



Peer Community Journal

Section: Mathematical & Computational Biology

Research article

Published
2026-06-17

Cite as
Igor Martayan, Léa Vandamme,
Bede Constantinides, Bastien
Cazaux, Charles Paperman and
Antoine Limasset (2026)
*Accelerating k-mer-based
sequence filtering*, Peer
Community Journal, 6: e55.

Correspondence
igor.martayan@univ-lille.fr

Peer-review
Peer reviewed and
recommended by
PCI Mathematical &
Computational Biology,
[https://doi.org/10.24072/pci.
mcb.100920](https://doi.org/10.24072/pci.mcb.100920)



This article is licensed
under the Creative Commons
Attribution 4.0 License.

Accelerating k -mer-based sequence filtering

Igor Martayan ¹, Léa Vandamme ¹, Bede
Constantinides ², Bastien Cazaux ¹, Charles
Paperman ³, and Antoine Limasset ¹

Volume 6 (2026), article e55

<https://doi.org/10.24072/pcjournal.735>

Abstract

Motivation. The exponential growth of global sequencing data repositories presents both analytical challenges and opportunities. While k -mer-based indexing has improved scalability over traditional alignment for identifying relevant documents, pinpointing the exact sequences matching numerous queries remains a hurdle. In particular, searching for numerous k -mers with a single large query or multiple distinct queries strains existing exact matching tools, whose performance scales poorly with an increasing number of patterns. At the same time, indexing entire vast datasets for infrequent or ad-hoc searches is often resource-prohibitive. Designing fast methods for matching a large number of k -mers without exhaustive pre-indexing is therefore critical. **Contributions.** We propose an efficient solution to the problem of k -mer-based sequence filtering: given a set of k -mers of interests and a threshold, quickly evaluate whether an arbitrary sequence has a number of k -mer matches above or below the threshold. Our approach demonstrates how minimizer-based sketching, alongside SIMD acceleration, can enhance the performance of streaming searches, and is implemented as a Rust tool named *K2Rmini*. On a consumer laptop, *K2Rmini* is able to filter long reads at 2 Gbp/s. **Availability.** <https://github.com/Malfoy/K2Rmini>.

¹Univ. Lille, CNRS, Centrale Lille, UMR 9189 CRISTAL, F-59000 Lille, France, ²University of Birmingham, Birmingham, United Kingdom, ³Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 CRISTAL, France

Peer Community Journal is a member of the
Centre Mersenne for Open Scientific Publishing
<http://www.centre-mersenne.org/>

e-ISSN 2804-3871



CENTRE
MERSENNE

Introduction

The immense volume of available sequencing data presents both significant opportunities and challenges in bioinformatics. Assembled datasets have reached the Petabase scale, as seen in the 2 Petabytes of unitigs from the Logan project (Chikhi et al., 2024), while even curated databases like GenBank are nearing 6 Terabases (Sayers et al., 2024). Although holistic analyses of these large-scale datasets, such as those on viral (Edgar et al., 2022), soil (Ma et al., 2023), gut (Nayfach et al., 2019), and environmental metagenomes (Parks et al., 2017) are promising, they remain computationally expensive. Identifying matches to a query sequence in large collections, a fundamental task exemplified by BLAST (Altschul et al., 1990), has evolved. Traditional alignment tools were overwhelmed by data growth, leading to the use of strings of fixed length k , or k -mers (Zielezinski et al., 2017). To improve scalability, indexing structures allowing false positives were adopted, making shared k -mer fingerprints a necessary, though not sufficient, condition for a match (Lemane et al., 2024; Marchet et al., 2021; Marchet and Limasset, 2023). Some methods further enhance scalability by indexing only a fraction of k -mers, relying on detecting substantial matches from partial signals (Agret et al., 2021; Darvish et al., 2022; Rahman Hera and Koslicki, 2025).

Large-scale indexes can thus efficiently identify a small, relevant subset of documents (e.g., read sets, unitigs, contigs). However, this initial filtering isn't final. Due to false positives from probabilistic indexes, or sparse/non-collinear matches, queried k -mers may not be truly present or usable for mapping. Therefore, recovering specific sequences within identified documents that genuinely share k -mers with the query is crucial. This refinement, moving from relevant documents to relevant sequences suitable for alignment, mitigates downstream scalability issues and can reduce data volume by orders of magnitude, such as identifying a few relevant sequences from billions.

The problem of *k -mer-based sequence filtering* can be formalized as follows: given a set of k -mers of interest $\mathcal{Q} = \{q_1, \dots, q_m\}$ (the “queries”), a set of sequences $\mathcal{S} = \{S_1, \dots, S_n\}$ and a threshold T , output for each $S_i \in \mathcal{S}$ whether $\sum_{q \in \mathcal{Q}} \# \text{occ}(q, S_i) \geq T$, i.e. whether S_i contains at least T occurrences of k -mers of interest. In practice, the threshold can either be set as an absolute value $T \in \mathbb{N}$ or as a relative value proportional to the size of the sequence: given $t \in [0, 1]$, $T_i = \lceil t \cdot (|S_i| - k + 1) \rceil$.

Related work

Although sequence-level indexes have been proposed (Karasikov et al., 2022; Vandamme et al., 2025), indexing and storing entire dataset collections at this resolution remains expensive and challenging at scale, thus favoring lightweight filtering methods directly from unindexed data. Efficient multi-pattern matching is a highly studied topic from both the theoretical (Crochemore et al., 1999) and practical point of view (Faro and Lecroq, 2013). Practical implementation can be found in classical data processing with grep-like efficient tools such as Ripgrep (Gallant, 2025) or Hyperscan (Wang et al., 2019). Some tools perform pattern matching adapted to genomic data, such as Seqkit (Shen et al., 2016, 2024), fqgrep (Homer et al., 2025) or grepq (Crosbie, 2025). However, these solutions tend to scale poorly as the number of query k -mers increases and fail to leverage the fixed length of the search patterns. This is a critical issue, as querying numerous k -mers is often necessary for large queries (contigs, long reads, whole genomes), searching for many distinct sequences (variant or gene collections), or optimizing batched queries. The recent

BackToSequences (Baire et al., 2024) tool, used in the Logan project, addresses these issues by simply indexing query k -mers in a hash table to search for target documents. Very recently, both Deacon (Constantinides et al., 2025) and Cleanifier (Zentgraf et al., 2025) were proposed as scalable solutions for contaminant depletion, that is removing sequences originating from a given genome, which is a direct application of sequence filtering. In particular, Deacon effectively uses many of the techniques presented in this article. Finally, while to our knowledge it has not been used for sequence filtering, the spectral Burrows-Wheeler transform (SBWT) (Alanko et al., 2023) is also a suitable method for this problem thanks to efficient matching statistics (Mäklin et al., 2025).

