

P&S Processing-in-Memory

Benchmarking and Workload Suitability
on a Real-World PIM Architecture

Dr. Juan Gómez Luna

Prof. Onur Mutlu

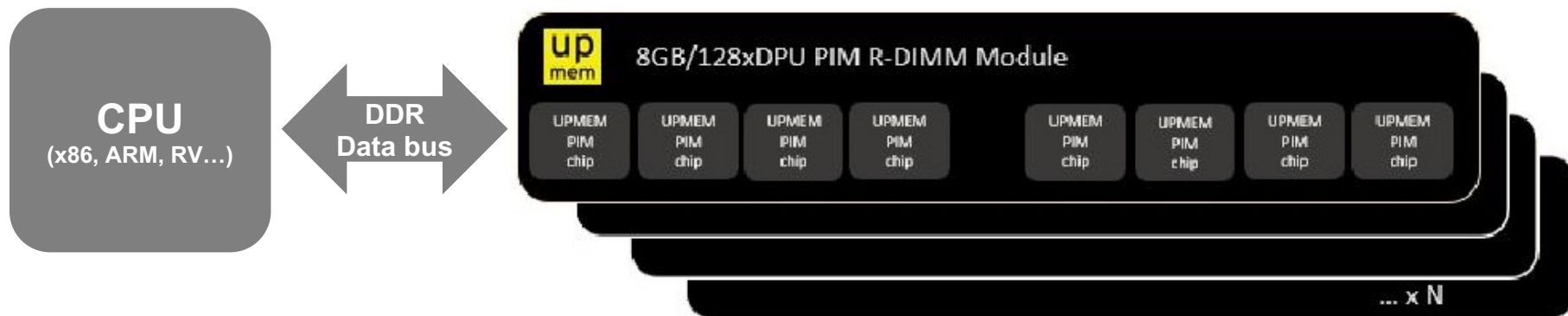
ETH Zürich

Fall 2022

13 December 2022

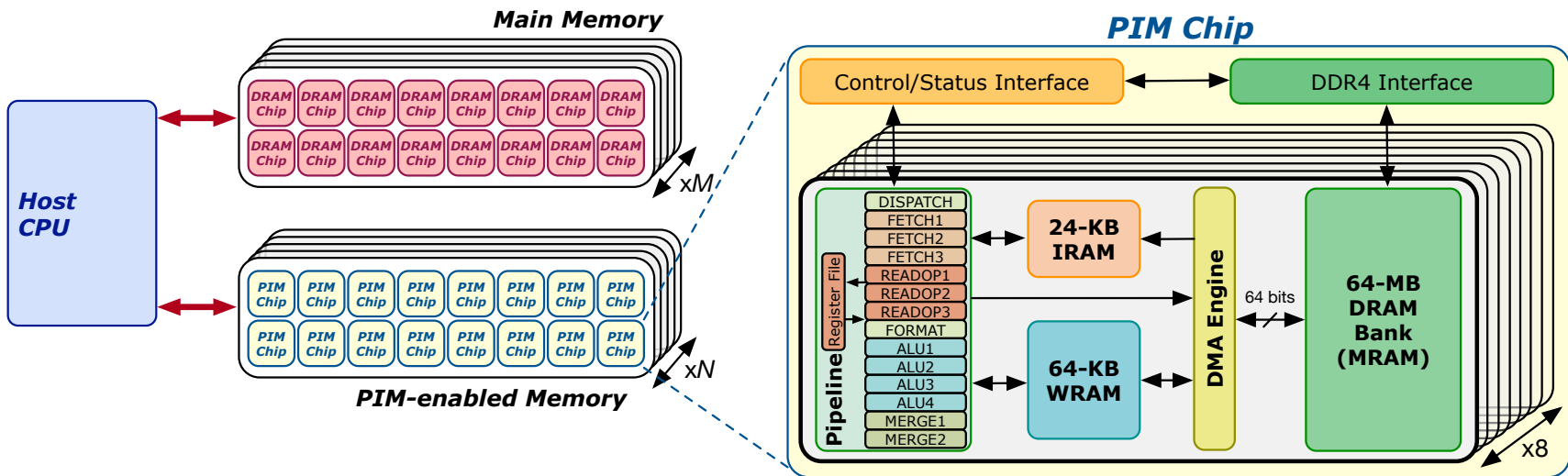
UPMEM Processing-in-DRAM Engine (2019)

- **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



System Organization

- 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



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

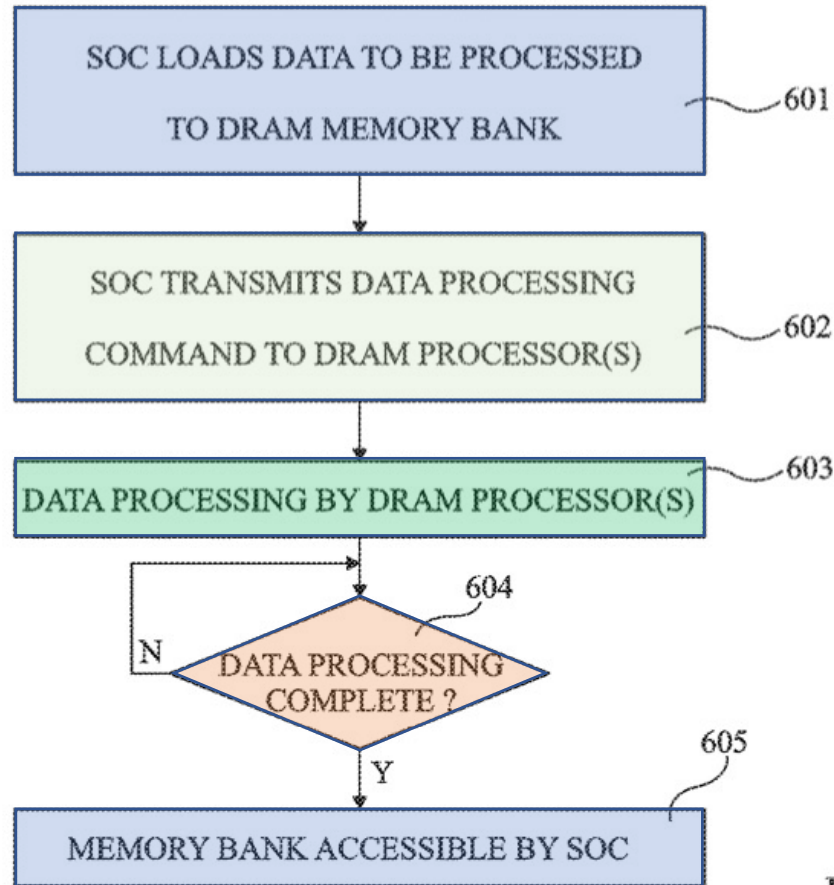
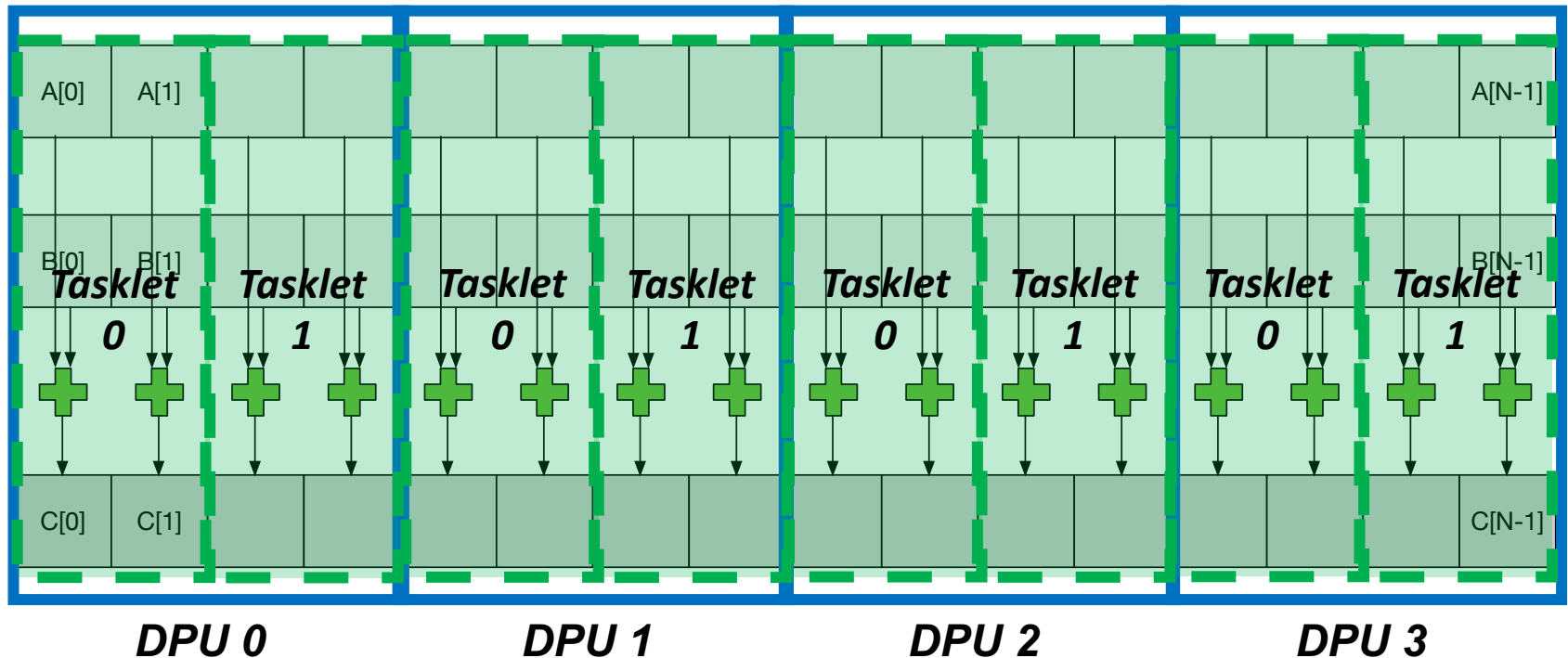


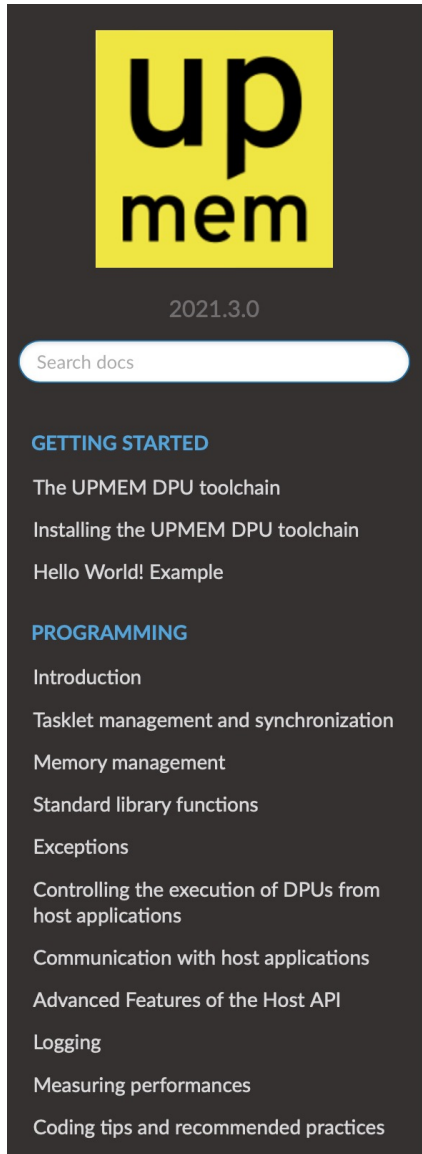
Fig 6

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



The screenshot shows the left sidebar of the UPMEM SDK documentation website. At the top is the UPMEM logo, which consists of the letters 'up' stacked above 'mem' in a bold, lowercase font, set against a yellow square background. Below the logo is the version number '2021.3.0'. A search bar with the placeholder text 'Search docs' is positioned below the version number. The sidebar is divided into two main sections: 'GETTING STARTED' and 'PROGRAMMING'. Under 'GETTING STARTED', there are links for 'The UPMEM DPU toolchain', 'Installing the UPMEM DPU toolchain', and 'Hello World! Example'. Under 'PROGRAMMING', there are links for 'Introduction', 'Tasklet management and synchronization', 'Memory management', 'Standard library functions', 'Exceptions', 'Controlling the execution of DPUs from host applications', 'Communication with host applications', 'Advanced Features of the Host API', 'Logging', 'Measuring performances', and 'Coding tips and recommended practices'.

[Home](#) » [User Manual](#)

User Manual

Getting started

- [The UPMEM DPU toolchain](#)
 - [Notes before starting](#)
 - [The toolchain purpose](#)
 - [dpu-upmem-dpurte-clang](#)
 - [Limitations](#)
 - [The DPU Runtime Library](#)
 - [The Host Library](#)
 - [dpu-ldb](#)
- [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.

