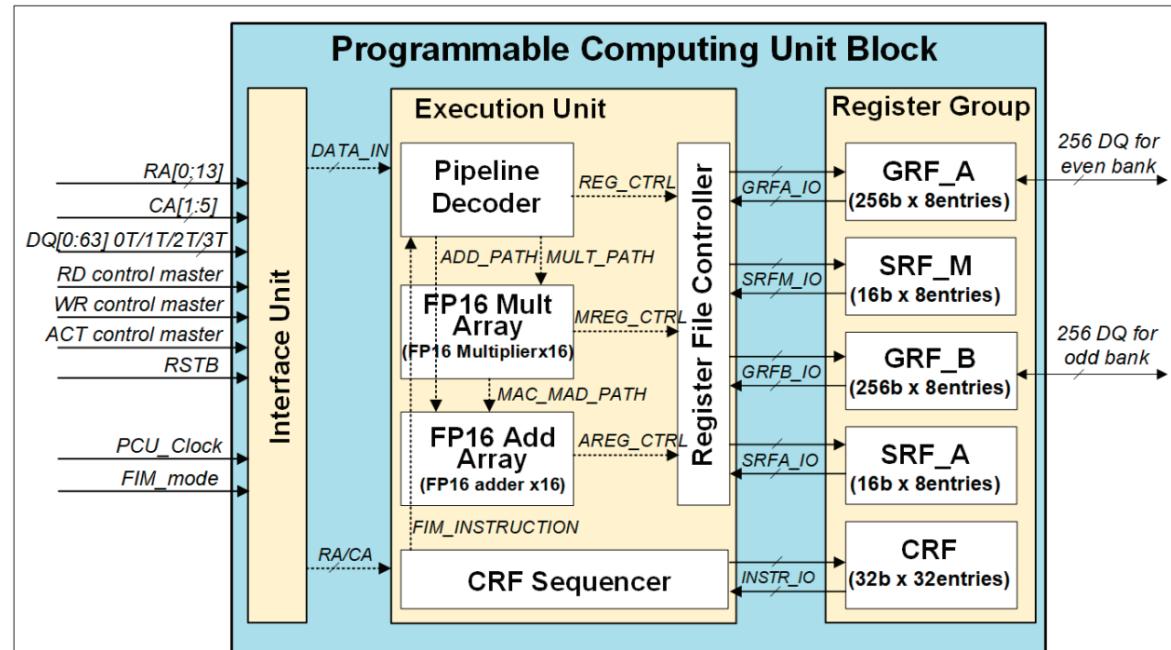
- PIM units respond to standard DRAM column commands (RD or WR)
 - Compliant with unmodified JEDEC controllers
- They execute one wide-SIMD operation commanded by a PIM instruction with deterministic latency in a lock-step manner
- A PIM unit can get 16 16-bit operands from IOSAs, a register, and/or the result bus



Samsung Function-in-Memory DRAM (2021)

Programmable Computing Unit

- Configuration of PCU block
 - Interface unit to control data flow
 - Execution unit to perform operations
 - Register group
 - 32 entries of CRF for instruction memory
 - 16 GRF for weight and accumulation
 - 16 SRF to store constants for MAC operations



[Block diagram of PCU in FIMDRAM]

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 Kwon¹, Suk Han Lee¹, Jaehoon Lee¹, Sang-Hyuk Kwon¹,
Je Min Ryu¹, Jong-Pil Son¹, Seengil O¹, Hak-Soo Yu¹, Haesuk Lee¹,
Soo Young Kim¹, Youngmin Cho¹, Jin Guk Kim¹, Jongyoon Choi¹,
Hyun-Sung Shin¹, Jin Kim¹, BengSeng Phua², HyoungMin Kim¹,
Myeong Jun Song¹, Ahn Choi¹, Daeho Kim¹, SooYoung Kim¹, Eun-Bong Kim¹,
David Wang¹, Shinhwa Kang¹, Yuhwan Ro¹, Seungwoo Seo¹, JoonHo Song¹,
Jayoun Youn¹, Kyomin Sohn¹, Nam Sung Kim¹

¹Samsung Electronics, Hwaseong, Korea

²Samsung Electronics, San Jose, CA

³Samsung Electronics, Suwon, Korea

Samsung Function-in-Memory DRAM (2021)

[Available instruction list for FIM operation]

Type	CMD	Description
Floating Point	ADD	FP16 addition
	MUL	FP16 multiplication
	MAC	FP16 multiply-accumulate
	MAD	FP16 multiply and add
Data Path	MOVE	Load or store data
	FILL	Copy data from bank to GRFs
Control Path	NOP	Do nothing
	JUMP	Jump instruction
	EXIT	Exit instruction

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 Kwon¹, Suk Han Lee¹, Jaehoon Lee¹, Sang-Hyuk Kwon¹,
Je Min Ryu¹, Jong-Pil Son¹, Seengil Oh¹, Hak-Soo Yu¹, Haesuk Lee¹,
Soo Young Kim¹, Youngmin Cho¹, Jin Guk Kim¹, Jongyoon Choi¹,
Hyun-Sung Shin¹, Jin Kim¹, BengSeng Phua², HyoungMin Kim¹,
Myeong Jun Song¹, Ahn Choi¹, Daeho Kim¹, SooYoung Kim¹, Eun-Bong Kim¹,
David Wang¹, Shinhaeng Kang¹, Yuhwan Ro¹, Seungwoo Seo¹, JoonHo Song¹,
Jayoun Youn¹, Kyomin Sohn¹, Nam Sung Kim¹

¹Samsung Electronics, Hwaseong, Korea

²Samsung Electronics, San Jose, CA

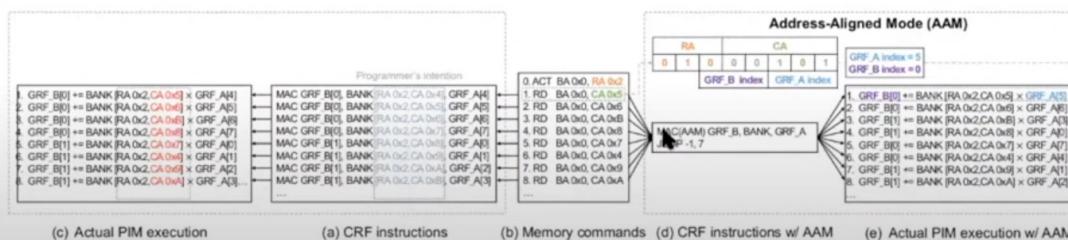
³Samsung Electronics, Suwon, Korea

Longer Lecture on HBM-PIM/FIMDRAM



FIMDRAM: Instruction Ordering

- One challenge is that DRAM commands may be re-ordered, and using fences is costly performance-wise
- Solution: **Address Aligned Mode (AAM)**
 - 8 MAC operations with 2 PIM instructions



A screenshot of a Zoom video call. The slide is titled "FIMDRAM: Instruction Ordering" and contains the diagrams and text from the previous slide. The video player controls are visible at the bottom, along with the Zoom logo. The URL "https://youtu.be/67oMjDmeUvo" is displayed at the bottom of the slide.

PIM Course: Lecture 5: Real-world PIM: Samsung HBM-PIM - Fall 2022



Subscribed

16



Share

Clip

Save



926 views 4 months ago Livestream - P&S Data-Centric Architectures: Fundamentally Improving Performance and Energy (Fall 2022)

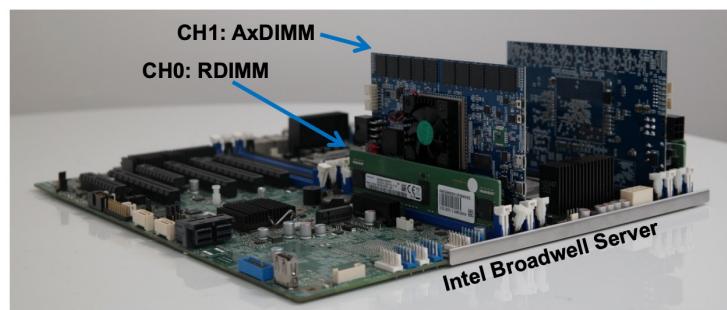
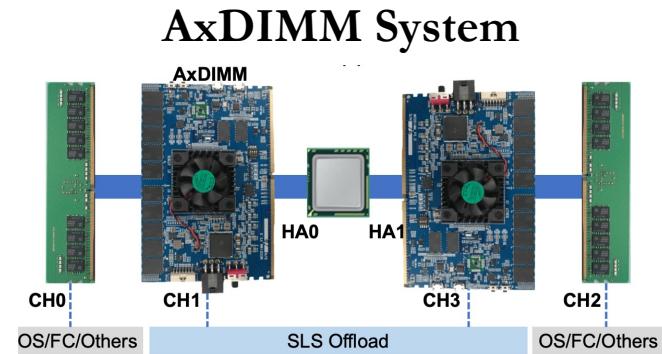
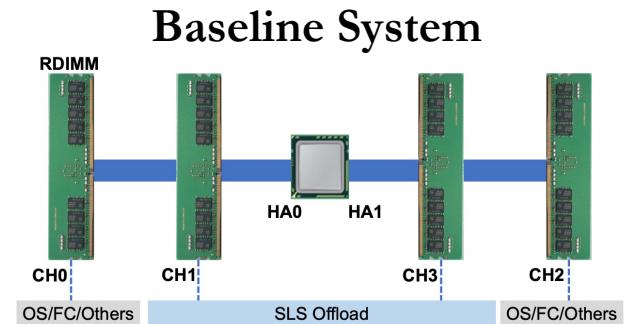
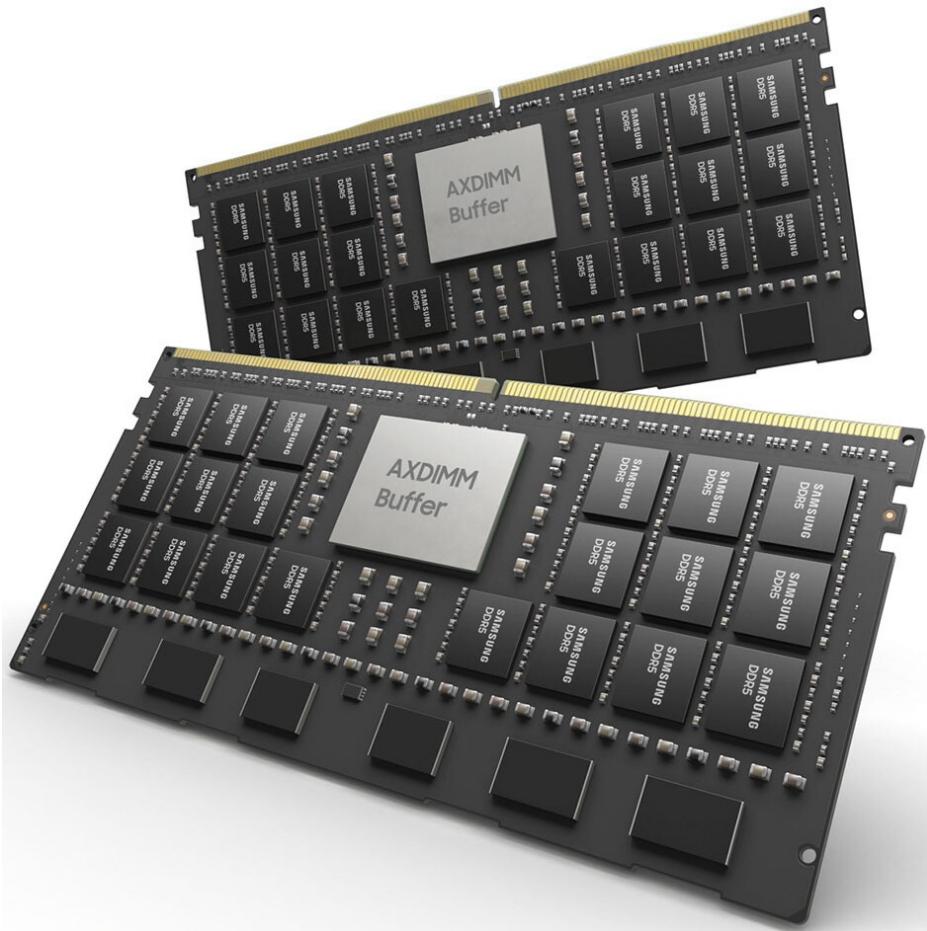
Projects & Seminars, ETH Zürich, Fall 2022

Data-Centric Architectures: Fundamentally Improving Performance and Energy

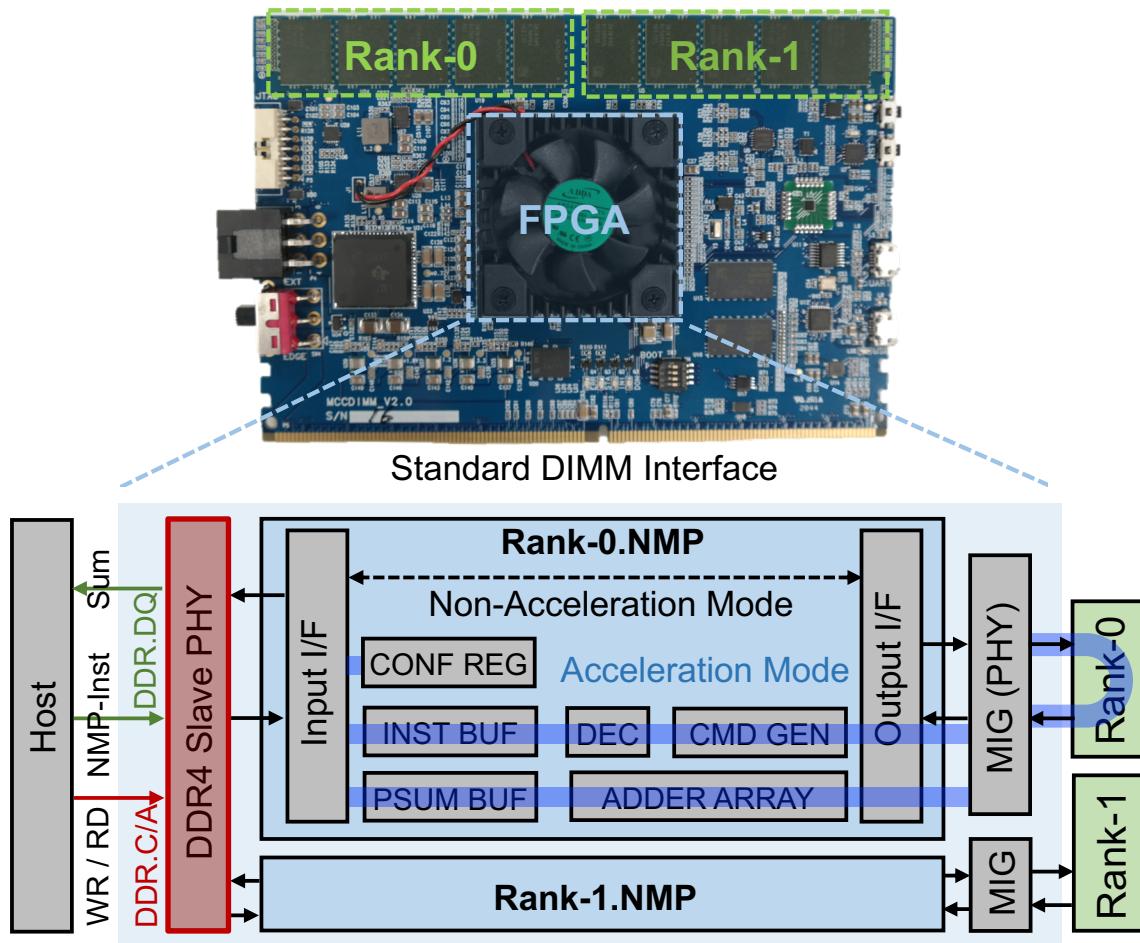
(https://safari.ethz.ch/projects_and_s...) Show more

Samsung AxDIMM (2021)

- DIMM-based PIM
 - DLRM recommendation system

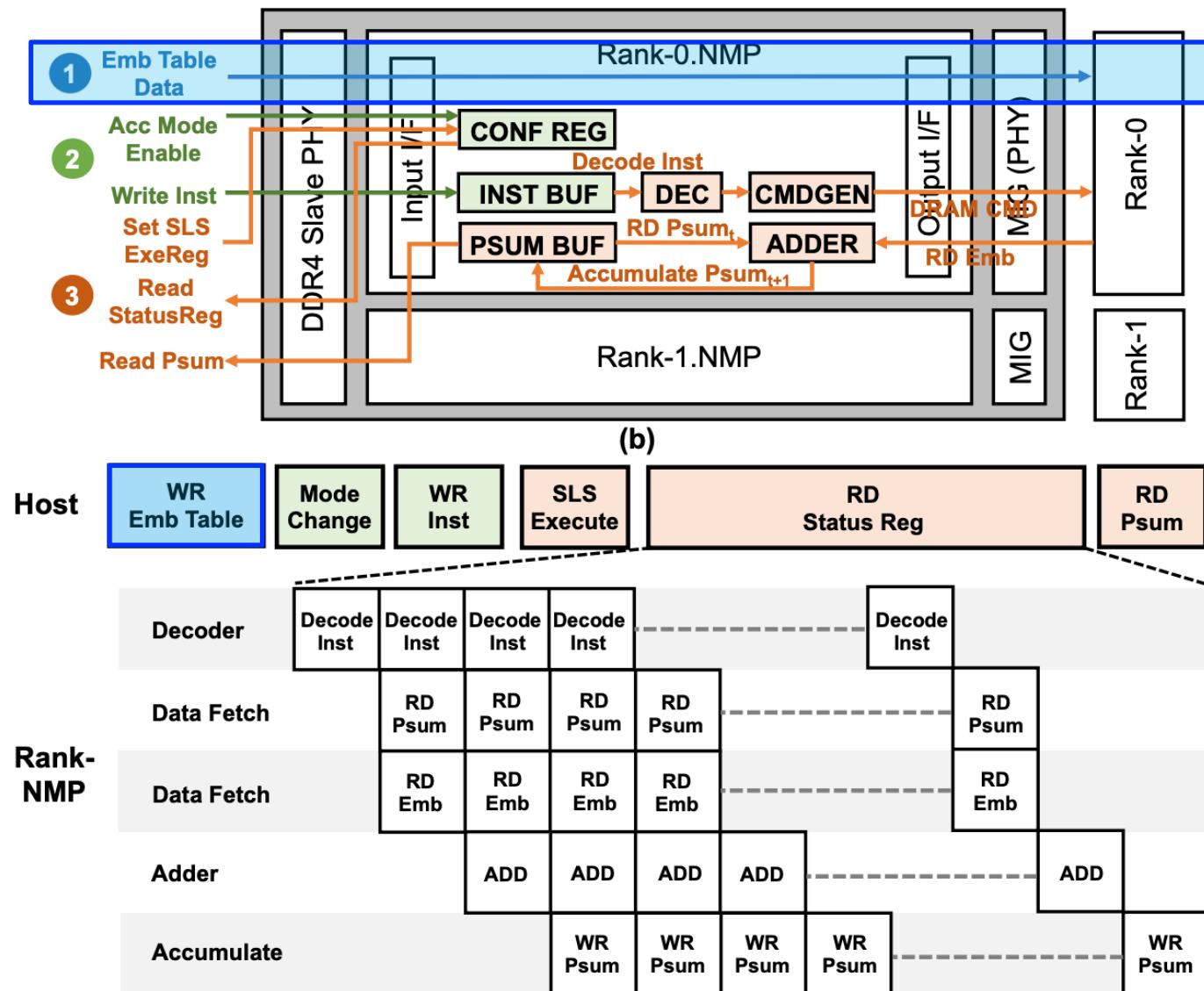


AxDIMM Design: Hardware Architecture



DDR4 slave PHY receives **DRAM commands and NMP instructions**
(via DQ pins) from the host side

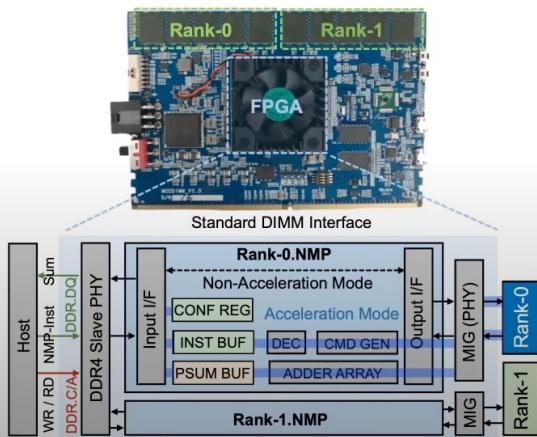
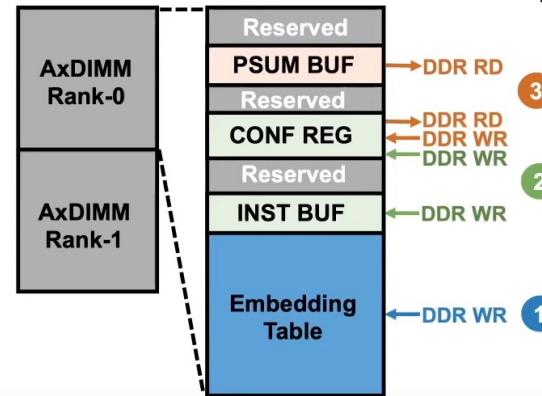
AxDIMM Design: Execution Flow



Longer Lecture on AxDIMM

AxDIMM Design: Address Map

- Memory map of AxDIMM



PIM Course: Lecture 7: Real-world PIM: Samsung AxDIMM - Fall 2022



Onur Mutlu Lectures

32.4K subscribers



Subscribed

21



Share



Clip

Save



846 views 4 months ago Livestream - P&S Data-Centric Architectures: Fundamentally Improving Performance and Energy (Fall 2022)

Projects & Seminars, ETH Zürich, Fall 2022

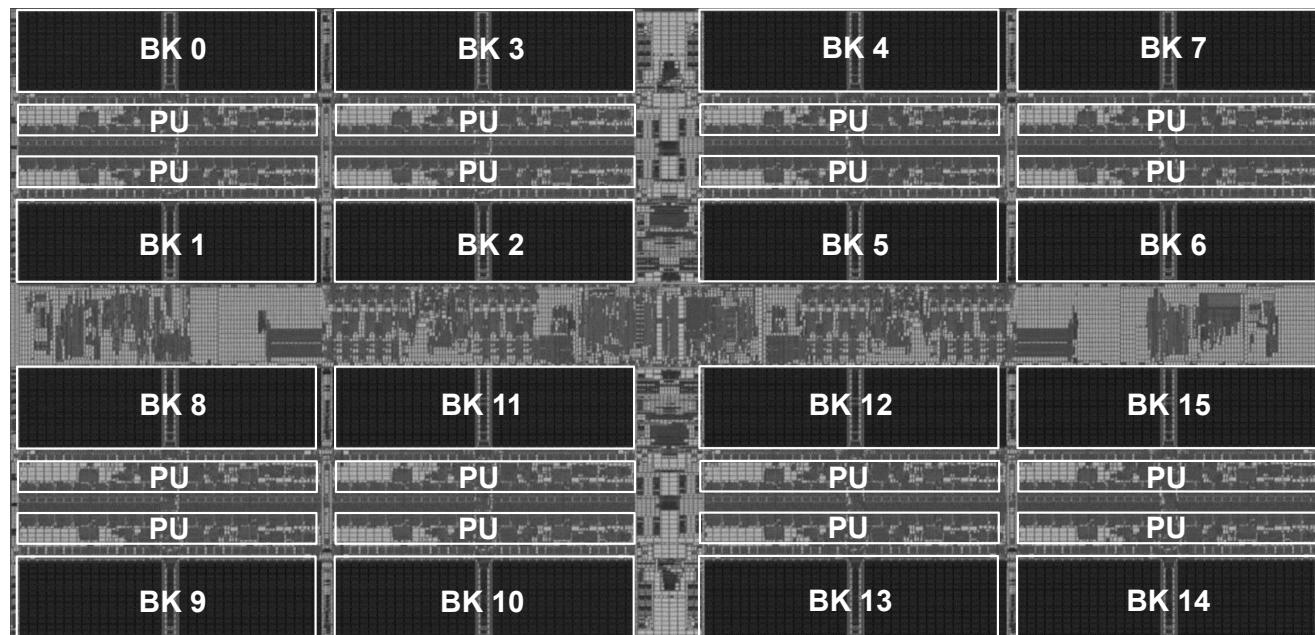
Data-Centric Architectures: Fundamentally Improving Performance and Energy

(https://safari.ethz.ch/projects_and_s...) Show more

SK Hynix AiM: Chip Implementation (2022)

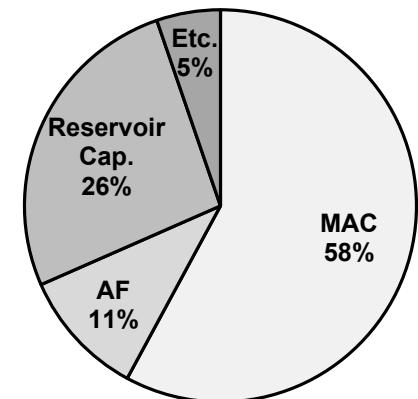
- 4 Gb AiM die with 16 processing units (PUs)

AiM Die Photograph



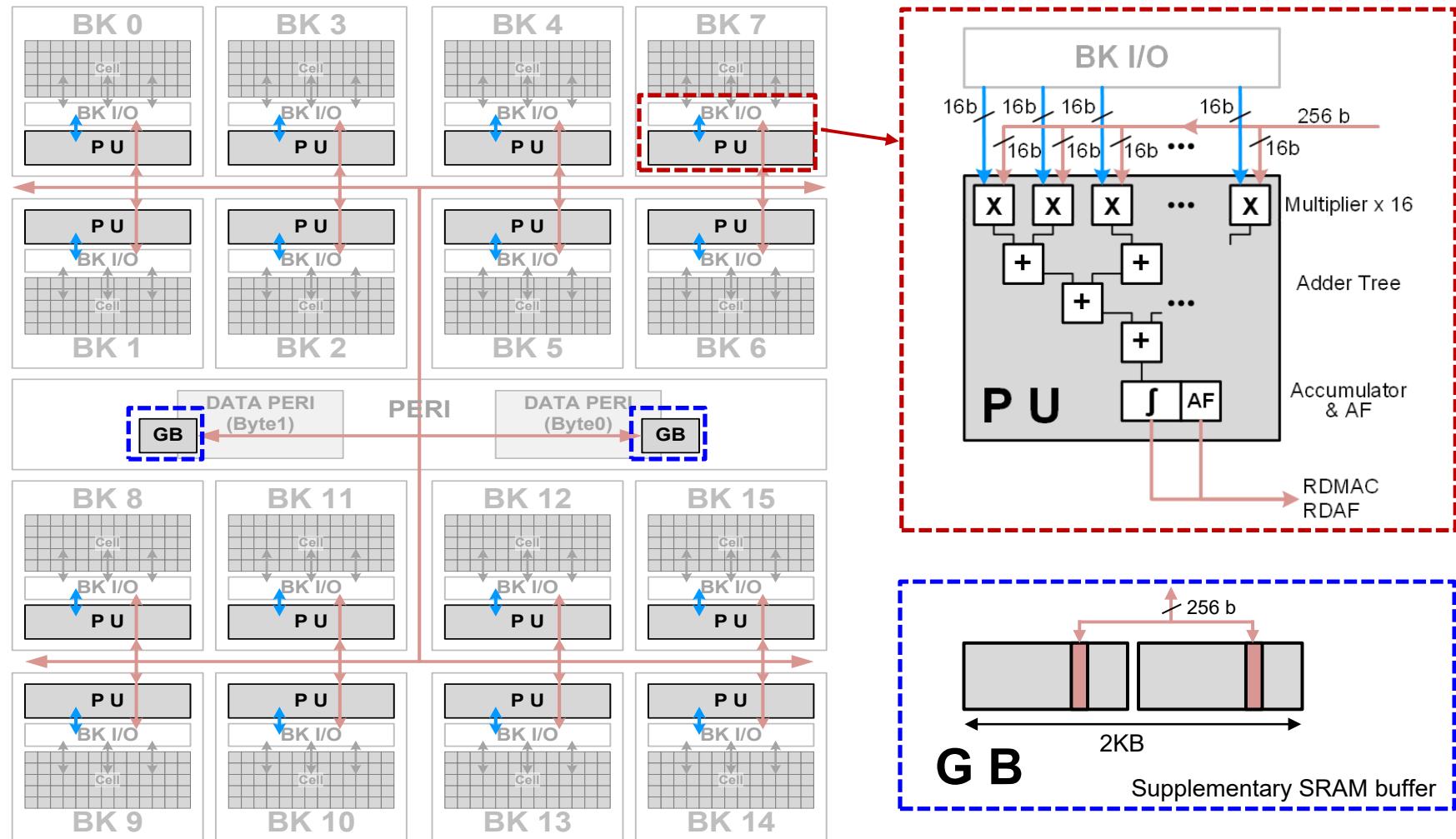
1 Process Unit (PU) Area

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



SK Hynix AiM: System Organization (2022)

GDDR6-based AiM architecture



Lecture on Accelerator-in-Memory

AiM: MAC Circuit

- 16 multipliers, adder tree, and accumulator
 - Bfloat16 (BF16) format

The diagram illustrates the AiM MAC Circuit architecture. At the top, 16 multipliers (MUL) receive weights (W_i) and vectors (V_j) as inputs. The results of these multiplications are then processed by an EXP Comp. block and a MAX EX block. The outputs of the EXP Comp. block are fed into an adder tree consisting of 16 MAU (Multi-Adder Unit) blocks. The adder tree's output is then combined with a bias and the result of the MAX EX block to produce the final MAC Result. Below this, the Bfloat16 (BF16) Format is shown, divided into Sign (S) 1b, Exponent (EX) 8b, and Mantissa (MA) 7b. A detailed view at the bottom shows a 16-bank weight vector (Weight) being multiplied by a vector (Vector) to produce a Multiplier Output. This output is then processed by an Adder Tree Output (MA SUM, MAX EX) block and an Accumulation block to produce the final MAC Result.

