

## Efficient Memory Layout for Pre-Alignment Filtering of Long DNA Reads Using Racetrack Memory

Khan, Asif Ali; Hameed, Fazal; Shahroodi, Taha; Jones, Alex K.; Castrillon, Jeronimo

**DOI**

[10.1109/LCA.2024.3350701](https://doi.org/10.1109/LCA.2024.3350701)

**Publication date**

2024

**Document Version**

Final published version

**Published in**

IEEE Computer Architecture Letters

**Citation (APA)**

Khan, A. A., Hameed, F., Shahroodi, T., Jones, A. K., & Castrillon, J. (2024). Efficient Memory Layout for Pre-Alignment Filtering of Long DNA Reads Using Racetrack Memory. *IEEE Computer Architecture Letters*, 23(1), 129-132. <https://doi.org/10.1109/LCA.2024.3350701>

**Important note**

To cite this publication, please use the final published version (if applicable).  
Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights.  
We will remove access to the work immediately and investigate your claim.

***Green Open Access added to TU Delft Institutional Repository***

***'You share, we take care!' - Taverne project***

**<https://www.openaccess.nl/en/you-share-we-take-care>**

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.

# Efficient Memory Layout for Pre-Alignment Filtering of Long DNA Reads Using Racetrack Memory

Asif Ali Khan<sup>1</sup>, Fazal Hameed<sup>2</sup>, Taha Shahroodi<sup>3</sup>, Alex K. Jones<sup>4</sup>, *Fellow, IEEE*,  
and Jeronimo Castrillon<sup>5</sup>, *Senior Member, IEEE*

**Abstract**—DNA sequence alignment is a fundamental and computationally expensive operation in bioinformatics. Researchers have developed *pre-alignment* filters that effectively reduce the amount of data consumed by the alignment process by discarding locations that result in a poor match. However, the filtering operation itself is memory-intensive for which the conventional Von-Neumann architectures perform poorly. Therefore, recent designs advocate compute near memory (CNM) accelerators based on stacked DRAM and more exotic memory technologies such as *racetrack memories* (RTM). However, these designs only support small DNA reads of circa 100 nucleotides, referred to as *short reads*. This letter proposes a CNM system for handling both long and short reads. It introduces a novel data-placement solution that significantly increases parallelism and reduces overhead. Evaluation results show substantial reductions in execution time ( $1.32\times$ ) and energy consumption (50%), compared to the state-of-the-art.

**Index Terms**—Domain wall memory, near memory computing, racetrack memory, sequence alignment.

## I. INTRODUCTION

HIGH-THROUGHPUT sequencing (HTS) technologies produce an extensive volume of sequencing data, encompassing both *short reads* (consisting of hundreds of nucleotides) and *long reads* (extending to thousands of nucleotides). Subsequently, the generated raw data is commonly used in diverse analyses, such as disease diagnostics and genetic evolution, among others. Sequence alignment represents a fundamental and computationally demanding stage within these analysis pipelines [1].

Several algorithms are employed for sequence alignment. Dynamic programming based solutions are accurate but scale poorly with the problem size [2]. Seed-and-extension algorithms are more scalable [3], [4] but still exhaust considerable effort on genome locations that eventually do not align. Many of these poor alignment locations can be easily discarded using pre-alignment filters [5], [6], [7].

Manuscript received 2 November 2023; accepted 1 December 2023. Date of publication 19 January 2024; date of current version 30 May 2024. This work was partially supported in part by the German Research Council DFG under Grant 502388442 and Grant 450944241, in part by US NSF under Grant 1822085 and Grant 2133267, in part by LPS, and in part by NSA. (*Corresponding authors: Asif Ali Khan.*)

Asif Ali Khan and Jeronimo Castrillon are with TU Dresden, 01069 Dresden, Germany (e-mail: asif\_ali.khan@tu-dresden.de; jeronimo.castrillon@tu-dresden.de).

Fazal Hameed is with the American University of Sharjah (AUS), Sharjah, UAE (e-mail: fazal.hameed@tu-dresden.de).

Alex K. Jones is with the University of Pittsburgh, Pittsburgh, PA 15261 USA (e-mail: akjones@pitt.edu).

