Sequence analysis

A near-tight lower bound on the density of forward sampling schemes

Bryce Kille (**b**^{1,*,†}, Ragnar Groot Koerkamp (**b**^{2,*,†}, Drake McAdams¹, Alan Liu¹, Todd J. Treangen^{1,3}

¹Department of Computer Science, Rice University, Houston, TX 77005, United States

²Department of Computer Science, ETH Zurich, Zurich 8092, Switzerland

³Ken Kennedy Institute, Rice University, Houston, TX 77005, United States

*Corresponding authors. Department of Computer Science, Rice University, 6100 Main Street, Houston, TX 77005, United States. E-mail: blk6@rice.edu (B.K.); Department of Computer Science, ETH Zurich, Rämistrasse 101, Zurich 8092, Switzerland. E-mail: ragnar.grootkoerkamp@inf.ethz.ch (R.G.K.)

[†]Equal contribution.

Associate Editor: Yann Ponty

Abstract

Motivation: Sampling *k*-mers is a ubiquitous task in sequence analysis algorithms. Sampling schemes such as the often-used random minimizer scheme are particularly appealing as they guarantee at least one *k*-mer is selected out of every *w* consecutive *k*-mers. Sampling fewer *k*-mers often leads to an increase in efficiency of downstream methods. Thus, developing schemes that have low density, i.e. have a small proportion of sampled *k*-mers, is an active area of research. After over a decade of consistent efforts in both decreasing the density of practical schemes and increasing the lower bound on the best possible density, there is still a large gap between the two.

Results: We prove a near-tight lower bound on the density of forward sampling schemes, a class of schemes that generalizes minimizer schemes. For small w and k, we observe that our bound is tight when $k \equiv 1 \pmod{w}$. For large w and k, the bound can be approximated by $\frac{1}{w+k} \left\lceil \frac{w+k}{w} \right\rceil$. Importantly, our lower bound implies that existing schemes are much closer to achieving optimal density than previously known. For example, with the current default minimap2 HiFi settings w=19 and k=19, we show that the best known scheme for these parameters, the double decycling-set-based minimizer of Pellow *et al.* is at most 3% denser than optimal, compared to the previous gap of at most 50%. Furthermore, when $k \equiv 1 \pmod{w}$ and the alphabet size σ goes to ∞ , we show that mod-minimizers introduced by Groot Koerkamp and Pibiri achieve optimal density matching our lower bound.

Availability and implementation: Minimizer implementations: github.com/RagnarGrootKoerkamp/minimizers ILP and analysis: github.com/ treangenlab/sampling-scheme-analysis.

1 Introduction

For over a decade, k-mer sampling schemes have served as a ubiquitous first step in many classes of bioinformatics tasks. By sampling k-mers in a way which ensures that two similar sequences will have similar sets of sampled k-mers, sampling schemes enable methods to bypass the need to compare entire sequences at the base level and instead allow them to work more efficiently using the sampled k-mers.

Local sampling schemes satisfy a window guarantee that at least one k-mer is selected out of every window of w consecutive k-mers. Most schemes used in practice, such as the random minimizer scheme (Schleimer et al. 2003, Roberts et al. 2004), are forward schemes that additionally guarantee that k-mers are sampled in the order in which they appear in the original sequence. These properties are particularly appealing since they guarantee that no region is left unsampled.

As the purpose of these schemes is to reduce the computational burden of downstream methods while upholding the window guarantee, the primary goal of most new schemes is to minimize the *density*, i.e. the expected proportion of sampled k-mers. Over the past decade, many new schemes have been proposed that obtain significantly lower densities than the original random minimizer scheme.

For example, there are schemes based on *hitting sets* (Orenstein *et al.* 2016, Marçais *et al.* 2017, 2018, DeBlasio *et al.* 2019, Ekim *et al.* 2020, Pellow *et al.* 2023, Golan *et al.* 2024), schemes that focus on sampling positions rather than *k*-mers (Loukides and Pissis 2021, Loukides *et al.* 2023), schemes that use an ordering on *t*-mers (t < k) to decide which *k*-mer to sample (Zheng *et al.* 2020, Groot Koerkamp and Pibiri 2024), and schemes that aim to minimize density on specific input sequences (Zheng *et al.* 2021b, Hoang *et al.* 2022). All of these improvements notwithstanding, it is still unknown how close these schemes are to achieving minimum density.

A trivial lower bound on density given by the window guarantee is $\frac{1}{w}$, and recently Groot Koerkamp and Pibiri (2024) improved the bound of Marçais *et al.* (2018) from $\frac{1.5 + \frac{1}{2w}}{w+k}$ to $\frac{1.5}{w+k-0.5}$. However, for many practical values of w and k, there is a sizeable gap between these lower bounds and the density of existing schemes. This raises the question

Received: 9 September 2024; Revised: 16 November 2024; Editorial Decision: 5 December 2024; Accepted: 10 December 2024 © The Author(s) 2024. Published by Oxford University Press.

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

whether schemes with density much closer to $\frac{1}{w}$ exist, but have not been found yet, or whether existing schemes are already very close to optimal and it is the lower bound that needs improvement. Our new lower bound closes most of the gap, and thus answers this question: Indeed, especially for $k \ge w$, the best existing schemes have near-optimal density in many cases. This allows future research to focus on improving other sampling scheme metrics, such as the *conservation* described by Edgar (2021) and Shaw and Yu (2022).

