
LISA: Towards Learned DNA Sequence Search

Darryl Ho
MIT

darryl140@mit.edu

Jialin Ding
MIT

jialind@mit.edu

Sanchit Misra
Intel Labs

sanchit.misra@intel.com

Nesime Tatbul
Intel Labs and MIT

tatbul@csail.mit.edu

Vikram Nathan
MIT

vikramn@mit.edu

Vasimuddin Md
Intel Labs

vasimuddin.md@intel.com

Tim Kraska
MIT

kraska@mit.edu

Abstract

Next-generation sequencing (NGS) technologies have enabled affordable sequencing of billions of short DNA fragments at high throughput, paving the way for population-scale genomics. Genomics data analytics at this scale requires overcoming performance bottlenecks, such as searching for short DNA sequences over long reference sequences. In this paper, we introduce LISA (Learned Indexes for Sequence Analysis), a novel learning-based approach to DNA sequence search. As a first proof of concept, we focus on accelerating one of the most essential flavors of the problem, called *exact search*. LISA builds on and extends FM-index, which is the state-of-the-art technique widely deployed in genomics toolchains. Initial experiments with human genome datasets indicate that LISA achieves up to a factor of $4\times$ performance speedup against its traditional counterpart.

1 Introduction

Rapid advances in high-throughput next-generation sequencing (NGS) technologies have enabled affordable sequencing of billions of short DNA fragments (called “reads”) at an unprecedented scale. For example, the Illumina NovaSeq 6000 Sequencer can sequence 20 billion reads of length 150 each in less than 2 days, generating 6 Terabases of data at a low cost of about \$1000 per human genome [1]. Already today, a growing number of public and private sequencing centers with hundreds of NGS deployments pave the way for population-level genomics. However, realizing this vision in practice heavily relies on building scalable systems for high-speed genomics data analysis.

DNA sequence alignment plays a critical role in genome analysis. In its simplest form, an aligner tries to piece together the short reads by mapping each individual read to a long reference genome (e.g., the human genome consisting of 3 billion bases). The key operation that has been shown to constitute a significant performance bottleneck during this mapping process is the search for *exact* or *inexact* matches of read substrings over the given reference sequence [2–7]. In this paper, we focus on the *exact search* variant of this search problem.

The state-of-the-art technique to perform exact search is based on building an FM-index over the reference genome [8]. The key idea behind an FM-index is that, in the lexicographically sorted order of all suffixes of the reference sequence, all matches of a short DNA sequence (a.k.a., a “query”) will fall in a single region matching the prefixes of contiguously located suffixes. Over the years, many improvements have been made to make the FM-index more efficient, leading to several state-of-the-art

implementations that are highly cache- and processor-optimized [3–7, 9–15, 10, 15]. Hence, it becomes increasingly more challenging to further improve this critical step in the genomics pipeline to scale with increasing data growth.

In this paper, we propose a different approach to improving the sequence search performance: LISA (Learned Indexes for Sequence Analysis). The core idea behind LISA, which enables a new ML-enhanced algorithm for DNA search, is to speed up the process of finding the right region of suffixes by learning the distribution of suffixes in the reference. We do this in a way similar to how learned indexes capture value distributions through models learned from data [16].

When evaluated on an Intel® Core™ i9-9900K 3.6 GHz processor, despite being single-threaded and not yet fully optimized to the underlying hardware architecture, our current implementation achieves nearly $4\times$ speedup against a state-of-the-art single-threaded, CPU-optimized version of the FM-index based algorithm [15], for a workload of 50 million queries matched against the human genome. This early result shows that learned DNA sequence search is a promising idea ¹.

To the best of our knowledge, this is the first work exploring how ML-enhanced algorithms can improve the process of DNA sequence search, while providing identical semantic guarantees as the traditional algorithms. This work is a preliminary proof of concept that can be used as a building block towards fully-optimized learning-based tools for DNA sequence search. In summary, this paper makes the following contributions:

- enhancements to the FM-index that enable the application of learning-based search,
- a new search algorithm that applies the learning-based approach to the enhanced FM-index to find all exact matches,
- an experimental comparison of our approach against a highly-tuned baseline using realistic workloads on the human genome.

In the rest of this paper, we first provide some brief background on the traditional FM-index based exact search algorithm as well as the idea of learned index structures which inspired this work; then we present our new approach LISA, along with results from our experimental study.

