

# Minmers are a generalization of minimizers that enable unbiased local Jaccard estimation

Bryce Kille<sup>1,\*</sup>, Erik Garrison<sup>2</sup>, Todd J Treangen<sup>1</sup>, and Adam M Phillippy<sup>3,\*</sup>

<sup>1</sup>Department of Computer Science, Rice University, Houston, TX, USA

<sup>2</sup>Department of Genetics, Genomics and Informatics, University of Tennessee Health Science Center, Memphis, TN, USA

<sup>3</sup>Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA

\*To whom correspondence should be addressed.

## 1 Abstract

**Motivation:** The Jaccard similarity on  $k$ -mer sets has shown to be a convenient proxy for sequence identity. By avoiding expensive base-level alignments and comparing reduced sequence representations, tools such as MashMap can scale to massive numbers of pairwise comparisons while still providing useful similarity estimates. However, due to their reliance on minimizer winnowing, previous versions of MashMap were shown to be biased and inconsistent estimators of Jaccard similarity. This directly impacts downstream tools that rely on the accuracy of these estimates.

**Results:** To address this, we propose the *minmer* winnowing scheme, which generalizes the minimizer scheme by use of a rolling minhash with multiple sampled  $k$ -mers per window. We show both theoretically and empirically that minmers yield an unbiased estimator of local Jaccard similarity, and we implement this scheme in an updated version of MashMap. The minmer-based implementation is over 10 times faster than the minimizer-based version under the default ANI threshold, making it well-suited for large-scale comparative genomics applications.

**Availability:** MashMap3 is available at <https://github.com/marbl/MashMap>

**Contact:** [blk6@rice.edu](mailto:blk6@rice.edu), [adam.phillippy@nih.gov](mailto:adam.phillippy@nih.gov)

## 2 Introduction

The recent deluge of genomic data accelerated by population-scale long-read sequencing efforts has driven an urgent need for scalable long-read mapping and comparative genomics algorithms. The completion of the first Telomere-to-Telomere (T2T) human genome Nurk *et al.* (2022) and the launch of the Human Pangenome Project Wang *et al.* (2022a) have paved the way to mapping genomic diversity at unprecedented scale and resolution. A key goal when comparing a newly sequenced human genome to a reference genome or pangenome is to accurately identify homologous sequences, that is, DNA sequences that share a common evolutionary source.

Algorithms for pairwise sequence alignment, which aim to accurately identify homologous regions between two sequences, have continued to advance in recent years Marco-Sola *et al.* (2021). While a powerful and

ubiquitous computational tool in computational biology, exact alignment algorithms are typically reserved for situations where the boundaries of homology are known *a priori*, due to their quadratic runtime costs and inability to model nonlinear sequence relationships such as inversions, translocations, and copy number variants. Because of this, long-read mapping or whole-genome alignment methods must first identify homologous regions across billions of nucleotides, after which the exact methods can be deployed to compute a base-level “gapped” read alignment for each region. To efficiently identify candidate mappings, the prevailing strategy is to first sample  $k$ -mers and then identify consecutive  $k$ -mers that appear in the same order for both sequences: known as “seeding” and “chaining”, respectively.

For many use cases, an exact gapped alignment is not needed and only an estimate of sequence identity is required. As a result, methods have been developed which can predict sequence identity without the cost of computing a gapped alignment. Jaccard similarity, a metric used for comparing the similarity of two sets, has found widespread use for this task, especially when combined with locality sensitive hashing of  $k$ -mer sets Ondov *et al.* (2016); Brown and Irber (2016); Ondov *et al.* (2019); Jain *et al.* (2017, 2018a); Baker and Langmead (2019); Shaw and Yu (2023). By comparing only  $k$ -mers, the Jaccard can be used to estimate the average nucleotide identity (ANI) of two sequences without the need for an exact alignment Ondov *et al.* (2016, 2019); Blanca *et al.* (2022).

To accelerate mapping and alignment,  $k$ -mers from the input sequences are often down-sampled using a “winnowing scheme” in a way that reduces the input size while still enabling meaningful comparisons. For example, both MashMap Jain *et al.* (2017, 2018a) and Minimap Li (2018) use a minimizer scheme Roberts *et al.* (2004), which selects only the *smallest*  $k$ -mer from all  $w$ -length substrings of the genome. Of relevance to this study, MashMap2 then uses these minimizers to approximate the Jaccard similarity between the mapped sequences, and these estimates have been successfully used by downstream methods such as FastANI Jain *et al.* (2018b) and MetaMaps Diltney *et al.* (2019).

However, a recent investigation noted limitations of the “winnowed minhash” scheme introduced by MashMap Belbasi *et al.* (2022). Although the origi-

92 nal MashMap paper notes a small, but negligible bias  
93 in its estimates Jain *et al.* (2017), Belbasi *et al.* proved  
94 that no matter the length of the sequences, the bias  
95 of the minimizer-based winnowed minhash estimator is  
96 never zero Belbasi *et al.* (2022).

97 To address this limitation, we propose a novel win-  
98 nowing scheme, the “minmer” scheme, which is a gen-  
99 eralization of minimizers that allows for the selection  
100 of multiple  $k$ -mers per window. We define this scheme,  
101 characterize its properties, and provide an implementa-  
102 tion in MashMap3. Importantly, we show that minmers,  
103 unlike minimizers, enable an unbiased prediction of the  
104 local Jaccard similarity.

### 105 3 Preliminaries

106 Let  $\Sigma$  be an alphabet and  $\mathcal{S}_k(S) : \Sigma^+ \rightarrow \{\Sigma^k\}^+$  be  
107 a function which maps a sequence  $S$  to the set of all  
108  $k$ -mers in  $S$ . Similarly, given a sequence  $S$ , we define  
109  $W_i^{(w)}(S)$  as the sequence of  $w$   $k$ -mers in  $S$  starting at  
110 the  $i$ th  $k$ -mer. When  $w$  and  $S$  are clear from context,  
111 we use  $W_i$ . We use the terms sequence and string  
112 interchangeably.

#### 113 3.1 Jaccard similarity and the minhash 114 approximation

115 Given two sets  $A$  and  $B$ , their Jaccard similarity is  
116 defined as  $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ . The Jaccard similarity  
117 between two sequences  $R$  and  $Q$  can be computed as  
118  $J(\mathcal{S}_k(R), \mathcal{S}_k(Q))$  for some  $k$ -mer size  $k$ .

However, computing the exact Jaccard for  $\mathcal{S}_k(R)$   
and  $\mathcal{S}_k(Q)$  is not an efficient method for determining  
similarity for long reads and whole genomes. Instead,  
the minhash algorithm provides an estimator for the  
Jaccard similarity while only needing to compare a  
fraction of the two sets. Assuming  $U$  is the universe  
of all possible elements and  $\pi : U \rightarrow |U|$  is a function  
which imposes a randomized total order on the universe  
of elements, we have that

$$J(A, B) = \Pr(\min_{x \in A}(\pi(x)) = \min_{x \in B}(\pi(x)))$$

119 This equivalency, proven by Broder (1997), is key to  
120 the minhash algorithm and yields an unbiased and con-  
121 sistent Jaccard estimator  $\hat{J}$  with the help of a sketching  
122 function  $\pi_s$ . Let  $\pi_s$  return the lowest  $s$  items from the  
123 input set according to the random total order  $\pi$ . Then  
124 we define the minhash as

$$\hat{J}(A, B) = \frac{|\pi_s(A \cup B) \cap \pi_s(A) \cap \pi_s(B)|}{|\pi_s(A \cup B)|}$$

125 Importantly, this Jaccard estimator has an expected  
126 error that scales with  $\mathcal{O}(1/\sqrt{s})$  and is therefore inde-  
127 pendent of the size of the original input sets. While  
128 there are a number of variants of minhash which pro-  
129 vide the same guarantee Cohen (2016), we will be using  
130 the “bottom- $s$  sketch” (as opposed to the  $s$ -mins and  
131  $s$ -partition sketch) since it ensures a consistent sketch  
132 size regardless of the parameters and requires only a  
133 single hash computation per element of  $\mathcal{S}_k$ . Addition-  
134 ally, the simplicity of the bottom- $s$  sketch leads to a  
135 streamlined application of the sliding window model,  
136 which we describe next.