PIM Course: Lecture 6: Real-world PIM: SK Hynix AiM (Spring 2023)



Onur Mutlu Lectures
33.4K subscribers

Subscribed

18



Share

Clip

Save



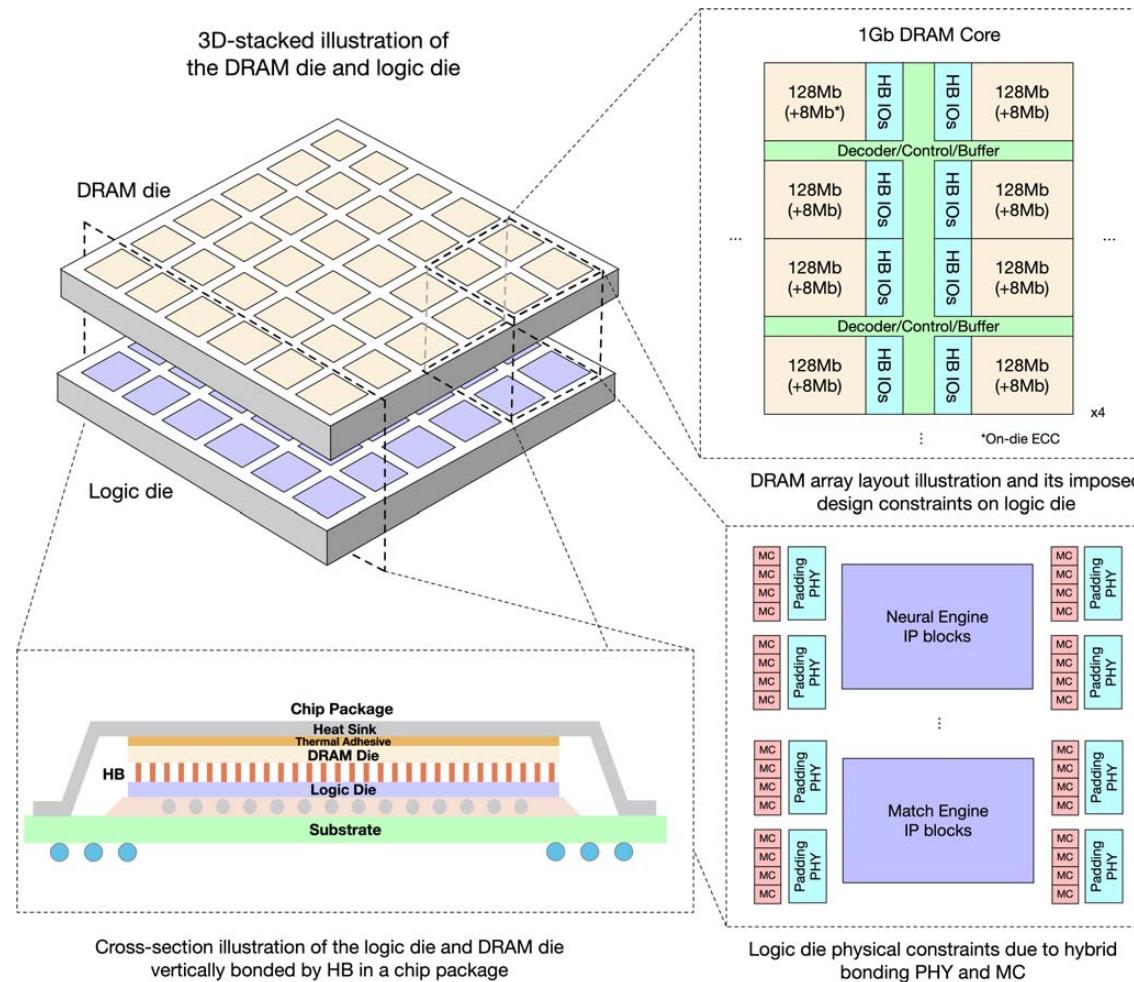
569 views 1 month ago Livestream - Data-Centric Architectures: Fundamentally Improving Performance and Energy (Spring 2023)

Projects & Seminars, ETH Zürich, Spring 2023

Data-Centric Architectures: Fundamentally Improving Performance and Energy

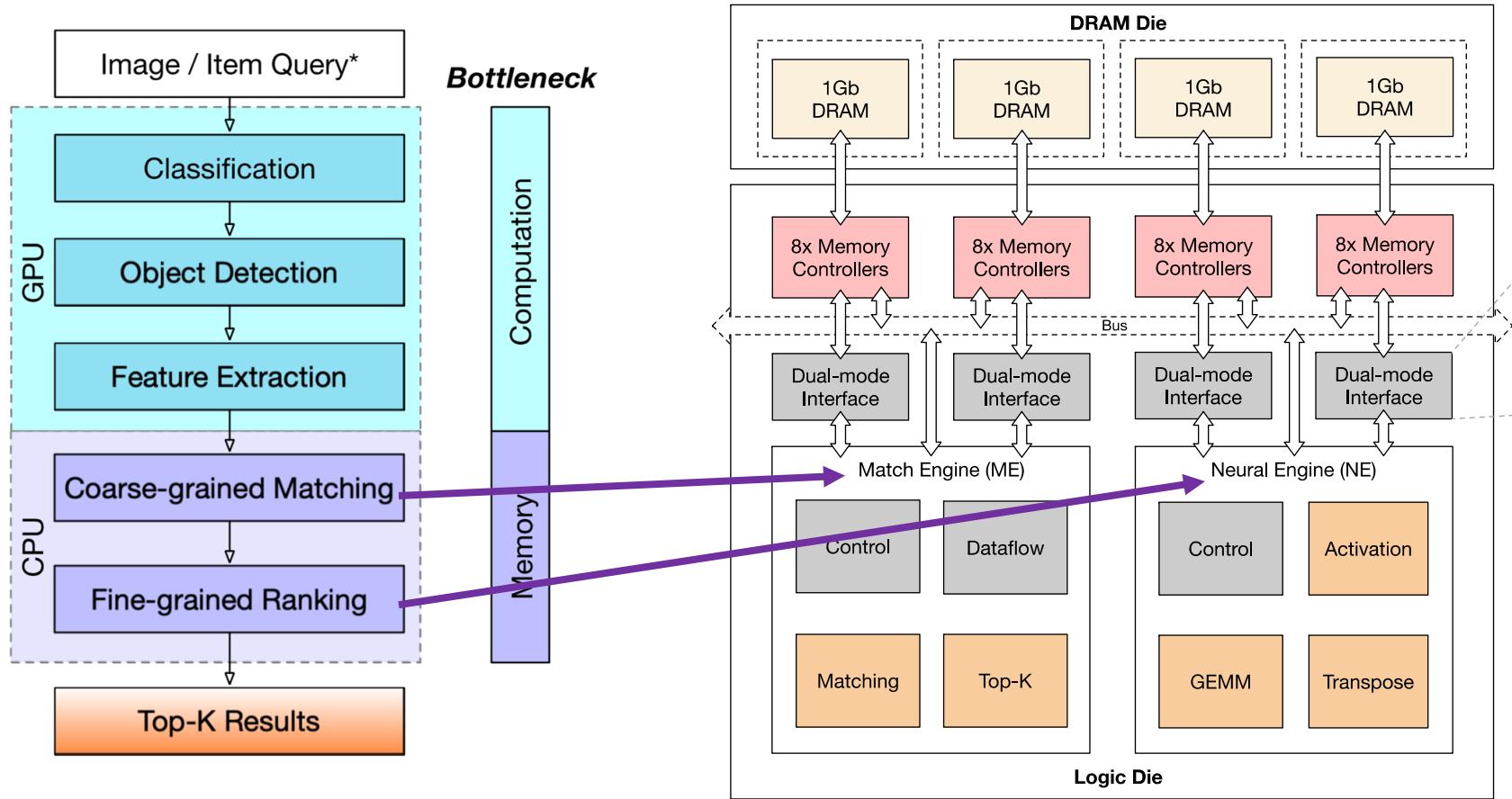
Alibaba HB-PNM: Overall Architecture (2022)

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

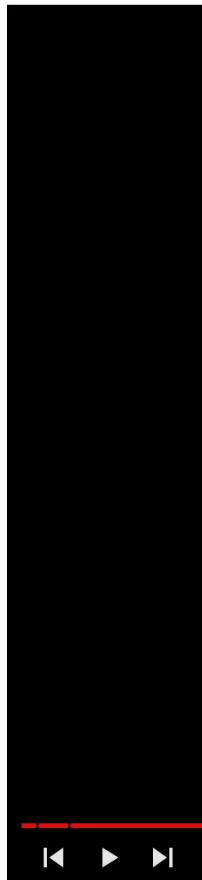


HB-PNM: Overall Architecture

- Match engine and neural engine for matching and ranking in a recommendation system

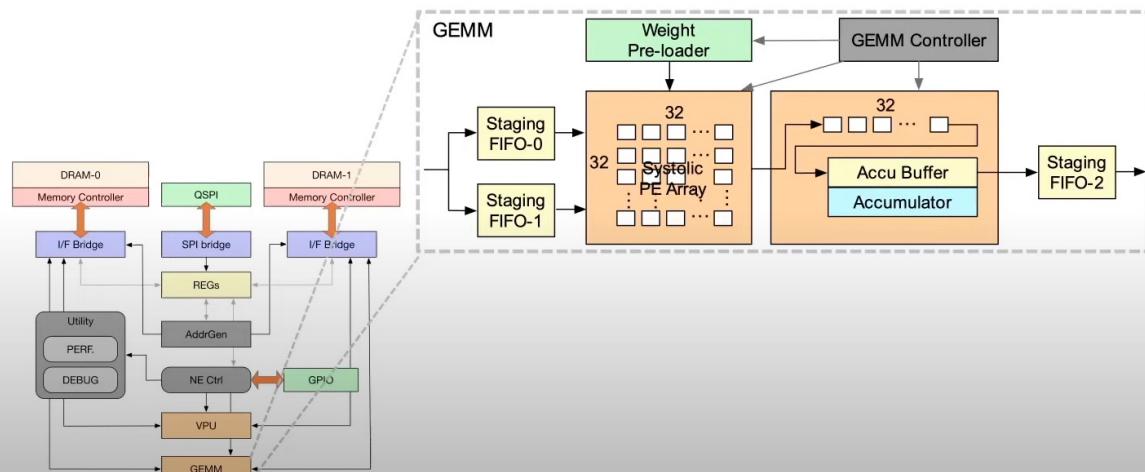


Longer Lecture on HB-PNM



Neural Engine: GEMM

- 32x32 INT8 fully-pipelined **systolic array**
 - Partial sums accumulated in INT32 accumulator



PIM Course: Lecture 8: Real-world PIM: Alibaba HB-PNM - Fall 2022



Onur Mutlu Lectures

32.4K subscribers

Subscribed

17



Share

Clip

Save



438 views 4 months ago Livestream - P&S Data-Centric Architectures: Fundamentally Improving Performance and Energy (Fall 2022)

Projects & Seminars, ETH Zürich, Fall 2022

Data-Centric Architectures: Fundamentally Improving Performance and Energy

(https://safari.ethz.ch/projects_and_s...) Show more

FPGA-based Processing Near Memory

- Gagandeep Singh, Mohammed Alser, Damla Senol Cali, Dionysios Diamantopoulos, Juan Gómez-Luna, Henk Corporaal, and Onur Mutlu,
"FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications"
IEEE Micro (IEEE MICRO), to appear, 2021.

FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications

Gagandeep Singh[◊] Mohammed Alser[◊] Damla Senol Cali[✉]

Dionysios Diamantopoulos[▽] Juan Gómez-Luna[◊]

Henk Corporaal^{*} Onur Mutlu^{◊✉}

[◊]*ETH Zürich* [✉]*Carnegie Mellon University*

^{*}*Eindhoven University of Technology* [▽]*IBM Research Europe*

Why In-Memory Computation Today?

- **Huge problems with Memory Technology**
 - Memory technology scaling is not going well (e.g., RowHammer)
 - Many scaling issues demand intelligence in memory
- **Huge demand from Applications & Systems**
 - Data access bottleneck
 - Energy & power bottlenecks
 - Data movement energy dominates computation energy
 - Need all at the same time: performance, energy, sustainability
 - We can improve all metrics by minimizing data movement
- **Designs are squeezed in the middle**

We Need to Think Differently
from the Past Approaches

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 near Memory
 - Processing using Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

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)]

A PIM Taxonomy

- **Nature (of computation)**
 - **Using**: Use operational properties of memory structures
 - **Near**: Add logic close to memory structures
 - **Technology**
 - Flash, DRAM, SRAM, RRAM, MRAM, FeRAM, PCM, 3D, ...
 - **Location**
 - Sensor, Cold Storage, Hard Disk, SSD, Main Memory, Cache, Register File, Memory Controller, Interconnect, ...
 - A tuple of the three determines “PIM type”
 - One can combine multiple “PIM types” in a system
-

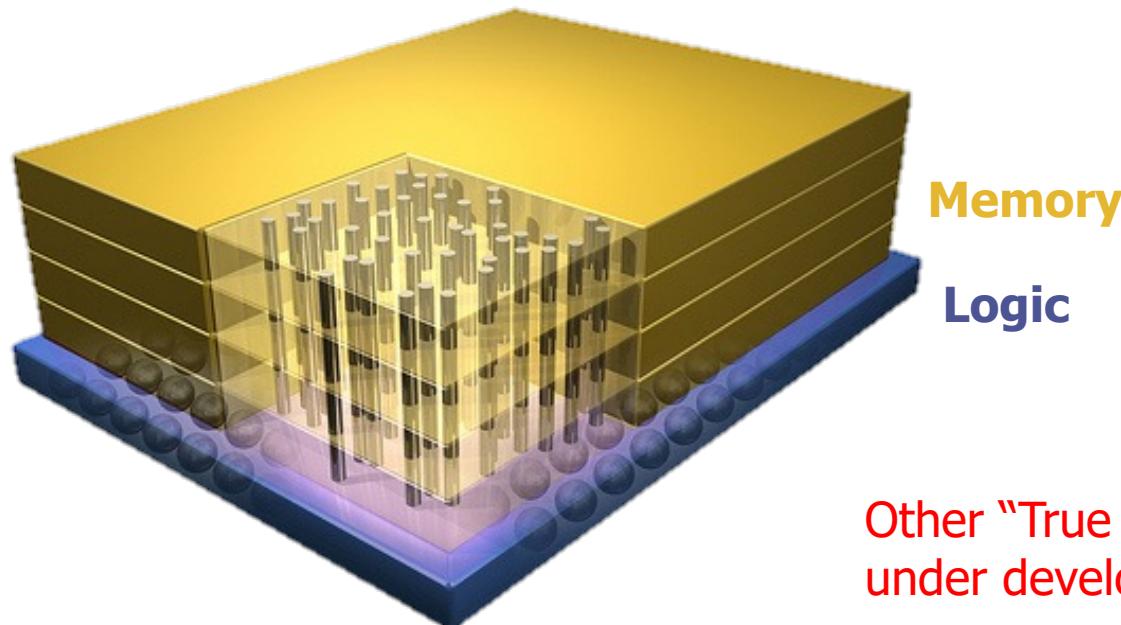
Processing in Memory: Two Approaches

- 1. Processing near Memory**
2. Processing using Memory

Opportunity: 3D-Stacked Logic+Memory



Hybrid Memory Cube
C O N S O R T I U M



Other “True 3D” technologies
under development

DRAM Landscape (circa 2015)

<i>Segment</i>	<i>DRAM Standards & Architectures</i>
Commodity	DDR3 (2007) [14]; DDR4 (2012) [18]
Low-Power	LPDDR3 (2012) [17]; LPDDR4 (2014) [20]
Graphics	GDDR5 (2009) [15]
Performance	eDRAM [28], [32]; RLDRAM3 (2011) [29]
3D-Stacked	WIO (2011) [16]; WIO2 (2014) [21]; MCDRAM (2015) [13]; HBM (2013) [19]; HMC1.0 (2013) [10]; HMC1.1 (2014) [11]
Academic	SBA/SSA (2010) [38]; Staged Reads (2012) [8]; RAIDR (2012) [27]; SALP (2012) [24]; TL-DRAM (2013) [26]; RowClone (2013) [37]; Half-DRAM (2014) [39]; Row-Buffer Decoupling (2014) [33]; SARP (2014) [6]; AL-DRAM (2015) [25]

Table 1. Landscape of DRAM-based memory

Kim+, "Ramulator: A Flexible and Extensible DRAM Simulator", IEEE CAL 2015.

Several Questions in 3D-Stacked PIM

- What are the performance and energy benefits of using 3D-stacked memory as a coarse-grained accelerator?
 - By changing the entire system
 - By performing simple function offloading

- What is the minimal processing-in-memory support we can provide?
 - With minimal changes to system and programming

Another Example: In-Memory Graph Processing

- Large graphs are everywhere (circa 2015)



36 Million
Wikipedia Pages



1.4 Billion
Facebook Users

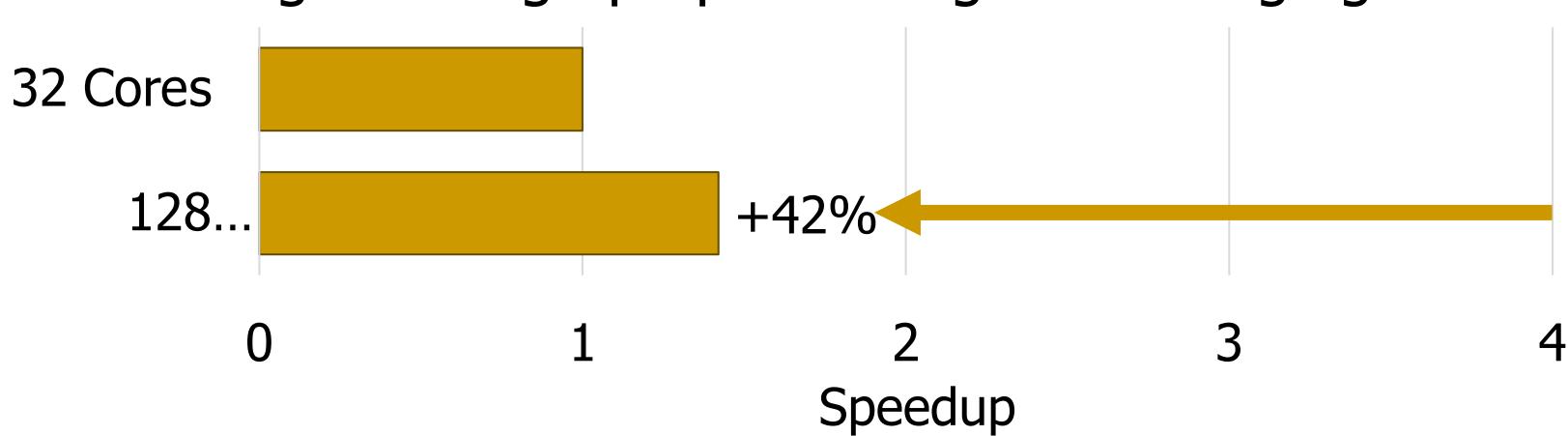


300 Million
Twitter Users



30 Billion
Instagram Photos

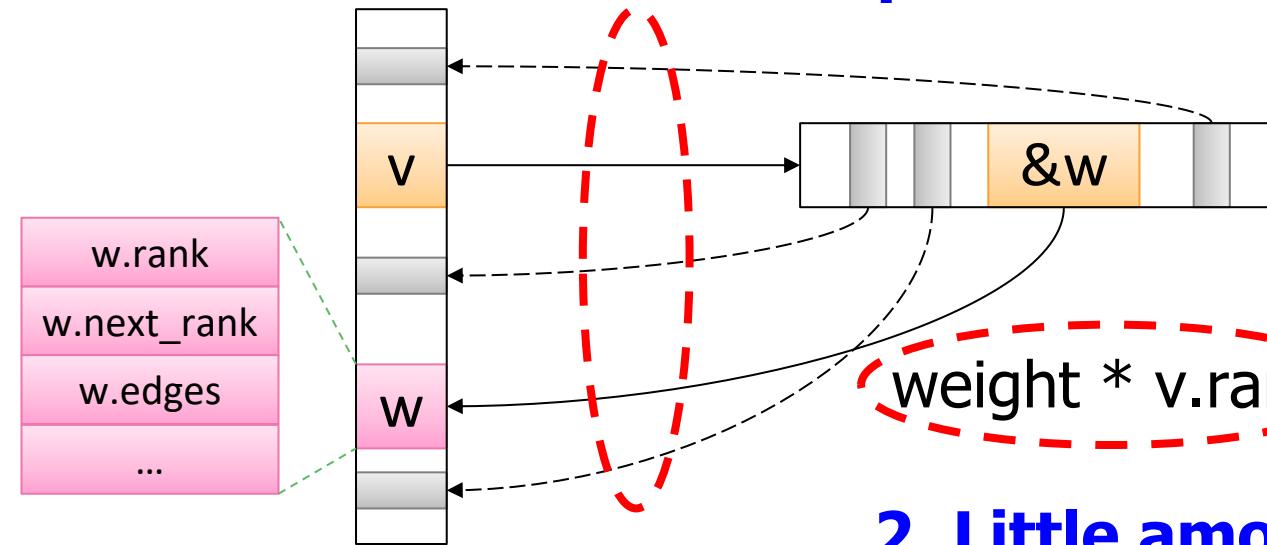
- Scalable large-scale graph processing is challenging



Key Bottlenecks in Graph Processing

```
for (v: graph.vertices) {  
    for (w: v.successors) {  
        w.next_rank += weight * v.rank;  
    }  
}
```

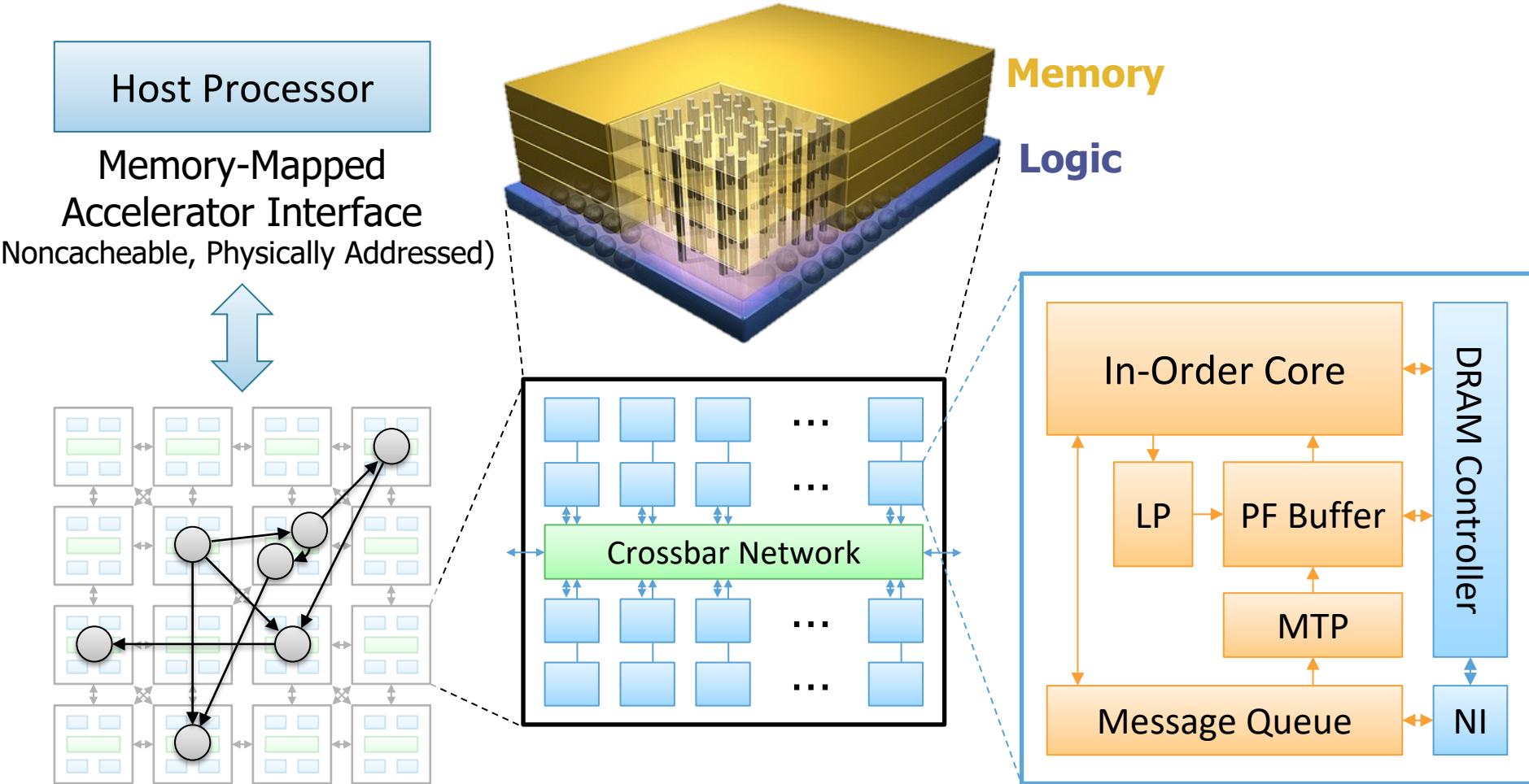
1. Frequent random memory accesses



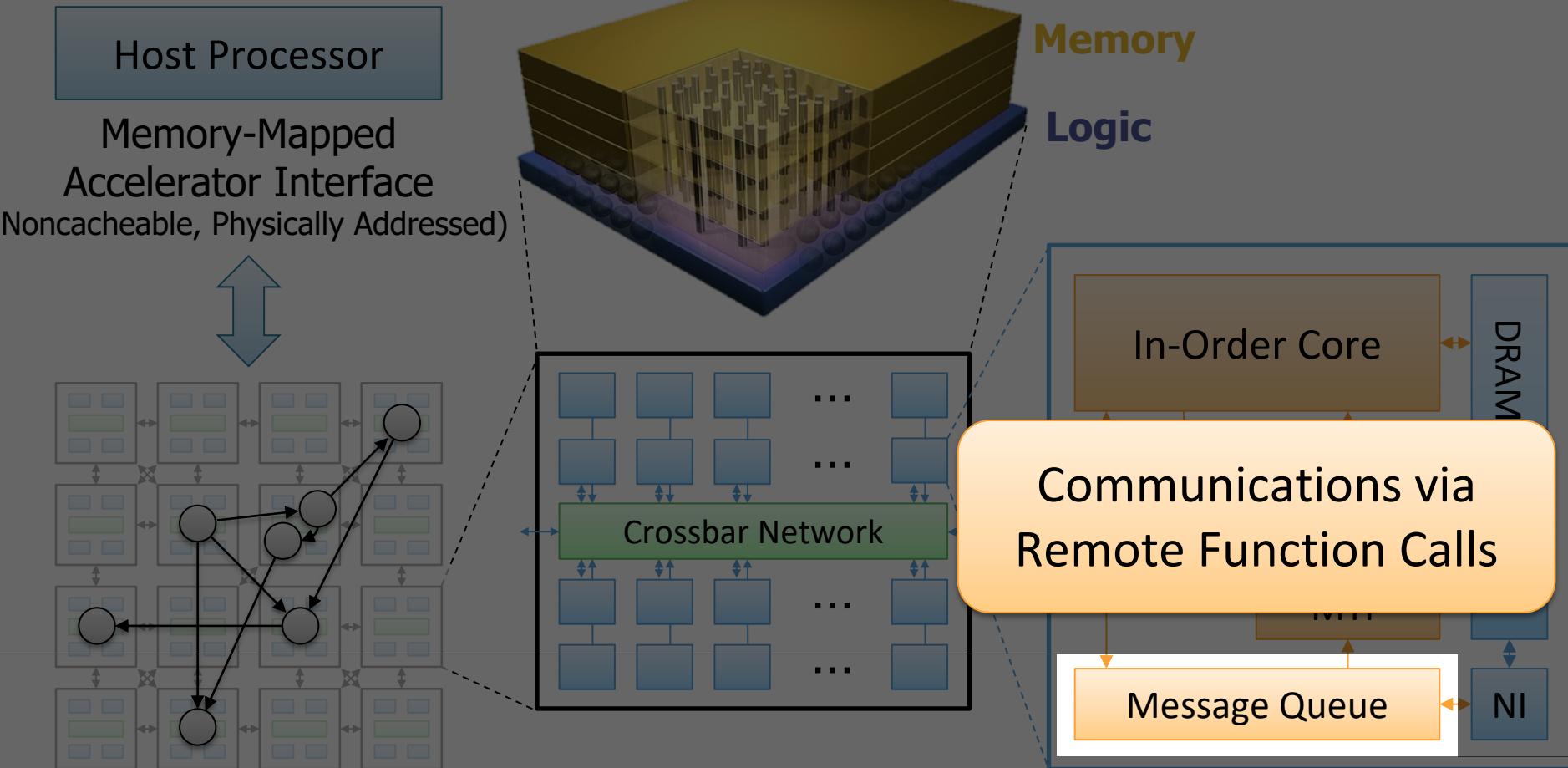
2. Little amount of computation

Tesseract System for Graph Processing

Interconnected set of 3D-stacked memory+logic chips with simple cores

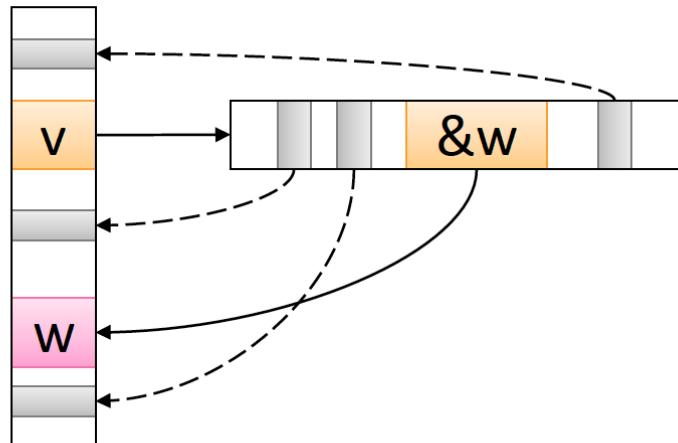


Tesseract System for Graph Processing



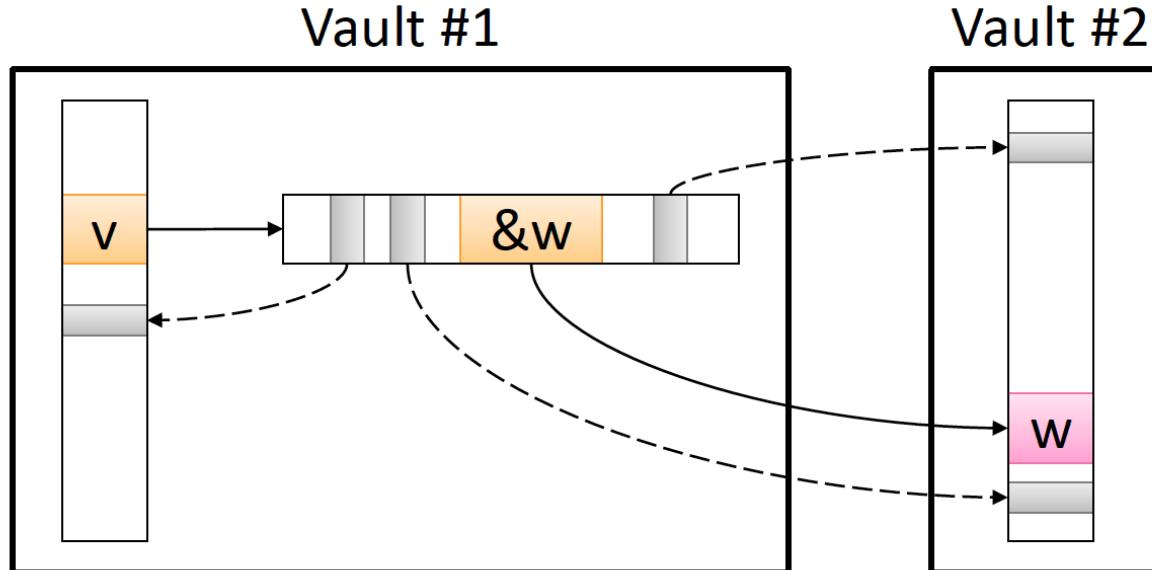
Communications In Tesseract (I)

```
for (v: graph.vertices) {  
    for (w: v.successors) {  
        w.next_rank += weight * v.rank;  
    }  
}
```



Communications In Tesseract (II)

```
for (v: graph.vertices) {  
    for (w: v.successors) {  
        w.next_rank += weight * v.rank;  
    }  
}
```

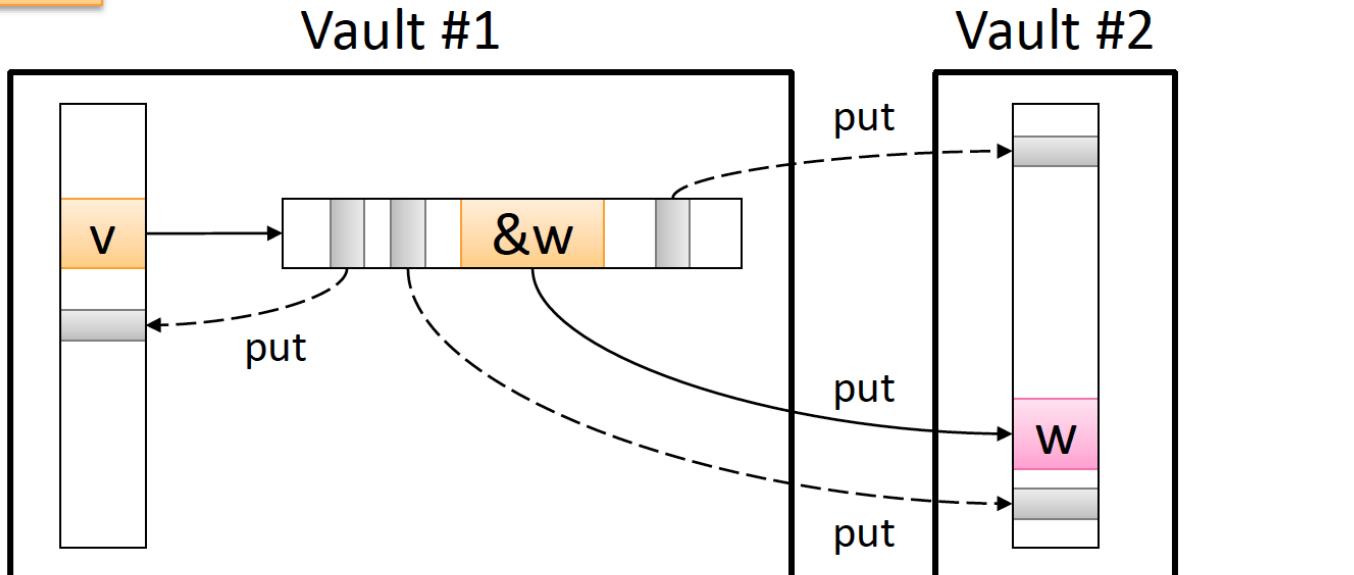


Communications In Tesseract (III)

```
for (v: graph.vertices) {  
    for (w: v.successors) {  
        put(w.id, function() { w.next_rank += weight * v.rank; });  
    }  
}  
barrier();
```