2 Background

A DNA sequence is a string over the alphabet, $\Sigma = \{A, C, G, T\}$, representing the four bases. For the rest of this paper, we use the terms “sequence” and “string” interchangeably, as well as the terms “base” and “character”. Exact DNA sequence search is a key kernel in many genomics tools, including the widely-used sequence mapping tool Bowtie 2 [4]. Given a reference sequence R and a query sequence Q , the goal of exact sequence search is to find exact end-to-end matches of Q in R . Typically, $|R| \approx 10^9$ bases; e.g., the length of the human genome is around 3 billion bases. On the other hand, $|Q|$ is typically less than 200 bases; e.g., the default query length in Bowtie 2 is 21.

FM-index: Fig. 1 depicts the construction of the FM-index for an example reference sequence R . First, we append R with the character $\$ \notin \Sigma$ which is lexicographically smaller than all characters in Σ . Subsequently, we obtain all the rotations of R (Rotations(R)). The lexicographically sorted order of the rotations forms the BW-matrix. The BWT (B) is the last column of the BW-matrix. The original positions in R of the first bases of these rotations constitute the suffix array (S).

All the exact matches of a query can be found as prefixes of the rotations in the BW-matrix. Since the BW-matrix is lexicographically sorted, these matches are located in contiguous rows of the BW-matrix. Therefore, for a query, all the matches can be represented as a range of rows of the BW-matrix. This range is called the *SA interval* of the query. For example, in Fig. 1, the *SA interval* of query “AC” is $[1, 2]$. The values of the suffix array in the SA interval are 5 and 2. Indeed, the sequence “AC” is found at positions 5 and 2 in the reference sequence.

The FM-index is used to expedite search for the SA interval [18]. It consists of the suffix array S and the BWT B , as well as D and O data structures. $D(x)$ is the count of bases in $R[0, |R| - 1]$ that are lexicographically smaller than $x \in \Sigma$. $O(x, i)$ is the count of occurrences of base x in $B[0, i]$. Note that the BW-matrix is not stored.

¹Intel Xeon and Intel Xeon Phi are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries. Other names and brands may be claimed as the property of others. ©Intel Corporation

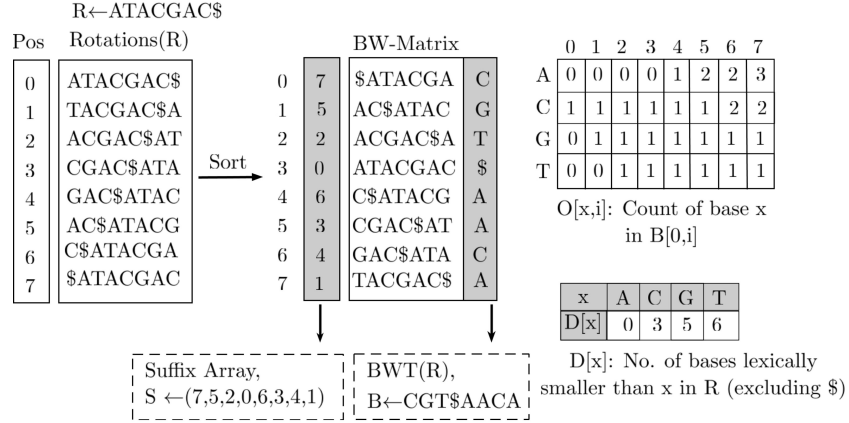


Figure 1: FM-index (S, B, D, O) and BW-matrix for sample reference sequence $R \leftarrow \text{ATACGAC}\$$. The lexicographical ordering is $\$ < A < C < G < T$ [17].

Using the FM-index, finding the SA interval is done using the *backward-search algorithm*. For a reference sequence of length n , the algorithm initializes the SA interval to $[0, n)$. The SA interval is updated over the course of the algorithm. Specifically, the algorithm processes the query sequence backwards, starting from the last character, prepending one character at a time and updating the SA interval after each prepended character in $O(1)$ time. To update the SA interval, the FM-index supports a function $f : (\text{char}, \text{int}) \rightarrow \text{int}$ that takes the prepended character c and an integer location i , and finds in $O(1)$ time the lower-bound location i' in the BW-matrix of string $\text{concat}(c, \text{BW-matrix}[i])$. (The lower bound of a string s is the first entry in the BW-matrix which does not compare less than s .) Given an SA interval $[l, u)$, after prepending the character c , the SA interval is updated to $[f(c, l), f(c, u))$. Since prepending each character takes $O(1)$ time, the overall algorithm takes $O(|Q|)$ time. For further information on the FM-index, see [8].