### 137 3.2 Winnowing

138 While sequences can be reduced into their correspond-  
139 ing sketch via the method described above, this is a  
140 *global* sketch and it is difficult to determine *where* two  
141 sequences share similarity. In order to perform local  
142 mapping, Schleimer *et al.* (2003) and Roberts *et al.*  
143 (2004) independently introduced the concept of *win-*  
144 *nowing* and *minimizers*. In short, given some total  
145 ordering on the  $k$ -mers, a window of length  $w$  is slid  
146 over the sequence and the element with the lowest rank  
147 in each window (the *minimizer*) is selected, using the  
148 left-most position to break ties Roberts *et al.* (2004).  
149 By definition, winnowing ensures that at least one el-  
150 ement is sampled per window and therefore there is  
151 never a gap of more than  $w$  elements between sampled  
152 positions. Here, we extend the winnowing concept to al-  
153 low the selection of more than one element per window  
154 (the *minmers*), and we refer to the set of all minmers  
155 and/or their positions as the *winnowed* sequence.

#### 156 3.2.1 Winnowing scheme characteristics

**Definition 3.1.** *A winnowing scheme has a  $(w, s)$ -*  
*window guarantee if for every window of  $w$   $k$ -mers,*  
*there are at least  $\max(\#_{distinct}, s)$   $k$ -mers sampled from*  
*the window, where  $\#_{distinct}$  is the number of distinct*  
 *$k$ -mers in the window.*

162 This definition is more general than the commonly  
163 used  $w$ -window guarantee, which is equivalent to the  
164  $(w, 1)$ -window guarantee. While not all winnowing  
165 schemes must have such a guarantee, this ensures that  
166 no area of the sequence is under-sampled. Shaw and  
167 Yu (2022) recently provided an analytical framework  
168 for winnowing schemes and showed that mapping sen-  
169 sitivity is related to the distribution of distances (or  
170 *spread*) between sampled positions, and precision is  
171 related to the proportion of unique values relative to  
172 the total number of sampled positions. As the over-  
173 arching goal of winnowing is to reduce the size of the  
174 input while preserving as much information as possi-  
175 ble, winnowing schemes typically aim to optimize the  
176 precision/sensitivity metrics given a particular density.

**Definition 3.2.** *The density  $d$  of a winnowing scheme*  
*is defined as the expected frequency of sampled positions*  
*from a long random string, and the density factor  $d_f$  is*  
*defined as the expected number of sampled positions in*  
*a window of  $w + 1$   $k$ -mers.*

182 There has been significant work on improving the  
183 performance of minimizers by identifying orderings that  
184 reduce the density factor Marçais *et al.* (2017). Mini-  
185 mizer schemes which use a uniformly random ordering  
186 have a density factor of  $d_f = 2$  and recent schemes like  
187 Miniception Zheng *et al.* (2020) and PASHA Ekim *et al.*  
188 (2020) are able to obtain density factors as low as 1.7  
189 for certain values of  $w$  and  $k$ .

190 For the remainder of this work, we will assume that  
191  $w \ll 4^k$ , i.e. the windows are not so large that we expect  
192 duplicate  $k$ -mers in a random string. This ensures that  
193 each  $k$ -mer in a window has probability  $s/w$  of being  
194 in the sketch for that window.

### 3.2.2 Winnowing scheme hierarchies

Recent winnowing methods have focused on schemes that select at most a single position per window, which simplifies analyses but restricts the universe of possible schemes. Minimizers belong to the class of *forward* winnowing schemes, where the sequence of positions sampled from adjacent sliding windows is non-decreasing Marçais *et al.* (2018). More general is the concept of a  $w$ -local scheme Shaw and Yu (2022), defined on windows of  $w$  consecutive  $k$ -mers but without the forward requirement. Non-forward schemes are more powerful and are not limited by the same density factor bounds as forward schemes. While the need of non-forward schemes to “jump back” in order to obtain lower sampling densities is acknowledged by Marçais *et al.* (2018), there are currently no well-studied, non-forward,  $w$ -local schemes.

## 3.3 MashMap

MashMap is a minimizer-based tool for long-read and whole-genome sequence homology mapping that is designed to identify all pairwise regions above some sequence similarity cutoff Jain *et al.* (2017, 2018a). Specifically, for a reference sequence  $R$  and a query sequence  $Q$  comprised of  $w$   $k$ -mers, MashMap aims to find all positions  $i$  in the reference such that  $J(A, B_i) \geq c$ , where  $A = \mathcal{S}_k(Q)$  and  $B_i = W_i^{(w)}(R)$ , and  $c$  is the sequence similarity cutoff. For ease of notation, we will use  $B$  to refer to the sequence of  $k$ -mers from the reference sequence  $R$ . Importantly, MashMap only requires users to specify a minimum segment length and minimum sequence identity threshold, and the algorithm will automatically determine the parameters needed to return all mappings that meet this criteria with parameterized confidence under a binomial mutation model.

Here we replace the minimizer-based approach of prior versions of MashMap with minmers. While the problem formulation remains the same, our method for computing the reference index and filtering candidate mappings is novel. We will first introduce the concept of minmers, which enable winnowing the input sequences while still maintaining the  $k$ -mers necessary to compute an unbiased Jaccard estimation between any two windows of length at least  $w$ . We will then discuss the construction of the reference index and show how query sequences can be efficiently mapped to the reference such that their expected ANI is above the desired threshold.

## 4 The minmer winnowing scheme

Minmers are a generalization of minimizers that allow for the selection of more than one minimum value per window. The relationship between minmers and minimizers was noted by Berlin *et al.* (2015) but as a global sketch and without the use of a sliding window. Here we formalize a definition of the minmer winnowing scheme.

**Definition 4.1.** Given a tuple  $(w, s, k, \pi)$ , where  $w, k$

and  $s$  are integers and  $\pi$  is an ordering on the set of all  $k$ -mers, a  $k$ -mer in a sequence is a minmer if it is one of the smallest  $s$   $k$ -mers in any of the subsuming windows of  $w$   $k$ -mers.

Similar to other  $w$ -local winnowing schemes, ties between  $k$ -mers are broken by giving priority to the leftmost  $k$ -mer. From the definition, it follows that by letting  $s = 1$  we obtain the definition of the minimizer scheme. Compared to minimizers with the same  $w$  value, minmers guarantee that at least  $s$   $k$ -mers will be sampled from each window. However, as a non-forward scheme, a minmer may be one of the smallest  $s$   $k$ -mers in two non-adjacent windows, yet not one of the smallest  $s$   $k$ -mers in an intervening window (Figure 1). To account for this and simplify development of this scheme, we define a *minmer interval* to be the interval for which the  $k$ -mer at position  $i$  is a minmer for all windows starting within that interval. Thus, a single  $k$ -mer may have multiple minmer intervals starting at different positions.

**Definition 4.2.** A tuple  $(i, a, b)$  is a minmer interval for a sequence  $S$  if the  $k$ -mer at position  $i$  is a minmer for all windows  $W_j$  where  $j \in [a, b)$ , but not  $W_{a-1}$  or  $W_b$ .

Any window  $W_j$  may contain more than  $s$  minmers, and so to naively compute the Jaccard between a query and  $W_j$  would require identification of the  $s$  smallest  $k$ -mers in  $W_j$ . Minmer intervals are convenient because for any window start position  $j$ , the  $s$  smallest  $k$ -mers in  $W_j$  are simply the ones whose minmer intervals contain  $j$ . Thus, indexing  $S$  with minmer intervals enables the efficient retrieval of the smallest  $s$   $k$ -mers for any window without additional sorting or comparisons.