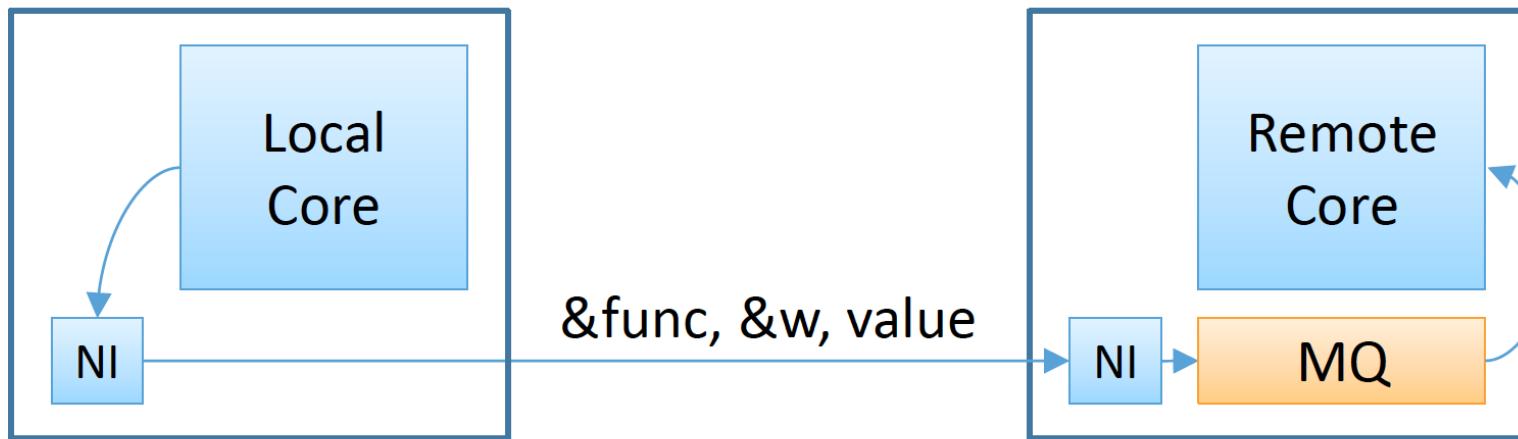
Non-blocking Remote Function Call

Can be **delayed**
until the nearest barrier



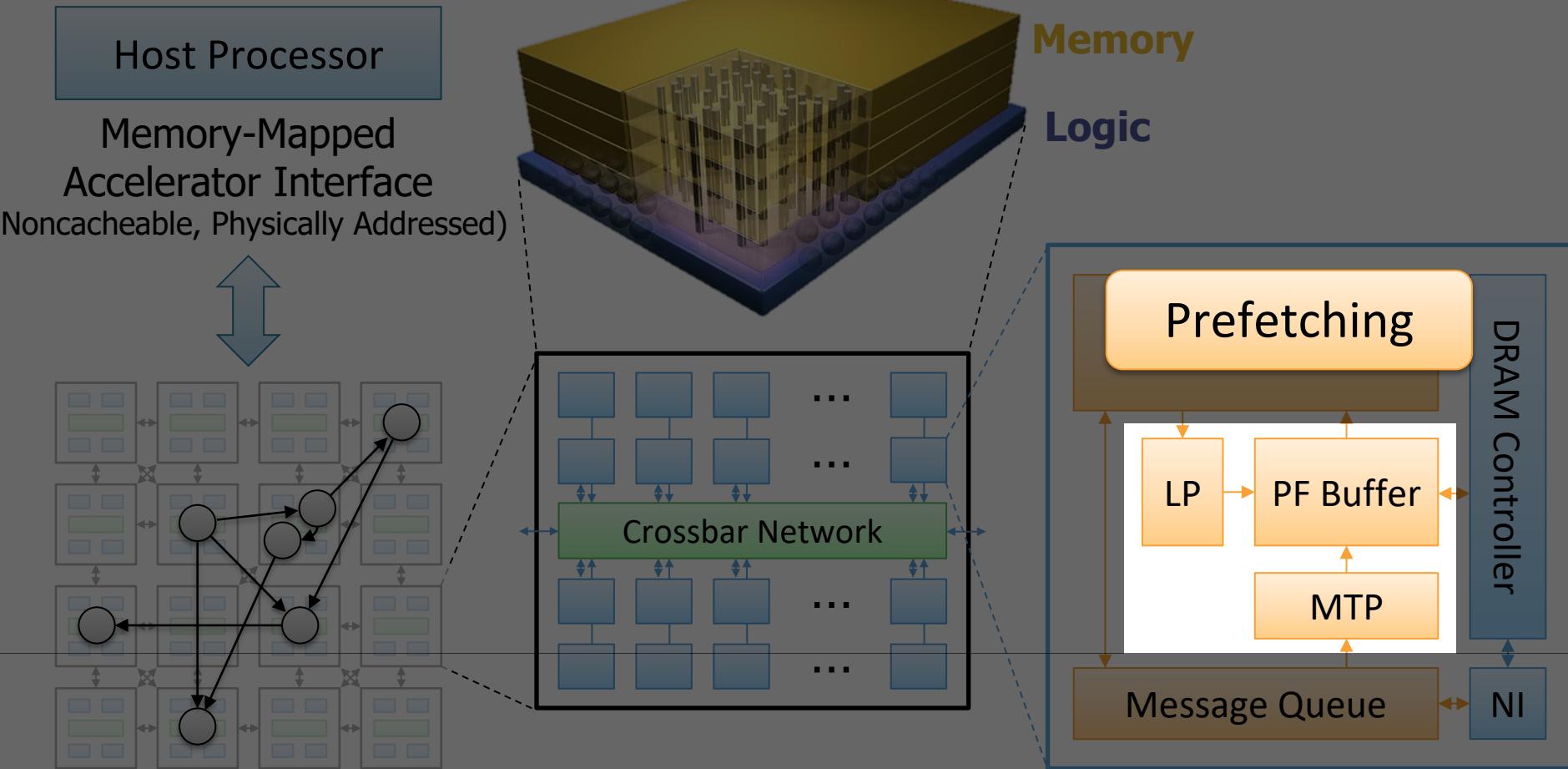
Remote Function Call (Non-Blocking)

1. Send function address & args to the remote core
2. Store the incoming message to the message queue
3. Flush the message queue when it is full or a synchronization barrier is reached



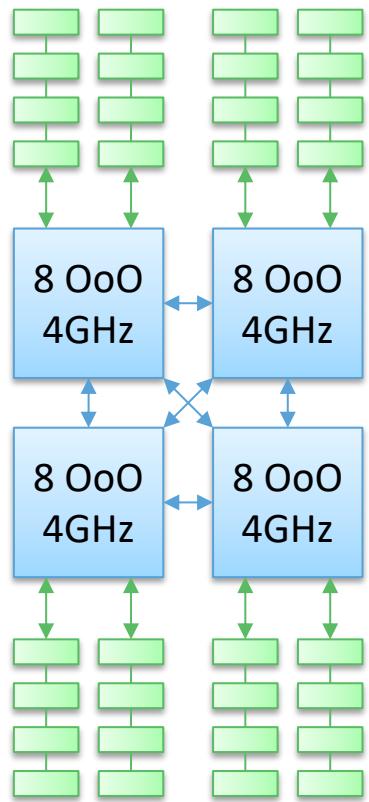
```
put(w.id, function() { w.next_rank += value; })
```

Tesseract System for Graph Processing

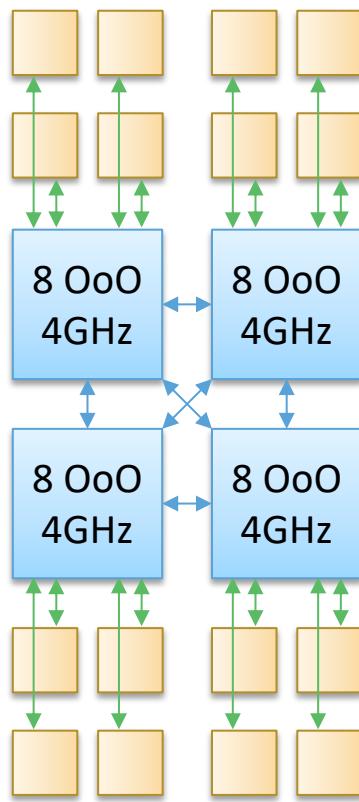


Evaluated Systems

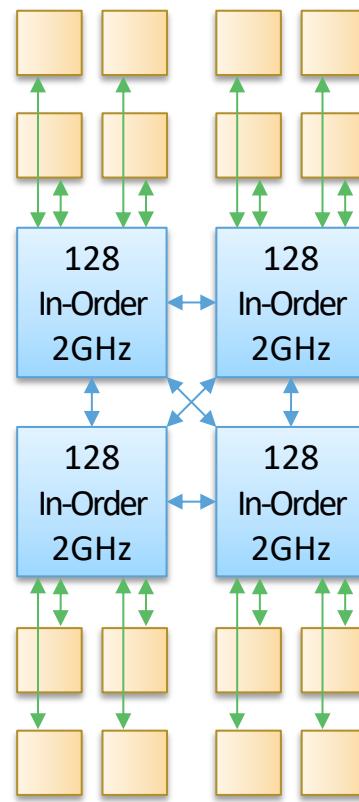
DDR3-OoO



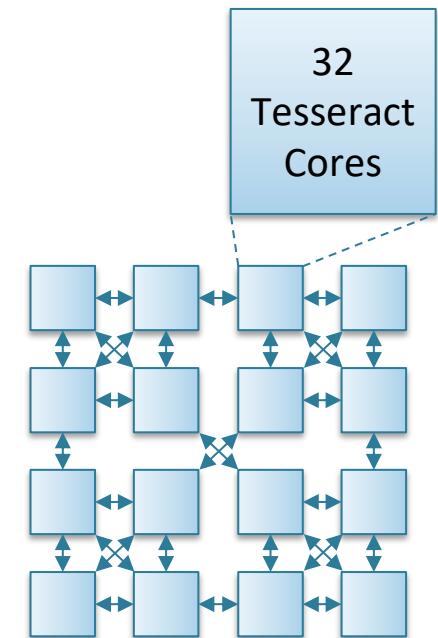
HMC-OoO



HMC-MC



Tesseract



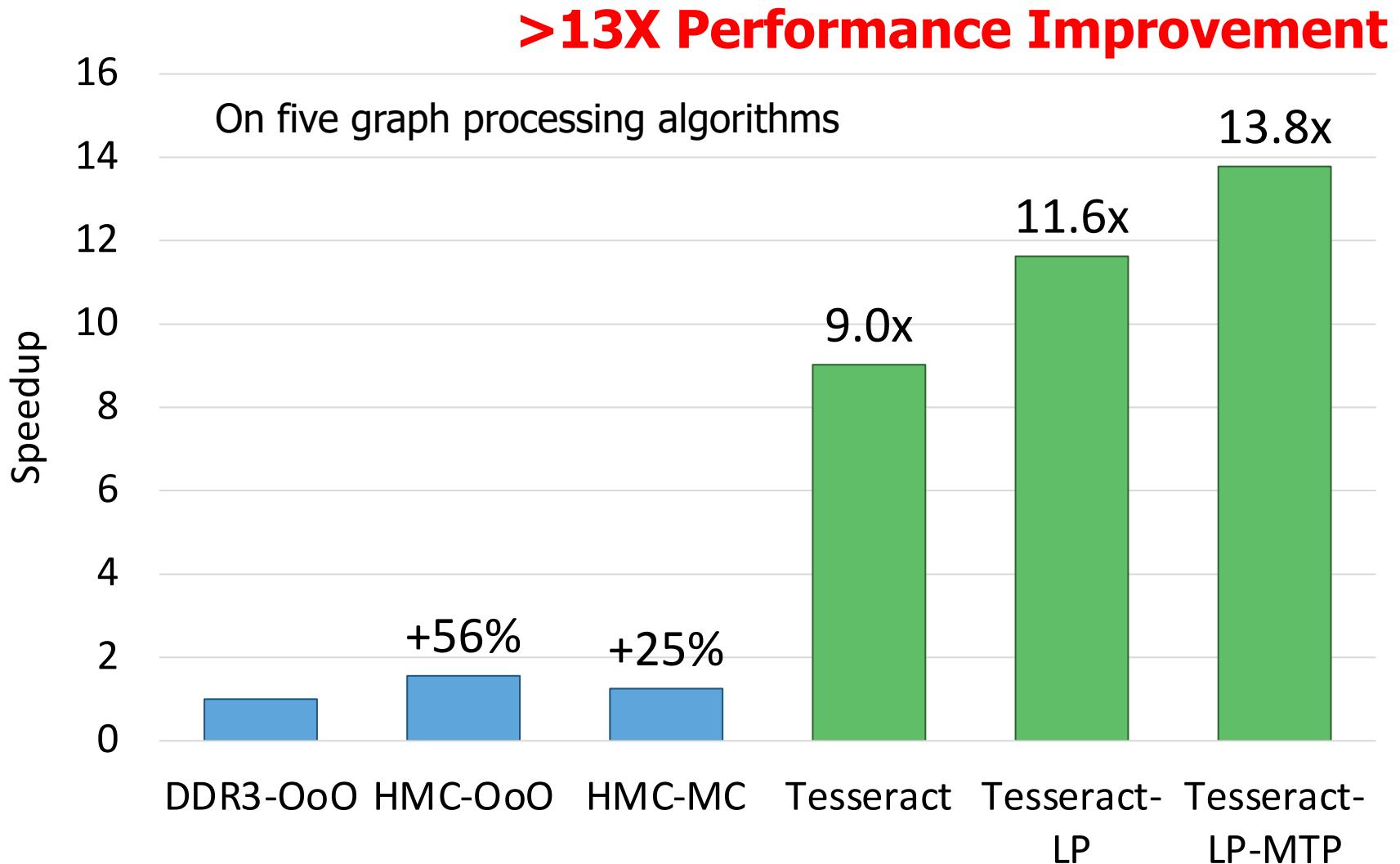
102.4GB/s

640GB/s

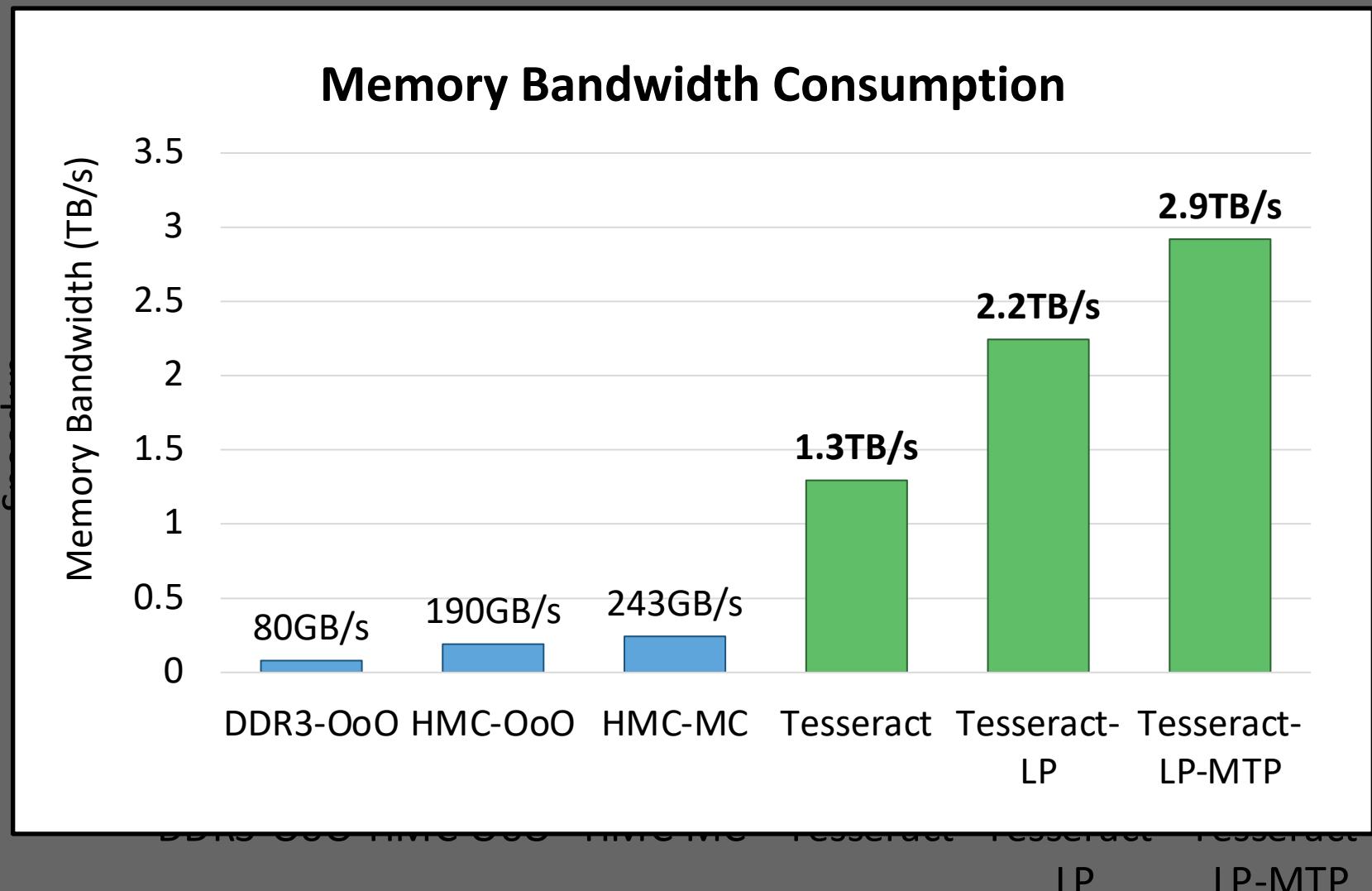
640GB/s

8TB/s

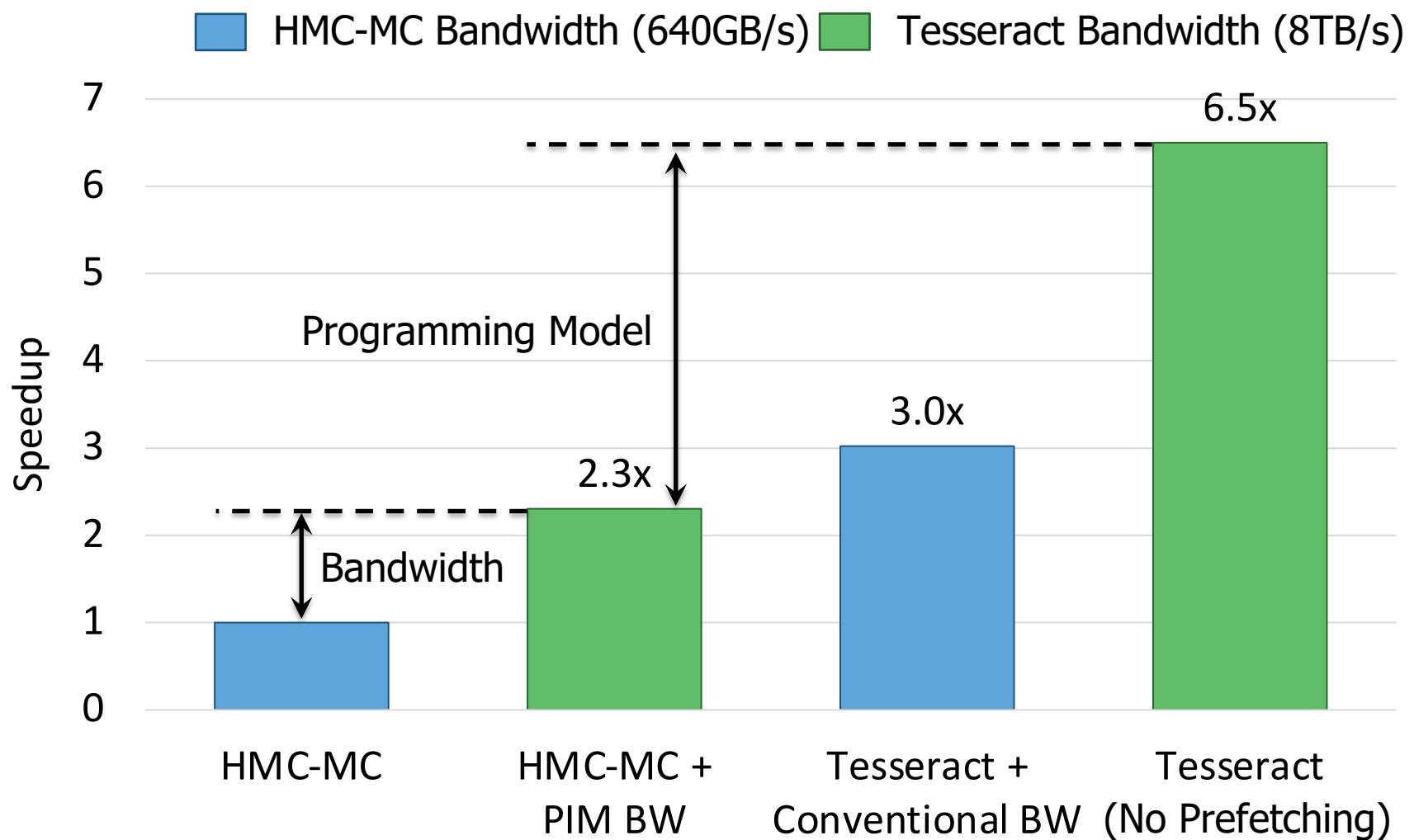
Tesseract Graph Processing Performance



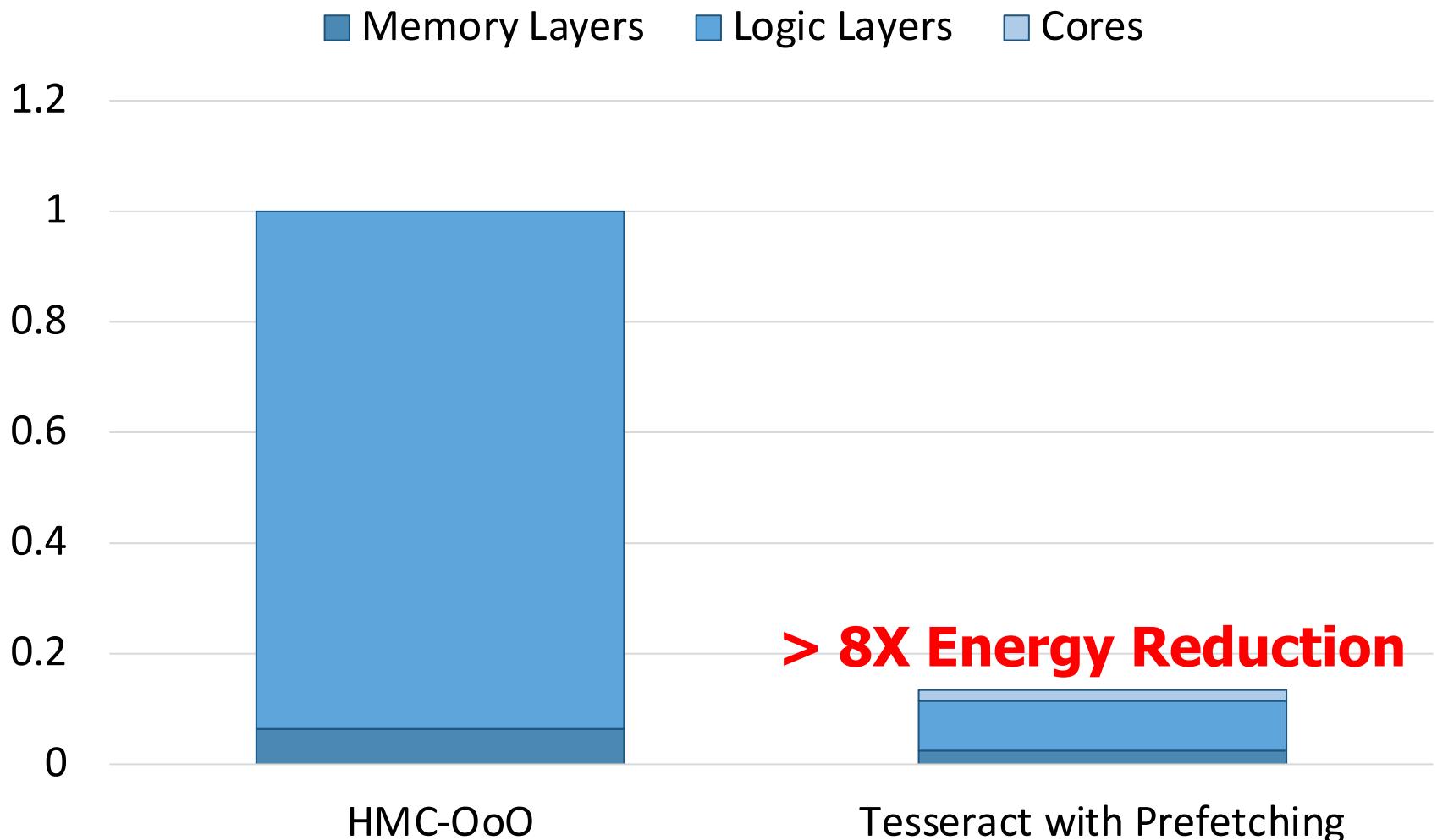
Tesseract Graph Processing Performance



Effect of Bandwidth & Programming Model



Tesseract Graph Processing System Energy



Tesseract: Advantages & Disadvantages

- **Advantages**
 - + Specialized graph processing accelerator using PIM
 - + Large system performance and energy benefits
 - + Takes advantage of 3D stacking for an important workload
 - + More general than just graph processing

- **Disadvantages**
 - Changes a lot in the system
 - New programming model
 - Specialized Tesseract cores for graph processing
 - Cost
 - Scalability limited by off-chip links or graph partitioning

More on Tesseract

- Junwhan Ahn, Sungpack Hong, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,

"A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing"

Proceedings of the 42nd International Symposium on Computer Architecture (ISCA), Portland, OR, June 2015.

[Slides (pptx) (pdf)] [Lightning Session Slides (pptx) (pdf)]

Top Picks Honorable Mention by IEEE Micro.

Selected to the ISCA-50 25-Year Retrospective Issue covering 1996-2020 in 2023 (Retrospective (pdf) Full Issue).

A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing

Junwhan Ahn Sungpack Hong[§] Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi

junwhan@snu.ac.kr, sungpack.hong@oracle.com, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr

Seoul National University

[§]Oracle Labs

[†]Carnegie Mellon University

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 near Memory
 - Processing using Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

Several Questions in 3D-Stacked PIM

- What are the performance and energy benefits of using 3D-stacked memory as a coarse-grained accelerator?
 - By changing the entire system
 - By performing simple function offloading

- What is the minimal processing-in-memory support we can provide?
 - With minimal changes to system and programming

3D-Stacked PIM on Mobile Devices

- Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu,
"Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks"

Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Williamsburg, VA, USA, March 2018.

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

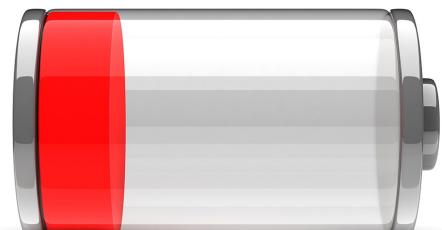
Amirali Boroumand¹ Saugata Ghose¹ Youngsok Kim²
Rachata Ausavarungnirun¹ Eric Shiu³ Rahul Thakur³ Daehyun Kim^{4,3}
Aki Kuusela³ Allan Knies³ Parthasarathy Ranganathan³ Onur Mutlu^{5,1}

Consumer Devices



Consumer devices are everywhere!

**Energy consumption is
a first-class concern in consumer devices**

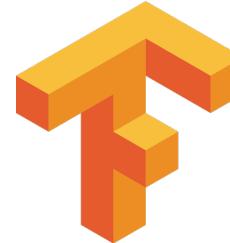


Four Important Workloads



Chrome

Google's web browser



TensorFlow Mobile
Google's machine learning
framework

VP9



Video Playback

Google's **video codec**

VP9

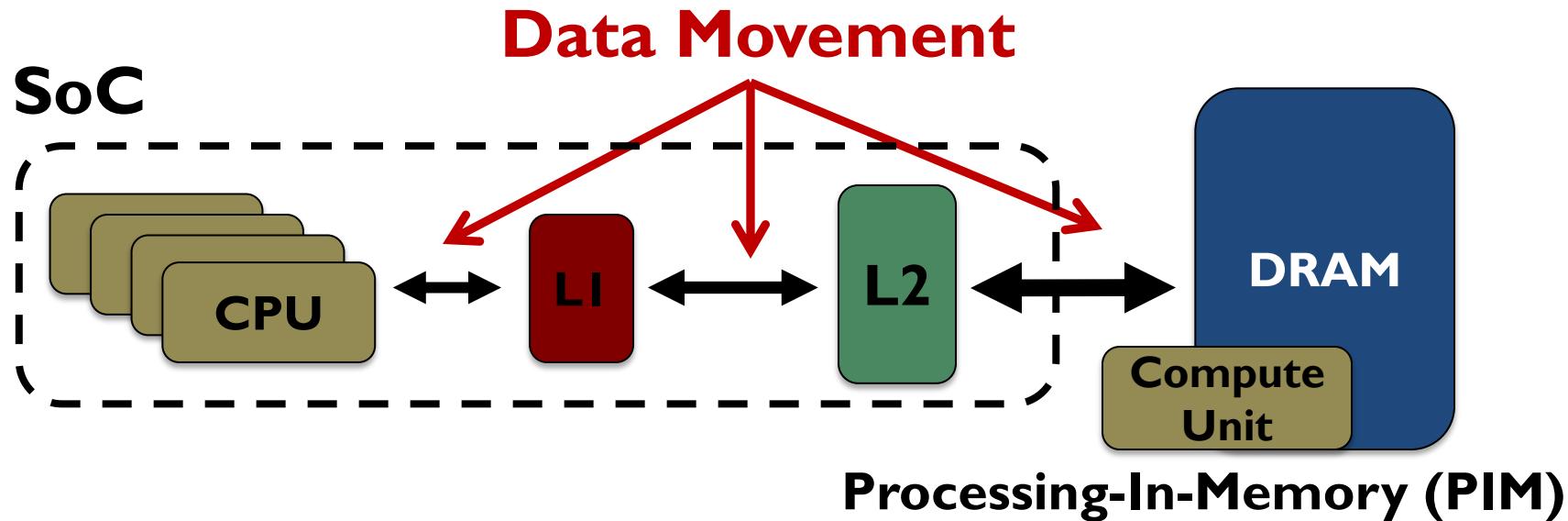


Video Capture

Google's **video codec**

Energy Cost of Data Movement

1st key observation: **62.7% of the total system energy is spent on data movement**



Potential solution: move computation **close to data**

Challenge: limited area and energy budget

Using PIM to Reduce Data Movement

2nd key observation: a significant fraction of the **data movement** often comes from **simple functions**

We can design lightweight logic to implement these simple functions in **memory**

**Small embedded
low-power core**



**Small fixed-function
accelerators**



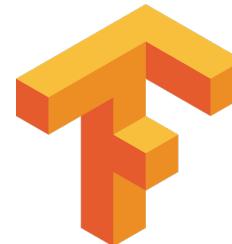
Offloading to PIM logic reduces energy and improves performance, on average, by 55.4% and 54.2%

Workload Analysis



Chrome

Google's web browser



TensorFlow Mobile

Google's machine learning
framework

VP9



Video Playback

Google's **video codec**

VP9

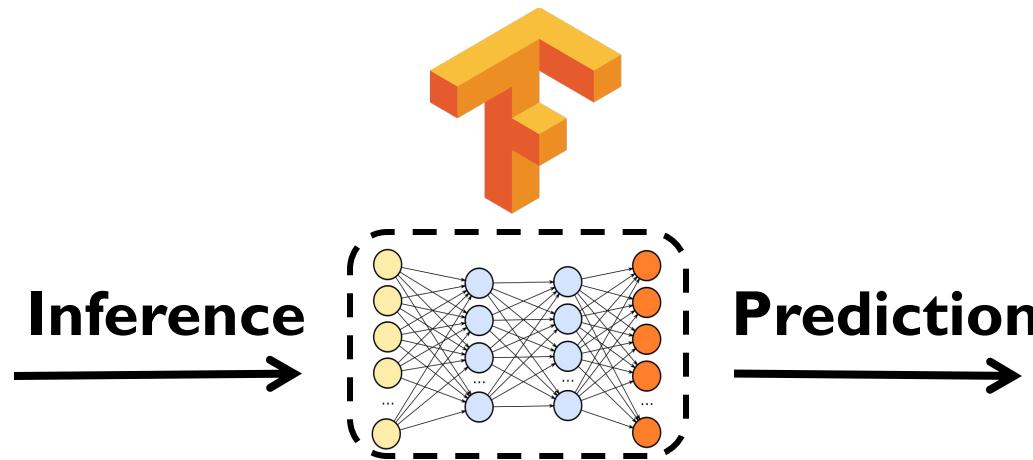


Video Capture

Google's **video codec**

SAFARI

TensorFlow Mobile



**57.3% of the inference energy is spent on
data movement**



**54.4% of the data movement energy comes from
packing/unpacking and quantization**

Packing



Reorders elements of matrices to minimize
cache misses during **matrix multiplication**



Up to **40%** of the
inference **energy** and **31%** of
inference **execution time**



Packing's data movement
accounts for up to
35.3% of the inference **energy**

A simple **data reorganization** process
that requires **simple arithmetic**

Quantization



Converts 32-bit floating point to 8-bit integers to improve inference execution time and energy consumption