Learned Indexes: Recent work on learned index structures has introduced the idea that indexes are essentially models that map inputs to positions and, therefore, can be replaced by other types of models, such as machine learning models [16]. For example, a B-tree index maps a given key to the position of that key in a sorted array. Kraska et al. show that using knowledge of the distribution of keys can produce a learned model, called the recursive model index (RMI), that outperforms B-trees in query time and memory footprint.

Taking a similar perspective, the FM-index can be seen as a model that maps a given query sequence to the SA interval for that query sequence. Based on this insight, in LISA, we use knowledge of the distribution of subsequences within the reference sequence to create a learned index structure that enables faster exact search queries.

3 LISA

Backward-search algorithm performs exact search of a query Q using the FM-index by iterating through the query sequence in backwards order, one base at a time, thereby consuming $O(|Q|)$ time. The key idea of LISA is to iterate backwards through the query sequence in chunks of K bases at a time, so that exact search takes $O(|Q|/K)$ time. To do this, LISA requires two components: (1) a new data structure called the IP-BWT that enables processing K bases at a time, and (2) a method to efficiently search through the IP-BWT, for which we use an RMI. Similar to the *backward-search algorithm*, the goal of LISA is to output the SA interval of a query sequence.

IP-BWT: In the *backward-search algorithm*, in order to update the SA interval after prepending a character, FM-index supports a function $f : (\text{char}, \text{int}) \rightarrow \text{int}$ that takes the character c and an integer location i , and finds in $O(1)$ time the lower-bound location i' in the BW-matrix of string $\text{concat}(c, \text{BW-matrix}[i])$. For LISA, we need to support a similar function that takes in a length- K string s and a location i , and returns another location i' (i.e. $f_K : (\text{char}^K, \text{int}) \rightarrow \text{int}$) representing the lower-bound location in the BW-matrix of string $\text{concat}(s, \text{BW-matrix}[i])$.

Algorithm 1 Exact Search Algorithm using IP-BWT

Input: Q , a query string; $f_K : (\text{char}^K, \text{int}) \rightarrow \text{int}$, a function that finds the lower-bound location of $(\text{char}^K, \text{int})$ in the IP-BWT

Output: $[low, high)$, an SA interval.

```
1:  $low, high \leftarrow 0, n$ 
2: split  $Q$  into  $\lceil |Q|/K \rceil$  chunks, each of length  $K$ , with the final chunk possibly shorter than  $K$ 
3: for  $\mathcal{C}$  in reversed order of chunks do
4:   if  $|\mathcal{C}| < K$  then
5:     /* Special case for when final chunk has length less than  $K$  */
6:      $\mathcal{C}_{low} \leftarrow \mathcal{C} + "\$"$  +  $"A" \times (K - |\mathcal{C}| - 1)$ 
7:      $\mathcal{C}_{high} \leftarrow \mathcal{C} + "T" \times (K - |\mathcal{C}|)$ 
8:      $low, high \leftarrow f_K((\mathcal{C}_{low}, low)), f_K((\mathcal{C}_{high}, high))$ 
9:   else
10:     $low, high \leftarrow f_K((\mathcal{C}, low)), f_K((\mathcal{C}, high))$ 
11:   end if
12: end for
13: return  $[low, high)$ 
```

For this purpose, we introduce the Index-Paired BWT (IP-BWT) array. Each entry of IP-BWT consists of a $(\text{char}^K, \text{int})$ pair. The first part is the first K characters of the corresponding BW-matrix row. The second part is the BW-matrix location of the string with the first K and the last $n - K$ characters swapped. Our desired function $f_K : (\text{char}^K, \text{int}) \rightarrow \text{int}$ is now equivalent to finding the lower-bound location of the input $(\text{char}^K, \text{int})$ pair in the IP-BWT array. We are free to choose any implementation for how to find that lower-bound location; for example, since the IP-BWT is sorted, we could do a binary search over the entries of the IP-BWT. Fig. 2 shows how to create an IP-BWT with $K = 3$.

Alg. 1 shows the exact search algorithm using IP-BWT. For example, using the reference sequence and IP-BWT from Fig. 2, let the query sequence be ATTA. We split this into two chunks: ATT and A (line 2). We first use the RMI to find the lower bound locations of $(A\$A, 0)$ and (ATT, n) , which are 1 and 5, respectively. We then use the RMI to find the lower bound locations of $(ATT, 1)$ and $(ATT, 5)$, which are 3 and 5. Our algorithm gives the interval $[3, 5)$. We can confirm that ATTA can be found in position 3 and 4 of the BW-matrix.