1.1 Contributions

Main lower bound theorem. We prove a novel lower bound on the density of forward schemes that is strictly tighter than all previously established lower bounds for all w, k, and alphabet size σ :

Theorem 1. Let f be a (w, k)-forward sampling scheme and $M_{\sigma}(p)$ count the number of aperiodic necklaces of length p over an alphabet of size σ . Then, the density of f is at least

$$g_{\sigma}(w,k) := \frac{1}{\sigma^{w+k}} \sum_{p \mid (w+k)} M_{\sigma}(p) \left\lceil \frac{p}{w} \right\rceil \ge \frac{\lceil \frac{w+k}{w} \rceil}{w+k} \ge \frac{1}{w}, \quad (1)$$

where the middle inequality is strict for w > 1.

We prove that this bound can be extended to work for more general classes of sampling schemes, such as the local schemes described by Marçais *et al.* (2018) and the multi-local schemes described by Kille *et al.* (2023).

Comparison with optimal schemes for small parameters. We show that our lower bound is tight for some small w, k, and σ by using an integer linear program to construct schemes whose density matches our lower bound. This marks the first time that there is an analytical description of a tight minimum density of any forward scheme. We conjecture that when $k \equiv 1 \pmod{w}$, there exist schemes with density matching our lower bound.

Comparison with practical schemes for large parameters. To show that our bound is significantly closer to the density achieved by existing schemes compared to previous lower bounds, we replicate the benchmark from Groot Koerkamp and Pibiri (2024) for a selection of w and k (Fig. 3). For example, with the default minimap2 (Li 2018) HiFi settings w = 19 and k = 19, the lower bound goes up from 50% of the density achieved by the double decycling based method to 97% of the achieved density (Table 1).

Analysis of the mod-minimizer. Finally, our new lower bound implies that the mod-minimizer scheme (Groot Koerkamp and Pibiri 2024) is optimal when $k \equiv 1 \pmod{w}$ and σ is large. Indeed, for the ASCII alphabet (σ =256), the mod-minimizer scheme density is consistently within 1% of the lower bound when $k \equiv 1 \pmod{w}$ (Supplementary D, Fig. S4).

2 Background

Notation. We begin by defining some necessary notation, as well as definitions of mathematical concepts that will be used throughout the work. We use [n] to refer to the set $\{0, 1, ..., n-1\}$. The alphabet is denoted by Σ and has size $\sigma := |\Sigma|$, with $\sigma = 4$ for DNA. The expression a|b indicates that a divides b. The summation $\sum_{a|b}$ is over all positive divisors a of b. We use $a \mod m$ for the remainder (in [m]) of a after dividing by m and we use $a \equiv b \pmod{m}$ to indicate that a and b have the same remainder modulo m. Given a string W, W[i,j)) refers to the substring of W containing the characters at 0-based positions i up to j - 1 inclusive. For two strings X and Y, XY represents the concatenation of X and Y.

Classes of sampling schemes. There are multiple established classes of sampling schemes. We begin by drawing a distinction between schemes with and without a window guarantee that guarantees that at least one every w k-mers is sampled. While schemes without a window guarantee, such as fracminhash (Irber *et al.* 2022), are often efficient to compute, the lack of a guarantee on the distance between sampled k-mers makes them ineffective or inefficient for certain tasks such as indexing and alignment. Indeed, we only consider schemes with a window guarantee:

Definition 1. A (w, k)-local scheme with window guarantee w and k-mer size k on an alphabet Σ corresponds to a sampling function $f : \Sigma^{w+k-1} \to [w]$.

In other words, given a window of w + k - 1 characters (w consecutive k-mers), the output of the sampling function f (W) is an integer in [w] which represents the index of the sampled k-mer in W. Recently, Kille *et al.* (2023) proposed a generalization of (w, k)-local schemes which samples at least s k-mers out of every w instead of at least 1 and we extend our results to these more general schemes in Supplementary A.

Local schemes have no restrictions on which of the w kmers can be selected for each window, but *forward schemes* are a subset of local schemes that enforce the restriction that they never select a k-mer which occurs before a previously selected k-mer.

Table 1. Minimum densities achieved by existing sampling schemes for default parameters of frequently-used bioinformatics methods ($\sigma = 4$)^a

Application	(w, k)	$\begin{array}{c} Random \\ 2/(w\!+\!1) \end{array}$	Best		Lower bound		Gap (%)	
			Scheme	Density	1/w	g′	1/w	\mathbf{g}'
Kraken2	(5, 31)	0.333	Mod-mini	0.226	0.200	0.222	12.8	1.6
SSHash	(12, 20)	0.154	Mod-mini	0.120	0.083	0.108	43.9	10.9
minimap2, hifi	(19, 19)	0.100	dbl decycling	0.079	0.053	0.077	50.1	2.7

^a The gap percentage describes the how much larger the lowest achieved density is than the lower bound and is calculated as $100 \cdot \frac{d(f) - LB(w,k)}{LB(w,k)}$, where $LB(w,k) = \frac{1}{w}$ for the old gap and $LB(w,k) = g'_4(w,k)$ for the new gap. While Groot Koerkamp and Pibiri (2024) showed that $\frac{1.5}{w+k-0.5}$ is also a lower bound, $\frac{1}{w}$ is tighter for all of the parameter choices in the table. For SSHash (Pibiri 2022), we show parameters used for indexing a single human genome.

