

Estimation of Substitution and Indel Rates via k -mer Statistics

Mahmudur Rahman Hera ✉ 

School of Electrical Engineering and Computer Science, The Pennsylvania State University,
University Park, PA, USA

Paul Medvedev ✉ 

Department of Computer Science and Engineering, The Pennsylvania State University,
University Park, PA, USA

Department of Biochemistry and Molecular Biology, The Pennsylvania State University,
University Park, PA, USA

Huck Institutes of the Life Sciences, The Pennsylvania State University, University Park, PA, USA

David Koslicki ✉ 

School of Electrical Engineering and Computer Science, The Pennsylvania State University,
University Park, PA, USA

Department of Biology, The Pennsylvania State University, University Park, PA, USA

Huck Institutes of the Life Sciences, The Pennsylvania State University, University Park, PA, USA

Antonio Blanca ✉ 

School of Electrical Engineering and Computer Science, The Pennsylvania State University,
University Park, PA, USA

Abstract

Methods utilizing k -mers are widely used in bioinformatics, yet our understanding of their statistical properties under realistic mutation models remains incomplete. Previously, substitution-only mutation models have been considered to derive precise expectations and variances for mutated k -mers and intervals of mutated and non-mutated sequences. In this work, we consider a mutation model that incorporates insertions and deletions in addition to single-nucleotide substitutions. Within this framework, we derive closed-form k -mer-based estimators for the three fundamental mutation parameters: substitution, deletion rate, and insertion rates. We provide theoretical guarantees in the form of concentration inequalities, ensuring accuracy of our estimators under reasonable model assumptions. Empirical evaluations on simulated evolution of genomic sequences confirm our theoretical findings, demonstrating that accounting for insertions and deletions signals allows for accurate estimation of mutation rates and improves upon the results obtained by considering a substitution-only model. An implementation of estimating the mutation parameters from a pair of fasta files is available here: github.com/KoslickiLab/estimate_rates_using_mutation_model.git. The results presented in this manuscript can be reproduced using the code available here: github.com/KoslickiLab/est_rates_experiments.git.

2012 ACM Subject Classification Applied computing → Computational biology; Theory of computation → Theory and algorithms for application domains; Mathematics of computing → Probabilistic inference problems

Keywords and phrases k -mers, mutation rate, indel, alignment-free, estimation, substitution, insertion, deletion

Digital Object Identifier 10.4230/LIPIcs.WABI.2025.16

Related Version *Full version with proofs:* <https://doi.org/10.1101/2025.05.14.653858>

Supplementary Material

Software: https://github.com/KoslickiLab/estimate_rates_using_mutation_model.git

Workflow: https://github.com/KoslickiLab/est_rates_experiments.git



© Mahmudur Rahman Hera, Paul Medvedev, David Koslicki, and Antonio Blanca;
licensed under Creative Commons License CC-BY 4.0

25th International Conference on Algorithms for Bioinformatics (WABI 2025).

Editors: Broňa Brejová and Rob Patro; Article No. 16; pp. 16:1–16:15

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

Funding This material is based upon work supported by the National Science Foundation under Grant No. DBI2138585. Research reported in this publication was supported by the National Institute Of General Medical Sciences of the National Institutes of Health under Award Number R01GM146462. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

1 Introduction

Estimating the mutation rate between two evolutionarily related sequences is a classical question in molecular evolution, with roots that pre-date the genomics era [20]. Early quantitative efforts focused on amino-acid substitution: the seminal PAM matrices of Dayhoff et al. converted curated alignments of close homologues into an evolutionary time-scale [3], while the BLOSUM series by Henikoff and Henikoff mined ungapped blocks of conserved proteins to improve sensitivity for more diverged sequences [6]. These approaches and later profile-based HMM models [4] were derived from pairwise or multiple alignments and remain the gold standard when accurate alignments are available.

Over the last decade, however, high-throughput sequencing has shifted the scale of comparative genomics from dozens to millions of genomes, rendering high computational cost pipelines (e.g. quadratic-time) increasingly impractical. Consequently, alignment-free techniques that summarize sequences by inexpensive statistics have become indispensable [17, 21]. These approaches most commonly utilize k -mer sets and sketches thereof. Popular tools such as Mash [10], Skmer [13] and more recent sketch-corrected frameworks like Sylph [16] and FracMinHash-based methods [7, 8, 11, 15] can build whole-genome phylogenies, screen metagenomic samples, and compute millions of pairwise point-mutation rate estimates in minutes rather than days.

Despite their empirical success, theoretical understanding of alignment-free estimators has lagged behind practice. Nearly all existing models treat evolution as a pure-substitution process, ignoring insertions and deletions (indels), or else their performance in the presence of indels is often not thoroughly evaluated [12]. When indels are frequent, substitution-only estimators systematically inflate divergence and can misplace taxa – even in otherwise well-resolved trees of primates constructed from k -mer Jaccard similarities [10]. Recent work has quantified how k -mer-based statistics are affected by substitutions and are also used to estimate substitution-only mutation rates [8], yet a principled treatment that jointly infers substitution and indel parameters from k -mer statistics is still absent. This omission is particularly significant because indels represent a substantial fraction of genomic variation and play crucial roles in evolution [19]. Such indels cause substitution-only k -mer-based methods to underperform, as just like with substitutions, disruption of k -mer content by indel events affects at least k k -mers, often leading to overestimates of mutation rates [2, 8].

In this paper we introduce the first closed-form, alignment-free estimators for the three fundamental mutation parameters: substitution rate p_s , deletion rate p_d , and mean insertion length d under a model that explicitly incorporates single-nucleotide substitutions, deletions, and geometrically-distributed insertions. Starting from elementary counts of unmutated and single-deletion k -mers, we derive algebraic expressions for p_s , p_d , and d and prove a sub-exponential concentration bound that guarantees a strong form of consistency as the sequence length grows as detailed in our main contribution: Theorem 9. Simulations on synthetic and real bacterial genomes demonstrate that modeling indels yields markedly more accurate distance estimates than substitution-only approaches. The remainder of this paper is organized as follows: Section 2 defines our mutation model, Section 3 derives basic statistical

properties, Section 4 presents our estimators, Section 5 provides theoretical guarantees about these estimators, Section 6 details our implementation, Section 7 presents experimental results, and Section 8 concludes with a discussion of implications and future directions.