Faster Chunk Processing using RMI: Using the IP-BWT, we are able to process the query sequence in chunks of K bases at a time. However, when processing each chunk, we must evaluate the function f_K , which involves a binary search over the IP-BWT. This takes $O(\log n)$ time, where n is the number of entries in the IP-BWT, which is equivalent to the length of the reference sequence. Therefore, the overall runtime of exact search using IP-BWT and evaluating f_K with binary search is $O(|Q| \log n / K)$. For large reference sequences, this might be slower than *backward-search* using the FM-index.

In order to avoid paying the cost of a binary search for each evaluation of f_K , we use a learned approach to support $O(1)$ evaluation of f_K . In particular, f_K is a model that maps input keys $((\text{char}^K, \text{int})$ pairs) to their positions in the sorted IP-BWT. We model f_K using the RMI, which is a hierarchy of models that is quick to evaluate [16]; the RMI conceptually resembles a hierarchical mixture of experts [19]. Fig. 3 shows an example of using a 3-layer RMI to evaluate f_K in three steps: (1) since the RMI only accepts numbers as inputs, we first convert the input $(\text{char}^K, \text{int})$ into a number. Since the alphabet Σ only has 4 characters, any character can be represented in 2 bits. Therefore, we convert char^K into a number with $2K$ bits by concatenating the bits of the individual characters together. We then append the bits of the int . Note that we have a special case for handling the sentinel character $\$$ while maintaining this 2-bit encoding. (2) We give the encoded input to the RMI and traverse down the layers of the RMI to a leaf model. The leaf model predicts the position in the IP-BWT where it expects to find the input pair. (3) If the predicted position does not contain the input pair, we use linear search over the IP-BWT starting from the predicted position to find the actual position of the pair. Note that this

BW-Matrix		IP-BWT
00	SCATTATTAGGA	00 SCA, 11
01	A\$CATTATTAGG	01 A\$C, 04
02	AGGA\$CATTATT	02 AGG, 01
03	ATTAGG\$CATT	03 ATT, 02
04	ATTATTAGGAS\$C	04 ATT, 03
05	CATTATTAGGAS	05 CAT, 09
06	GAS\$CATTATTAG	06 GA\$, 05
07	GGA\$CATTATTA	07 GGA, 00
08	TAGGAS\$CATTAT	08 TAG, 06
09	TATTAGG\$CAT	09 TAT, 08
10	TTAGGAS\$CATT	10 TTA, 07
11	TTATTAGGAS\$CA	11 TTA, 10

Figure 2: Conversion of a BW-Matrix to IP-BWT on reference sequence CATTATTAGGA, where $K = 3$.

learning-based approach to modeling f_K guarantees correctness; LISA will produce exactly the same results as using backward search with FM-index.

Unlike the RMI proposed in [16], which constructs the model hierarchy top-down according to the user-selected number of models at each layer, we construct the RMI bottom-up according to desired bounds on the average error between the predicted position and the actual position. We use these desired bounds to determine the number of models at each layer. Bounds on the average prediction error are useful for limiting the time spent on step 3 of the RMI evaluation workflow described above, because they directly reflect the average number of iterations of linear search. Given a desired bound α on the average error, we begin building the RMI at the leaf layer by partitioning the IP-BWT entries into contiguous blocks that can each be modeled with average error no more than α . We find this partitioning by starting with one partition that comprises of the entire IP-BWT, then recursively splitting equally in two until each partition achieves the α bound. A leaf model is built on the entries of each partition. The smallest entry in each partition is used to repeat this procedure in order to find the partitions in the layer above the leaf layer, and so forth until we have one model at the root layer. Since the non-leaf models are allowed to have prediction error, we may need to perform linear search at each layer of the RMI, instead of only the leaf layer. Note that we can set different values of α for each layer of the RMI. Also, note that at the root layer, there is only one partition; therefore, we do not set an α bound for the root model. Since we have no guarantee on average error for the root model, the root model uses exponential search instead of linear search.