Definition 2. A (w, k)-local scheme is also (w, k)-forward if for all strings $W \in \Sigma^{w+k}$ representing two adjacent windows,

$$f(W[0, w+k-1)) \le f(W[1, w+k)) + 1.$$

Definition 3. The *density* d(f) of a sampling scheme *f* is defined as the expected proportion of sampled positions from an infinite, uniformly random string.

For a further background on types of sampling schemes, we refer to Shaw and Yu (2022), Zheng *et al.* (2023), Groot Koerkamp and Pibiri (2024), and Ndiaye *et al.* (2024).

De Bruijn graphs. Let $B_{n,\sigma} = (V, E)$ denote the complete De Bruijn graph of order *n*, which has as vertices all strings of length *n*, $V = \Sigma^n$, and edges between vertices that overlap in *n* - 1 positions, $E = \{(X, X[1, n)c) | X \in V, c \in \Sigma\}$. When σ is clear from the context or irrelevant for a particular discussion, it is omitted. It is worth noting that the vertices of B_{n+1} correspond to edges of B_n .

For each string *s* of length *n*, the *n* rotations of *s* induce a *pure cycle* in B_n consisting of (up to) *n* vertices cyclically connected by edges. Note that when *s* is repetitive, e.g. a single repeated character or some other repeated string, the length of the cycle will be a divisor of *n*. These pure cycles are also called *necklaces*. The set of necklaces of length *n* corresponds to a partitioning of the vertices of B_n into a vertex-disjoint set of pure cycles. We use C_n to refer to this set of pure cycles of B_n , and for $c \in C_n$, we write |c| for the number of vertices in the cycle.

When a string of length n has n unique rotations, the corresponding necklace is said to be *aperiodic*. The total number of necklaces and the number of aperiodic necklaces of length n are given by Moreau (1872) (and see also Riordan (1957)) as, respectively,

$$N_{\sigma}(n) = rac{1}{n} \sum_{p|n} arphi(n/p) \cdot \sigma^p, \qquad M_{\sigma}(n) = rac{1}{n} \sum_{p|n} \mu(n/p) \cdot \sigma^p,$$

where $\varphi(p)$ is Euler's totient function that counts the number of integers in [*p*] coprime to *p*. The formula $M_{\sigma}(n)$ counting *aperiodic* necklaces follows from the formula for $N_{\sigma}(n)$ via *Möbius inversion* (Möbius 1832), where μ is the *Möbius function* defined to be 0 if *n* is divisible by a square (>1) and $\mu(n) = (-1)^q$ otherwise, where *q* is the number of prime factors of *n*.

Charged contexts. The context of a window of length w+k-1 in a sequence is the set of preceding windows that influences whether the current window samples a new position.

For a local scheme to select a new position, none of the previous w - 1 windows may have selected the same k-mer as the current window. As a result, the context for local schemes consists of 2w + k - 2 characters: the current window of w kmers as well as the w - 1 windows preceding the current window.

For a forward scheme, however, as soon as a window samples a different position than the preceding window, this position must be a new position. Thus, one needs only to consider the context of two consecutive windows of w kmers, for a total of w + k characters.

When a sampling scheme selects a new position for the last window in a context, the context is *charged*. Marçais *et al.* (2017) showed that the density of a scheme *f* can be defined as the proportion of contexts which are charged. In the case of forward schemes, each edge in B_{w+k-1} represents a context, and the charged contexts are the edges (u, v) for which $f(u) \neq f(v) + 1$.

Universal hitting sets. In 2021, Zheng et al. (2021a) related the density of forward and local schemes to the concept of universal hitting sets (UHS). A (w, ℓ) -UHS is defined as a set of ℓ -mers U such that any sequence of w adjacent ℓ -mers must contain at least one ℓ -mer from U. Theorem 1 of Zheng et al. (2021a) showed that when k = 1, one can use the minimum size of a $(w, \ell = w + k)$ -UHS to bound the density of (w, k = 1)-forward schemes, and the minimum size of a $(w, \ell = 2w + k - 2)$ -UHS to bound the density of a (w, k = 1)-local scheme.

3 Theoretical results

In this section, we prove our main result: an improved lower bound on the density of forward sampling schemes. We first generalize some existing theorems to arbitrary w and k(Sections 3.1 and 3.2), after which our main theorem follows in Section 3.3.

3.1 A lower bound on the size of a (w, ℓ) -UHS

We begin by considering a $(w = 2, \ell)$ -UHS. A $(2, \ell)$ -UHS is equivalent to a vertex cover in B_{ℓ} , i.e., a subset of vertices such that each edge in B_{ℓ} is adjacent to at least one vertex in the subset. Lichiardopol (2006) used the fact that for every cycle *C*, at least $\lceil |C|/2 \rceil$ of its vertices must be in a vertex cover, and obtained a lower bound on the size of a vertex cover by partitioning B_{ℓ} into its pure cycles. We naturally extend this argument to obtain a lower bound on the cardinality of a (w, ℓ) -UHS for any $w \ge 2$.