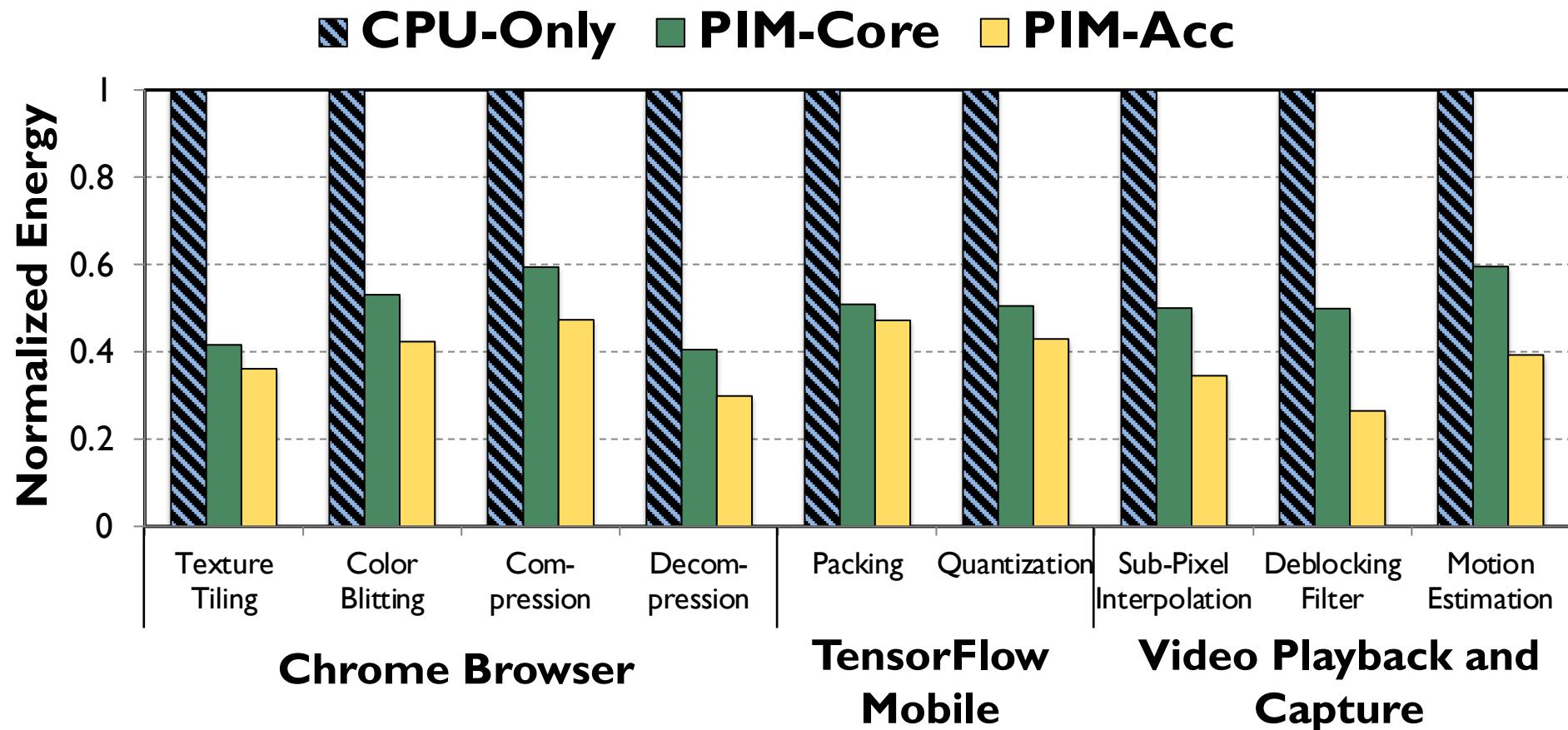
Up to **16.8%** of the inference **energy** and **16.1%** of inference **execution time**



Majority of **quantization energy** comes from **data movement**

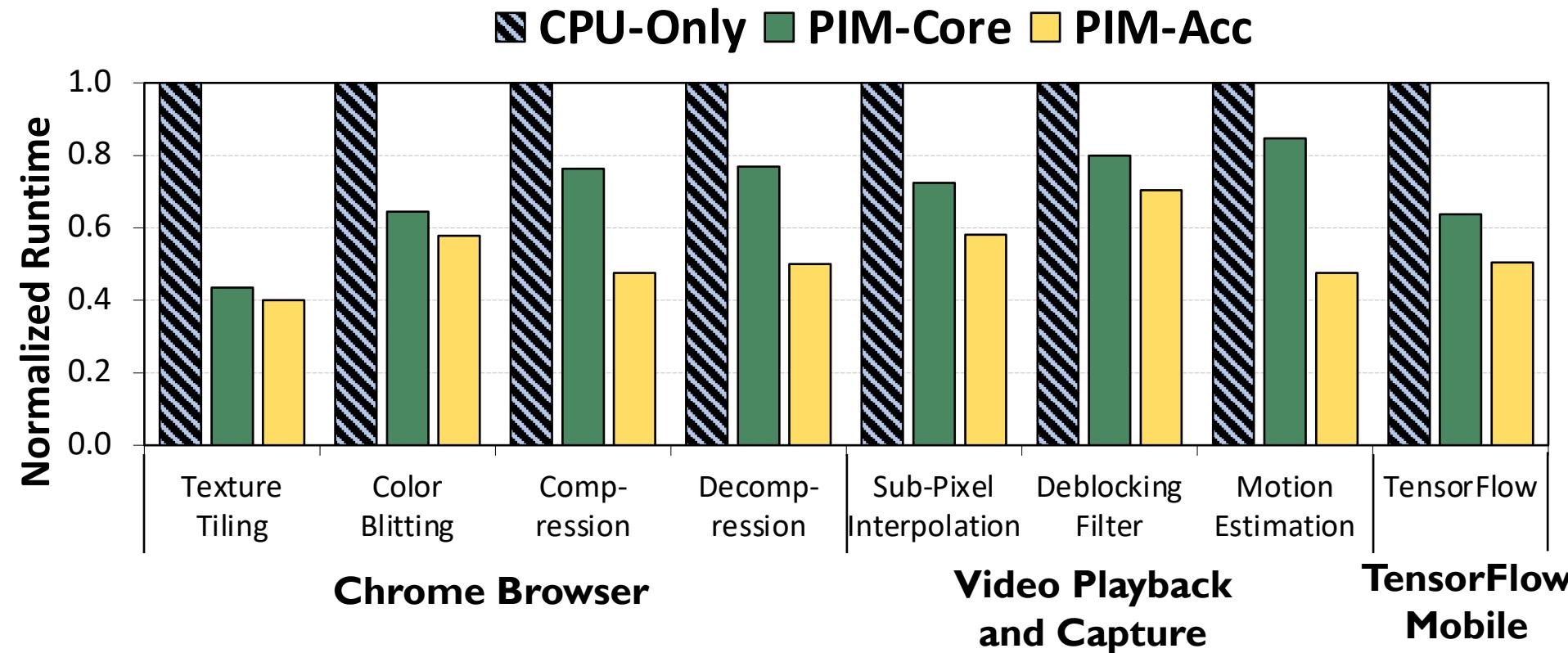
A simple **data conversion** operation that requires **shift**, **addition**, and **multiplication** operations

Normalized Energy



**PIM core and PIM accelerator reduce
energy consumption on average by 49.1% and 55.4%**

Normalized Runtime



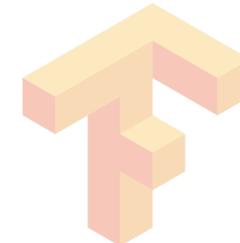
Offloading these kernels to **PIM core** and **PIM accelerator** improves performance on average by **44.6%** and **54.2%**

Workload Analysis



Chrome

Google's web browser



TensorFlow

Google's machine learning
framework



Video Playback

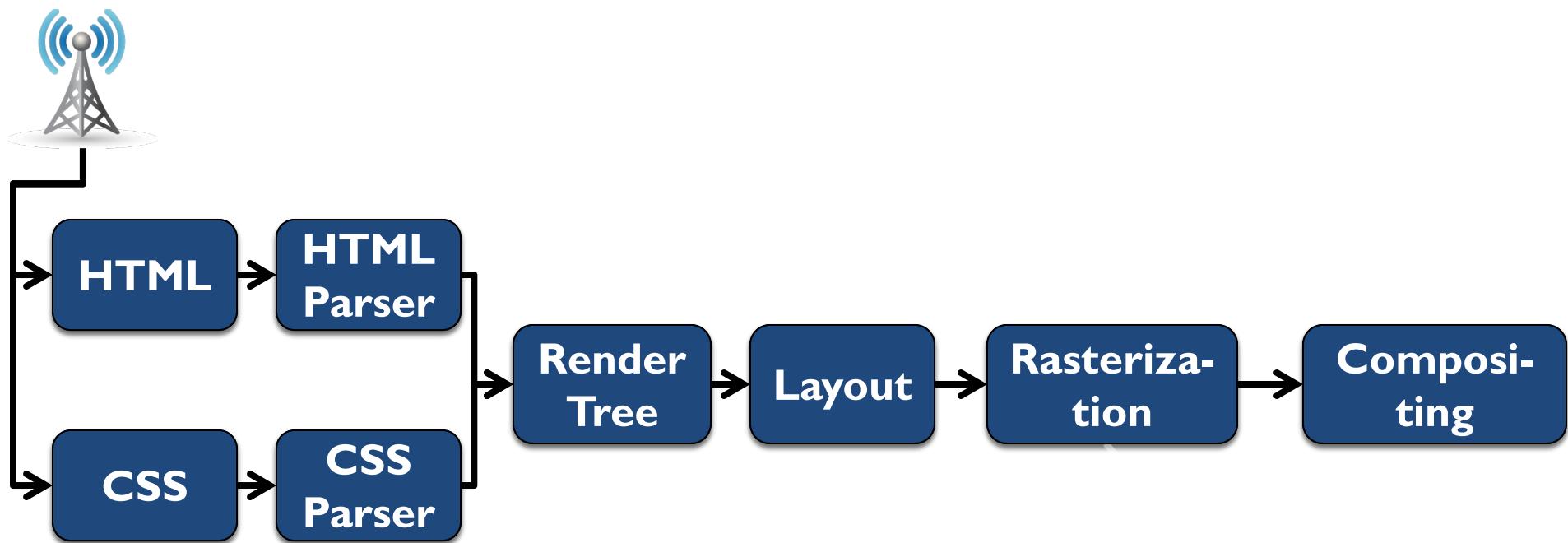
Google's video codec



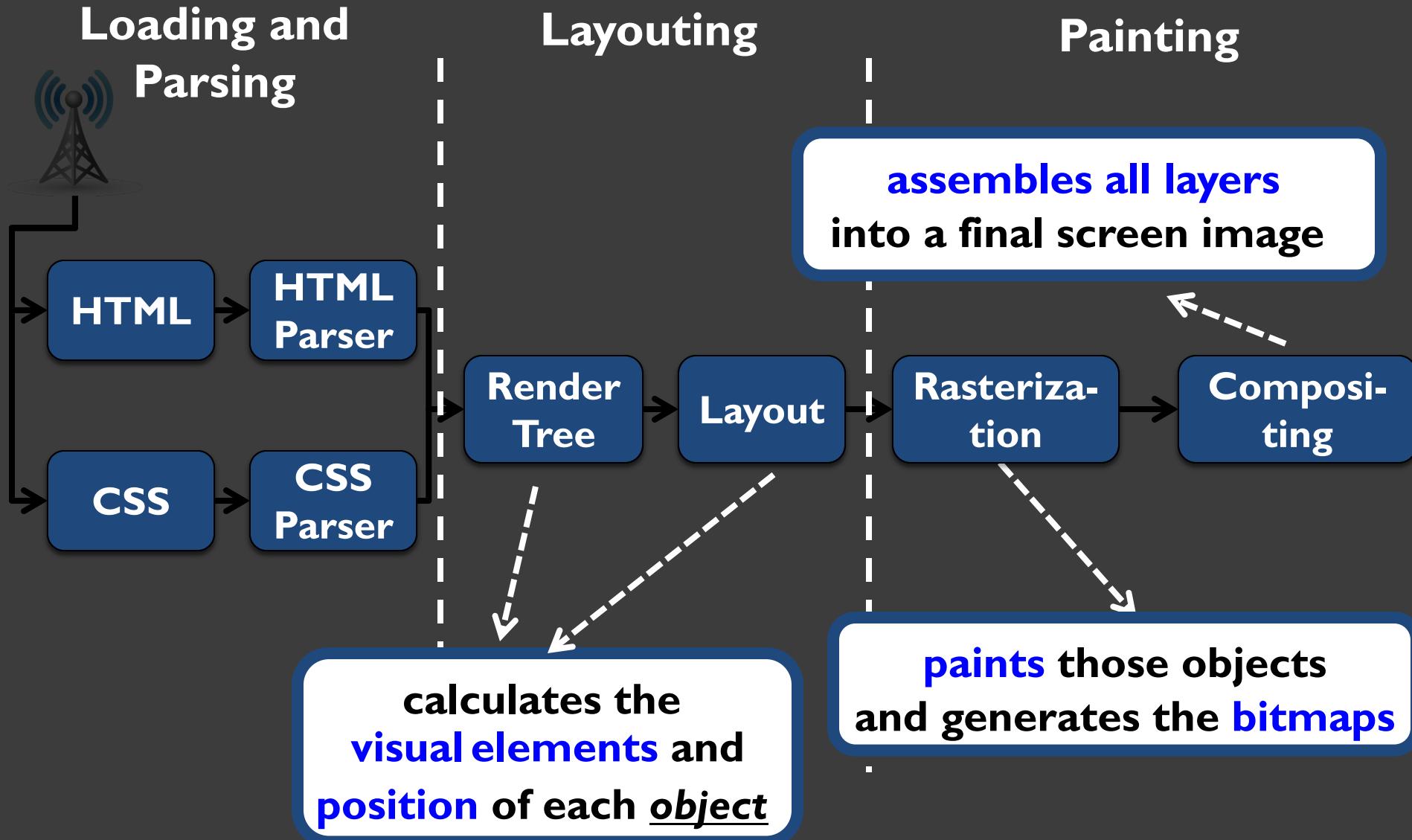
Video Capture

Google's video codec

How Chrome Renders a Web Page



How Chrome Renders a Web Page



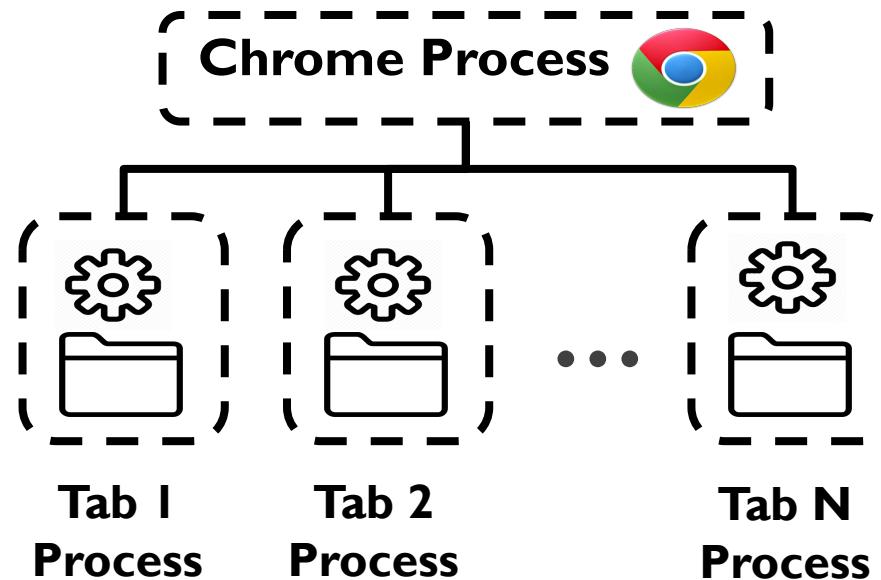
Browser Analysis

- To satisfy user experience, the browser must provide:
 - Fast **loading** of webpages
 - Smooth **scrolling** of webpages
 - Quick **switching** between browser tabs
- We focus on two important user interactions:
 - 1) **Page Scrolling**
 - 2) **Tab Switching**
 - Both include page loading

Tab Switching

What Happens During Tab Switching?

- Chrome employs a **multi-process architecture**
 - Each tab is a separate process

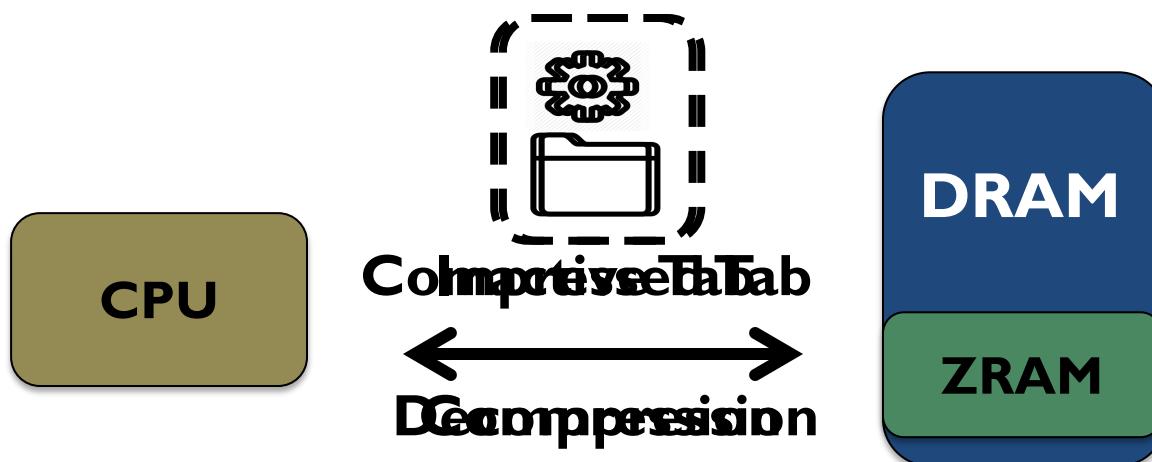


- Main operations during tab switching:
 - Context switch
 - Load the new page

Memory Consumption

- Primary concerns during tab switching:
 - How fast a new tab loads and becomes interactive
 - Memory consumption

Chrome uses compression to reduce each tab's memory footprint



Data Movement Study

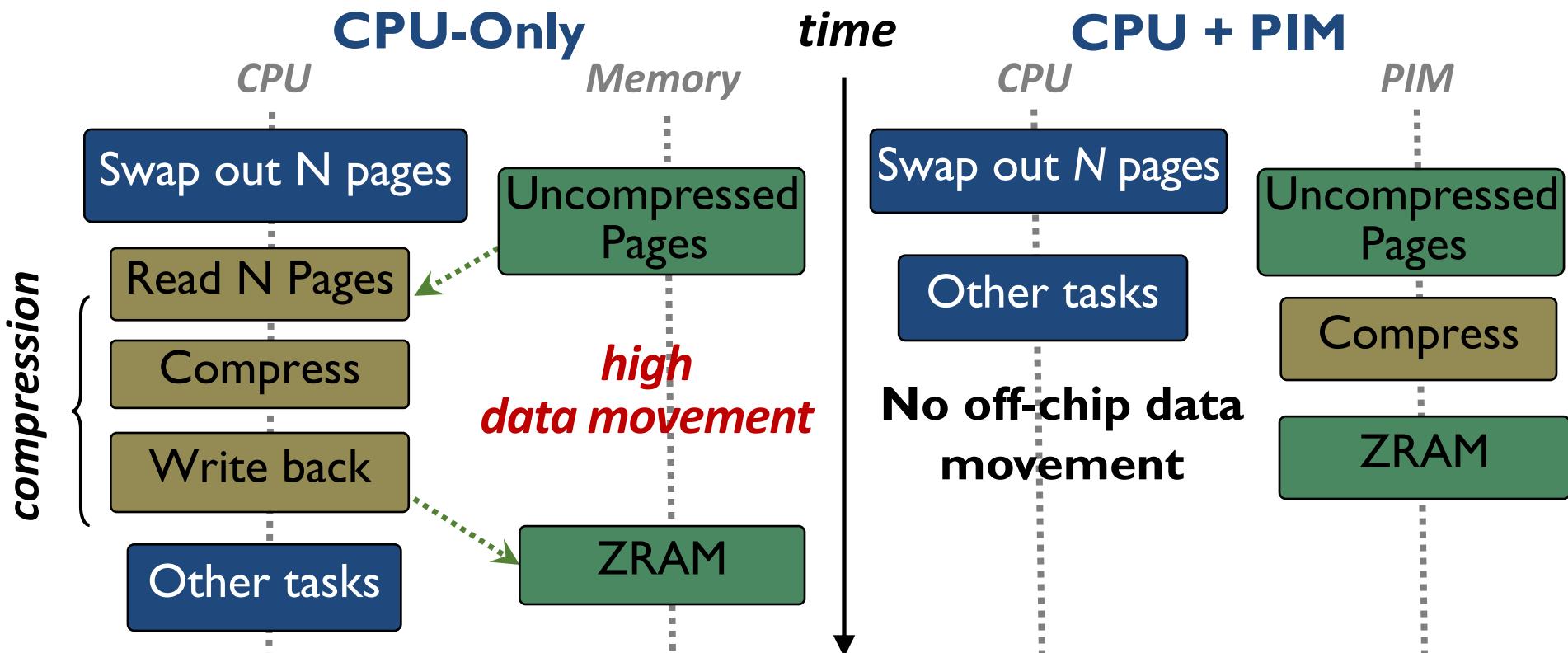
- To study **data movement** during tab switching, we emulate a user switching through 50 tabs

We make two **key observations**:

1 **Compression and decompression contribute to 18.1% of the total system energy**

2 **19.6 GB of data moves between CPU and ZRAM**

Can We Use PIM to Mitigate the Cost?



PIM core and PIM accelerator are feasible to implement in-memory compression/decompression

Tab Switching Wrap Up

A large amount of **data movement** happens during **tab switching** as Chrome attempts to **compress** and **decompress** tabs

Both functions can benefit from PIM execution and can be implemented as PIM logic

More on PIM for Mobile Devices

- Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu,
"Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks"

Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Williamsburg, VA, USA, March 2018.

**62.7% of the total system energy
is spent on data movement**

Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand¹

Rachata Ausavarungnirun¹

Aki Kuusela³

Saugata Ghose¹

Eric Shiu³

Allan Knies³

Youngsok Kim²

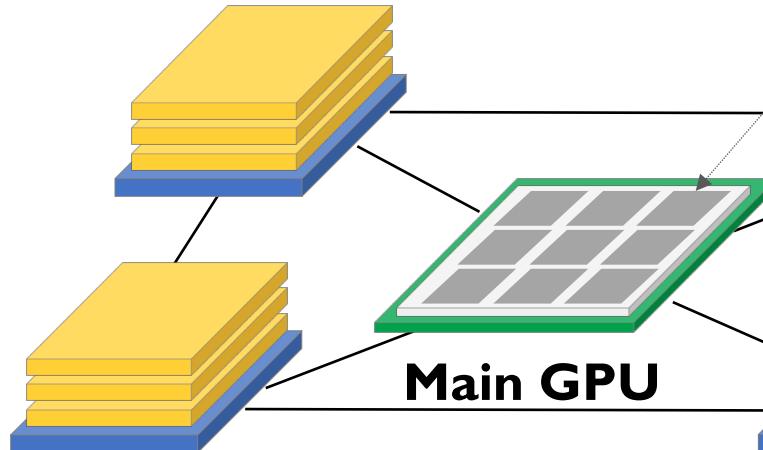
Rahul Thakur³

Daehyun Kim^{4,3}

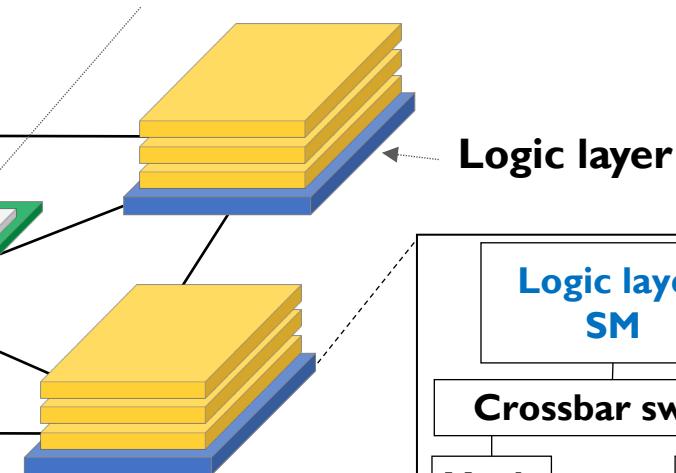
Onur Mutlu^{5,1}

Truly Distributed GPU Processing with PIM?

**3D-stacked memory
(memory stack)**



SM (Streaming Multiprocessor)

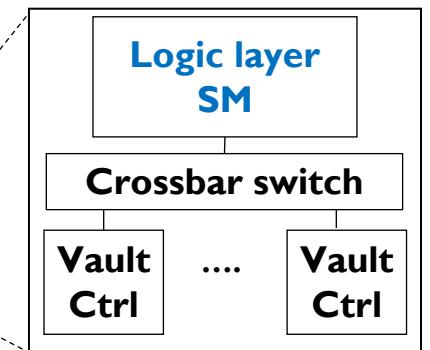


```
__global__
void applyScaleFactorsKernel( uint8_T * const out,
    uint8_T const * const in, const double *factor,
    size_t const numRows, size_t const numCols )
{
    // Work out which pixel we are working on.
    const int rowIdx = blockIdx.x * blockDim.x + threadIdx.x;
    const int colIdx = blockIdx.y;
    const int sliceIdx = threadIdx.z;

    // Check this thread isn't off the image
    if( rowIdx >= numRows ) return;

    // Compute the index of my element
    size_t linearIdx = rowIdx + colIdx*numRows +
        sliceIdx*numRows*numCols;
```

Logic layer



Accelerating GPU Execution with PIM (I)

- Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladrish Chatterjee, Mike O'Connor, Nandita Vijaykumar, Onur Mutlu, and Stephen W. Keckler,

"Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems"

Proceedings of the 43rd International Symposium on Computer Architecture (ISCA), Seoul, South Korea, June 2016.

[Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems

Kevin Hsieh[†] Eiman Ebrahimi[†] Gwangsun Kim^{*} Niladrish Chatterjee[†] Mike O'Connor[†]
Nandita Vijaykumar[†] Onur Mutlu^{§‡} Stephen W. Keckler[†]

[†]Carnegie Mellon University [†]NVIDIA ^{*}KAIST [§]ETH Zürich

Accelerating GPU Execution with PIM (II)

- Ashutosh Pattnaik, Xulong Tang, Adwait Jog, Onur Kayiran, Asit K. Mishra, Mahmut T. Kandemir, Onur Mutlu, and Chita R. Das,
"Scheduling Techniques for GPU Architectures with Processing-In-Memory Capabilities"

Proceedings of the 25th International Conference on Parallel Architectures and Compilation Techniques (PACT), Haifa, Israel, September 2016.

Scheduling Techniques for GPU Architectures with Processing-In-Memory Capabilities

Ashutosh Pattnaik¹ Xulong Tang¹ Adwait Jog² Onur Kayıran³
Asit K. Mishra⁴ Mahmut T. Kandemir¹ Onur Mutlu^{5,6} Chita R. Das¹

¹Pennsylvania State University ²College of William and Mary

³Advanced Micro Devices, Inc. ⁴Intel Labs ⁵ETH Zürich ⁶Carnegie Mellon University

Accelerating Linked Data Structures

- Kevin Hsieh, Samira Khan, Nandita Vijaykumar, Kevin K. Chang, Amirali Boroumand, Saugata Ghose, and Onur Mutlu,

"Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation"

Proceedings of the 34th IEEE International Conference on Computer Design (ICCD), Phoenix, AZ, USA, October 2016.

Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation

Kevin Hsieh[†] Samira Khan[‡] Nandita Vijaykumar[†]
Kevin K. Chang[†] Amirali Boroumand[†] Saugata Ghose[†] Onur Mutlu^{§†}
[†]*Carnegie Mellon University* [‡]*University of Virginia* [§]*ETH Zürich*

Accelerating Dependent Cache Misses

- Milad Hashemi, Khubaib, Eiman Ebrahimi, Onur Mutlu, and Yale N. Patt,
"Accelerating Dependent Cache Misses with an Enhanced Memory Controller"

Proceedings of the 43rd International Symposium on Computer Architecture (ISCA), Seoul, South Korea, June 2016.

[Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

Accelerating Dependent Cache Misses with an Enhanced Memory Controller

Milad Hashemi*, Khubaib†, Eiman Ebrahimi‡, Onur Mutlu§, Yale N. Patt*

*The University of Texas at Austin †Apple ‡NVIDIA §ETH Zürich & Carnegie Mellon University

Accelerating Runahead Execution

- Milad Hashemi, Onur Mutlu, and Yale N. Patt,
"Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads"
Proceedings of the 49th International Symposium on Microarchitecture (MICRO), Taipei, Taiwan, October 2016.
[Slides (pptx) (pdf)] [Lightning Session Slides (pdf)] [Poster (pptx) (pdf)]

Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads

Milad Hashemi*, Onur Mutlu[§], Yale N. Patt*

*The University of Texas at Austin §ETH Zürich

Accelerating Climate Modeling

- Gagandeep Singh, Dionysios Diamantopoulos, Christoph Hagleitner, Juan Gómez-Luna, Sander Stuijk, Onur Mutlu, and Henk Corporaal,

"NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling"

Proceedings of the 30th International Conference on Field-Programmable Logic and Applications (FPL), Gothenburg, Sweden, September 2020.

[[Slides \(pptx\)](#) ([pdf](#))]

[[Lightning Talk Slides \(pptx\)](#) ([pdf](#))]

[[Talk Video](#) (23 minutes)]

Nominated for the Stamatis Vassiliadis Memorial Award.

NERO: A Near High-Bandwidth Memory Stencil Accelerator for Weather Prediction Modeling

Gagandeep Singh^{a,b,c}

Dionysios Diamantopoulos^c

Christoph Hagleitner^c

Juan Gómez-Luna^b

Sander Stuijk^a

Onur Mutlu^b

Henk Corporaal^a

^aEindhoven University of Technology

^bETH Zürich

^cIBM Research Europe, Zurich

Accelerating Approximate String Matching

- Damla Senol Cali, Gurpreet S. Kalsi, Zulal Bingol, Can Firtina, Lavanya Subramanian, Jeremie S. Kim, Rachata Ausavarungnirun, Mohammed Alser, Juan Gomez-Luna, Amirali Boroumand, Anant Nori, Allison Scibisz, Sreenivas Subramoney, Can Alkan, Saugata Ghose, and Onur Mutlu,

["GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis"](#)

Proceedings of the [53rd International Symposium on Microarchitecture \(MICRO\)](#), Virtual, October 2020.

- [[Lightning Talk Video](#) (1.5 minutes)]
- [[Lightning Talk Slides \(pptx\)](#) ([pdf](#))]
- [[Talk Video](#) (18 minutes)]
- [[Slides \(pptx\)](#) ([pdf](#))]

GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali^{†✉} Gurpreet S. Kalsi[✉] Zülal Bingöl[▽] Can Firtina[◊] Lavanya Subramanian[‡] Jeremie S. Kim^{◊†}
Rachata Ausavarungnirun[○] Mohammed Alser[◊] Juan Gomez-Luna[◊] Amirali Boroumand[†] Anant Nori[✉]
Allison Scibisz[†] Sreenivas Subramoney[✉] Can Alkan[▽] Saugata Ghose^{★†} Onur Mutlu^{◊†▽}

[†]*Carnegie Mellon University* [✉]*Processor Architecture Research Lab, Intel Labs* [▽]*Bilkent University* [◊]*ETH Zürich*

[‡]*Facebook* [○]*King Mongkut's University of Technology North Bangkok* [★]*University of Illinois at Urbana-Champaign*

Accelerating Time Series Analysis

- Ivan Fernandez, Ricardo Quislant, Christina Giannoula, Mohammed Alser, Juan Gómez-Luna, Eladio Gutiérrez, Oscar Plata, and Onur Mutlu,
"NATSA: A Near-Data Processing Accelerator for Time Series Analysis"
Proceedings of the 38th IEEE International Conference on Computer Design (ICCD), Virtual, October 2020.