Taha Shahroodi is with TU Delft, 2628 Delft, The Netherlands (e-mail: t.shahroodi@tudelft.nl).

Digital Object Identifier 10.1109/LCA.2024.3350701

Pre-alignment filters significantly reduce the number of DNA fragments that are passed on to the computationally expensive alignment phase by eliminating the apparent mismatches. This improves the end-to-end performance of the aligners by circa  $2\times$  [5]. These filters employ lightweight algorithms to compute a similarity score for each candidate location. Detailed alignment operations are only performed at locations that pass a set threshold on the similarity score.

For pre-filtering to accelerate the overall processing time effectively, it must be considerably faster than the actual alignment operation. However, considering the larger genomics data sets, this can become challenging. To reduce their computational time, pre-alignment filters have been implemented in CPUs, GPUs, FPGAs, ASICs, and computing-near-memory (CNM)<sup>1</sup> solutions [5], [6], [7], [8]. The CNM solutions dramatically reduce the data movement on the external channel between the memory and the processor.

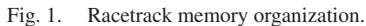
Pre-alignment filtering is particularly beneficial for long reads, which are predicted to revolutionize genomics [9] but require considerably more resources for processing. Unfortunately, they have received no attention from the community. SneakySnake is the only available filtering algorithm for long reads but is not compatible with CNM acceleration [7]. On the contrary, GRIM [5], ALPHA [8] and FIRM [6] are CNM-accelerated but do not support long reads.

This letter proposes an extension to the FIRM filter to support long reads. Concretely, we introduce algorithmic modifications to the FIRM filter to support long reads. Similar to FIRM [6], our CNM-accelerator leverages the unprecedented density, energy, and performance efficiencies of *racetrack memory* (RTM) [10]. However, the FIRM memory layout necessitates substantial data duplication, resulting in significant performance degradation when handling long reads. To address this challenge, we present an innovative RTM layout for our CNM accelerator that supports both short and long reads. Our layout minimizes RTM shifts to an absolute minimum and maximizes parallelism for varying numbers of parallel working bin units in the accelerator. To our best knowledge, this is the first CNM solution that efficiently handles long reads.

## II. BACKGROUND AND RELATED WORK

RTM is a spintronic nonvolatile memory. A single cell in RTM is a magnetic nanowire, or *track* that can be split into circa 100 magnetic regions, or *domains* separated by *domain walls* (see Fig. 1, top right). During an RTM access, a domain must be *shifted* and aligned to an access port (AP) position. RTM

<sup>1</sup>This is also sometimes referred to as processing near memory (PNM) or processing in memory (PIM).



RTM exhibits several promising properties, including SRAM-class access speed, lower energy requirements, SRAM/DRAM-class endurance, and exceptional density. For DNA pre-alignment filtering, similar to FIRM, we favor RTM over other nonvolatile memories (NVMs) because of two primary considerations: (i) the remarkable density of RTMs, essential for accommodating entire reference genomes, and (ii) the inherently sequential memory access pattern of the application which perfectly aligns with the sequentiality in RTMs.

Pre-alignment filters proactively skip poor alignment mapping positions, minimizing the number of extend operations in read aligners [5], [6], [8]. Among others, the GRIM filter has demonstrated effective utilization of 3D-stacked memory devices, leading to significant performance improvements in read mappers through substantial parallelism [5].

### B. State-of-The-Art Pre-Alignment Filters

[illegible]

bin units in the accelerator (shown in green on the left). Let us call these simultaneously executed bins a *binset*. For a fair comparison with existing filters, we assume the same binset and the memory access granularity, i.e., 4k bits.

ALPHA [8] exploits the tokens repetition in query reads to minimize the number of memory accesses in the GRIM filter. Concretely, if the set of distinct tokens in  $q_i$  is represented by  $\Theta$ , i.e.,  $\Theta(q_i) = \{\text{distinct tokens in } q_i\}$ , ALPHA maintains each token count value in a dedicated buffer called *CountBuffer*. In the filtering processing, instead of accessing a row of bitvectors multiple times for repeated tokens in a read, ALPHA only accesses rows storing unique tokens to compute the sum value. For instance, for a particular bin  $r_k \in \mathbb{R}$ , the accumulated sum is computed as  $\sum_{j \in \Theta(q_i)} \text{CountBuffer}(j) \cdot \vec{r}_k[f(j)]$ .

