

# Supplementary Information - Entropy predicts sensitivity of pseudo-random seeds.

Benjamin D. Maier<sup>1</sup>, Kristoffer Sahlin<sup>1</sup>.

<sup>1</sup> Department of Mathematics, Science for Life Laboratory, Stockholm University, 106 91, Stockholm, Sweden.

## Contents

<b>S.1 Construction of mixedstrokes and altstrokes</b>	<b>SI-2</b>
S.1.1 Construction of mixedstrokes . . . . .	SI-2
S.1.2 Construction of altstrokes . . . . .	SI-2
S.1.2.1 Minimum $k_s$ for altstrokes . . . . .	SI-5
S.1.2.2 Selecting altstrokes parametrization for experiments . . . . .	SI-5
S.1.2.3 Altstrokes implementation details . . . . .	SI-5
S.1.3 Construction of multistrokes . . . . .	SI-5
<b>S.2 Example computations of <math>P(X_i Y_j)</math></b>	<b>SI-8</b>
S.2.1 $k$ -mers . . . . .	SI-8
S.2.2 Randstrokes . . . . .	SI-8
S.2.3 Altstrokes . . . . .	SI-8
S.2.4 Mixedstrokes . . . . .	SI-9
S.2.5 Multistrokes . . . . .	SI-9
<b>S.3 Empirically estimating <math>P(N_m &gt; 0)</math></b>	<b>SI-9</b>
<b>S.4 Experiments</b>	<b>SI-9</b>
S.4.1 Evaluation metrics . . . . .	SI-9
S.4.2 Data simulation . . . . .	SI-10
S.4.3 Simulated results . . . . .	SI-11
S.4.4 <i>E. coli</i> Oxford Nanopore Technology Reads . . . . .	SI-11
S.4.5 Biological Data . . . . .	SI-12
<b>S.5 Minimap2 implementation details</b>	<b>SI-12</b>
S.5.1 Minimap2 analysis . . . . .	SI-12
<b>S.6 Figures and tables</b>	<b>SI-12</b>

## 1 S.1 Construction of mixedstrokes and altstrokes

### 2 S.1.1 Construction of mixedstrokes

3 Mixedstrokes consists out of a specified fraction of  $k$ -mers and strobemers and may be  
 4 sampled with either of the three strobemers seeding methods (minstrokes, hybridstrokes  
 5 and randstrokes), but we will only consider randstrokes here.

6 Whether a strobemer or a  $k$ -mer is seeded depends on the hash value of the first  
 7 strobe  $h(S[i : i + \ell])$  and the user-defined strobe fraction  $q$ . For instance, for mixed-  
 8 strokes with  $q = 0.8 = 80\%$ , strobemers should be generated when  $h(S[i : i + \ell]) \%$   
 9  $100 < 80$ ; otherwise a  $k$ -mer should be sampled instead. Strobemers are sampled fol-  
 10 lowing the routine [SR2]. When sampling a  $k$ -mer  $(1 - q)$ ,  $n\ell$  consecutive nucleotides  
 11 are taken starting from the start position of the first strobe ( $S[i : i + n\ell]$ ) and converted  
 12 to its respective hash value (Python). By taking  $n\ell$  nucleotides, we obtain the same  
 13 subsequence lengths as strobemers consisting out of  $n$  strobes of length  $\ell$ . In C++,  $k$ -  
 14 mers are constructed by summing the hash values of  $S[i + j * \ell]$ , which are each divided  
 15 by  $(n + j)$  to guarantee non-commutativity, whereby  $j \in [1, n]$ .

### 17 S.1.2 Construction of altstrokes

18 Altstrokes are modified randstrokes where the strobe length is alternating between  
 19 shorter ( $k_s$ ) and longer strobes ( $k_l$ ) with  $|k_s| + |k_l| = k$ . Whether the first strobe  
 20 is of length  $|k_s|$  or  $|k_l|$  is decided based on the hash value of the substring of length  
 21  $|k_s|$  (*i.e. the potential first strobe*). Note that it is highly advised against making the  
 22 decision based on the hash value of  $k_l$  as this may lead to unnecessarily many seeds  
 23 being destroyed where  $(k_s, k_l)$ -altstrokes should have been sampled and there is are  
 24 mutations within the positions  $[k_s, k_l]$  downstream from the start position of the seed.  
 25 Furthermore, one should not seed altstrokes where the length of the shorter strobe is  
 26 smaller than 5 to avoid seed repetitiveness (low uniqueness).

27 In case, that we are sampling mixed-altstrokes, the hash value is divided by 100,  
 28 whereby the fractional part (remainder) is discarded. This integer division is required  
 29 as we have to make two independent decisions using the hash value: whether to sample  
 30 an altstrobe or a  $k$ -mer and in case of an altstrokes whether to sample a short or a long  
 31 strobe first.

32 The following downstream strobes are selected from a window  $W$  by alternatively  
 33 sampling a short and a long strobe using the randstrobe linking routine as described  
 34 in [SR2]. We adjust the offset of  $(w_{min}, w_{max})$  depending on if it is the long or  
 35 short strobe we sample. Specifically, we let  $k_l$  in altstrobe  $(k_s, k_l)$  be sampled from  
 36  $[w_{min} - (k_l - k_s)/2, w_{max} - (k_l - k_s)/2]$  and  $k_s$  in altstrobe  $(k_l, k_s)$  be sampled from  
 37  $[w_{min} + (k_l - k_s)/2, w_{max} + (k_l - k_s)/2]$ . This guarantees that the maximum length of  
 38 the altstrobe seed remains the same as randstrokes ( $w = w_{max} + k/2$ ) which is important  
 39 for benchmarking.