NATSA: A Near-Data Processing Accelerator for Time Series Analysis

Ivan Fernandez[§]

Ricardo Quislant[§]

Christina Giannoula[†]

Mohammed Alser[‡]

Juan Gómez-Luna[‡]

Eladio Gutiérrez[§]

Oscar Plata[§]

Onur Mutlu[‡]

[§]*University of Malaga*

[†]*National Technical University of Athens*

[‡]*ETH Zürich*

Several Questions in 3D-Stacked PIM

- What are the performance and energy benefits of using 3D-stacked memory as a coarse-grained accelerator?
 - By changing the entire system
 - By performing simple function offloading

- What is the minimal processing-in-memory support we can provide?
 - With minimal changes to system and programming

PIM-Enabled Instructions

- Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,
"PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture"

Proceedings of the 42nd International Symposium on Computer Architecture (ISCA), Portland, OR, June 2015.
[Slides (pdf)] [Lightning Session Slides (pdf)]

PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture

Junwhan Ahn Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi

junwhan@snu.ac.kr, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr

Seoul National University

[†]Carnegie Mellon University

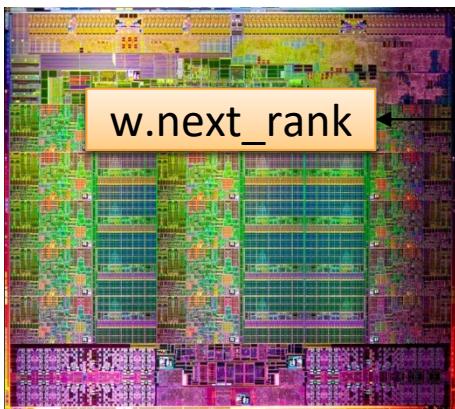
PEI: PIM-Enabled Instructions (Ideas)

- **Goal:** Develop mechanisms to get the most out of near-data processing with **minimal cost**, **minimal changes to the system**, no changes to the programming model
- **Key Idea 1:** Expose each PIM operation as a **cache-coherent, virtually-addressed host processor instruction** (called PEI) that operates on **only a single cache block**
 - e.g., `__pim_add(&w.next_rank, value) → pim.add r1, (r2)`
 - No changes sequential execution/programming model
 - No changes to virtual memory
 - Minimal changes to cache coherence
 - No need for data mapping: Each PEI restricted to a single memory module
- **Key Idea 2: Dynamically decide where to execute a PEI** (i.e., the host processor or PIM accelerator) based on simple locality characteristics and simple hardware predictors
 - Execute each operation at the location that provides the best performance

Simple PIM Operations as ISA Extensions (II)

```
for (v: graph.vertices) {  
    value = weight * v.rank;  
    for (w: v.successors) {  
        w.next_rank += value;  
    }  
}
```

Host Processor



Main Memory



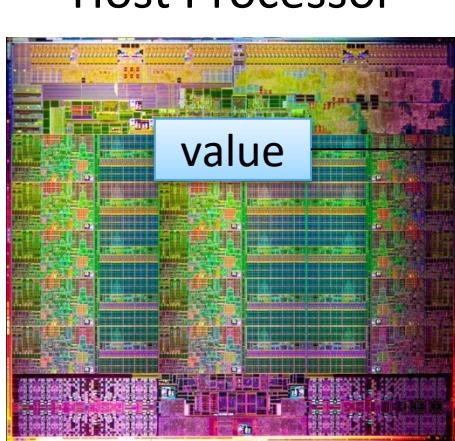
64 bytes in
64 bytes out

Conventional Architecture

Simple PIM Operations as ISA Extensions (III)

```
for (v: graph.vertices) {  
    value = weight * v.rank;  
    for (w: v.successors) {  
        __pim_add(&w.next_rank, value);  
    }  
}
```

pim.add r1, (r2)



Host Processor

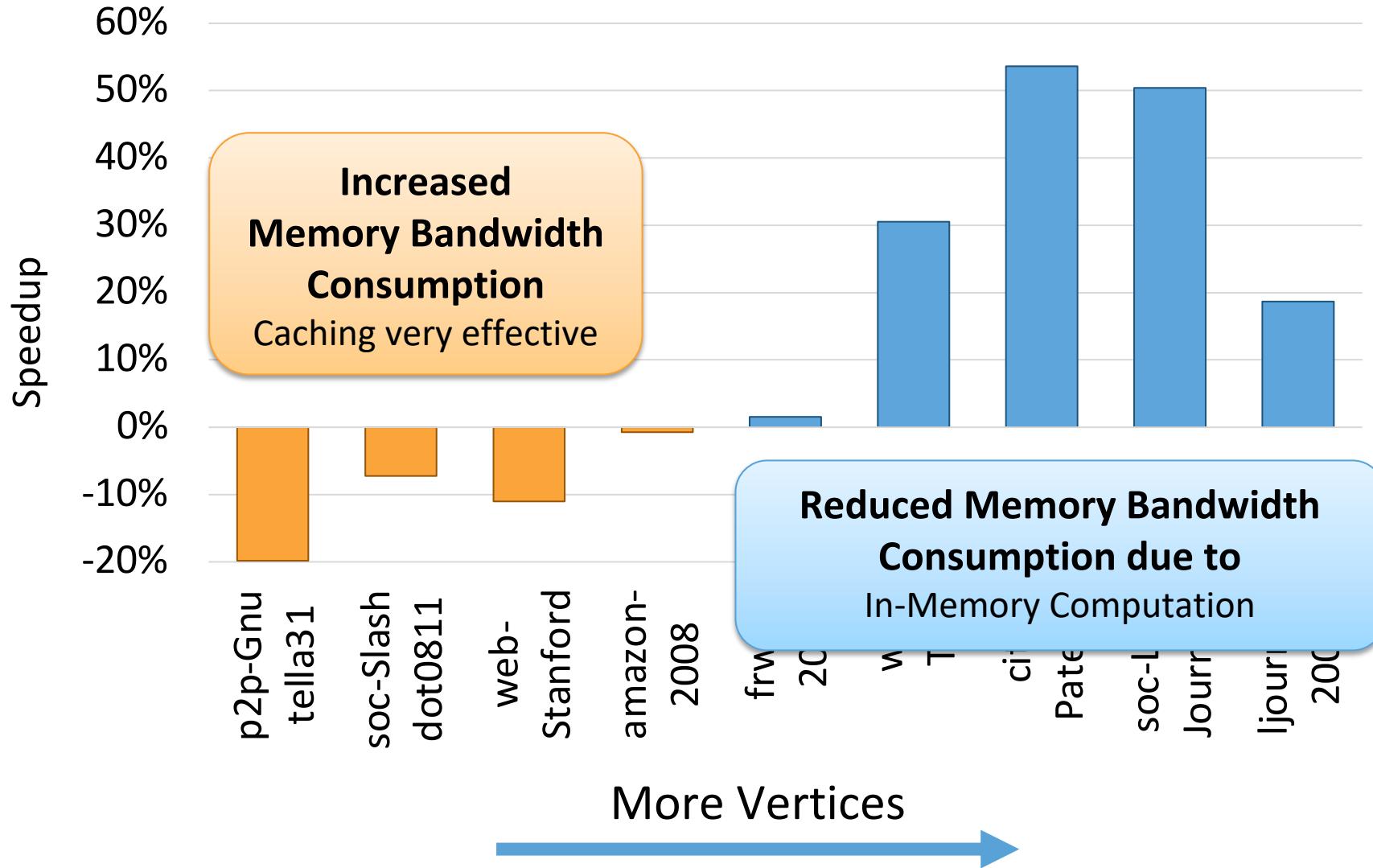
8 bytes **in**
0 bytes **out**



Main Memory

In-Memory Addition

Always Executing in Memory? Not A Good Idea



PEI: PIM-Enabled Instructions (Example)

```
for (v: graph.vertices) {  
    value = weight * v.rank;  
    for (w: v.successors) {  
        __pim_add(&w.next_rank, value);  
    }  
}  
pfence();
```

pim.add r1, (r2)

Table 1: Summary of Supported PIM Operations

Operation	R	W	Input	Output	Applications
8-byte integer increment	O	O	0 bytes	0 bytes	AT
8-byte integer min	O	O	8 bytes	0 bytes	BFS, SP, WCC
Floating-point add	O	O	8 bytes	0 bytes	PR
Hash table probing	O	X	8 bytes	9 bytes	HJ
Histogram bin index	O	X	1 byte	16 bytes	HG, RP
Euclidean distance	O	X	64 bytes	4 bytes	SC
Dot product	O	X	32 bytes	8 bytes	SVM

- Executed either in memory or in the processor: dynamic decision
 - Low-cost locality monitoring for a single instruction
- Cache-coherent, virtually-addressed, single cache block only
- Atomic between different PEIs
- Not* atomic with normal instructions (use *pfence* for ordering)

PIM-Enabled Instructions

- Key to practicality: single-cache-block restriction
 - Each PEI can access *at most one last-level cache block*
 - Similar restrictions exist in atomic instructions
- Benefits
 - **Localization:** each PEI is bounded to one memory module
 - **Interoperability:** easier support for cache coherence and virtual memory
 - **Simplified locality monitoring:** data locality of PEIs can be identified simply by the cache control logic

PEI: Initial Evaluation Results

- Initial evaluations with **10 emerging data-intensive workloads**
 - Large-scale graph processing
 - In-memory data analytics
 - Machine learning and data mining
 - Three input sets (small, medium, large) for each workload to analyze the impact of data locality
- Pin-based cycle-level x86-64 simulation
- **Performance Improvement and Energy Reduction:**
 - 47% average speedup with large input data sets
 - 32% speedup with small input data sets
 - 25% avg. energy reduction in a single node with large input data sets

Table 2: Baseline Simulation Configuration

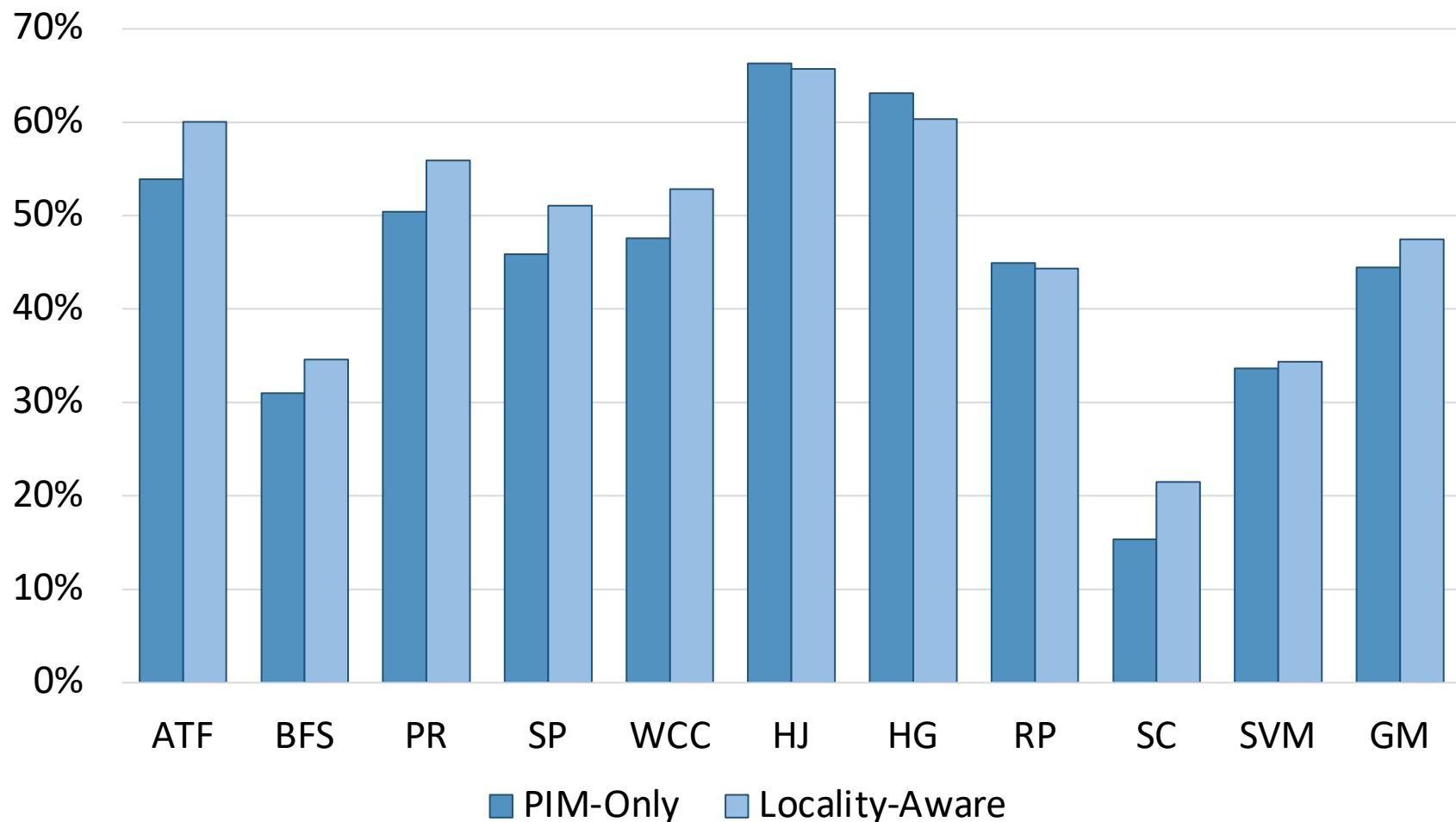
Component	Configuration
Core	16 out-of-order cores, 4 GHz, 4-issue
L1 I/D-Cache	Private, 32 KB, 4/8-way, 64 B blocks, 16 MSHRs
L2 Cache	Private, 256 KB, 8-way, 64 B blocks, 16 MSHRs
L3 Cache	Shared, 16 MB, 16-way, 64 B blocks, 64 MSHRs
On-Chip Network	Crossbar, 2 GHz, 144-bit links
Main Memory	32 GB, 8 HMCs, daisy-chain (80 GB/s full-duplex)
HMC	4 GB, 16 vaults, 256 DRAM banks [20]
– DRAM	FR-FCFS, tCL = tRCD = tRP = 13.75 ns [27]
– Vertical Links	64 TSVs per vault with 2 Gb/s signaling rate [23]

Evaluated Data-Intensive Applications

- **Ten emerging data-intensive workloads**
 - Large-scale graph processing
 - Average teenage follower, BFS, PageRank, single-source shortest path, weakly connected components
 - In-memory data analytics
 - Hash join, histogram, radix partitioning
 - Machine learning and data mining
 - Streamcluster, SVM-RFE
- Three input sets (small, medium, large) for each workload to show the impact of data locality

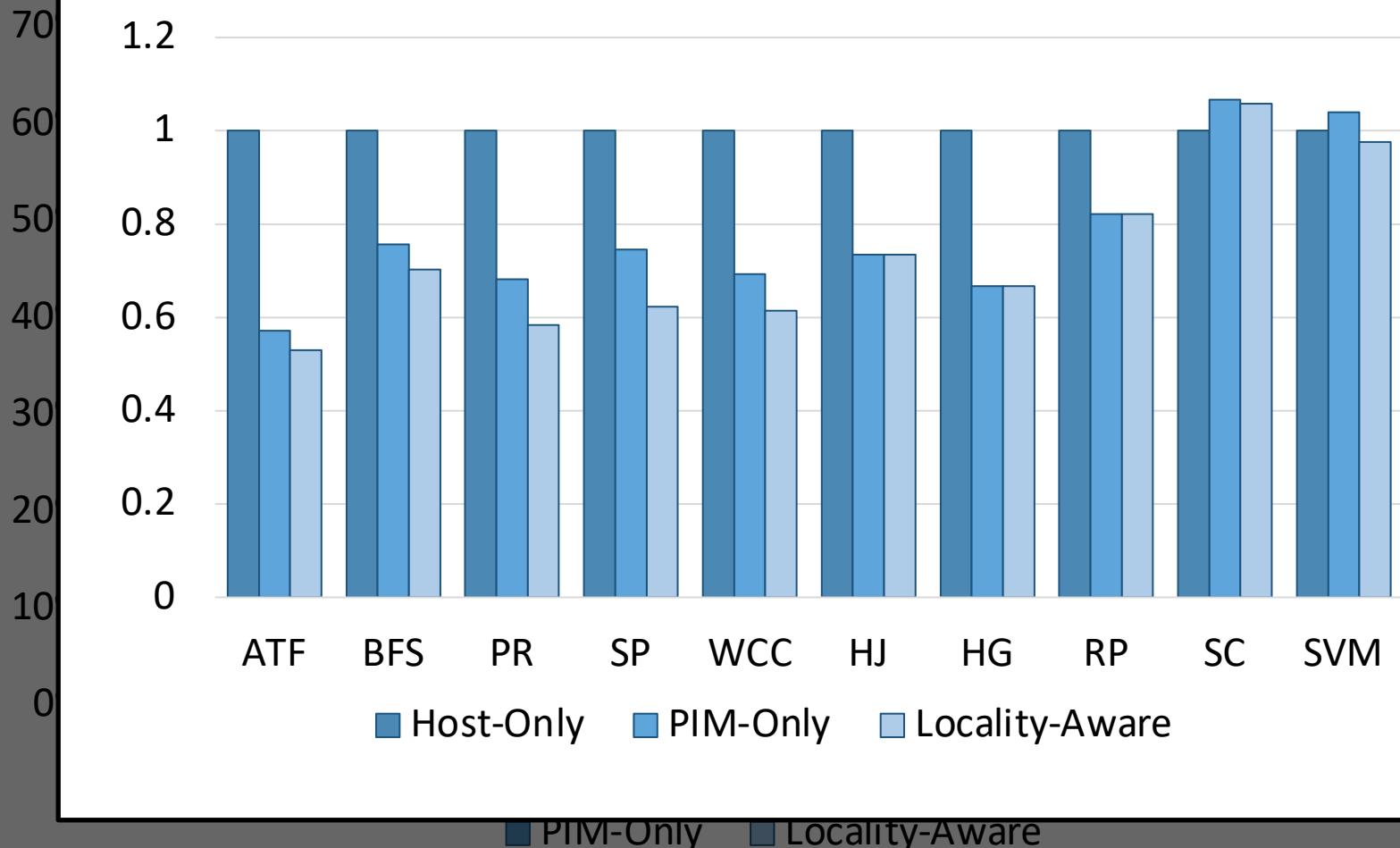
PEI Performance Delta: Large Data Sets

(Large Inputs, Baseline: Host-Only)



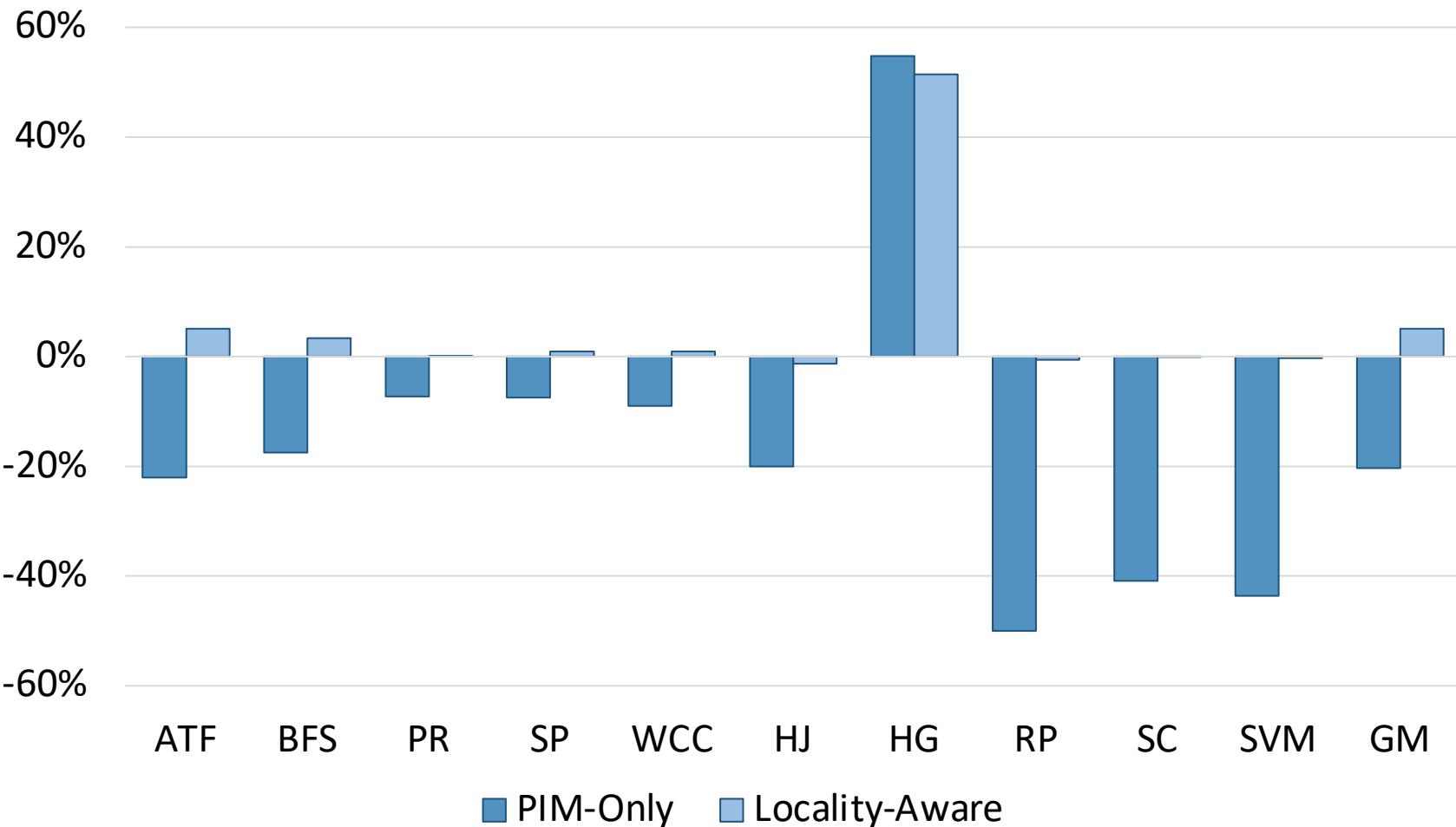
PEI Performance: Large Data Sets

Normalized Amount of Off-chip Transfer

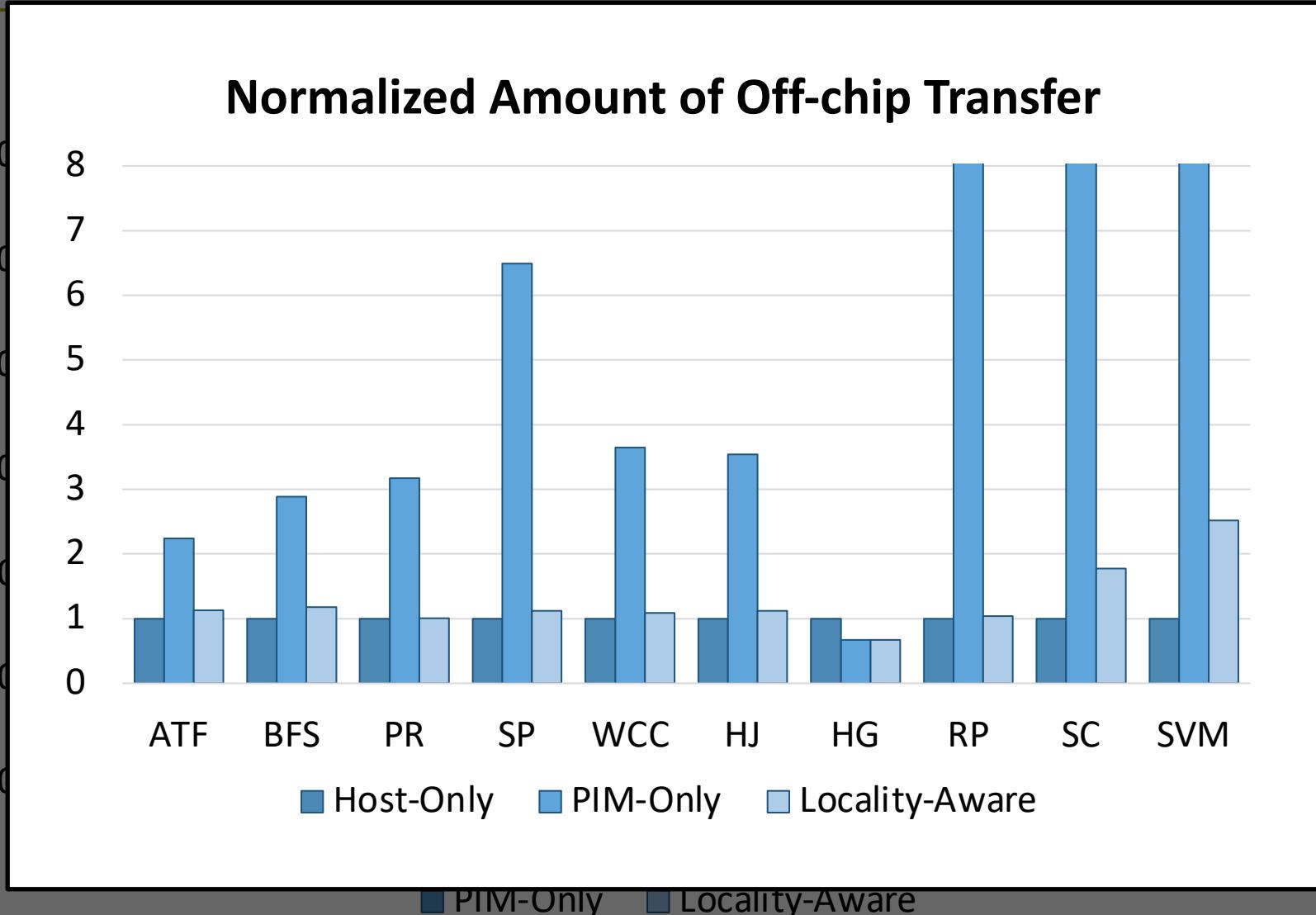


PEI Performance Delta: Small Data Sets

(Small Inputs, Baseline: Host-Only)

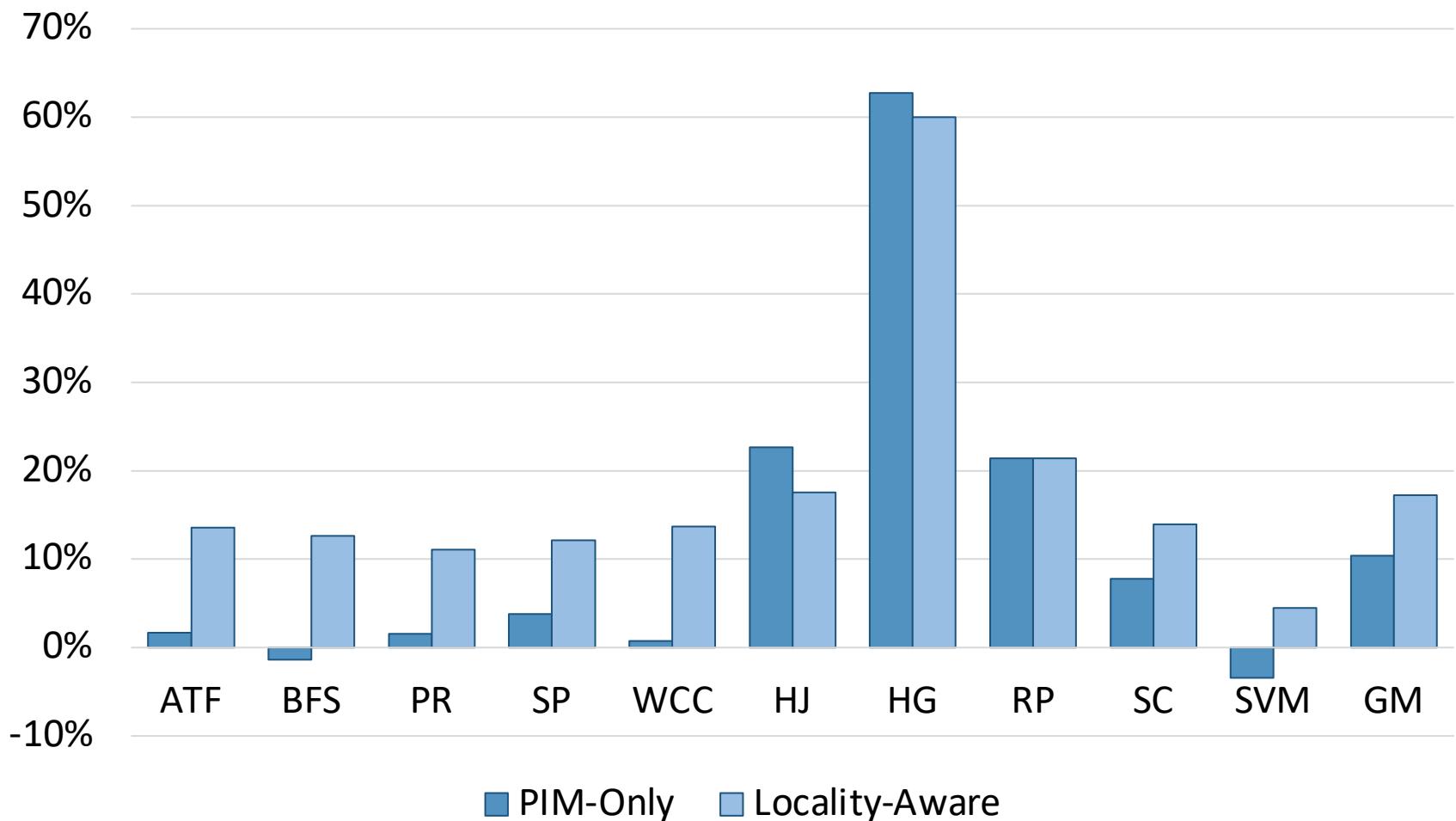


PEI Performance: Small Data Sets

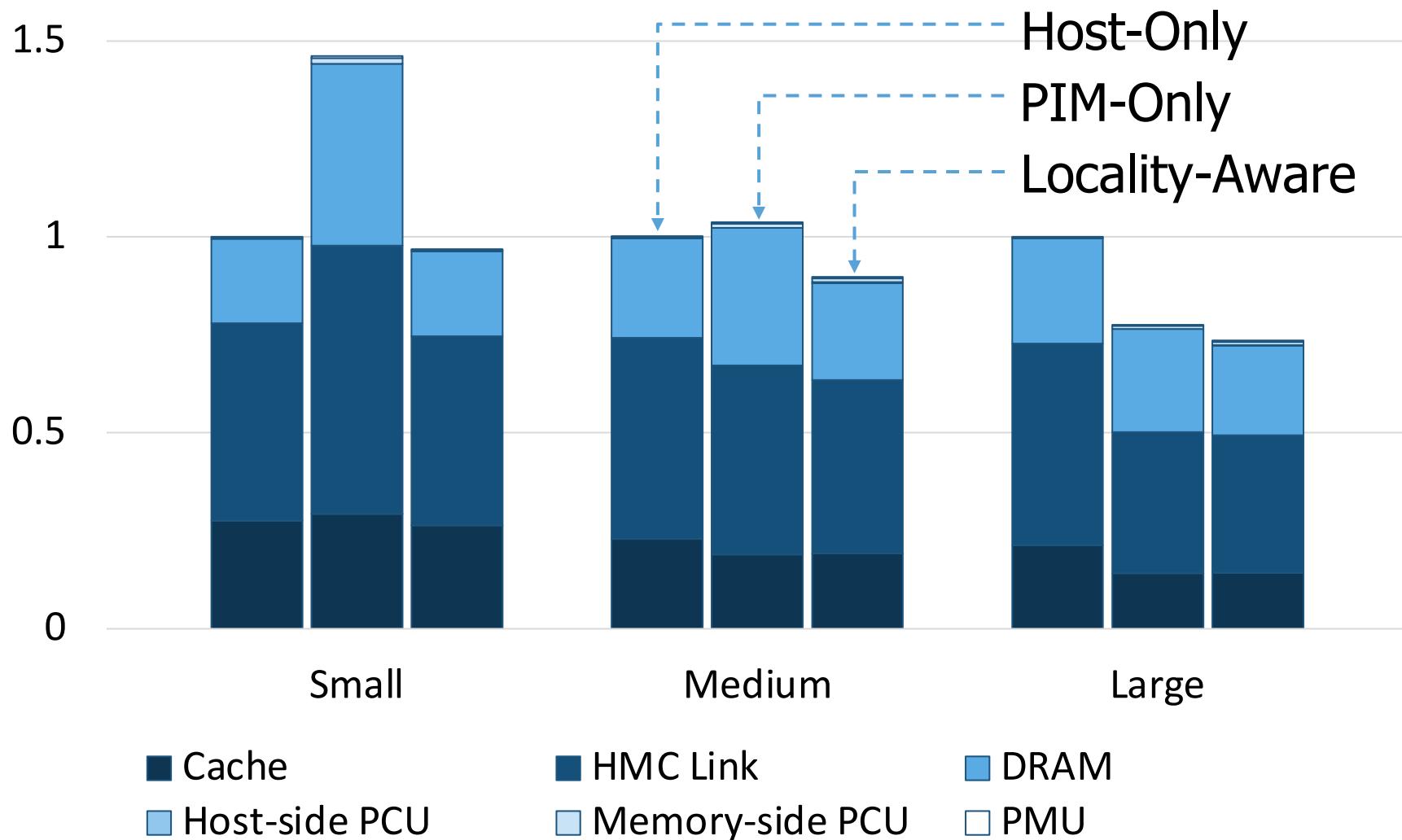


PEI Performance Delta: Medium Data Sets

(Medium Inputs, Baseline: Host-Only)