Another benefit of minmer intervals is that the smallest  $s$   $k$ -mers for any window of length  $w' > w$  are guaranteed to be a subset of the combined  $(w, s)$ -minmers contained in that window. This subset can be easily computed with minmer intervals, since the set of  $(w, s)$ -minmer intervals that overlap with the range  $[i, i + w' - w]$  are also guaranteed to include the  $s$  smallest  $k$ -mers of the larger window, and the overlapping minmer intervals can be inspected to quickly identify them.

### 4.1 Constructing the rolling minhash index

In this section, we will describe our rolling bottom- $s$  sketch algorithm for collecting minmers and their corresponding minmer intervals. Popic and Batzoglou (2017) proposed a related rolling minhash method for short-read mapping, but using an  $s$ -mins scheme without minmer intervals. For the remainder of the section, we will assume no duplicate  $k$ -mers in a window and an ideal uniform hash function which maps to  $[0, 1]$ . Duplicate  $k$ -mers are handled in practice by keeping a counter of the number of active positions for a particular  $k$ -mer, similar to the original MashMap implementation Jain *et al.* (2017). Minmer intervals longer than the window length sometimes arise due to duplicate  $k$ -mers and are

split into adjacent windows of length at most  $w$ . This bound on the minmer interval length is necessary for the mapping step.

For ease of notation, we now consider  $B$  as a sequence of  $k$ -mer hash values  $x_0, x_1, \dots, x_n$  where each  $x_i \in [0, 1]$  and refer to these elements as hashes and  $k$ -mers interchangeably. We use a min-heap  $H$  and a sorted map  $M$ , both ordered on the hash values, to keep track of the rolling minhash index. As the window slides across  $B$ ,  $M$  will contain the minmer intervals for the lowest  $s$  hashes in the window and  $H$  will contain the remaining hashes in the window. We denote the minmer interval of a hash  $x$  in  $M$  by  $M[x]^{(start)}$  and  $M[x]^{(end)}$ . In practice,  $H$  may contain “expired”  $k$ -mers which are no longer part of the current window, however by storing the  $k$ -mer position as well, we can immediately discard such  $k$ -mers whenever they appear at the top of the heap. To prevent expired  $k$ -mers from accumulating, all expired  $k$ -mers from the heap are pruned whenever the heap size exceeds  $2w$ . After initialization of  $H$  and  $M$  with the first  $w$   $k$ -mers of  $B$ , we begin sliding the window for each consecutive position  $i$  and collect the minmer intervals in an index  $I$ . For each window  $B_i$ , there will be a single “exiting”  $k$ -mer  $x_{i-1}$  and a single “entering”  $k$ -mer  $x_{i+w-1}$ , each of which may or may not belong to the lowest  $s$   $k$ -mers. Therefore, we have four possibilities, examples of which can be seen in Figure 1.

1.  $x_{i-1} > \max(M)$  and  $x_{i+w-1} > \max(M)$   
Neither the exiting nor entering  $k$ -mer is in the sketch. Insert  $x_{i+w-1}$  into  $H$ .
2.  $x_{i-1} > \max(M)$  and  $x_{i+w-1} < \max(M)$   
The exiting  $k$ -mer was not in the sketch, but the entering  $k$ -mer will be. Since the incoming  $k$ -mer  $x_{i+w-1}$  enters the sketch, the largest element in the sketch must be removed. Therefore,  $M[\max(M)]^{(end)}$  is set to  $i$  and the minmer interval is appended to the index  $I$ .  $\max(M)$  is then removed from  $M$  and the new  $k$ -mer  $x_{i+w-1}$  is inserted to  $M$ , marking  $M[x_{i+w-1}]^{(start)} = i$ .
3.  $x_{i-1} \leq \max(M)$  and  $x_{i+w-1} > \max(M)$   
The exiting  $k$ -mer was in the sketch, but the entering  $k$ -mer will not be. Since the exiting  $k$ -mer  $x_{i-1}$  was a member of the sketch, set  $M[x_{i-1}]^{(end)} = i$ , remove  $M[x_{i-1}]$  from  $M$  and append it to  $I$ , and insert  $x_{i+w-1}$  into  $H$ . At this point,  $|M| = s - 1$ , as we removed an element from the sketch but did not replace it. To fill the empty sketch position,  $k$ -mers are popped from  $H$  until a  $k$ -mer  $x$  which has not expired is obtained. This  $k$ -mer is added to  $M$ , setting  $M[x]^{(start)} = i$ .
4.  $x_{i-1} \leq \max(M)$  and  $x_{i+w-1} \leq \max(M)$   
Both the exiting and entering  $k$ -mers are in the sketch. As before, set  $M[x_{i-1}]^{(end)} = i$  and remove  $M[x_{i-1}]$  from  $M$  and append it to  $I$ . The entering  $k$ -mer belongs in the sketch, so set  $M[x_{i+w-1}]^{(start)} = i$ .

Our implementation of  $M$  uses a balanced binary tree and  $H$  is pruned in  $\mathcal{O}(w)$  time at most every  $w$   $k$ -mers and therefore the amortized time complexity of each sliding window update is  $\mathcal{O}(\log(w))$ . In order to efficiently use the index for mapping, we sort  $I$  based on the start positions of the minmers. In addition to  $I$ , we compute a reverse lookup table  $T$  which maps hash values to ordered lists of start and end points of minmer intervals for that hash value. Overall, the indexing time requires  $\mathcal{O}(n \log(w) + |I| \log(|I|))$ , where  $|I|$  is estimated to be  $1 - \binom{w-1}{s} / \binom{w+1}{s}$ , as shown in section 5.1.2.

## 4.2 Querying the rolling minhash index

MashMap computes mappings in a two-stage process. In the first stage, all regions within the reference that may contain a mapping satisfying the desired ANI constraints are obtained. In the second stage, the minhash algorithm is used to estimate the Jaccard for each candidate mapping position  $i$  produced by the first stage. As the second stage is the most computationally intensive step, we introduce both a new candidate region filter and a more efficient minhash computation to improve overall runtime. We assume here that query sequences are  $w$   $k$ -mers long. In practice, sequences longer than  $w$  are split into windows of  $w$   $k$ -mers, mapped independently, and then chained and filtered as described in Jain *et al.* (2018a).

### 4.2.1 Stage 1: Candidate region filter

First, the query sequence  $A$  is winnowed using a min-heap to obtain the  $s$  lowest hash values. All  $m$  minmer intervals in the reference with matching hashes are obtained from  $T$  and a sorted list  $L$  is created in  $\mathcal{O}(m \log(s))$  time, where  $L$  consists of all minmer start and end positions. In this way, we can iterate through the list and keep a running count of the overlapping minmer intervals by incrementing the count for each start-point and decrementing the count for each end-point.

Unlike the previous versions of MashMap that look for all mappings above a certain ANI threshold, MashMap3 provides the option to instead filter out all mappings which are not likely to be within  $\Delta_{ANI}$  of the best predicted mapping ANI. This significantly reduces the number and size of the candidate regions passed on to the more expensive second stage.

