

Sequence analysis

FroM Superstring to Indexing: a space-efficient index for unconstrained k -mer sets using the Masked Burrows-Wheeler Transform (MBWT)

Ondřej Sladký^{1,2,*}, Pavel Veselý^{2,*}, Karel Brinda^{3,*}

¹ETH Zurich, 8092 Zurich, Switzerland

²Computer Science Institute of Charles University, 118 00 Prague, Czechia

³INRIA, IRISA, Univ. Rennes, 35042 Rennes, France

*Corresponding authors. Ondřej Sladký, ETH Zurich, Zurich, Switzerland. E-mail: ondra.sladky@gmail.com; Pavel Veselý, Computer Science Institute of Charles University, Malostranské nám. 25, 118 00 Prague, Czechia. E-mail: vesely@iuuk.mff.cuni.cz; Karel Brinda, INRIA, IRISA, Univ. Rennes, Campus de Beaulieu, 35042 Rennes, France. E-mail: karel.brinda@inria.fr.

Associate Editor: Thomas Lengauer

Abstract

Motivation: The growing volumes and heterogeneity of genomic data call for scalable and versatile k -mer-set indexes. However, state-of-the-art indexes such as SBWT and SShash depend on long non-branching paths in de Bruijn graphs, which limits their efficiency for small k , sampled data, or high-diversity settings.

Results: We introduce FMSI, a superstring-based index for arbitrary k -mer sets that supports efficient membership and compressed dictionary queries with strong theoretical guarantees. FMSI builds on recent advances in k -mer superstrings and uses the Masked Burrows-Wheeler Transform, a novel extension of the classical Burrows-Wheeler Transform that incorporates position masking. Across a range of k values and dataset types—including genomic, pangenomic, and metagenomic—FMSI consistently achieves superior query space efficiency, using up to 2–3× less memory than state-of-the-art methods, while maintaining competitive query times. Only a space-optimized version of SBWT can match the FMSI's footprint in some cases, but then FMSI is 2–3× faster. Our results establish superstring-based indexing as a robust, scalable, and versatile framework for arbitrary k -mer sets across diverse bioinformatics applications.

Availability and implementation: FMSI is developed in C++ and released under the MIT license, with source code provided at <https://github.com/OndrejSladky/fmsi> and an installable package available through Bioconda. The datasets used in the experiments are deposited at Zenodo (<https://doi.org/10.5281/zenodo.14722244>).

1 Introduction

The exponential growth of DNA sequencing data calls for efficient approaches for their compression and search (Loh *et al.* 2012, Stephens *et al.* 2015). To address the growing complexity of data, modern bioinformatics methods increasingly rely on k -mer tokenization for data storage and analysis. This strategy enables a unified representation of various types of genomic data, including sequencing reads, transcripts, assemblies, or species pangenomes. Notable applications of k -mer-based methods include large-scale data search (Bradley *et al.* 2019, Bingmann *et al.* 2019, Brinda *et al.* 2025, Karasikov *et al.* 2025), metagenomic classification (Wood and Salzberg 2014, Brinda *et al.* 2017), rapid diagnostics (Bradley *et al.* 2015, Brinda *et al.* 2020), and transcript abundance quantification (Bray *et al.* 2016, Patro *et al.* 2017). An essential low-level functionality is the storage, and membership and dictionary queries on single k -mer sets (Marchet *et al.* 2021a).

A key challenge is to design k -mer set data structures that can scale with the set size and simultaneously adapt to the structural characteristics of different sets. Comprehensive

databases now reach hundreds of billions of k -mers (Karasikov *et al.* 2025) and continue to grow rapidly. At the same time, their structural complexity is increasing due to denser sampling of genetic polymorphisms and the accumulation of technological artifacts. Although advances in k -mer sampling (Wood and Salzberg 2014) and sketching (Ondov *et al.* 2016) techniques can alleviate size issues in some applications, the resulting k -mer sets are then less amenable to standard compression techniques.

The efficiency of k -mer set data structures is tightly connected to the compressive properties of their support representations. Real-world k -mer sets are largely non-random, and a wisely chosen indexable representation allows one to exploit redundancies, going beyond the traditional information-theoretic lower bounds (Conway and Bromage 2011). Of particular benefit has been the so-called *spectrum-like property*, stating that k -mer sets are typically generated by a small number of long strings (Chikhi *et al.* 2022). As these strings correspond to paths in the associated de Bruijn graphs (dBG), indexing dBG paths enables efficient indexes for the original k -mer sets.

Received: September 26, 2025; Accepted: October 30, 2025

© The Author(s) 2025. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

State-of-the-art methods for indexing dBG paths typically use one of two approaches: either they proceed implicitly by indexing successor k -mers in the graph (Bowe et al. 2012, Alanko et al. 2023), or they enumerate a path cover of the dBG explicitly and then index it using either hash-based (Marchet et al. 2021b, Pibiri, 2022) or full-text (Chikhi et al. 2015, Rahman and Medvedev 2021, Brinda et al. 2021) techniques. We refer to such explicit path-cover representations as *Spectrum-Preserving String Sets* (SPSS), and refer to the literature for more details (Brinda et al. 2021, Rahman and Medvedev 2021, Schmidt and Alanko 2023, Schmidt et al. 2023). Efficient SPSS computation is now supported by a wide range of software tools (Chikhi et al. 2016, Minkin et al. 2017, Holley and Melsted 2020, Khan and Patro 2021, Brinda et al. 2021, Rahman and Medvedev 2021, Khan et al. 2022, Schmidt and Alanko 2023, Cracco and Tomescu 2023, Schmidt et al. 2023, Khan et al. 2025).

However, the compression power of dBG-based methods is fundamentally limited by the existence of $(k-1)$ -long overlaps between the indexed k -mers. Even when a fraction of k -mers do not share $(k-1)$ -long overlaps with other k -mers, long paths in the graph cannot be formed, and the corresponding k -mer pairs cannot be jointly compressed, impacting all these indexing methods. The resulting fragmentation causes time and space inefficiencies at all levels, especially in applications involving complex metagenomic or pangenomic k -mer sets, or when the sets are subsampled.

Masked superstrings (Sladký et al. 2023) eliminate the dependency on $(k-1)$ -long overlaps. A masked superstring consists of a superstring of the original k -mer set, along with an associated mask that encodes the introduced “false positive” k -mers. Masked superstrings can represent arbitrary k -mer sets in text form while leveraging overlaps of all possible lengths. Conceptually, this corresponds to replacing dBGs with overlap graphs to construct the path cover. Masked superstrings thus unify all SPSS-based representations while enabling further data reduction by using shorter k -mer overlaps as well. Importantly, approximately shortest masked superstrings can be efficiently computed using greedy approaches implemented in KmerCamel (Sladký et al. 2023) (<https://github.com/ondrejsladky/kmercamel>), starting from either the original k -mer set or a precomputed SPSS.

However, masked superstrings are not indexable using existing methods. While the superstring itself can be indexed with full-text indexing, this yields an inexact data structure with false positives. A straightforward inclusion of masks would require recovering the positions of k -mers in the support superstring, which in turn necessitates auxiliary data structures, such as the subsampled suffix array, thereby increasing both the space and time requirements of queries (Salikhov et al. 2018, Brinda et al. 2021). Furthermore, efficient query processing in real-world applications often relies on streaming consecutive k -mers; yet, masked superstrings are incompatible with existing BWT-based streaming techniques (Salikhov 2017, Salikhov et al. 2018, Alanko et al. 2023) due to complications arising from k -mer canonicity when minimizing index size. Finally, in downstream applications requiring index associativity, to retain high space efficiency even for subsampled k -mer sets, the index must support the full functionality of compressed k -mer dictionaries (Pibiri 2022).

Here, we introduce FMSI, a scalable data structure and software for space-efficient membership and dictionary

queries over k -mer sets with arbitrary structure, with support for accelerated streamed queries. FMSI combines compact k -mer set representation using approximately shortest masked superstrings (Sladký et al. 2023) with indexing via the Burrows-Wheeler Transform (BWT) (Burrows and Wheeler 1994), extended to accommodate for masks. The resulting index provides strong theoretical guarantees on space for arbitrary k -mer sets, including those resulting from sampling and sketching. We implement this in a program called FMSI (<https://github.com/ondrejsladky/fmsi>) and demonstrate its superior space efficiency across a broad range of datasets and use cases compared to current state-of-the-art k -mer-set indexing methods.