Contributions

This work presents three main contributions. First, we introduce a new algorithm that uses random minimizers to quickly filter out sequences that contain too few k -mers of interest, effectively reducing the cost of negative matches by a factor $w/2$. Second, we propose an optimized implementation named *K2Rmini*, that takes advantage of vectorized instructions to parse sequences, hash k -mers and compute minimizers efficiently. Third, we compare our approach to many different tools and evaluate their scalability for a large number of patterns, highlighting the benefits of different methods. Our tool, *K2Rmini*, is able to filter 2 Gbp/s on a consumer laptop and is available online at <https://doi.org/10.5281/zenodo.20715935> (Martayan and Limasset, 2026a).

Methods

Classical multi-pattern matching algorithms, which generalize the Knuth-Morris-Pratt approach, construct deterministic finite automata to achieve linear time complexity relative to the pattern and input data sizes. While these methods offer strong theoretical guarantees, they often exhibit poor cache locality. Furthermore, they address the more general problem of matching variable-length patterns, whereas our work focuses exclusively on fixed-size patterns. This constraint creates an opportunity for hash-based techniques, which are better suited for large pattern sets and more amenable to hardware acceleration through vectorized instructions.

One of the state-of-the-art solutions, BackToSequences (Baire et al., 2024), uses an efficient hash table containing all k -mers of interest and skips irrelevant fields. However, its main limitation is its necessity to perform a hash table lookup for every k -mer in the data stream. This process becomes computationally expensive, especially when the hash table is large and does not fit into low-level caches. To address this bottleneck, we propose *K2Rmini* (Martayan and Limasset, 2026a), a method that integrates minimizer-based filtering and Single Instruction, Multiple Data (SIMD) operations.

Using minimizers to upper bound the number of k -mer matches

Minimizers (Roberts et al., 2004; Schleimer et al., 2003) are a widely used k -mer sampling method that selects the smallest m -mer out of every w consecutive one, usually using a random ordering. The set of sampled minimizers is typically much smaller, by a factor $2/(w + 1)$ called the *density* (Pellow et al., 2023), than the whole sequence while preserving a lot of information, making it a key component of large-scale sequence analysis. An important property of minimizers is that they are *locally consistent*: two sequences sharing contexts of w m -mers will always sample

the same minimizers from these contexts. Thus, comparing minimizers can be used as a proxy to identify similar sequences (Darvish et al., 2022; Ekim et al., 2021; Li, 2018).

One core idea of our approach is to leverage this property to avoid exhaustive k -mer scanning. Given a set Q of k -mers of interest, we compute the set of associated minimizers $\mathcal{M}(Q) = \{\text{minimizer}(q) : q \in Q\}$. Then, instead of checking every k -mer from a sequence S , we only compare the minimizers of S against $\mathcal{M}(Q)$. We know that each minimizer match implies that up to w k -mers of S are in Q . Thus, if we have ℓ minimizer matches, the number of k -mer matches is upper-bounded by $\ell \times w$, so having $\ell < \lceil \frac{T}{w} \rceil$ is sufficient to discard a sequence.

While this upper bound is quite accurate for sparse matches, it is very loose (by a factor close to 2) for dense matches. One way to refine this upper bound is to observe that if two minimizers are d positions apart from each other, they cover at most $w + d$ k -mers. Therefore, if we have a chain of consecutive minimizer matches at positions p_1, \dots, p_r , the number of k -mer matches is upper-bounded by $w + p_r - p_1$.

In practice, we often have access to the exact number of k -mers covered by each minimizer as a byproduct of the sliding window algorithm (Groot Koerkamp and Martayan, 2025). In that case, we can simply sum them to get a tight upper bound.

Minimizer-based filtering algorithm

Our algorithm, detailed in Algorithm 1, works in two passes. The first pass computes an upper bound on the number of k -mer matches, as detailed in the previous section, by performing $\mathcal{O}\left(\frac{|S|}{w}\right)$ queries to \mathcal{T}_M (the hash table of minimizers). This has two main advantages: not only do we reduce the number of lookups, but querying \mathcal{T}_M is also faster than \mathcal{T}_Q (the hash table of k -mers) because it is roughly $w/2$ times smaller when k -mers from Q are consecutive.

The second pass only happens for sequences with an upper bound on k -mer matches above the threshold. In that case, we must count the exact number of k -mer matches to identify true hits. Deacon (Constantinides et al., 2025), another tool tailored for decontamination developed by some of the authors, uses a simplified version of the method proposed here. In particular, it omits this second pass, reducing computational overhead at the cost of potential false positives.

Algorithm 1 Minimizer-based filtering

```

1: function Filter(set of  $k$ -mers  $Q$ , set of sequences  $S$ , threshold  $T$ )
2:    $\mathcal{T}_Q \leftarrow \text{HashSet}(Q)$  // hash table of  $k$ -mers of interest
3:    $\mathcal{T}_M \leftarrow \text{HashSet}(\mathcal{M}(Q))$  // hash table of minimizers of interest
4:   for  $S \in \mathcal{S}$  do
5:      $u \leftarrow 0$  // upper bound of  $k$ -mer matches
6:     for each minimizer  $m$  in  $S$ , covering  $c$   $k$ -mers do
7:       if  $m \in \mathcal{T}_M$  then
8:          $u \leftarrow u + c$ 
9:       if  $u < T$  then
10:        yield false
11:        continue
12:      $s \leftarrow 0$  // exact count of  $k$ -mer matches
13:     for each  $k$ -mer  $x$  in  $S$  do
14:       if  $x \in \mathcal{T}_Q$  then
15:          $s \leftarrow s + 1$ 
16:     yield  $s \geq T$ 

```

Vectorization

Multiple parts of our algorithm are accelerated using vectorized instructions. First of all, we implemented a library called `helicase` (Martayan et al., 2026) that vectorizes the parsing of the sequence files and outputs bitpacked representations for faster processing. Then, we build upon `SimdMinimizers` (Groot Koerkamp and Martayan, 2025) to have a vectorized computation of minimizer positions, along with the number of k -mers that they cover. `SimdMinimizers` relies on a completely branchless sliding window minimum algorithm using two stacks and processes 8 independent chunks of the sequence within a single SIMD register. Finally, we use a vectorized rolling hash adapted from `NtHash` (Mohamadi et al., 2016) to accelerate the k -mer lookups in the second pass.