Proposition 4. Let $M_{\sigma}(p)$ count the number of aperiodic necklaces of length p. For any (w, ℓ) -UHS U,

$$|U| \ge \sum_{p|\ell} M_{\sigma}(p) \Big\lceil \frac{p}{w} \Big\rceil.$$

Proof. The pure cycles of C_{ℓ} partition the vertices of B_{ℓ} . For any simple cycle of size p in B_{ℓ} , a (w, ℓ) -UHS must contain at least $\lceil p/w \rceil \ell$ -mers. As there is a one-to-one correspondence between the pure cycles of length $p | \ell$ in B_{ℓ} and the $M_{\sigma}(p)$ aperiodic necklaces of length p, we have

$$|U| \geq \sum_{c \in C_{\ell}} \left\lceil \frac{|c|}{w} \right\rceil = \sum_{p \mid \ell} M_{\sigma}(p) \left\lceil \frac{p}{w} \right\rceil.$$

Figure 1b provides a depiction of a minimum (2, 4)-UHS as well as the pure-cycle partitioning of B_4 on a binary alphabet. Notably, the pure cycle (0011, 0110, 1100, 1001) has three vertices in the UHS, even though the lower bound given by Proposition 4 only requires it have 2. This is an example where the lower bound is not tight.

For certain values of w and ℓ , such as when ℓ is prime or w=2 and ℓ is odd, Proposition 4 can be simplified to remove the summation and ceil function (Supplementary B).

Proposition 4 is the core of the proof of Theorem 1 and already has the right structure. The remainder of this section translates this result on universal hitting sets to a result on the density of sampling schemes.

3.2 A connection between sampling scheme density and UHS size

Zheng *et al.* (2021a, Theorem 1) showed a connection between universal hitting sets and the density of sampling schemes when k = 1. We naturally extend their result to $k \ge 1$ for both local schemes (Lemma 5) and forward schemes (Corollary 6).

Lemma 5. Let f be a (w, k)-local scheme, and let C_f be its corresponding set of charged contexts defined as the set of strings W of length 2w + k - 2 for which the last window W[w - 1, 2w + k - 2) selects a position w - 1 + f(W[w - 1, 2w + k - 1)) not selected by any previous window:

$$C_f := \{ W \in \Sigma^{2w+k-2} | \forall 0 \le i \le w-2, \\ f(W[w-1, 2w+k-2)) + (w-1) \ne f(W[i, i+w+k-1)) + i \}$$

Then, C_f is a (w, 2w + k - 2)-UHS.

Proof. For the sake of a contradiction, suppose there is a walk of length w in the De Bruijn graph of order (2w+k-2), say (W_0, \ldots, W_{w-1}) , that avoids *cf* Let *S* be the *spelling* of the walk, i.e., the sequence of length 3w+k-3 such that $S[i,i+2w+k-2) = W_i$. Since $W_{w-1} \notin C_f$ and *S* contains W_{w-1} , this implies that on the last (w+k-1)-mer of W_{w-1} (i.e. S[2w-2, 3w+k-3)), *f* selects an index $j \ge 2w-2$ in *S* which has already been picked.

Since $0 \le f(\cdot) \le w - 1$ and $j \ge 2w - 2$, the first window that selects position *j* must begin at an index $m \ge w - 1$. Therefore, the context $W_{m-w+1} = S[m - (w-1), m+w+k-1)$ is charged, as *f* selects a previously unselected position when applied to its last (w+k-1)-mer. By definition, $W_{m-w+1} \in C_f$, contradicting our supposition and therefore C_f is a (w, 2w+k-2)-UHS.

Identically, one can consider contexts for a (w, k)-forward scheme f, which requires only verifying that the selection for a window of length w+k-1 is distinct from the selection for the previous window. Therefore, the length of a context for forward f is only w+k. As above, every w contexts must have at least one charged context, leading to the following conclusion:

Corollary 6. If f is a (w, k)-forward scheme and C_f is its corresponding set of charged contexts, defined as $C_f = \{W \in \Sigma^{w+k} | f(W[0, w+k-1)) \neq f(W[1, w+k)) + 1\}$, then C_f is a (w, w+k)-UHS.

As all contexts of a particular length ℓ are equally likely to occur in an infinite, uniform random string, the proportion of charged contexts corresponds to the density of the sampling scheme (Marçais *et al.* 2017), i.e. $d(f) = |C_f|/\sigma^{\ell}$, where $\ell = w + k$ for forward schemes and $\ell = 2w + k - 2$ for local schemes. An example of the charged contexts of a (2, 2)-



Figure 1. (a) A De Bruijn graph B_3 corresponding to a minimum density (w = 2, k = 2)-forward scheme. The underlined characters in each vertex represent the 2-mer that is selected for that window. The solid edges represent the charged contexts and the edge colors represent the pure cycles in B_4 (not in B_3 itself). For characters c_i , each edge $(c_0c_1c_2, c'_0c'_1c'_2)$ in B_3 corresponds to the vertex $c_0c_1c_2c'_2$ in B_4 . (b) The corresponding ($w = 2, \ell = 4$)-UHS in B_4 . The vertices are partitioned by color, representing the pure-cycles. The 2-mer(s) selected in each context are underlined. The vertices with a double border represent the charged edges in B_3 in (a) and the corresponding (2, 4)-UHS. Each pure cycle *c* has at least [|c|/w] vertices in the UHS.

forward scheme and the corresponding UHS is depicted in Fig. 1.

3.3 Lower bounds on local and forward scheme density

We are now ready to state and prove our main theorem.