2 The mutation model

We first describe the mutation model under consideration. We use the model from [14]. This model is a type of indel channel and various variations of it have been used to model sequence evolution (e.g. [5]). Let S be a string over the alphabet $\{A, C, G, T\}$, and let L be the number of characters in S . Let S_i denote the i -th character in S where $1 \leq i \leq L$. We define four operations on a string S as follows:

- **Substitute(i)**: select a character $c \in \{A, C, G, T\} \setminus \{S_i\}$ uniformly at random and replace S_i by c .
- **Delete(i)**: remove the character S_i .
- **Stay(i)**: do nothing.
- **Insert(i, s)**: insert string s between S_{i-1} and S_i if $i > 0$. If $i = 0$, prepend s to the start of S instead.

The mutation process takes as input a string S and three parameters p_s , p_d , and d , where $0 \leq p_s, p_d < 1$, $p_s + p_d < 1$, and $d \geq 0$. Then,

1. For each i , let a_i be the operation to be performed at position i . Then, $a_i = \text{Sub}$ with probability p_s , $a_i = \text{Del}$ with probability p_d , or $a_i = \text{Stay}$ with the remaining probability $1 - p_s - p_d$. (These operations are not performed at this point, but only recorded.)
2. Let $track$ be a function mapping from a position in the original string S to its position in the modified S . Initially, $track(i) = i$ for all i . We assume that $track$ is updated accordingly whenever an insertion or deletion operation is performed.
3. For each i that is a substitute action, apply **Substitute(i)**.
4. For each i , let $l_i \geq 0$ be a sample from a geometric distribution with mean d . Then generate a random string Q_i of length exactly l_i by drawing each character in Q_i from $\{A, C, G, T\}$ independently and uniformly at random. This is equivalent to sampling a uniformly random string among all strings of length l_i with characters in $\{A, C, G, T\}$. If $l_i = 0$, Q_i is the empty string.
5. For every i , execute **Insert($track(i), Q_i$)**.
6. For every i such that $a_i = \text{Del}$, execute **Delete($track(i)$)**.
7. Return the resulting S' .

3 Preliminary statistics of the mutation process

In this section, we define several quantities whose statistics we use to design an estimator of the parameters of the mutation process. For the remainder of the manuscript, we assume that the mutation process described in Section 2 is applied to string S of length L , with the unknown parameters p_s , p_d , and d , and a mutated string S' is returned. The theoretical results we present are centered around the concept of k -spans. We define \mathcal{K}_i , the k -span at position $1 \leq i \leq L - k + 1$ as the range of integers $[i, i + k - 1]$ (inclusive of the endpoints of the range). In simpler terms, a k -span captures the interval of a k -mer. We assume that the string S has at least k nucleotides, and therefore, at least one k -mer (and at least one k -span). For the sake of theoretical rigor, we use k -spans to develop our results, and later discuss how our practical implementation uses k -mers to estimate the substitution and indel rates.

► **Lemma 1.** *Let L' be the length of S' . Then*

$$\mathbf{E}[L'] = L(1 + d - p_d), \quad (1)$$

$$\mathbf{Var}[L'] = L(d(d+1) + p_d(1 - p_d)).$$

► **Lemma 2.** *Let f_A and f_A' be the number of “A”s in S and S' , respectively. Then,*

$$\mathbf{E}[f_A'] = f_A(1 - p_s - p_d) + \frac{p_s(L - f_A)}{3} + \frac{dL}{4}. \quad (2)$$

We can analogously obtain expressions for the expectations of the corresponding counts of “C”s, “G”s, and “T”s in S' , denoted by f_C' , f_G' , and f_T' , respectively.

Next, we say that a k -span has no mutations if (a) for all k positions in the k -span, the mutation process picks a Stay(\cdot) action, and if (b) for all $k - 1$ intermediate positions in the k -span, the mutation process does not insert anything. Let \mathcal{N} be the number of such k -spans. Let \mathcal{K} be the number of all k -spans in S .

► **Lemma 3.** $\mathbf{E}[\mathcal{N}] = \mathcal{K} (1 - p_s - p_d)^k \frac{1}{(d+1)^{k-1}}.$

The last quantities we will use correspond to the number of k -spans with a single kind of mutation. Let \mathcal{S} be the number of k -spans with a single substitution, and no other mutations, let \mathcal{D} be the number of k -spans with a single deletion and no other mutations, and let \mathcal{I} be the number of k -spans with a single insertion and no other mutations.

► **Lemma 4.** *The expectations of \mathcal{S} , \mathcal{D} , and \mathcal{I} are given by:*

$$\mathbf{E}[\mathcal{S}] = \mathcal{K} k (1 - p_s - p_d)^{k-1} p_s \frac{1}{(d+1)^{k-1}}. \quad (3)$$

$$\mathbf{E}[\mathcal{D}] = \mathcal{K} k (1 - p_s - p_d)^{k-1} p_d \frac{1}{(d+1)^{k-1}}. \quad (4)$$

$$\mathbf{E}[\mathcal{I}] = \mathcal{K} (k-1) (1 - p_s - p_d)^k \frac{d}{(d+1)^k}. \quad (5)$$

4 Estimating p_s , p_d and d

We describe next our estimators for the parameters of the mutation model p_s , p_d , and d . We derive estimators for these rates based on the statistics defined in Section 3. Recall that S is the known input string to the mutation process, and S' is the resulting random string or observation.

Our mutation model assumes $p_s + p_d < 1$, and $d \geq 0$. In addition, we assume that $L \geq k$ which implies that $\mathbf{E}[\mathcal{N}] > 0$. In reality, strings are typically much longer than k -mer sizes and therefore this is a reasonable assumption.

From Lemmas 3 and 4 we get

$$\frac{\mathbf{E}[\mathcal{D}]}{\mathbf{E}[\mathcal{N}]} = \frac{k p_d}{1 - p_s - p_d}. \quad (6)$$

From (1), (2) and (6), we obtain the following system of linear equations with variables p_s , p_d , and d .

$$-p_d + d = \frac{\mathbf{E}[L']}{L} - 1 \quad (7)$$

$$\frac{L - 4f_A}{3} p_s - f_A p_d + \frac{L}{4} d = \mathbf{E}[f_A'] - f_A \quad (8)$$

$$\mathbf{E}[\mathcal{D}]p_s + (k\mathbf{E}[\mathcal{N}] + \mathbf{E}[\mathcal{D}])p_d = \mathbf{E}[\mathcal{D}]. \quad (9)$$