Parallelization

On top of the vectorized computation of the matches, we adopt a parallelization strategy for the main loop based on a producer-consumer model. A producer thread parses the sequences and sends batches of sequences (up to a desired length) to multiple consumer threads, which in turn compute the matches on their batch and output the results. One approach that we have not implemented yet but could be worth trying would be to let each consumer thread directly parse a chunk of the input file in parallel, thus removing the need for a producer thread and reducing the communication overhead.

Results

All experiments were run with the benchmarking workflow available at <https://doi.org/10.5281/zenodo.20715943> (Martayan and Limasset, 2026b). The benchmarking framework is easily extensible, since each additional program only requires a dedicated command wrapper to be added to a common interface before it can be executed and plotted with the other tools. All reported times include both the cost of indexing the queried patterns and the filtering step. Benchmarks were performed on a dual-socket Intel Xeon Gold 6430 machine with 64 cores and 512 GB of RAM running Ubuntu 24.04.2. Because no ready-to-use SBWT-based sequence filtering implementation was available for our setting, we implemented the SBWT baseline ourselves¹ using the Rust `sbwt` crate (<https://docs.rs/sbwt/latest/sbwt/>).

State-of-the-art analysis

To evaluate state-of-the-art tools, we benchmarked a selection of specialized and general-purpose programs while varying the number of queried k -mers. The specialized tools included `Seqkit` 2.13.0 (Shen et al., 2016, 2024), `fqgrep` 1.1.1 (Homer et al., 2025), `grepq` 1.6.5 (Crosbie, 2025), `BackToSequences` 0.8.3 (Baire et al., 2024), `Cleanifier` 1.2.0 (Zentgraf et al., 2025), `Deacon` 0.14.0 (Constantinides et al., 2025), and our SBWT-based implementation. For comparison, we also included `grep` 3.11, `ripgrep` 15.1.0 (Gallant, 2025), and `Hyperscan` 5.4.2 (Wang et al., 2019) (or `Vectorscan` 5.4.12 on ARM). As shown in [Figures 1 and 2](#), we benchmarked these tools on 2.6 Gbp of PacBio HiFi reads from HG002² while increasing the number of queried k -mers.

¹The SBWT baseline is included in the archived benchmarking workflow under `sbwt_filter`: <https://doi.org/10.5281/zenodo.20715943>.

²https://s3-us-west-2.amazonaws.com/human-pangenomics/NHGRI_UCSC_panel/HG002/hpp_HG002_NA24385_son_v1/PacBio_HiFi/15kb/m54328_180928_230446.Q20.fastq