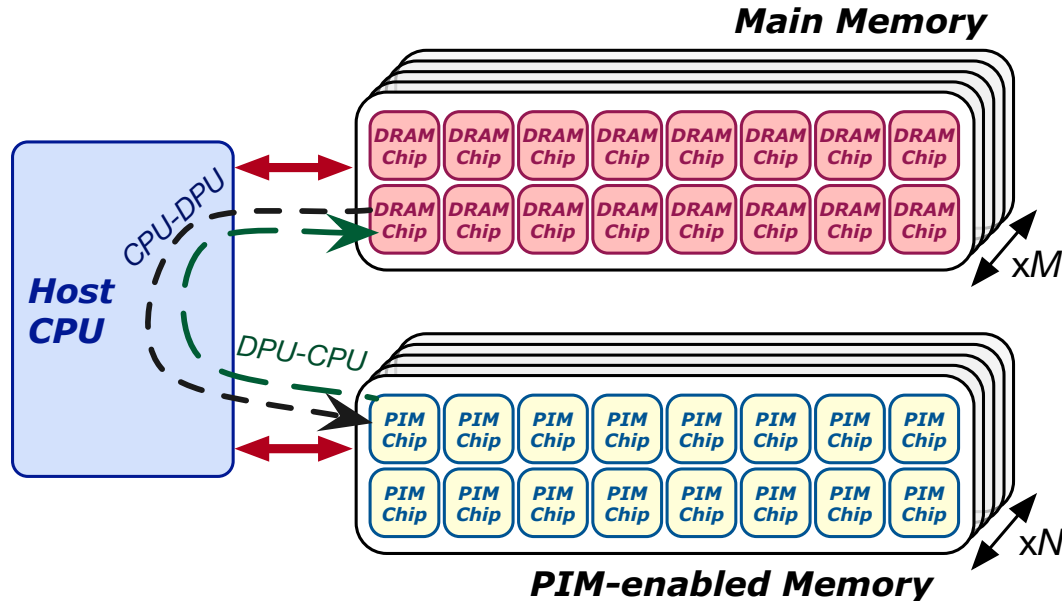
* <https://sdk.upmem.com/2021.1.1/index.html>

★ 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

CPU-DPU/DPU-CPU Data Transfers

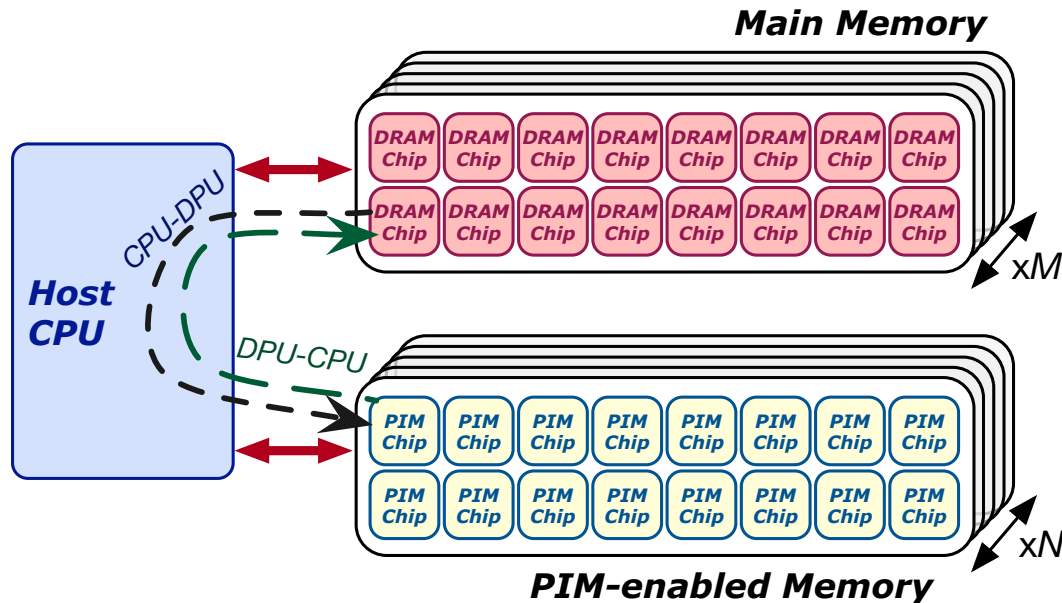
- CPU-DPU and DPU-CPU transfers
 - Between host CPU's main memory and DPUs' MRAM banks



- **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

Inter-DPU Communication

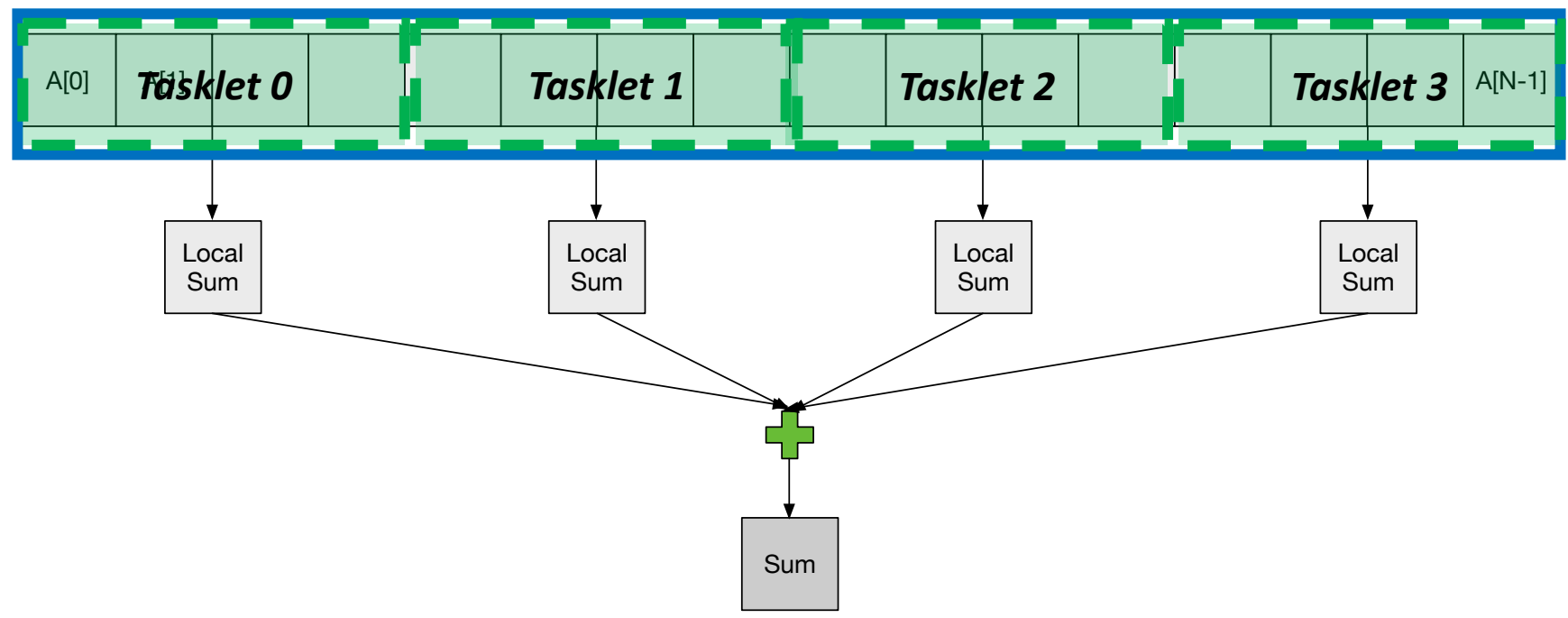
- There is **no direct communication channel between DPUs**



- 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

Parallel Reduction (II)

- Each tasklet computes a local sum



Benchmark Selection and Diversity

PrIM Benchmarks

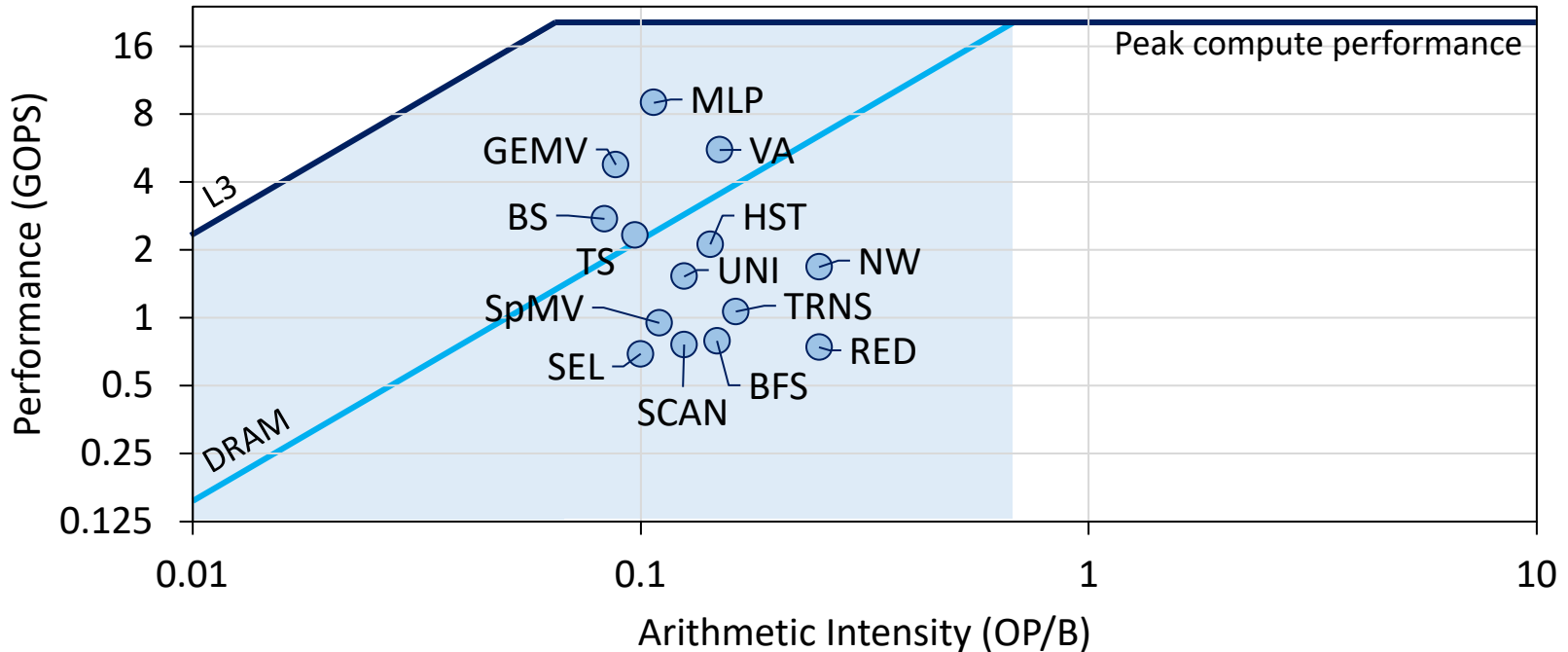
- Goal
 - A **common set of workloads** that can be used to
 - evaluate the UPMEM PIM architecture,
 - compare software improvements and compilers,
 - compare future PIM architectures and hardware
- Two key selection criteria:
 - Selected workloads from **different application domains**
 - **Memory-bound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

PrIM Benchmarks: Application Domains

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

Roofline Model

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



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

PrIM Benchmarks: Diversity

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

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

PrIM Benchmarks: Inter-DPU Communication

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

PrIM Benchmarks: Inter-DPU Communication

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

- Inter-DPU communication

- Result merging:

- SEL, UNI, HST-S, HST-L, RED
- Only DPU-CPU transfers

- Redistribution of intermediate results:

- BFS, MLP, NW, SCAN-SSA, SCAN-RSS
- DPU-CPU and CPU-DPU transfers

Benchmark Evaluation

Evaluation Methodology

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

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

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

Evaluation Methodology

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

Datasets

- Strong and weak scaling experiments

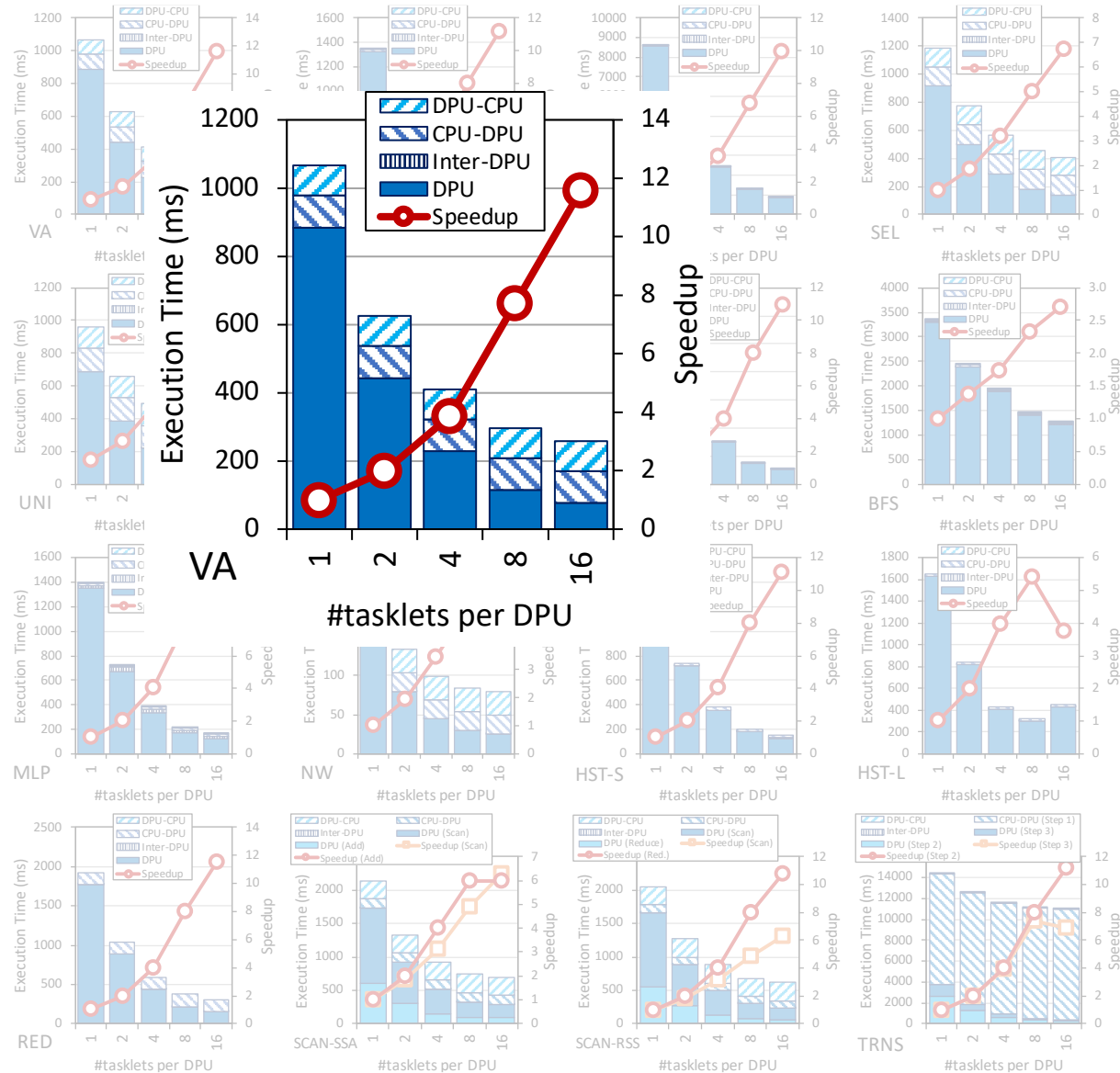
Benchmark	Strong Scaling Dataset	Weak Scaling Dataset	MRAM-WRAM Transfer Sizes
VA	1 DPU-1 rank: 2.5M elem. (10 MB) 32 ranks: 160M elem. (640 MB)	2.5M elem./DPU (10 MB)	1024 bytes
GEMV	1 DPU-1 rank: 8192 × 1024 elem. (32 MB) 32 ranks: 163840 × 4096 elem. (2.56 GB)	1024 × 2048 elem./DPU (8 MB)	1024 bytes
SpMV	<i>bcstk30</i> [253] (12 MB)	<i>bcstk30</i> [253]	64 bytes
SEL	1 DPU-1 rank: 3.8M elem. (30 MB) 32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
UNI	1 DPU-1 rank: 3.8M elem. (30 MB) 32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
BS	2M elem. (16 MB). 1 DPU-1 rank: 256K queries. (2 MB) 32 ranks: 16M queries. (128 MB)	2M elem. (16 MB). 256K queries./DPU (2 MB).	8 bytes
TS	256 elem. query. 1 DPU-1 rank: 512K elem. (2 MB) 32 ranks: 32M elem. (128 MB)	512K elem./DPU (2 MB)	256 bytes
BFS	<i>loc-gowalla</i> [254] (22 MB)	<i>rMat</i> [255] (≈100K vertices and 1.2M edges per DPU)	8 bytes
MLP	3 fully-connected layers. 1 DPU-1 rank: 2K neurons (32 MB) 32 ranks: ≈160K neur. (2.56 GB)	3 fully-connected layers. 1K neur./DPU (4 MB)	1024 bytes
NW	1 DPU-1 rank: 2560 bps (50 MB), large/small sub-block= $\frac{2560}{\#DPU_s}/2$ 32 ranks: 64K bps (32 GB), l./s.=32/2	512 bps/DPU (2MB), l./s.=512/2	8, 16, 32, 40 bytes
HST-S	1 DPU-1 rank: 1536 × 1024 input image [256] (6 MB) 32 ranks: 64 × input image	1536 × 1024 input image [256]/DPU (6 MB)	1024 bytes
HST-L	1 DPU-1 rank: 1536 × 1024 input image [256] (6 MB) 32 ranks: 64 × input image	1536 × 1024 input image [256]/DPU (6 MB)	1024 bytes
RED	1 DPU-1 rank: 6.3M elem. (50 MB) 32 ranks: 400M elem. (3.1 GB)	6.3M elem./DPU (50 MB)	1024 bytes
SCAN-SSA	1 DPU-1 rank: 3.8M elem. (30 MB) 32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
SCAN-RSS	1 DPU-1 rank: 3.8M elem. (30 MB) 32 ranks: 240M elem. (1.9 GB)	3.8M elem./DPU (30 MB)	1024 bytes
TRNS	1 DPU-1 rank: 12288 × 16 × 64 × 8 (768 MB) 32 ranks: 12288 × 16 × 2048 × 8 (24 GB)	12288 × 16 × 1 × 8/DPU (12 MB)	128, 1024 bytes