2 Methods

2.1 Problem formulation

Let us consider the nucleotide alphabet $\{A, C, G, T\}$. For a fixed $k > 0$, k -mers are strings of length k over this alphabet. We assume the *bidirectional model* (formalized, e.g. in Chikhi et al. (2016), online suppl. and Schmidt and Alanko (2023)), where a k -mer and its reverse complement (RC) are considered equal, unless the *unidirectional model* is explicitly indicated. We are interested in providing the following indexing functionality for a k -mer set K .

- (P1) **Membership.** Single $O(k)$ -time operation:
 - ▷ MEMBER(Q): Determine if a given k -mer Q is in K .
- (P2) **Dictionary.** Two $O(k)$ -time operations enabling a bijective mapping between k -mers $Q \in K$ and a discrete interval of hash values $H = \{0, \dots, |K| - 1\}$:
 - ▷ LOOKUP(Q): Return the hash $h \in H$ of Q , or -1 if $Q \notin K$.
 - ▷ ACCESS(h): Return the k -mer Q corresponding to $h \in H$.
- (P3) **Streaming.** MEMBER/LOOKUP in $O(1)$ time for a k -mer Q following a query for k -mer Q' when the $(k-1)$ -long prefix of Q is a suffix of Q' and both $Q, Q' \in K$.

The dictionary index (P2) corresponds to evaluating a *minimal perfect hash function* with rejection and its inverse. (P2) is more general than (P1), but (P1) can be implemented more efficiently. (P3) corresponds to extending the core data structure for isolated queries by a component for rapid FWD or BWD operations (as defined in (Chikhi et al. 2022)), thus trading a small additional space for improving the time efficiency of positive streamed queries.

Our goal is to develop a data structure that simultaneously addresses four challenges. First, it should minimize space requirements while maintaining competitive query speed. Second, it must make no assumptions about the structure in the data, such as about the presence of many $(k-1)$ -long overlaps. Third, when structural regularities such as overlaps are available, the method should be able to exploit them. Fourth, the data structure should be available in four optimized variants, each targeting a specific functionality: a membership index (P1), a compressed dictionary (P2), and two extended versions with streaming support (P1 + P3 and P2 + P3, respectively).

2.2 Two core ingredients of FMSI

Our approach is based on indexing approximately shortest superstrings of the input k -mers, using a technique based on the BWT (Burrows and Wheeler 1994) (Fig. 1). Depending on the desired functionality (Fig. 1a), we compute a

particular variant of the masked superstring representation (Sladký *et al.* 2023), obtain the BWT image of superstring together with an appropriate mask transformation, and complete the data structure with supporting rank and select data structures into an exact index termed FMSI (Fig. 1b). FMSI shares similarities but is different from existing BWT-based approaches such as dbgFM (Chikhi *et al.* 2015, Rahman and Medvedev 2021), BWA (Li 2012, Brinda *et al.* 2021)/ProPhex (Salikhov 2017, Salikhov *et al.* 2018), and SBWT (Alanko *et al.* 2023) (Fig. 1c; Table 1).

2.2.1 Approximately shortest masked superstrings of k -mers

A *masked superstring* representation (Sladký *et al.* 2023) of a given k -mer set K consists of two parts: a superstring S of the k -mers, i.e. a single string containing all the k -mers in K as substrings, and a binary mask M determining which of the k -mers appearing in S are represented. A k -mer is considered represented by a masked superstring (S, M) if at least one of its occurrences in S has the corresponding mask symbol set to 1, and the set of represented k -mers should be equal to K . Since it is sufficient to mask

each k -mer as present only once, this makes the mask nonunique, allowing for additional mask optimization. Masked superstrings unify all SPSS representations and provide a more general framework with better compression properties (Sladký *et al.* 2023).

Here, we consider uniquely masked superstrings obtained via greedily approximated shortest superstrings (Sladký *et al.* 2023) and, based on the target functionality, masked to have either the maximum or the minimum number of ones, termed max-one or min-one masked superstrings, respectively (Steps 1, 2 in Fig. 1b). Given input k -mer data, such a superstring can be efficiently approximated by the global greedy algorithm as implemented in KmerCamel (Sladký *et al.* 2023), with possible acceleration via parallelized SPSS precomputation. The output min-one masked superstring can further be converted to max-one using a two-pass algorithm (Sladký *et al.* 2023).

2.2.2 Masked Burrows-Wheeler Transform

We index the resulting masked superstring using the *Masked Burrows Wheeler Transform*, an extension of the traditional

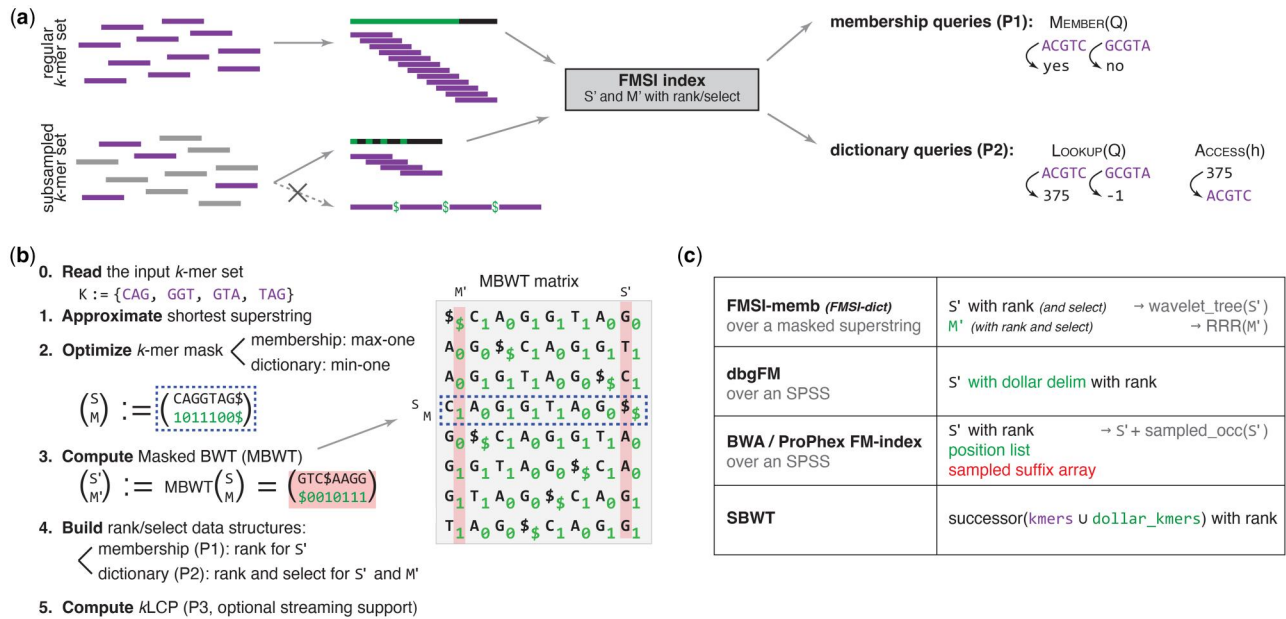


Figure 1. Overview of the FMSI method. Concepts related to input k -mers and masks are shown in violet and green, respectively. (a) Core FMSI functionality. FMSI indexes regular or subsampled k -mer sets and supports both membership or dictionary queries. Unlike other approaches, its representation remains compact even without $(k - 1)$ -long overlaps. (b) FMSI index construction. A greedy global algorithm (Sladký *et al.* 2023) approximates the shortest superstring of the input k -mer set; the resulting mask is optimized for the target functionality (P1 or P2). The masked superstring is transformed via the Masked BWT (Definition 1 and Algorithm 1) and equipped with rank/select data structures. Optionally, the k LCP array (Salikhov 2017) can be computed to accelerate streamed queries. (c) Comparison with other BWT-based indexes. dbgFM (Chikhi *et al.* 2015, Rahman and Medvedev 2021) indexes $\$$ -separated SPSS using an FM index without a sampled inverse suffix array. BWA/ProPhex (Li 2012, Salikhov 2017, Salikhov *et al.* 2018, Brinda *et al.* 2021) index SPSS concatenations, avoiding false matches at simplitig borders by checking simplitigs' start positions. SBWT (Alanko *et al.* 2023) indexes dBG paths via k -mer successors. Only FMSI removes redundancy among k -mers with overlaps shorter than $k - 1$.