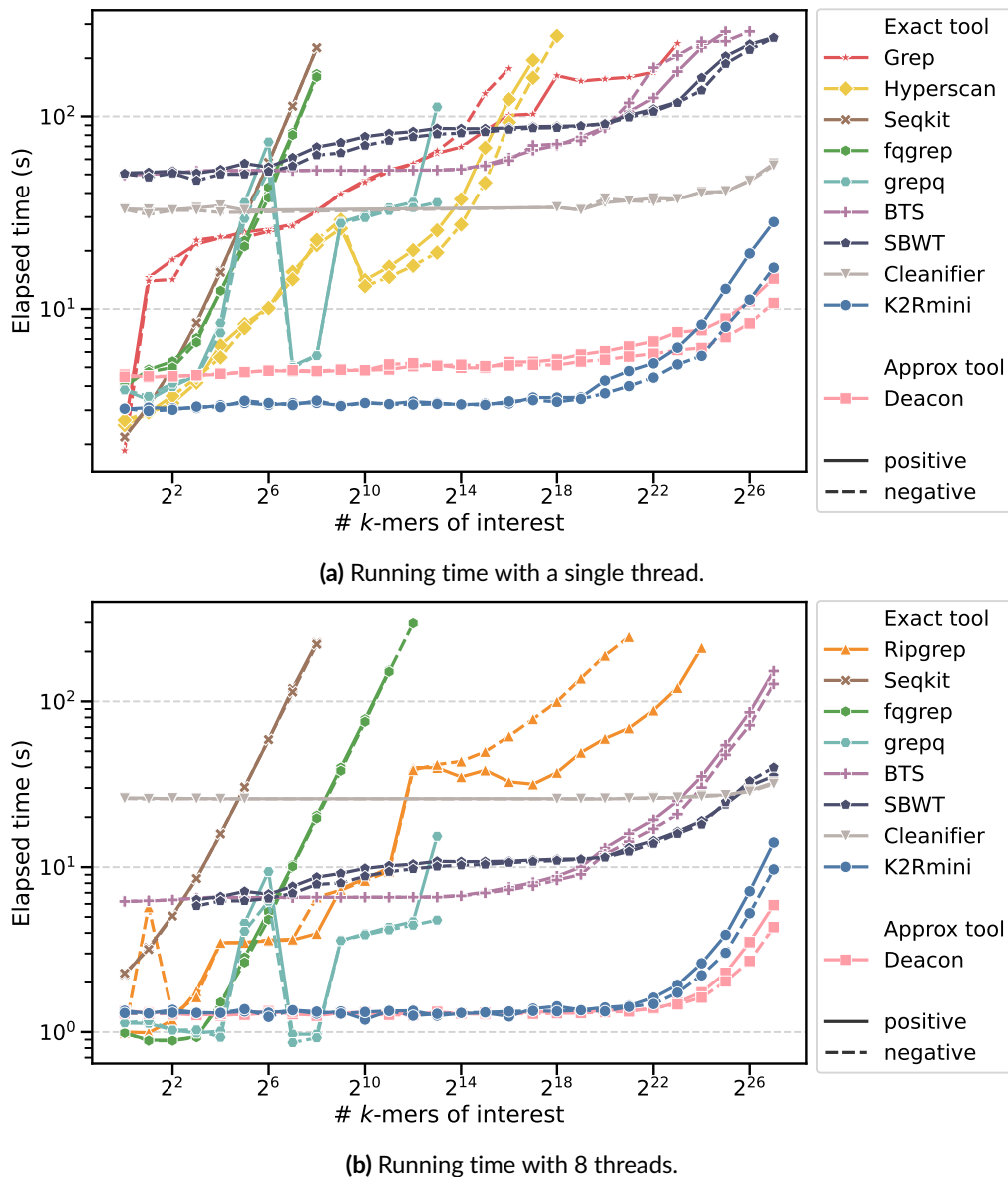


Figure 1 – Running time of state-of-the-art tools when filtering 2.6 Gbp of PacBio HiFi reads from HG002 (m54328_180928_230446) with an increasing number of queried k -mers for $k = 31$. Positive k -mers are extracted from the reads while negative k -mers are random sequences.

The results separate the tested methods into two broad groups. Classical multiple-pattern matching tools such as grep, ripgrep, Hyperscan, Seqkit, fqgrep, and grepq scale poorly as the number of queried k -mers increases. Their running time rises rapidly and they become impractical once the query set reaches the large regimes that matter for genomic filtering. In contrast, methods that explicitly index the queried patterns, or a sketch derived from them, remain effective over the full range of tested query sizes.

Among these scalable approaches, Deacon and K2Rmini are the fastest overall. Cleanifier is also largely insensitive to the number of queried k -mers, but with a clearly higher absolute running time. SBWT and BackToSequences remain much more scalable than general-purpose pattern matchers, although both become slower than Deacon and K2Rmini as the number of

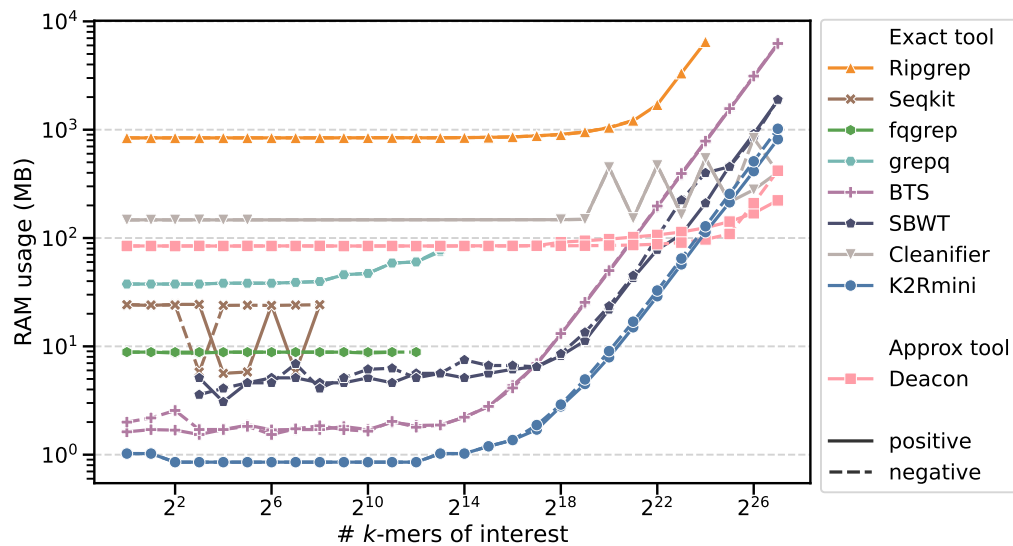


Figure 2 – Memory usage of state-of-the-art tools when filtering 2.6 Gbp of PacBio HiFi reads from HG002 (m54328_180928_230446) with an increasing number of queried k -mers for $k = 31$. Positive k -mers are extracted from the reads while negative k -mers are random sequences.

queried k -mers grows. In the multithreaded setting, SBWT remains consistently ahead of BackToSequences for the largest query sets, while both are clearly outperformed by the minimizer-based approaches.

The gap between positive and negative queries is especially informative for K2Rmini. For negative k -mers, most reads are rejected during the minimizer pass, which keeps the running time close to flat across a broad range of query counts. For positive k -mers, more reads pass this first filter and trigger the exact counting phase, so the running time increases more substantially. This behavior directly reflects the design of K2Rmini, where minimizers are used to avoid unnecessary exact searches, but exact verification is still required whenever the upper bound is above threshold.

This comparison must be interpreted more carefully for Deacon. Deacon only indexes minimizers and does not perform an exact second pass on the surviving reads. Its excellent scalability therefore comes from solving a relaxed version of the problem, which may introduce false positives. By contrast, K2Rmini uses minimizers as a lossless prefilter and performs exact k -mer counting whenever necessary.

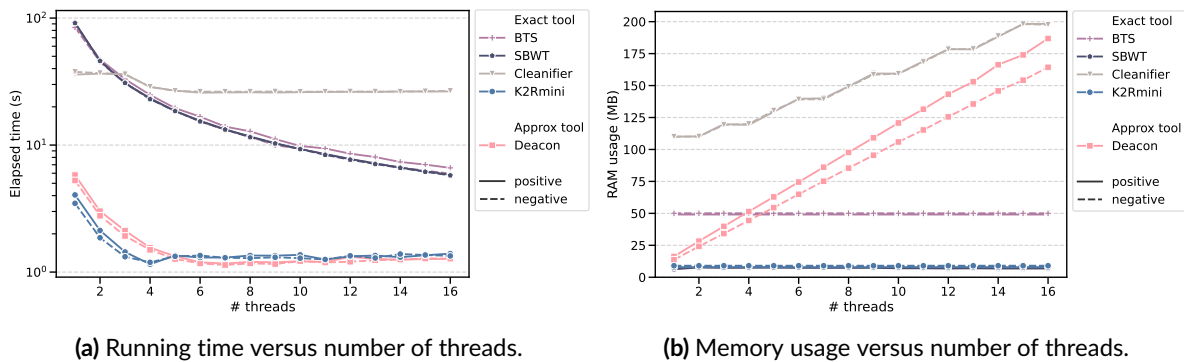
A practical observation is that at some point, pattern indexing itself is a visible part of the total running time and filtering is no longer the only bottleneck.