Solving this system of equations, we obtain that:

$$p_s = 3 \frac{k \mathbf{E}[\mathcal{N}] \left(\mathbf{E}[L'] - L + 4f_A - 4\mathbf{E}[f_A'] \right) + \mathbf{E}[\mathcal{D}] \left(\mathbf{E}[L'] - 4\mathbf{E}[f_A'] \right)}{(4f_A - L) (\mathbf{E}[\mathcal{D}] + 4k\mathbf{E}[\mathcal{N}])},$$

$$p_d = \frac{\mathbf{E}[\mathcal{D}] (4f_A + 12\mathbf{E}[f_A'] - L - 3\mathbf{E}[L'])}{(4f_A - L) (\mathbf{E}[\mathcal{D}] + 4k\mathbf{E}[\mathcal{N}])},$$

$$d = \frac{\mathbf{E}[L']}{L} - 1 + \frac{\mathbf{E}[\mathcal{D}] (-4f_A - 12\mathbf{E}[f_A'] + L + 3\mathbf{E}[L'])}{(-4f_A + L) (\mathbf{E}[\mathcal{D}] + 4k\mathbf{E}[\mathcal{N}])}.$$

Given S and S' , we can compute L' and f_A' . Let us also assume that we know \mathcal{N} and \mathcal{D} (we will discuss how to find these in Section 6). By replacing the expectations above with these observations, we obtain our estimators for the parameters of the mutation model. That is,

$$\hat{p}_s = 3 \frac{k \mathcal{N} (L' - L + 4f_A - 4f_A') + \mathcal{D} (L' - 4f_A')}{(4f_A - L) (\mathcal{D} + 4k\mathcal{N})}, \quad (10)$$

$$\hat{p}_d = \frac{\mathcal{D} [3(L' - 4f_A') + (L - 4f_A)]}{(L - 4f_A) (\mathcal{D} + 4k\mathcal{N})}, \quad (11)$$

$$\hat{d} = \frac{L'}{L} - 1 + \frac{\mathcal{D} [3(L' - 4f_A') + (L - 4f_A)]}{(L - 4f_A) (\mathcal{D} + 4k\mathcal{N})}. \quad (12)$$

We briefly comment on our choice of estimators for the mutation model parameters, as various statistical approaches based on a different set of observables could yield a different set of estimators. For example, we considered a variant of the estimators based on the counts of k -spans with a single insertion, deletion, or substitution (i.e., \mathcal{I} , \mathcal{S} , and \mathcal{D}). These quantities contain enough information to estimate the mutation parameters, specifically, by solving the non-linear system of equations given by (3), (4), and (5); the resulting estimators performed quite well in real data. However, establishing theoretical guarantees for these estimators proved challenging, as they were defined as roots of degree- k polynomials. Our current estimators address this theoretical limitation as they involve solely linear equations. As we shall see in Section 7, the performance of our estimators in real data is strong, and they strike a more favorable balance by offering both reasonable accuracy and rigorous theoretical guarantees.

5 Concentration results

In this section, we provide a theoretical guarantee for the estimator (10) of the substitution rate. In particular, we show that our estimate of p_s is not only a consistent estimator, but also that it is tightly concentrated in a symmetric interval around the true value of p_s . Similar techniques can likely be used to prove the consistency and concentration of p_d and d , though we do not do so in this paper. We provide the bias analysis for \hat{p}_s by proving asymptotically tight concentration bounds for \mathcal{N} , \mathcal{D} , and $L' - 4f_A'$.

The concentration bounds for \mathcal{N} and \mathcal{D} stem from the fact that, from the perspective of the mutation process k -spans located more than k positions apart are independent of each other.

► **Lemma 5.** For any $\delta \in (0, 1)$:

$$\Pr \left[\left| \mathcal{N} - \mathbf{E}[\mathcal{N}] \right| \geq \delta \mathbf{E}[\mathcal{N}] \right] \leq 3k \exp \left\{ -\frac{\delta^2 \mathbf{E}[\mathcal{N}]}{3k} \right\}.$$

► **Lemma 6.** For any $\delta \in (0, 1)$:

$$\Pr \left[\left| \mathcal{D} - \mathbf{E}[\mathcal{D}] \right| \geq \delta \mathbf{E}[\mathcal{D}] \right] \leq 3k \exp \left\{ -\frac{\delta^2 \mathbf{E}[\mathcal{D}]}{3k} \right\}.$$

We prove a similar result for $P = L' - 4f_A'$ by noting that P can be expressed as a sum of independent sub-exponential random variables. Our concentration bound for P is as follows.

► **Lemma 7.** Let $J_1 = \ln 2 / \min\{d + 1, 8\}$. Then, there exist absolute constants $c_1, c_2 > 0$ such that the following holds for any $\delta > 0$:

$$\Pr \left[|P - \mathbf{E}[P]| \geq 3\delta \right] \leq 2 \exp \left\{ -\frac{\delta^2}{8f_A} \right\} + 2 \exp \left\{ -\frac{\delta^2}{8(L - f_A)} \right\} + 2 \exp \left\{ -c_1 \min \left(\frac{\delta^2}{c_2^2 J_1^2}, \frac{\delta}{c_2 J_1} \right) \right\}.$$

Finally, we prove that L' is also strongly concentrated around $\mathbf{E}[L']$.

► **Lemma 8.** For any $\delta \in (0, 1)$:

$$\Pr \left[\left| L' - \mathbf{E}[L'] \right| \leq \delta(Ld + Lp_d) \right] \geq 2 \exp \left\{ -\frac{L\delta^2 d^2}{2(d + 1 - \delta d)(d + 1)} \right\} + 2 \exp \left\{ -\frac{Lp_d \delta^2}{3} \right\}.$$

We can piece the results of these lemmas together and prove the following result.

► **Theorem 9.** Suppose $4f_A < L$ and $\frac{4}{3}p_s + p_d < 1$. Then, for sufficiently small $\delta > 0$, there exists constants $c_1, c_2 > 0$ such that

$$\begin{aligned} \Pr[|\hat{p}_s - p_s| \geq 12\delta] &\leq 8k \exp \left\{ -\frac{\delta^2 \mathbf{E}[\mathcal{N}]}{3k} \right\} + 6k \exp \left\{ -\frac{\delta^2 \mathbf{E}[\mathcal{D}]}{3k} \right\} + 2 \exp \left\{ -\frac{\delta^2 \mathbf{E}[P]^2}{72f_A} \right\} \\ &\quad + 2 \exp \left\{ -\frac{\delta^2 \mathbf{E}[P]^2}{72(L - f_A)} \right\} + 2 \exp \left\{ -c_1 \min \left(\frac{\delta^2 \mathbf{E}[P]^2}{c_2^2}, \frac{\delta \mathbf{E}[P]}{c_2} \right) \right\}. \end{aligned}$$

The requirement that $4f_A < L$ in this theorem does not restrict generality: aside from equal nucleotide frequency (where the estimators are naturally invalid), at least one character $c \in \{A, C, T, G\}$ must satisfy $4f_c \leq L$. In addition, the assumption that $\frac{4}{3}p_s + p_d < 1$ holds when p_s and p_d are small (e.g., $p_s < 1/4$ and $p_d < 1/4$ suffices) which is the case most frequently encountered in practice.