Table 1. Functional comparison of FMSI with state-of-the-art single-set k -mer indexes.

	P1: Membership	P2: Dictionary	P3: Streaming	Scalable	Dynamic
BWA	✓	✗	ProPhex (Salikhov 2017) (Salikhov <i>et al.</i> 2018)	✓	✗
dbgFM	✓	ABySS (Simpson <i>et al.</i> 2009)	✗	✗	✗
SSHash	✓	✓	✓	✓	✗
CBL	✓	✗	✓	✓	✓
SBWT	✓	✗	✓	✓	✗
FMSI	✓	✓	✓	✓	f -MS (Sladký <i>et al.</i> 2025)

BWA Fastmap (Li 2012) on simplitigs (Brinda *et al.* 2021), dbgFM (Chikhi *et al.* 2015, Rahman and Medvedev 2021), SSHash (Pibiri 2022), CBL (Martayan *et al.* 2024), and SBWT (Alanko *et al.* 2023).

BWT (Burrows and Wheeler 1994) with support for masks in the suffix array (SA) order (Fig. 1b).

Definition 1 (MBWT). Let (S, M) be a masked superstring and let S' be the BWT image of S . The *SA-transformed mask* is a string M' of the same length as M such that for all i , $M'[i] = M[j]$ where j is the starting position of the i -th lexicographically smallest suffix of S . We call the pair (S', M') the *Masked Burrows-Wheeler Transform* (MBWT) of (S, M) .

The SA-transformed mask can be obtained alongside the standard BWT via a small tweak: we start by rotating the mask by one position to the left, “glue” the mask symbols to the superstring symbols, and then compute the BWT of the superstring with respect to the lexicographical order based solely on the superstring characters (i.e. disregarding the mask bits in the string comparisons). We formalize the whole process in [Algorithm 1](#) and prove its correctness in [Lemma 1](#).

Lemma 1 (SA-transformed mask construction).

Let (S, M) be a masked superstring, M' the SA-transformed mask, S' the BWT of S , and $BWT_S^{-1}(i)$ the original coordinates in S of the i -th character from S' . Then $M'[i] = M[BWT_S^{-1}(i) + 1 \bmod |S|]$.

Proof. Consider the i -th lexicographically smallest rotation of S . The last character of this rotation has position $BWT_S^{-1}(i)$ in S and since the first character cyclically succeeds it in S , it has coordinate $BWT_S^{-1}(i) + 1 \bmod |S|$ in S . Therefore, $M'[i]$ must be $M[BWT_S^{-1}(i) + 1 \bmod |S|]$. \square

Therefore, MBWT can be constructed with any method for BWT computation that allows alphabets with eight symbols, to represent all pairs (s, m) of superstring and mask symbols, while considering pairs (s, m) and (s, \hat{m}) with the same superstring symbol as lexicographically equivalent even if $m \neq \hat{m}$, that is, using only s for comparisons.

2.3 The FMSI membership and dictionary data structures

The FMSI data structure combines the two ingredients with suitable supporting data structures. FMSI indexes an approximately shortest masked superstring (S, M) via its MBWT image (S', M') , which is complemented by the rank data structure for S' . In the case of the dictionary variant of FMSI, it additionally requires select for S' and rank and select for M' (Fig. 1b).

The memory footprint of FMSI is tightly connected with the (S', M') representation, whose size should be minimized. Since approximately shortest k -mer superstrings are typically not compressible beyond two bits per k -mer, a natural representation for S' is a plain bit-vector. In contrast, the masks tend to be highly imbalanced: towards all-zero vectors for subsampled k -mer sets, and towards all-one vectors with the spectrum-like property (Sladký et al. 2023). Therefore, S' can be represented using rank and select succinct data structures, and M' can be encoded via zeroth-order entropy compression in data structures such as the *succinct indexable dictionary* (Raman et al. 2007; Feigenblat et al. 2016), with support for constant-time rank and select operations.

Definition 2 (FMSI). Given a k -mer set K , a corresponding masked superstring (S, M) , and its associated MBWT image (S', M') , we define the following data structures:

- **FMSI-memb** = $(S'$ with rank, M'), assuming that M maximizes the number of 1's in the mask, for membership queries (P1).
- **FMSI-dict** = $(S'$ with rank and select, M' with rank and select), assuming that M minimizes the number of 1's in the mask, for dictionary queries (P2).

Algorithm 1 MBWT computation.

```

1: function COMPUTEMBWT( $S, M$ )
2:    $M = \text{ROTATE\_TO\_LEFT}(M, 1)$ 
3:    $Z = \text{zip}(S, M)$ 
4:    $Z' = \text{BWT}(Z, \text{compare} = x \rightarrow x.\text{first})$ 
5:    $S', M' = \text{unzip}(Z')$ 
6:   return  $S', M'$ 

```

Algorithm 2 k -mer search in FMSI

```

1: function SINGLESEARCH( $Q$ ) ▷ Backward search
2:    $i = 0, j = |S'|$ 
3:   for all  $c \in \text{reversed}(Q)$  do
4:      $i, j = |S|_{<c} + S'.\text{rank}(c, i), |S|_{<c} + S'.\text{rank}(c, j)$ 
5:     if  $i = j$  then return 0, 0
6:   return  $i, j$ 
7: function STREAMINGSEARCH( $Q, i, j$ )
8:   if  $i = j$  then return  $i, j = \text{SINGLESEARCH}(Q)$ 
9:    $i = \text{klcp.select}(0, \text{klcp.rank}(0, i))$ 
10:   $j = \text{klcp.select}(0, 1 + \text{klcp.rank}(0, j - 1))$ 
11:   $c = Q[0]$ 
12:  return  $i, j = |S|_{<c} + S'.\text{rank}(c, i), |S|_{<c} + S'.\text{rank}(c, j)$ 

```

Algorithm 3 Query operations (given SA range $[i, j]$ or hash h)

```

1: function MEMBER( $i, j$ ) ▷ requires max-one mask
2:   return  $i \neq j \ \& \ M'[i] = 1$ 
3: function LOOKUP( $i, j$ ) ▷ requires min-one mask
4:    $r_i = M'.\text{rank}(1, i), r_j = M'.\text{rank}(1, j)$ 
5:   if  $r_i = r_j$  then return -1 else return  $r_i$ 
6: function ACCESS( $h$ )
7:    $Q = "", i = M'.\text{select}(h)$ 
8:   for  $k$  steps do
9:      $c = \max\{c \in \Sigma : |S|_{\leq c} < i\}$ 
10:     $i = S'.\text{select}(c, i - |S|_{<c})$ 
11:     $Q = Q + S'[i]$ 
12:   return  $Q$ 

```

2.3.1 FMSI index theoretical properties

We now analyze the theoretical guarantees of the FMSI index construction and its space requirements. First, we analyze index construction time, starting from a suitable (max-one or min-one) masked superstring (S, M) of the indexed k -mer set K .

Lemma 2. *Given a masked superstring (S, M) of a k -mer set K , both FMSI-memb and FMSI-dict data structures can be constructed in expected time $O(|S|)$.*

Proof. The MBWT image (S', M') can be computed in linear time with respect to the superstring length (Algorithm 1). All aforementioned data structures for rank and select can be constructed in $O(|S|)$ time as well; only for the compressed representation of M' using the succinct indexable dictionary, the construction runs in *expected* linear time (Feigenblat et al. 2016). \square

Second, we analyze the storage requirement of FMSI—first, in the general case of arbitrary k -mer sets, and then when the underlying superstring size approaches the number of distinct canonical k -mers (e.g. only 4% difference for the human genome and $k = 31$; Tab. 3).