Let  $Y_i$  be a random variable representing the numerator of the minhash formula for  $A$  and  $B_i$ . Given  $c_i = |\pi_s(A) \cap \pi_s(B_i)|$ , we observe that  $Y_i$  is distributed hypergeometrically, where we have  $s$  success states in a population of  $2s - c_i$  states (proof in Supplementary Materials). Let  $z = \arg \max_i c_i$  be a position with the maximum intersection size over all  $B_i$ , i.e. the position in  $B$  that overlaps with the most selected minmer intervals. We can now find a minimum intersection size  $\tau$  such that for any  $c_i < \tau$ ,

$$\Pr(\hat{J}(A, B_i) > \hat{J}(A, B_z) - \Delta_J) < 1 - \delta$$

where  $\Delta_J$  is the difference in the Jaccard that corresponds to an ANI value  $\Delta_{ANI}$  less than the ANI value

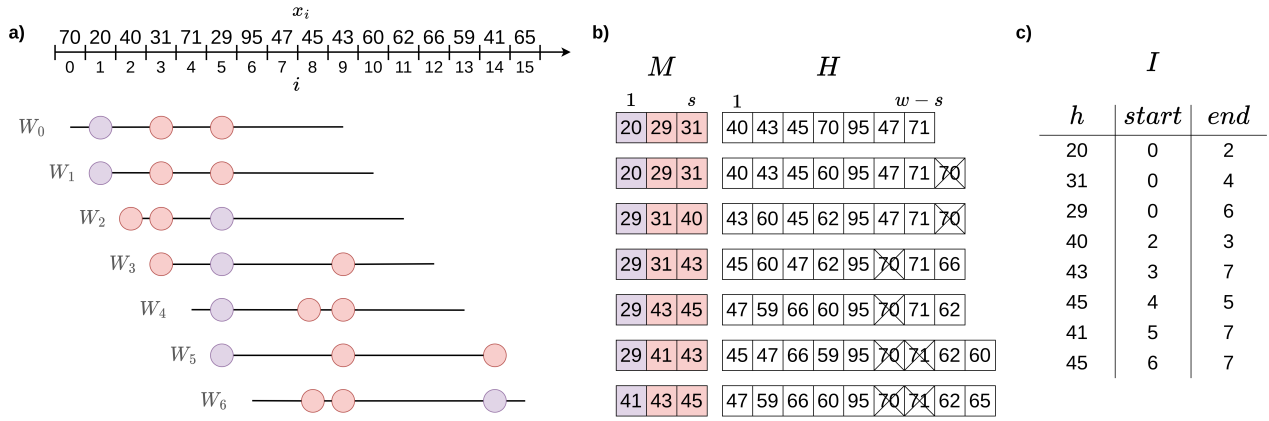


Figure 1: **Constructing the rolling minhash index.** (a) A sliding window  $B_i$  of length  $w = 10$  is moved over the hashes of all  $k$ -mers. At each position  $i$  of the sliding window, the positions with the  $s = 3$  lowest hash values are marked as minmers. The 3 minmers for each window are highlighted with colored circles, with the smallest hash in each window (the minimizer) highlighted in purple. (b) The values of the hashes in the map  $M$  and heap  $H$  as the window slides over the sequence. The expired  $k$ -mers in the heap are crossed out. (c) The final sorted minmer interval index  $I$ .

426 predicted by  $\hat{J}(A, B_z)$  and  $\delta$  is a desired confidence level.  
 427 To calculate this probability, we can use the following  
 428 summation

$$\Pr(\hat{J}(A, B_i) > \hat{J}(A, B_z) - \Delta_J) = \sum_{y=0}^s \Pr(Y_i = y | c_i) \Pr(Y_z < y + \Delta_J | c_z)$$

429 For each intersection size, we can identify a cutoff in  
 430  $\mathcal{O}(s \log(s))$  time. As a preprocessing step, we compute  
 431 cutoffs for each of the  $s$  possible intersection sizes at  
 432 the indexing stage. Candidate regions that are unlikely  
 433 to have an ANI within  $\Delta_{ANI}$  of the best predicted ANI  
 434 are then pruned. The default  $\Delta_{ANI}$  and  $\delta$  confidence  
 435 parameters of MashMap3 are 0 and 0.999, respectively,  
 436 as in many cases the lower scoring mappings for a  
 437 segment are filtered out by the plane-sweep filtering  
 438 method of MashMap described in Jain *et al.* (2018a).

439 We compute two passes over the interval endpoints  
 440 in  $L$ . In the first pass of stage 1, the maximum intersec-  
 441 tion size  $c_z$  is obtained. In the second pass, candidate  
 442 mappings whose intersection is above the cutoff derived  
 443 from  $c_z$  are obtained. Consecutive candidate mappings  
 444 are grouped into candidate regions and passed to stage  
 445 2.

#### 4.2.2 Stage 2: Efficiently computing the rolling minhash

446 Given a candidate region  $[a, z]$ , the goal of stage 2 is to  
 447 calculate the minhash for all  $A, B_i$  pairs for  $i \in [a, z]$ .  
 448 In order to track the minhash of  $A$  and  $B_i$  for each  $i$ ,  
 449 MashMap2 previously used a sorted map to track all  
 450 active seeds in each window. We improve upon this by  
 451 observing that the minhash can be efficiently tracked  
 452 using only  $\pi_s(A)$ ,  $\pi_s(A) \cap \pi_s(B_i)$ , and the number of  
 453 minmers from  $\pi_s(B_i)$  in-between each consecutive pair  
 454 of minmers from  $\pi_s(A)$ . To do so, MashMap3 uses  
 455 an array  $V = (-1, 0, 0), (x_1, \alpha_1, \beta_1), (x_2, \alpha_2, \beta_2), \dots,$   
 456  $(x_s, \alpha_s, \beta_s)$  where each  $x_j$  represents one of the  $s$  minmer  
 457 hash values from  $\pi_s(A)$  in increasing order and for each  
 458

460  $i \in [a, z]$ , the values  $\alpha_j$  and  $\beta_j$  are

- 461 •  $\alpha_j = 1$  if  $x_j \in \pi_s(B_i)$  else 0
- 462 •  $\beta_j = 1 + |\{x \in \pi_s(B_i) \text{ s.t. } x_{j-1} < x < x_j\}|$

463 We can imagine  $V$  as a set of  $s$  buckets labeled by  
 464 the  $s$  corresponding hash values of  $A$  and sorted in  
 465 increasing order. At each position  $i \in [a, z]$ , each  
 466 bucket  $j$  holds  $x_j$  and all  $\beta_j - 1$  reference minmers in  
 467  $\pi_s(B_i)$ , which are between  $x_j$  and  $x_{j-1}$ . A bucket is  
 468 marked “good” ( $\alpha_j \rightarrow 1$ ) if  $x_j \in \pi_s(B_i)$ . It remains  
 469 to find the largest integer  $p_i$  such that the number of  
 470 minmers in the first  $p_i$  buckets is at most  $s$ . Given  
 471  $p_i$ , the numerator of the minhash formula,  $Y_i$ , is the  
 472 number of “good” buckets in the first  $p_i$  buckets.

For a candidate region  $[a, z]$ , we initialize  $V$  by insert-  
 ing all of the minmers from the reference index whose  
 intervals overlap with  $a$  and set

$$p_a = \max_q \left( \sum_{j=0}^{j \leq q} \beta_j \leq s \right)$$

It follows that  $Y_a = \sum_{j=1}^{j \leq p_a} \alpha_j$

473 In order to keep track of intervals which overlap with  
 474 the current position, we use a min-heap  $H$  sorted on  
 475 interval endpoints. We then continue to iterate through  
 476 minmer intervals from the reference in order based on  
 477 their start points, stopping once the intervals no longer  
 478 overlap with  $[a, z]$ . For each minmer interval starting  
 479 at  $i \in [a + 1, z]$ , we pop intervals from  $H$  that end at or  
 480 before  $i$ . For each interval popped from  $H$ , we update  $V$   
 481 in  $\mathcal{O}(\log(s))$  time through a binary search, decrementing  
 482 the corresponding  $\beta_j$  and setting  $\alpha_j = 0$  if the interval  
 483 represents a shared minmer. The new interval is added  
 484 in a similar manner and the necessary  $\alpha$  and  $\beta$  values  
 485 are updated. After  $V$  is updated,  $p_i$  is updated from  
 486  $p_{i-1}$  by incrementing or decrementing until it is the  
 487 maximal value such that  $p_i = \max_q \left( \sum_{j=0}^{j \leq q} \beta_j \leq s \right)$ . By  
 488 keeping track of  $p_{i-1}$  and the sums  $\sum_{j=0}^{j \leq p_{i-1}} \beta_j$  and  
 489

490  $\sum_{j=0}^{p_i-1} \alpha_j$ , the new  $p_i$  and corresponding sums are  
 491 updated in constant time per window.

