

# FiltPIM: In-Memory Filter for DNA Sequencing

Marcel Khalifa\*, Rotem Ben-Hur, Ronny Ronen, Orian Leitersdorf, Leonid Yavits, and Shahar Kvatinsky  
*Viterbi Faculty of Electrical and Computer Engineering, Technion – Israel Institute of Technology, Haifa, Israel*  
 marcelkh@campus.technion.ac.il\*, rotembenhur@campus.technion.ac.il, ronny.ronen@technion.ac.il,  
 orianl@campus.technion.ac.il, yavits@technion.ac.il, shahar@ee.technion.ac.il

**Abstract**—Aligning the entire genome of an organism is a compute-intensive task. Pre-alignment filters substantially reduce computation complexity by filtering potential alignment locations. The base-count filter successfully removes over 68% of the potential locations through a histogram-based heuristic. This paper presents FiltPIM, an efficient design of the base-count filter that is based on memristive processing-in-memory. The in-memory design reduces CPU-to-memory data transfer and utilizes both intra-crossbar and inter-crossbar memristive stateful-logic parallelism. The reduction in data transfer and the efficient stateful-logic computation together improve filtering time by 100x compared to a CPU implementation of the filter.

## I. INTRODUCTION

DNA sequencing, the process of reading the genome of a given organism, is the foundation of many scientific and medical discoveries. For example, human genome sequencing enables personalized medicine and early diagnosis of genetic diseases [1]. A genome is a sequence of the bases [A, T, G, C], with the human genome being composed of approximately 3.2 billion bases. As reading the entire DNA sequence at once is infeasible for large-scale genomes, DNA sequencers typically extract sub-sequences called *reads*. These reads are orders of magnitude shorter than the whole genome sequence, ranging from dozens of bases for *short reads* to thousands of bases for *long reads* [2]. This work focuses on short reads.

DNA sequencing typically extracts numerous short reads, ranging from hundreds of millions to billions for the human genome. Then, through *read mapping*, a step towards the genome assembly, the reads are aligned to a reference genome of a similar organism. Many previous works have accelerated read mapping on CPU [3]–[5], GPU [6], [7], and FPGA [8], [9]. Standard hardware, however, suffers from the massive amount of data transfer between the memory and the processing unit (memory wall bottleneck). Therefore, *processing-in-memory* (PIM) platforms have gained interest. Such platforms inherently support parallel logic on sets of data residing in memory without the need to read the data, thereby reducing the memory wall bottleneck [10].

Read mapping can be divided into the *indexing*, *pre-alignment filtering*, and *sequence alignment* stages, as shown in Figure 1. To guide the mapping process, a full reference genome of a similar organism is utilized. Initially, the indexing stage generates potential locations on the reference genome for each read. We focus on the pre-alignment stage, where the similarity between a DNA read and reference fragments

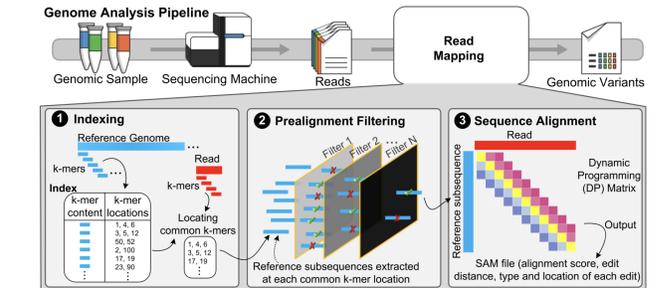


Fig. 1. DNA sequencing pipeline, demonstrating the *indexing*, *prealignment filtering*, and *sequence alignment* stages. Figure adapted from [1].  
 © 2020, IEEE.

is evaluated via heuristics. This stage filters false-positives for the computationally-intensive optimal sequence alignment.

Pre-alignment filtering aims to remove potential locations with edit distance above a certain predefined edit threshold (*eth*) compared to the read. Edits may be substitutions, insertions, or deletions. They are allowed due to genome variations among different organisms as well as sequencing errors. The base count filter [1] is a pre-alignment filtering algorithm that discards more than 68% of the potential locations by using a histogram-based heuristic for similarity comparison. Crucially, the base-count filter does not harm the sensitivity of the mapper as it does not discard true locations.