Lemma 3 (FMSI storage requirements). *Given a masked superstring (S, M) of a k -mer set K , the corresponding FMSI-memb and FMSI-dict indexes require $2|S| + \log \binom{|S|}{|M|_1} + o(|S|)$ bits, which is at most $(3 + o(1)) \cdot |S|$. Furthermore, if $|S| = |K| + o(|K|)$, the index requires $2 + o(1)$ bits per distinct canonical k -mer.*

Proof. S' stored as a plain bit-vector-based wavelet tree (Grossi et al. 2003) requires $2 + o(1)$ bits per character, including the rank and select support. M' stored using a zeroth-order entropy compression (Raman et al. 2007, Feigenblat et al. 2016) requires $\log \binom{|S|}{|M|_0} = \log \binom{|S|}{|M|_1}$ bits. Both rank and select for the compressed vector M' occupy only sublinear space (Raman et al. 2007; Feigenblat et al. 2016). Last, suppose that $|S| = |K| + \Delta$ with $\Delta \in o(|K|)$. Letting $\Delta = \frac{|K|}{f(|K|)}$ for $f(|K|) \rightarrow \infty$ as $|K| \rightarrow \infty$, the space complexity of M' is $\log \binom{|S|}{|M|_1} \leq \log \binom{|K| + \Delta}{\Delta} \leq \Delta \log \left(\frac{2e|K|}{\Delta} \right) = |K| \frac{\log(2e \cdot f(|K|))}{f(|K|)} = o(|K|)$ and, thus, the index requires $2 + o(1)$ bits per distinct k -mer. \square

2.3.2 FMSI membership and dictionary queries

FMSI can act as two slightly different data structures, a membership index and a dictionary, with membership queries being slightly more space- and time-efficient (Table 2). In both cases, the search starts with the standard backward search, as used, e.g. in the FM index, resulting in a suffix array (SA) interval corresponding to occurrences of the query k -mer in the superstrings (Algorithm 2). The backward search uses the “count” values for characters c , denoted $|S|_{<c}$ and equal to the total number of occurrences of all characters $c' < c$ in S , which can be computed during the FMSI construction in constant space. The candidate k -mer matches in the SA intervals, if non-empty, are subsequently validated against M' to exclude false positives or to obtain the k -mer’s hash (Algorithm 3). In the bidirectional model, everything works the same way with only one difference: If the forward k -mer is not found in the superstring, we subsequently query also the RC of the k -mer.

Table 2. Query time and memory usage for membership vs. dictionary queries.^a

Dataset	FMSI-memb	FMSI-dict
ec-pg-hq	17.2 3.3	20.5 3.4
mtg-ilm	7.7 3.6	14.6 3.7
hg-t2t	11.4 2.4	14.8 2.4

^a Comparison of query time (per query k -mer, μ s) and memory (bits per distinct canonical k -mer, in italics) on three datasets with $k = 31$. See Section 3.

Table 3. Index construction resources on the human genome.^a

Method	SPSS	MS	Index	Overall
CBL	–	–	38.4 101.1	38.4 101.1
SBWT-small	–	–	21.5 ^b 49.9	21.5 ^b 49.9
SBWT-fast	–	–	39.5 ^b 93.1	39.5 ^b 93.1
SSHash ^c	30.5 ^b 21.6	–	14.6 10.9	45.1 21.6
FMSI-memb	30.5 ^b 21.6	25.3 39.0	15.4 44.2	71.2 44.2
FMSI-dict ^c	30.5 ^b 21.6	14.1 1.7	18.8 44.2	63.4 44.2

^a CPU time of construction (usr + sys, minutes) and memory usage (GB, in italics) for the human genome dataset and $k = 31$. The greedily computed superstring had length $|S| = 2.602$ G, with $|K| = 2.512$ G. See Section 3.

^b With eight threads (1 otherwise).

^c Indexes supporting dictionary queries.

Membership queries. In FMSI-memb, MEMBER is evaluated by inspecting the first bit of the SA-transformed mask within the interval $[i, j)$ obtained via backward search (Algorithm 3). This is possible thanks to the max-one property of the mask, which ensures that all bits in the interval are set to the same value.

Dictionary queries. FMSI-dict supports the LOOKUP operation via additional rank data structure, required for two additional invocations of M' .rank to compute $r_i = M'.rank(1, i)$ and $r_j = M'.rank(1, j)$ (Algorithm 3). If $r_i = r_j$, the k -mer is not in K and is reported as absent from the dictionary; otherwise, r_i is returned as its unique minimal hash. This is possible thanks to the min-one property of the mask, which ensures that r_i equals the number of distinct represented k -mers appearing before the queried k -mer in the SA order, making r_i both unique and minimal.

In the other direction, ACCESS requires select data structures for both S' and M' . It is evaluated via $M'.select(h)$, which returns the position i of the h -th occurrence of 1 in M' . This position marks the start of the SA interval corresponding to the queried k -mer Q , which we then reconstruct in k steps by iteratively applying select on S' (Algorithm 3).

Queries with unoptimized masks. Membership operations are possible with masks other than max-one, but at the cost of additional time overhead due to two invocations of rank on M' to determine the presence of 1 in the SA interval via $M'.rank(1, i) < M'.rank(1, j)$. Likewise, dictionary operations are possible with masks other than min-one, but at the cost of losing the contiguity and surjectivity of hash values, as their range length increases by $|M|_1 - |K|$.

2.4 Streaming acceleration

In many practical scenarios, k -mers are queried in batches consisting of all overlapping k -mers of a given query string, such as a sequencing read or a genomic element. Streamed queries are often predominantly positive, e.g. when a query gene differs only in one nucleotide due to a single mutation. By *positive streamed queries*, we mean a query for (all

occurrences of k -mers of) a string T such that all the k -mers of T are present in the represented k -mer set.

To accelerate positive streamed queries, we use the k LCP array technique, originally developed in the context of k -mer indexing with FM-indexes (Salikhov 2017) and previously implemented in ProPhex (Salikhov *et al.* 2018). This technique is based on the bit-vector of the same length as S such that the i th bit is set to 1 if and only if the longest common prefix between the i -th and $(i+1)$ th suffixes of S in the lexicographic order has length at least $k-1$. This allows us to obtain a range corresponding to a $(k-1)$ -mer obtained after deleting the last character. In the unidirectional model, this immediately gives fast positive streamed queries, in the bidirectional model, we need to ensure that we query most k -mers first in the direction in which they appear in the superstring.

Unidirectional streaming via k LCP. Without taking RCs into account, we traverse the query sequence in the reverse direction (Algorithm 3). Starting from the last k -mer, after we have queried a *present* k -mer, we obtain the range for its prefix of size $k-1$ by extending the range in the SA coordinates in both directions to the first 0 in the k LCP array using two rank and select operations. To get the range for the next k -mer, we perform an update of the range with its first character, i.e. one step of the backward search, which only costs $O(1)$ time in total. Only when the current queried k -mer does not appear in the superstring, we perform the backward search for the next k -mer in time $O(k)$.

Streaming bidirectionality via saturating strand counters. In the bidirectional model, consecutive queried k -mers may change directionality in the superstring; however, the k LCP extension technique supports only one fixed direction. Nevertheless, streaming can still be accelerated locally if we can split queries into locally consistent parts and query them in the right direction. As a simple but effective heuristic, we split queries into blocks of size B , and for each block we estimate the likely dominating directionality of its k -mers via a saturating counter. Specifically, whenever a k -mer is found on the forward strand, we increment the counter, and decrement otherwise (and vice versa for the reverse strand). Then, the directionality of the next block can be estimated based on the statistics stored in this counter, informing whether the block should be first reverse-complemented or not.

This technique provides favorable theoretical time guarantees for positive streamed queries, both in the unidirectional and bidirectional model (Lemma 4). For the unidirectional model, we have slightly better guarantees; however, the difference vanishes as the size of the query sequence grows. The speedup in the bidirectional model assumes that consecutive k -mers do not change strand in the superstring too often; for example, this readily holds for greedy superstrings built from k -mer sets corresponding to DBGs with limited branching.