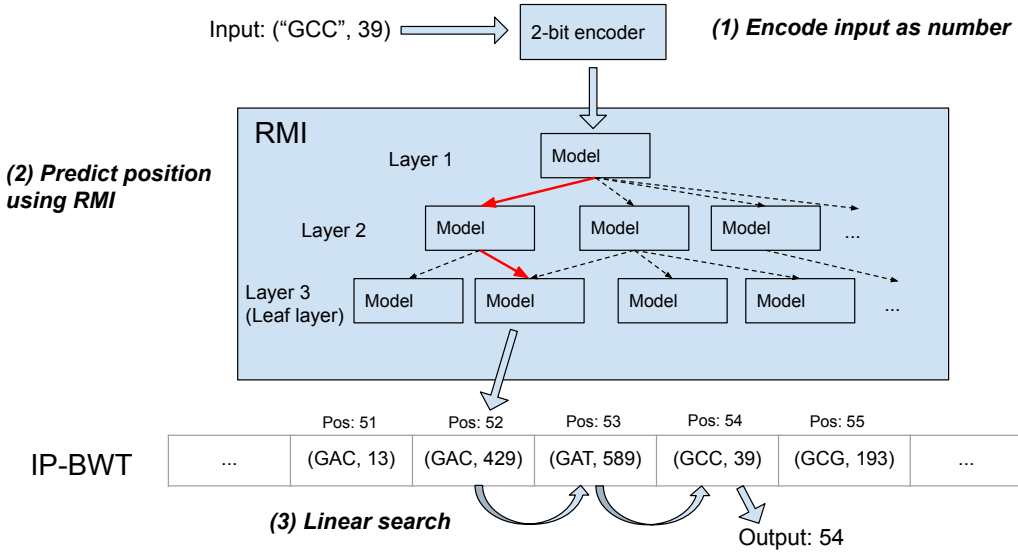


Figure 3: Using a 3-layer RMI to evaluate f for an example input ($GCC, 39$). The RMI predicts position 52, and linear search finds the correct position. The solid red lines show the traversal path down the RMI.

Optimization for Batched Queries: So far, we have presented our method of processing a single query. However, in practice, queries arrive in large batches. For large query batches, we can adopt an additional optimization. Because the final layer of the RMI (i.e., RMI leaves), which are responsible for predicting locations, are arranged in sorted order, if the queries also arrive in sorted order, we can use the *double-pointer technique* (Alg. 2) to find for each query its corresponding RMI leaf. This resembles the merge step of merge-sort. When the number of queries is Ω (number of RMI leaves), the amortized time for each query to find its RMI leaf is $O(1)$. This avoids the cost of traversing down the RMI for each query, and also improves cache locality.

Note that this optimization for batches cannot be easily applied to the FM-index. This optimization requires that inputs to f_K are sorted; since IP-BWT processes in chunks of K characters, we only need to sort queries $O(|Q|/K)$ times. However, since backward search with the FM-index prepends only one character at a time, it would need to sort the inputs to f after each prepended character. Sorting inputs $O(|Q|)$ times would impose a significant performance bottleneck, which makes this optimization impractical when using FM-index.

Algorithm 2 The RMI-based exact search algorithm for large query batch

Input: $qs[]$, a list of query strings encoded as numbers.

Output: a list of SA intervals.

```
1: sort  $qs$ 
2:  $leafPtr \leftarrow \&rmi.leaves[0]$ 
3: for  $q$  in  $qs$  do
4:   while  $q \geq leafPtr.ipbwt\_range\_upper\_bound$  do
5:     increment  $leafPtr$ 
6:   end while
7:    $prediction \leftarrow leafPtr.predict(q)$ 
8:   perform linear search around  $prediction$  to find the SA interval of  $q$ 
9: end for
```

Discussion: LISA’s speed advantage over FM-index comes from two components: (1) the IP-BWT, which allows LISA to process K -character chunks of the query at a time, whereas FM-index processes one character at a time, and (2) using an RMI to process each K -character chunk in $O(1)$ time, whereas a naive binary search would take $O(\log n)$ time per chunk. Therefore, an exact search query using FM-index takes $O(|Q|)$ time, whereas LISA takes $O(|Q|/K)$ time.

To discuss memory consumption, we measure in bytes in terms of n , the length of the reference sequence. The reference sequence itself therefore takes $0.25n$ space. In the state-of-the-art implementation of *backward-search algorithm*, the space of suffix array is $4n$, a compressed structure that combines BWT and O is $2n$, and D is negligible, for a total space consumption of $6n$. For LISA, the space of the suffix array is $4n$, IP-BWT is $(0.25K + 4)n$, and the RMI is usually around $0.5n$, for a total space consumption of $8.5n + 0.25Kn$. For example, in Bowtie 2 the default value of $|Q|$ is 21, so using an IP-BWT with $K = 21$, LISA takes around $13.75n$ space, which is larger than the FM-index. However, slightly larger space consumption is usually not a concern in practice, and if necessary LISA can use smaller K or compress the IP-BWT. Though the space usage of LISA would increase with larger values of K , we find through experiments that LISA maintains good performance using an IP-BWT with $K = 21$, even for large query lengths. Therefore, the space of LISA does not need to grow beyond around $13.75n$.