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

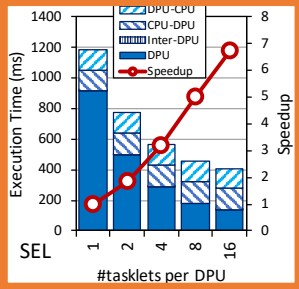
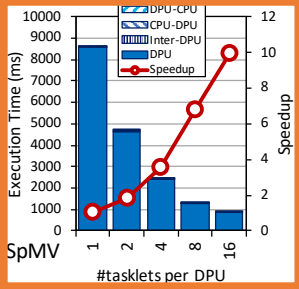
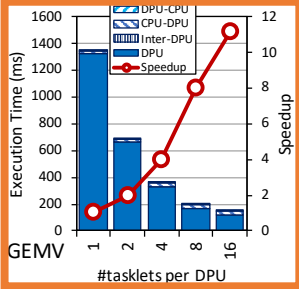
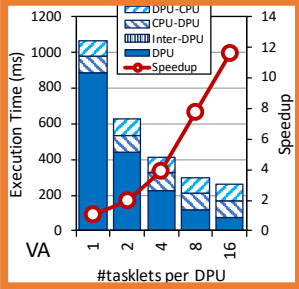
Strong Scaling: 1 DPU (I)

- Strong scaling experiments on 1 DPU

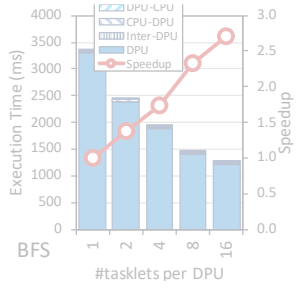
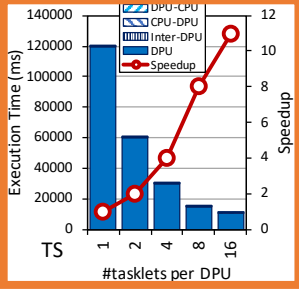
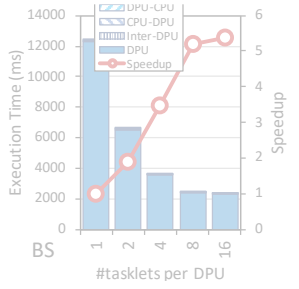
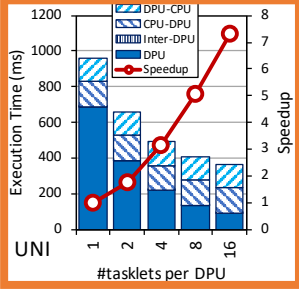
- We set the number of tasklets to 1, 2, 4, 8, and 16
- We show the breakdown of execution time:
 - **DPU**: Execution time on the DPU
 - **Inter-DPU**: Time for inter-DPU communication via the host CPU
 - **CPU-DPU**: Time for CPU to DPU transfer of input data
 - **DPU-CPU**: Time for DPU to CPU transfer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



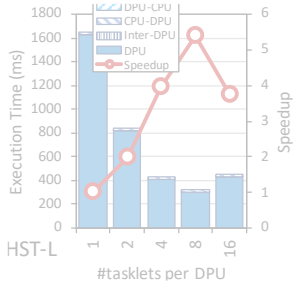
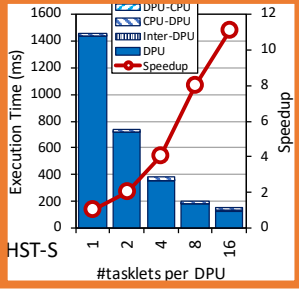
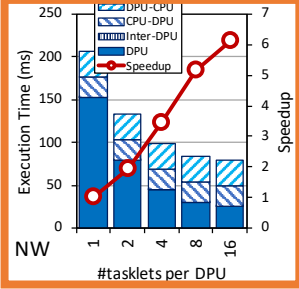
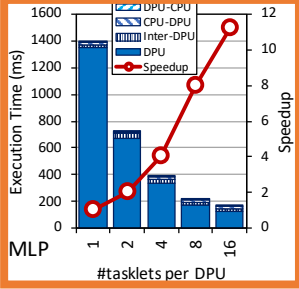
Strong Scaling: 1 DPU (II)



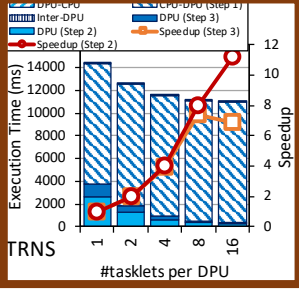
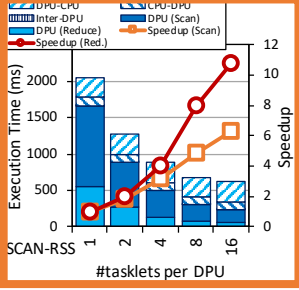
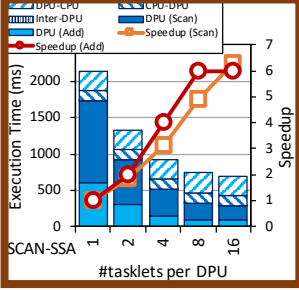
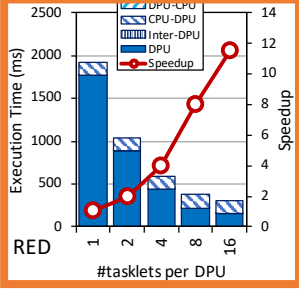
VA, GEMV, SpMV, SEL, UNI, TS, MLP, NW, HST-S, RED, SCAN-SSA (Scan kernel), SCAN-RSS (both kernels), and TRNS (Step 2 kernel), the best performing number of tasklets is 16



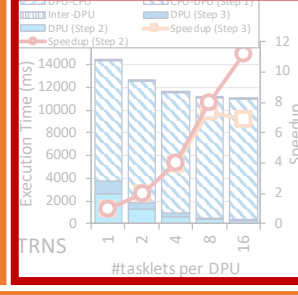
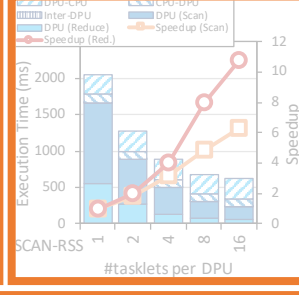
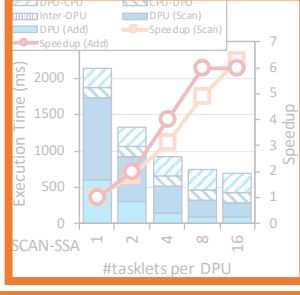
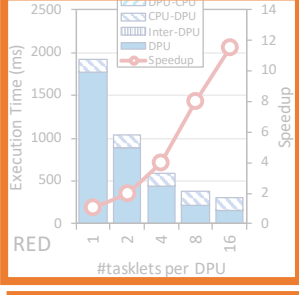
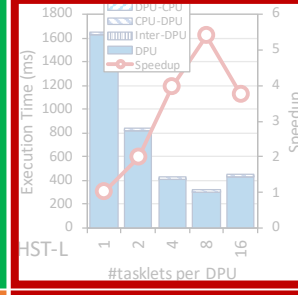
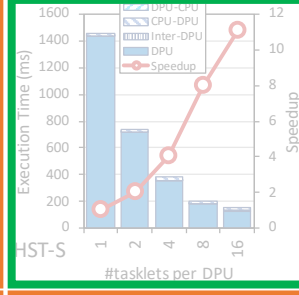
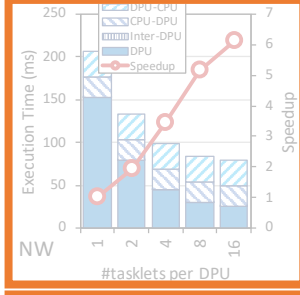
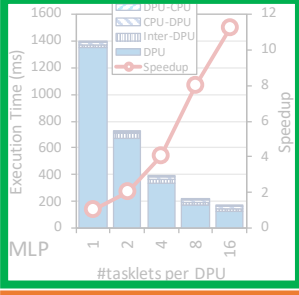
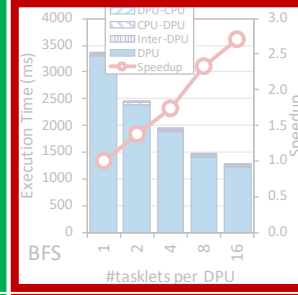
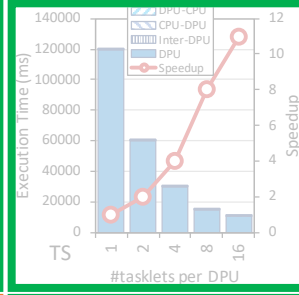
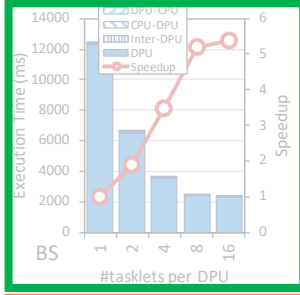
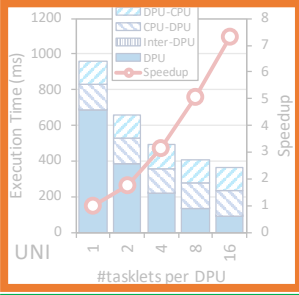
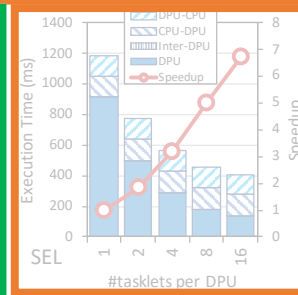
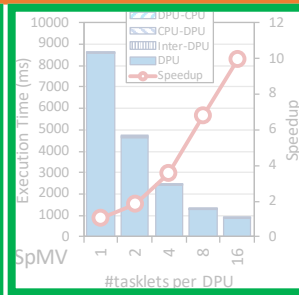
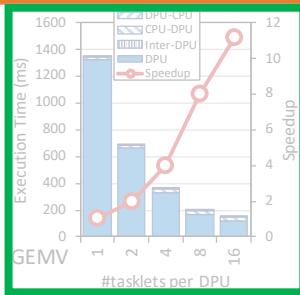
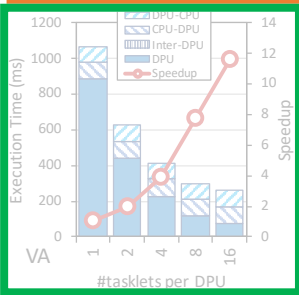
Speedups 1.5-2.0x as we double the number of tasklets from 1 to 8. Speedups 1.2-1.5x from 8 to 16, since the pipeline throughput saturates at 11 tasklets



KEY OBSERVATION 10
A number of tasklets greater than 11 is a good choice for most real-world workloads we tested (16 kernels out of 19 kernels from 16 benchmarks), as it fully utilizes the DPU's pipeline.