Theorem 9 is our central theoretical result, establishing that the estimators developed in Section 4 are sound under the reasonable conditions. Before discussing our implementation and presenting the experimental results, we comment on the error probability in Theorem 9. This probability is small when each of the terms in the sum are small. Since each of these terms decays (at least) exponentially with δ times an expectation that grows linearly with the length of the string (assuming the mutation parameters are fixed), they will all generally be small.

6 Implementation details

Our estimators for the three rates require counting the number of k -spans with single deletion, \mathcal{D} , and the number of k -spans with no mutation, \mathcal{N} . Counting \mathcal{D} and \mathcal{N} can be challenging, particularly because k -mers do not contain the contextual information, and so we do not have access to their corresponding k -spans. An additional layer of complexity comes into play from the fact that identifying a k -mer with no mutation (or a single kind of mutation) is more difficult, considering many edge cases that may arise from inserting the same character that has been deleted. These challenges are circumvented when we use k -spans, and therefore, counting \mathcal{D} and \mathcal{N} solely from k -mers is not trivial and can be considered an interesting problem in and of itself. We therefore implemented an ad hoc solution to estimate \mathcal{D} and \mathcal{N} given the two strings S and S' . The steps for estimating \mathcal{D} and \mathcal{N} are as follows.

We start by extracting all k -mers in S , and building a de Bruijn graph using these k -mers using the cuttlefish tool [9]. We then extract the unitigs from this graph. Let the set of unitigs computed from S be U . We also compute the unitigs in S' in a similar manner and call this U' . We next take an arbitrary unitig u from U , and align every unitig in U' with u . To allow for partial overlap, we use semi-global alignment by using the infix option in edlib [18], which makes sure that gaps at the beginning and at the end of the alignment are not penalized. For a particular $u \in U$, we align every $u' \in U'$ to make sure all relevant alignments are considered. We use these alignments to look at all windows of length k , and count \mathcal{D} and \mathcal{N} accordingly for u . We repeat this for all $u \in U$, and accumulate the measurements from individual u 's into a single global count.

The motivation behind using unitigs is that if there is an isolated mutation, and if the mutation is in the first or the last position of a k -mer, then there is no way to understand if the mutation is a substitution, an insertion, or a deletion only from the k -mers. The only way to resolve this ambiguity (and other similar ambiguities) is to scan beyond the context of k characters – and unitigs are a natural way to do this. The core goal of our implementation of these steps described above was not to make it efficient, but rather to obtain a working solution. We found that executing these steps estimates \mathcal{D} and \mathcal{N} reasonably well, and the estimated rates are also acceptable. As such, we leave finding an efficient way to compute \mathcal{D} and \mathcal{N} as an open research question.

7 Experiments and results

In this section, we present a series of experiments to evaluate the performance of the estimators detailed in Section 4. As discussed earlier, these estimators are sensitive to several input parameters, including k -mer size, sequence length, and the fraction of “A” characters in the sequence. Sections 7.1 through 7.3 explore the sensitivity of the estimators with respect to these parameters. In Sections 7.4 and 7.5, we estimate mutation rates across a wide range of known rate combinations. And finally, in Section 7.6, we demonstrate that our estimated substitution rate outperforms estimates obtained under a substitution-only mutation model. For the experiments in Sections 7.1 through 7.4, the original sequence is a randomly generated synthetic sequence. In these cases, we compute the number of k -spans containing a single deletion and the number of k -spans with no mutation directly from the known mutation process. For the experiments in Sections 7.5 and 7.6, we use real reference genomes as the original sequences. In these cases, the two types of k -span counts are estimated using the steps described in Section 6.

7.1 Sensitivity of the estimators to k -mer lengths

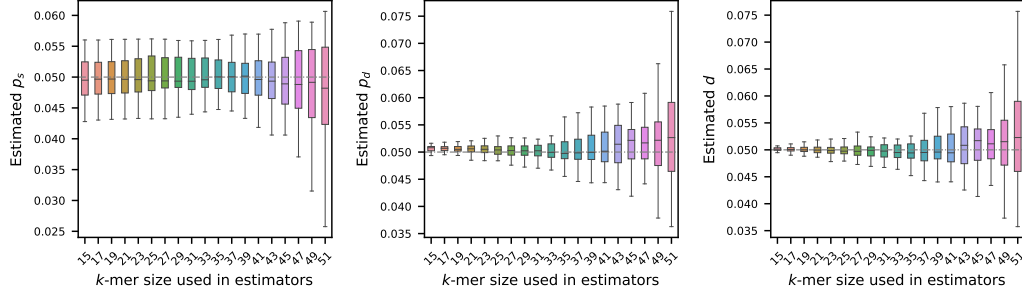


Figure 1 Effect of k -mer length on mutation rate estimation (true rates set to 0.05). A synthetic genome of 1 million nucleotides, mutated genomes were generated by setting p_s , p_d , and $d = 0.05$ – shown by the gray dashed horizontal line. Estimated rates were then computed using a range of k -mer sizes. Each boxplot shows the variability in estimation across 20 simulations, with error bars showing one standard deviation. The plots show that the estimators become more accurate and more precise for shorter k -mers.

We begin our analysis by examining how the choice of k -mer length affects the estimation of the mutation rates. To investigate this effect, we first generated a synthetic reference sequence of 1 million nucleotides, randomly sampling bases with fixed frequencies: 30% “A”, and equal proportions of “C”, “G”, and “T” – making sure total frequency is 100%. From this reference, we simulated 20 mutated sequences, independently from each other, using the mutation model described in Section 2. For each of these mutated sequences, we estimated mutation rates using the estimators defined in Section 4 using a range of values for k . The results of this analysis are summarized in Figure 1.

As illustrated, the choice of k has a substantial impact on the stability of the estimators. In particular, longer k -mers tend to produce estimates with higher variability. This behavior is consistent with the known sensitivity of k -mers to mutations: since a single mutation can disrupt up to k consecutive k -mers, the longer the k -mer, the more susceptible it becomes to such perturbations. Our theoretical result in Theorem 9 also captures this: with a larger k , the error probabilities become larger, and the probabilistic guarantee for the estimators’ performances decreases accordingly.

Interestingly, for the estimator of substitution probability p_s , we observe that the variability in the estimated values does not change significantly from 15 to 39. The reason behind this behavior is not immediately clear and warrants further investigation. It is possible that incorporating the number of k -spans with a single substitution \mathcal{S} into the estimators may correct this behavior, but additional analyses are required to substantiate this hypothesis.