The base count filter operates on massive amounts of data, which makes it suitable for PIM. In this paper, we present FiltPIM, a memristive processing-in-memory accelerator that efficiently implements this filter. The PIM capabilities of FiltPIM are based on memristive stateful logic [11]. Stateful logic provides inherent inter-crossbar and intra-crossbar parallelism. By reducing data transfer and exploiting both types of parallelism, FiltPIM reduces filtering time by 100x compared to the CPU implementation of the filter.

## II. MEMRISTOR-AIDED LOGIC (MAGIC)

Memristive stateful logic is a PIM concept for representing data with resistance and performing logic operations on the resistances using the same cells in a memristive memory crossbar array [12]. One of the classic stateful logic techniques is MAGIC [11]. MAGIC implements logic gates using memristor cells within the same row (or within the same column) of a memristive crossbar array, see Figure 2. The inputs of the MAGIC gate are the states (resistance) of the input memristors prior to the computation, and the output is the resistance of the output memristor after computation. MAGIC NOR

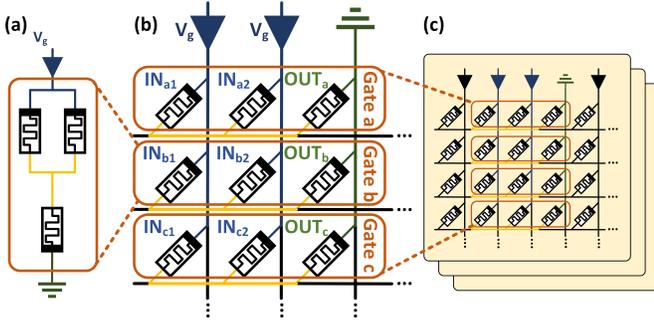


Fig. 2. (a) MAGIC NOR gate. (b) Parallel mapping of the MAGIC NOR gate to crossbar array rows, and (c) parallel computation across crossbars.

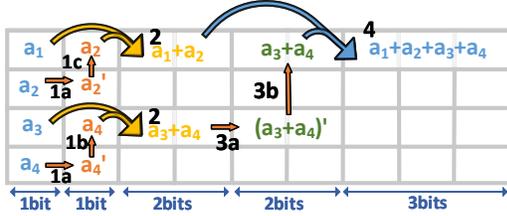


Fig. 3. Calculating the popcount of a 4-bit column using MAGIC NOR/NOT.

is executed in two steps (clock cycles): (1) initializing the output memristor to logical '1' (low resistance), (2) applying a voltage  $V_g$  across the gate.

MAGIC supports inherent parallelism as the same in-row gate can be performed in parallel across multiple rows (see Figure 2(b)) and across multiple crossbars (see Figure 2(c)). Table I lists different operations that are implemented using MAGIC NOR gates and specifies the number of cycles they require per single bit operands (initialization cycles are not included, as different initialization can be executed in parallel in the same clock cycle [13]). All operations, except *Popcount* (count 1's in a given column), are straightforward.

Popcount is a special case of the reduction algorithm introduced by Ronen *et al.* [14]. The algorithm is based on a recursive-tree technique that pairs up numbers and accumulates them in parallel, as seen in Figure 3. The first iteration is to add each two vertically adjacent numbers by aligning them in the same row (steps 1a, 1b, 1c in Figure 3), and then summing them (step 2). We iteratively continue with summing the vertically adjacent results (steps 3a, 3b, 4). In the last iteration, we will the sum of all bits residing at the end of the first row. A popcount of 100 bits requires 414 cycles.