Strong Scaling: 1 DPU (III)

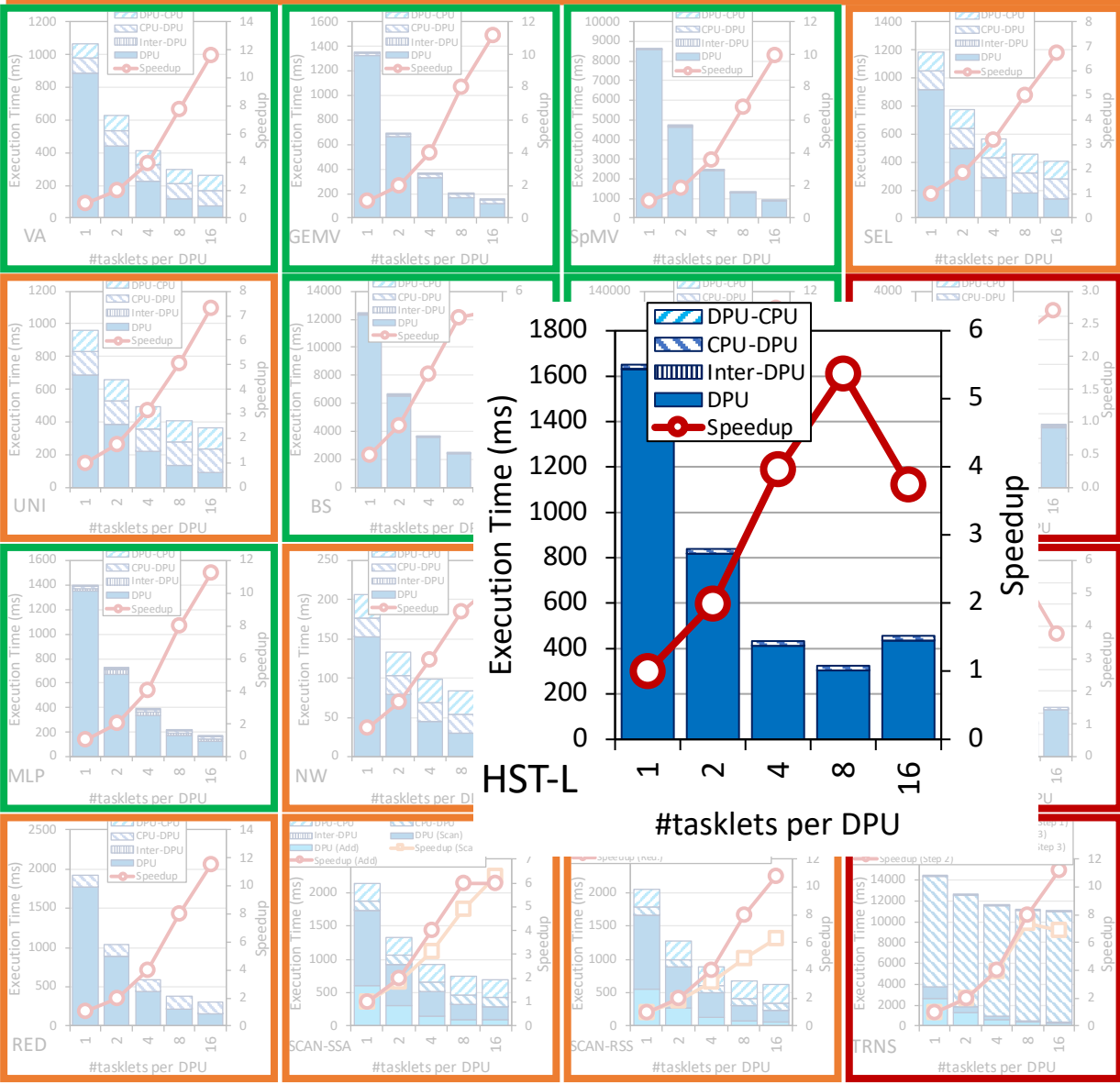


VA, GEMV, SpMV, BS, TS, MLP, HST-S do not use intra-DPU synchronization primitives

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

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

Strong Scaling: 1 DPU (IV)



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

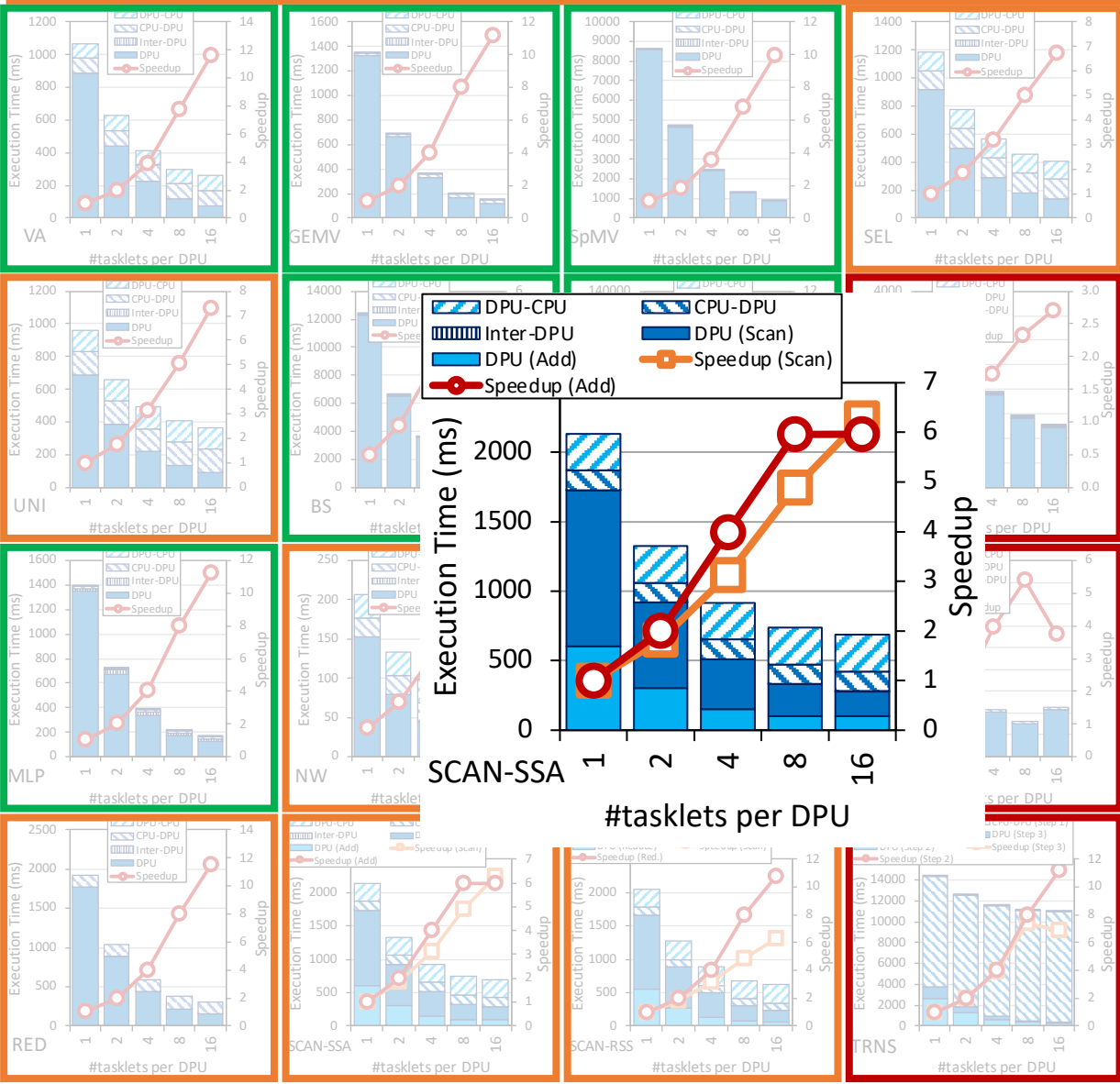
In SEL, UNI, NW, RED, SCAN-SSA (Scan kernel), SCAN-RSS (both kernels), synchronization is lightweight

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

KEY OBSERVATION 11

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

Strong Scaling: 1 DPU (V)

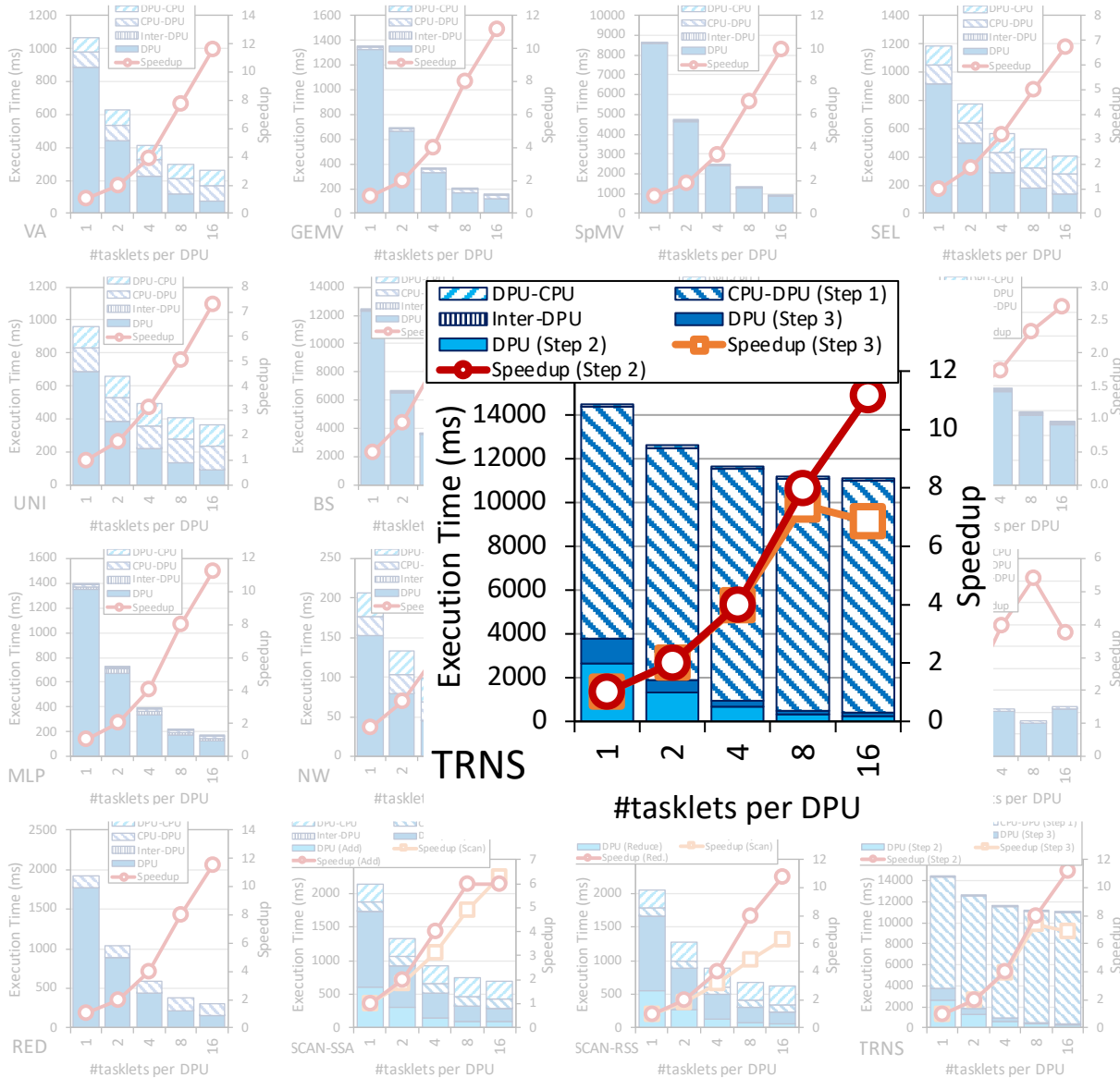


SCAN-SSA (Add kernel) is **not compute-intensive**. Thus, performance saturates with less than 11 tasklets (recall STREAM ADD). BS shows similar behavior

KEY OBSERVATION 12

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

Strong Scaling: 1 DPU (VI)