7.2 Sensitivity of the estimators to sequence length

To investigate how the length of the original sequence S influences estimation of the mutation rates, we simulated synthetic genomes ranging from 10K to 1M nucleotides in length. For each genome length, we generated 10 independent synthetic sequences to capture variability due to random sampling. The nucleotide composition of each sequence was fixed, with the frequency of “A” set to 30% and frequencies of “C”, “G”, and “T” set equally – making sure total frequency is 100%. For each synthetic sequence, we generated its mutated version by

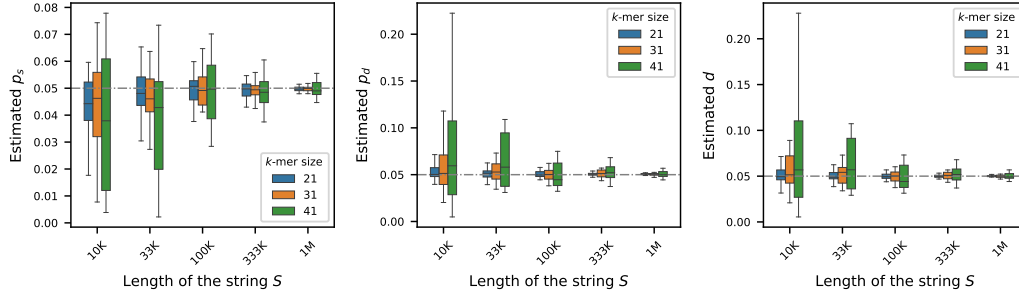


Figure 2 Effect of sequence length on mutation rate estimation (true rates set to 0.05). For synthetic genomes of varying lengths, mutated genomes were generated by setting p_s , p_d , and $d = 0.05$ – shown by the gray dashed horizontal line. Estimated rates were computed using three k -mer sizes: 21, 31, and 41. Each boxplot shows the variability in estimation across 20 simulations. The plots show that the estimators become more accurate and more precise for longer sequences.

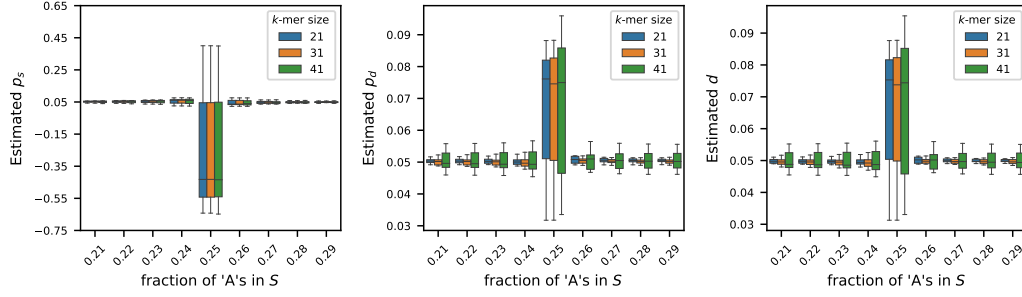
running the mutation process described in Section 2, setting each of p_s , p_d , and d to 0.05. We then estimated the mutation rates using the estimators outlined in Section 4 for three k -mer sizes: $k = 21, 31, 41$.

Figure 2 displays the estimated rates across the varying sequence lengths. As shown, the estimators are less stable for shorter sequences. However, with longer sequences, the estimators yield more accurate results – a trend expected from our core theoretical result in Theorem 9, which states that the associated error is asymptotically vanishing in L , the length of the string S : as L increases, the number of k -spans \mathcal{K} increases, and therefore the probability of error decreases, leading to a more precise estimation.

7.3 Sensitivity of the estimators to base composition

We generated synthetic genomes of 1M nucleotides to investigate how the fraction of “A” characters affects the estimation of the mutation rates. We varied the fraction of “A”s from 21% to 29% in increments of 1%. For each fraction of “A”s, we set the frequency of “C”s, “G”s, and “T”s equally. For each preset fraction of “A”s, we generated 10 random genomes to capture stochastic variation. For each of these genomes, we generated its mutated version using the mutation process described in Section 2, setting each of p_s , p_d , and d to 0.05. We then estimated the mutation rates using the estimators described in Section 4 using three k -mer sizes: $k = 21, 31$, and 41. Figure 3 shows the sensitivity of the estimators to the fraction of “A”s in the original string S .

We observe that the estimators \hat{p}_s , \hat{p}_d , and \hat{d} work reasonably well to estimate the true rates when the frequency of “A” characters, f_A is not $L/4$. On the contrary, when the fraction of “A”s is exactly 25% in the original string S , the estimator gives inaccurate values, some of which are even negative (see estimated values of p_s). This behavior is captured in Theorem 9: when $L - 4f_A = 0$, $\mu_P = \mathbf{E}[P] = 0$, and the probabilistic guarantees become unbounded. We only get meaningful guarantees of consistency when $4f_A$ is strictly smaller than L . While Theorem 9 does not guarantee consistency when $4f_A > L$, this does not restrict generality, as explained in Section 4, and as demonstrated by the estimators’ performances when the fraction of “A”s is larger than $1/4$.



■ **Figure 3** Effect of nucleotide composition on mutation rate estimation (true rates set to 0.05). For synthetic genomes of length $L = 1\text{M}$, the fraction of “A”s is varied from 0.21 to 0.29, and the frequencies of “G”s, “C”s, and “T”s are set equally. For each setting, the mutated string S' was generated by setting $p_s, p_d, d = 0.05$. Estimated rates were computed for three k -mer sizes: 21, 31, and 41. Each boxplot shows the variability in estimation across 20 simulations. The results show that the estimators generally work well for all three k -mer sizes, except when $f_A \approx L/4$, in which case the estimators become unstable – as predicted by Theorem 9.