PEI Energy Consumption



PEI: Advantages & Disadvantages

- **Advantages**
 - + Simple and low cost approach to PIM
 - + No changes to programming model, virtual memory
 - + Dynamically decides where to execute an instruction

- **Disadvantages**
 - Does not take full advantage of PIM potential
 - Single cache block restriction is limiting

Simpler PIM: PIM-Enabled Instructions

- Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,
**"PIM-Enabled Instructions: A Low-Overhead,
Locality-Aware Processing-in-Memory Architecture"**
*Proceedings of the 42nd International Symposium on
Computer Architecture (ISCA)*, Portland, OR, June 2015.
[Slides (pdf)] [Lightning Session Slides (pdf)]

PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture

Junwhan Ahn Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi

junwhan@snu.ac.kr, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr

Seoul National University

[†]Carnegie Mellon University

Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand

Ravi Narayanaswami

Eric Shiu

Saugata Ghose

Geraldo F. Oliveira

Onur Mutlu

Berkin Akin

Xiaoyu Ma

PACT 2021

SAFARI

UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN

Carnegie Mellon



Google

ETH Zürich

Executive Summary

Context: We extensively analyze a state-of-the-art edge ML accelerator (**Google Edge TPU**) using 24 Google edge models

- Wide range of models (CNNs, LSTMs, Transducers, RCNNs)

Problem: The Edge TPU accelerator suffers from three challenges:

- It operates significantly below its peak throughput
- It operates significantly below its theoretical energy efficiency
- It inefficiently handles memory accesses

Key Insight: These shortcomings arise from the monolithic design of the Edge TPU accelerator

- The Edge TPU accelerator design does not account for layer heterogeneity

Key Mechanism: A new framework called Mensa

- Mensa consists of heterogeneous accelerators whose dataflow and hardware are specialized for specific families of layers

Key Results: We design a version of Mensa for Google edge ML models

- Mensa improves performance and energy by 3.0X and 3.1X
- Mensa reduces cost and improves area efficiency

Google Edge Neural Network Models

We analyze inference execution using 24 edge NN models



Speech Recognition

6 RNN
Transducers



2 LSTMs



Language Translation



Face Detection

13 CNN

Google Edge TPU

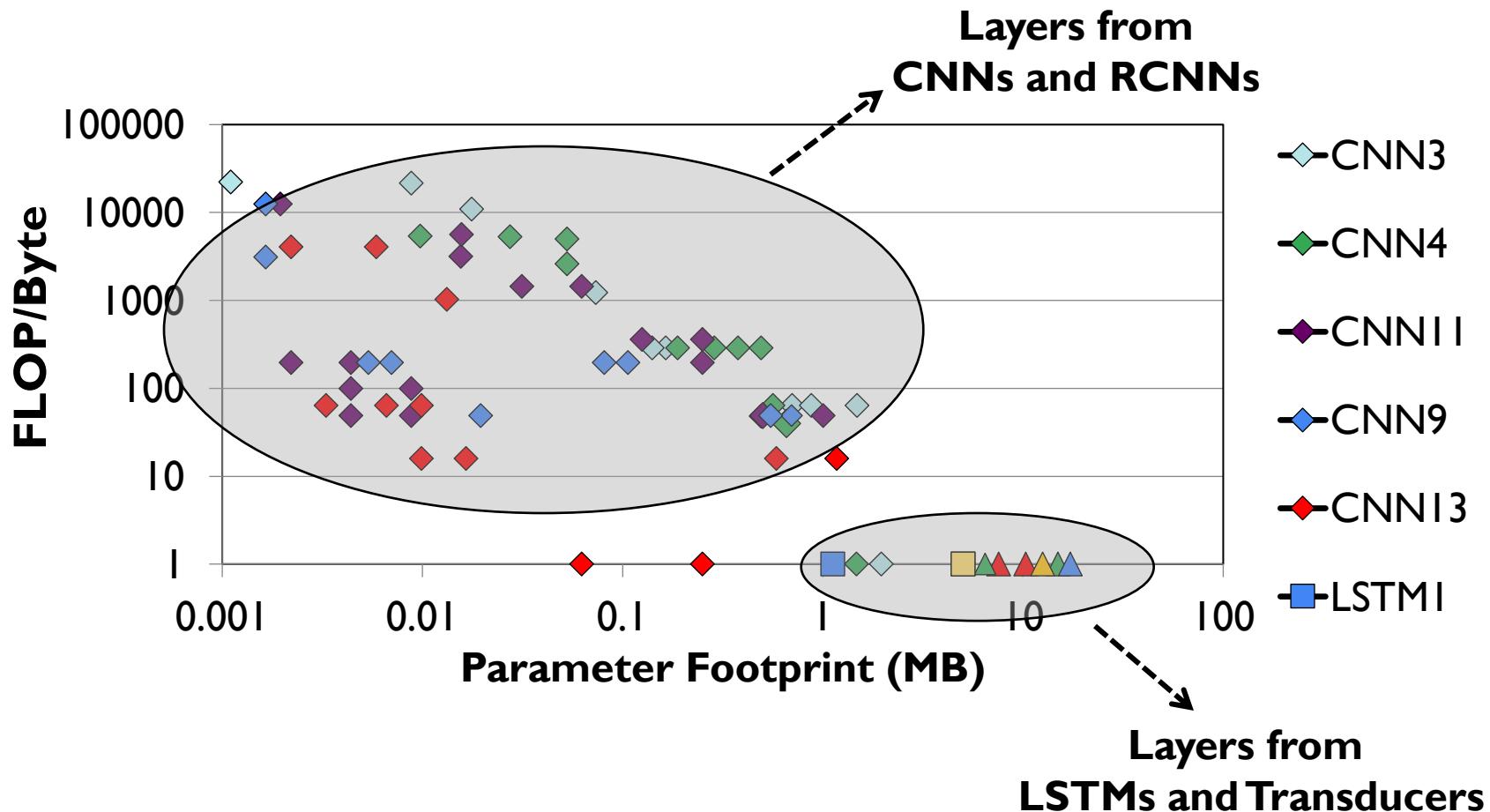
3 RCNN



Image Captioning

Diversity Across the Models

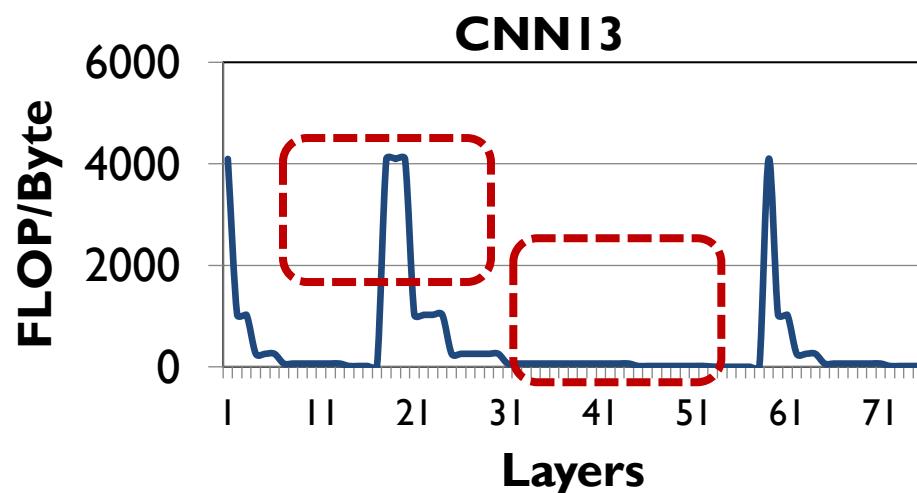
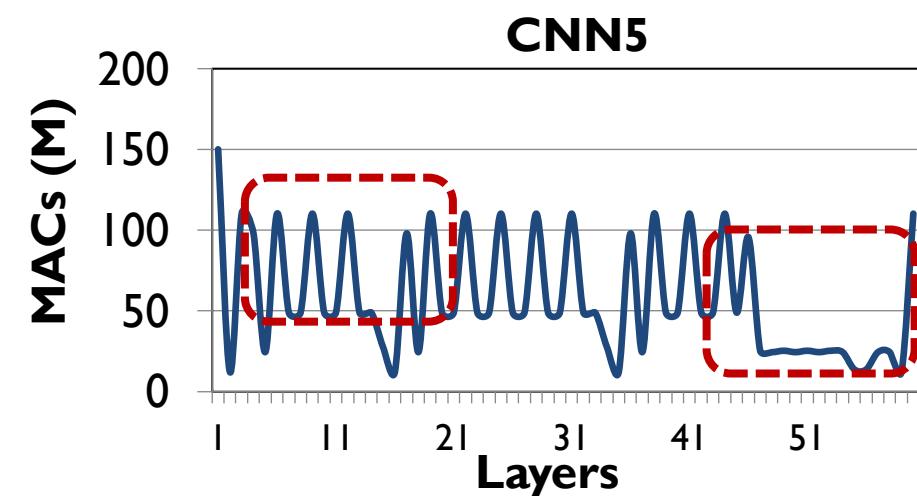
Insight I: there is **significant variation** in terms of layer characteristics across the models



Diversity Within the Models

Insight 2: even **within each model, layers exhibit significant variation in terms of layer characteristics**

For example, our analysis of edge CNN models shows:



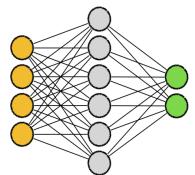
Variation in MAC intensity: up to 200x across layers

Variation in FLOP/Byte: up to 244x across layers

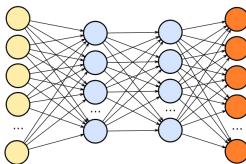
Mensa High-Level Overview

Edge TPU Accelerator

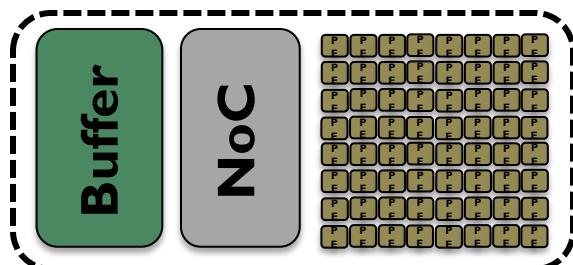
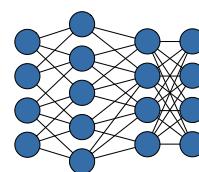
Model A



Model B



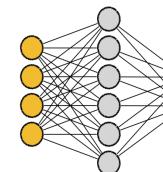
Model C



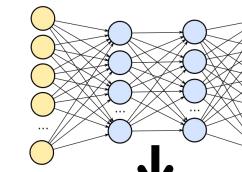
Monolithic Accelerator

Mensa

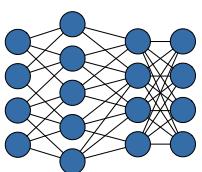
Model A



Model B

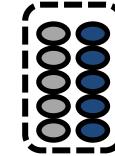


Model C

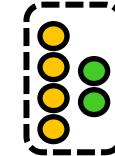


Runtime

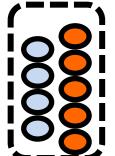
Family 1



Family 2

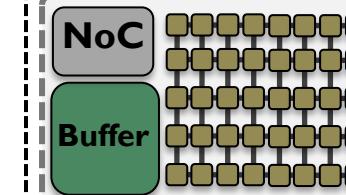


Family 3

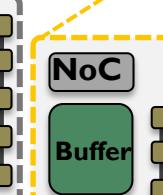


CPU

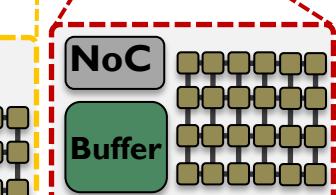
3D-Stacked DRAM



Acc. 1



Acc. 2

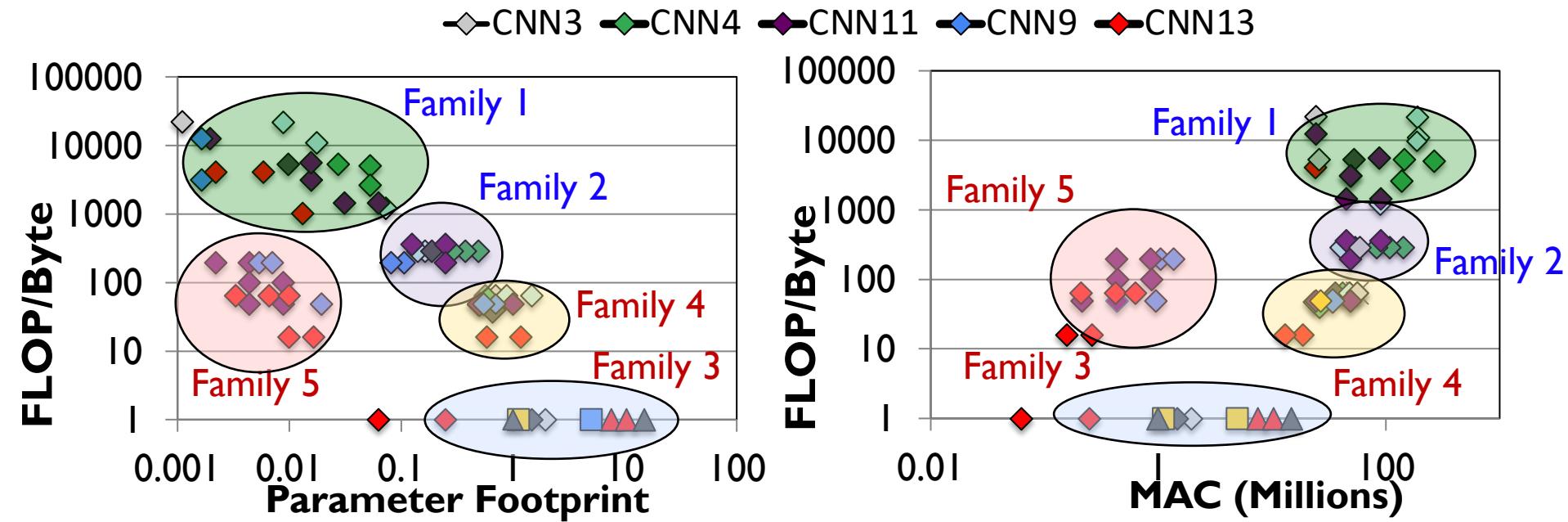


Acc. 3

Heterogeneous Accelerators

Identifying Layer Families

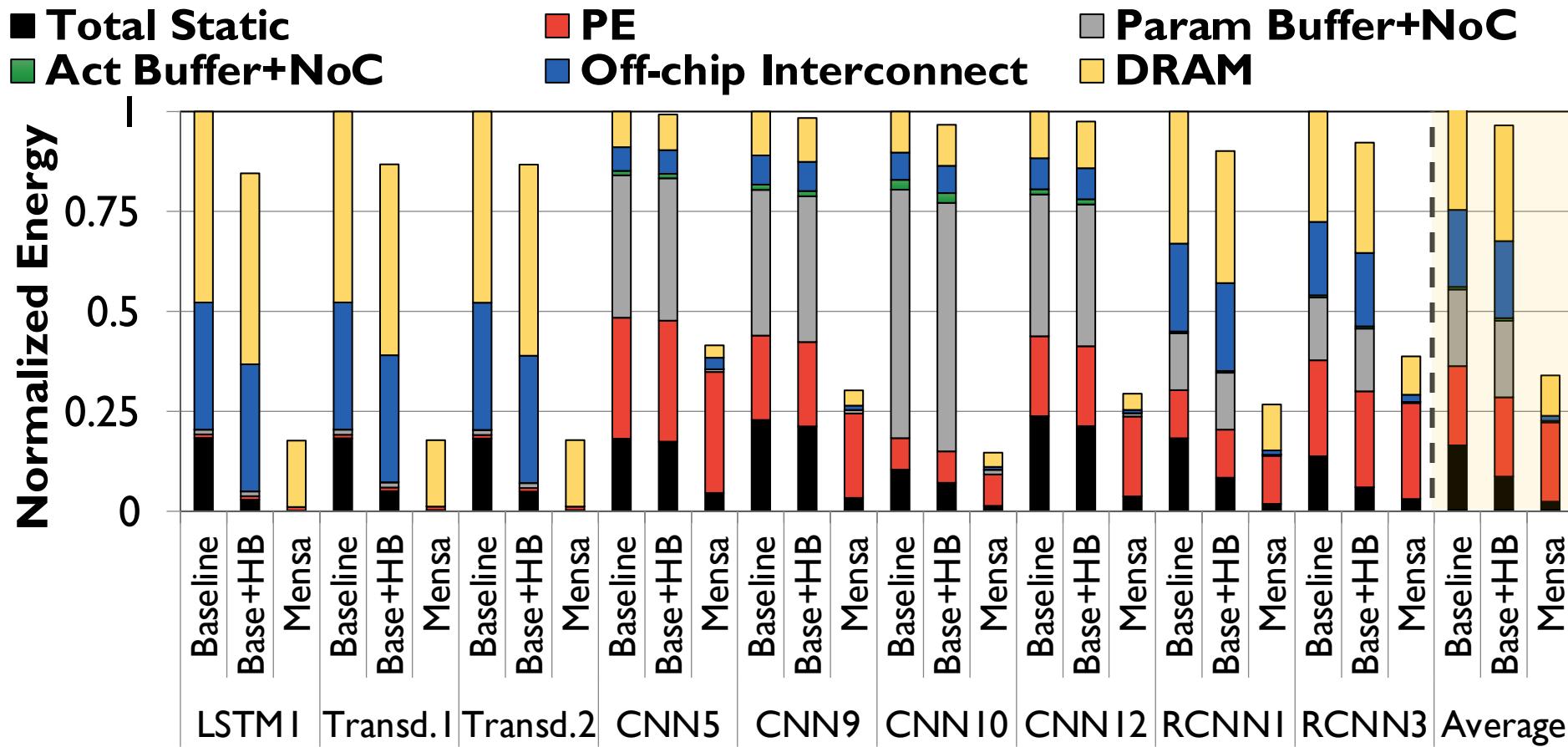
Key observation: the majority of layers group into a small number of layer families



Families 1 & 2: low parameter footprint, high data reuse and MAC intensity
→ compute-centric layers

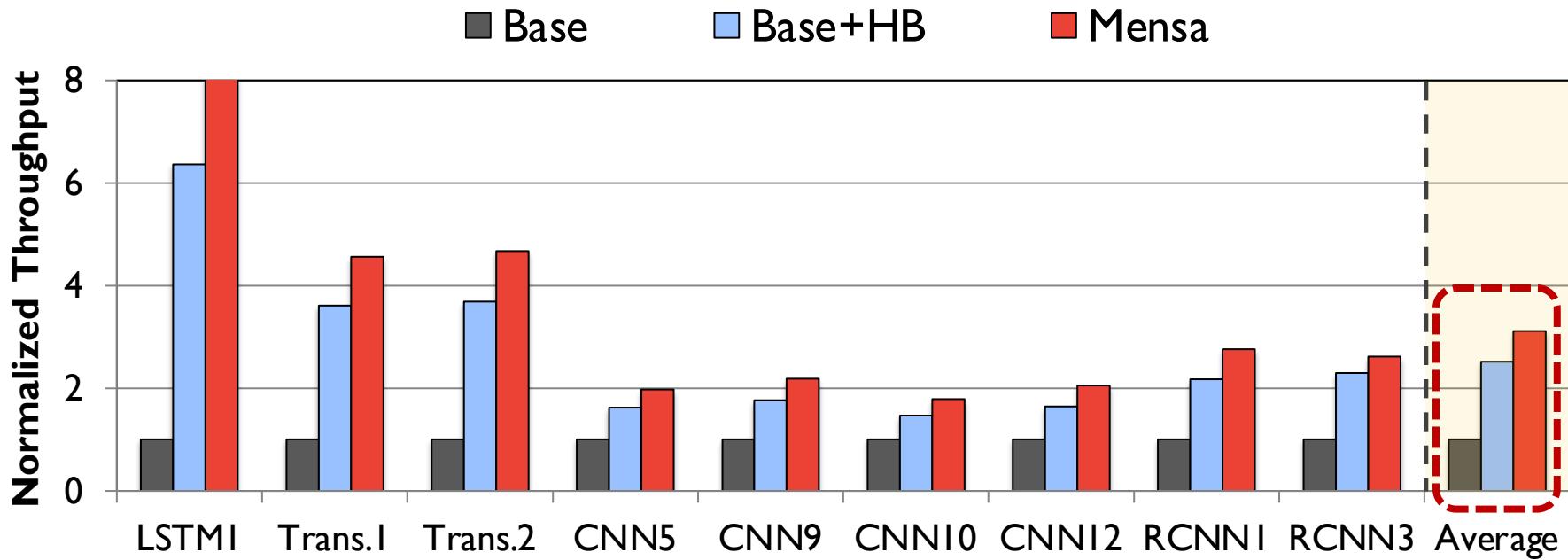
Families 3, 4 & 5: high parameter footprint, low data reuse and MAC intensity
→ data-centric layers

Mensa: Energy Reduction



**Mensa-G reduces energy consumption by 3.0X
compared to the baseline Edge TPU**

Mensa: Throughput Improvement



**Mensa-G improves inference throughput by 3.1X
compared to the baseline Edge TPU**

Mensa: Highly-Efficient ML Inference

- Amirali Boroumand, Saugata Ghose, Berkin Akin, Ravi Narayanaswami, Geraldo F. Oliveira, Xiaoyu Ma, Eric Shiu, and Onur Mutlu,

"Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks"

Proceedings of the 30th International Conference on Parallel Architectures and Compilation Techniques (PACT), Virtual, September 2021.

[Slides (pptx) (pdf)]

[Talk Video (14 minutes)]

Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks

Amirali Boroumand^{†◊}

Geraldo F. Oliveira*

Saugata Ghose[‡]

Xiaoyu Ma[§]

Berkin Akin[§]

Eric Shiu[§]

Ravi Narayanaswami[§]

Onur Mutlu^{*†}

[†]*Carnegie Mellon Univ.*

[◊]*Stanford Univ.*

[‡]*Univ. of Illinois Urbana-Champaign*

[§]*Google*

^{*}*ETH Zürich*

Automatic Code and Data Mapping

- Kevin Hsieh, Eiman Ebrahimi, Gwangsun Kim, Niladrish Chatterjee, Mike O'Connor, Nandita Vijaykumar, Onur Mutlu, and Stephen W. Keckler,

"Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems"

Proceedings of the 43rd International Symposium on Computer Architecture (ISCA), Seoul, South Korea, June 2016.

[Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems

Kevin Hsieh[†] Eiman Ebrahimi[†] Gwangsun Kim^{*} Niladrish Chatterjee[†] Mike O'Connor[†]
Nandita Vijaykumar[†] Onur Mutlu^{§‡} Stephen W. Keckler[†]

[†]Carnegie Mellon University [†]NVIDIA ^{*}KAIST [§]ETH Zürich

Automatic Offloading of Critical Code

- Milad Hashemi, Khubaib, Eiman Ebrahimi, Onur Mutlu, and Yale N. Patt,
"Accelerating Dependent Cache Misses with an Enhanced Memory Controller"

Proceedings of the 43rd International Symposium on Computer Architecture (ISCA), Seoul, South Korea, June 2016.

[Slides (pptx) (pdf)]

[Lightning Session Slides (pptx) (pdf)]

Accelerating Dependent Cache Misses with an Enhanced Memory Controller

Milad Hashemi*, Khubaib†, Eiman Ebrahimi‡, Onur Mutlu§, Yale N. Patt*

*The University of Texas at Austin †Apple ‡NVIDIA §ETH Zürich & Carnegie Mellon University

Automatic Offloading of Prefetch Mechanisms

- Milad Hashemi, Onur Mutlu, and Yale N. Patt,
"Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads"
Proceedings of the 49th International Symposium on Microarchitecture (MICRO), Taipei, Taiwan, October 2016.
[Slides (pptx) (pdf)] [Lightning Session Slides (pdf)] [Poster (pptx) (pdf)]

Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads

Milad Hashemi*, Onur Mutlu[§], Yale N. Patt*

**The University of Texas at Austin* [§]*ETH Zürich*

Adoption: How to Keep It Simple?

- Junwhan Ahn, Sungjoo Yoo, Onur Mutlu, and Kiyoung Choi,
"PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture"
Proceedings of the 42nd International Symposium on Computer Architecture (ISCA), Portland, OR, June 2015.
[Slides (pdf)] [Lightning Session Slides (pdf)]

PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture

Junwhan Ahn Sungjoo Yoo Onur Mutlu[†] Kiyoung Choi

junwhan@snu.ac.kr, sungjoo.yoo@gmail.com, onur@cmu.edu, kchoi@snu.ac.kr

Seoul National University

[†]Carnegie Mellon University

Adoption: How to Maintain Coherence? (I)

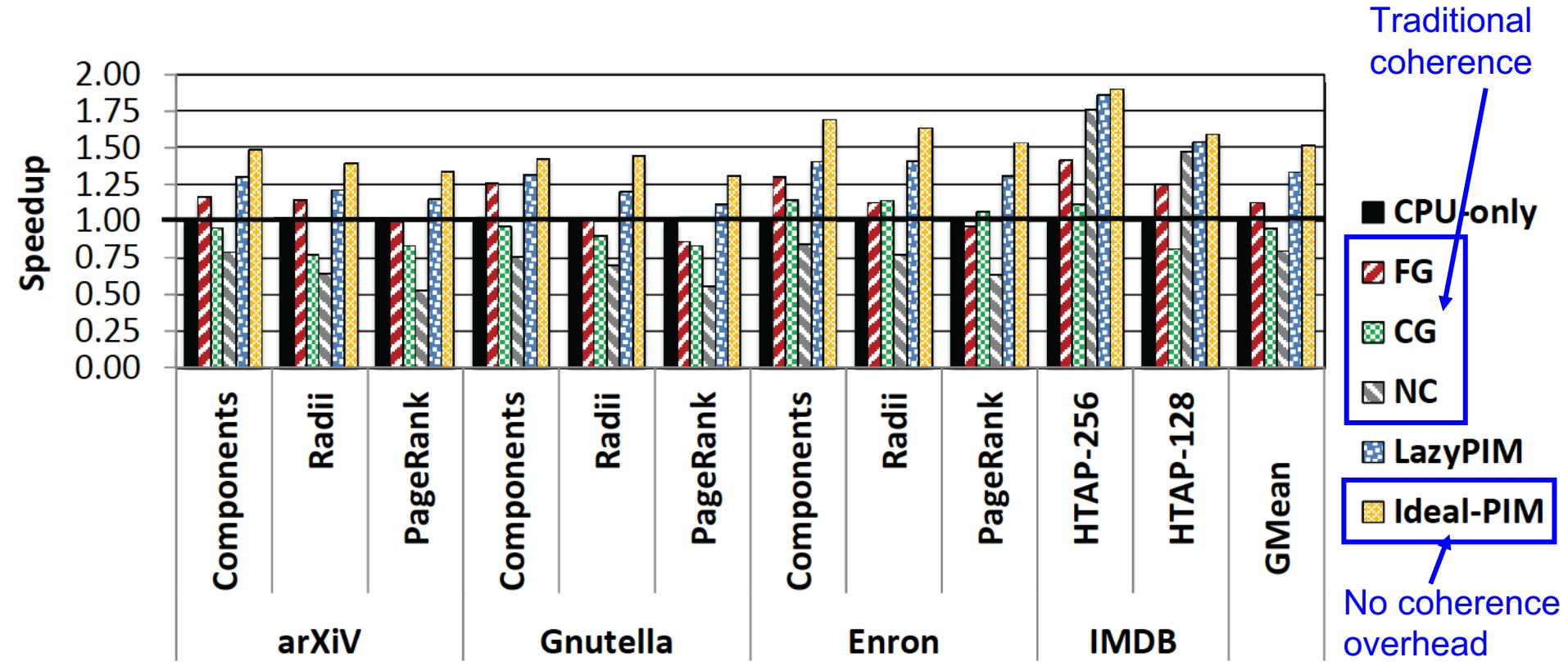
- Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu,
"LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory"
IEEE Computer Architecture Letters (CAL), June 2016.

LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory

Amirali Boroumand[†], Saugata Ghose[†], Minesh Patel[†], Hasan Hassan^{†§}, Brandon Lucia[†],
Kevin Hsieh[†], Krishna T. Malladi^{*}, Hongzhong Zheng^{*}, and Onur Mutlu^{‡†}

[†]*Carnegie Mellon University* ^{*}*Samsung Semiconductor, Inc.* [§]*TOBB ETÜ* [‡]*ETH Zürich*

Challenge: Coherence for Hybrid CPU-PIM Apps



Adoption: How to Maintain Coherence? (II)

- Amirali Boroumand, Saugata Ghose, Minesh Patel, Hasan Hassan, Brandon Lucia, Kevin Hsieh, Krishna T. Malladi, Hongzhong Zheng, and Onur Mutlu,
"CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators"

Proceedings of the 46th International Symposium on Computer Architecture (ISCA), Phoenix, AZ, USA, June 2019.

CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators

Amirali Boroumand[†]

Brandon Lucia[†]

Nastaran Hajinazar^{◦†}

Saugata Ghose[†]

Rachata Ausavarungnirun^{†‡}

Krishna T. Malladi[§]

Minesh Patel^{*}

Kevin Hsieh[†]

Hongzhong Zheng[§]

Hasan Hassan^{*}

Onur Mutlu^{★†}

[†]Carnegie Mellon University

[◦]Simon Fraser University

^{*}ETH Zürich

[‡]KMUTNB

[§]Samsung Semiconductor, Inc.

Adoption: How to Support Synchronization?

- Christina Giannoula, Nandita Vijaykumar, Nikela Papadopoulou, Vasileios Karakostas, Ivan Fernandez, Juan Gómez-Luna, Lois Orosa, Nectarios Koziris, Georgios Goumas, Onur Mutlu,

"SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures"

Proceedings of the 27th International Symposium on High-Performance Computer Architecture (HPCA), Virtual, February-March 2021.

[Slides (pptx) (pdf)]

[Short Talk Slides (pptx) (pdf)]

[Talk Video (21 minutes)]

[Short Talk Video (7 minutes)]

SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures

Christina Giannoula^{†‡} Nandita Vijaykumar^{*‡} Nikela Papadopoulou[†] Vasileios Karakostas[†] Ivan Fernandez^{§‡}
Juan Gómez-Luna[‡] Lois Orosa[‡] Nectarios Koziris[†] Georgios Goumas[†] Onur Mutlu[‡]

[†]*National Technical University of Athens* [‡]*ETH Zürich* ^{*}*University of Toronto* [§]*University of Malaga*

Adoption: How to Support Virtual Memory?

- Kevin Hsieh, Samira Khan, Nandita Vijaykumar, Kevin K. Chang, Amirali Boroumand, Saugata Ghose, and Onur Mutlu,

**"Accelerating Pointer Chasing in 3D-Stacked Memory:
Challenges, Mechanisms, Evaluation"**

*Proceedings of the 34th IEEE International Conference on Computer
Design (ICCD), Phoenix, AZ, USA, October 2016.*

Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation

Kevin Hsieh[†] Samira Khan[‡] Nandita Vijaykumar[†]
Kevin K. Chang[†] Amirali Boroumand[†] Saugata Ghose[†] Onur Mutlu^{§†}
[†]*Carnegie Mellon University* [‡]*University of Virginia* [§]*ETH Zürich*

Eliminating the Adoption Barriers

Processing-in-Memory in the Real World

Fundamentally Energy-Efficient (Data-Centric) Computing Architectures

Fundamentally High-Performance (Data-Centric) Computing Architectures

Computing Architectures with Minimal Data Movement

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 near Memory
 - Processing using Memory
- How to Enable Adoption of Processing in Memory
- Conclusion

Eliminating the Adoption Barriers

How to Enable Adoption of Processing in Memory

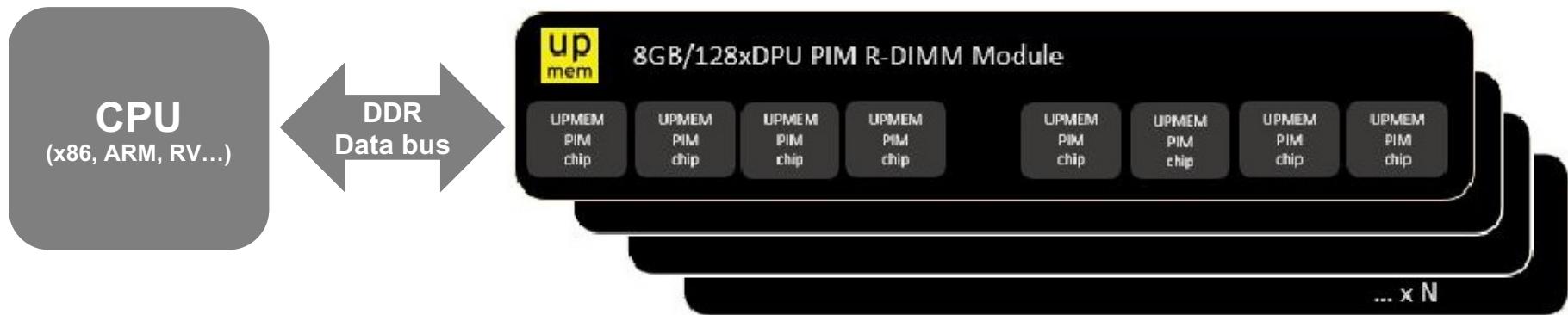
Potential Barriers to Adoption of PIM

1. **Applications & software** for PIM
2. Ease of **programming** (interfaces and compiler/HW support)
3. **System** and **security** support: coherence, synchronization, virtual memory, isolation, communication interfaces, ...
4. **Runtime** and **compilation** systems for adaptive scheduling, data mapping, access/sharing control, ...
5. **Infrastructures** to assess benefits and feasibility

All can be solved with change of mindset

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-LUNA¹, IZZAT EL HAJJ², IVAN FERNANDEZ^{1,3}, CHRISTINA GIANNOULA^{1,4},
GERALDO F. OLIVEIRA¹, AND ONUR MUTLU¹

¹ETH Zürich

²American University of Beirut

³University of Malaga

⁴National 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 page includes a header with repository statistics (Unwatch 2, Star 2, Fork 1), a navigation bar with links like Code, Issues, Pull requests, Actions, Projects, Wiki, Security, Insights, and Settings, and a main content area for the README.md file. The README.md content discusses the PrIM benchmark suite, its purpose, and its characteristics.