Theorem 1. Let f be a (w, k)-forward sampling scheme and $M_{\sigma}(p)$ count the number of aperiodic necklaces of length p over an alphabet of size σ . Then, the density of f is at least

$$g_{\sigma}(w,k) := \frac{1}{\sigma^{w+k}} \sum_{p \mid (w+k)} M_{\sigma}(p) \left\lceil \frac{p}{w} \right\rceil \ge \frac{\lceil \frac{w+k}{w} \rceil}{w+k} \ge \frac{1}{w}, \quad (1)$$

where the middle inequality is strict for w > 1.

Proof. Due to Corollary 6 and Marçais *et al.* (2017), we can see that a (w, k)-forward sampling scheme of density d(f) implies a $(w, \ell = w + k)$ -UHS of size $\sigma^{w+k} \cdot d(f)$. By Proposition 4, this implies that every forward sampling scheme has a density of at least $g_{\sigma}(w, k)$, and hence $d(f) \ge g_{\sigma}(w, k)$ follows.

For any p that divides w+k, we have $\lfloor \frac{p}{w} \rfloor \ge \frac{p}{w+k} \lfloor \frac{w+k}{w} \rfloor$, with strict inequality when p=1 and w>1. Substituting this in $g_{\sigma}(w,k)$, the middle inequality follows directly using the identity $\sum_{p|w+k} p \cdot M_{\sigma}(p) = \sigma^{w+k}$ that counts the number of strings of length w+k partitioned by their shortest period. The last inequality follows directly from $\frac{1}{w+k} \left[\frac{w+k}{w}\right] \ge \frac{1}{w+k} \frac{w+k}{w} = 1/w.$ \Box As shown in Section 4, $g_{\sigma}(w,k)$ is a tight bound for many small cases. Since its formula is somewhat unwieldy, $\frac{1}{w+k} \left[\frac{w+k}{w}\right]$ can be used as an approximation that quickly approaches $g_{\sigma}(w,k)$ (Fig. 2). Simple arithmetic shows that both $g_{\sigma}(w,k)$ and $\frac{1}{w+k} \left[\frac{w+k}{w}\right]$ improve the previous lower bound of $\frac{1.5}{w+k-0.5}$ of Groot Koerkamp and Pibiri (2024).

Given a (w, k)-local scheme f_k , we can construct a $(w,k' \ge k)$ -local scheme $f_{k'}$ of the same density by ignoring the last k'-k characters in each window, i.e. $f_{k'}(W) = f_k(W[0...(w+k))$. This directly implies $d(f_k) = d(f_{k'})$ (Zheng *et al.* 2021a). It follows that the minimum density of a (w, k)-local or forward scheme is monotonically decreasing as k increases. However, as can be seen in Fig. 2, $g_{\sigma}(w,k)$ is not a monotonically decreasing function. The local maxima appear to be at $k \equiv 1 \pmod{w}$, which motivates the following improved lower bound.

Theorem 2. For any (w, k)-forward scheme f, an improved lower bound g' is given by

$$d(f) \ge g'_{\sigma}(w,k) := \max(g_{\sigma}(w,k), g_{\sigma}(w,k'))$$
$$\ge \max\left(\frac{1}{w+k} \left\lceil \frac{w+k}{w} \right\rceil, \frac{1}{w+k'} \left\lceil \frac{w+k'}{w} \right\rceil\right),$$

where k' is the smallest integer $\geq k$ such that $k' \equiv 1 \pmod{w}$.

5

Remark 7. Similar to Theorem 1, Lemma 5 implies that any (w, k)-local scheme f has density at least $d(f) \ge g_{\sigma}(w, w+k-2)$. As this bound is in terms of g_{σ} , the improved bound in Theorem 2 can be applied to local schemes as well, i.e., for any (w, k)-local scheme f, an improved lower bound is given by

$$\mathbf{d}(f) \ge \mathbf{g'}_{\sigma}(w, w+k-2).$$

4 Empirical tightness of our bounds

Here, we compare our bounds g_{σ} and g'_{σ} to existing lower bounds. Further, we show how tight these bounds are for small w, k, and σ by searching for optimal schemes via an integer linear programming (ILP) formulation. We also show how close existing sampling scheme densities are to g'_{σ} for practical choices of w, k, and σ . Finally, we show when the recently described mod-minimizer scheme (Groot Koerkamp and Pibiri 2024) achieves optimal density as $\sigma \to \infty$.

ILP description. We use an ILP to search for minimum density forward sampling schemes. In short, we use a single integer variable $x_W \in [w]$ for every window W of length w + k - 1 (corresponding to a vertex in B_{w+k-1}) that indicates the position of the chosen k-mer, and a single boolean variable $y_{(W,W')}$ for each edge in B_{w+k-1} that indicates whether the corresponding context is charged. On each edge, we require that the scheme be forward. The objective function is to minimize the number of charged edges. To reduce the search space, we add an additional constraint corresponding to our lower bound g_{σ} by requiring that for each pure



Figure 2. Comparison of forward scheme lower bounds and optimal densities for small w, k, and σ . Optimal densities were obtained via the ILP and are plotted as black circles that are solid when the optimal density matches our lower bound, g'_{σ} , and hollow otherwise. Each column corresponds to a parameter being fixed to the lowest non-trivial value, i.e., $\sigma = 2$ in the first column, w = 2 in the second column, and k = 1 in the third column. Note that the *x*-axis in the third column corresponds to w, not k.

cycle of length |c| in B_{w+k} , at least $\lceil |c|/w \rceil$ of the corresponding edges in B_{w+k-1} are charged. Further details, including the ILP formulation for local schemes, can be found in Supplementary C.