7.4 Estimating rates from a randomly generated synthetic sequence

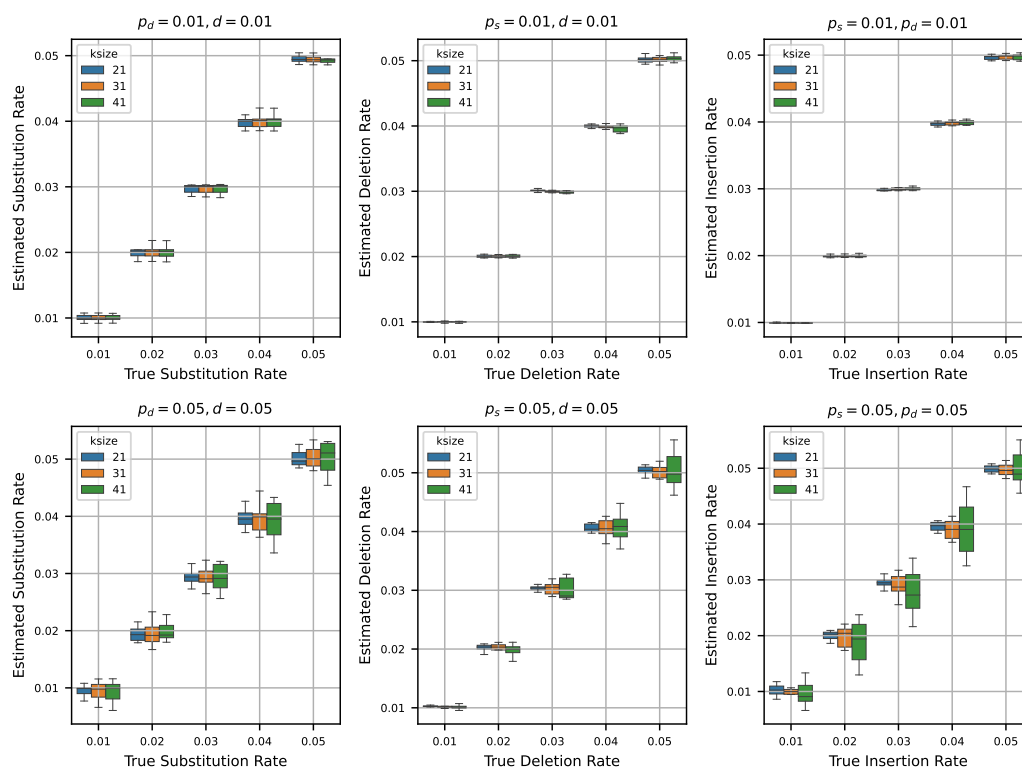
After testing our estimators for varying k -mer lengths, sequence lengths, and base compositions, we next turn to estimating mutation rates by varying the true rates across a range of values. To do this, we generated a synthetic reference genome of 1 million base pairs, fixing the base composition at 30% “A”, and equal proportions of “C”, “G”, and “T” – making sure total frequency is 100%. Using the mutation model described in Section 2, we then simulated mutated genomes from the synthetic reference by varying the mutation rates p_s , p_d , and d across the values $\{0.01, 0.02, 0.03, 0.04, 0.05\}$. For every parameter combination, we generated 10 independent mutated genomes to capture stochastic variability. We then estimated the mutation rates using the estimators detailed in Section 4 for each of these mutated genomes.

In Figure 4, we show two sets of results:

- **Fixed low rates (0.01):** Fixing two of the rates at 0.01, we show the estimates of the third rate as the true rate varies from 0.01 to 0.05. We repeat this process independently for p_s , p_d , and d .
- **Fixed high rates (0.05):** Fixing two of the rates at a higher value of 0.05, we also show the estimates of the third rate as the true rate varies from 0.01 to 0.05. Again, we repeat this for all three rates.

When we set the other two rates to 0.01 and estimate the third rate, we observe that the estimated rates are highly accurate across all trials. In many cases, the boxplots of the estimates nearly vanish due to minimal variance, indicating tight clustering around the true values. This trend remains consistent across multiple k -mer sizes, suggesting that the estimators are robust at low rates of mutation.

In contrast, when we fix the other two rates at 0.05 and estimate the third rate, the accuracy of the estimation decreases slightly. While the estimates still track the true values reasonably well, the variance increases, and the boxplots become more prominent. Notably, the median estimate remains close to the true rate in most settings, which indicates that the estimators retain their central tendency even under higher mutation rates. However, for



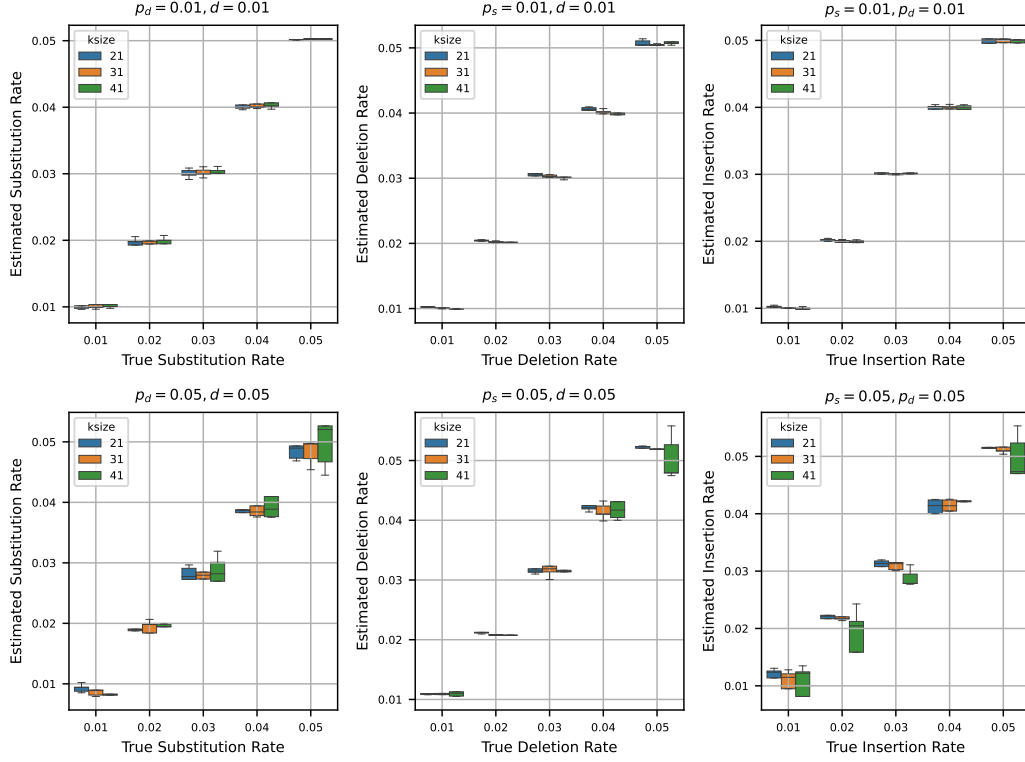
■ **Figure 4** Estimated mutation rates versus true values, where the original string is a synthetic sequence. Each subplot corresponds to a case where two mutation rates are fixed (either at 0.01 or 0.05) and the third is varied from 0.01 to 0.05. Each boxplot shows the variability in estimation across 10 simulations. The results show that the rate estimation is very accurate when the other two rates are small, and is reasonably accurate when the other two rates are larger.

larger k -mer sizes, we observe increased variability in the estimates – an effect that mirrors our earlier observations in Section 7.1, where longer k -mers resulted in decreased precision of the estimators.