### III. BASE COUNT FILTERING

The base-count filter introduced by Wendi *et al.* [17], aims to discard locations that have more than  $eth$  (edit-threshold) edits. The filter is based on comparing the base-count histograms of the read and the potential location. For each base type  $B \in [A, T, G, C]$ , the base error is defined as the absolute value of the count differences for the base. These base errors are accumulated to receive a single error. If that error is greater than  $2 \cdot eth$ , then the read and the potential location must have an edit distance of at least  $eth$

TABLE I  
MAGIC OPERATIONS

Operation	Cycles/bit	Notes
NOR/NOT	1	NOT is a 1-input MAGIC NOR
Copy	2	NOT(NOT(cell))
Half Adder	5	[15]
$N$ -bit Adder	$9 \cdot N + 1$	[14]
$N$ -bit Subtractor	$9 \cdot N + 1$	[16]
Popcount for 100 bits	414	Section II
$N$ -bit MUX( $X, Y, sel$ )	$4 \cdot N$	$((x_i + sel)^i + (y_i + sel)^i)^i$

and thus the potential location is discarded. As histograms are not influenced by permutation, then the heuristic does not remove all locations with more than  $eth$  errors.

CPU implementations of the filter, such as GASSST [18], limit the comparison to a sub-sequence of the read to reduce computation complexity. In FiltPIM, we exploit the parallelism of the memristive crossbars to compare base count for the entire read, thus ensuring higher precision (higher number of discarded locations) while reducing the execution time.

To evaluate the efficiency of the filter, we incorporate a CPU implementation of the algorithm into *mrFAST* [19], a state-of-the-art read mapping tool. Various human-genome data-sets with 100-base reads were considered: *ERR240726\_I*, *ERR240727\_I*, and *ERR240730\_I*<sup>1</sup>. The filter successfully discards more than 68% of *all potential locations*, while not affecting the sensitivity of the mapper as no true locations were discarded.

### IV. BASE-COUNT FILTERING WITHIN A MEMRISTIVE MEMORY ARRAY

In this section, we demonstrate an implementation of the filtering algorithm within a single memristive crossbar array. The array is structured to support reads of length up to 100-bases, yet it can be expanded to support other lengths.

We consider 128x256 arrays. 128 rows allow pre-storage of 100 bases of a reference genome within a column while reserving the remaining cells for temporary storage. The 256 columns enable pre-storage of a significant portion of the reference genome within each array:  $65 \cdot 100 = 6500$  bases, implying 130 columns as two adjacent bits are used to represent each base. An array is selected to participate in the filtering only when a fragment of the sequence containing the potential location fully resides within that array. The selected array receives the base-count histogram of the read (the number of occurrences of each base type [A, G, C, T]) and the exact potential location to be checked (a total of  $4 \times 7 + 32 = 60$  bits), and returns *true* if the location is to be discarded (returns 32 bits location + 1 bit result). To check the location, a series of operations inside the array are executed (see Figure 4):

- *Step 1:* For each base type  $B$ , the number of occurrences for  $B$  in the read,  $R_B$ , is stored to reserved cells in the array. 2 cycles to store each value. **8 cycles in total.**

<sup>1</sup>Available at [www.ebi.ac.uk/ena/browser](http://www.ebi.ac.uk/ena/browser).