Could we replace the RMI with some other index structure that can evaluate f_K over the IP-BWT even faster? [16] showed that RMIs perform better than binary search and B trees. A lookup table (implemented as an array) that stores the output of f_K for every possible pair $(char^K, int)$ would also allow $O(1)$ evaluations but would far exceed memory capacity, even if we use an IP-BWT with very small K . However, it is possible to combine a downsampled version of the lookup table with binary search, which we discuss in the evaluation. Another idea is to use a hash table to map all n existing pairs $(char^K, int)$ to their positions; however, this fails because we almost always need to evaluate f_K on pairs that do not exist. For example, in the example attached to Fig. 2, we find the lower bounds of four pairs— $(A\$A, 0)$, (ATT, n) , $(ATT, 1)$, and $(ATT, 5)$ —none of which exist in the IP-BWT. Also, note that since the entries of the IP-BWT must be sorted in order to maintain one contiguous SA interval, the IP-BWT itself cannot be replaced with a hash table to enable faster searches.

4 Evaluation

We present preliminary results² for LISA. Experiments use SIMD, running a single-thread implementation on an Ubuntu system with Intel® Core™ i9-9900K 3.6GHz processor and 64GB RAM. As our baseline, we use a highly CPU-optimized implementation of backward search algorithm that is

²Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more information go to www.intel.com/benchmarks.

Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

significantly faster than its alternatives [15]. We also compare with doing backwards search using IP-BWT and binary search (i.e., without the RMI).

Workload Scenario: We evaluate LISA on a real-world scenario: we use whole human genome as the reference sequence and for query sets, we use large sets of short query sequences of various lengths as would be found in several prominent sequence aligners. We train LISA using an IP-BWT with $K = 21$. The RMI has three layers, only uses linear regression models, and we construct the RMI using $\alpha = 14$ for the second layer models and $\alpha = 6$ for the leaf layer models. α is tuned once for optimal performance on the reference sequence; the RMI is not re-trained or re-tuned for each experiment. We evaluate on four different query sequence lengths: (1) $|Q| = 21$, so that LISA processes exactly one chunk, (2) $|Q| = 42$, so that LISA processes multiple whole chunks, (3) $|Q| = 32$, so that LISA processes a chunk shorter than K , and (4) a very long query, $|Q| = 200$. For each of these query sequence lengths, we generate a batch of 50 million query sequences randomly from the human genome to be our query set, and run them together with LISA, using the batched-query optimization from §3.

Tab. 1 shows that LISA is $2.73\times$ to $3.97\times$ faster than the optimized FM-index baseline on these query lengths. Query lengths that are a perfect multiple of K (i.e., 21 and 42) perform the best. The query length of 32 has slightly lower relative performance, as LISA must still process 2 chunks per query, as if the query has length 42. The very long query sequence has lower relative performance, as LISA’s performance does not scale linearly: for longer queries, the time spent on sorting grows super-linearly. LISA achieves around a $2\times$ performance boost from using the RMI instead of binary search.

	$ Q =21$	$ Q =32$	$ Q =42$	$ Q =200$
Optimized FM-index	1509	2414	3284	15411
IP-BWT with binary search	785 (1.92 \times)	1424 (1.70 \times)	1779 (1.85 \times)	10414 (1.48 \times)
LISA	383 (3.94 \times)	762 (3.17 \times)	827 (3.97 \times)	5646 (2.73 \times)

Table 1: Average query time (measured in ticks per query) of exact search on 50 million queries, while varying query length. Relative speedup to FM-index is shown in parentheses.

Batch Size: In order to measure the effect of batch size on query times, we fix the query sequence length to $|Q| = 21$ and vary the batch size from 1 to 500 million. Fig. 4 shows that as batch size increases, LISA’s time per query decreases. For batches larger than 10 thousand, LISA starts to outperform the FM-index.