### III. PRE-ALIGNMENT FILTERING FOR LONG READS

To filter long reads without significant overhead, we keep the bin and token sizes unchanged and instead split the long reads into *chunks* with sizes comparable to the reference bins. Each

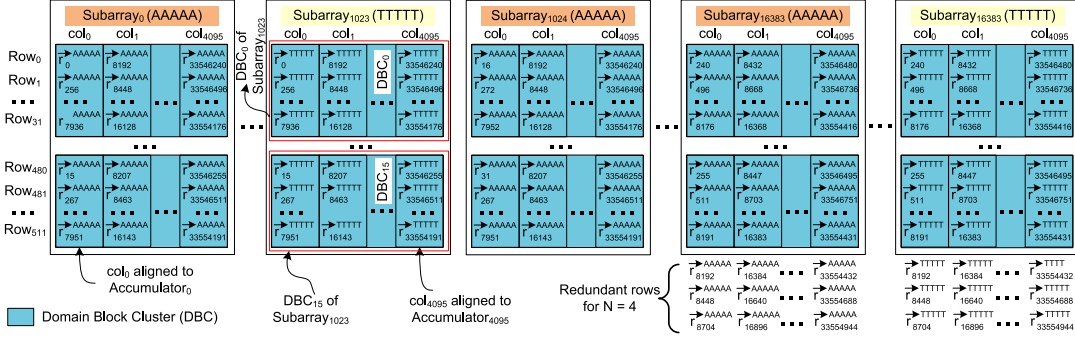


Fig. 3. Our optimized data placement for long reads, with backward compatibility for short reads.

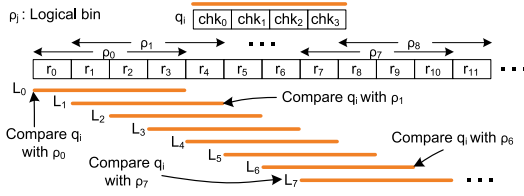


Fig. 4. Overview of the pre-filtering for long reads.

chunk is compared with successive bins in the reference genome, and the similarity score of individual chunks is accumulated to compute the overall similarity score. For instance, in Fig. 4, the long read  $q_i$  is divided into  $N = 4$  chunks:  $chk_0, chk_1, chk_2$ , and  $chk_3$ . In order to find the mapping similarity score of  $q_i$  at any location in the reference genome, referred to as a *logical* bin,  $p_j$ , all chunks of  $q_i$  having separate count buffers are compared to successive bins in R starting at the location  $r_j$ . For instance, if we want to compute the similarity of  $q_i$  with  $\rho_0$  at location  $r_0$ ,  $chk_0$  is compared to  $r_0$ ,  $chk_1$  to  $r_1$ ,  $chk_2$  to  $r_2$ , and  $chk_3$  to  $r_3$  and the similarity score of the comparisons is accumulated:  $\sum_{i=0}^{N-1} \sum_{j \in \Theta(chk_i)} \text{CountBuffer}^{chk_i}(j) * r_{k+i}^v[f(j)]$  where  $k$  is the starting position in the reference genome (0 in this case).

Intuitively, the proposed *split-and-compare* mechanism for long reads should not affect the filtering accuracy as we retain the small token window of short reads. Nonetheless, we validate the correctness of our filtering algorithm by comparing its output with the BWA-mem and minimap2 aligners. We ensure that the bin positions filtered out are not contained in the output of the two aligners.