40 In the default altstrobe protocol where  $2k_s = k_l$ , we store  $S[i]$  for short strobes  
 41 and  $S[i]$  as well as  $S[i] + k_s$  for long strobes. This allows us to keep track of whether  
 42  $(|k_s|, |k_l|)$  or  $(|k_l|, |k_s|)$  was sampled and facilitates analysis (e.g. sequence coverage) as  
 43 now all strobes are of equal length ( $k_s$ ). If generalized altstrokes are sampled where  
 44  $2k_s \neq k_l$  and information about the exact strobe composition is required, we store all  
 45 strobe positions.

---

**Algorithm 1** Mixedstrokes

---

**Function:** Mixedstrokes ( $S, n, \ell, w_{min}, w_{max}, s$ )**Input:** Sequence  $S$ , number of strokes  $n$ , stroke lengths  $\ell$ , stroke window  $W[w_{min}, w_{max}]$ , stroke fraction  $s$ **Output:** Mixedstrokes of order  $n$  and their positions from  $S$  $s = \frac{N}{D}$  // Split stroke fraction  $s$  in numerator  $N$  and denominator  $D$  $O = []$  // Initialize array of strobemers and their positions**for**  $i \in [1, |S| - (n * \ell + 1)]$  **do** // Iterate over all positions     $m_1 = S[i : i + \ell]$      $P = [i, ]$      $H = h(m_1)$     **if**  $h(m_1) \% D < N$  **then** // Sample strobemer at position  $i$          $w_u = \min(w_{max}, ([S] - i) / (n - 1))$  // Second argument only active at end of  $S$          $w_l = \max(w_{min} - (w_{max} - w_u), \ell)$         **for**  $j \in [2, n]$  **do**             $w' = [i + w_l + (j - 2)w_u, i + (j - 1)w_u]$  // Window to look for current stroke             $p = \arg \min_p \{p : h(m \oplus S[p : p + \ell]) \leq h(m \oplus S[i' : i' + \ell]), \forall i' \in w'\}$  // Selecting strokes based on strobemer protocol, here exemplified using randstrokes)             $P += p$              $H += (-1)^j * j * h(S[p : p + \ell])$  // Non-commutative hash value combination        **end**    **else** // Sample  $k$ -mer of length  $n\ell$  at position  $i$         **for**  $j \in [2, n]$  **do**             $P += i + j * \ell$              $H = h(S[i : i + n\ell])$         **end**    **end**     $O += (P, H)$     **return**  $O$ **end**

---

---

**Algorithm 2** Altstrobes

---

**Function:** Altstrobes ( $S, n, k_s, k_l, w_{min}, w_{max}$ )**Input:** Sequence  $S$ , number of strobes  $n$ , length of short and long strobes  $k_s$  and  $k_l$ , strobe window  $W[w_{min}, w_{max}]$ **Output:** Altstrobes of order  $n$  and their positions from  $S$ **Require :**  $n \in 2\mathbb{Z}$  $O = []$  // Initialize array of altstrobes and their positions**for**  $i \in [1, |S| - (n * \ell + 1)]$  **do** // Iterate over all positions     $P = [i, ]$     **if**  $h(S[i : i + k_s]) \% 2 = 0$  **then** // Sample altstrobe( $k_s, k_l$ ) at position  $i$          $H = h(S[i : i + k_s])$          $strobes = [k_s, k_l]$     **else** // Sample altstrobe( $k_l, k_s$ ) at position  $i$          $H = h(S[i : i + k_l])$          $strobes = [k_l, k_s]$     **end**     $w_u = \min(w_{max}, ([S] - i) / (n - 1))$  // Second argument only active at end of  $S$      $w_l = \max(w_{min} - (w_{max} - w_u), k_l)$     **for**  $j \in [2, n]$  **do**         $k' = strobes[(j + 1) \% 2]$  // Retrieve length of next strobe         $offset = (k_s + k_l) / 2 - k'$  // Adjusting window offsets         $w' = [i + w_l + (j - 2)w_u + offset, i + (j - 1)w_u + offset]$  // Window to look for current strobe         $p = \arg \min_p \{p : h(m \oplus S[p : p + k']) \leq h(m \oplus S[i' : i' + k']), \forall i' \in w'\}$   
        // Selecting strobes based on randstrobe protocol         $P += p$          $H += (-1)^j * j * h(S[p : p + k'])$  // Non-commutative hash value combination    **end**     $O += (P, H)$     **return**  $O$ **end**

---

### 46 S.1.2.1 Minimum $k_s$ for altstrobes

47 The combination of strobe lengths  $k_s$  and  $k_l$  matter in practice. A too short  $k_s$  leads  
 48 to degenerate seed constructs for two reasons. Firstly, the number of possible hash  
 49 values is limited for very short strobes (e.g., 4, 16, and 64 hash values for 1-, 2-, and  
 50 3-mers respectively). Thus, it cannot be guaranteed that an even fraction of  $(k_s, k_l)$   
 51 and  $(k_l, k_s)$  is computed for such parametrizations, which decreases randomness in  
 52 seed selection. Secondly,  $|k_s|$  needs to be long enough to avoid random repetitions in  
 53 substrings of length  $w_{max} - w_{min}$ . To see why this is the case, consider using a short  
 54 strobe size of 2. If we have a sequence ACGACTACA..., where AC hash to an even  
 55 number (short strobe selected), we will very likely link AC (short strobe) with the same  
 56  $k_l$  multiple times as the downstream selection window of size  $w_{max} - w_{min}$  will share  
 57 many strobes for reasonably large windows. This deteriorates both randomness and  
 58 uniqueness of seeds. We show uniqueness per  $k_s$  length in Suppl. Fig. S4, where for  
 59 random simulated sequences has a window size of 25nt. For this window size,  $|k_s| \geq 6$   
 60 is needed to be competitive with other parametrizations in terms of uniqueness.