7.5 Estimating rates from real sequences

Having the estimators tested for a synthetic reference, we next estimate rates from a real genome sequence. For this set of experiments, we used the reference assembly of *Staphylococcus aureus* (subspecies: aureus USA300_TCH1516), which has 2.8 million nucleotides. We simulated mutated sequences from this reference by running the mutation process described in Section 2 by varying the mutation rates p_s , p_d , and d from 0.01 to 0.05. Similar to Section 7.4, we generated 10 independent mutated sequences for each combination of p_s , p_d , and d to capture stochastic variability. We then estimated the mutation rates using the estimators outlined in Section 4 for each of these simulated sequences.

Figure 5 shows the estimated mutation rates plotted against the true rates that were used to run the mutation process. We observe that the results shown in Figure 5 are consistent with previously discussed results. Specifically, when estimating a given mutation rate while keeping the other two rates low (0.01), the corresponding estimator performs with high precision, closely tracking the true value. On the other hand, when the estimation is carried



■ **Figure 5** Estimated mutation rates versus true values, where the original string is the reference genome of *Staphylococcus* (length 2.8 million). Each subplot corresponds to a case where two mutation rates are fixed (either at 0.01 or 0.05) and the third varies from 0.01 to 0.05. Each boxplot shows the variability in estimation across 10 simulations. Estimated mutation rates closely match true values when other rates are low, but estimation becomes less precise under high total mutation rates ($>10\%$) due to increased difficulty in real genomic contexts.

out with the other two mutation rates set to higher values (0.05), the estimates appear more confounded. This is likely due to the increased difficulty of accurately estimating the number of k -spans with a single deletion (\mathcal{D}) or no mutation (\mathcal{N}) in a real genomic context. Notably, these experiments involve challenging conditions, with total mutation rates exceeding 10%. Despite this, the estimators yield reasonably accurate results, indicating potential for practical effectiveness.

7.6 Comparison with substitution rates estimated using simple mutation model

We conclude the experiments section by contrasting the substitution rates estimated using (10) with substitution rates estimated considering a simple mutation model. We use the statistics of k -mers developed in a recent work [1] to estimate substitution rates under a simple mutation model. The simple mutation model captures only substitutions, and no insertions or deletions. Consequently, we can only compute substitution rates considering this simple model. Henceforth, we refer to the simple mutation model as SMM.

We estimated the substitution rates using the SMM for the same simulated mutated sequences described in Section 7.5. In Figure 6, we show the estimated substitution rates using our estimators in (10), and using SMM. The results highlight that the substitution

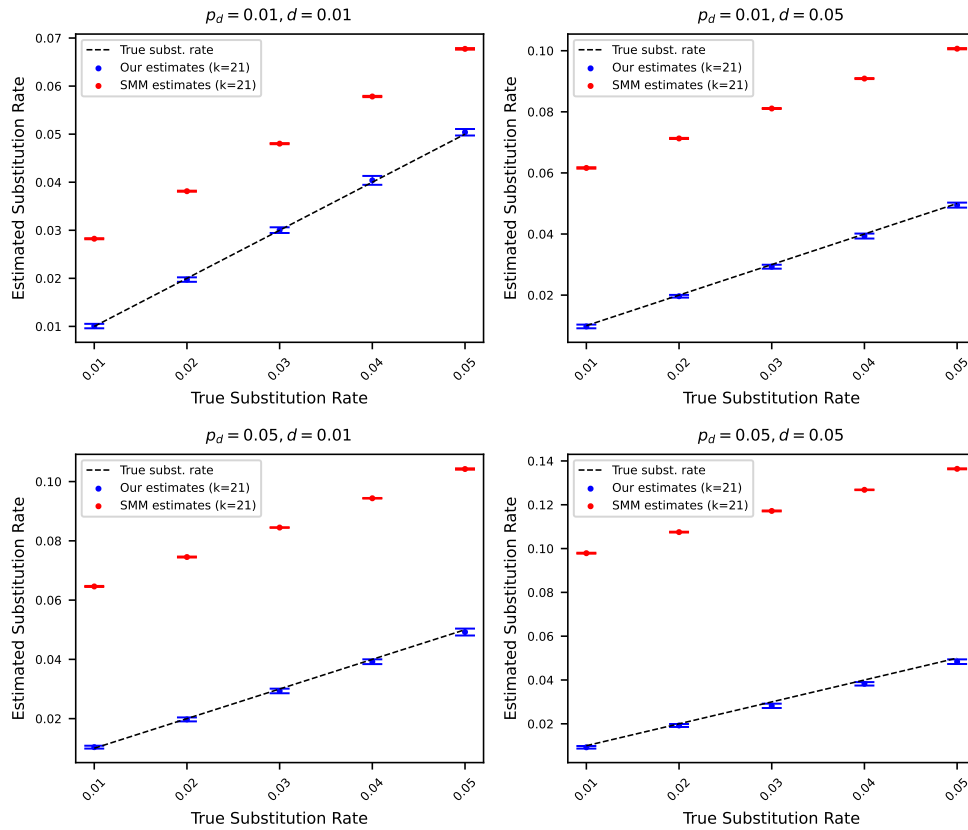


Figure 6 Estimated substitution rates versus true substitution rates, where the original string is the reference genome of *Staphylococcus* (length 2.8 million). Rates were estimated using (10) and using a simple mutation model (SMM) that only considers substitution. Each subplot corresponds to a case where p_d and d are fixed at either 0.01 or 0.05. The points show the average of 10 estimates, and the error bars show one standard deviation. The dashed black line corresponds to the true substitution rates. Estimated substitution rates using our method closely match true rates, whereas SMM overestimates due attributing to substitutions mismatching k -mers originating from insertions and deletions.

rates estimated using the estimator we developed track the true substitution rates accurately. On the other hand, the substitution rates estimated using SMM make a gross overestimation. This is because the SMM does not consider indels, and therefore, the effects of all three mutation rates are subsumed in the single substitution rate we get using the SMM. As such, a simple mutation model cannot disentangle the distinct contributions of substitution, insertion, and deletion rates. In contrast, the mutation model we introduce effectively decomposes these components, enabling more accurate and meaningful estimation of individual mutation rates.