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

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

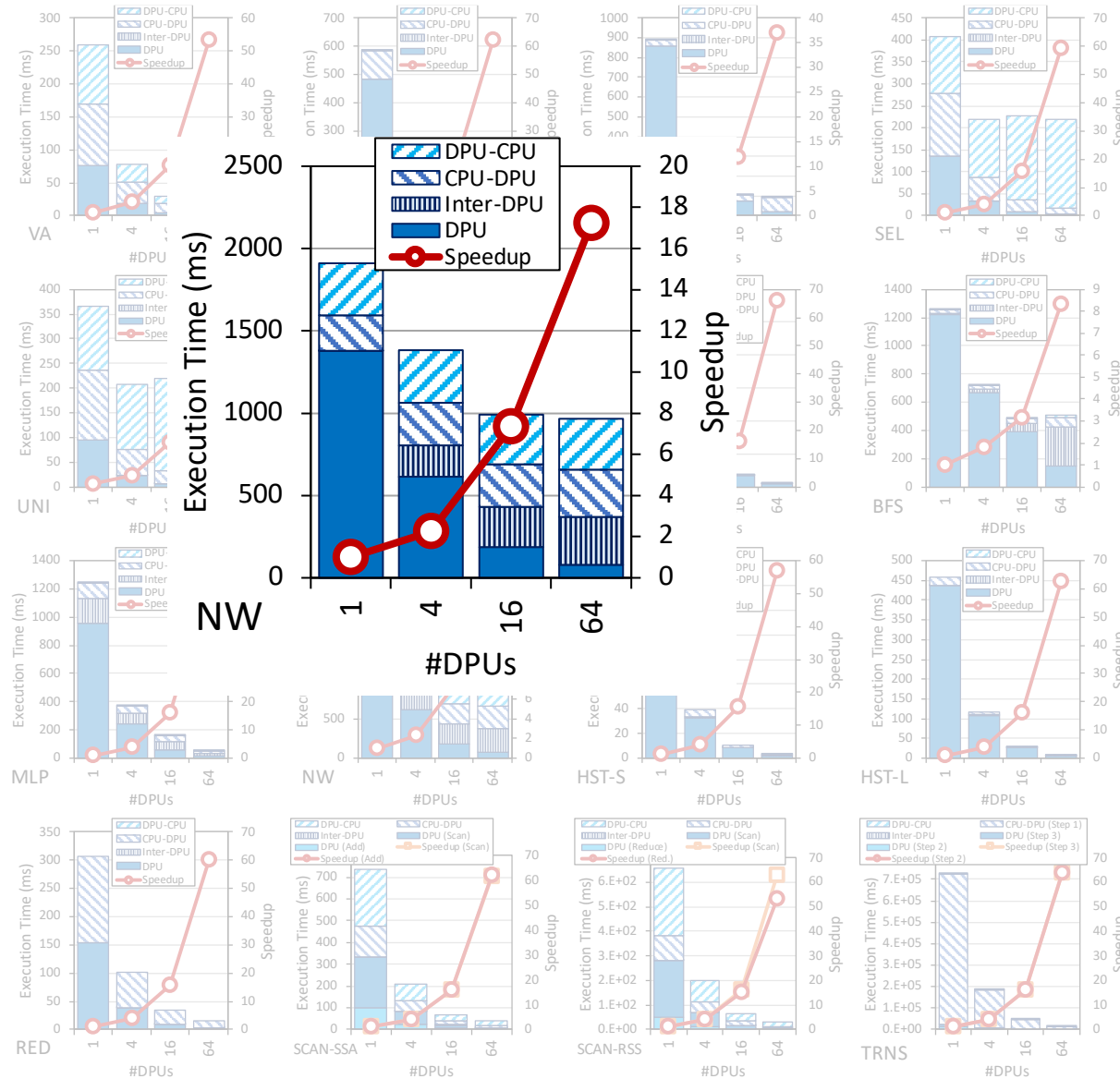
KEY OBSERVATION 13

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

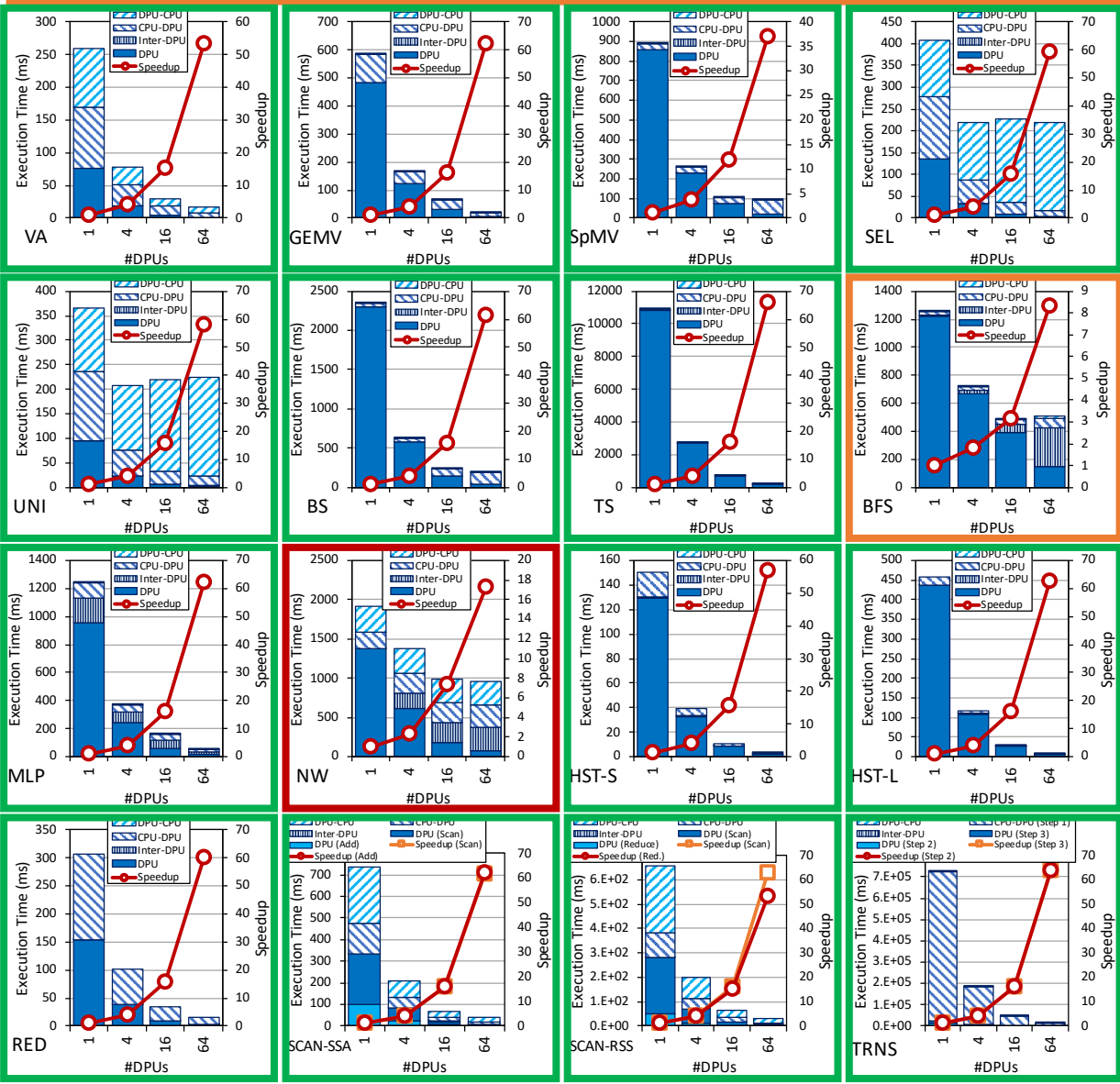
Strong Scaling: 1 Rank (I)

- Strong scaling experiments on 1 rank

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



Strong Scaling: 1 Rank (II)



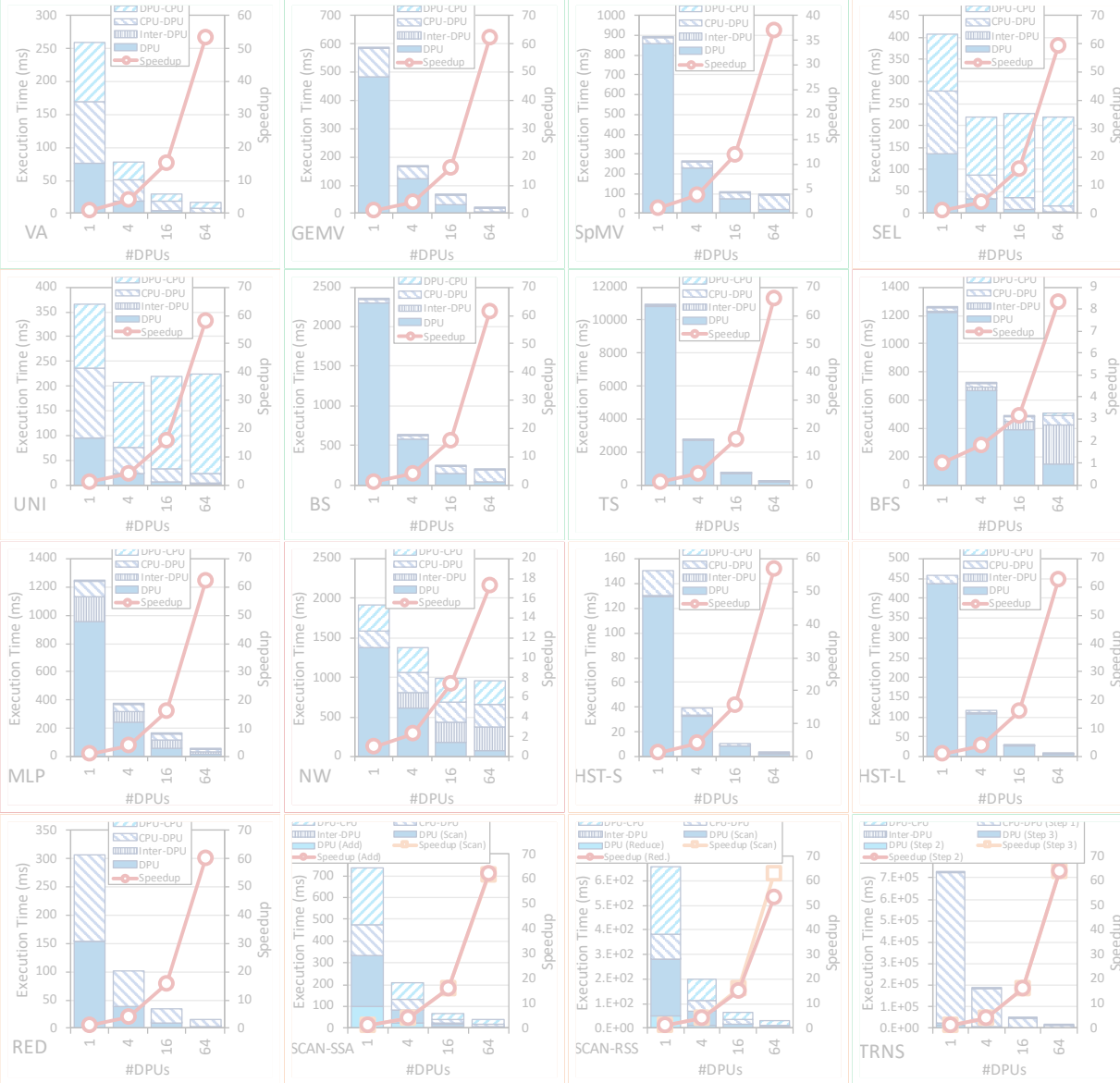
VA, GEMV, SpMV, SEL, UNI, BS, TS, MLP, HST-S, HSTS-L, RED, SCAN-SSA (both kernel), SCAN-RSS (both kernels), and TRNS (both kernels) scale linearly with the number of DPUs

Scaling is sublinear for BFS and NW

BFS suffers load imbalance due to irregular graph topology

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

Strong Scaling: 1 Rank (III)

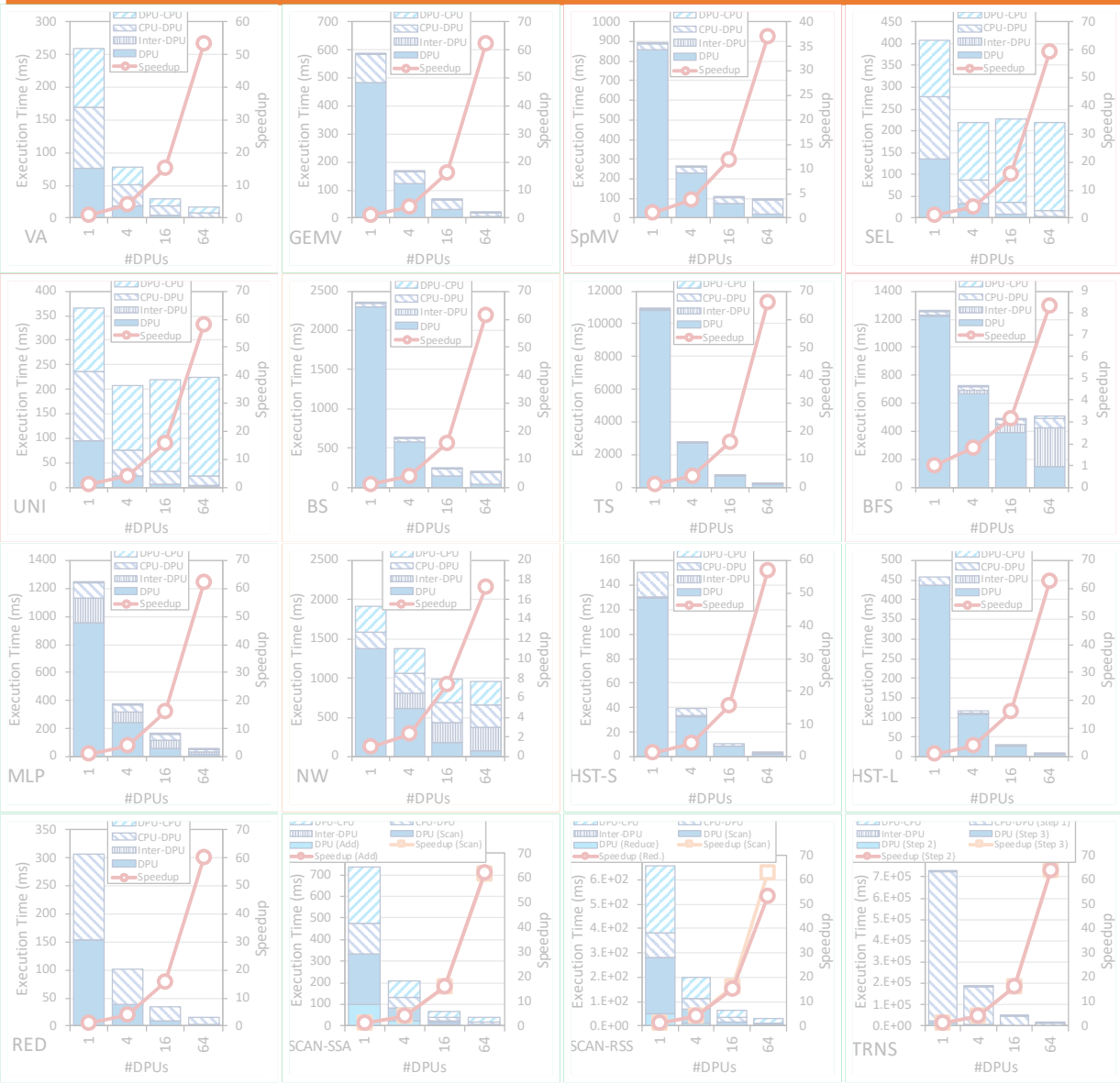


VA, GEMV, SpMV, BS, TS, TRNS **do not need inter-DPU synchronization**

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

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

Strong Scaling: 1 Rank (IV)