### 61 S.1.2.2 Selecting altstrobes parametrization for experiments

62 Several combinations of  $(k_s, k_l)$  in altstrobes could lead to lower correlation and poten-  
 63 tially outperform randstrobes, as was shown in Figure S5. We performed simulations  
 64 on random sequences for all possible altstrobe combinations of combined strobe length  
 65 of 30. We found that sequence coverage, match coverage and expected island size was  
 66 best for altstrobes with  $k_s$  between 7 and 10 (Fig. S5). Since simulated sequences are  
 67 less repetitive than biological sequences, we opt for the longest possible  $k_s$  with good  
 68 metrics.

### 69 S.1.2.3 Altstrobes implementation details

70 Despite altstrobes orders being multiples of 2, we output the results as multiples of  
 71 order 3 to keep track on whether a short-long or a long-short combination was sampled.  
 72 As the long strobe  $k_l$  is exactly double the size of  $k_s$  we can store altstrobe seeds of  
 73 order 2 as  $(x_1, x_1 + k_s, x_2)$  and  $(x_1, x_2, x_2 + k_s)$  with  $x_1$  and  $x_2$  being the start positions  
 74 of the strobes and  $k_s$  being the shorter strobe length.

## 75 S.1.3 Construction of multistrobes

76 Multistrobes are generalized altstrobes which allows to sample a range of possible strobe  
 77 lengths ranging from  $k_s$  to  $k_l$  with  $|k_s| + |k_l| = k$  within the same seeding pass.

78 As the strobe length combination is decided based on the hash value of a subsequence  
 79 of length  $k_s$  ( $h(S[i : i + |k_s|])$ ), we have to deal with uniformity issues. For instance,  
 80 for  $|k_s| = 2$  and  $|k_l| = 28$ , there are only  $4^2 = 16$  hash values  $h(S[i : i + |k_s|])$  making  
 81 it impossible to select all 27 options from  $|k_l - k_s + 1|$ . However, there are less strobe  
 82 length combinations from  $[k_s, k/2]$  than the 16 hash values. Therefore, if we first pick  
 83 a strobe size from  $[k_s, k/2]$  (14 options), we may in a second step choose whether it  
 84 should be the first or the second strobe. In more detail, the approach is performed as  
 85 follows.

86 First, a strobe length  $k'$  is sampled from  $[k_s, k/2]$  based on the hash value of  $h(S[i :$   
 87  $i + |k_s|])$  as

$$k' \doteq |k_s| + (h(S[i : i + |k_s|]) \% (|k_l| - |k|/2 + 1)).$$

88 Second, we decide if the strobe of length  $k'$  should come first or second based on

$$f(i, k, w, S, *) = \begin{cases} \text{Sample multistrobe } (k', k - k') & \text{if } ((h(S[i : i + |k'|]) // 100) \% 2 = 0 \\ \text{Sample multistrobe } (k - k', k'), & \text{otherwise} \end{cases}$$

89 This approach can, at best, double the possible hash value selection options and also  
90 improve the uniformity of the number of short-long and long-short combinations when  
91  $k_s$  is small. Note however that our two-step implementation described above is not fully  
92 uniform when  $|k_l - k_s + 1|$  is uneven, as the option  $(k/2, k/2)$  is appearing with double  
93 probability as their short-long and long-short combinations are identical. The python  
94 implementation implements this two-step procedure to cope with very small  $k_s$  for our  
95 sensitivity benchmarks while the C++ implements the sampling as described in the  
96 main paper for runtime. While the number of possible hash values is largely exceeding  
97 the number of strobe length options for  $k_s > 3$  ensuring a near-uniform distribution  
98 of strobe sizes, we observe a high variability of the simulation results for very small  
99  $k_s$  at each new experiment, since Python's hash function starts with a different seed  
100 value. This has an effect on the distribution of strobe sizes. Especially for  $k_s = 1$ ,  
101 it may happen that a poor hash value distribution cause a non-uniform strobe sizes  
102 distorting simulation results. Hence, we decided to avoid degenerate hash functions for  
103 the smallest  $k_s \in [1, 2]$ . For  $k_s \leq 2$  we use a customized hash function based on the first  
104 two nucleotides at the start position of the first strobe, using the following strobe size  
105 assignments for our sensitivity estimation in Fig. 2 and Suppl. Figs. S1 - S3. When  
106  $k_s = 1$ , we use the mapping  
107 ("AA": 1, "AC": 9, "AG": 19, "AT": 27, "CA": 14, "CC": 3, "CG": 11, "CT": 21,  
108 "GA": 23, "GC": 16, "GG": 5, "GT": 13, "TA": 29, "TC": 25, "TG": 17, "TT": 7).  
109 When  $k_s = 2$ , we use the mapping  
110 ("AA": 2, "AC": 9, "AG": 19, "AT": 26, "CA": 14, "CC": 4, "CG": 11, "CT": 21,  
111 "GA": 23, "GC": 16, "GG": 6, "GT": 13, "TA": 28, "TC": 24, "TG": 17, "TT": 7).

