

In-Memory Accelerator Architectures for Machine Learning and Bioinformatics

Roman Kaplan

PhD Advisor: Prof. Ran Ginosar

Resistive Content Addressable Memory for Machine Learning and Bioinformatics

Storage ↔ Memory Bottleneck Problem

Storage (TBs-PBs) → Limited Bandwidth
Storage to memory data transfers are the most time & energy consuming parts

Big Data Input → PIM performs the bulk of computation in main memory
✓ Reduces Mem → CPU transfers

Bulk of data-intensive computation → Storage → Mem data transfers are the main bottleneck

Short, non-data-intensive tasks → Processing in-storage can increase performance and reduce energy by 100x

Resistive Content Addressable Memory (ReCAM): A PRiNS Device

Registers → Data row 256/512 bit

Resistive Bitcells → Based on resistive elements (e.g., memristors)
✓ High density & non-volatile

Full ReCAM-Based Storage System → Divided to multiple ReCAM ICs
✓ E.g., IC can be 256MB = 8M rows
✓ Linearly scalable performance
→ More ICs = More parallelism

Logic next to data: *in-situ* processing
→ Processing unit (PU) per row
→ All PUs can operate in parallel
→ Instructions are based on bit-serial associative processing

PRiNS Application DNA Local Sequence Alignment

Smith-Waterman Algorithm

- Finds regions with "highest similarity" between two sequences (DNA/protein)
- Proven to be optimal
- Every "match", "mismatch" & "gap" has a score
- S-W is based on dynamic programming
- Fills a $m \times n$ matrix
- Has a quadratic computational complexity $O(m \cdot n)$
- We focus on the computationally-heavy scoring step

Matrix-fill order is on the main diagonal
Entire anti-diagonal is calculated in parallel
Anti-diagonal per group of ReCAM columns
Search for the maximal score
Only 3 anti-diagonals are required in each iteration

* This work was done with prof. Uri Weiser

Machine Learning Application K-Means and K-Nearest Neighbors

K-Means Clustering Algorithm

Common Machine Learning algorithm
For each sample, finds its group among possible K
Widely used in many fields, including:
Image processing
Anomaly detection
Data intensive
Multiple iterations over entire dataset

Two Main Loops:
1. For each Sample: For each Center: Calculate |Sample - Center|
2. For each Center: For each Sample: Calculate Sum(Sample)

$O(n)$ in a von Neumann machine
 $O(1)$ in ReCAM

PRiNS Implementation & Performance Comparison

In-Memory Computation

IC = Integrated Circuit
Shift operations move data between ICs

Iteration Start → Iteration End

IC = 1 → IC = 2

Performance Comparison

Cycle-accurate simulator: 8GB of storage running at 500MHz

Compared to multi-accelerator state-of-the-art solutions: FPGA, Xeon Phi and GPU

ReCAM shows 4.7x higher performance than a 384-GPU solution

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

PRiNS Implementation & Performance Comparison

In-Storage Computation

In large datasets: #Data Samples >> #Clusters
Computation over ALL samples in one instruction
Each sample attributes requires multiple temp fields
Difference between center's attribute
Squared difference
Total distance between sample and center

Comparison to other Works

Work	Platform	Samples	Attributes	Size on disk	Clusters
[1]	FPGA	1M	1	4MB	128
[2]	FPGA	2M	4	33.6MB	4
[3]	Intel i7	2.5M	68	318.8MB	10000
[4]	GPU	1.4M	5	21.3MB	240
[5]	10-GPU Cluster	18	40	152.7GB	120

Improved Architecture: Batch-Write NAND ReCAM

Batch-Write: Write One Output Value in One Cycle

Associative processor: one write cycle for truth table row
With Batch-Write: one write cycle for output value
All outputs with the same value are batched
New TAG logic to accumulate consecutive compares
In Bioinformatics: BLOSUM protein scoring matrix is 23x23
Most scores repeat multiple times

NAND CAM Bitcell
Match: discharge → no mismatch
Mismatch: no discharge → mismatch

NAND Data Row

K-Nearest Neighbors Algorithm

Common classification algorithm

Given a query sample, find the k most closest dataset samples
Large datasets cannot fit on memory of one chip
Requires communication between chips
On PRiNS, no data transfer between chip