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

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

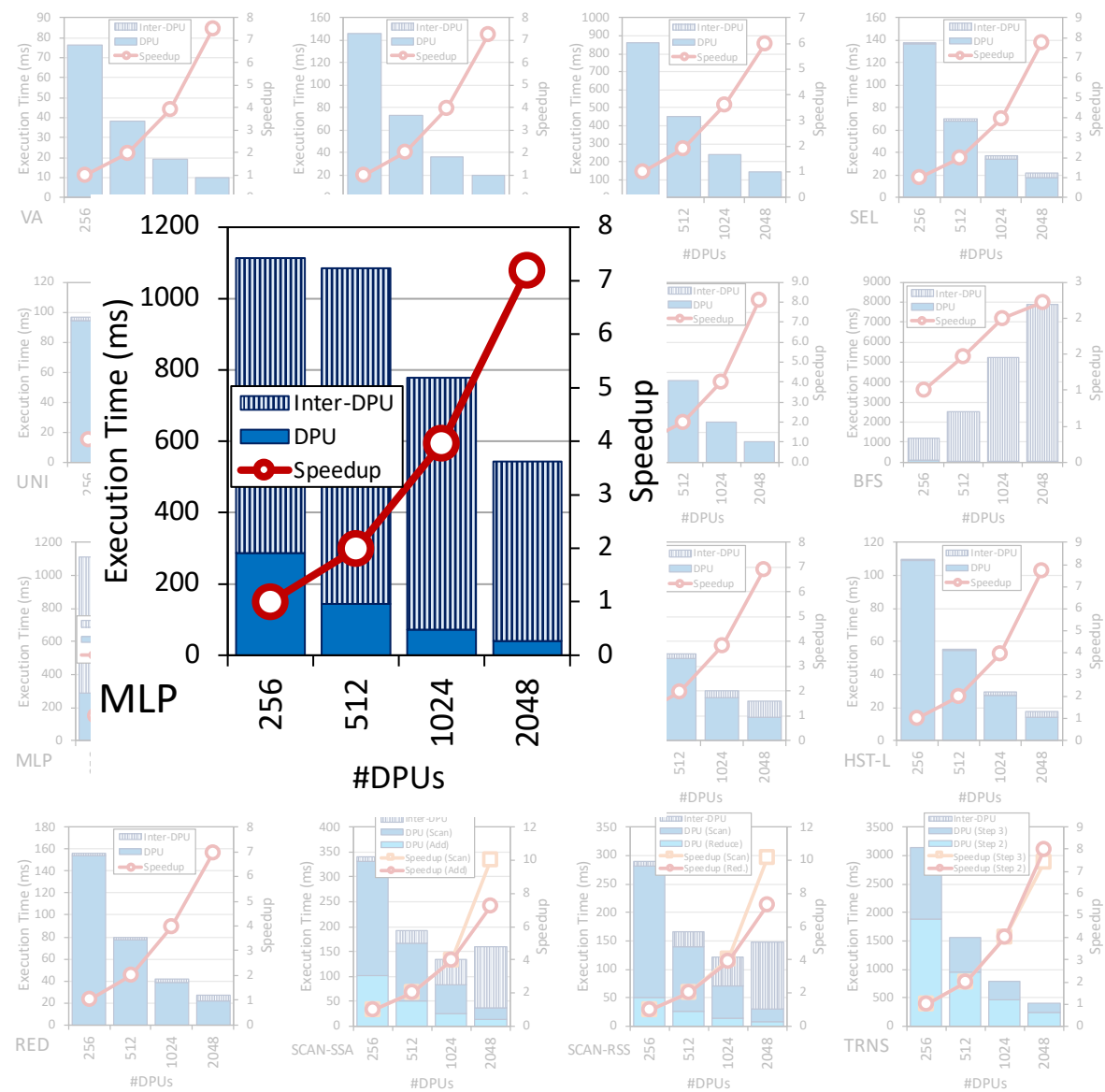
- BS transfers a complete array to all DPUs.
- NW does not use all DPUs in all iterations

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

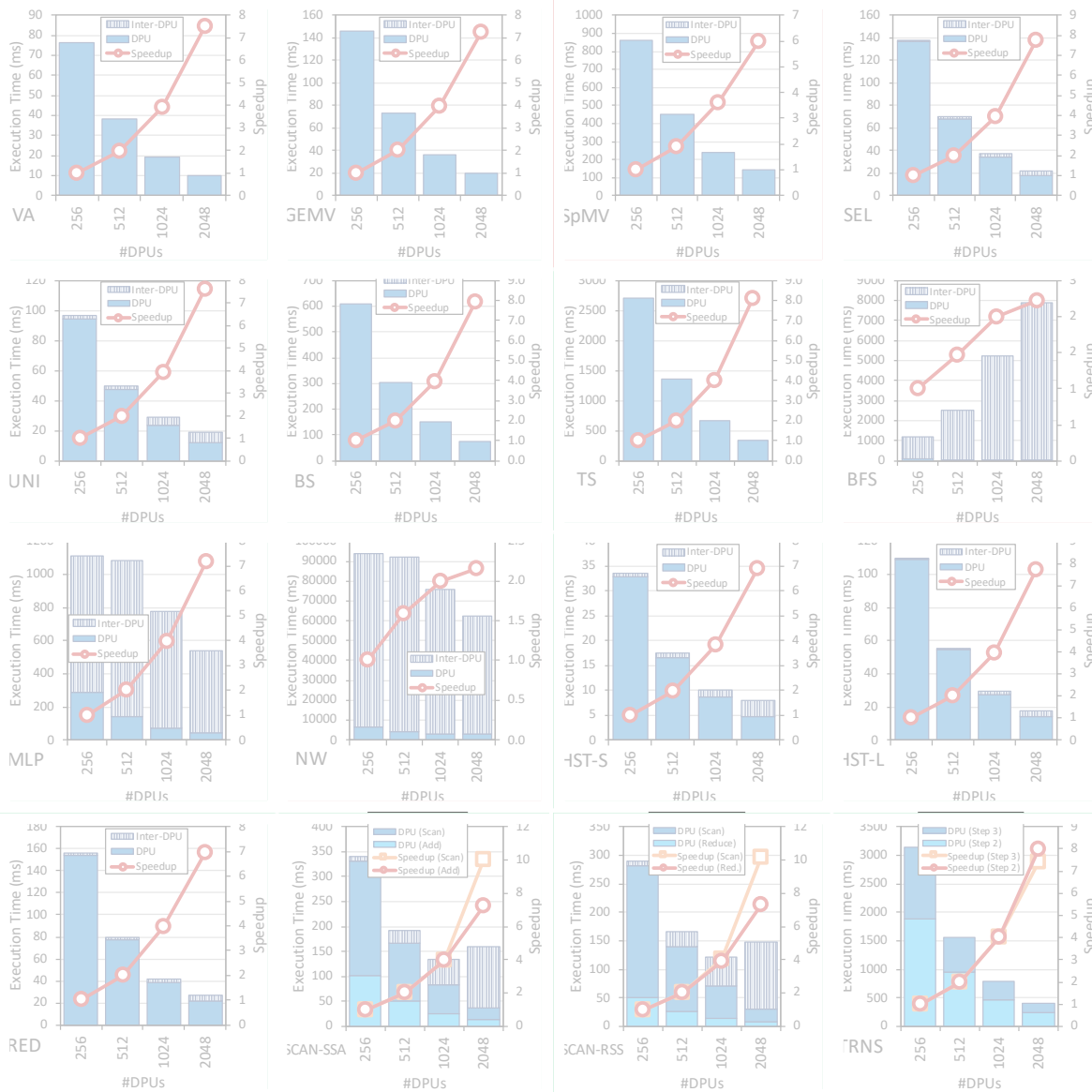
PROGRAMMING RECOMMENDATION 5
Parallel CPU-DPU/DPU-CPU transfers inside a rank of DPUs are recommended for real-world workloads when all transferred buffers are of the same size.

Strong Scaling: 32 Ranks (I)

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



Strong Scaling: 32 Ranks (II)

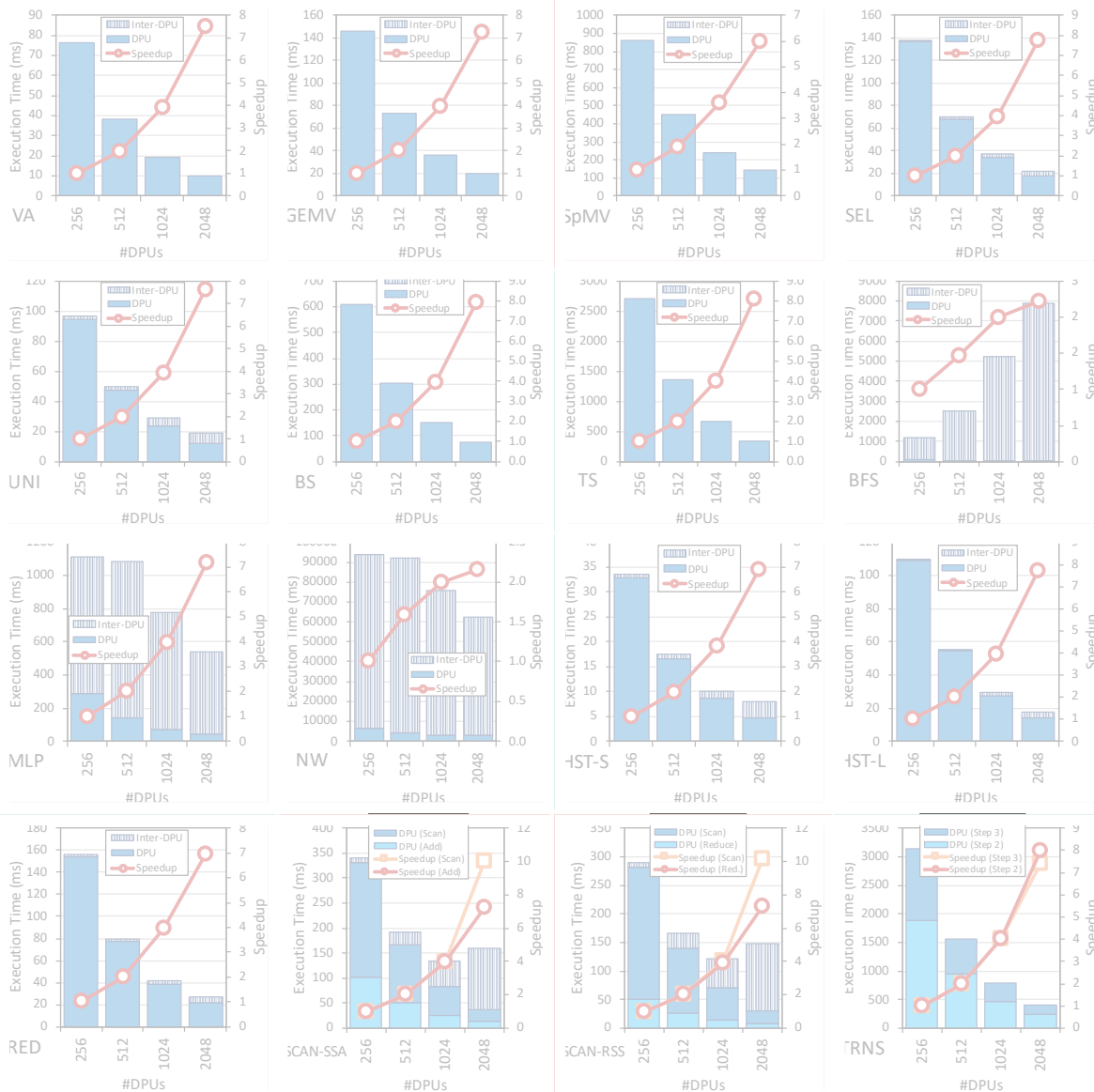


VA, GEMV, SEL, UNI, BS, TS, MLP, HST-S, HSTS-L, RED, SCAN-SSA (both kernel), SCAN-RSS (both kernels), and TRNS (both kernels) **scale linearly with the number of DPUs**

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

KEY OBSERVATION 14
Load balancing across DPUs ensures linear reduction of the execution time spent on the DPUs for a given problem size, when all available DPUs are used (as observed in strong scaling experiments).

Strong Scaling: 32 Ranks (III)



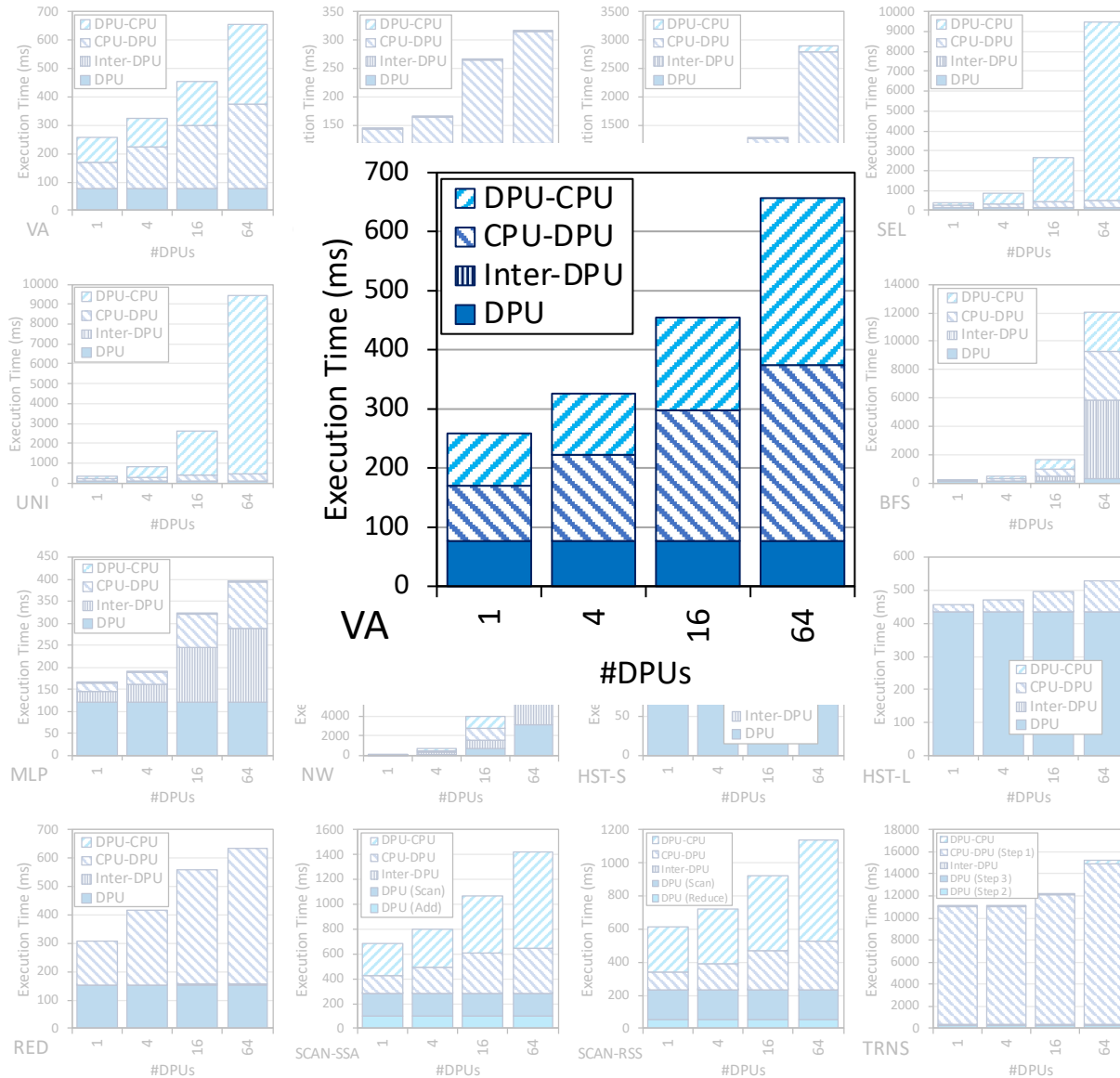
SEL, UNI, HST-S, HST-L, RED only need to merge final results

KEY OBSERVATION 15
 The overhead of merging partial results from DPUs in the host CPU is tolerable across all PRIM benchmarks that need it.