Lemma 4 *Given an FMSI index over a masked superstring (S, M) of a k -mer set K , and a query sequence T containing $t = |T| - k + 1$ occurrences of k -mers from K , streamed queries can be processed such that:*

- In the unidirectional model, querying T takes time $O(t+k) = O(|T|)$, or $O(1 + \frac{k}{\sqrt{t}})$ per k -mer.
- In the bidirectional model, if there are at most $O(t/k^2)$ pairs of consecutive k -mers Q_i, Q_{i+1} in T , such that Q_{i+1} does not appear on the same strand as Q_i in S , then

querying k -mers in T takes time $O(t+k\sqrt{t})$, or $O(1 + \frac{k}{\sqrt{t}})$ per k -mer.

Proof sketch. Querying the last k -mer of T takes $O(k)$ time. In the unidirectional model, since all queries are positive, every other k -mer of T can be queried in $O(1)$ time as explained above.

In the bidirectional model, let $B = \min(\sqrt{t}, k)$. We split T into blocks of size $\Theta(B)$. For each of these blocks, we predict whether the k -mer or its RC are in S using the saturating counter with small constant saturation computed on previous blocks, which predicts to start with the strand with which the previous block ended with. We process the first block in time at most $O(B \cdot k)$. For other blocks, we spend time $O(k)$ on the first k -mer. Then, if we predict the strand correctly and the strand does not change during the block, we spend time $O(k+B)$ per block, and otherwise, we spend time $O(k)$ for each k -mer. The total time spend on blocks in the former case is thus $O((t/B) \cdot (B+k)) = O(t+k \cdot \sqrt{t})$. The latter can happen if consecutive k -mers switch strands either inside the block, in the previous block, or between blocks which may cause an incorrect prediction. However, such switching happens at most $O(t/k^2)$ times in total and hence, the total time complexity is bounded by $O(B \cdot k + t + k \cdot \sqrt{t} + (t/k^2) \cdot B \cdot k)$, which is always bounded by $O(t+k\sqrt{t})$.

2.5 Implementation in the FMSI program

We implemented the FMSI index in a tool called FMSI, an abbreviation of ‘From Superstring to Indexing’. FMSI is developed in C++, provided on Github under the MIT license (<https://github.com/OndrejSladky/fmsi>), deposited on Zenodo (<https://doi.org/10.5281/zenodo.13905020>), and distributed via Bioconda (Grüning *et al.* 2018). The version of FMSI used in this paper is 0.4.0. FMSI can calculate an index for user-provided masked superstring (S, M) for an arbitrary k , and then query it with Membership or Lookup queries, possibly streamed. The input masked superstring can be computed by KmerCamel, or alternatively inferred from any SPSS representation.

The index in the FMSI software consists of three components: S' as a BWT image of S (stored via a wavelet tree over non-compressed plain bit-vectors), M' as the corresponding SA-transformed mask (an RRR-compressed bit-vector (Raman *et al.* 2007)), and an optional k LCP array (a plain bit-vector), all based on the sds-lite library (Gog *et al.* 2014). MBWT computation proceeds via suffix array construction using QSufSort (Larsson and Sadakane 2007).

For single- k -mer queries, we also utilize the strand prediction technique we described for streamed queries. If a k -mer is not found in the superstring (Algorithm 2 and 3), its RC is also queried, starting from a predicted strand (using a counter with saturation 7). FMSI uses two additional counters, to accumulate and predict only if the previous k -mer was found in the forward or reverse direction respectively. Streaming (Algorithm 2) additionally simplifies the range extension: the k LCP array is traversed in both directions until 0s are found, avoiding the additional rank data structures. Bidirectionality is supported via splitting each queried sequence of length ℓ into parts of size $\lfloor 2\sqrt{\ell} \rfloor$. The initial query strand of each part is predicted from the saturating counters of the preceding parts, with only the first strand queried entirely and the other one only for the not-yet-found k -mers.

3 Experimental evaluation

Experimental setup. We benchmarked FMSI against three state-of-the-art indexes for single k -mer sets: Spectral Burrows-Wheeler Transform (SBWT) (Alanko *et al.* 2023, commit b795178), SShash (Pibiri 2022), version 3.0.0, Conway-Bromage-Lyndon (CBL) (Martayan *et al.* 2024, commit 328bcc6), each configured as recommended by their authors. Specifically, SShash used minimizer sized as $\min\{\lceil \log_4 |K| \rceil + 1, k - 2\}$ and indexes based on optimal simplitigs (eulertigs) computed by GGCAT (Cracco and Tomescu 2023). CBL was independently compiled per k -mer size, with SIMD support and 28 prefix bits. SBWT was evaluated in two distinct modes: a standard mode optimized for speed (*SBWT-fast*, with a plain matrix and index with RCs) and mode tuned for minimum memory (*SBWT-small*, with RRR-split and index without RCs), both considered with and without streaming support. *DbgFM* (Chikhi *et al.* 2015) was excluded due to scalability issues and recurrent index construction failures.

We selected eight datasets representative of diverse tradeoffs between k -mer set sizes and sets' structural properties common across bioinformatics applications: SARS-CoV-2 pangenome (*sc2-pg*; $n = 17\text{M}$; GISAID; Shu and McCauley (2017)), *S. pneumoniae* pangenome (*sp-pg*; $n = 616$; Croucher *et al.* (2015)), *Escherichia coli* 661k pangenomes (Blackwell *et al.* 2021) with (*ec-pg-hq*; $n = 86\text{K}$) and without (*ec-pg-all*; $n = 89\text{K}$) quality filtering, metagenomic Illumina reads of a human tongue microbiome (*mtg-ilm*; SRX023459), RNA-seq Illumina reads (*rna-ilm*; SRX348811), human T2T assembly (*hg-t2t*; chm13v2.0; Nurk *et al.* (2022)), and human WGS Illumina reads (*hg-ilm*; SRX016231). Each dataset was preprocessed with GGCAT (Cracco and Tomescu 2023) to compute an SPSS, followed by KmerCamel (Sladký *et al.* 2023) to compute greedy masked superstrings for FMSI. Illumina raw-read datasets (**-ilm*) were cleaned via k -mer filtering ($k = 32$, min-freq 2). An example summary of construction resources for *hg-t2t* and $k = 31$ is provided in Table 3. Benchmarks were performed on an AMD EPYC 7302 (3 GHz) server with 251 GB RAM and SSD storage. Comprehensive protocols and pre-processed datasets are provided in Supplementary material, available as supplementary data at *Bioinformatics Advances* online.

FMSI consistently achieves best space efficiency. We systematically evaluated the memory usage of FMSI, SBWT-small, SBWT-fast, SShash, and CBL queries across dataset types (without subsampling) and various k -mer sizes (Table 4). FMSI consistently showed the lowest memory footprint. Specifically, it was by 1–2 orders of magnitude more space-efficient than CBL, over $2.5\times$ more space-efficient than SShash, and over $3\times$ more space-efficient than SBWT-fast. The only method competitive in space with FMSI was SBWT-small, but only on large k -mer sets with *DBG* dominated by long paths (e.g. the human genome assembly *hg-t2t* and the *E. coli* pangenome without contamination filtering *ec-pg-all*, both only at large k). Since index serialization on disk enables additional compression and omits recomputable index parts and buffers, we also evaluated on-disk index sizes (smaller values in Table 4). The results were consistent, except for *sc2-pg*, where the SBWT-small serialized index was marginally smaller than FMSI's at large k .

FMSI is robust to k . We next examined how memory usage varied with k (Fig. 2). FMSI and SBWT-fast exhibited stable memory consumption per k -mer, remaining within a narrow range of the minimum observed value (except *rna-ilm*). SShash' performance varied largely with k , with the best space requirements achieved for large values of k ; otherwise, it required up to 2–4 \times more space. SBWT-small remained space efficient at large k , but performance was degraded by a factor of 1.6–3.3 at smaller k . In contrast, CBL showed the opposite trend—beyond $k > 23$, its memory exceeded 40 bits per k -mer. We note SShash and CBL performance may be further tuned via targeted (dataset, k)-specific optimization of the minimizer-size and prefix-bit parameters, respectively.