Pseudocode of KNN on PRiNS

- Denotes the number of nearest neighbors. Every sample i if it may be stored in several consecutive ReCAM rows, the row associated with i per sample line simplicity.
- Each sample is characterized by A_i attributes.
- Calculate distance of each dataset sample from query Q : $Dist_i = \sum_{j=1}^n |A_{ij} - Q_j|^2$
- Choose k i 's in parallel, on entire ReCAM
- $dist_{min} = \min_i Dist_i$
- $dist_{max} = \max_i Dist_i$
- On microcontroller, $dist_{min} < dist_{max} + 1$
- (Assume unique $dist_{min}$ values)
- If $dist_{min} < dist_{max} + 1$, all classes maintained by microcontroller (Start with all samples connected)
- Long k time
- Find all unclassified samples
- Find each next row with min value of $dist_{min}$
- Continue until all samples are classified
- On microcontroller, $dist_{min} < dist_{max} + 1$
- Classification: Class with highest histogram

NAND Array: Improved Energy Efficiency

In associative processing, most compares result in mismatch
NAND CAM bitcell: 2T2R bitcells that discharge on match (instead of mismatch)
No discharge on mismatch: charge remains on mismatch
Less precharge energy is required for next compare operation

PRiNS Implementation & Performance Comparison

Distance of each dataset sample from query sample is computed

Temp fields store intermediate results
Distance over each dimensions is accumulated
Once distance of each dataset sample is found, closest sample is searched serially
In all cases: N is usually small (<100)
FindMin is performed in $O(1)$ on ReCAM

Work	Platform	Samples	Attributes	Size on disk	Clusters
[1]	GPU	1.4M	5	21.3MB	240
[2]	FPGA	2M	4	33.6MB	4
[3]	FPGA	1M	1	4MB	128
[4]	Intel i7	2.5M	68	318.8MB	10000
[5]	10-GPU Cluster	18	40	152.7GB	120

Performance and Energy Efficiency Results

BioPRiNS: PRiNS for Biological Sequence Search

BioPRiNS compared to five large-scale platforms, including PRiNS
Performance measured in Tera CUPS (TCUPS)
Comparisons over multiple datasets, DNA and protein
Results taken from publications (below)
BioPRiNS achieves 2.2-5.5x higher perf. than other solutions, with 2-156x better energy efficiency

Other works vs BioPRiNS
Performance (TCUPS) vs Energy Efficiency (TCUPS/Watt)

PRiNS Application: In-Storage Deduplication

Traditional (RAM+CPU) Systems

New block write:
1. Hash (create key)
2. Search in key table
3. Write to three tables in RAM (A, B, C)

In-ReCAM Deduplication

Use CAM operations:
1. CAM search → No need to hash
2. Write block + 1 pointer

ReCAM has 100x higher throughput than deduplication with RAM+CPU
Energy consumption is similar or lower for the common block sizes (4 & 8KB)

DNA Long Read Mapping

3rd Generation Sequencing: Real-Time Single-Molecule

DNA is composed of 4 nucleotides: 'A', 'G', 'C', 'T' (base pairs, bps)
Reading the content of an entire DNA strand at once (e.g., Chromosome) isn't possible
DNA sequencers output strands of the DNA (reads)
3rd generation sequencers vs. 2nd generation:
Produce outputs faster (hours vs. days)
Reads are longer (10k+ bps vs. 100-300 bps) → called long reads
Simpler preparation (less lab work)

The downsides: many errors, up to 15%
Errors can be insertions, deletions and substitutions

DNA Long Read Mapping

The problem: "stitching" together all the DNA reads
When an organism was previously sequenced (e.g., human), this sequence is used as a reference to construct the new organism genome
Requires mapping 1M+ reads against a reference sequence (e.g., human is 3Gbps)
Existing read mapping tools use technology-specific heuristics, complex data structures and large memory requirements
Hardware solutions only addressed short reads with few errors (from 2nd generation)
Long reads contain many errors and pose a challenge for an effective hardware acceleration solution

The Solution

RASSA: Resistive Approximate Similarity Search Accelerator

Processing-in-Memory Using Memristors

Basic cell: 1 memristor, 2 transistors (2T1R)
The memristor serves as a programmable non-volatile switch
Low resistance (R_{on}) - open switch
High resistance (R_{off}) - closed switch
One-hot encoding to store DNA base pairs
4 cells encode one base pair