492 While the MashMap3 implementation of the second  
 493 filtering stage still requires  $\mathcal{O}(\log(s))$  time to update the  
 494 minhash for each sliding window within the candidate  
 495 region, it is significantly more efficient than MashMap2's  
 496 ordered map in practice due to  $V$  being a static data  
 497 structure in contiguous memory, only requiring updates  
 498 to counters.

### 499 4.2.3 Early termination of stage 2

500 Instead of computing the stage 2 step for each can-  
 501 didate region obtained in the first stage, we aim to  
 502 terminate the second stage once we have confidently  
 503 identified all mappings whose predicted ANI is within  
 504  $\Delta_{\text{ANI}}$  of the best predicted ANI. We do this by sorting  
 505 the candidate regions in decreasing order of their max-  
 506 imum interval overlap size obtained in stage 1. The  
 507 stage 2 minhash calculation is then performed on each  
 508 candidate region in order, keeping track of the best  
 509 predicted ANI value seen. Let  $\kappa$  be numerator of the  
 510 minhash that corresponds to an ANI value  $\Delta_{\text{ANI}}$  less  
 511 than the best predicted ANI value seen so far. Then,  
 512 given a candidate region with a maximum overlap size  
 513 of  $c_i < \kappa$ , we know that  $\Pr(Y_i \geq \kappa) = 0$  and therefore  
 514 no more candidate regions can contain mappings whose  
 515 predicted ANI is within  $\Delta_{\text{ANI}}$  of the predicted ANI of  
 516 the best mapping.

## 517 5 Results

### 518 5.1 Characteristics of the minmer 519 scheme

520 Here we provide formulas for the density of minmers and  
 521 minmer intervals and an approximation for the distance  
 522 between adjacent minmers. Proofs of the formulas are  
 523 presented in the Supplementary Materials. We then  
 524 compare these formulas to results on both simulated  
 525 and empirical sequences. For the simulated dataset, we  
 526 generated a sequence of 1 million uniform random hash  
 527 values. For the empirical dataset, we used MurmurHash  
 528 to hash the sequence of  $k$ -mers in the recently-completed  
 529 human Y-chromosome Rhie *et al.* (2022) with  $k = 18$ .

#### 530 5.1.1 Minmer density

531 To obtain the formula for the minmer density, we con-  
 532 sider how the rank of a random  $k$ -mer changes with  
 533 each consecutive window that contains it. As a result,  
 534 we have a distribution of the rank of a random  $k$ -mer  
 535 throughout consecutive sliding windows. This distribu-  
 536 tion enables us to not only obtain the density (Figure  
 537 2), but also determine other characteristics such as the  
 538 likelihood of being a minmer given some initial rank  $r_1$   
 539 or given a hash value  $z$ .

**Theorem 5.1.** Let  $d_{(w,s)}$  be the expected density of  
 $(w, s)$ -minmers in a random sequence. Then,

$$d_{(w,s)} = \frac{1}{w} \sum_{r_1, r_w \in \{1 \dots w\}} \Pr(C = 1 | r_1, r_w) \Pr(R_w = r_w | r_1)$$

540 where  $R_w | r_1 \sim \text{BetaBinomial}(r_1, w - r_1 + 1)$  and

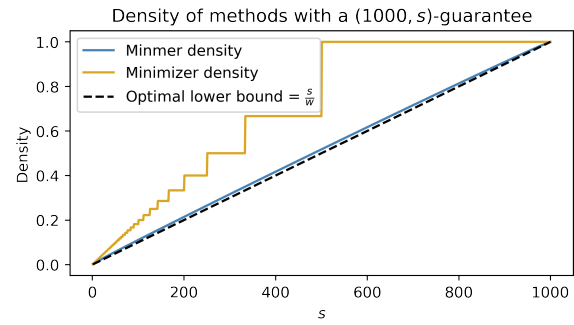


Figure 2: The density of a  $(1000, s)$ -minmer scheme compared to a  $w'$ -minimizer scheme which also yields a  $(1000, s)$ -window guarantee. To ensure that the minimizer scheme satisfies the  $(1000, s)$  window guarantee, the minimizer scheme is set with  $w' = \lfloor 1000/s \rfloor$ .

$$\Pr(C = 1 | r_1, r_w) = \begin{cases} \sum_{u=0}^{\delta} \Pr(U = u) \binom{2u+r_w-r_1}{u+r_w-s} \binom{2u+r_w-r_1}{u}^{-1} & r_1, r_w > s \\ 1 & \text{otherwise} \end{cases}$$

541 where  $U \sim \text{Hypergeometric}(w - 1, r_1 - 1, w - r_w)$  and  
 542  $\delta = \min(r_1 - 1, w - r_w)$ .

#### 543 5.1.2 Minmer interval density

**Theorem 5.2.** Let  $d_{(w,s)}^*$  be the density of  $(w, s)$ -  
 544 minmer intervals in a random sequence. Then,  
 545

$$d_{(w,s)}^* = 1 - \frac{(w - s + 1)(w - s)}{w(w + 1)}$$

546 As expected, letting  $s = 1$  yields the same density as  
 547 minimizers,  $2/(w + 1)$ , and a similar formula appears  
 548 when determining the probability of observing  $s$  con-  
 549 secutive unsampled  $k$ -mers under the the minimizer  
 550 scheme Spouge (2022). As the number of minmers is a  
 551 strict lower bound on the number of minmer intervals,  
 552 this result also gives an upper bound on the density of  
 553  $(w, s)$ -minmers.

#### 554 5.1.3 Minmer window guarantee

555 As the main difference between minimizers and minmers  
 556 is the window guarantee, it is important to observe the  
 557 difference in the density of the minmer scheme compared  
 558 to a minimizer scheme which also satisfies the  $(w, s)$ -  
 559 window guarantee. In Figure 2, we consider the case  
 560 where we have a  $(1000, s)$ -minmer scheme and a  $w'$ -  
 561 minimizer scheme, where  $w'$  is set to obtain the same  
 562 window guarantee of the minmer scheme by letting  $w' =$   
 563  $\lfloor 1000/s \rfloor$ . We observe that for sketch sizes other than 1  
 564 and 1000, for which the density of the schemes are equal,  
 565 the density of the minmer scheme is strictly less than  
 566 the density of the corresponding minimizer scheme. For  
 567 some values of  $s$ , the density of the  $\lfloor 1000/s \rfloor$ -minimizer  
 568 scheme is over 70% larger than the  $(1000, s)$ -minmer  
 569 scheme.

#### 570 5.1.4 Minmer spread

571 Let  $G_i$  be the distance between the  $i^{\text{th}}$  selected minmer  
 572 and the  $(i + 1)^{\text{th}}$  selected minmer. For a  $(w, s)$ -minmer  
 573 scheme with a density factor  $d_f$ , we have that

$$\Pr(G_i = d) \approx \frac{\binom{w-d}{d_f-2}}{\binom{w}{d_f-1}}$$

574 To see how well this approximation holds, we plot  
575 the results on both empirical and simulated data in  
576 Supplemental Figure 2.

## 577 5.2 ANI prediction ideal sequences

578 We replicated the experiments for Table 1 of Belbasi  
579 *et al.* (2022) using the minmer-based MashMap3 (com-  
580 mit 0b47608), with the exception that we report the  
581 mean predicted sequence divergence as opposed to the  
582 median. For each divergence rate  $r \in \{0.01, 0.05, 0.10\}$ ,  
583 100 random windows of 10,000 base pairs were selected  
584 from the *Escherichia coli* genome and 10,000r positions  
585 were selected at random and mutated, ensuring that  
586 no duplicate  $k$ -mers were generated. The reads were  
587 mapped back to the reference *E. coli* genome and the  
588 predicted divergence was compared to the ground truth  
589 (Figure 3).