LISA’s performance is poor for small batch sizes, where the double pointer technique actually hurts performance, as the merge step will skip over many leaves. However, in real-world scenarios, it is rare to have such small batch sizes where this would matter. Next-generation sequencing technologies produce billions of DNA fragments, and downstream applications such as DNA sequence alignment need to process all or most of the produced fragments at once. Therefore, the typical use case for LISA is to process large batches of queries. If high performance on small batches is absolutely necessary, it is possible to achieve comparable/better performance than FM-index for small batch sizes by not using the double-pointer technique.

Alternative to RMI: To evaluate the effectiveness of RMI against alternative index structures, we compared LISA to a version which evaluates f_K using a combination of a lookup-table and binary search, instead of using the RMI. If we treat a $(char^K, int)$ pair as a $(2K + 32)$ -bit number, we maintain a “downsampled” lookup-table (implemented as an array) with 2^p entries, where $p < 2K + 32$. Let $q = 2K + 32 - p$. The i -th entry of the lookup-table holds the lower-bound location in the IP-BWT of the pair represented by the number $i \cdot q$. Essentially, our downsampled lookup-table holds every q -th entry of a “full” lookup-table that contains every possible $(char^K, int)$ pair. To evaluate f_K on an input represented as a $2K + 32$ -bit number, we do a lookup in the lookup-table using the first p bits to

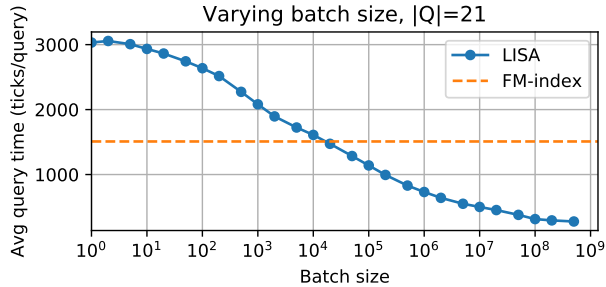


Figure 4: LISA’s performance advantage increases with batch size.

obtain lower and upper bounds on the location of the input pair in the IP-BWT, then do a binary search with the remaining q bits on the IP-BWT starting from those bounds to find the location of the input. This is our best effort at an alternative index structure that is most competitive with RMI.

For $K = 21$, $|Q| = 21$, and batch size of 50 million, replacing the RMI with a downsampled lookup-table of the same memory size results in 40% slower query times. More performant lookup-tables use significantly more memory: the lookup-table itself has space 2^{p+2} bytes, so even with only $p = 34$, the lookup table takes 69GB, which exceeds the memory of our machine. Therefore, the lookup-table is impractical and we do not include it in this paper.

5 Conclusions

In this work, we introduced a preliminary version of our learned indexing approach for DNA sequence search, LISA, which produces promising initial results for the exact search problem when tested on workloads of realistic queries against the human genome. In particular, LISA achieves up to nearly $4\times$ faster query times than an extensively optimized version of the FM-index based approach, which has been the common practice in sequence search. We believe that the core ideas behind LISA can be extended to other types of DNA sequence search problems. In particular, we are currently working on using learned indexes to speed up searching for super maximal exact matches (SMEMs) for a query in the reference. For any position in query, an SMEM is the longest substring of the query through that position that has an exact match in the reference [17, 20, 21]. We are also working with the Broad Institute to integrate LISA into applications that are widely used by the genomics community.

Acknowledgments

We thank Tony Peng, Ashwath Thirumalai, and Elizabeth Wei for their contributions to the original design of LISA; and Pradeep Dubey and Heng Li for their valuable feedback. This research has been funded in part by affiliate members and supporters of DSAIL (Data Systems and AI Lab) at MIT – Google, Intel, and Microsoft.