Memory usage reveals a similar separation between approaches. General-purpose tools such as ripgrep already have a large baseline memory footprint and can reach several gigabytes for the largest query sets. By contrast, the specialized genomic tools initially use much less memory, but their scaling behavior differs substantially. Among the scalable exact methods, K2Rmini has the lowest memory footprint over most of the explored range and remains below BackToSequences and SBWT at large query sizes. BackToSequences starts from a very small footprint on small query sets, but its memory usage grows steeply and eventually becomes the highest among the indexed genomic methods. SBWT scales more smoothly than BackToSequences, but still requires substantially more memory than K2Rmini for large pattern sets. Cleaner and Deacon exhibit a

higher baseline memory usage, yet their growth is more moderate, which is consistent with their good scalability in running time. Overall, these results show that K2Rmini provides the best end-to-end compromise among exact methods, combining low running time with low memory usage over a broad range of query sizes.

Influence of parameters

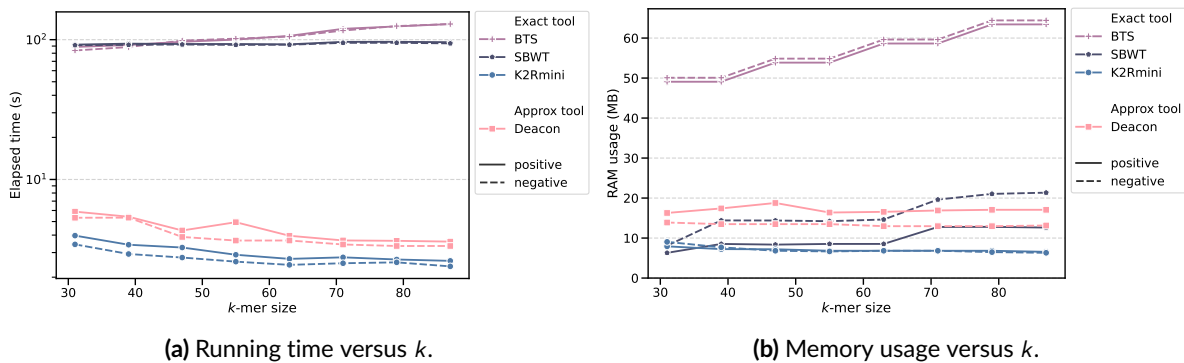
To better understand the behavior of the most scalable approaches, we compared BackToSequences, SBWT, Deacon, and K2Rmini while varying the number of threads and the k -mer size. The results are shown in Figures 3 and 4.



(a) Running time versus number of threads.

(b) Memory usage versus number of threads.

Figure 3 – Influence of the number of threads when filtering 2.6 Gbp of reads with 1M patterns of size $k = 31$.



(a) Running time versus k .

(b) Memory usage versus k .

Figure 4 – Influence of the k -mer size when filtering 2.6 Gbp of reads with 1M patterns and a single thread. Cleanifier is not included because it does not support $k > 32$.

When increasing the number of threads, K2Rmini remains the fastest method throughout the benchmark. Most of the speedup is obtained within the first four threads, after which the running time stabilizes around 1.3–1.4 seconds. Deacon shows a very similar trend, with slightly higher running times and the same early saturation. This suggests that, for both minimizer-based approaches, the main bottleneck quickly shifts away from the matching procedure itself toward parsing and other shared overheads. In contrast, BackToSequences and SBWT continue to benefit from additional threads over the full range, with their running times decreasing steadily from about 100 seconds on one thread to around 6 seconds on 16 threads. Cleanifier shows almost no parallel speedup in this experiment, remaining close to 35 seconds regardless of the thread count.

The memory trends differ markedly from the time trends. K2Rmini has the lowest memory footprint and remains nearly constant as the number of threads increases, staying around 8–10 MB. SBWT also remains essentially flat, at about 8 MB. BackToSequences has a larger but still stable footprint, around 50 MB, independent of the number of threads. By contrast, Deacon and Cleanifier both show a substantial increase in memory usage with additional threads. Deacon grows roughly linearly from about 15 MB on one thread to around 180 MB on 16 threads, while Cleanifier increases from about 110 MB to about 195 MB. Therefore, among the tested methods, K2Rmini combines the best running time with the smallest and most stable memory usage in the multithreaded setting.

We also evaluated the effect of the k -mer size with a single thread and 1M queried patterns. As k increases, K2Rmini becomes slightly faster, decreasing from about 4 seconds at $k = 31$ to about 2.3 seconds for the largest tested values. This trend is consistent with our minimizer-based design. Keeping the minimizer size fixed while increasing k enlarges the minimizer window, reducing minimizer density, and therefore decreasing the number of lookups required during the first filtering pass. Deacon exhibits a similar, though less pronounced, decrease in running time. In contrast, BackToSequences becomes slower as k increases, rising from about 80 seconds to more than 120 seconds, while SBWT remains roughly constant around 90 seconds. These results again highlight that minimizer-based filtering benefits from larger k -mer sizes, whereas exact indexed lookup of every k -mer does not.

The memory usage as a function of k remains moderate for all methods, but the trends are again distinct. K2Rmini stays almost constant around 7–8 MB across the full range of tested values. Deacon remains below 20 MB, with only a small increase as k grows. SBWT increases more noticeably, specifically after each power of two (32 and 64 in the plot), from about 7–8 MB at small k to about 20 MB for the largest tested values, especially on negative queries. BackToSequences has the highest memory usage throughout this experiment, steadily increasing from about 50 MB to more than 60 MB. Overall, these experiments confirm that K2Rmini is not only the fastest of the exact methods considered here, but also the most memory-efficient and the least sensitive to both thread count and k -mer size.