8 Conclusions

We have presented a mutation model that accommodates single-nucleotide substitutions, as well as insertions and deletions while retaining enough mathematical structure to admit closed-form rate estimators derived solely from k -mer statistics. From this model, we

obtained algebraic estimators for the three elementary mutation rates: p_s , p_d , and d ; we also proved a relatively tight sub-exponential concentration bound on p_s that guarantees consistency as sequence length grows. We also identified regimes in which the estimation becomes ill-conditioned (i.e. large k , $p_d = 0$, or sequence composition with 25% “A”). These results establish a bridge between sequence evolution and combinatorial word statistics, thus providing additional tools for theoretical algorithmic computational biology.

In our prototype implementation, we demonstrated that our estimates remain accurate on simulated evolution of real genomes, and outperforms a substitution-only simple mutation model by avoiding spurious attribution of indel signals. While naive counting of unmutated and single-deletion k -mers sufficed to show practical accuracy of our estimators, this raises an interesting open problem: estimating the number of these unmutated and single-deletion k -mers efficiently for large scale data sets.

Several directions invite further investigation. First, incorporating the count \mathcal{S} of single-substitution k -spans may illuminate why p_s remains relatively stable even for moderately large k . Second, our framework can extend to heterogeneous or context-dependent rates by replacing global expectations with position-specific covariates. Third, coupling our estimators with sketch-based distance measures (such as in [8]) may provide a theory-backed avenue for larger scale applications such as phylogenetic placement in the presence of high indel activity. Finally, a more thorough investigation on real genomic data (where the unitig-based approach we used in the practical implementation starts to become infeasible) will be necessary to understand the utility of the mutation estimates in practice.

In summary, by utilizing probabilistic modeling and concentration inequalities, we provide a theoretical foundation and initial practical implementation for quantifying the parameters of a relatively complex mutation process directly from k -mers. We anticipate that these ideas will continue to inform new alignment-free computational biology tasks, particularly relevant as sequencing data continues to outpace traditional alignment-based paradigms.

References

- 1 Antonio Blanca, Robert S Harris, David Koslicki, and Paul Medvedev. The statistics of k -mers from a sequence undergoing a simple mutation process without spurious matches. *Journal of Computational Biology*, 29(2):155–168, 2022. doi:10.1089/CMB.2021.0431.
- 2 Ke Chen, Vinamratha Pattar, and Mingfu Shao. Sequence similarity estimation by random subsequence sketching. *bioRxiv*, pages 2025–02, 2025.
- 3 Margaret Oakley Dayhoff. Computer analysis of protein evolution. *Scientific American*, 221(1):86–95, 1969.
- 4 Richard Durbin, editor. *Biological sequence analysis*. Cambridge Univ. Press, Cambridge [u.a.], 17. print. edition, 2013.
- 5 Arun Ganesh and Aaron Sy. Near-Linear Time Edit Distance for Indel Channels. In Carl Kingsford and Nadia Pisanti, editors, *20th International Workshop on Algorithms in Bioinformatics (WABI 2020)*, volume 172 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 17:1–17:18. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2020. doi:10.4230/LIPIcs.WABI.2020.17.
- 6 Steven Henikoff and Jorja G Henikoff. Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences*, 89(22):10915–10919, 1992.
- 7 Mahmudur Rahman Hera, Shaopeng Liu, Wei Wei, Judith S Rodriguez, Chunyu Ma, and David Koslicki. Metagenomic functional profiling: to sketch or not to sketch? *Bioinformatics*, 40(Supplement_2):ii165–ii173, 2024. doi:10.1093/BIOINFORMATICS/BTAE397.

- 8 Mahmudur Rahman Hera, N Tessa Pierce-Ward, and David Koslicki. Deriving confidence intervals for mutation rates across a wide range of evolutionary distances using FracMinHash. *Genome Research*, 33(7):1061–1068, 2023.
- 9 Jamshed Khan and Rob Patro. Cuttlefish: fast, parallel and low-memory compaction of de bruijn graphs from large-scale genome collections. *Bioinformatics*, 37(Supplement_1):i177–i186, 2021.
- 10 Brian D Ondov, Gabriel J Starrett, Anna Sappington, Aleksandra Kostic, Sergey Koren, Christopher B Buck, and Adam M Phillippy. Mash Screen: high-throughput sequence containment estimation for genome discovery. *Genome Biology*, 20:1–13, 2019.
- 11 Mahmudur Rahman Hera and David Koslicki. Estimating similarity and distance using fracminhash. *Algorithms for Molecular Biology*, 20(1):1–13, 2025.
- 12 Sophie Röhling, Alexander Linne, Jendrik Schellhorn, Morteza Hosseini, Thomas Dencker, and Burkhard Morgenstern. The number of k-mer matches between two dna sequences as a function of k and applications to estimate phylogenetic distances. *Plos one*, 15(2):e0228070, 2020.
- 13 Shahab Sarmashghi, Kristine Bohmann, M Thomas P Gilbert, Vineet Bafna, and Siavash Mirarab. Skmer: assembly-free and alignment-free sample identification using genome skims. *Genome Biology*, 20(1):1–20, 2019.
- 14 Tizian Schulz and Paul Medvedev. ESKEMAP: exact sketch-based read mapping. *Algorithms for Molecular Biology*, 19(1):19, 2024. doi:10.1186/S13015-024-00261-7.
- 15 Jim Shaw and Yun William Yu. Fast and robust metagenomic sequence comparison through sparse chaining with skani. *Nature Methods*, 20(11):1661–1665, 2023.
- 16 Jim Shaw and Yun William Yu. Rapid species-level metagenome profiling and containment estimation with sylph. *Nature Biotechnology*, pages 1–12, 2024.
- 17 Kai Song, Jie Ren, Gesine Reinert, Minghua Deng, Michael S Waterman, and Fengzhu Sun. New developments of alignment-free sequence comparison: measures, statistics and next-generation sequencing. *Briefings in bioinformatics*, 15(3):343–353, 2014. doi:10.1093/BIB/BBT067.
- 18 Martin Šošić and Mile Šikić. Edlib: a c/c++ library for fast, exact sequence alignment using edit distance. *Bioinformatics*, 33(9):1394–1395, 2017. doi:10.1093/BIOINFORMATICS/BTW753.
- 19 Shiheng Tao, Yanhui Fan, Wenjuan Wang, Guoji Ma, Lijing Liang, and Qi Shi. Patterns of insertion and deletion in mammalian genomes. *Current Genomics*, 8(6):370–378, 2007.
- 20 Kun Wu, Danqi Qin, Yang Qian, and Haoxuan Liu. A new era of mutation rate analyses: Concepts and methods. *Zoological Research*, 45(4):767, 2024.
- 21 Andrzej Zielezinski, Susana Vinga, Jonas Almeida, and Wojciech M Karlowski. Alignment-free sequence comparison: benefits, applications, and tools. *Genome biology*, 18:1–17, 2017.