References

- [1] Illumina, Inc. Novaseq 6000 sequencing system. <https://www.illumina.com/content/dam/illumina-marketing/documents/products/datasheets/novaseq-6000-system-specification-sheet-770-2016-025.pdf>, 2019. Accessed: September 2019.
- [2] Vasimuddin Md, Sanchit Misra, and Srinivas Aluru. Identification of Significant Computational Building Blocks through Comprehensive Investigation of NGS Secondary Analysis Methods. *bioRxiv*, April 2018. URL <https://www.biorxiv.org/content/early/2018/07/25/301903>.
- [3] Ben Langmead, Cole Trapnell, Mihai Pop, and Steven L Salzberg. Ultrafast and Memory-efficient Alignment of Short DNA Sequences to the Human Genome. *Genome Biology*, 10(3):1, 2009.
- [4] Ben Langmead and Steven L Salzberg. Fast Gapped-read Alignment with Bowtie 2. *Nature Methods*, 9(4):357–359, 2012.
- [5] Ruiqiang Li, Chang Yu, Yingrui Li, Tak-Wah Lam, Siu-Ming Yiu, Karsten Kristiansen, and Jun Wang. SOAP2: An Improved Ultrafast Tool for Short Read Alignment. *Bioinformatics*, 25(15):1966–1967, 2009.
- [6] Ruibang Luo, Thomas Wong, Jianqiao Zhu, Chi-Man Liu, Xiaoqian Zhu, Edward Wu, Lap-Kei Lee, Haoxiang Lin, Wenjuan Zhu, David W Cheung, et al. SOAP3-dp: Fast, Accurate, and Sensitive GPU-Based Short Read Aligner. *PLOS ONE*, 8(5):e65632, 2013.
- [7] Heng Li and Richard Durbin. Fast and Accurate Short Read Alignment with Burrows–Wheeler Transform. *Bioinformatics*, 25(14):1754–1760, 2009.
- [8] Paolo Ferragina and Giovanni Manzini. Opportunistic Data Structures with Applications. In *IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 390–398, 2000.

- [9] Alejandro Chacón, Juan Carlos Moure, Antonio Espinosa, and Porfidio Hernández. n-step FM-Index for Faster Pattern Matching. *Procedia Computer Science*, 18:70–79, 2013.
- [10] Alejandro Chacón, Santiago Marco-Sola, Antonio Espinosa, Paolo Ribeca, and Juan Carlos Moure. Boosting the FM-Index on the GPU: Effective Techniques to Mitigate Random Memory Access. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 12(5):1048–1059, September 2015.
- [11] J. Pantaleoni and N. Subtil. NVBIO: A Library of Reusable Components Designed by NVIDIA Corporation to Accelerate Bioinformatics Applications using CUDA. <http://nvlabs.github.io/nvbio/>. Accessed: November 2017.
- [12] Jing Zhang, Heshan Lin, Pavan Balaji, and Wu-chun Feng. Optimizing Burrows-Wheeler Transform-based Sequence Alignment on Multi-core Architectures. In *IEEE/ACM International Symposium on Cluster, Cloud, and Grid Computing (CCGrid)*, pages 377–384, 2013.
- [13] E. Fernandez, W. Najjar, and S. Lonardi. String Matching in Hardware Using the FM-Index. In *IEEE Annual International Symposium on Field-Programmable Custom Computing Machines*, pages 218–225, May 2011.
- [14] Szymon Grabowski, Marcin Raniszewski, and Sebastian Deorowicz. FM-index for Dummies. In *International Conference on Beyond Databases, Architectures, and Structures (BDAS)*, pages 189–201, 2017.
- [15] Sanchit Misra, Tony C Pan, Kanak Mahadik, George Powley, Priya N. Vaidya, Md Vasimuddin, and Srinivas Aluru. Performance Extraction and Suitability Analysis of Multi- and Many-core Architectures for Next Generation Sequencing Secondary Analysis. In *International Conference on Parallel Architectures and Compilation Techniques (PACT)*, pages 3:1–3:14, 2018.
- [16] Tim Kraska, Alex Beutel, Ed H. Chi, Jeffrey Dean, and Neoklis Polyzotis. The Case for Learned Index Structures. In *ACM International Conference on Management of Data (SIGMOD)*, pages 489–504, June 2018.
- [17] Vasimuddin Md, Sanchit Misra, Heng Li, and Srinivas Aluru. Efficient Architecture-Aware Acceleration of BWA-MEM for Multicore Systems. In *IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2019.
- [18] Paolo Ferragina and Giovanni Manzini. An Experimental Study of an Opportunistic Index. In *Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 269–278, 2001.
- [19] Michael I. Jordan and Robert A. Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural Comput.*, 6(2):181–214, March 1994. ISSN 0899-7667. doi: 10.1162/neco.1994.6.2.181. URL <http://dx.doi.org/10.1162/neco.1994.6.2.181>.
- [20] Heng Li. Exploring Single-Sample SNP and INDEL Calling with Whole-Genome De Novo Assembly. *Bioinformatics*, 28(14):1838–1844, 2012.
- [21] Heng Li. Aligning Sequence Reads, Clone sequences, and Assembly Contigs with BWA-MEM. *arXiv preprint arXiv:1303.3997*, 2013.