Comparison against optimal schemes for small k. We used Gurobi (Gurobi Optimization, LLC 2024) to solve the ILP for all combinations of w, k, and σ such that $1 \le w \le 12, 1 \le k \le 12$, and $2 \le \sigma \le 4$ for both forward and local schemes and limited the runtime for each instance to 12h on 128 threads. All results are reported in in Supplementary D, Table S2. While the additional constraint on pure cycles corresponding to g_{σ} significantly sped up the search, for most large w, k, and σ , the ILP failed to terminate with an optimal solution in the allotted time. As a result, we restrict most of our analysis to the following three cases: fixed alphabet size σ = 2, fixed window size w=2, and fixed k-mer size k=1 (Fig. 2).

For all (w,k,σ) where $k \equiv 1 \pmod{w}$ (including when k=1), the minimum density exactly matches our lower bound $g_{\sigma}(w,k)$. Additionally, when $\sigma = 2$ and w=2, the minimum density was equal to $g'_{\sigma}(w,k)$.

Comparison against existing schemes for large k. Using a sequence of 10 million random characters over alphabet size $\sigma = 4$, we approximated the density of recent sampling schemes using the benchmarking implementation from Groot Koerkamp and Pibiri (2024). To compare each density to the particular proportion of selected k-mers on a genomic sequence, we also ran all sampling schemes on the human Y chromosome (Rhie *et al.* 2023) after removing all non-ACTG characters. The densities of the best performing methods, Miniception (Zheng *et al.* 2023), and mod-minimizers (Groot Koerkamp and Pibiri 2024) are plotted in Fig. 3 along with random minimizers and lower bounds.

The ratio between the minimum achieved densities and lower bounds for a selection of (w, k) pairs used by existing *k*-mer-based methods are presented in Table 1. Additional results for $\sigma \in \{2, 256\}$ and $w \in \{2, 50\}$ are provided in in Supplementary D, Fig. S4.

The mod-minimizer has optimal density for large σ when $w \equiv k \pmod{1}$. When w and k are constant and $\sigma \to \infty$, the probability of duplicate characters in a window goes to 0. This implies that we can use t=1 for the mod-minimizer. When $k \equiv t = 1 \pmod{w}$, the density of the mod-minimizer (Theorem 10 of Groot Koerkamp and Pibiri, 2024) is given by

$$\frac{\lfloor \frac{w+k-2}{w} \rfloor + 2}{w+k} + o(1/\ell).$$

The $o(1/\ell)$ term only accounts for duplicate *t*-mers, and hence disappears when $\sigma \to \infty$. We get

$$\frac{\lfloor \frac{w+k-2}{w} \rfloor+2}{w+k} = \frac{\lfloor \frac{k+3w-2}{w} \rfloor}{w+k} = \frac{\lceil \frac{k+2w-1}{w} \rceil}{w+k} \stackrel{k=1 \pmod{w}}{=} \frac{\lceil \frac{w+k}{w} \rceil}{w+k}.$$

Thus, the mod-minimizer has density equal to the lower bound provided by Theorem 1 when σ goes to ∞ and w and $k \equiv 1 \pmod{w}$ are fixed.

In practice, for $\sigma = 256$ the mod-minimizer scheme is within 1% from optimal when $k \equiv 1 \pmod{w}$ (Supplementary D, Fig. S4). When $\sigma = 4$ (Fig. 3), a t > 1must be used, causing the density plot to 'shift right' Existing schemes vs. lower bounds ($\sigma = 4$)



Figure 3. Comparison of existing schemes to lower bounds with practical parameters. Densities are calculated by applying each scheme to a random sequence of 10 million characters over an alphabet of size $\sigma = 4$ (dotted lines) and are compared with the corresponding proportion of sampled *k*-mers on the human Y chromosome (Rhie *et al.* 2023) (soft lines). The mod-minimizer uses parameter r = 4, and miniception uses parameter max(4, k - w). The window sizes 5 and 19 are the default window sizes for Kraken2 (Wood *et al.* 2019) and minimap2 (-ax hifi) (Li 2018), respectively. For SSHash, w = 12 was the window size used when indexing the human genome (Pibiri 2022).

compared to the lower bound. Because of that, the modminimizer does not quite match the lower bound for practical values of σ .

5 Discussion

5.1 Conjecture on when our lower bound is tight

Analytically, it is clear that $g'_{\sigma}(w,k)$ is much larger than $\frac{1}{w}$. In all cases, $g'_{\sigma}(w,k)$ is nearly tight, if not completely. In particular, our bound is tight for all 40 tested parameter sets where $k \equiv 1 \pmod{w}$, leading us to our conjecture:

Conjecture 1. For any w and k satisfying $k \equiv 1 \pmod{w}$, there exists a (w, k)-forward sampling scheme f such that $d(f) = g_{\sigma}(w, k)$.

While the minimum size of a decycling set, i.e., a $(w = \infty, \ell)$ -UHS, is well known to be $N_{\sigma}(\ell)$ (Mykkeltveit 1972), very little is known about the minimum size of a (w, ℓ) -UHS for finite w. In addition to providing the minimum density of a (w, k)-forward scheme for $k \equiv 1 \pmod{w}$, proving Conjecture 1 would also determine the minimum size of a $(w, \ell = w + k)$ -UHS when $k \equiv 1 \pmod{w}$.