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

KEY OBSERVATION 16
 Complex synchronization across DPUs (i.e., inter-DPU synchronization involving two-way communication with the host CPU) imposes significant overhead, which limits scalability to more DPUs.

Weak Scaling: 1 Rank



KEY OBSERVATION 17

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

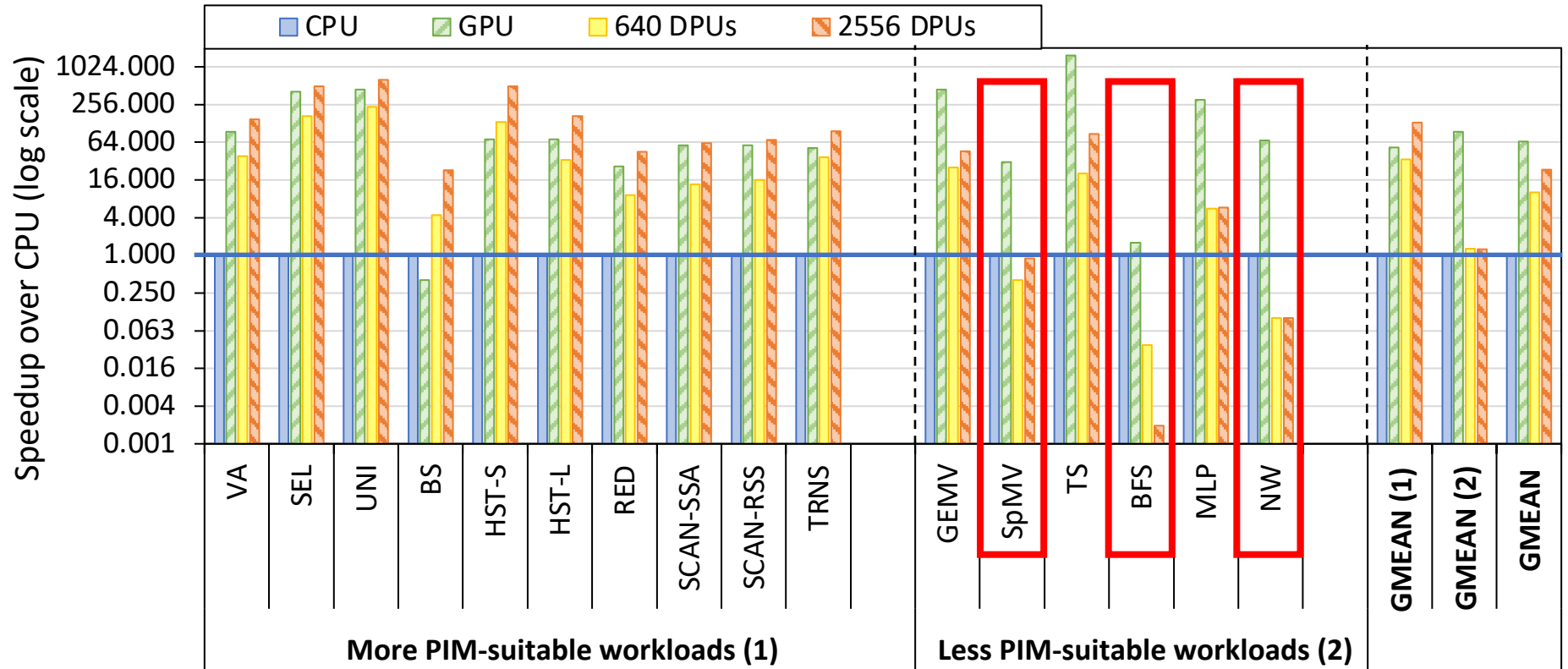
KEY OBSERVATION 18

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

CPU/GPU: Evaluation Methodology

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

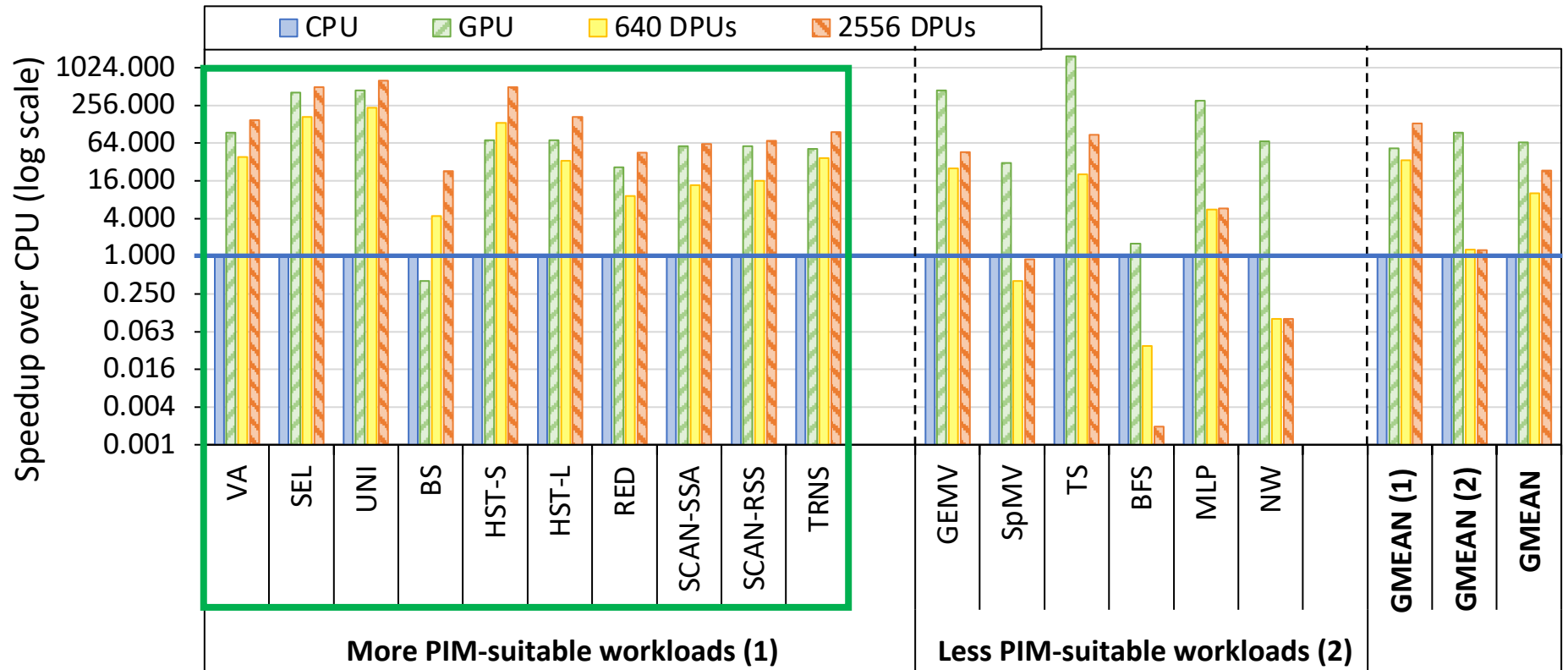
CPU/GPU: Performance Comparison (I)



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

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

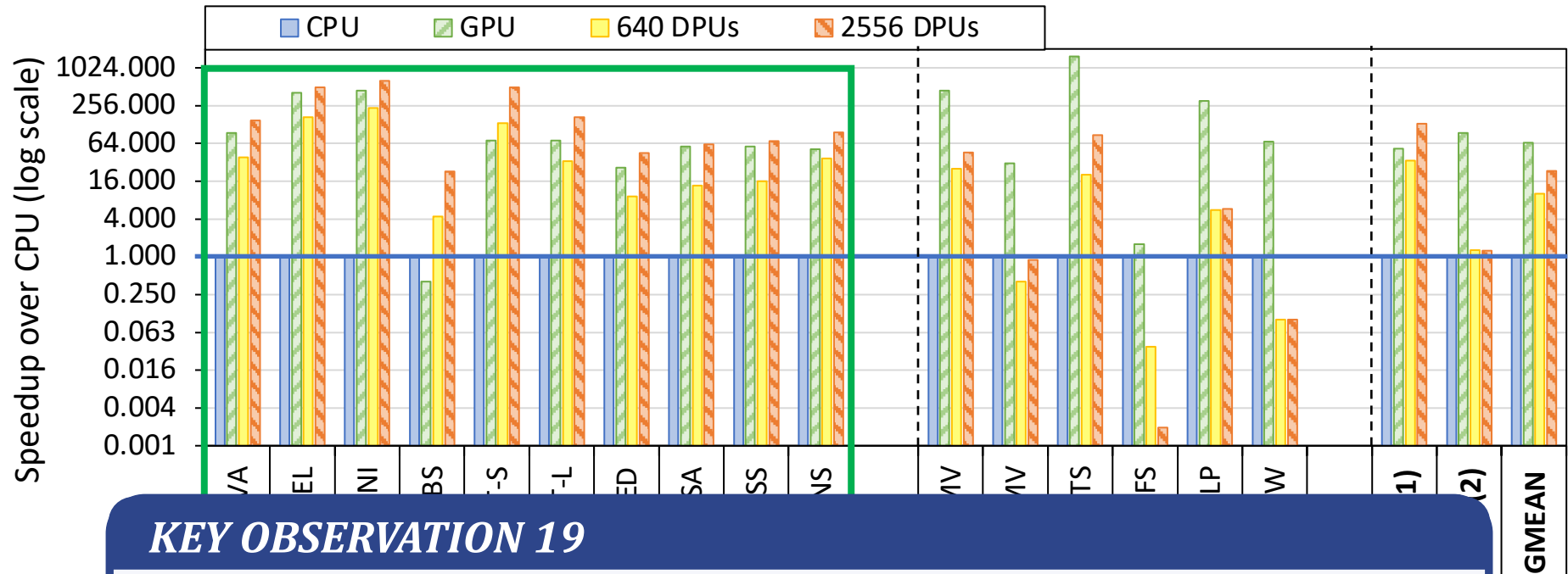
CPU/GPU: Performance Comparison (II)



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

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

CPU/GPU: Performance Comparison (III)