Memristors (Memory Resistor)

- Change resistance if applied voltage higher than a threshold
- High density
- Non-volatile

Comparing a pattern

Compared pattern is applied on all bitcells
 R_{on} allows charge to flow
Match: no charge flow → N.C. in Match line
Mismatch: charge flows through R_{on} memristor → Match line voltage drops

The RASSA Bottom-Up Architecture

Sub-Word: Contains 60 bitcells → encodes 15 DNA base pairs (0 through 15 matches = 4 bit)
Analog-to-Digital converter translates match line voltage level to number of matches

Word Row: Contains 16 Sub-Words (240 bps)
Connected to top and bottom Word Rows
All mismatch values are summed (mismatch score)
Two options for a mismatch score:
1) Add mismatch score from top Word Row
2) Compare mismatch score to a threshold

DNA Read Mapping with RASSA

The Main Idea: High-Similarity Regions → Mapping Locations

A fixed-size chunk is compared against a long sequence
Counting mismatch score in every position approximates the correlation between the chunk and sequence
In the example, the read chunk is mismatch score is calculated for every position against the reference sequence
The mapping location for the chunk is the position of a mismatch score below the threshold

Find Long Read Mapping Locations with RASSA

Reads are divided to fixed-size chunks
For example: 100bps / 200bps
The threshold is set at 40-50% of the chunk length (determined empirically)
Every chunk is compared against the entire reference sequence (stored in RASSA)
Chunk with mismatch score below a threshold signals a mapping location
When chunk is split between two Word Rows, mismatch score for every part is found in separate cycles
Two consecutive Word Rows are needed
Mismatch score from top Word Row is transmitted to the bottom Word Row
Bottom Word Row sums and compares to threshold

Evaluations

Chip Parameters, Performance and Accuracy

Chip Parameters

A Sub-Word circuit was designed, placed and routed using the Global Foundries 28nm CMOS High-k Metal Gate library for:
Transistor sizing
Timing
Power analysis
Spectre simulations for the FF and SS corners at 70°C and nominal voltage

Performance and Accuracy Comparison with minimap2

Two organism reference sequences were used: e.coli and yeast
Input sequences from Pacific Biosciences (PacBio) and Oxford Nanopore Technologies (ONT)
For PacBio: regular and high-quality Circular Consensus Sequencing (CCS) reads
Accuracy and performance compared to minimap2, a state-of-art long read mapper
Executing platform: 16-core Intel Xeon E5-2650 @2GHz w/ 64GB of RAM
Minimap2 is was executed active SIMD extensions and multi-threading enabled
Sensitivity: % of reads found by RASSA from those found by minimap2

Organism	Platform	200bp chunk			100bp chunk		
		Sensitivity	False Positives	Speedup	Sensitivity	False Positives	Speedup
e.coli	PacBio CCS	79.3%	13.4%	25x	83.2%	13.6%	16x
	ONT	88.8%	10.5%	48x	87.6%	12.4%	31x
	RASSA	85.9%	34.9%	31x	85.1%	39.2%	49x
Yeast	PacBio CCS	69.8%	8.7%	77x	72%	11.8%	51x
	ONT	85.9%	34.9%	31x	85.1%	39.2%	49x
	RASSA	85.9%	34.9%	31x	85.1%	39.2%	49x

Performance Comparison with FPGA

Gatekeeper [1], a pre-alignment FPGA accelerator
Counts number of mismatches between short reads and a reference sequence
Implemented in a Virtex-7 FPGA using Xilinx VC709 board, running @250MHz
Host machine uses 3.6GHz Intel i7-3820 CPU w/ 8GB of RAM
Comparison of RASSA vs. GateKeeper throughput
Throughput measured in Billion Evaluated Mappings Locations per second (BEMLS)
GateKeeper results were taken from [1], RASSA results are normalized to 250MHz

RASSA vs. GateKeeper Throughput Comparison

Read Lengths	GateKeeper	RASSA @250MHz
100bp	1.7 BEMLS	226.8 BEMLS
200bp	-	175.2 BEMLS
300bp	0.2 BEMLS	142.8 BEMLS