5.2 Existing schemes are nearly optimal when $k \ge w$ or σ is large

A natural investigation which follows our proposed lower bound is to determine the gap between $g'_{\sigma}(w,k)$ and current forward scheme densities. Previously, the gap between known densities and lower bound was rather large, making it unclear how much more the density could be reduced.

In Table 1, we observe that existing schemes are already within 11% from the optimal density for practical values of w and k across different applications, and in many cases are even within 3% of the optimal density. In Fig. 3, we see that this difference holds not just for the specific (w, k) in Table 1, but for most $k \ge w$. This is much more informative than the previous lower bound of 1/w, which implied that most current schemes are at most 50% denser than optimal for many of the parameters in Fig. 3.

For alphabets much larger than DNA ($\sigma = 4$), such as the ASCII alphabet ($\sigma = 256$), we observe that when $k \equiv 1 \pmod{w}$, the mod-minimizer scheme recently proposed by Groot Koerkamp and Pibiri (2024) is at most 1% denser than optimal and furthermore, we show that it is optimal as $\sigma \to \infty$. This makes the mod-minimizer scheme the first practical scheme for which there exist finite parameters k and w for which it is close to optimal.

5.3 Tightening the bound for small k

Our new bound for forward schemes always improves over 1/w and appears tight when $k \equiv 1 \pmod{w}$. This leads to an increasingly close bound for $k \not\equiv 1 \pmod{w}$ as k increases, but leaves a large gap when 1 < k < w. A better understanding of these small cases will be necessary to obtain a tight lower bound for all w and k. Based on Supplementary D, Figs. S3 and S4, one might conjecture that the double decyling-setbased methods of Pellow *et al.* (2023) are near-optimal, but subsequent work such as the greedy minimizer (Golan *et al.* 2024) has shown better schemes are possible. From Fig. 2, we already know that our lower bound is not always tight, so this leaves the question:

Open problem 1. How close can practical sampling schemes get to the density given by our lower bound?

5.4 Extending the bound to local schemes

For local schemes, though, our bound appears much less tight. We identified eight sets of (w, k, σ) where local schemes can obtain lower densities than their forward counterparts. In all cases, however, the difference between the local and forward densities was minuscule, with the largest difference of being found for $(w = 4, k = 2, \sigma = 2)$ where the density decreased from 0.375 to 0.371 (Supplementary D, Table S2). Nevertheless, for some parameters, local schemes are able to achieve densities lower than our $g'_{\sigma}(w, k)$ lower bound for forward schemes. Given the trend observed in Supplementary D, Table S2, we arrive at our final open problem:

Open problem 2. How much lower can the density of a(w, k)-local scheme be compared to a(w, k)-forward scheme?

Acknowledgements

We thank Giulio Ermanno Pibiri, Nicolae Sapoval, and our reviewers for their suggestions which helped improve the manuscript.

Supplementary data

Supplementary data are available at *Bioinformatics* online.

Conflict of interest: None declared.

Funding

This work was supported by the National Library of Medicine Training Program in Biomed-ical Informatics and Data Science [T15LM007093 to B.K.]; in part, the National Institute of Allergy and Infectious Diseases [P01-AI152999 to B.K. and T.J.T.], National Science Foundation (NSF) awards [IIS-2239114 and EF-2126387 to T.J.T]. R.G.K. is supported by ETH Research Grant ETH-1721-1 to Gunnar Rätsch.

Data availability

No new data were generated or analysed in support of this research.

References

- DeBlasio D, Gbosibo F, Kingsford C et al. Practical universal k-mer sets for minimizer schemes. In: Proceedings of the 10th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics, Niagara Falls, NY, USA. New York, NY, USA: Association for Computing Machinery, 2019, 167–76. https://doi.org/10.1145/3307339.3342144
- Edgar R. Syncmers are more sensitive than minimizers for selecting conserved *k*-mers in biological sequences. *PeerJ* 2021;9:e10805. https:// doi.org/10.7717/peerj.10805
- Ekim B, Berger B, Orenstein Y. A randomized parallel algorithm for efficiently finding near-optimal universal hitting sets. In: *International Conference on Research in Computational Molecular Biology*, *Padua, Italy*. Cham, Switzerland: Springer International Publishing, 2020, 37–53. https://doi.org/10.1007/978-3-030-45257-5_3
- Golan S, Tziony I, Kraus M et al. Generating low-density minimizers. bioRxiv, 2024. https://doi.org/10.1101/2024.10.28.620726, preprint: not peer reviewed.
- Groot Koerkamp R, Pibiri GE. The mod-minimizer: a simple and efficient sampling algorithm for long k-mers. In: Pissis SP, Sung WK (eds), 24th International Workshop on Algorithms in Bioinformatics (WABI 2024), volume 312 of Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl—Leibniz-Zentrum für Informatik, 2024, 11:1–11:23. ISBN 978-3-95977-340-9. https://doi.org/10.4230/LIPIcs.WABI. 2024.11.
- Gurobi Optimization, LLC. *Gurobi Optimizer Reference Manual.* 2024. https://www.gurobi.com (September 2024, date last accessed).
- Hoang M, Zheng H, Kingsford C. Differentiable learning of sequencespecific minimizer schemes with DeepMinimizer. J Comput Biol 2022;29:1288–304. https://doi.org/10.1089/cmb.2022.0275
- Irber L, Brooks PT, Reiter T *et al.* Lightweight compositional analysis of metagenomes with FracMinHash and minimum metagenome covers. *bioRxiv*, 2022-01. 2022, https://doi.org/10.1101/2022.01. 11.475838, preprint: not peer reviewed.
- Kille B, Garrison E, Treangen TJ et al. Minmers are a generalization of minimizers that enable unbiased local Jaccard estimation. *Bioinformatics* 2023;39:btad512. https://doi.org/10.1093/bioinfor matics/btad512
- Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics 2018;34:3094–100. https://doi.org/10.1093/bioin formatics/bty191
- Lichiardopol N. Independence number of de Bruijn graphs. *Discrete* Math 2006;306:1145-60. https://doi.org/10.1016/j.disc.2005. 10.032
- Loukides G, Pissis S. Bidirectional string anchors: a new string sampling mechanism. In: ESA 2021-29th Annual European Symposium on Algorithms, Held online due to COVID-19, Vol. 204. Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2021, 1–21. https://doi.org/10.4230/lipics.esa.2021.64