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

PrIM (Processing-In-Memory Benchmarks)

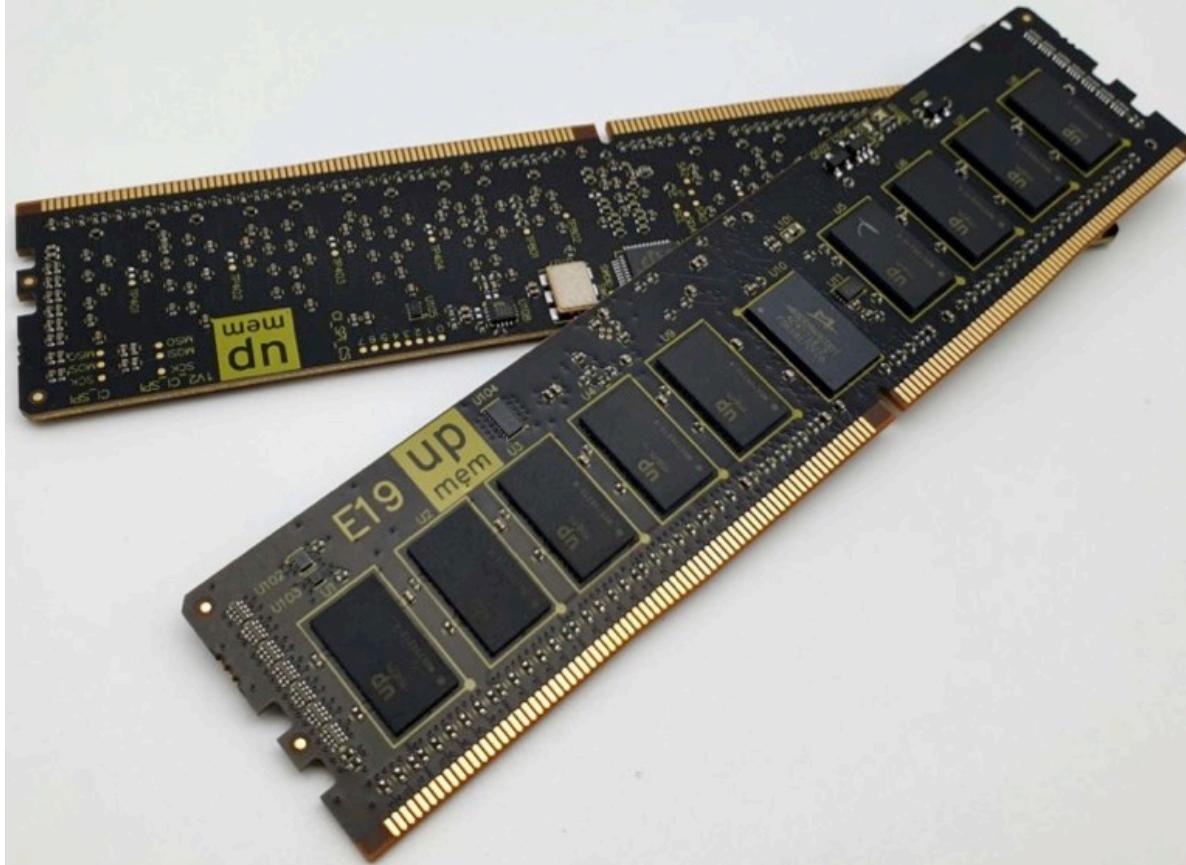
PrIM is the first benchmark suite for a real-world processing-in-memory (PIM) architecture. PrIM is developed to evaluate, analyze, and characterize the first publicly-available real-world processing-in-memory (PIM) architecture, the UPMEM PIM architecture. The UPMEM PIM architecture combines traditional 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	~700W	~300W
speed-up		~ x20	x1
energy gain		~ x10	x1
TCO gain		~ x10	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

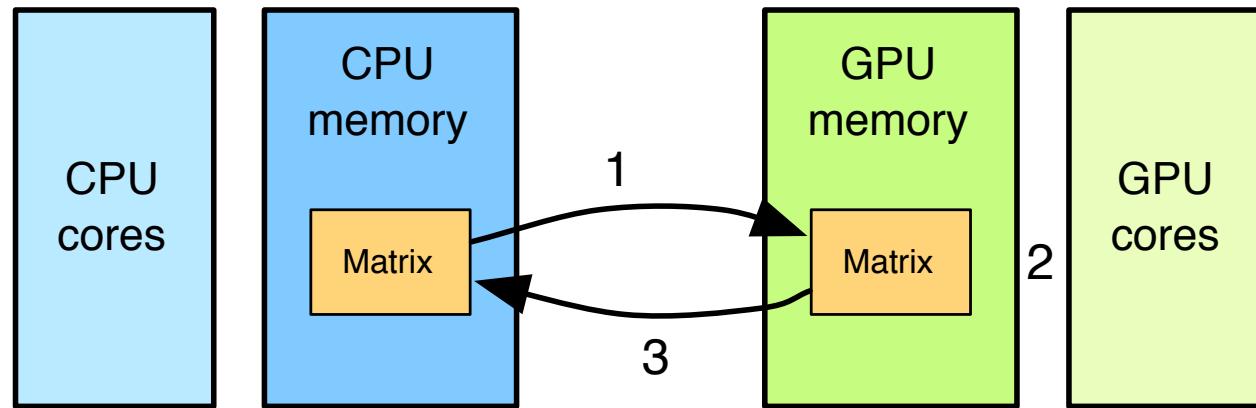
- Introduction
 - Accelerator Model
 - UPMEM-based PIM System Overview
- UPMEM PIM Programming
 - Vector Addition
 - CPU-DPU Data Transfers
 - Inter-DPU Communication
 - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
 - Arithmetic Throughput
 - WRAM and MRAM Bandwidth
- PrIM Benchmarks
 -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

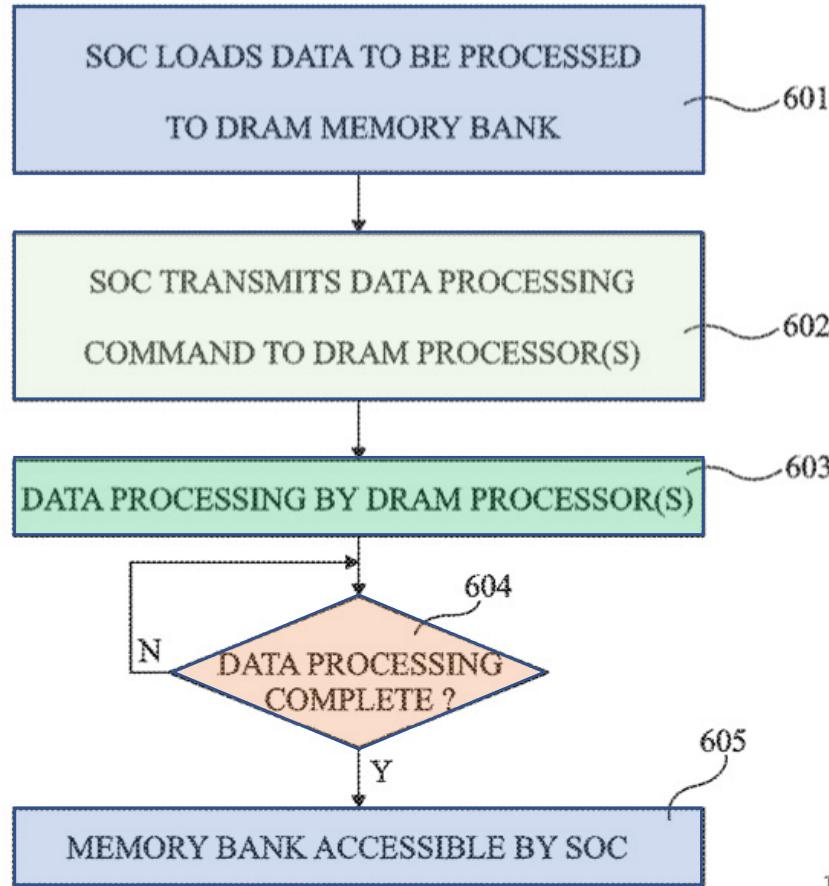


Fig 6

System Organization (I)

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

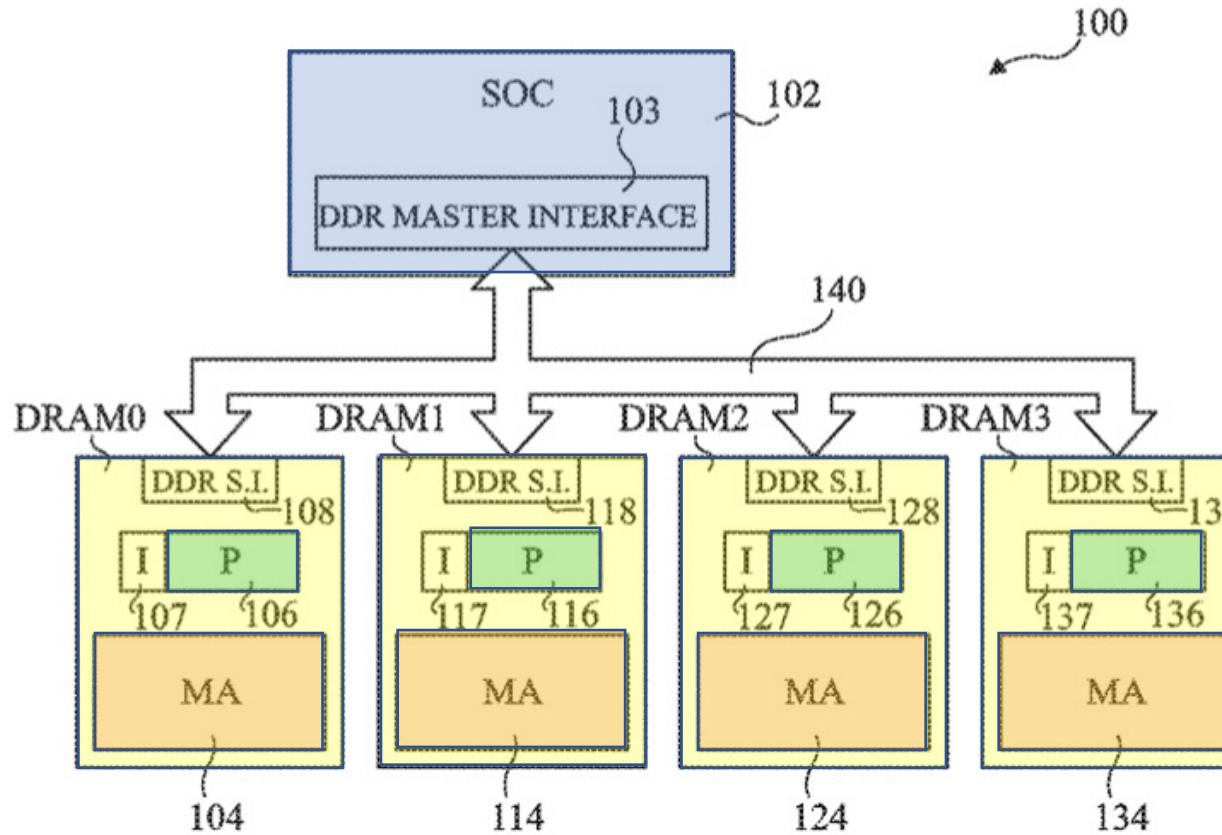
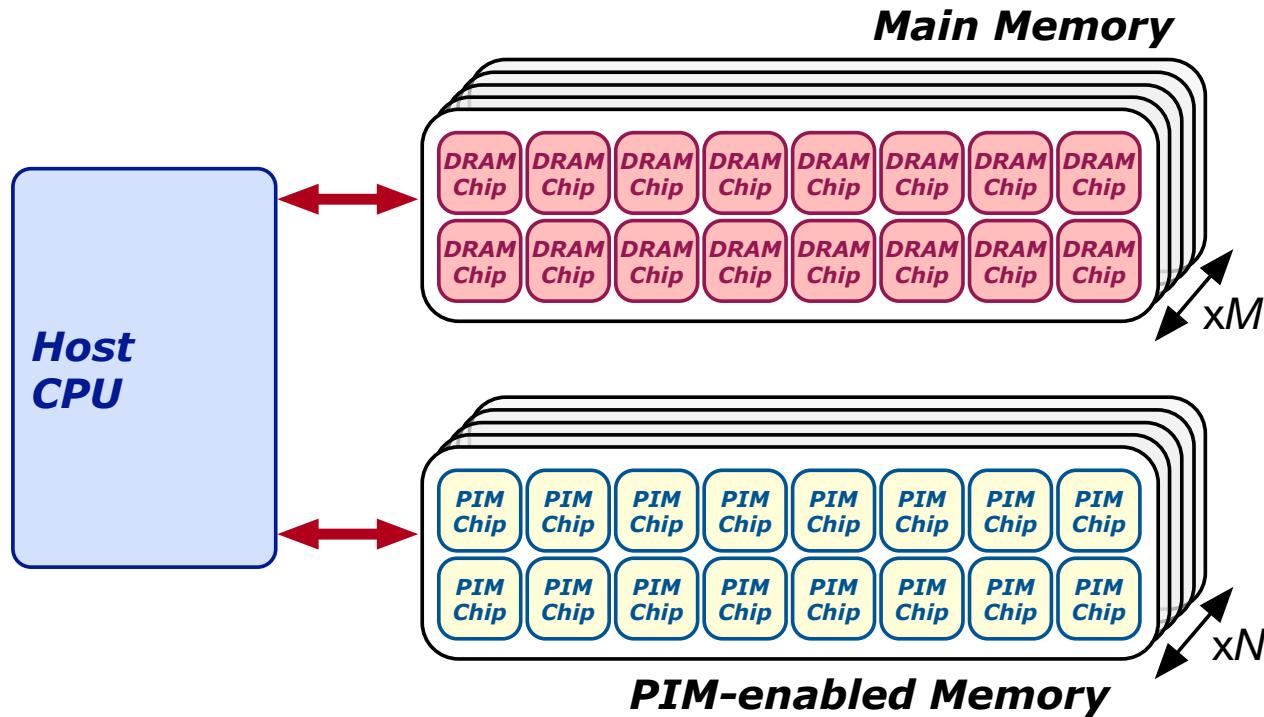


Fig 1

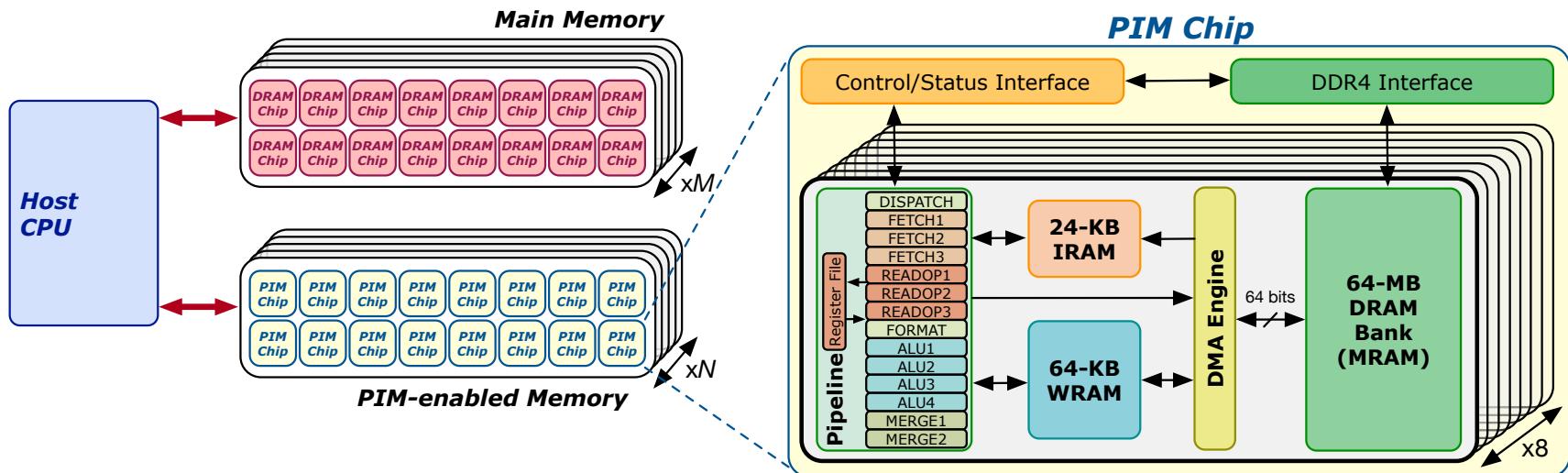
System Organization (II)

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



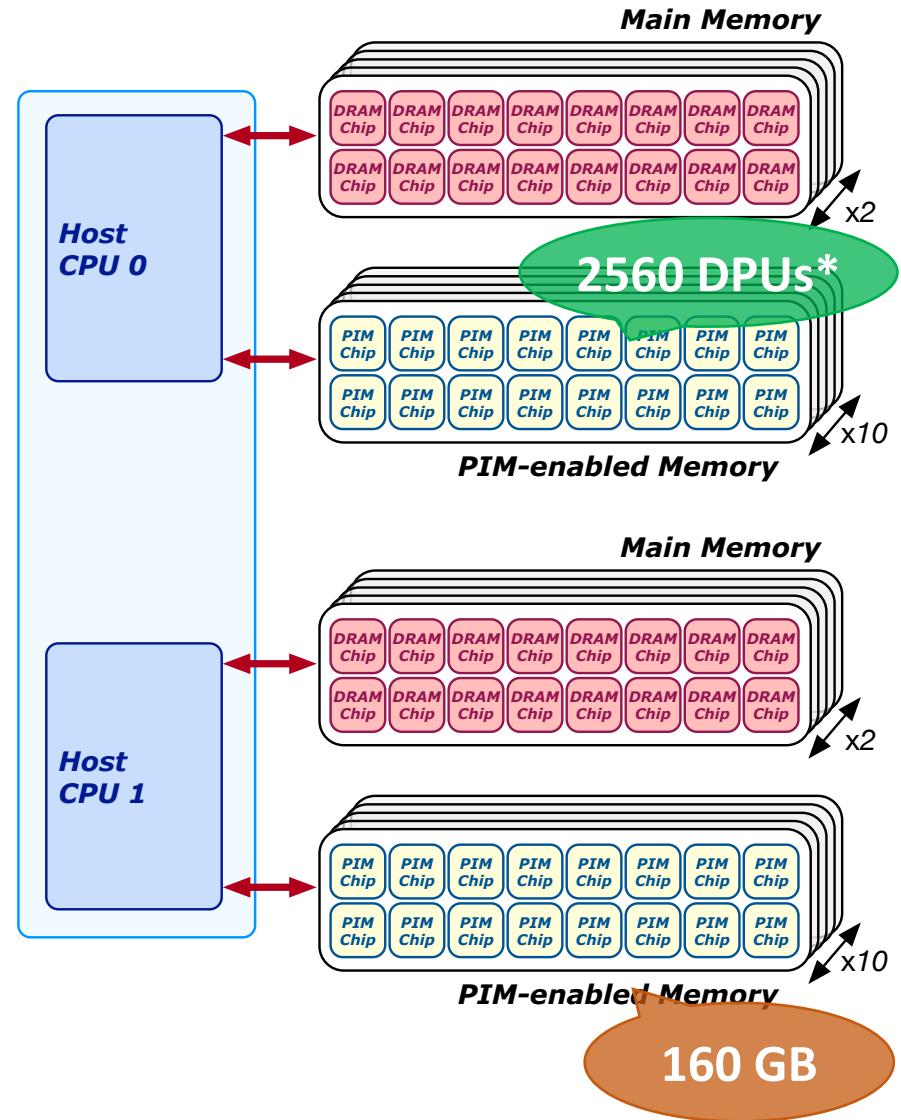
System Organization (III)

- 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

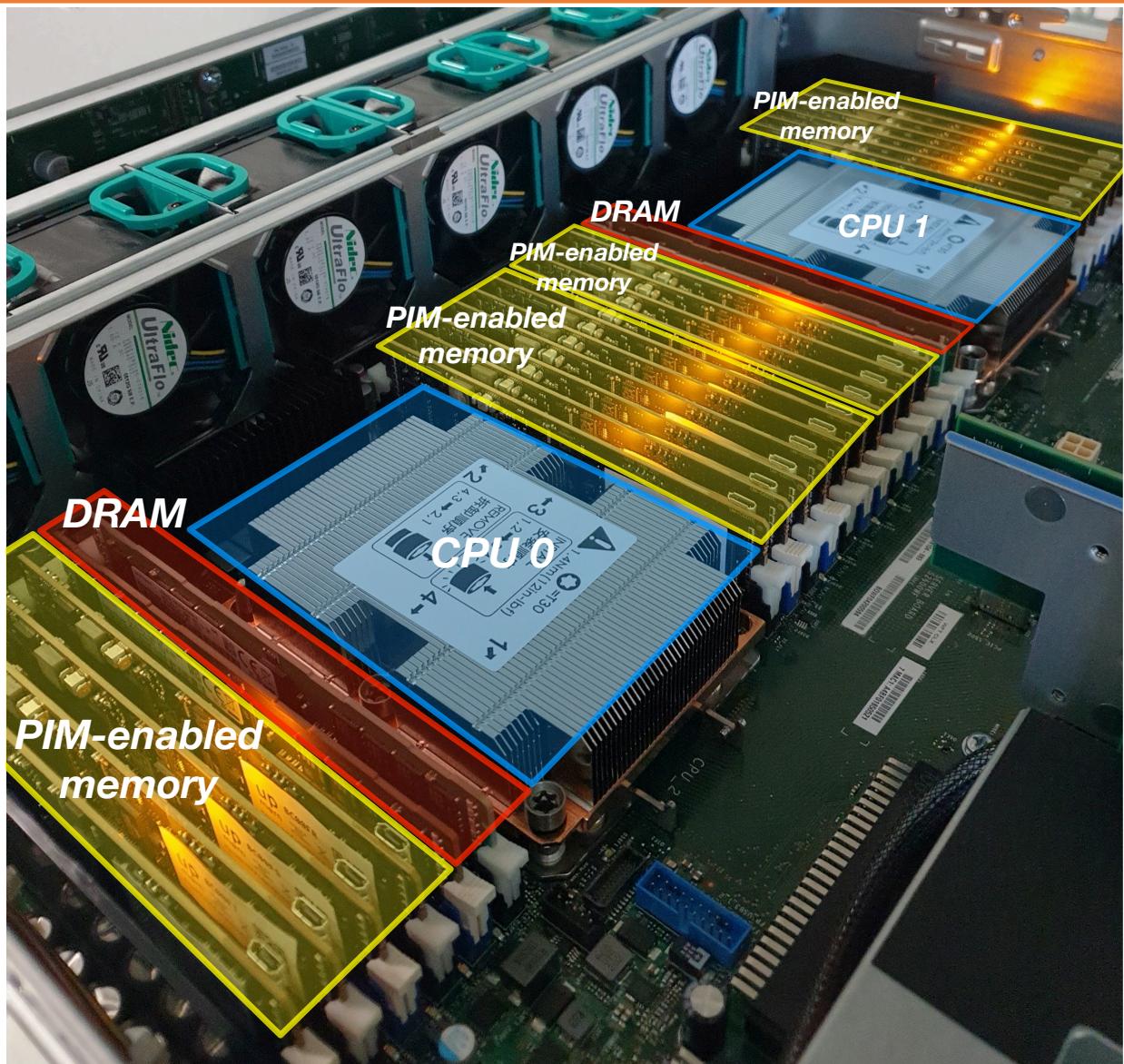
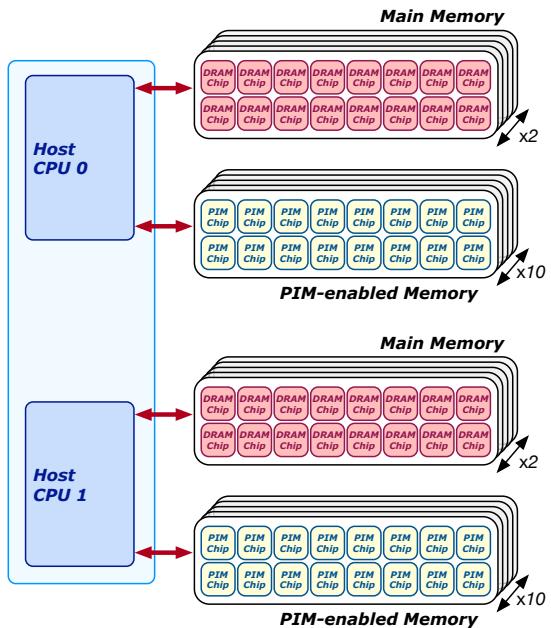


2,560-DPU System (I)

- 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



2,560-DPU System (II)

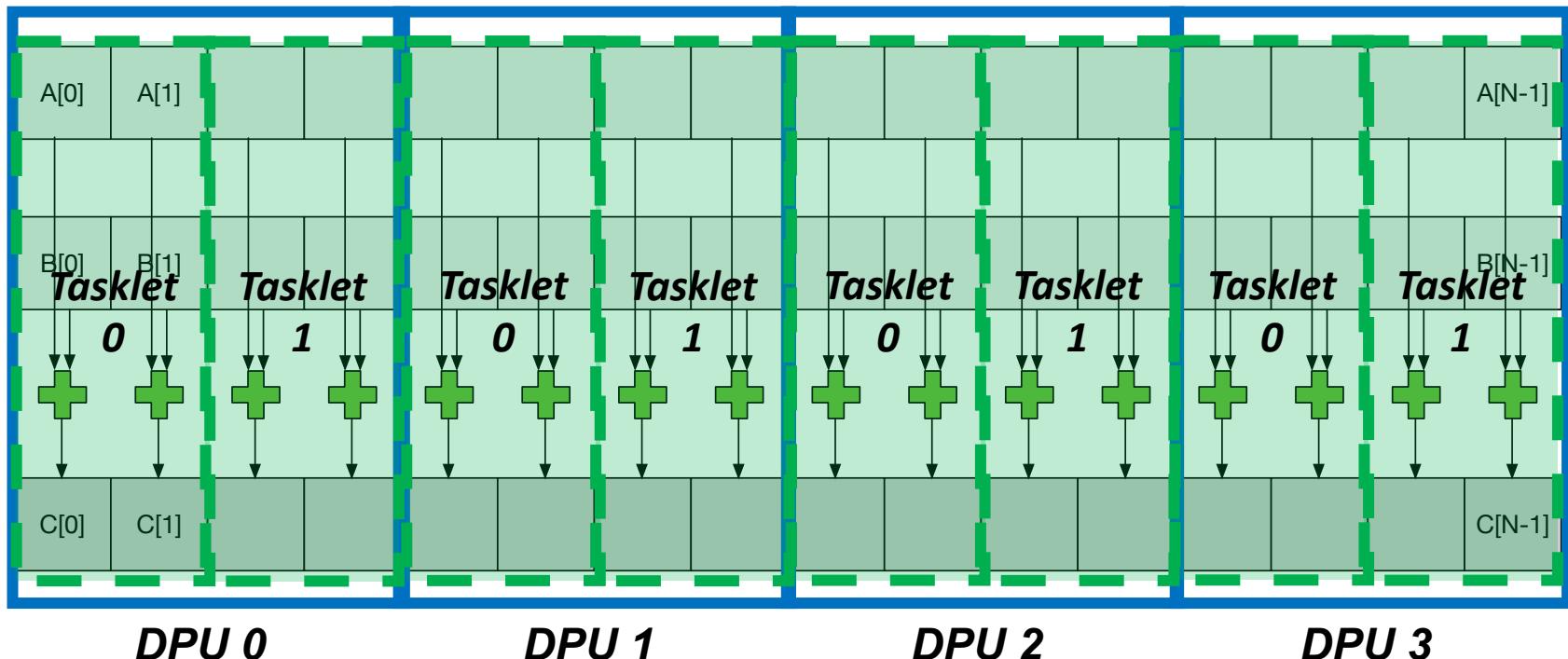


Outline

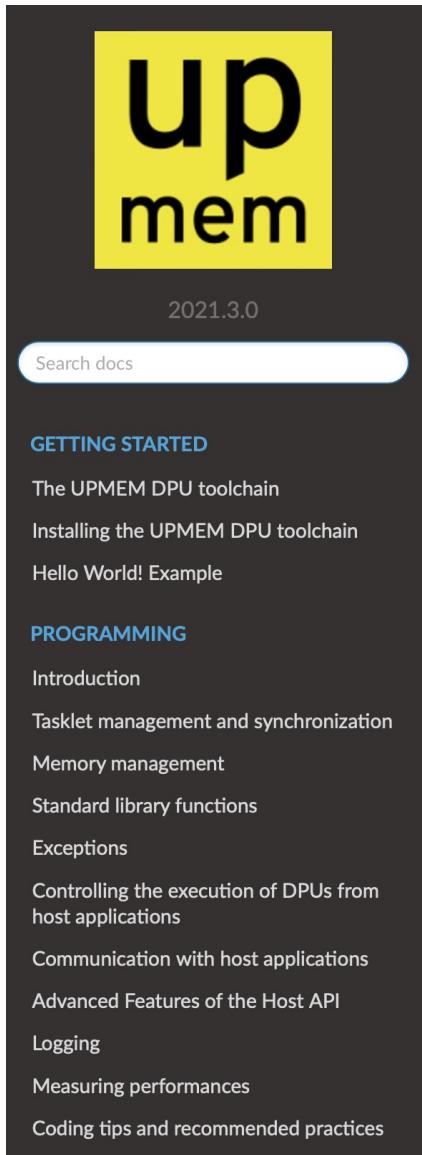
- Introduction
 - Accelerator Model
 - UPMEM-based PIM System Overview
- UPMEM PIM Programming
 - Vector Addition
 - CPU-DPU Data Transfers
 - Inter-DPU Communication
 - CPU-DPU/DPU-CPU Transfer Bandwidth
- DRAM Processing Unit
 - Arithmetic Throughput
 - WRAM and MRAM Bandwidth
- PrIM Benchmarks
 -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

Further Slides for Your Own Study
(May Be Covered in Future Lectures)

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
```

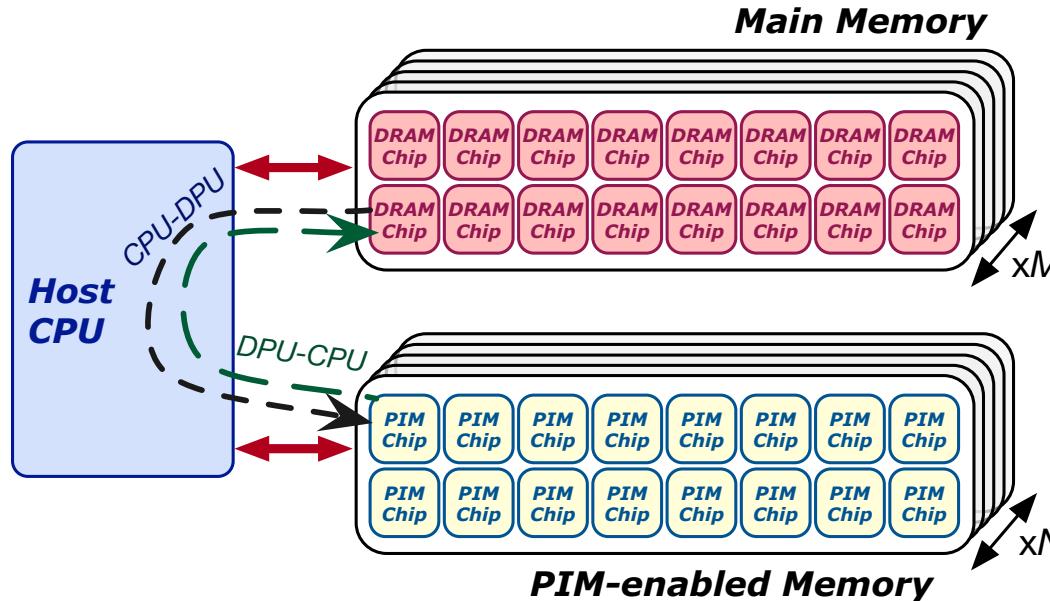
Is it possible to launch different kernels onto different DPUs?

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) { Pointer to main memory
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) { Offset within MRAM
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
```