---

**Algorithm 3** Multistrobes

---

**Function:** Multistrobes ( $S, n, k_s, k_l, w_{min}, w_{max}$ )**Input:** Sequence  $S$ , number of strobes  $n$ , minimum strobe length  $k_s$ , maximum strobe length  $k_l$ , strobe window  $W[w_{min}, w_{max}]$ **Output:** Multistrobes of order  $n$  and their positions from  $S$ **Require :**  $n \in 2\mathbb{Z}$  $O = []$  // Initialize array of multistrobes and their positions**for**  $i \in [1, |S| - (n * \ell + 1)]$  **do** // Iterate over all positions     $P = [i, ]$      $|k'| \doteq |k_s| + (h(S[i : i + |k_s|]) \% (|k_l| - |k|/2 + 1))$     **if**  $(h(S[i : i + |k'|]) // 100) \% 2 = 0$  **then** // Decide whether to start with a short strobe         $k_1 = k'$          $k_2 = k - k'$     **else**         $k_1 = k - k'$          $k_2 = k'$     **end**     $strobes = [k_1, k_2]$      $H = h(S[i : i + k_1])$      $w_u = \min(w_{max}, (|S| - i) / (n - 1))$  // Second argument only active at end of  $S$      $w_l = \max(w_{min} - (w_{max} - w_u), k_2)$     **for**  $j \in [2, n]$  **do**         $k' = strobes[(j + 1) \% 2]$  // Retrieve length of next strobe         $offset = k/2 - k'$  // Adjusting window offsets         $w' = [i + w_l + (j - 2)w_u + offset, i + (j - 1)w_u] + offset$  // Window to look for current strobe         $p = \arg \min_p \{p : h(m \oplus S[p : p + k']) \leq h(m \oplus S[i' : i' + k']), \forall i' \in w'\}$   
        // Selecting strobes based on randstrobe protocol         $P += p$          $H += (-1)^j * j * h(S[p : p + k'])$  // Non-commutative hash value combination    **end**     $O += (P, H)$     **return**  $O$ **end**

---

## 112 S.2 Example computations of $P(X_i|Y_j)$

113 For convenience, we will denote  $P(X_i|Y_j)$  for  $k$ -mers, randstrokes, mixedstrokes, alt-  
 114 strokes, and multistrokes as  $P_k$ ,  $P_r$ ,  $P_{mi}$ ,  $P_a$ ,  $P_{mu}$ , respectively.

### 115 S.2.1 $k$ -mers

116 We have

$$P_k = \begin{cases} 1 & \text{if } i + j + 2 \leq k \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

117 We will refer to this equation as  $E_6(k)$  taking the length of the  $k$ -mer as argument,  
 118 as it turns out handy to reuse for some of the other constructs.

### 119 S.2.2 Randstrokes

120 Let the stroke size be  $k' = \lfloor k/2 \rfloor$ . If  $Y_j$  is located on the second stroke, the probability  
 121 is  $E_6(k')$  since  $i$  can only be covered by the second stroke. If  $j$  is located on the first  
 122 stroke,  $i$  can be either covered by the first or second stroke. Hence, we need to structure  
 123 the probability up into cases. The possible second stroke coverings of position  $i$  given  $Y_j$   
 124 depends on the downstream window location, and is further restricted by the window  
 125 size ( $B = w_{max} - w_{min} + 1$ ) and the stroke length. The second stroke coverings is  
 126 computed as

$$A = \min(k', B, i + j + 2 - w_{min}, w - (i + j + 2)).$$

127 Under the assumption that  $w_{min} > k'$ , we have

$$E_7(k') = \begin{cases} 1 & \text{if } i + j + 2 \leq k' \\ A/B & \text{if } w_{min} - j \leq i \leq w - (j + 1) \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

128 Let  $I_j$  be an indicator variable with  $I_j = 1$  if  $j$  is placed on first stroke and 0  
 129 otherwise. Then,

$$P_r = E_7(k')I_j + E_6(k')(1 - I_j)$$

### 130 S.2.3 Altstrokes

For altstrokes, we have the same cases as for randstrokes ( $E_6$  and  $E_7$ ), but we need  
 to sum over the scenario that either the first or the second stroke is short (with equal  
 probabilities 0.5). We have

$$P_a = (0.5E_7(k_s) + 0.5E_7(k_l))I_j + (0.5E_6(k_l) + 0.5E_6(k_s))(1 - I_j).$$

131 To guarantee the same maximum seed length we employed a window offset based on  
 132 if the first stroke as short or long (as described in section 2.8), this window adjustment  
 133 is computed before  $E_7$  is computed.



## 134 S.2.4 Mixedstrokes

For mixedstrokes, we need to condition on the probability  $q$  that either a randstroke or a  $k$ -mer is sampled. We have

$$P_{mi} = (qE_7(k') + (1 - q)E_6(k))I_j + qE_6(k')(1 - I_j).$$

135 Where we have included in the event  $I_j = 1$  that a  $k$ -mer is sampled, as it can be  
136 seen as a stroke of length  $k$  with the second stroke being of length 0.

## 137 S.2.5 Multistrokes

138 Finally, multistrokes samples a stroke length  $\ell \in [k_b, k - k_b]$  and sets  $k - \ell$  as length  
139 for the second stroke. Let  $p_\ell = \frac{1}{k - 2k_b + 1}$  denote this probability (which is uniform over  
140 the considered stroke lengths). Let the shorter of the two be denoted by  $k_s$  and the  
141 longer be denoted by  $k_l$ , similarly to altstrokes. Then the probability is similar to  $P_a$   
142 but summed over possible stroke lengths. We have

$$P_{mu} = \sum_{\ell=k_b}^{k-k_b} p_\ell \left( (0.5E_7(k_s) + 0.5E_7(k_l))I_j + (0.5E_6(k_l) + 0.5E_6(k_s))(1 - I_j) \right).$$

143 In case of  $P_r$ ,  $P_{mi}$ ,  $P_a$ , and  $P_{mu}$ , the above probabilities assume that the second  
144 stroke is chosen uniformly at random over the window, meaning that they assume a  
145 perfect random hash function. When the first stroke in altstrokes or multistrokes is  
146 very short it is not possible to sample uniformly from the window, as was discussed and  
147 shown in the construction of altstrokes section.