590 The parameters of the minmer-based MashMap3 were  
591 set to obtain a similar numbers of sampled  $k$ -mers as  
592 the minimizer-based MashMap2 under MashMap2’s  
593 default settings, resulting in a density of 0.009 for both  
594 tools. As expected, the results show that the ANI  
595 values predicted by the minmer scheme are significantly  
596 closer to the ground truth than those predicted by the  
597 minimizer scheme. Notably, in the case where the true  
598 divergence was 1%, the relative error is reduced from  
599 29% to 2% (Figure 3).

## 600 5.3 ANI prediction on simulated reads

601 In addition to the ANI prediction measurements from  
602 Belbasi *et al.* (2022), we also simulated reads from  
603 the human T2T-CHM13 reference genome Nurk *et al.*  
604 (2022) at varying error rates to determine the accuracy  
605 of the ANI predictions. We compared the minmer-based  
606 MashMap3 against the minimizer-based MashMap2  
607 with similar densities for each run as well as against  
608 Minimap2 Li (2018). Minimap2 was run in its de-  
609 fault mode with `-x map-ont` set which, like MashMap,  
610 computes approximate mappings and estimates the  
611 alignment identity. MashMap2 was modified to use  
612 the binomial model for estimating the ANI from the  
613 Jaccard estimator which has been shown to be more  
614 accurate Belbasi *et al.* (2022).

615 We used Pbsim Ono *et al.* (2013) to simulate three  
616 datasets: “ONT-95”, “ONT-98”, and “ONT-99”, where  
617 the number following the dash represents the average  
618 ANI across reads. The standard deviation of the error  
619 rates was set to 0, and the ratio of matches, insertions,  
620 and deletions was set to 20:40:40, respectively, to en-  
621 sure that mapped regions would, on average, be the  
622 same length as the reads. For each dataset, 5,000bp  
623 reads were generated with the CLR profile at a depth  
624 of 2, resulting in 1.25 million reads for each dataset.  
625 The mappings output by the different methods were  
626 parsed and the predicted ANI was compared to the gap-  
627 compressed ANI of the ground-truth mapping. The  
628 results of the simulations can be seen in Table 1.

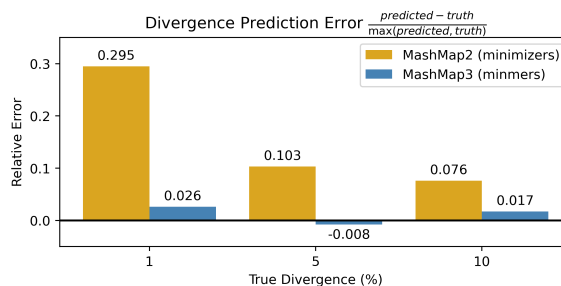


Figure 3: **Eliminating the bias in MashMap.** The experiments from Table 1 of Belbasi *et al.* (2022) were replicated. Divergence, defined as 1-ANI, was predicted across 100 sequences for both MashMap2 and MashMap3 using a density of 0.009 ( $L = 10,000$ ,  $s = 78$ ). In the case of identifying sequence divergence of very closely related genomes, the minmer scheme reduces relative prediction error from 29% to 2%.

629 For MashMap2 and MashMap3, we used a  $k$ -mer size  
630 of 19 and set the MashMap2 minimizer  $w$  to 89 and min-  
631 mer  $s$  to  $s = 100$  to obtain a density of 0.0222 for both  
632 tools. The ANI cutoff was set to 94%, 93%, and 90%  
633 for the ONT-99, ONT-98, and ONT-95 datasets, respec-  
634 tively. The indexing times for Minimap2, MashMap2,  
635 and MashMap3 were 1.7, 2.8, and 9.8 minutes, respec-  
636 tively.

## 637 5.4 ANI prediction on mammalian 638 genome alignments

639 To test the performance of MashMap3 at the genome-  
640 mapping scale, we computed mappings between the  
641 T2T human reference genome and reference genomes for  
642 chimpanzee Kronenberg *et al.* (2018) and macaque War-  
643 ren *et al.* (2020). In absence of ground truth ANI values,  
644 we used wfmask Guarracino *et al.* (2021) to compute the  
645 gap-compressed ANI of the segment mappings output  
646 by MashMap and report the results of the mappings  
647 with  $\geq 80\%$  complexity in Table 2. For a small pro-  
648 portion of segment mappings output by MashMap2  
649 and MashMap3, wfmask did not produce an alignment.  
650 When the ANI threshold is 85%, these cases accounted  
651 for 0.07% of chimpanzee mappings and 0.3% macaque  
652 mappings. When the ANI threshold was 90% or 95%,  
653 less than 0.01% of mappings were not aligned with wf-  
654 mask for both chimpanzee and macaque. We consider  
655 these mappings as false positives. For the ANI thresh-  
656 olds of 95%, 90%, and 85%, the winnowing scheme  
657 densities were set to 0.043, 0.053, and 0.064, respec-  
658 tively.

659 To isolate the effect of the new seeding method, we  
660 turned chaining off for both tools. As the Jaccard es-  
661 timator is known to perform poorly in the presence of  
662 many degenerate  $k$ -mers, results for query regions above  
663 and below 80% complexity are reported separately,  
664 where complexity is defined as the ratio of observed  
665 distinct  $k$ -mers in a region to  $w$ . Low-complexity map-  
666 pings make up for at most 1% and 3% of the mappings  
667 for chimpanzee and macaque genomes, respectively. We  
668 show the table of the metrics for the low-complexity  
669 mappings in Supplementary Table 1.

Dataset	Minimap2				MashMap2				MashMap3			
	CPU time (m)	Memory (Gb)	ME	MAE	CPU time (m)	Memory (Gb)	ME	MAE	CPU time (m)	Memory (Gb)	ME	MAE
ONT-99	154.20	<b>9.89</b>	-0.25	0.34	80.27	9.92	-0.27	0.29	<b>33.64</b>	13.07	<b>0.03</b>	<b>0.17</b>
ONT-98	147.29	<b>9.89</b>	-0.36	0.52	82.46	9.92	-0.33	0.39	<b>35.13</b>	13.09	<b>0.06</b>	<b>0.29</b>
ONT-95	96.35	<b>9.89</b>	-0.46	0.81	106.81	9.92	-0.25	<b>0.59</b>	<b>42.81</b>	13.10	<b>0.21</b>	0.62

Table 1: **Metrics for simulated Nanopore read mapping to the human genome.** Minmer and minimizer-based MashMap implementations as well as Minimap2 were used to map simulated reads from the human reference genome using Pbsim Ono *et al.* (2013).

Query Species	ANI Threshold	MashMap2					MashMap3				
		Basepairs mapped (Gbp)	CPU time (m)	Memory (Gb)	ME	MAE	Basepairs mapped (Gbp)	CPU time (m)	Memory (Gb)	ME	MAE
Chimpanzee	95%	2.80	39.76	<b>19.95</b>	-0.25	0.29	<b>2.81</b>	<b>32.76</b>	27.07	<b>0.01</b>	<b>0.22</b>
Chimpanzee	90%	2.82	118.31	<b>24.55</b>	-0.22	0.29	2.82	<b>51.12</b>	36.20	<b>0.01</b>	<b>0.25</b>
Chimpanzee	85%	2.83	787.44	44.96	-0.18	0.27	2.83	<b>64.48</b>	<b>39.47</b>	<b>0.02</b>	<b>0.25</b>
Macaque	95%	0.38	30.0	<b>20.83</b>	<b>0.29*</b>	<b>0.46</b>	<b>1.08</b>	<b>28.67</b>	28.97	0.57*	0.66
Macaque	90%	2.54	40.49	<b>23.04</b>	-0.30	<b>0.69</b>	<b>2.56</b>	<b>34.87</b>	35.91	<b>0.01</b>	0.74
Macaque	85%	2.60	446.71	<b>38.13</b>	-0.24	<b>0.74</b>	<b>2.61</b>	<b>43.74</b>	39.49	<b>0.05</b>	0.87

Table 2: **Comparison of MashMap2 and MashMap3 for identifying mappings between pairs of mammalian genomes.** MashMap2 and MashMap3 were used to align the human reference genome to chimpanzee and macaque genomes. The ME and MAE metrics shown are for query segments with at least 80%  $k$ -mer complexity. Corresponding metrics for low-complexity mappings can be found in Supplementary Table 1. \*Sampling bias leads to ANI over-estimation. See discussion for details.