Direction

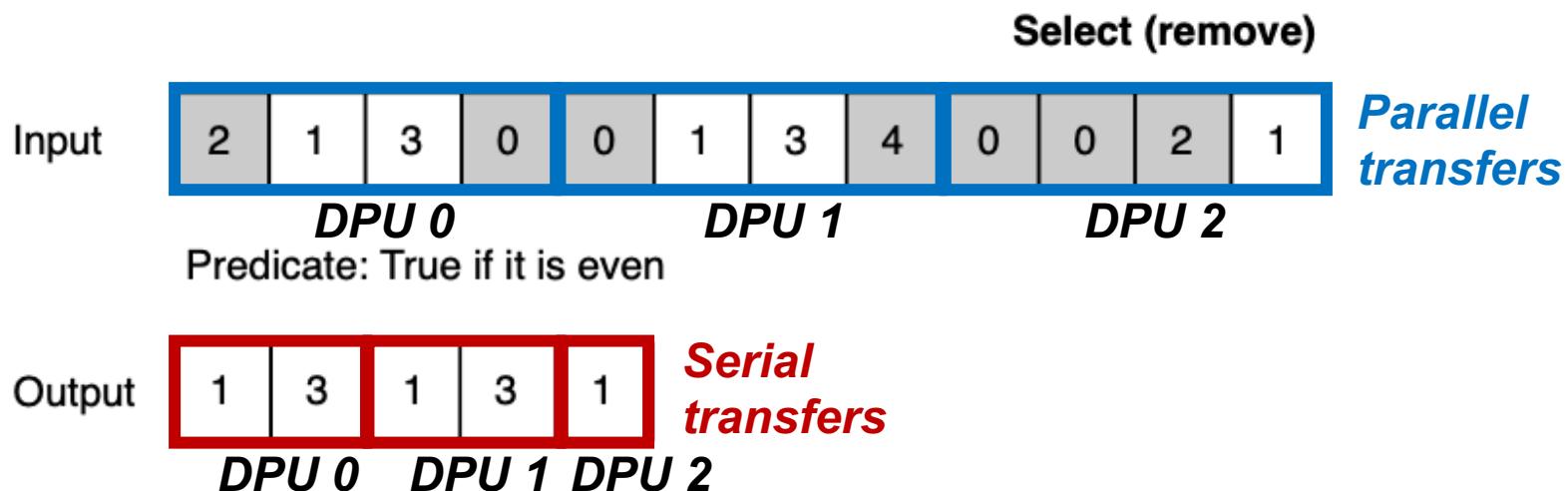
Broadcast Transfers

- `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
```

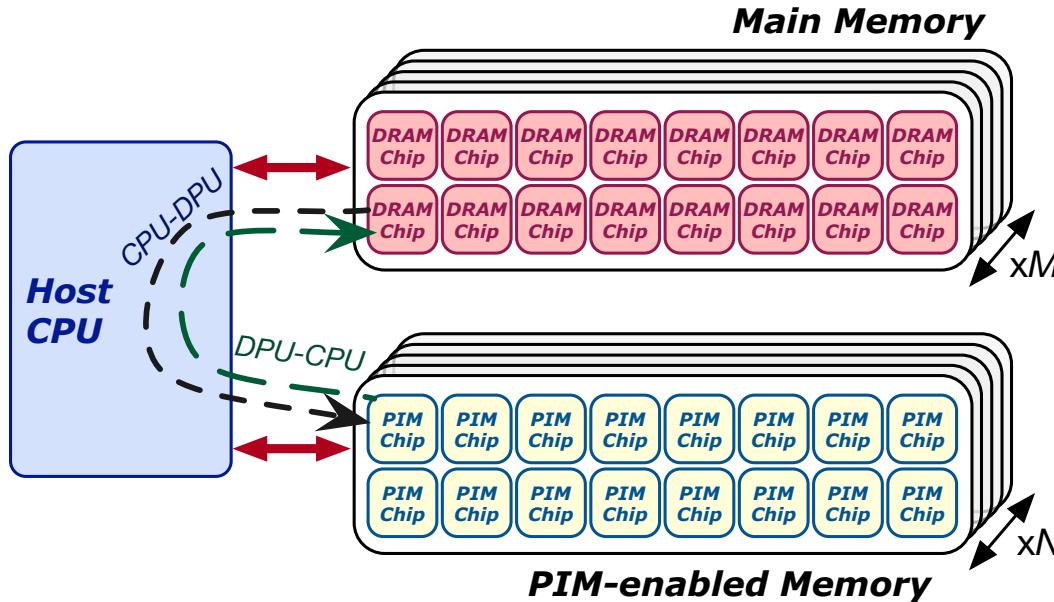
Different Types of Transfers in a Program

- 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?

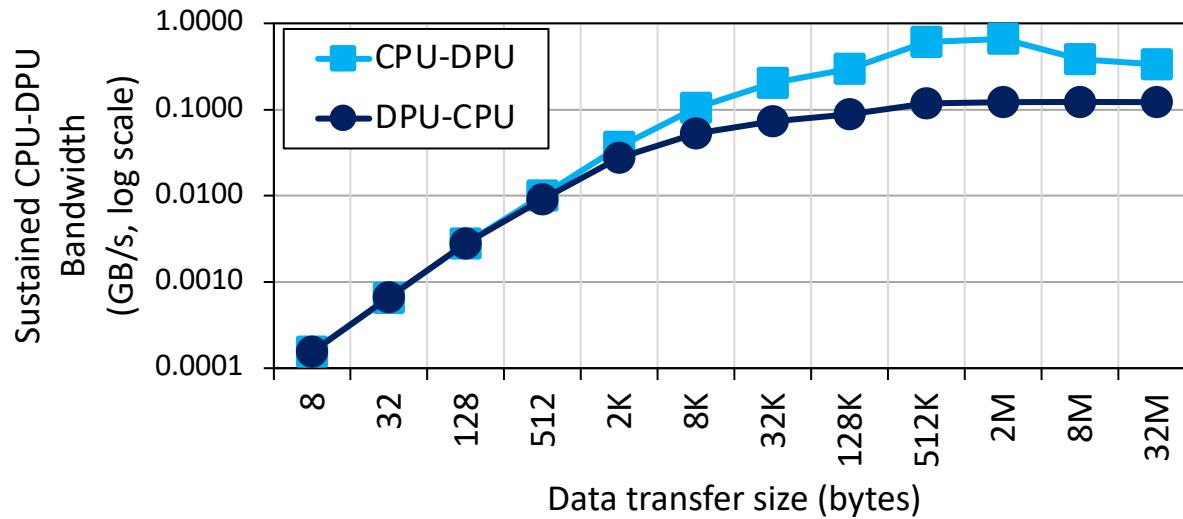
- With a microbenchmark, we obtain the **sustained bandwidth of all types of CPU-DPU and DPU-CPU transfers**
- Two experiments:
 - 1 DPU: variable CPU-DPU and DPU-CPU transfer size (**8 bytes to 32 MB**)
 - 1 rank: 32 MB CPU-DPU and DPU-CPU transfers to/from a set of **1 to 64 MRAM banks** within the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

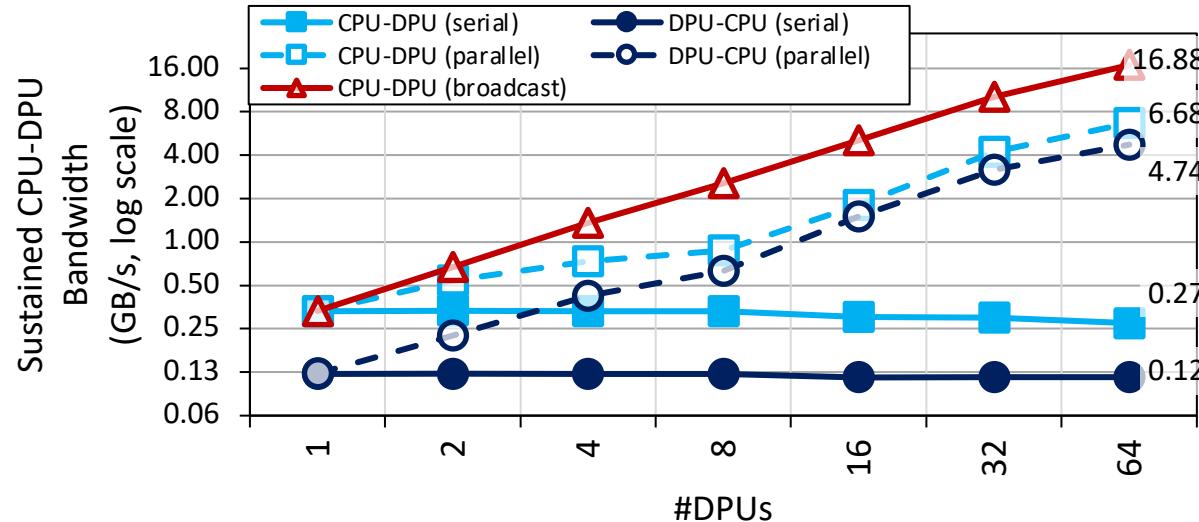


KEY OBSERVATION 7

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

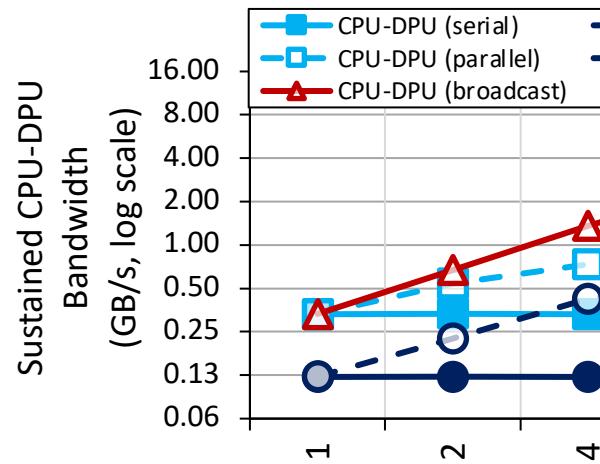


KEY OBSERVATION 8

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

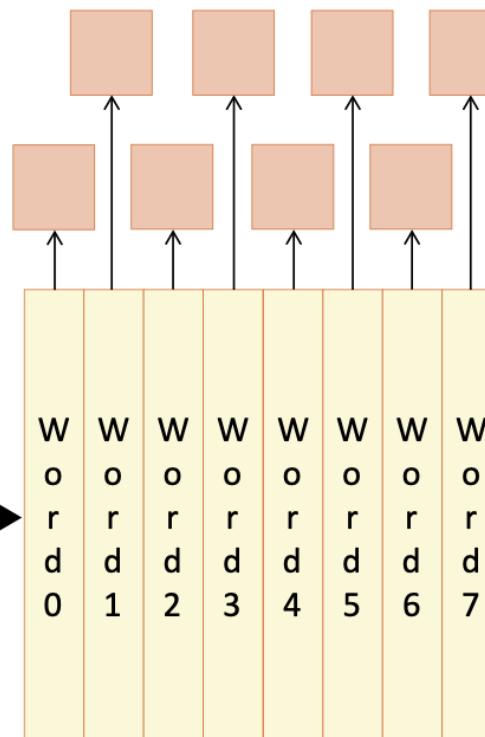
The library feeds DPUs with correct data

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)

The screenshot shows a GitHub repository page for **CMU-SAFARI / prim-benchmarks**. The repository has 2 stars and 0 forks. The main navigation tabs are Code, Issues, Pull requests, Actions, Projects, Wiki, Security, Insights, and Settings. Below the tabs, there's a breadcrumb navigation: **main** > **prim-benchmarks** / **Microbenchmarks** / **CPU-DPU** / . The commit details are as follows:

Commit	Author	Date	Actions
Juan Gomez Luna PrIM -- first commit	Juan Gomez Luna	3de4b49 7 days ago	History
..			
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?

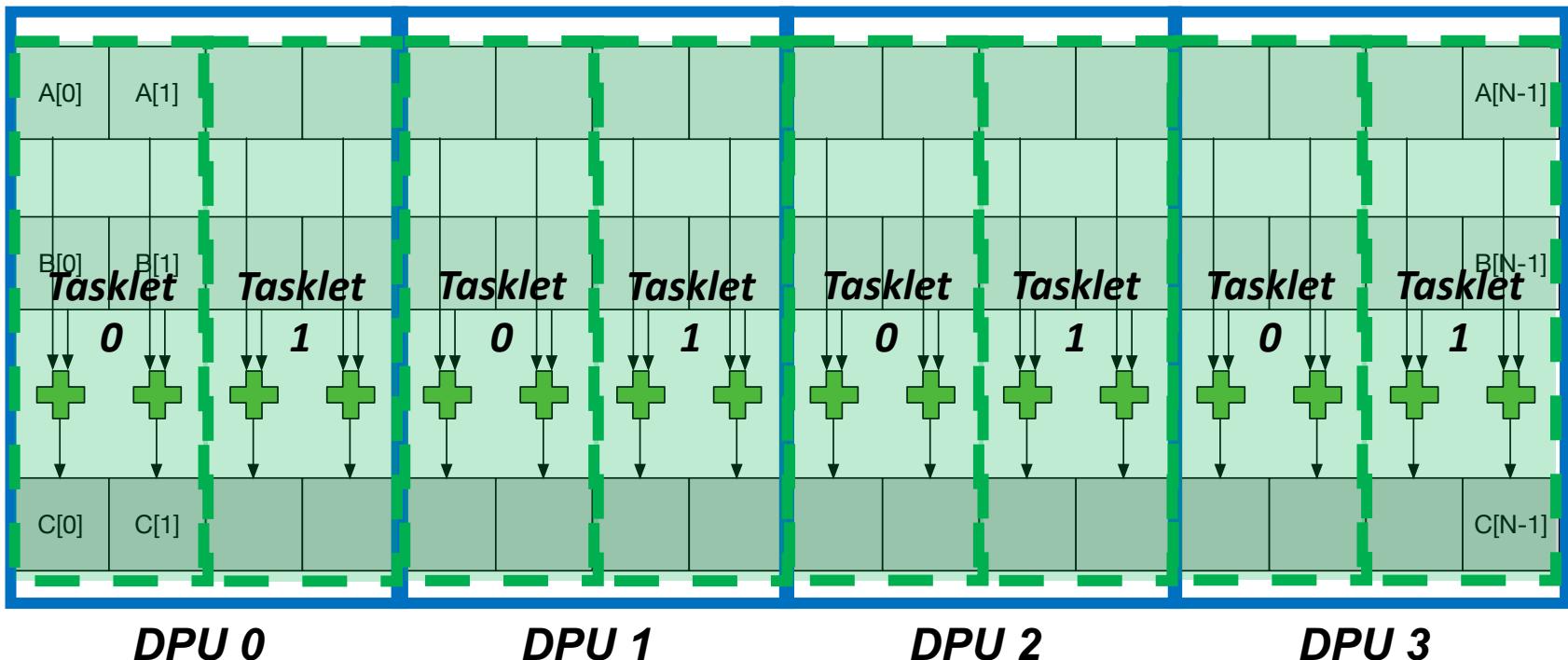
- We can use serial and parallel transfers
- We pass them directly to the scratchpad memory of the DPU
 - Working RAM (WRAM): We introduce it in the next slides
- This is useful for **input parameters and some results**

```
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
```

```
1 // Host code (host/app.c)
2 #ifdef SERIAL
3     DPU_FOREACH(dpu_set, dpu) {
4         DPU_ASSERT(dpu_copy_to(dpu, "DPU_INPUT_ARGUMENTS", 0, (const void *)&input_arguments[i], sizeof(input_arguments[0])));
5         i++;
6     }
7 #else
8     DPU_FOREACH(dpu_set, dpu, i) {
9         DPU_ASSERT(dpu_prepare_xfer(dpu, &input_arguments[i]));
10    }
11 DPU_ASSERT(dpu_push_xfer(dpu_set, DPU_XFER_TO_DPU, "DPU_INPUT_ARGUMENTS", 0, sizeof(input_arguments[0]), DPU_XFER_DEFAULT));
12
13#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);                         WWRAM 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-WWRAM 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);           WWRAM-MRAM DMA transfer
29
30    }
31
32    return 0;
}
```

Programming a DPU Kernel (II)

- Vector addition

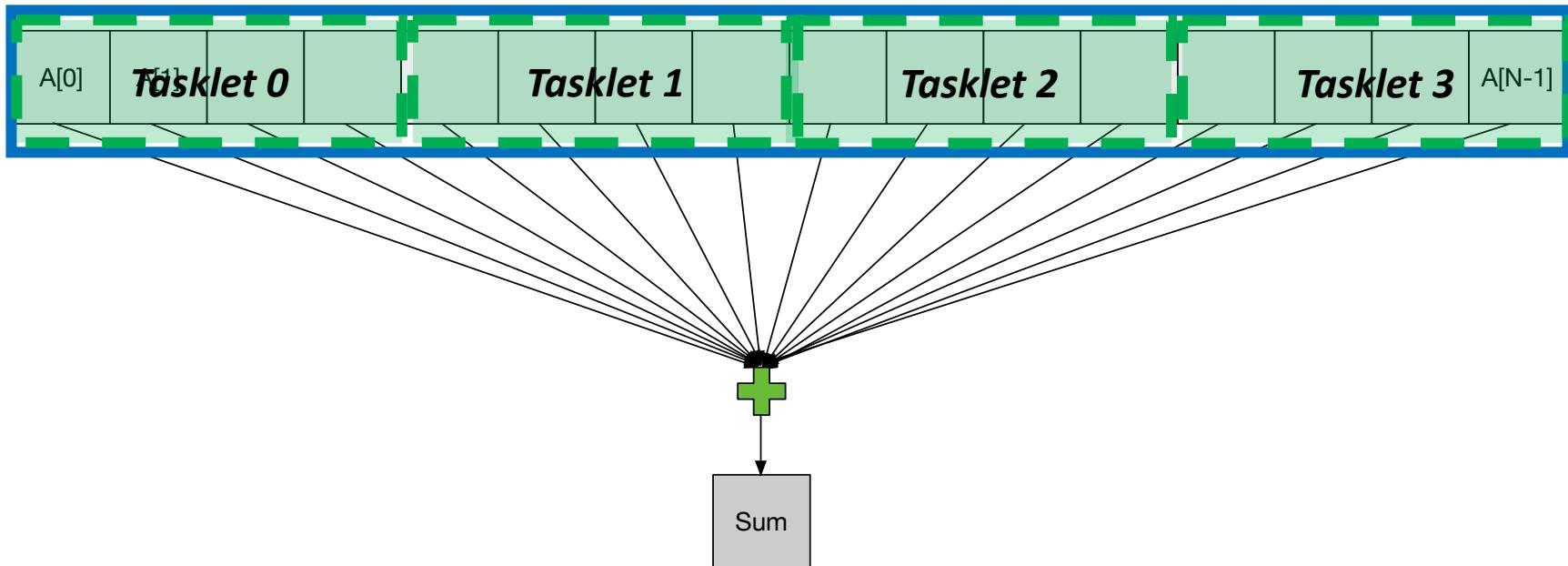
```
1 // vector_addition: Computes the vector addition of a cached block
2 ▼ static void vector_addition(T *bufferB, T *bufferA, unsigned int l_size) {
3
4 ▼     for (unsigned int i = 0; i < l_size; i++){
5         bufferB[i] += bufferA[i];
6 ▲     }
7
8 ▲ }
```

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();`

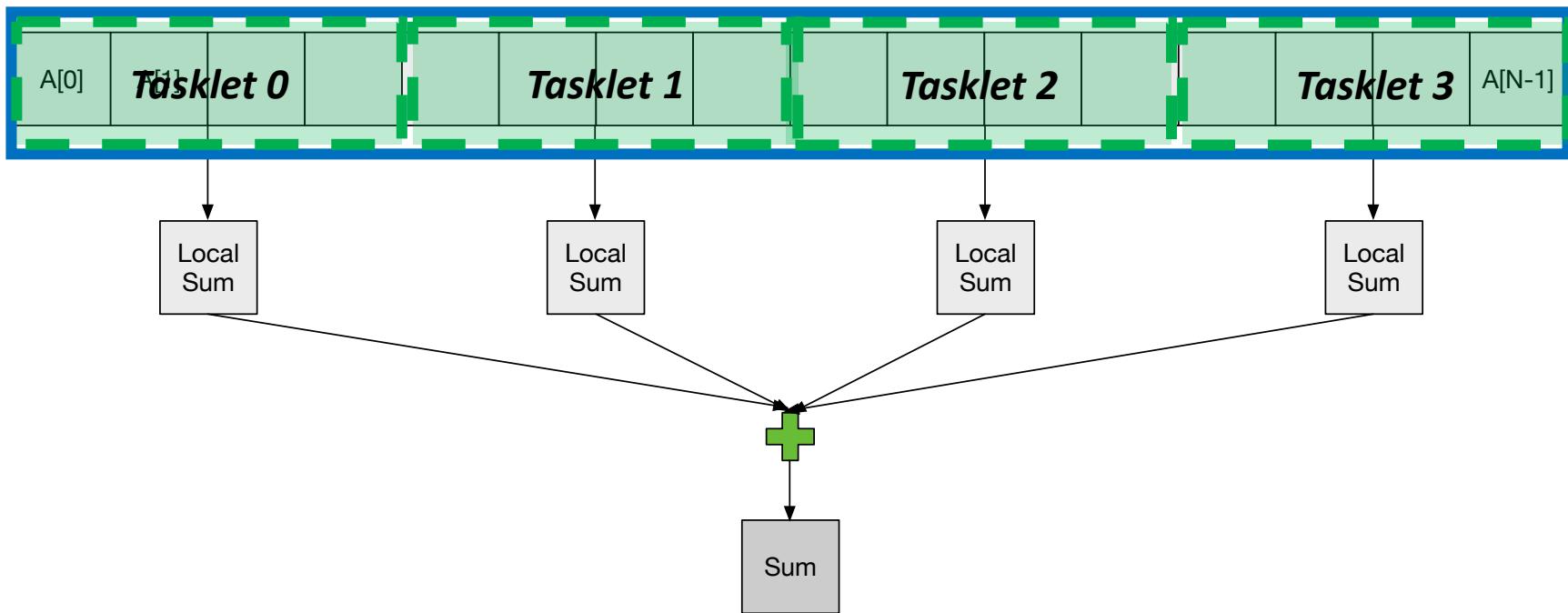
Parallel Reduction (I)

- 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

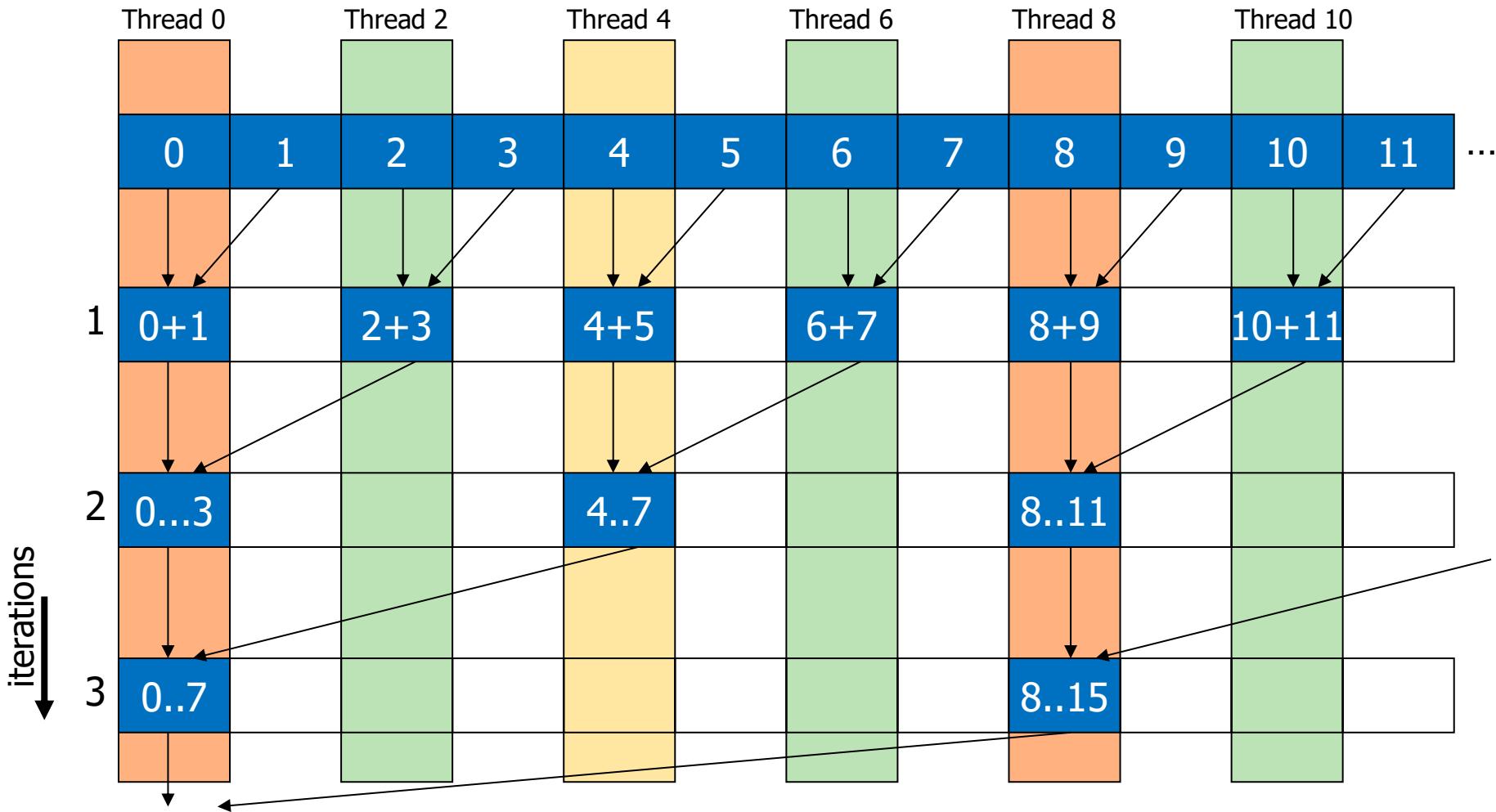
```
1▼ for(unsigned int byte_index = base_tasklet; byte_index < input_size_dpu_bytes; byte_index += BLOCK_SIZE * NR_TASKLETS){  
2  
3    // Bound checking  
4    uint32_t l_size_bytes = (byte_index + BLOCK_SIZE >= input_size_dpu_bytes) ? (input_size_dpu_bytes - byte_index) : BLOCK_SIZE;  
5  
6    // 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

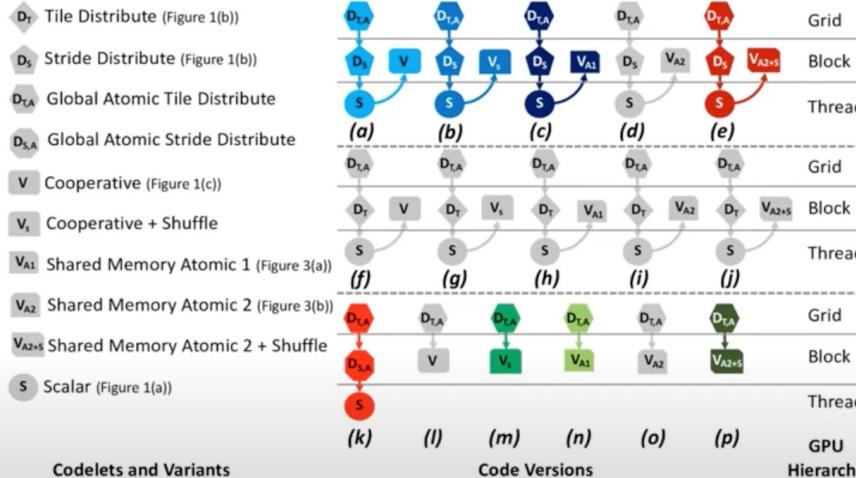
- Multiple tasklets can perform a tree-based reduction
 - After every iteration tasklets synchronize with a barrier
 - Half of the tasklets retire at the end of an iteration

```
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!

Garcia de Gonzalo et al., "Automatic Generation of Warp-Level Primitives and Atomic Instructions for Fast and Portable Parallel Reduction on GPUs," CGO 2019

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

201 views • Premiered Apr 19, 2022

15 DISLIKE SHARE CLIP SAVE ...



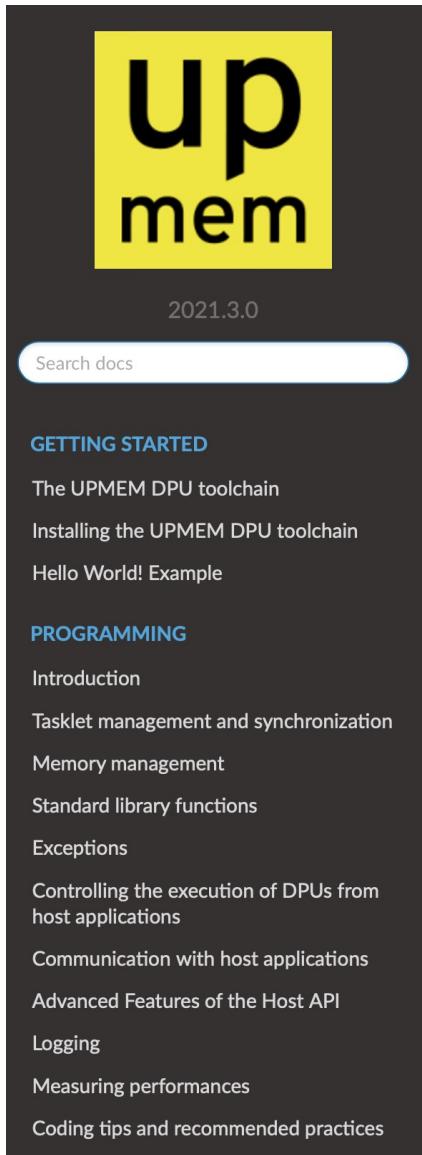
Onur Mutlu Lectures
24.1K subscribers

SUBSCRIBED



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

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](#)

Lecture on Programming UPMEM PIM

“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

DRAM chip have 8-bit data bus

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

Library →

W	W	W	W	W	W	W	W
o	o	o	o	o	o	o	o
r	r	r	r	r	r	r	r
d	d	d	d	d	d	d	d
0	1	2	3	4	5	6	7

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

Copyright UPMEM® 2019

HOT CHIPS 31

up mem

25:46 / 46:42 • "Transposing" Library > True Processing In Memory accelerator, HotChips 2019. doi: 10.1109/HOTCHIPS.2019.887568 CC BY-SA Share Clip Save

PIM Course: Lecture 9: Programming PIM Architectures - Fall 2022



Onur Mutlu Lectures

32.6K subscribers



Subscribed

10



Share

Clip

Save

...

424 views 4 months ago Livestream - P&S Data-Centric Architectures: Fundamentally Improving Performance and Energy (Fall 2022)

Projects & Seminars, ETH Zürich, Fall 2022

Data-Centric Architectures: Fundamentally Improving Performance and Energy

Another Lecture on PIM Programming

The image shows a YouTube video player interface. At the top, the title reads "Computer Architecture" in red, followed by "Lecture 10: Programming a Real-world PIM Architecture" in green. Below the title, the speakers are listed as "Dr. Juan Gómez Luna" and "Prof. Onur Mutlu" from "ETH Zürich" in "Fall 2022". The date "28 October 2022" is also mentioned. The video progress bar indicates it's at 0:34 / 2:45:58. In the bottom right corner of the video frame, there is a small thumbnail of a person speaking, with the text "Juan Gomez L..." below it. The YouTube interface includes standard controls like play/pause, volume, and a zoom logo.

Livestream - Computer Architecture - ETH Zürich (Fall 2022)

Computer Architecture - Lecture 10: Real Processing in Memory Systems: UPMEM Case Study (Fall 2022)



Onur Mutlu Lectures
29.4K subscribers

Subscribed



18



Share

Clip

Save



830 views Streamed live on Oct 28, 2022

Computer Architecture, ETH Zürich, Fall 2022 (<https://safari.ethz.ch/architecture/f...>)

Lecture 10: Real Processing in Memory Systems: UPMEM Case Study

Computer Architecture

Lecture 3: Processing near Memory

Prof. Onur Mutlu
ETH Zürich
Fall 2023
5 October 2023