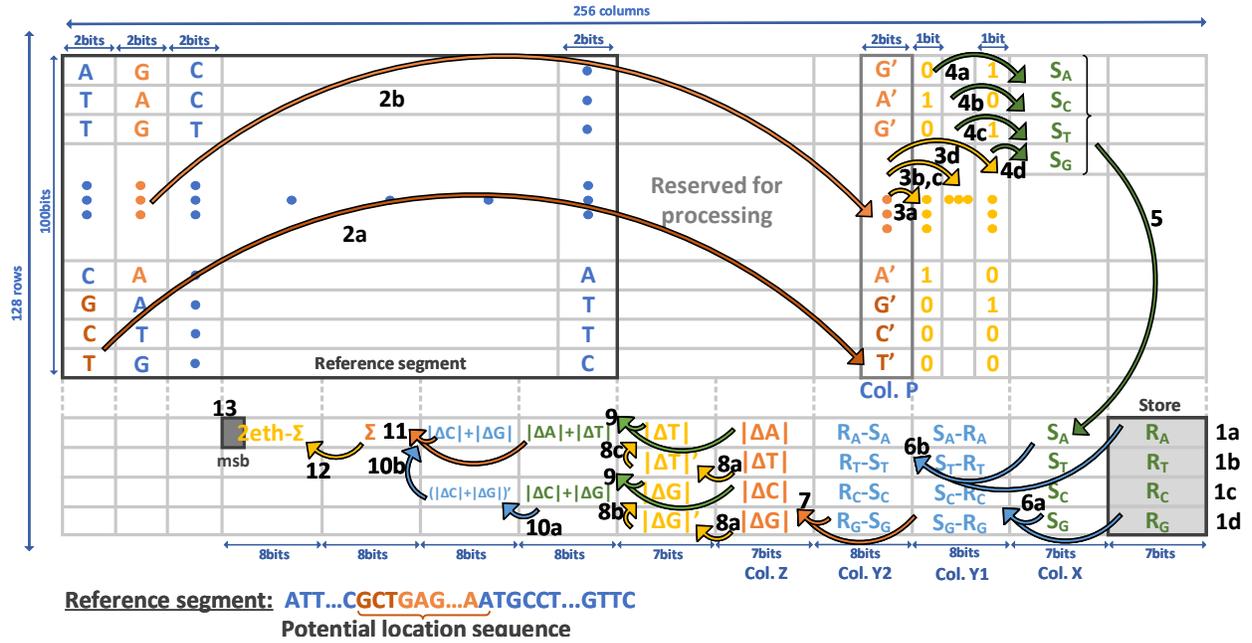


Fig. 4. Base-count filtering within a memristive memory array. The steps in black are explained in detail in Section IV.

- *Step 2:* A NOT operation is performed on the reference genome fragment of length  $read\_length$ , starting at  $potential\_location$ , to copy it (inverted) to column  $P$ . It is done in two steps, as the bases of the sub-sequence may span over two different 2-bit columns. Note that the bases will not overlap horizontally (as the number of bases in each column is each to the read length). **4 cycles in total.**
- *Step 3:* For each base type  $B$ ,  $P$  is compared against  $NOT(B)$ . Each base is represented by two bits, as follows: A="00", T="01", G="10", C="11". For each type  $B$ , a different function is performed to compare a base against it, as listed in Table II. **6 cycles in total.**
- *Step 4:* For each base type  $B$ , the matches of  $B$  in  $P$  are counted by popcounting the 100-bit columns computed in step 3. **1656 cycles in total.**
- *Step 5:* For all base types, the base counts in a potential location's sequence and in the read are aligned by copying the results from step 4 to column  $X$ . Two NOTs are performed to copy each result. **8 cycles in total.**
- *Step 6:* For all base types, in parallel, perform  $S_B - R_B$  and then compute their 2's complement ( $R_B - S_B$ ) by inverting all bits followed by performing Half Adder (HA) 8 times along the bits. **111 cycles in total.**
- *Step 7:* For all base types, in parallel, calculate  $|S_B - R_B|$  by choosing the non-negative value between both values calculated in step 6. To do this, we apply a MUX operation (see Table I), where the *most significant bit* (msb) of the result in column  $Y_1$  is the selector. The result in column  $Y_2$  is  $X$  and the result in column  $Y_1$  is  $Y$ . If, for example, the selector is '0' ( $msb_{(S_B - R_B)} = '0'$ ), we conclude that  $(S_B - R_B) \geq 0$  and that this value shall be chosen. **28 cycles in total.**

TABLE II  
BASE COMPARISONS

Compare "ab" against	#Cycles	Function	Notes
NOT(A)	3	NOR(NOR(a), NOR(b))	
NOT(T)	1	NOR(NOR(a), b)	NOR(a) already computed
NOT(G)	1	NOR(a, NOR(b))	NOR(b) already computed
NOT(C)	1	NOR(a, b)	

- *Steps 8-11:* Addition is performed between all differences in counts (four values in column  $Z$ ). This is performed in two iterations. In each iteration each two vertically adjacent values are aligned in the same row, and then full-adder (FA) is performed. **153 cycles in total.**
- *Step 12:* The sum of all differences in count is subtracted from the pre-defined constant ( $2 \cdot eth$ ). **72 cycles in total.**
- *Step 13:* Read out the msb of the final result of step 12. **1 cycle in total.** If the msb of the result is '1', meaning the sum of differences in counts is greater than twice the error threshold, then the potential location is discarded.