## 6 Discussion

Minmers are a novel “non-forward” winnowing scheme with a  $(w, s)$ -window guarantee. Similar to what has been done for other proposed schemes, we have derived formulas (approximate and exact) that describe the scheme’s characteristics. We have replaced minimizers with minmers in MashMap3 and demonstrated that minmers eliminate Jaccard estimator bias and enable new methods to reduce mapping runtime compared to MashMap2. In addition, we show that minmers require substantially less density than minimizers when a  $(w, s)$ -window guarantee is required.

### The minmer scheme enables sparser sketches

The minimizer winnowing scheme has long been the dominant method for winnowing due to its  $(w, 1)$ -window guarantee, simplicity, and performance. Other 1-local methods such as strobemers Sahlin (2021) and syncmers Edgar (2021) remove the window guarantee and rely on a random sequence assumption to provide probabilistic bounds on the expected distance between sampled  $k$ -mers.

Minmers represent a novel class of winnowing schemes that extend the window guarantee of minimizers. Unlike strobemers, syncmers, and other 1-local methods, the minmer scheme guarantees the desired number of  $k$ -mers will be sampled from every window, so long as it contains at least  $s$  distinct  $k$ -mers. This is particularly desirable for accurate Jaccard estimation and the winnowing of low-complexity sequence where the density of sampled  $k$ -mers from 1-local schemes can vary significantly.

### Minmers yield an unbiased estimator at lower computational costs

Indexing minmers rather than minimizers removes the Jaccard estimator bias present in earlier versions of MashMap. For any window, the set of sampled  $k$ -mers is guaranteed to be a superset of the bottom- $s$  sketch of that window. Therefore, running the minhash algorithm on minmers yields the same estimator as running the minhash algorithm on the full set of  $k$ -mers.

In addition to the experiments from Belbasi *et al.* (2022), which focus on “ideal” sequences with no repetitive  $k$ -mers, we also measured the performance of the ANI prediction for different levels of divergence on the human genome across mappings of simulated reads and a sample of mammalian genomes. Our results showed that MashMap3 with minmers not only produced unbiased and more accurate predictions of the ANI than Minimap2 and MashMap2, but it did so in a fraction of the time.

We replicated the behavior of minimizers to under-predict ANI as seen in Belbasi *et al.* (2022) across all experiments. At the same time, in both the simulated reads and empirical genome alignment results, we see that MashMap3 slightly over-predicts the ANI at larger divergences. Further inspection reveals that this is due to indels in the alignment, which are not modeled by the binomial mutation model used to convert the Jaccard to ANI (Supplementary Table 2).

The optimizations to the second stage of mapping combined with the minmer interval indexing leads to significantly better mapping speeds in MashMap3. Relative to Minimap2 and MashMap2, MashMap3 spends



733 a significant amount of time indexing the genome. This, 732  
734 however, serves as an investment for the mapping phase 733  
735 which is significantly faster than MashMap2, particu- 734  
736 larly at lower ANI thresholds. As an additional feature, 735  
737 MashMap3 provides the option to save the reference 736  
738 index so that users can leverage the increased mappings 737  
739 speeds for previously indexed genomes. 738

740 Similar to MashMap2, MashMap3 by default uses the 739  
741 plane-sweep post-processing algorithm described in Jain 740  
742 *et al.* (2018a) to filter out redundant segment mappings. 741  
743 We show that by using the probabilistic filtering method 742  
744 described in Section 4.2.1, we can discard many of these 743  
745 mappings at the beginning of the process as opposed 744  
746 to the end, yielding significant runtime improvements. 745

747 MashMap3 is significantly more efficient at lower ANI 746  
748 thresholds, which is helpful for detecting more distant 747  
749 homologies. For example, in our human-chimpanzee 748  
750 mapping, we recovered an additional 50 Mbp of mapped 749  
751 sequence by reducing the ANI threshold from 95% 750  
752 to 85% while also completing over 10x quicker than 751  
753 MashMap2. It is also worth noting that the default 752  
754 ANI of MashMap2 and MashMap3 is 85%, and often 753  
755 the ANI of homologies between genomes is not known 754  
756 a priori. 755

757 Further motivating the improved efficiency of low ANI 756  
758 thresholds is the fact that thresholds above the true ANI 757  
759 can lead to recovering mappings which over-predict the 758  
760 ANI while discarding those which accurately or under- 759  
761 predict the ANI. This sampling bias leads to an increase 760  
762 in the ANI estimation bias. We see this behavior in 761  
763 the human-macaque alignment with a threshold 95% 762  
764 ANI (Table 2). At lower ANI thresholds, we observe 763  
765 that the majority of mappings are in the 90%-95% ANI 764  
766 range. 765

## 767 Limitations and future directions

768 MashMap’s Jaccard-based similarity method tends to 767  
769 overestimate ANI in low-complexity sequences. For 768  
770 downstream alignment applications, the resulting false- 769  
771 positive mappings can be pruned using a chaining or 770  
772 exact alignment algorithm to validate the mappings. 771  
773 Unreliable ANI estimates could also be flagged by using 772  
774 the bottom-*s* sketch to determine the complexity of a 773  
775 segment as described in Cohen and Kaplan (2007), but 774  
776 a sketching method and distance metric that better 775  
777 approximates ANI across all sequence and mutational 776  
778 contexts would be desirable. 777

779 An important characteristic of MashMap is the rela- 778  
780 tively few parameter settings necessary to tune across 779  
781 different use cases. Building on this, we aim to de- 780  
782 velop a methodology that can find maximal homologies 781  
783 without a pre-determined segment size, similar to the 782  
784 approach of Wang *et al.* (2022b). 783

## 785 7 Conclusion

786 In this work, we proposed and studied the charac- 785  
787 teristics of the minmer scheme and showed that they 786  
788 belong to the unexplored class of non-forward local 787  
789 schemes, which have the potential to achieve lower den- 788  
790 sities under the same locality constraints as forward 789  
791 schemes Marçais *et al.* (2018). We derived formulas

792 for the density and approximate spread of minmers, 792  
793 enabling them to be objectively compared to other 793  
794 winnowing schemes. 794

795 By construction, minmers, unlike minimizers, enable 795  
796 an unbiased estimation of the Jaccard. We replaced 796  
797 the minimizer winnowing scheme in MashMap2 with 797  
798 minmers and showed that minmers significantly reduce 798  
799 the bias in both simulated and empirical datasets. 799

800 Through leveraging the properties of the minmers, 800  
801 we implemented a number of algorithmic improvements 801  
802 in MashMap3. In our experiments, these improvements 802  
803 yielded significantly lower runtimes, particularly in the 803  
804 case when the ANI threshold of MashMap is set to the 804  
805 default of 85%. With the improvements in MashMap3, 805  
806 it is no longer necessary to estimate the ANI of ho- 806  
807 mologies a priori to avoid significantly longer runtimes, 807  
808 making it an ideal candidate for a broad range of com- 808  
809 parative genomics applications. 809

## 810 Acknowledgements

811 We would like to thank Chirag Jain for helpful discus- 811  
812 sions and his implementation of the original MashMap 812  
813 software, as well Andrea Guarracino for improvements 813  
814 and discussions. We would also like to thank Nicolae 814  
815 Sapoval and Fritz Sedlazeck for their feedback on the 815  
816 proofs and primate alignments, respectively. 816