FMSI maintains competitive query times. Next, we evaluated whether FMSI maintains competitive query times compared to other methods (Fig. 3a). We selected $k = 23$ as a representative setting where all indexes achieve reasonable space usage and compared their positive and negative query times, with and without streaming acceleration. Overall, FMSI's query speed was in the mid-range relative to the other methods. SShash and SBWT-fast were generally the fastest, outperforming FMSI by a factor of 2–4. At the same time,

Table 4. Memory and disk usage across k -mer sizes for membership query indexes.^a

Dataset	k	31-mers	CBL	SShash	SBWT-fast	SBWT-small	FMSI
sc2-pg	15, 17, ... 43	12.7 M	153.8–654.1 ↳46.6–78.7	14.0–25.0 ↳10.7–19.3	22.2–56.0 ↳9.4–10.3	16.0–50.8 ↳3.2 ^b –5.0	5.5 ^b –9.0 ^b ↳3.4–3.7 ^b
sp-pg	15, 17, ... 43	9.62 M	150.3–300.9 ↳47.5–88.9	8.2–22.6 ↳5.5–18.0	29.9–45.9 ↳9.2–9.8	24.5–41.4 ↳3.6–5.1	6.1 ^b –7.4 ^b ↳2.8 ^b –3.4 ^b
ec-pg-hq	15, 17, ... 43	551 M	34.3–173.0 ↳20.7–80.9	6.6–22.7 ↳6.6–22.6	9.0–9.7 ↳8.5–8.6	3.4–7.8 ↳3.1–6.6	3.1 ^b –3.5 ^b ↳2.8 ^b –3.3 ^b
ec-pg-all	15, 23, ... 31	1.32 G	25.5–69.0 ↳18.7–48.9	7.6–22.0 ↳7.5–21.9	8.8–9.2 ↳8.5–8.6	3.1–7.8 ↳2.9–7.1	2.9 ^b –3.0 ^b ↳2.7 ^b –2.8 ^b
mtg-ilm	15, 17, ... 31	355 M	37.1–65.5 ↳21.0–50.9	3.8–22.7 ↳7.7–22.6	9.5–10.8 ↳8.6–10.1	3.7–7.3 ↳3.1–6.1	3.0 ^b –3.7 ^b ↳2.8 ^b –3.4 ^b
rna-ilm	15, 23, ... 31	686 M	39.6–130.4 ↳20.5–48.8	14.9–22.8 ↳14.9–22.6	9.7–14.4 ↳8.7–14.0	4.8–7.8 ↳4.4–6.7	3.3 ^b –6.2 ^b ↳3.0 ^b –5.9 ^b
hg-ilm	15, 23, ... 31	3.56 G	25.1–341.1 ↳18.6–48.5	10.8–20.6 ↳10.8–20.6	9.2–10.5 ↳8.5–10.4	3.5–8.5 ↳3.4–7.8	2.9 ^b –4.1 ^b ↳2.6 ^b –3.9 ^b
hg-t2t	15, 17, ... 43	2.51 G	21.1–597.0 ↳19.0–73.1	5.9–24.2 ↳5.9–24.2	8.6–9.2 ↳8.5–8.5	2.4 ^b –8.0 ↳2.3–7.3	2.4 ^b –3.0 ^b ↳2.2 ^b –2.9 ^b
Overall			21.1–654.1 ↳18.6–88.9	5.9–25.0 ↳5.5–24.2	8.6–56.0 ↳8.5–14.0	2.4 ^b –50.8 ↳2.3–7.8	2.4 ^b –9.0 ^b ↳2.2 ^b –5.9 ^b

^a Reported values are in bits per canonical k -mer, measured across all benchmark datasets and indicated k -mer sizes. Space usage was measured in memory (via GNU time; larger numbers) and on disk (smaller, italicized numbers below memory usage, preceded by ↳). Due to memory constraints (>251 GiB), CBL could not be evaluated for *hg-t2t* at $k > 39$.

^b The best lower and upper value for each dataset.

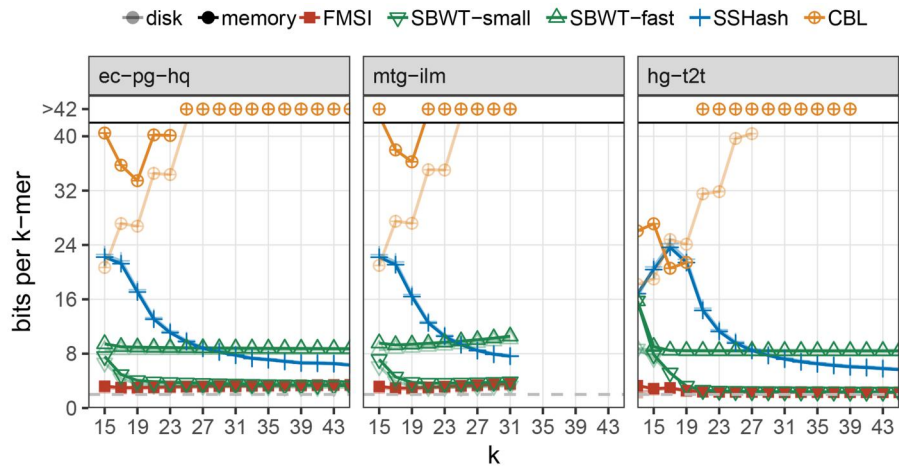


Figure 2. Memory usage of indexes across k -mer lengths. Memory usage is reported in bits per canonical k -mer, shown as a function of k . Corresponding on-disk sizes are indicated in lighter color. For *mtg-ilm*, evaluation is limited to $k \leq 32$ due to initial k -mer filtering. The gray dashed line indicates two bits per k -mer.

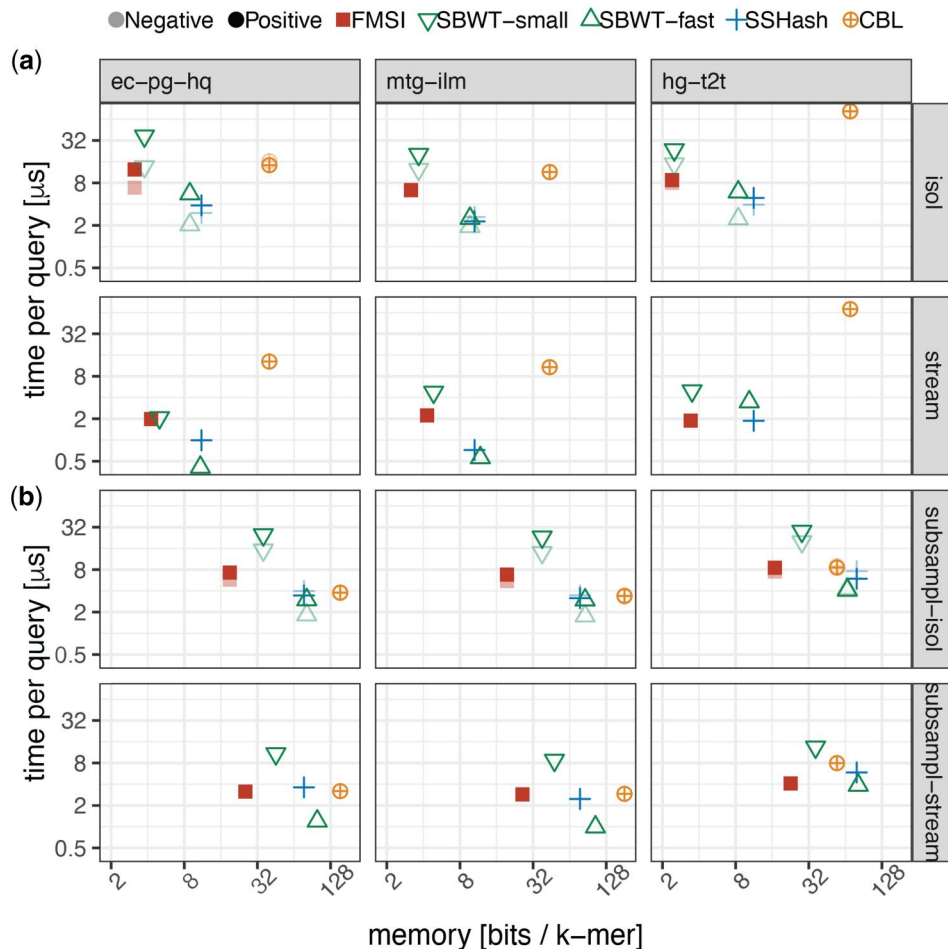


Figure 3. Time-memory trade-offs for membership queries for (a) regular and (b) subsampled k -mer sets (10%). Query time (per k -mer) and memory usage are shown for selected datasets with $k = 23$, with and without streaming acceleration, and with and without reference set subsampling.