In summary, using in-crossbar array operations, the validity of a potential location for a given read can be checked in **less than 3000 MAGIC-NOR cycles** including initialization cycles (less than 2000 cycles without initialization cycles).

## V. EVALUATION

To evaluate FiltPIM, we compare its performance against a CPU-based implementation of an equivalent base-count filter. The advantage of FiltPIM originates from its massive intra-crossbar and inter-crossbar parallelism, and from the reduced CPU-to-memory data transfer.

TABLE III  
EVALUATION RESULTS FOR FILTPIM

	<i>Compute</i>	Notes	FiltPIM
a	Cycle time (ns)	[14]	10
b	PIM Crossbars		500,000
c	Latency per iteration (cycles)		3,000
d	Latency per iteration (ns)	$a \times c$	30,000
e	# Potential locations		$46 \cdot 10^9$
f	# Iterations	$5 \times (e/b)^1$	460,000
g	Total latency (sec)	$d \times f$	13.8
<i>Transferred Data</i>			
h	Transfer per potential location (B)	data in&out	$8 + 5 = 13$
i	Data transfer rate (GB/sec)		10
j	Total data transferred (GB)	$h \times e$	598
k	Data transfer latency (sec)	$j/i$	59.8
<i>Total Time (sec)</i>			$k + g$
			73.6

<sup>1</sup> Iterating 5 times the average locations per crossbar ( $5 \cdot 46G/500K = 460K$ ) in each crossbar, covers over 99% of all locations, decreasing the filter's efficiency by less than 0.7% (inferred from the locations distribution among the crossbars).

The human genome consists of approximately 3.2 billion bases. We pre-store the genome in 500,000 crossbars, each containing  $64+1=65$  read-size (100 bases) *fragments*. The genome portion within a crossbar has dimensions 100 *rows*  $\times$  130 *columns*. The first and the last fragments in each crossbar overlap with the neighboring crossbars to guarantee that each potential-location fragment resides in a single crossbar.

A set of human-genome reads (*ERR240727\_1*) is considered. When aligning the reads against a human genome reference (*humanG1Kv37*), the indexing stage of mrFast provides a total of 46 billion potential locations.

To evaluate the improvement of FiltPIM over CPU, we developed an optimized standalone implementation for the base-count filter, and measured the latency for checking 46 billion potential locations. The tool was executed on an Intel(R) Xeon(R) CPU E5-2683 v4 at 2.1 GHz, with 256GB of DRAM, 2400MHz DDR4, 2x 1TB HD and 480GB SSD. We observed total latency of approximately 7360 seconds, of which 70% are for data transfer and 30% for computation.

We analyzed FiltPIM performance using the algorithm from Section IV. We assumed the same CPU-to-memory interface (2400MHz DDR4). The actual data transfer rate may depend on the full architecture of FiltPIM, which is out of scope for this paper. Therefore, we assumed 10 GB/s, about half of the peak performance rate for DDR4 2400 (19.2 GB/s). Table III summarizes the results.

The results show that FiltPIM is faster than the CPU on both the computation time (160x) and the data transfer time (86x). In total, it reaches a speedup of 100x over the CPU.

We use the Bitlet model [14] to evaluate the power consumption of FiltPIM. According to Bitlet, simultaneous operation of 1000 arrays, each using up to 100 rows, consumes 1W. If we set the power budget limit at 100W, only 100K arrays can operate simultaneously. This increases the compute (total) time by 5x (1.75x). Figure 5 shows how increasing the number of arrays working in parallel in FiltPIM affects its performance. For over 200 working arrays, FiltPIM outperforms the CPU.

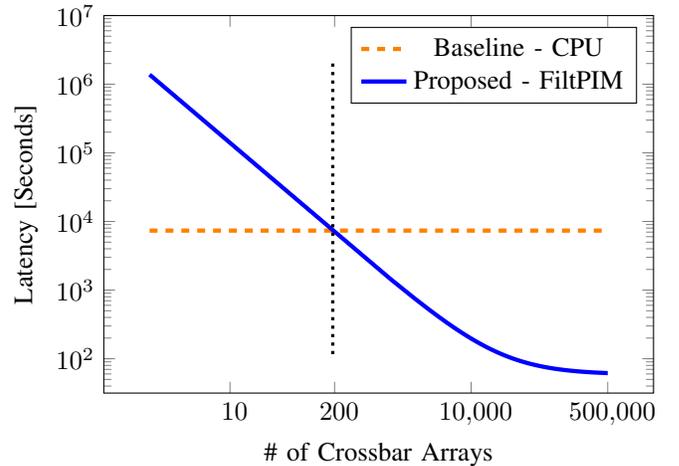


Fig. 5. Execution time of two platforms running the same filter, the CPU and FiltPIM. Lower latency means better performance.