## 148 S.3 Empirically estimating $P(N_m > 0)$

149 We empirically estimated  $P(N_m > 0)$  for each of the seed construct parametrizations  
150 as follows. A string  $S$  (letters A, C, G, and T) of length  $4w$  is simulated at random,  
151 and a second string  $T$  is simulated by copying  $S$  and randomly performing exactly  $m$   
152 mutations within the first  $2w$  nucleotides of  $S$ , creating a fixed mutation rate within the  
153 segment of  $2w$  on  $S$ . As for the mutation profile, we uniformly draw the substitution  
154 rate  $s \in [0, 1]$  and set insertions and deletions each to  $(1 - s)/2$ . We then construct seeds  
155 from  $S$  and  $T$  with each seed construct and parametrization. We store the event that  
156 seeds from the first  $w$  seeds on  $S$  and  $T$  has at least one match. We estimate  $P(N_m > 0)$   
157 from the fraction of experiments with at least one match out of  $U$  experiment replicates,  
158 where  $U = 10,000$  for  $w_{max} = 50$  and  $U = 1,000$  for  $w_{max} = 100$  and 200. Finally,  
159 we obtain the summed sensitivity  $\sum_{m=1}^M P(N_m > 0)$  to capture the sensitivity over a  
160 range of error rates.  $M$  is chosen such that it corresponds to 30% error rate within the  
161 segment of length  $2w$ . We chose 30% as all seed constructs returned no matches with  
162 the given seed lengths for this error rate. We chose  $w_{max}$  of 50, 100, and 200, giving  
163 rise to  $w$  of 64, 114, and 214 respectively. We chose these values to study the effect of  
164 the window size, where 50 and 100 are also consistent with previous study [SR2].

## 165 S.4 Experiments

### 166 S.4.1 Evaluation metrics

167 We also evaluated our best performing seed constructs from the sensitivity analysis in  
168 a more general sequence matching scenario using previously designed metrics in [SR2],

169 namely fraction of matches, match coverage, sequence coverage, and expected island  
170 size:

- 171 • **Fraction of Matches:** proportion of the query seeds that matched the reference
- 172 • **Match Coverage:** proportion of nucleotides covered by the  $k$ -mers and strobem-  
173 mers from end-to-end including potential gaps
- 174 • **Sequence Coverage:** proportion of nucleotides covered by the strobemers of matches,  
175 so it distinguished from match coverage by disregarding the gaps between the  
176 strobemers
- 177 • **Expected Island Size:** An island is the maximal interval of consecutive nu-  
178 cleotides without matches. If a random location from the reference genome is  
179 selected, they may either be covered by matches (size of island = 0) or islands  
180 of various length. For a sequence  $S$  and a set of islands length  $X$ , the expected  
181 island size  $E$  is computed as follows:

$$E = \frac{1}{|S|} \sum_{x \in X} x^2 \quad (3)$$

182 .

183 We used the same parameters and simulation setup as in [SR2] (details in Suppl.  
184 Section S2). Our seed sensitivity measure investigated in section 3.1 is most related to  
185 the metrics sequence coverage and island E-size, which measures two important metrics  
186 for sequence matching. While some sequence comparison algorithms may require a  
187 high fraction of matches for accurate similarity estimation and therefore optimize for  
188 number of matches, it is not typically needed for, e.g., read mapping applications. For  
189 example, it was shown in [SR2] that  $k$ -mers produce the highest fraction of matches  
190 under a random error distribution, but the high fraction of matches often occur because  
191 of several consecutive overlapping matches. In [SR1], the authors argued that this is  
192 bad as many of the matches are redundant, and they aim to select combinations of seeds  
193 that yield matches which overlap as little as possible which does not optimize for a high  
194 fraction of matches.

## 195 S.4.2 Data simulation

196 The performance of mixedstrobemers, altstrobemers and multistrobemers was benchmarked and  
197 compared to (spaced)  $k$ -mers and strobemers on simulated sequencing data in a similar  
198 scenario to the match evaluation in [SR2], described here for convenience. Random  
199 reference sequences of 10,000 nucleotides were simulated, whereby the probability of each  
200 of the four nucleotides was 25% for each position. To create the corresponding query  
201 sequence, 1%, 5% and 10% of the nucleotides of the reference string were mutated for the  
202 different experimental conditions. Insertions, deletions and substitutions were hereby  
203 added with equal probability of 1/3. To reduce sample variation bias, these simulations  
204 were repeated 1,000 times. The randomly-simulated references and queries were now  
205 used as inputs to seed  $k$ -mers, spaced  $k$ -mers, randstrobemers, minstrobemers, hybridstrobemers,  
206 mixedstrobemers, and altstrobemers. All  $k$ -mer and strobemer parameters were set as in [SR2],  
207 namely  $k = 30$  and  $(2,15,25,50)$ , altstrobemers with  $(2,10,20,25,50)$  and multistrobemers with  
208  $(2,5,25,25,50)$ , all yielding valid  $(30, 64)$ -seeds. Then, mixedstrobemers was sampled with  
209  $(2,15,25,50,q)$ ,  $q \in [0, 0.1, 0.2, \dots, 1.0]$ . We sampled mixedstrobemers by combining each  
210 strobemer type (randstrobemers, minstrobemers, hybridstrobemers, altstrobemers) with  $k$ -mers as it  
211 was easy with our experiment setup.

### 212 S.4.3 Simulated results

213 Our results demonstrate that  $k$ -mers performed best regarding the fraction of matches  
 214 (Fig. S6A), which decreased linearly with increasing fractions of strobemers. Contrary,  
 215 the match coverage (Fig. S6B) and the sequence coverage (Fig. S6C) increased with  
 216 higher strobemer fractions, even though it leveled off towards the highest strobemer  
 217 fractions. Sequence and match coverage were shown to be best for multistobes and  
 218 pure altstobes (100%) followed by randstobes and hybridstobes with a strobemer  
 219 content of 70-80% (Fig. S6C), largely outperforming  $k$ -mers and even being slightly  
 220 better than the pure strobemer implementations. We observed similar results when  
 221 looking at the expected island size which displayed an exponential decay behavior when  
 222 adding more strobemers, whereby the benefit of going beyond 80% strobemers was either  
 223 non-existent or very low (Fig. S6D).