Real data

To complement the synthetic benchmarks, we compared BackToSequences and K2Rmini on several real datasets with heterogeneous sequence lengths and structures, summarized in [Table 1](#). The evaluation includes Oxford Nanopore ultra-long reads (SRR23365080), PacBio HiFi reads of 10–20 kb (SRX7897685–8), and Illumina 250 bp reads (ERR3239454). We also evaluated assembled sequences from the Logan project (Chikhi et al., 2024), including both contigs and unitigs (SRR7853572), and the complete human T2T reference genome (GCF_009914755.1), which represents very long contiguous sequences. These experiments were run on a laptop equipped with a 13th Gen Intel Core i9-13950HX processor, 64 GB of RAM, and an SSD, using 8 threads. For the positive case, the queried patterns were sequences of length 1000 selected from each dataset, using 100 patterns for short reads and unitigs so as to total 10^6 queried bases. For the negative case, the queried patterns were random sequences of length 1000. Across all datasets and in both positive and negative settings, K2Rmini is consistently faster than BackToSequences in both CPU time and wall-clock time.

The largest gains are observed on ONT and HiFi reads. On ONT data, K2Rmini is $10.55\times$ faster than BackToSequences in CPU time and $4.94\times$ faster in elapsed time for positive queries,

while the speedup increases to 17.74 \times and 16.09 \times on negative queries. On HiFi reads, the gains are even larger for positive queries, reaching 26.94 \times in CPU time and 27.41 \times in elapsed time, and remain very high on negative queries with speedups of 17.23 \times and 16.55 \times . On Illumina reads, the improvement is more moderate but still substantial, ranging from 5.53 \times to 5.98 \times in CPU time and from 3.97 \times to 4.48 \times in elapsed time. Overall, these results confirm that the minimizer-based prefilter is especially effective on raw reads, with the strongest gains on long-read datasets.

The same trend holds on assembled and reference sequences, although with more variable gains. On the T2T genome, K2Rmini is 4.00 \times faster in CPU time and 4.19 \times faster in elapsed time for positive queries, while negative queries lead to much larger speedups of 15.19 \times and 11.10 \times . On Logan contigs, the gains are stable between positive and negative cases, around 7.3 \times to 7.4 \times in CPU time and 5.8 \times to 5.9 \times in elapsed time. Logan unitigs show the smallest difference, with speedups around 2.5 \times in CPU time and 1.9 \times in elapsed time in both settings. More generally, negative queries tend to benefit more from K2Rmini than positive ones, which is consistent with the minimizer filter rejecting most sequences early and therefore avoiding the exact counting phase. Overall, K2Rmini consistently outperforms BackToSequences across all tested datasets, with the strongest improvements on long reads and on negative queries.

Table 1 – Runtime comparison between K2Rmini and BackToSequences (BTS) using 8 threads.

Experiment	CPU time (s)			Elapsed time (s)		
	K2Rmini	BTS	Speedup	K2Rmini	BTS	Speedup
ONT positive	1850.65	19529.04	$\times 10.55$	531.83	2626.34	$\times 4.94$
ONT negative	959.72	17021.16	$\times 17.74$	125.13	2013.64	$\times 16.09$
HiFi positive	279.46	7527.53	$\times 26.94$	33.11	907.56	$\times 27.41$
HiFi negative	273.75	4716.76	$\times 17.23$	33.70	557.88	$\times 16.55$
Illumina positive	278.24	1537.61	$\times 5.53$	43.56	173.02	$\times 3.97$
Illumina negative	267.35	1597.82	$\times 5.98$	40.41	180.89	$\times 4.48$
T2T positive	55.85	223.67	$\times 4.00$	8.88	37.19	$\times 4.19$
T2T negative	9.81	148.98	$\times 15.19$	2.31	25.63	$\times 11.10$
Contigs positive	14.05	102.10	$\times 7.27$	2.18	12.66	$\times 5.81$
Contigs negative	13.83	102.29	$\times 7.40$	2.15	12.70	$\times 5.91$
Unitigs positive	76.56	191.45	$\times 2.50$	13.87	26.49	$\times 1.91$
Unitigs negative	75.29	192.07	$\times 2.55$	14.07	26.65	$\times 1.89$

Discussion

This work has demonstrated that a search strategy based on SIMD-accelerated minimizer filtering is fast and resource-efficient across a broad spectrum of applications. Similar tools are being developed for screening sequence repositories for antimicrobial resistance mutations, emerging pathogen surveillance, or sequencing contaminant filtration (Constantinides et al., 2025). Currently, the main limitation of our implementation is that indexing the patterns is single-threaded and relies on a generic hash table. One promising direction to improve the indexing step would be to implement a concurrent hash table that relies on minimizers for grouping keys, such as the

one recently described in kache-hash (Khan et al., 2026). This would be especially relevant since we already have efficient methods to compute both minimizers and k -mer hashes, and we could thus easily batch insertions and queries. A second direction for improvement would be to parallelize the parsing step by attributing independent chunks of the input file to each thread, which would fully benefit from our vectorized parsing library and reduce the communication overhead.

Future work could explore alternative search strategies, including using sparser minimizer schemes to further reduce the number of minimizers to query (Alanko et al., 2025; Golan et al., 2025; Groot Koerkamp et al., 2025; Groot Koerkamp and Pibiri, 2024; Pellow et al., 2023). We could also imagine a middle ground between the costly exact search used in the second pass and doing no search at all at the cost of false positives, possibly by using small filters for faster queries with a lower false positive rate.

Finally, as the throughput of the core algorithm increases, performance bottlenecks shift to other components of the pipeline, mostly I/O-bound. While uncompressed FASTA/Q files are unnecessarily wasteful, reading input from compressed FASTA/Q files is a major bottleneck, highlighting the need for more performant alternatives (Patro et al., 2025; Teysier and Dobin, 2025).

Acknowledgements

Preprint version 2 of this article has been peer-reviewed and recommended by Peer Community In Mathematical and Computational Biology (<https://doi.org/10.24072/pci.mcb.100920>; Pibiri, 2026).

Funding

This work is funded by the French National Research Agency AGATE ANR-21-CE45-0012 and SxC ANR-24-CE25-2874-01. BC is supported by Wellcome Trust Award 313694/Z/24/Z.

Conflict of interest disclosure

The authors declare they comply with the PCI rule of having no financial conflicts of interest.