#### A. CNM Architecture for Short and Long Reads

FIRM [6] only handles short reads. In order to use the FIRM accelerator for long reads, it requires data duplication and column shifting. To explain it further, let us assume we want to compare the first 4k logical bins (i.e.,  $\rho_0$  to  $\rho_{4095}$ ) in R with a given long read  $q_i$  having  $N = 4$ . For the comparison to work, the last  $N - 1 = 3$  logical bins (i.e.,  $\rho_{4093}, \rho_{4094}$ , and  $\rho_{4095}$ ) requires three redundant columns for  $r_{4096}, r_{4097}$ , and  $r_{4098}$  in each subarray. Second, in order to use the same accumulator for all  $N$  chunks, this data layout requires shifting each fetched row for alignment prior to summation, as in the accelerator design a single accumulator to use to compute the similarity score of one entire reference bin. The shifting problem can be avoided by using  $N$  accumulators per long read. However, this requires considerable changes and hardware overhead to the accelerator design to sum up all the partial results of the  $N$  accumulators. This also incurs energy overhead and reduces the parallelism degree from 4096 to  $4096/N$ .

Authorized licensed use limited to: TU Delft Library. Downloaded on June 20, 2024 at 11:48:07 UTC from IEEE Xplore. Restrictions apply.

Alternatively, placing successive bins (belonging to the same logical bin) in adjacent rows also solves the accumulator alignment problem. However, in RTM-based systems, it may result in a notable increase in RTM shifts. In the following, we introduce a novel data placement mechanism that resolves the accumulator alignment problem while reducing the number of RTM shifts. In the following, we propose a novel data placement mechanism that solves the accumulator alignment problem while reducing the number of RTM shifts.

1) *Optimized Data Placement for Long and Short Reads:* Fig. 3 shows our novel data layout that supports both long and short reads. Note that the subarray layout (e.g., Subarray<sub>0</sub>) is similar to Fig. 1, i.e., each subarray has 16 DBCs, each DBC has 4096 nanowires (number of columns) and each nanowire has 32 bits (e.g., DBC<sub>0</sub> contains Row<sub>0-31</sub>). In the proposed mapping, presence bits of the reference bins are stored in a subarray-interleaved fashion, as indicated by the subarray labels (e.g., Subarray<sub>0</sub> (AAAAA), meaning this subarray exclusively stores presence bits for token AAAAA). This interleaving reduces the number of active subarrays and minimizes within-subarray movements (RTM shifts), as subarrays dedicated to tokens that are not present in the query read are never accessed. Presence bits of successive bins (e.g.,  $r_{0-3}$ ) in a logic bin ( $\rho_0$ ) are placed in the same column (col<sub>0</sub>) and the same index of adjacent DBCs (DBC<sub>0-3</sub>). The row and column offsets (see Fig. 1) in our layout are parameterized so that they can work with different memory sizes and configurations. For instance, for the memory configuration in Fig. 3, column offset is determined as  $S \times W/1024 = 8192$ , where  $S$  is the number of subarrays and  $W$  is the rows per subarray. The row offset is  $o = D \times S/1024 = 256$ , where  $D$  is DBCs per subarray. This layout does not require column shifting as all bins belonging to logical bin  $p_j$  are already aligned to the same accumulator.

To demonstrate how this layout solves the accumulator alignment problem, consider the example in Fig. 5 where all bins of  $\rho_0$  (i.e.,  $r_0, r_1, r_2$ , and  $r_3$ ) are aligned to Acc<sub>0</sub>. Similarly all bins pertinent to  $\rho_{33546240}$  are aligned to Acc<sub>4095</sub>. Compared to FIRM's layout where  $(N - 1) \times S$  redundant columns were needed for long reads, which accumulates to a considerable 19.3% storage overhead for a subarray with  $W = 512$ , our new data placement for long reads requires  $(N - 1) \times 1024$  redundant rows (required for the last 1024 subarrays) for the entire memory, which is an overhead of 1.2% for a 8GB memory. Note that this data layout also supports short reads and produces comparable results to the layout in FIRM [6].