224 Overall, mixedrandstobes with a strobe fraction of 70-80% were found to be superior  
 225 to the pure strobemer implementations as higher fraction of matches and a lower island  
 226 size was observed while maintaining high sequence and match coverage. Furthermore  
 227 our analysis suggests that altstobes are strictly better than the currently best known  
 228 strobemer construct (randstobes) as the fraction of matches, sequence coverage and  
 229 match coverage are higher while the number of expected island sizes remained lower.  
 230 The generalized implementation of altstobes, multistobes, performed even better as it  
 231 is outperforming altstobes in all matching metrics (Suppl. Table 1).

232 As the spaced  $k$ -mer protocol performed worse than  $k$ -mers on all metrics (fewer  
 233 matches, lower match coverage, and larger expected island size), their results are not  
 234 displayed for better visualisation. However, full data including spaced  $k$ -mers is pre-  
 235 sented in Suppl. Table 1.

236 We also looked at strobemers of orders 2 to 5. We observe that in practice,  $n \geq 4$   
 237 does not offer benefit over  $n = 2$  and 3 for the application of strobemers that we consider  
 238 in this study (Fig S6). Note that  $k = 30$ , can be divided by 2,3, and 5 without leaving  
 239 a remainder, guaranteeing same number of positions sampled. Also, altstobes and  
 240 multistobes were only computed for  $n = 2, 4$  to guarantee same number of positions  
 241 sampled as the other seeds. For strobemers of order 4, we observe that altstobes and  
 242 multistobes perform very similar which can be explained by their shortest possible  
 243 stobes being of identical length ( $k_s = 5$ ) and the considerably lower number of strobe  
 244 combinations (6 vs. 21) for multistobes of order 4.

### 245 S.4.4 *E.coli* Oxford Nanopore Technology Reads

246 In contrast to simulated data where we compare query sequences and their correspond-  
 247 ing references individually, biological data may contain spurious matches, which requires  
 248 measurements of the sequence matching metrics on biological data to reaffirm our pre-  
 249 vious simulation results.

250 To this end, the one thousand *E.coli* Oxford Nanopore Technology reads used in  
 251 [SR2] ranging from 17,360nt to 52,197nt (median 19,601nt) were split up in disjoint  
 252 segments of 2,000bp before computing the collinear chain solution of the raw hits for  
 253 each of the segments. The collinear chaining algorithm determines chains of matches  
 254 that are in identical order in both sequences. Hence, the collinear solution can be  
 255 viewed as a proxy for finding the true location of the hits by taking only the longest  
 256 collinear chain of the hits into account, as in [SR2], to avoid overcounting “spurious”  
 257 hits caused by local matches in repetitive regions that occur throughout the genome.  
 258 Raw unmerged hits were assessed rather than non-overlapping approximately matchings  
 259 (NAMs) to make the analysis identical to the simulated experiment. Subsequently, for  
 260 each read, the number of matches, the match coverage, the sequence coverage and the  
 261 expected island size were computed for the collinear solution of hits (Fig. 3 and Suppl.

262 Figs. S7 and S8).

### 263 S.4.5 Biological Data

264 The thousand *E.coli* reads were downloaded from Sequence Read Archive with Run  
 265 ID SRR13893500 (available here: [https://trace.ncbi.nlm.nih.gov/Traces/?view=](https://trace.ncbi.nlm.nih.gov/Traces/?view=run_browser&acc=SRR13893500&display=download)  
 266 [run\\_browser&acc=SRR13893500&display=download](https://trace.ncbi.nlm.nih.gov/Traces/?view=run_browser&acc=SRR13893500&display=download)). We selected the 1,000 longest  
 267 reads that aligned to the *E.coli* genome with more than 95% of the total read length  
 268 as in [SR2]. The reads were mapped to the *E.coli* genome assembly GCA\_003018135.1  
 269 ASM301813v1 (available here: [https://www.ncbi.nlm.nih.gov/genome/167?genome\\_](https://www.ncbi.nlm.nih.gov/genome/167?genome_assembly_id=368373)  
 270 [assembly\\_id=368373](https://www.ncbi.nlm.nih.gov/genome/167?genome_assembly_id=368373)). For the E-hits and uniqueness analysis we used human chro-  
 271 mosome 21 (NC\_000021.9) from GRCh38.p13 Primary Assembly. For the human data  
 272 analysis with minimap2, we used assembly CHM13 v2.2 (available here: [https://](https://s3-us-west-2.amazonaws.com/human-pangenomics/T2T/CHM13/assemblies/analysis_set/chm13v2.0.fa.gz)  
 273 [s3-us-west-2.amazonaws.com/human-pangenomics/T2T/CHM13/assemblies/analysis\\_](https://s3-us-west-2.amazonaws.com/human-pangenomics/T2T/CHM13/assemblies/analysis_set/chm13v2.0.fa.gz)  
 274 [set/chm13v2.0.fa.gz](https://s3-us-west-2.amazonaws.com/human-pangenomics/T2T/CHM13/assemblies/analysis_set/chm13v2.0.fa.gz)).

## 275 S.5 Minimap2 implementation details

276 We implemented randstrokes, mixedstrokes, altstrokes and multistrokes in the fast se-  
 277 quence mapping and alignment program minimap2. First, we find minimizer  $k$ -mers  
 278 based on the shorter strobe  $k_s$  (altstrokes and multistrokes) or the first strobe  $k/2$   
 279 (randstrokes and mixedstrokes). To decide whether a  $k$ -mer or a randstroke (mixed-  
 280 strokes) should be sampled, a modulo operation is performed on the hash value of the  
 281 first strobe (see mixedstroke implementation). Analogously, a decision about  $k_s$  and  $k_l$ ,  
 282 and  $k_1$  and  $k_2$  is made for altstrokes and multistrokes, respectively (see Section 2.6 of  
 283 the paper).