FMSI was usually 2–4 \times faster than the second best space-efficient tool, SBWT-small, although ties occurred, e.g. for *E. coli* pangenomic datasets with streaming support. FMSI's use of k LCP for accelerated streaming improved queries 2–4 \times (and similarly for SBWT), but at the cost of an additional bit per superstring character. In contrast, SSHash and CBL

accelerate streaming without additional memory overhead. Importantly, the relative performance of methods under streaming remained consistent with that observed in isolated query mode. Interestingly, for *hg-t2t* with accelerated streaming, FMSI achieved comparable query times as SSHash, both being approximately 2 \times faster than SBWT-fast.

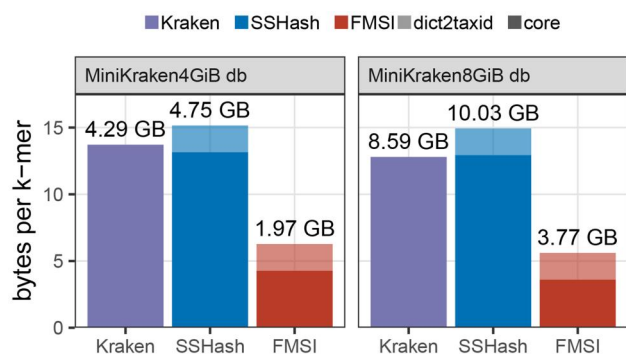


Figure 4. Disk usage of LCA indexes on MiniKraken database. Comparison of disk size for Kraken, FMSI, and SSHash indexes on subsampled MiniKraken v20171019 ($k = 31$). FMSI and SSHash include an additional two bytes per k -mer to associate each hash with the lowest common ancestor in the taxonomic tree.

Membership on subsampled k -mer sets. Next, we sought to understand the effect of sampling on the query time and space of all indexes. We uniformly subsampled the k -mer sets to 10% on the previous k -mer sets (Fig. 3b). This resulted in a substantial increase in memory per k -mer of all tested tools, except for CBL on hg-t2t. The memory increase is primarily due to missing overlaps of length $k - 1$ in the dataset, resulting in substantially larger representation lengths, when normalized by the k -mer set size. FMSI typically used at most half the space of the second best tool, SBWT-small, while also being 2 \times faster. FMSI was approximately 4 \times more space-efficient than SBWT-fast and SSHash, while being 2–3 \times slower on isolated queries. FMSI was particularly efficient on streamed queries, similarly fast as SSHash while using about 4 \times less space, achieved thanks to the “ k -mer interpolation” nature of superstrings.

FMSI reduces the Minikraken index size 2 \times . Finally, we sought to demonstrate the utility of FMSI as a dictionary of sampled k -mer sets, and evaluated its compressive properties on the subsampled database of the Kraken metagenomic classifier (Wood and Salzberg 2014) termed the MiniKraken databases (Fig. 4). We built FMSI-dict and SSHash 31-mer dictionaries for both MiniKraken 31-mer databases available online (with an additional two bytes per k -mer for taxid arrays), and compared them to the size of the original indexes. We found that while SSHash moderately increased memory requirements compared to the Kraken index, FMSI provided more than a twofold memory improvement in both cases.

4 Discussion and conclusions

k -mer-based methods are central to modern bioinformatics, and efficient k -mer set indexing is therefore a critical objective. In this work, we introduced FMSI, a data structure and its implementation for space-efficient indexing of single k -mer sets. FMSI consistently achieved lower space usage than all other state-of-the-art methods across a broad operational range, while maintaining competitive query times. This was possible thanks to approximately shortest superstrings and BWT indexing—two techniques that had not previously been combined in this context.

Several performance tradeoffs in FMSI could be further improved. First, index construction currently depends on QSufSort (Larsson and Sadakane 2007), which limits its scalability, whereas modern BWT construction algorithms

already scale to terabyte-sized inputs (Li 2024). Second, query speed would benefit from parallelization, implemented similarly to BWT-based indexes (Li 2012, Alanko *et al.* 2023). Third, space usage may be further reduced by applying Elias-Fano (Fano 1971, Elias 1974) to the SA-transformed mask instead of RRR (Raman *et al.* 2007). Finally, FMSI could be re-optimized for speed instead of space—for instance, by adding a reverse complement superstring to the index or replacing wavelet trees with 2-bit encoding, in analogy to the fast modes in SBWT.

Despite decades of work and a rich palette of available methods, k -mer indexing remains an exciting problem with substantial room for improvement. This is particularly relevant in settings with subsampling, sketching, or high levels of genetic polymorphisms and technological artifacts, such as metagenomics and pangenomics. In these cases, superstrings offer an elegant means to bypass the limitations of traditional methods, representing a major weapon in the algorithmic arsenal for future generic, space-efficient, and parameter-free k -mer-based data structures. Their potential impact may extend well beyond bioinformatics, into domains such as indexing in large language models.

Acknowledgements

Portions of this research were conducted at the GenOuest bioinformatics core facility (<https://www.genouest.org>). We are grateful to Camille Marchet, Igor Martayan, Giulio Pibiri, Jarno Alanko, and Paul Medvedev for their constructive comments.

Author contributions

Ondřej Sladký, Pavel Veselý, and Karel Břinda developed the methodology, designed and conducted the experiments, analyzed the results, and wrote the manuscript. Ondřej Sladký developed the software.

Supplementary data

Supplementary data are available at *Bioinformatics Advances* online.

Conflict of interest

None declared.

Funding

This work was supported by the French National Research Agency (ANR-24-CE45-1226; REALL); Czech Science Foundation (22-22997S); Czech Ministry of Education, Youth and Sports (ERC-CZ LL2406); Charles University (UNCE 24/SCI/008); French Ministry for Europe and Foreign Affairs, French Ministry for Higher Education and Research, and Czech Ministry for Education, Youth and Sports (PHC BARRANDE 2025 grant no. 52374TC; EFFIMAS).