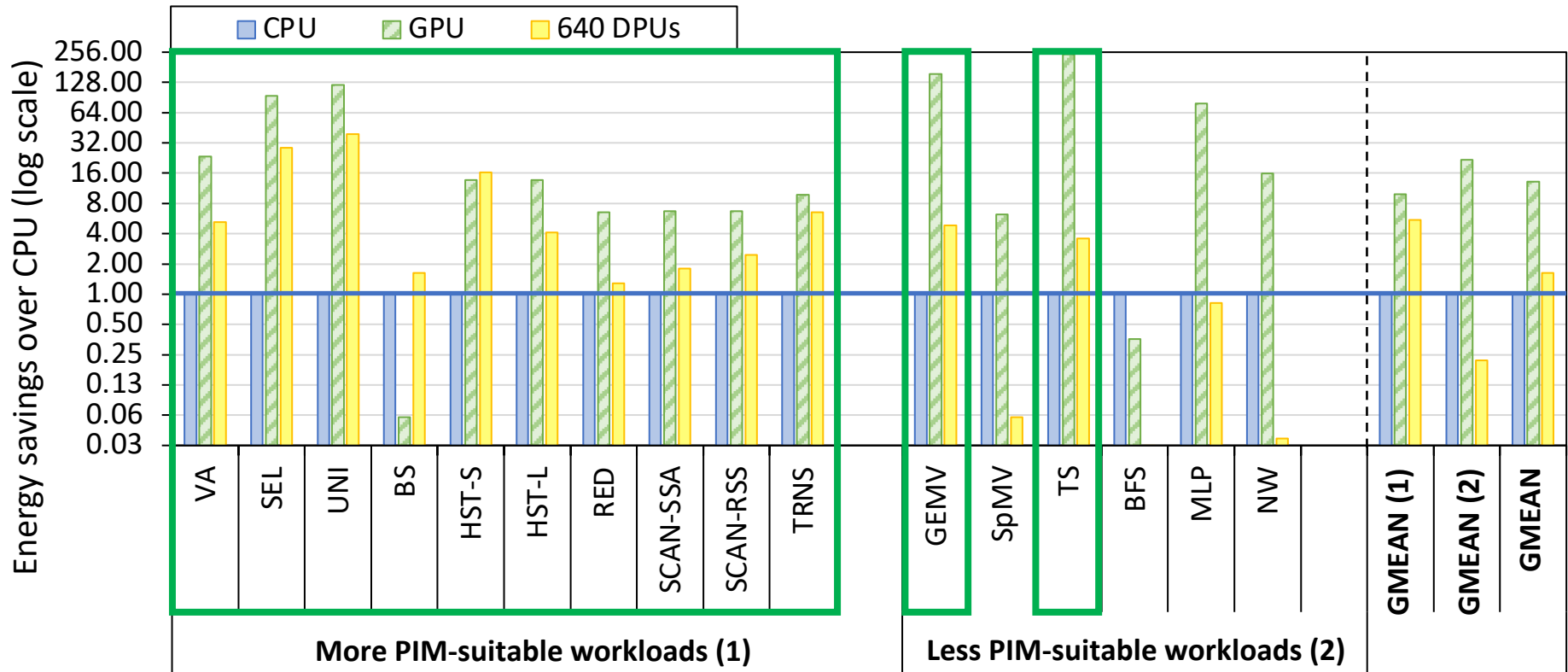
KEY OBSERVATION 19

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

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

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

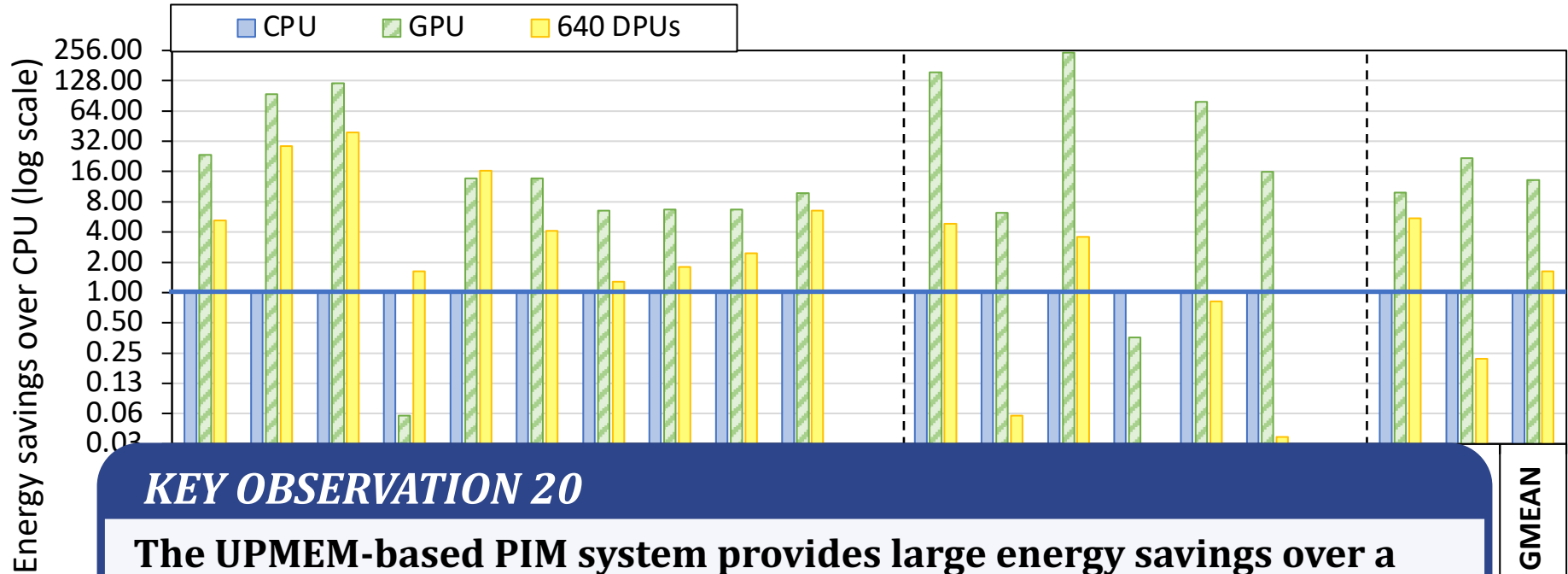
CPU/GPU: Energy Comparison (I)



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

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

CPU/GPU: Energy Comparison (II)



KEY OBSERVATION 20

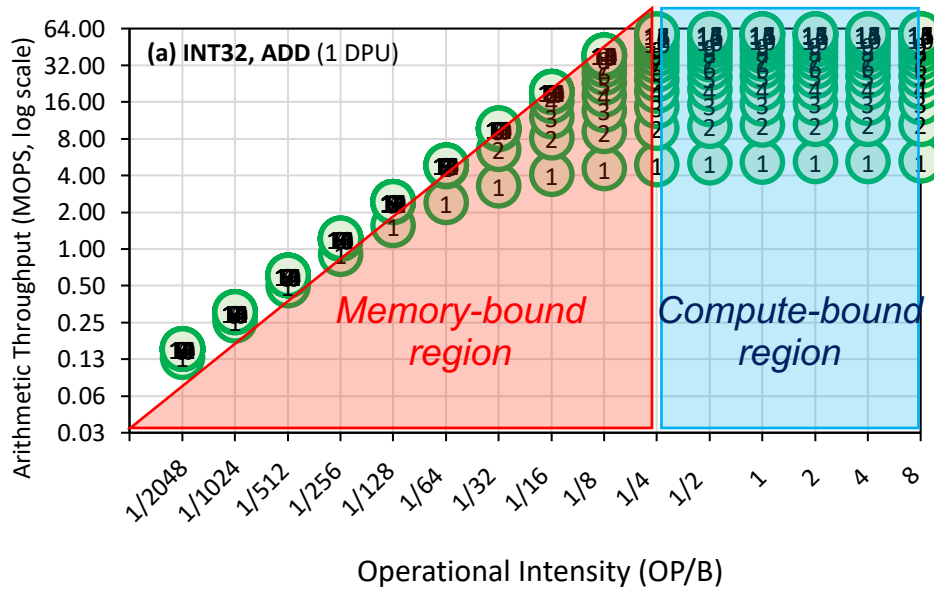
The UPMEM-based PIM system provides large energy savings over a state-of-the-art CPU due to higher performance (thus, lower static energy) and less data movement between memory and processors.

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

This is because the source of both performance improvement and energy savings is the same: **the significant reduction in data movement between the memory and the processor cores**, which the UPMEM-based PIM system can provide for PIM-suitable workloads.

Key Takeaways

Key Takeaway 1

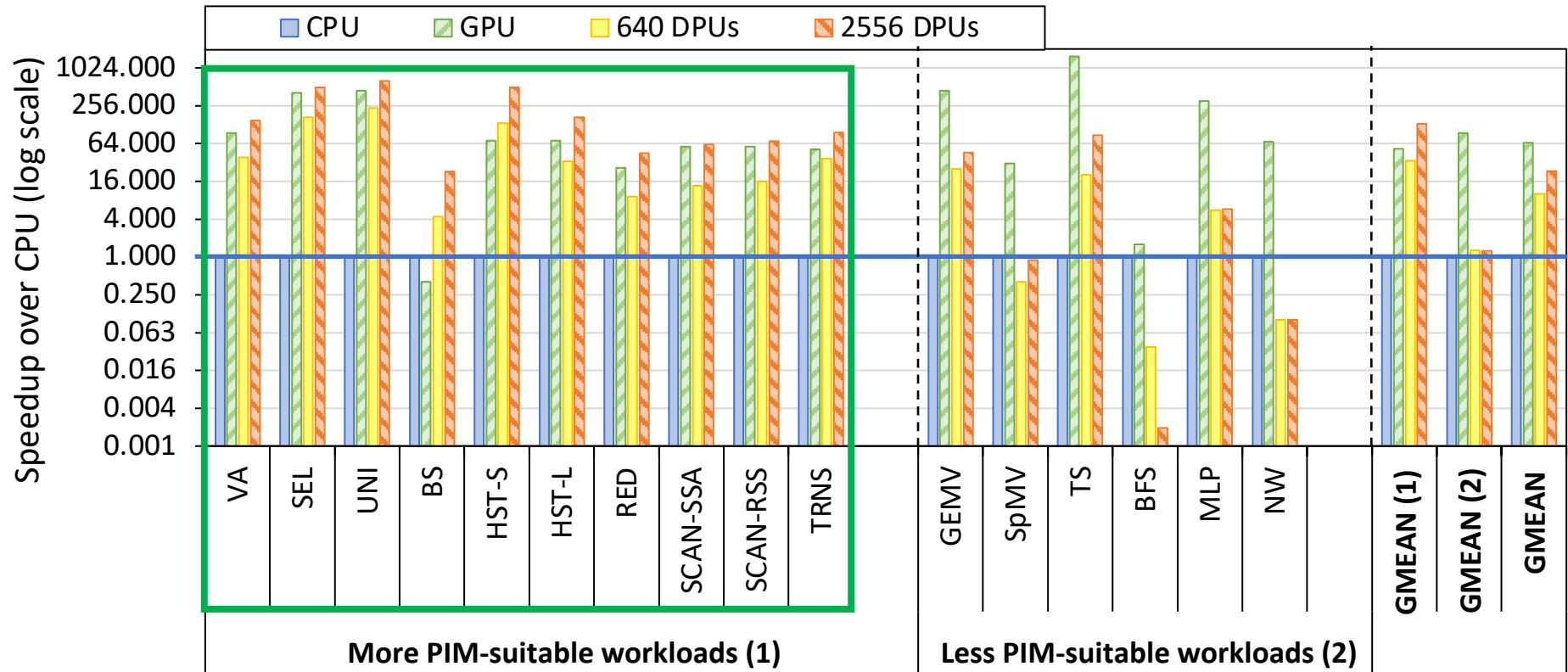


The throughput saturation point is as low as $\frac{1}{4}$ OP/B, i.e., 1 integer addition per every 32-bit element fetched

KEY TAKEAWAY 1

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

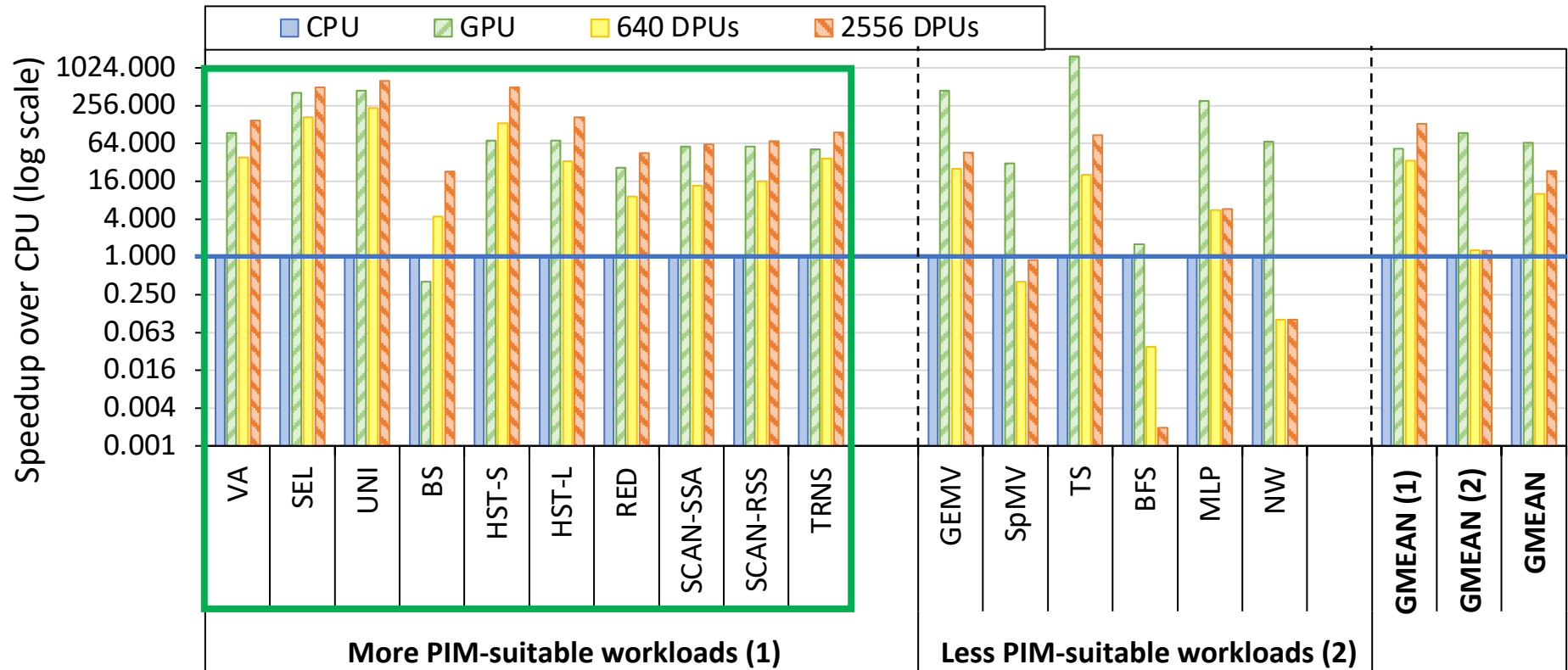
Key Takeaway 2



KEY TAKEAWAY 2

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

Key Takeaway 3



KEY TAKEAWAY 3

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

Key Takeaway 4

KEY TAKEAWAY 4

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

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

Juan Gómez-Luna
ETH Zürich

Izzat El Hajj
*American University
of Beirut*

Ivan Fernandez
*University
of Malaga*

Christina Giannoula
*National Technical
University of Athens*

Geraldo F. Oliveira
ETH Zürich

Onur Mutlu
ETH Zürich

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

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

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

Juan Gómez-Luna¹ Izzat El Hajj² Ivan Fernandez^{1,3} Christina Giannoula^{1,4}
Geraldo F. Oliveira¹ Onur Mutlu¹

¹ETH Zürich ²American University of Beirut ³University of Malaga ⁴National Technical University of Athens

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

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

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 Benchmarks

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

CMU-SAFARI / prim-benchmarks

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

main 1 branch 0 tags

Go to file Add file Code

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

Upcoming Lectures

- More real-world PIM architectures
- PUM architectures and prototypes
 - SIMD RAM: An end-to-end framework for bit-serial SIMD computing in DRAM

P&S Processing-in-Memory

Benchmarking and Workload Suitability
on a Real-World PIM Architecture

Dr. Juan Gómez Luna

Prof. Onur Mutlu

ETH Zürich

Fall 2022

13 December 2022