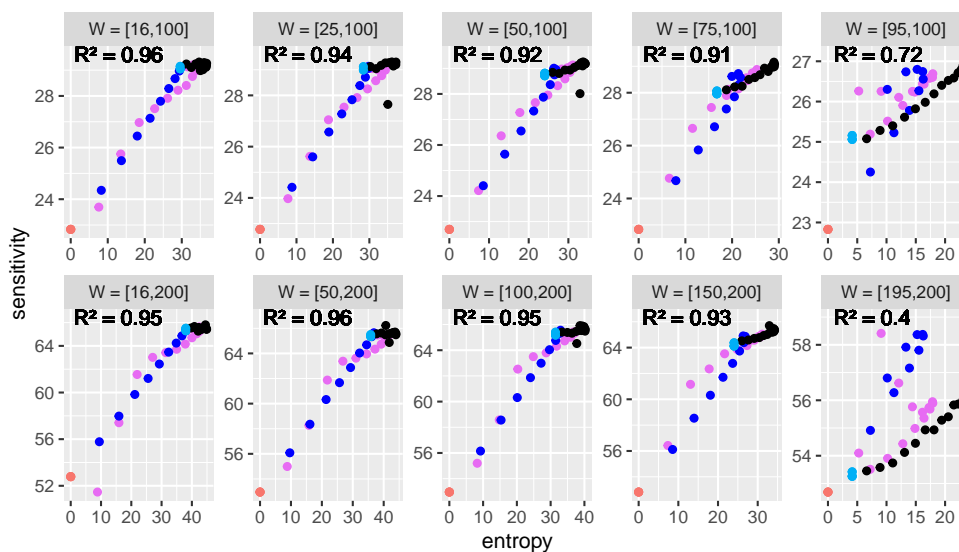
284 Next, the second strokes are selected from a downstream window [25,50] that min-  
 285 imizes the function  $\arg \min_p \{p : h(m \oplus S[p : p + \ell]) \leq h(m \oplus S[i' : i' + \ell]), \forall i' \in w'\}$ .  
 286 The hash values of the selected strokes are combined using non-commutative concate-  
 287 nation (hash values of linked  $k$ -mers) and are assigned to our fuzzy seeds of length  
 288  $w_{max} + k - 1$ . As it is not possible to sample seeds of these length from the reverse  
 289 strand at the beginning of the sequence, only forward seeds are sampled for the first  
 290  $w_{max} + k - 1$  nucleotides, and *vice versa* only reverse strand seeds at the end of the  
 291 sequence.

