

From Processing-in-Memory to Processing-in-Storage (PrinS)

Roman Kaplan

Advisor: Prof. Ran Ginosar

In collaboration with: Dr. Leonid Yavits

Processing-in-Memory (PiM): Idea & Problems

Storage (TBs-PBs)

Large Input (E.g., "big data")

Most of the computation is processed in-memory

CPU used for short non data-intensive tasks

Data transfer: The most energy- & time-consuming part

- PiM performs the bulk of computation in main memory
 - Reducing Data transfers
 - Saving time & energy
- Data transfers from storage to main memory are still required
- Storage ↔ main memory is the most energetically expensive link
- Moving large amounts of data takes a lot of time and energy
- To reach 100x lower energy and higher performance we **must bring processing closer to data**

Resistive Content Addressable Memory (ReCAM): A PrinS Device

Registers

Data row 128/256-bit

10⁶-10⁹ Rows
Row = Processing element

Resistive Bit-cells

- Based on resistive technology (e.g., memristors)
- High density & non-volatile

Logic next to storage

- Every data row also has a processing unit
- Instructions are performed bit-serially on all PUs in parallel

Entire ReCAM crossbar array can be GBs in size divided to separate ICs

- E.g., IC can be 256MB in size
- 256MB = 8M rows
- Future technology might allow TBs of storage

In-situ processing:

- Bit-wise logic
- Algebraic operations
- Column-shift
- Row-wise comparison
- CAM operations (explained in 'In-Storage Deduplication' below)

Bit-cell switching time can allow for a GHz operational frequency

PrinS Application: DNA Local Sequence Alignment

Smith-Waterman: Local Sequence Alignment

DNA comparison of human and fruit fly

"Highest similarity" region

- The Smith-Waterman (S-W) algorithm finds regions with "highest similarity" between two sequences (DNA or protein)
 - Proven to be optimal
- Every "match", "mismatch" & "gap" inside the pair of sequences has a score
- S-W is based on dynamic programming
 - Fills a $m \times n$ matrix
 - Has a quadratic computational complexity $O(m \cdot n)$
- Two main steps in the algorithm: (1) scoring and (2) traceback
 - We implemented the scoring step – more computationally demanding

Parallel S-W

- Matrix-fill order is on the main diagonal
- The entire antidiagonal is calculated in parallel
 - Antidiagonal per group of ReCAM columns
- We search for the maximal score
- Only 3 antidiagonals are stored in each iteration

In-ReCAM S-W

- Every score is a 32-bit integer = 32 ReCAM bitcells
- DNA sequence is stored in 2 columns (2 bits per base-pair)
- Sequence B shifts down at the beginning of an iteration

Dynamic Programming Matrix

	b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8
a_1				4,1	5,1	6,1		
a_2		4,2	5,2	6,2				
a_3	4,3	5,3	6,3					
a_4	4,4	5,4	6,4					
a_5	5,5	6,5						
a_6	6,6							
a_7								
a_8								

ReCAM (iteration start)

			A_1	B_1
6,1	5,1			
6,2	5,2	4,1	a_1	b_5
6,3	5,3	4,2	a_3	b_4
6,4	5,4	4,3	a_4	b_3
6,5	5,5	4,4	a_5	b_2
6,6	5,6		a_6	b_1
			a_7	
			a_8	

ReCAM (iteration end)

			A_1	B_1
6,1	5,1			
6,2	5,2	4,1	a_1	b_5
6,3	5,3	4,2	a_3	b_4
6,4	5,4	4,3	a_4	b_3
6,5	5,5	4,4	a_5	b_2
6,6	5,6		a_6	b_1
			a_7	
			a_8	

Simulations

- Cycle-accurate simulator: 8GB of storage @500MHz
- Compared to multi-accelerator state-of-the-art solutions: FPGA, Xeon Phi and GPU

Performance Compared to State-of-the-Art

Accelerator	FPGA	Xeon Phi	GPU	ReCAM
Performance (TCUPS)	6.02	0.23	11.08	52.68
# of ICs	128	4	384	32

Publications: [1] The PRIN-based High-Performance Computer Architecture for Applications in Bioinformatics, in Proceedings of the 2016 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2016, June 2016, pp. 383-392. Springer, 2016. [2] Y. Li and B. Schmitz, 2016. 256MB Dense Non-volatile Memory Array for Applications in Bioinformatics, in Proceedings of the 2016 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2016, June 2016, pp. 383-392. Springer, 2016. [3] S. Kim, E. C. O. et al., "100x Higher Performance than State-of-the-Art in GPU Clusters," IEEE Transactions on Parallel and Distributed Systems, vol. 28, no. 12, pp. 3581-3592, Dec 2017.

PrinS Application: In-Storage Deduplication

What is Deduplication?

- Deduplication is a technique for storing a single copy of each data block in storage
 - Can reach 10x reduction in data volume
- How it works:
 - Data is broken into fixed blocks
 - A fingerprint (FP) is calculated for each block
 - Only pointers are stored for identical blocks

Traditional (RAM+CPU) vs. In-ReCAM Deduplication

RAM+CPU Deduplication

New block write requires:

- FP calculation (hash)
- Search in a hash table
- Write to three different tables in RAM (A,B,C below)

ReCAM

Use CAM operations:

- Compare all storage content to a specific word → No need to hash
- Write to specific column in parallel
- Write to specific rows in parallel

Simulations: ReCAM vs. OpenDedup

Throughput vs. duplicate %

Energy vs. duplicate %

- ReCAM was simulated with a cycle-accurate simulator
 - ReCAM size = 256GB
 - Frequency = 1GHz
- Openedup was executed on a high-end server for comparison: 4x8 octa-core CPU, 64GB of RAM and 800GB SSD drive
- ReCAM has 100x higher throughput than deduplication with RAM+CPU**
- Energy consumption is similar or lower for the common block sizes (4 & 8KB)

Future: Deep Learning

Next Work

- Nowadays Deep Learning for large nets is done in cloud
- No energy figures reported, main concern is performance
- Deep Learning has two parts: (1) Feedforward (2) Backpropagation. Each requires different set of tools

Main Questions

- Which instruction set?
- Do we need new modules? (E.g., sum reduction)
- Will we be able to solve larger problems with PrinS?