## 817 Funding

818 B.K. was supported by the NLM Training Program 818  
819 in Biomedical Informatics and Data Science (Grant: 819  
820 T15LM007093). A.M.P. was supported by the Intramur- 820  
821 al Research Program of the National Human Genome 821  
822 Research Institute, National Institutes of Health. B.K. 822  
823 and T.T. were supported in part by the National Insti- 823  
824 tute of Allergy and Infectious Diseases (Grant# P01- 824  
825 AI152999). T.T. was supported in part by National Sci- 825  
826 ence Foundation grant EF-2126387. E.G. was supported 826  
827 by National Institutes of Health/NIDA U01DA047638, 827  
828 National Institutes of Health/NIGMS R01GM123489, 828  
829 NSF PPOSS Award #2118709, and the Tennessee Gov- 829  
830 ernor’s Chairs program. 830

## 831 References

- 832 Baker, D. N. and Langmead, B. (2019). Dashing: fast 832  
833 and accurate genomic distances with hyperloglog. 833  
834 *Genome biology*, **20**(1), 1–12. 834
- 835 Belbasi, M. *et al.* (2022). The minimizer jaccard esti- 835  
836 mator is biased and inconsistent. *Bioinformatics*. 836
- 837 Berlin, K. *et al.* (2015). Assembling large genomes 837  
838 with single-molecule sequencing and locality-sensitive 838  
839 hashing. *Nature biotechnology*, **33**(6), 623–630. 839
- 840 Blanca, A. *et al.* (2022). The statistics of k-mers from 840  
841 a sequence undergoing a simple mutation process 841  
842 without spurious matches. *Journal of Computational 842*  
843 *Biology*, **29**(2), 155–168. 843
- 844 Broder, A. Z. (1997). On the resemblance and contain- 844  
845 ment of documents. In *Proceedings. Compres-* 845

- 846 *sion and Complexity of SEQUENCES 1997 (Cat. No.*  
847 *97TB100171)*, pages 21–29. IEEE.
- 848 Brown, C. T. and Irber, L. (2016). sourmash: a library  
849 for minhash sketching of dna. *Journal of Open Source*  
850 *Software*, **1**(5), 27.
- 851 Cohen, E. (2016). Min-hash sketches.
- 852 Cohen, E. and Kaplan, H. (2007). Summarizing data  
853 using bottom-k sketches. In *Proceedings of the twenty-*  
854 *sixth annual ACM symposium on Principles of dis-*  
855 *tributed computing*, pages 225–234.
- 856 Dilthey, A. T. *et al.* (2019). Strain-level metagenomic  
857 assignment and compositional estimation for long  
858 reads with metamaps. *Nature communications*, **10**(1),  
859 1–12.
- 860 Edgar, R. (2021). Syncmers are more sensitive than  
861 minimizers for selecting conserved k-mers in biological  
862 sequences. *PeerJ*, **9**, e10805.
- 863 Ekim, B. *et al.* (2020). A randomized parallel algo-  
864 rithm for efficiently finding near-optimal universal  
865 hitting sets. In *International Conference on Research*  
866 *in Computational Molecular Biology*, pages 37–53.  
867 Springer.
- 868 Guarracino, A. *et al.* (2021). wfmash: a pangenome-  
869 scale aligner.
- 870 Jain, C. *et al.* (2017). A fast approximate algorithm for  
871 mapping long reads to large reference databases. In  
872 *International Conference on Research in Computa-*  
873 *tional Molecular Biology*, pages 66–81. Springer.
- 874 Jain, C. *et al.* (2018a). A fast adaptive algorithm for  
875 computing whole-genome homology maps. *Bioinform-*  
876 *atics*, **34**(17), i748–i756.
- 877 Jain, C. *et al.* (2018b). High throughput ani analy-  
878 sis of 90k prokaryotic genomes reveals clear species  
879 boundaries. *Nature communications*, **9**(1), 1–8.
- 880 Kronenberg, Z. N. *et al.* (2018). High-resolution com-  
881 parative analysis of great ape genomes. *Science*,  
882 **360**(6393), eaar6343.
- 883 Li, H. (2018). Minimap2: pairwise alignment for nu-  
884 cleotide sequences. *Bioinformatics*, **34**(18), 3094–  
885 3100.
- 886 Marçais, G. *et al.* (2017). Improving the performance of  
887 minimizers and winnowing schemes. *Bioinformatics*,  
888 **33**(14), i110–i117.
- 889 Marçais, G. *et al.* (2018). Asymptotically optimal mini-  
890 mizers schemes. *Bioinformatics*, **34**(13), i13–i22.
- 891 Marco-Sola, S. *et al.* (2021). Fast gap-affine pairwise  
892 alignment using the wavefront algorithm. *Bioinform-*  
893 *atics*, **37**(4), 456–463.
- 894 Nurk, S. *et al.* (2022). The complete sequence of a  
895 human genome. *Science*, **376**(6588), 44–53.
- Ondov, B. D. *et al.* (2016). Mash: fast genome 896  
and metagenome distance estimation using minhash. 897  
*Genome biology*, **17**(1), 1–14. 898
- Ondov, B. D. *et al.* (2019). Mash screen: high- 899  
throughput sequence containment estimation for 900  
genome discovery. *Genome biology*, **20**(1), 1–13. 901
- Ono, Y. *et al.* (2013). Pbsim: Pacbio reads simula- 902  
tor—toward accurate genome assembly. *Bioinformat-* 903  
*ics*, **29**(1), 119–121. 904
- Popic, V. and Batzoglou, S. (2017). A hybrid cloud 905  
read aligner based on minhash and kmer voting that 906  
preserves privacy. *Nature communications*, **8**(1), 1–7. 907
- Rhie, A. *et al.* (2022). The complete sequence of a 908  
human y chromosome. *bioRxiv*. 909
- Roberts, M. *et al.* (2004). Reducing storage require- 910  
ments for biological sequence comparison. *Bioinform-* 911  
*atics*, **20**(18), 3363–3369. 912
- Sahlin, K. (2021). Effective sequence similarity de- 913  
tection with strobemers. *Genome research*, **31**(11), 914  
2080–2094. 915
- Schleimer, S. *et al.* (2003). Winnowing: local algo- 916  
rithms for document fingerprinting. In *Proceedings* 917  
*of the 2003 ACM SIGMOD international conference* 918  
*on Management of data*, pages 76–85. 919
- Shaw, J. and Yu, Y. W. (2022). Theory of local k-mer 920  
selection with applications to long-read alignment. 921  
*Bioinformatics*, **38**(20), 4659–4669. 922
- Shaw, J. and Yu, Y. W. (2023). Fast and robust metage- 923  
nomic sequence comparison through sparse chaining 924  
with skani. *bioRxiv*, pages 2023–01. 925
- Spouge, J. L. (2022). A closed formula relevant to 926  
‘theory of local k-mer selection with applications to 927  
long-read alignment’ by jim shaw and yun william yu. 928  
*Bioinformatics*, **38**(20), 4848–4849. 929
- Wang, T. *et al.* (2022a). The human pangenome project: 930  
a global resource to map genomic diversity. *Nature*, 931  
**604**(7906), 437–446. 932
- Wang, Z. *et al.* (2022b). Txtalign: Efficient near- 933  
duplicate text alignment search via bottom-k sketches 934  
for plagiarism detection. In *Proceedings of the 2022* 935  
*International Conference on Management of Data*, 936  
pages 1146–1159. 937
- Warren, W. C. *et al.* (2020). Sequence diversity analyses 938  
of an improved rhesus macaque genome enhance its 939  
biomedical utility. *Science*, **370**(6523), eabc6617. 940
- Zheng, H. *et al.* (2020). Improved design and analysis 941  
of practical minimizers. *Bioinformatics*, **36**(Supple- 942  
ment\_1), i119–i127. 943