## VI. CONCLUSION

A base count filter improves the performance of DNA read mapping by reducing the work for the computation-intensive sequence alignment stage. This paper introduces FiltPIM, a novel PIM base-count filter that accelerates filtering by 100x compared to a CPU based base-count filter. This gain is due to reduction in CPU-to-memory data transfer and to the massive parallelism enabled by memristive crossbars. In future work, we aim to build a complete architecture for FiltPIM and measure its benefit to the entire read mapping pipeline.

## VII. ACKNOWLEDGMENT

This work was supported by the European Research Council through the European Union's Horizon 2020 Research and Innovation Program under Grant 757259.

## REFERENCES

- [1] M. Alser *et al.*, "Accelerating genome analysis: A primer on an ongoing journey," *IEEE Micro*, 2020.
- [2] S. Goodwin *et al.*, "Coming of age: ten years of next-generation sequencing technologies," *Nature Reviews Genetics*, 2016.
- [3] J. Daily, "Parasail: SIMD C library for global, semi-global, and local pairwise sequence alignments," *BMC bioinformatics*, 2016.
- [4] R. Li *et al.*, "SOAP: short oligonucleotide alignment program," *Bioinformatics*, 2008.
- [5] B. Langmead *et al.*, "Ultrafast and memory-efficient alignment of short dna sequences to the human genome," *Genome biology*, 2009.
- [6] N. Ahmed *et al.*, "GASAL2: a gpu accelerated sequence alignment library for high-throughput ngs data," *BMC bioinformatics*, 2019.
- [7] A. Zeni *et al.*, "LOGAN: High-performance gpu-based x-drop long-read alignment," in *IPDPS*, 2020.
- [8] M. Alser *et al.*, "GateKeeper: a new hardware architecture for accelerating pre-alignment in DNA short read mapping," *Bioinformatics*, 2017.
- [9] P. Chen *et al.*, "Accelerating the next generation long read mapping with the FPGA-based system," *IEEE/ACM TCBB*, 2014.
- [10] S. Kvatinisky, "Real processing-in-memory with memristive memory processing unit (mmpu)," in *ASAP*, vol. 2160-052X, 2019, pp. 142-148.
- [11] S. Kvatinisky *et al.*, "MAGIC—memristor-aided logic," *IEEE TCAS-II*, vol. 61, no. 11, pp. 895-899, 2014.
- [12] J. Reuben *et al.*, "Memristive logic: A framework for evaluation and comparison," in *PATMOS*, 2017.
- [13] R. Ben-Hur *et al.*, "SIMPLER MAGIC: synthesis and mapping of in-memory logic executed in a single row to improve throughput," *IEEE TCAD*, 2020.

- [14] R. Ronen *et al.*, "The bitlet model: A parameterized analytical model to compare pim and cpu systems," *ACM JETC*, 2021.
- [15] P. L. Thangkhiew *et al.*, "Efficient implementation of adder circuits in memristive crossbar array," in *TENCON*, 2017.
- [16] P. Saini, "Full subtractor." [Online]. Available: <https://gateoverflow.in/31359/minimum-nand-gates-realization-exor-exnor-adder-subtractor>
- [17] W. Wang, P. Zhang, and X. Liu, "Short read dna fragment anchoring algorithm," *BMC bioinformatics*, vol. 10, no. 1, pp. 1–11, 2009.
- [18] G. Rizk and D. Lavenier, "GASSST: global alignment short sequence search tool," *Bioinformatics*, 2010.
- [19] H. Xin *et al.*, "Accelerating read mapping with fasthash," in *BMC genomics*. Springer, 2013.