Data, script, code, and supplementary information availability

K2Rmini is available at <https://doi.org/10.5281/zenodo.20715935> (Martayan and Limasset, 2026a). Benchmarking scripts and wrappers are available at <https://doi.org/10.5281/zenodo.20715943> (Martayan and Limasset, 2026b). Public sequencing datasets used in the experiments are identified in the text.

References

- Agret C, Cazaux B, Limasset A (2021). Toward optimal fingerprint indexing for large scale genomics. *bioRxiv*, 2021.11.04.467355. <https://doi.org/10.1101/2021.11.04.467355>.
- Alanko JN, Biagi E, Puglisi SJ (2025). Finimizers: Variable-Length Bounded-Frequency Minimizers for k -mer Sets. *IEEE Transactions on Computational Biology and Bioinformatics* **22**, 899–910. <https://doi.org/10.1109/tcbbio.2025.3545285>.

- Alanko JN, Puglisi SJ, Vuohtoniemi J (2023). Small Searchable κ -Spectra via Subset Rank Queries on the Spectral Burrows-Wheeler Transform. In: *SIAM Conference on Applied and Computational Discrete Algorithms (ACDA23)*. Society for Industrial and Applied Mathematics, pp. 225–236. <https://doi.org/10.1137/1.9781611977714.20>.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990). Basic local alignment search tool. *Journal of molecular biology* **215**, 403–410. [https://doi.org/10.1016/s0022-2836\(05\)80360-2](https://doi.org/10.1016/s0022-2836(05)80360-2).
- Baire A, Marijon P, Andrace F, Peterlongo P (2024). Back to sequences: Find the origin of k-mers. *Journal of Open Source Software* **9**, 7066. <https://doi.org/10.21105/joss.07066>.
- Chikhi R, Raffestin B, Korobeynikov A, Edgar R, Babaian A (2024). Logan: planetary-scale genome assembly surveys life's diversity. *bioRxiv*, 2024.07.30.605881. <https://doi.org/10.1101/2024.07.30.605881>.
- Constantinides B, Lees J, Crook DW (2025). Deacon: fast sequence filtering and contaminant depletion. *bioRxiv*, 2025.06.09.658732. <https://doi.org/10.1101/2025.06.09.658732>.
- Crochemore M, Czumaj A, Gasieniec L, Lecroq T, Plandowski W, Rytter W (1999). Fast practical multi-pattern matching. *Information Processing Letters* **71**, 107–113. [https://doi.org/10.1016/S0020-0190\(99\)00092-7](https://doi.org/10.1016/S0020-0190(99)00092-7).
- Crosbie ND (2025). grepq: A Rust application that quickly filters FASTQ files by matching sequences to a set of regular expressions. *Journal of Open Source Software* **10**, 8048. <https://doi.org/10.21105/joss.08048>.
- Darvish M, Seiler E, Mehninger S, Rahn R, Reinert K (2022). Needle: a fast and space-efficient prefilter for estimating the quantification of very large collections of expression experiments. *Bioinformatics* **38**, 4100–4108. <https://doi.org/10.1093/bioinformatics/btac492>.
- Edgar RC, Taylor B, Lin V, Altman T, Barbera P, Meleshko D, Lohr D, Novakovsky G, Buchfink B, Al-Shayeb B, et al. (2022). Petabase-scale sequence alignment catalyses viral discovery. *Nature* **602**, 142–147. <https://doi.org/10.1038/s41586-021-04332-2>.
- Ekim B, Berger B, Chikhi R (2021). Minimizer-space de Bruijn graphs: Whole-genome assembly of long reads in minutes on a personal computer. *Cell systems* **12**, 958–968. <https://doi.org/10.1016/j.cels.2021.08.009>.
- Faro S, Lecroq T (2013). The exact online string matching problem: A review of the most recent results. *ACM Computing Surveys (CSUR)* **45**, 1–42. <https://doi.org/10.1145/2431211.2431212>.
- Gallant A (2025). ripgrep. GitHub repository, version 15.1.0, released 22 October 2025, accessed 15 June 2026, <https://github.com/BurntSushi/ripgrep/releases/tag/15.1.0>.
- Golan S, Tziony I, Kraus M, Orenstein Y, Shur A (2025). GreedyMini: generating low-density DNA minimizers. *Bioinformatics* **41**, i275–i284. <https://doi.org/10.1093/bioinformatics/btaf251>.
- Groot Koerkamp R, Liu D, Pibiri GE (2025). The open-closed mod-minimizer algorithm. *Algorithms for Molecular Biology* **20**, 4. <https://doi.org/10.1186/s13015-025-00270-0>.
- Groot Koerkamp R, Martayan I (2025). SimdMinimizers: Computing Random Minimizers, fast. In: *23rd International Symposium on Experimental Algorithms (SEA 2025)*. Ed. by Petra Mutzel and Nicola Prezza. Vol. 338. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 20:1–20:19. <https://doi.org/10.4230/LIPIcs.SEA.2025.20>.