References

- Alanko JN, Puglisi SJ, Vuohthoniemi J. Small searchable κ -spectra via subset rank queries on the spectral Burrows-Wheeler transform. In: *SIAM ACDA '23*. 2023.
- Bingmann T, Bradley P, Gauger F *et al.* COBS: a compact bit-sliced signature index. In: Brisaboa N, Puglisi S (eds), *String Processing and*

- Information Retrieval. SPIRE 2019*. Lecture Notes in Computer Science, Vol. 11811. Cham: Springer, 2019, 285–303.
- Blackwell GA, Hunt M, Malone KM *et al.* Exploring bacterial diversity via a curated and searchable snapshot of archived DNA sequences. *PLoS Biol* 2021;19:e3001421.
- Bowe A, Onodera T, Sadakane K *et al.* Succinct de Bruijn graphs. In: Raphael B, Tang J (eds), *Algorithms in Bioinformatics. WABI 2012*. Lecture Notes in Computer Science, Vol. 7534. LNCS WABI '12. Berlin, Heidelberg: Springer, 2012, 225–35.
- Bradley P, Den Bakker HC, Rocha EP *et al.* Ultrafast search of all deposited bacterial and viral genomic data. *Nat Biotechnol* 2019;37:152–9.
- Bradley P, Gordon NC, Walker TM *et al.* Rapid antibiotic-resistance predictions from genome sequence data for *Staphylococcus aureus* and *Mycobacterium tuberculosis*. *Nat Commun* 2015;6:10063.
- Bray NL, Pimentel H, Melsted P *et al.* Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol* 2016;34:888.
- Břinda K, Baym M, Kucherov G. Simplitigs as an efficient and scalable representation of de Bruijn graphs. *Genome Biol* 2021;22:96.
- Břinda K, Callendrello A, Ma KC *et al.* Rapid inference of antibiotic resistance and susceptibility by genomic neighbour typing. *Nat Microbiol* 2020;5:455–64.
- Břinda K, Lima L, Pignotti S *et al.* Efficient and robust search of microbial genomes via phylogenetic compression. *Nat Methods* 2025;22:692–7.
- Břinda K, Salikhov K, Pignotti S *et al.* 2017. Prophyle 0.3.1.0: An Accurate, Resource-Frugal and Deterministic DNA Sequence Classifier. <https://doi.org/10.5281/zenodo.1045429> (16 September 2025, date last accessed).
- Burrows M, Wheeler D. A block-sorting lossless data compression algorithm. Technical report 124. Digital Equipment Corporation, 1994.
- Chikhi R, Holub J, Medvedev P. Data structures to represent a set of k-long DNA sequences. *ACM Comput Surv* 2022;54:1–22.
- Chikhi R, Limasset A, Jackman S *et al.* On the representation of de Bruijn graphs. *J Comput Biol* 2015;22:336–52.
- Chikhi R, Limasset A, Medvedev P. Compacting de Bruijn graphs from sequencing data quickly and in low memory. *Bioinformatics* 2016;32:i201–i208.
- Conway TC, Bromage AJ. Succinct data structures for assembling large genomes. *Bioinformatics* 2011;27:479–86.
- Cracco A, Tomescu AI. Extremely fast construction and querying of compacted and colored de Bruijn graphs with ggcatt. *Genome Res* 2023;33:1198–207.
- Croucher NJ, Finkelstein JA, Pelton SI *et al.* Population genomic datasets describing the post-vaccine evolutionary epidemiology of *Streptococcus pneumoniae*. *Sci Data* 2015;2:150058.
- Elias P. Efficient storage and retrieval by content and address of static files. *J ACM* 1974;21:246–60.
- Fano RM. 1971. *On the Number of Bits Required to Implement an Associative Memory*. Cambridge, MA, United States: MIT, Project MAC.
- Feigenblat G, Porat E, Shiftan A. Linear time succinct indexable dictionary construction with applications. In: *2016 Data Compression Conference (DCC)*, Snowbird, UT, USA. p. 13–22. IEEE, 2016.
- Gog S, Beller T, Moffat A *et al.* From theory to practice: Plug and play with succinct data structures. In: *SEA '14*. 2014.
- Grossi R, Gupta A, Vitter JS. High-order entropy-compressed text indexes. In: Gudmundsson J, Katajainen J (eds), *Experimental Algorithms*. SEA 2014. Lecture Notes in Computer Science, Vol 8504. Cham: Springer, 2014.
- Grüning B, Dale R, Sjödin A *et al.*; Bioconda Team. Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nat Methods* 2018;15:475–6.
- Holley G, Melsted P. Bifrost: highly parallel construction and indexing of colored and compacted de Bruijn graphs. *Genome Biol* 2020;21:249.
- Karasikov M, Mustafa H, Danciu D *et al.* Efficient and accurate search in petabase-scale sequence repositories. *Nature* 2025;1–9.
- Khan J, Dhulipala L, Patro R. Fast and scalable parallel external-memory construction of colored compacted de Bruijn graphs with Cuttlefish 3. *bioRxiv*, 2025, preprint: not peer reviewed.
- Khan J, Kokot M, Deorowicz S *et al.* Scalable, ultra-fast, and low-memory construction of compacted de Bruijn graphs with cuttlefish 2. *Genome Biol* 2022;23:190.
- Khan J, Patro R. Cuttlefish: fast, parallel and low-memory compaction of de Bruijn graphs from large-scale genome collections. *Bioinformatics* 2021;37:i177–86.
- Larsson NJ, Sadakane K. Faster suffix sorting. *Theor Comput Sci* 2007;387:258–72.
- Li H. Exploring single-sample SNP and INDEL calling with whole-genome de novo assembly. *Bioinformatics* 2012;28:1838–44.
- Li H. BWT construction and search at the terabase scale. *Bioinformatics* 2024;40.
- Loh P-R, Baym M, Berger B. Compressive genomics. *Nat Biotechnol* 2012;30:627–30.
- Marchet C, Boucher C, Puglisi SJ *et al.* Data structures based on k-mers for querying large collections of sequencing data sets. *Genome Res* 2021a;31:1.
- Marchet C, Kerbiriou M, Limasset A. Blight: efficient exact associative structure for k-mers. *Bioinformatics* 2021b;37:2858–65.
- Martayan I, Cazaux B, Limasset A, Marchet C. Conway–Bromage–Lyndon (CBL): an exact, dynamic representation of k-mer sets. *Bioinformatics* 2024;40:i48–57.
- Minkin I, Pham S, Medvedev P. Twopaco: an efficient algorithm to build the compacted de Bruijn graph from many complete genomes. *Bioinformatics* 2017;33:4024–32.
- Nurk S, Koren S, Rhie A *et al.* The complete sequence of a human genome. *Science* 2022;376:44–53.
- Ondov BD, Treangen TJ, Melsted P *et al.* Mash: fast genome and meta-genome distance estimation using MinHash. *Genome Biol* 2016;17:132.
- Patro R, Duggal G, Love MI *et al.* Salmon provides fast and bias-aware quantification of transcript expression. *Nat Methods* 2017;14:417–9.
- Pibiri GE. Sparse and skew hashing of K-mers. *Bioinformatics* 2022;38:i185–i194.
- Rahman A, Medvedev P. Representation of k-mer sets using spectrum-preserving string sets. *J Comput Biol* 2021;28:381–94.
- Raman R, Raman V, Satti SR. Succinct indexable dictionaries with applications to encoding k-ary trees, prefix sums and multisets. *ACM Trans Algorithms* 2007;3:43.
- Salikhov K. Efficient algorithms and data structures for indexing DNA sequence data. Ph.D. Thesis, Université Paris, 2017.
- Salikhov K, Břinda K, Pignotti S *et al.* 2018. Prophex 0.1.1. <https://doi.org/10.5281/zenodo.1247431> (16 September 2025, date last accessed).
- Schmidt S, Alanko JN. Eulertigs: minimum plain text representation of k-mer sets without repetitions in linear time. *Algorithms Mol Biol* 2023;18.
- Schmidt S, Khan S, Alanko JN *et al.* Matchtigs: minimum plain text representation of k-mer sets. *Genome Biol* 2023;24:136.
- Shu Y, McCauley J. GISAID: global initiative on sharing all influenza data—from vision to reality. *Eurosurveillance* 2017;22.
- Simpson JT, Wong K, Jackman SD *et al.* ABySS: a parallel assembler for short read sequence data. *Genome Res* 2009;19:1117–23.
- Sladký O, Veselý P, Břinda K. Masked superstrings as a unified framework for textual k-mer set representations. *bioRxiv*, <https://doi.org/10.1101/2023.02.01.526717> 2023, preprint: presented at RECOMB-seq 2023.
- Sladký O, Veselý P, Břinda K. 2025. Towards efficient k-mer set operations via function-assigned masked superstrings. In: *Proceedings of the Prague Stringology Conference 2025*. p. 26–40. Prague: Czech Technical University.
- Stephens ZD, Lee SY, Faghri F *et al.* Big data: astronomical or genomics? *PLoS Biol* 2015;13:e1002195.
- Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol* 2014;15:R46.

© The Author(s) 2025. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Bioinformatics Advances, 2025, 00, 1–10

<https://doi.org/10.1093/bioadv/vbaf290>

Original Article