Loukides G, Pissis SP, Sweering M. Bidirectional string anchors for improved text indexing and top-k similarity search. IEEE Trans Knowl Data Eng 2023;35:11093–111. https://doi.org/10.1109/ tkde.2022.3231780

8

- Marçais G, Pellow D, Bork D et al. Improving the performance of minimizers and winnowing schemes. Bioinformatics 2017;33:i110–7. https://doi.org/10.1093/bioinformatics/btx235
- Marçais G, DeBlasio D, Kingsford C. Asymptotically optimal minimizers schemes. *Bioinformatics* 2018;34:i13–22. https://doi.org/10. 1093/bioinformatics/bty258
- Moreau C. Sur les permutations circulaires distinctes. *Nouvelles Ann Math* 1872;11:309–14.
- Mykkeltveit J. A proof of Golomb's conjecture for the de Bruijn graph. J. Comb Theory B 1972;13:40–5. https://doi.org/10.1016/0095-8956(72)90006-8
- Möbius A. Über eine besondere art von umkehrung der reihen. J. Comb Theory B 1832;1832:105–23. https://doi.org/10.1515/crll.1832.
 9.105
- Ndiaye M, Prieto-Baños S, Fitzgerald LM *et al.* When less is more: sketching with minimizers in genomics. *Genome Biol* 2024;25:270. https://doi.org/10.1186/s13059-024-03414-4
- Orenstein Y, Pellow D, Marçais G et al. Compact universal k-mer hitting sets. In: Algorithms in Bioinformatics: 16th International Workshop, WABI 2016, Aarhus, Denmark, August 22–24, 2016. Proceedings 16. Cham, Switzerland: Springer International Publishing, 2016, 257–68. https://doi.org/10.1007/978-3-319-43681-4_21
- Pellow D, Pu L, Ekim B et al. Efficient minimizer orders for large values of k using minimum decycling sets. Genome Res 2023;33:1154–61. https://doi.org/10.1101/gr.277644.123
- Pibiri GE. Sparse and skew hashing of k-mers. Bioinformatics 2022;38: i185–94. ISSN 1367-4811. https://doi.org/10.1093/bioinformat ics/btac245

- Rhie A, Nurk S, Cechova M *et al.* The complete sequence of a human Y chromosome. *Nature* 2023;621:344–54. https://doi.org/10.1038/s41586-023-06457-y
- Riordan J. The combinatorial significance of a theorem of Pólya. J Soc Industrial Appl Math 1957;5:225–37. ISSN 2168-3484. https://doi. org/10.1137/0105015.
- Roberts M, Hunt BR, Yorke JA et al. A preprocessor for shotgun assembly of large genomes. J Comput Biol 2004;11:734–52.
- Schleimer S, Wilkerson DS, Aiken A. Winnowing: local algorithms for document fingerprinting. In: Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data, San Diego, CA, USA. New York, NY, USA: Association for Computing Machinery, 2003, 76–85. https://doi.org/10.1145/872757.872770
- Shaw J, Yu YW. Theory of local k-mer selection with applications to long-read alignment. *Bioinformatics* 2022;38:4659–69. https://doi. org/10.1093/bioinformatics/btab790
- Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome Biol 2019;20:257. https://doi.org/10.1186/ s13059-019-1891-0
- Zheng H, Kingsford C, Marçais G. Improved design and analysis of practical minimizers. *Bioinformatics* 2020;36:i119–27. https://doi. org/10.1093/bioinformatics/btaa472
- Zheng H, Kingsford C, Marçais G. Lower density selection schemes via small universal hitting sets with short remaining path length. J Comput Biol 2021a;28:395–409. https://doi.org/10.1089/cmb. 2020.0432
- Zheng H, Kingsford C, Marçais G. Sequence-specific minimizers via polar sets. *Bioinformatics* 2021b;37:i187–95. https://doi.org/10.1093/ bioinformatics/btab313
- Zheng H, Marçais G, Kingsford C. Creating and using minimizer sketches in computational genomics. J Comput Biol 2023;30: 1251–76. https://doi.org/10.1089/cmb.2023.0094.

Downloaded from https://academic.oup.com/bioinformatics/article/41/1/btae736/7922553 by guest on 15 January 2025

ec

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited. Bioinformatics, 2024, 41, 1–8 https://doi.org/10.1093/bioinformatics/btae736 Original Paper