#### IV. EVALUATION

For evaluation, we use four paired-end input short read genomes of size 100, same as in [6], and four long read genomes [11] of size 10 K (L1-10 kb, L2-10 kb) and 15 K



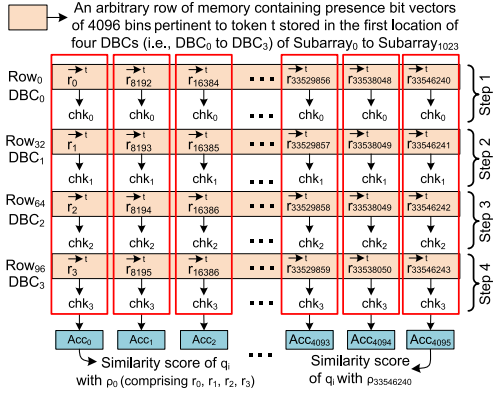


Fig. 5. An example showing how the novel data placement solves the accumulator misalignment problem for long reads.

TABLE I  
ENERGY AND LATENCY VALUES [12], [13], [14]

Operation	Value
Activation and precharge: DRAM/RTM	1964/1087 [pJ]
DRAM energy: Access / IO	1.25/0.40 [pJ]/bit
RTM energy: Read / shift	0.76/0.23 [pJ]/bit
Background power: DRAM/RTM	410/212 [mW]
RTM shift latency	1.87 (2 cycles)
$t_{RAS-tRCD-tRP-tCAS-tWR}$ (DRAM)	20-8-8-8-8 [cycles]
$t_{RAS-tRCD-tRS-tCAS-tWR}$ (RTM)	9-4-2S-4-4 [cycles]
Accel. power per bit (dynamic/leakage)	$1.77 \cdot 10^3$ / 11.16 [mW]
Preproc. unit per bit (dynamic/leakage)	11.6 / 1.13 [mW]

(L1-15 kb, L2-15 kb). For the sake of comparison to the GRIM and ALPHA filters, we assume bin and token sizes of 100 and 5, respectively.

We model an 8GB memory having 64 banks and 512 subarrays per bank using the cycle-accurate RTM simulator RTSIM [15]. The energy and latency numbers of the logical components and the memory subsystems shown in Table I are estimated using McPat [13], DRAMSpec [12], DESTINY [14], and Cadence RTL-Compiler Synthesis for the custom components on the logic layer, all presuming 32 nm fabrication technology.

We compare the pre-alignment filtering runtime and energy consumption of our design with the state-of-the-art *GRIM* [5], *ALPHA* [8] and *FIRM* [6] designs. We also explicitly compare the naive *FIRM* layout (*CNM-LC*) with our shift and performance aware layout (*CNM-LI*) and report the energy and performance results. All results include the query read preprocessing, i.e., populating the count buffers, loading the metadata of *R* and the data transfers between the memory and the logic layers.

#### A. Results and Discussions

The performance and energy efficiency of the proposed RTM-based CNM system largely depends on access parallelism and shift cost. Naively extending the short reads' optimized layout proposed in *FIRM* for long reads, noted as *CNM-LC*, leads to performance and energy inefficiencies. As discussed in Section III-A1, our *CNM-LI* avoids unnecessary shifts compared to *CNM-LC* data placement, which reduces the shifting overhead in RTM by  $11.8\times$  for long reads. This shift reduction is because our *CNM-LI* design requires a single shift per access. As a result, *CNM-LI* reduces the execution time on average by  $1.32\times$  compared to *CNM-LC*. It also reduces the energy consumption by 50% compared to the *CNM-LC* (see Fig. 6).

For short reads, our proposed data mapping produces comparable results to *FIRM* [6] while outperforming *GRIM* and *ALPHA* by a factor of  $3.6\times$  and  $3.4\times$  respectively. Moreover, it reduces energy consumption by  $2.5\times$  and  $2.3\times$  compared

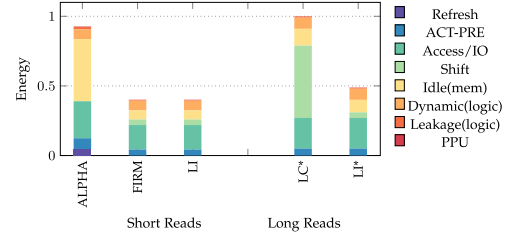


Fig. 6. Short and long reads energy comparison.

to *GRIM* and *ALPHA*, respectively (see Fig. 6). The energy results analysis shows that the background and the refresh energy account for a significant portion of the total energy in the DRAM-based designs, *GRIM* and *ALPHA*.