- Groot Koerkamp R, Pibiri GE (2024). The mod-minimizer: A Simple and Efficient Sampling Algorithm for Long k-mers. In: *24th International Workshop on Algorithms in Bioinformatics (WABI 2024)*. Ed. by Solon P. Pissis and Wing-Kin Sung. Vol. 312. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 11:1–11:23. <https://doi.org/10.4230/LIPIcs.WABI.2024.11>.
- Homer N, Stadick S, Lambert S, Stone M, Fennell T (2025). fulcrumgenomics/fqgrep: v1.1.1. <https://github.com/fulcrumgenomics/fqgrep>. Version 1.1.1. <https://doi.org/10.5281/zenodo.15034074>.
- Karasikov M, Mustafa H, Rättsch G, Kahles A (2022). Lossless indexing with counting de Bruijn graphs. *Genome Research* **32**, 1754–1764. <https://doi.org/10.1101/gr.276607.122>.
- Khan J, Patro R, Pandey P (2026). cache-hash: A dynamic, concurrent, and cache-efficient hash table for streaming k-mer operations. *bioRxiv*, 2026.02.13.705625. <https://doi.org/10.64898/2026.02.13.705625>.
- Lemane T, Lezsoche N, Lecubin J, Pelletier E, Lescot M, Chikhi R, Peterlongo P (2024). Indexing and real-time user-friendly queries in terabyte-sized complex genomic datasets with kmindex and ORA. *Nature Computational Science* **4**, 104–109. <https://doi.org/10.1038/s43588-024-00596-6>.
- Li H (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>.
- Ma B, Lu C, Wang Y, Yu J, Zhao K, Xue R, Ren H, Lv X, Pan R, Zhang J, et al. (2023). A genomic catalogue of soil microbiomes boosts mining of biodiversity and genetic resources. *Nature communications* **14**, 7318. <https://doi.org/10.1038/s41467-023-43000-z>.
- Mäklin T, Alanko JN, Biagi E, Puglisi SJ (2025). Sequence alignment with k-bounded matching statistics. *bioRxiv*, 2025.05.19.654936. <https://doi.org/10.1101/2025.05.19.654936>.
- Marchet C, Boucher C, Puglisi SJ, Medvedev P, Salson M, Chikhi R (2021). Data structures based on k-mers for querying large collections of sequencing data sets. *Genome research* **31**, 1–12. <https://doi.org/10.1101/gr.260604.119>.
- Marchet C, Limasset A (2023). Scalable sequence database search using partitioned aggregated Bloom comb trees. *Bioinformatics* **39**, i252–i259. <https://doi.org/10.1093/bioinformatics/btad225>.
- Martayan I, Limasset A (2026a). Malfloy/K2Rmini: Peer Community Journal version. <https://doi.org/10.5281/zenodo.20715935>.
- Martayan I, Limasset A (2026b). Malfloy/K2Rmini_experiments: Peer Community Journal version. <https://doi.org/10.5281/zenodo.20715943>.
- Martayan I, Lobet L, Marchet C, Paperman C (2026). Helicase: Vectorized parsing and bitpacking of genomic sequences. <https://doi.org/10.64898/2026.03.19.712912>.
- Mohamadi H, Chu J, Vandervalk BP, Birol I (2016). ntHash: recursive nucleotide hashing. *Bioinformatics* **32**, 3492–3494. <https://doi.org/10.1093/bioinformatics/btw397>.
- Nayfach S, Shi ZJ, Seshadri R, Pollard KS, Kyrpides NC (2019). New insights from uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505–510. <https://doi.org/10.1038/s41586-019-1058-x>.

- Parks DH, Rinke C, Chuvochina M, Chaumeil PA, Woodcroft BJ, Evans PN, Hugenholtz P, Tyson GW (2017). Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nature microbiology* **2**, 1533–1542. <https://doi.org/10.1038/s41564-017-0012-7>.
- Patro R, Bharti S, Singhanian P, Dhakal R, Dahlstrom TJ, Groot Koerkamp R (2025). mim: A light-weight auxiliary index to enable fast, parallel, gzipped FASTQ parsing. *bioRxiv*, 2025.11.24.690271. <https://doi.org/10.1101/2025.11.24.690271>.
- Pellow D, Pu L, Ekim B, Kotlar L, Berger B, Shamir R, Orenstein Y (2023). Efficient minimizer orders for large values of k using minimum decycling sets. *Genome Research* **33**, 1154–1161. <https://doi.org/10.1101/gr.277644.123>.
- Pibiri GE (2026). A lossless shortcut for k-mer-based sequence filtering. *Peer Community in Mathematical and Computational Biology*, 100920. <https://doi.org/10.24072/pci.mcb.100920>.
- Rahman Hera M, Koslicki D (2025). Estimating similarity and distance using FracMinHash. *Algorithms for Molecular Biology* **20**, 1–13. <https://doi.org/10.1186/s13015-025-00276-8>.
- Roberts M, Hayes W, Hunt BR, Mount SM, Yorke JA (2004). Reducing storage requirements for biological sequence comparison. *Bioinformatics* **20**, 3363–3369. <https://doi.org/10.1093/bioinformatics/bth408>.
- Sayers EW, Cavanaugh M, Clark K, Pruitt KD, Sherry ST, Yankie L, Karsch-Mizrachi I (2024). GenBank 2024 update. *Nucleic Acids Research* **52**, D134–D137. <https://doi.org/10.1093/nar/gkad903>.
- Schleimer S, Wilkerson DS, Aiken A (2003). Winnowing: local algorithms for document fingerprinting. In: *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*. SIGMOD '03. New York, NY, USA: Association for Computing Machinery, pp. 76–85. <https://doi.org/10.1145/872757.872770>. (Visited on 10/18/2024).
- Shen W, Le S, Li Y, Hu F (2016). SeqKit: a cross-platform and ultrafast toolkit for FASTA/Q file manipulation. *PloS one* **11**, e0163962. <https://doi.org/10.1371/journal.pone.0163962>.
- Shen W, Sipos B, Zhao L (2024). SeqKit2: A Swiss army knife for sequence and alignment processing. *Imeta* **3**, e191. <https://doi.org/10.1002/imt2.191>.
- Teyssier N, Dobin A (2025). BINSEQ: A Family of High-Performance Binary Formats for Nucleotide Sequences. *bioRxiv*, 2025.04.08.647863. <https://doi.org/10.1101/2025.04.08.647863>.
- Vandamme L, Cazaux B, Limasset A (2025). K2R: Tinted de Bruijn Graphs implementation for efficient read extraction from sequencing datasets. *Bioinformatics Advances* **5**, vbaf111. <https://doi.org/10.1093/bioadv/vbaf111>.
- Wang X, Hong Y, Chang H, Park K, Langdale G, Hu J, Zhu H (2019). Hyperscan: A Fast Multi-pattern Regex Matcher for Modern CPUs. In: *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*. USENIX Association, pp. 631–648. url: <https://www.usenix.org/conference/nsdi19/presentation/wang-xiang>.
- Zentgraf J, Schmitz JE, Rahmann S (2025). Cleanifier: contamination removal from microbial sequences using spaced seeds of a human pangenome index. *Bioinformatics* **42**, btaf632. <https://doi.org/10.1093/bioinformatics/btaf632>.
- Zielezinski A, Vinga S, Almeida J, Karlowski WM (2017). Alignment-free sequence comparison: benefits, applications, and tools. *Genome biology* **18**, 1–17. <https://doi.org/10.1186/s13059-017-1319-7>.