#### V. CONCLUSION

In this letter, we present an RTM based CNM system for energy-efficient pre-alignment filtering. We propose a novel data mapping that significantly reduces the data duplication and RTM shift operation overheads and maximizes parallelism. Our exploration of the design space and experimental evaluation show that our proposed layouts support both long and short reads. For short reads, our novel mapping produces comparable results to *FIRM* while for long reads, it reduces the execution time and energy consumption by  $1.32\times$  and 50%, compared to the data mapping presented in *FIRM*.

#### REFERENCES

- [1] R. C. Edgar et al., "Petabase-scale sequence alignment catalyses viral discovery," *Nature*, vol. 602, no. 7895, pp. 142–147, 2022.
- [2] T. F. Smith and M. S. Waterman, "Identification of common molecular sub sequences," *J. Mol. Biol.*, vol. 147, no. 1, pp. 195–197, Mar. 1981.
- [3] M. Vasimuddin, S. Misra, H. Li, and S. Aluru, "Efficient architecture-aware acceleration of BWA-MEM for multicore systems," in *Proc. IEEE Int. Parallel Distrib. Process. Symp.*, 2019, pp. 314–324.
- [4] H. Li, "Minimap2: Pairwise alignment for nucleotide sequences," *Bioinformatics*, vol. 34, no. 18, pp. 3094–3100, 2018.
- [5] J. Kim et al., "GRIM-Filter: Fast seed location filtering in DNA read mapping using processing-in-memory technologies," *BMC Genomic.*, vol. 19, no. 2, 2018, Art. no. 89.
- [6] F. Hameed, A. A. Khan, S. Ollivier, A. K. Jones, and J. Castrillon, "DNA pre-alignment filter using processing near racetrack memory," *IEEE Comput. Architecture Lett.*, vol. 21, no. 2, pp. 53–56, Jul.–Dec. 2022.
- [7] M. Alser et al., "SneakySnake: A fast and accurate universal genome pre-alignment filter for CPUs, GPUs and FPGAs," *Bioinformatics*, vol. 36, no. 22/23, pp. 5282–5290, 2020.
- [8] F. Hameed, A. A. Khan, and J. Castrillon, "ALPHA: A novel algorithm-hardware co-design for accelerating DNA seed location filtering," *IEEE Trans. Emerg. Topics Comput.*, vol. 10, no. 3, pp. 1464–1475, Jul.–Sep. 2022.
- [9] G. Logsdon et al., "Long-read human genome sequencing and its applications," *Nature Rev. Genet.*, vol. 21, no. 10, pp. 597–614, 2020.
- [10] R. Bläsing et al., "Magnetic racetrack memory: From physics to the cusp of applications within a decade," *Proc. IEEE*, vol. 108, no. 8, pp. 1303–1321, Aug. 2020.
- [11] "Genome in a bottle data indexes," Accessed: Aug. 9, 2023. [Online]. Available: [https://github.com/genome-in-a-bottle/giab\\_data\\_indexes/tree/master/AshkenazimTrio](https://github.com/genome-in-a-bottle/giab_data_indexes/tree/master/AshkenazimTrio)
- [12] O. Naji, C. Weis, M. Jung, N. Wehn, and A. Hansson, "A high-level DRAM timing, power and area exploration tool," in *Proc. Int. Conf. Embedded Comput. Syst.: Architectures, Model. Simul.*, 2015, pp. 149–156.
- [13] L. Sheng et al., "McPAT: An integrated power, area, and timing modeling framework for multicore and manycore architectures," in *Proc. Int. Symp. Microarch.*, 2009, pp. 469–480.
- [14] S. Mittal et al., "DESTINY: A comprehensive tool with 3D and multi-level cell memory modeling capability," *J. Low Power Electron. Appl.*, vol. 7, no. 3, 2017, Art. no. 23.
- [15] A. A. Khan, F. Hameed, R. Bläsing, S. Parkin, and J. Castrillon, "RT-Sim: A cycle-accurate simulator for racetrack memories," *IEEE Comput. Architecture Lett.*, vol. 18, no. 1, pp. 43–46, Jan.–Jun. 2019.