292

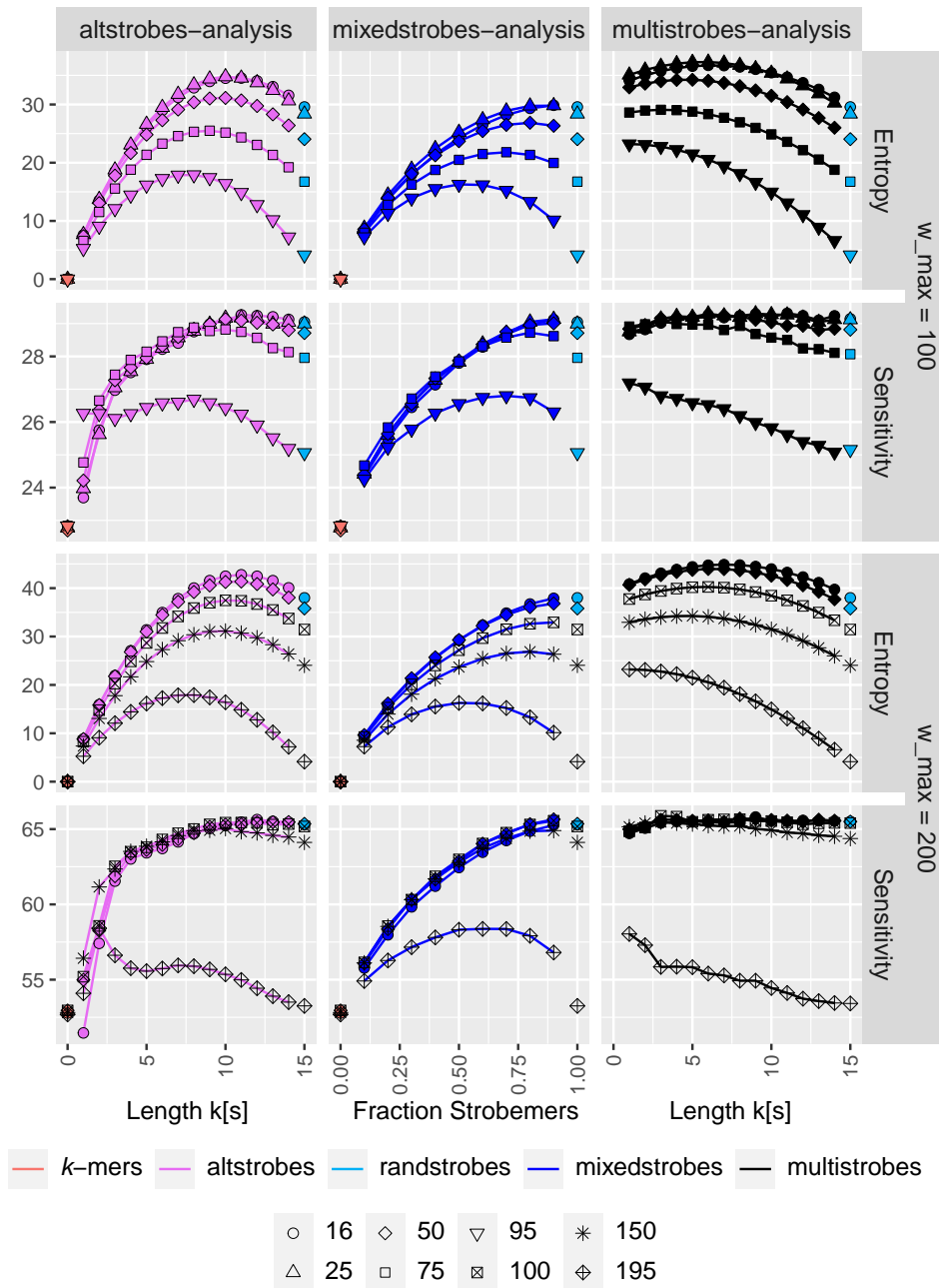
### 293 S.5.1 Minimap2 analysis

294 To measure speed and accuracy of our newly implemented seeding methods, we bench-  
 295 marked them against  $k$ -mers with  $k = 15$  (default setting) and  $k = 28$ . To this end,  
 296 we sampled 100,000 reads of length 10,000nt from random positions across the *Homo*  
 297 *sapiens* (human) genome assembly CHM13 that did not contain any incompletely spec-  
 298 ified bases. Next, we inserted, deleted, or substituted nucleotides with equal probability  
 299 of 1/3 each across the reads with the mutation rate 0.01, 0.05, 0.1, as in [SR2], and  
 300 converted the read sequence into its reverse-complement counterpart with probability  
 301 of 1/2. These mutated reads were mapped and aligned back to the CHM13 assembly  
 302 using minimap2 with default mapping and aligning settings. A read was considered to  
 303 be mapped correctly if at least 1nt was mapped to its correct location.

## 304 S.6 Figures and tables

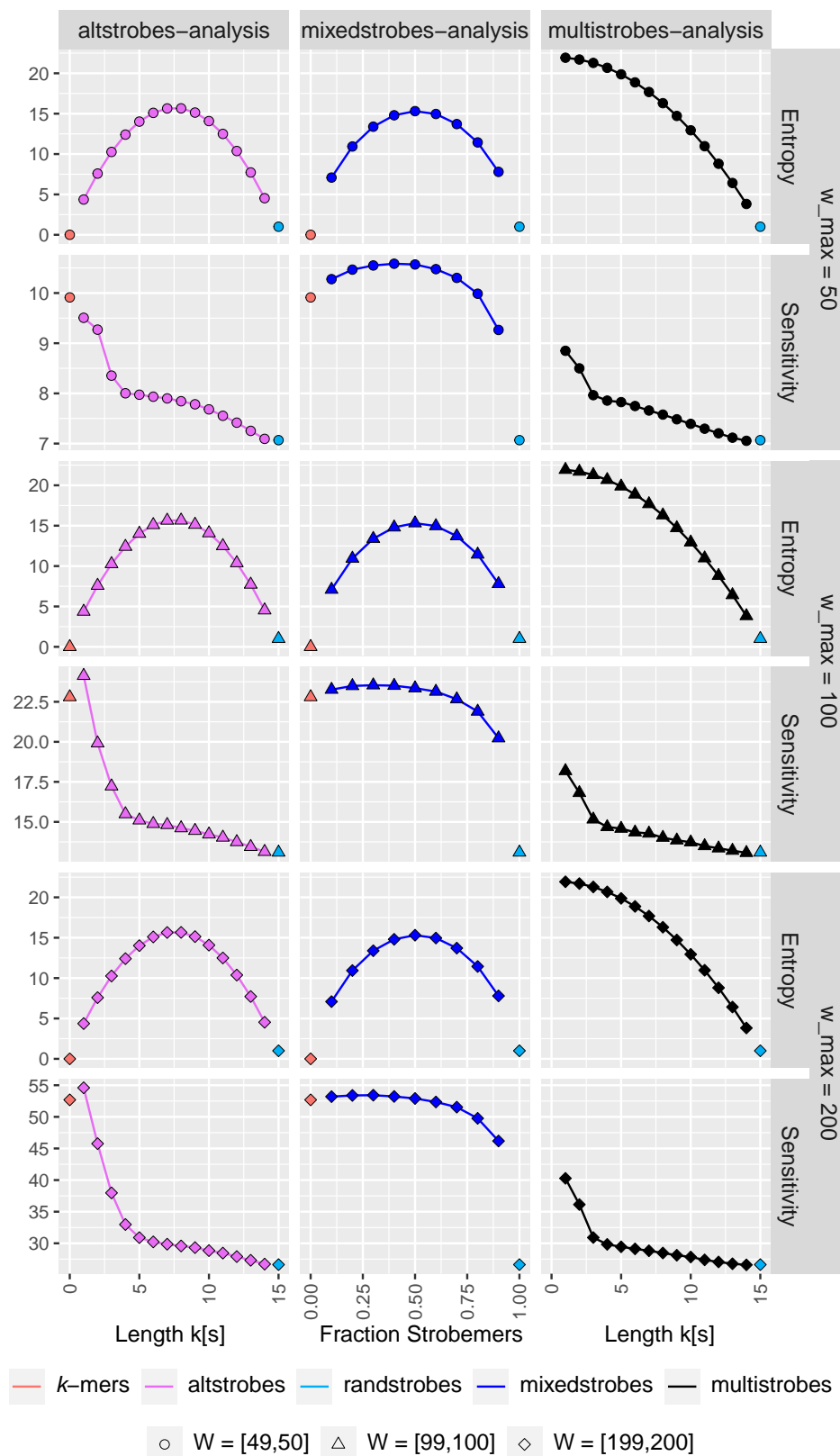


**Figure S1. Scatterplot showing the relationship between entropy and sensitivity.** The coefficient of determination (square of the Pearson correlation coefficient) to quantify whether entropy can predict sensitivity.

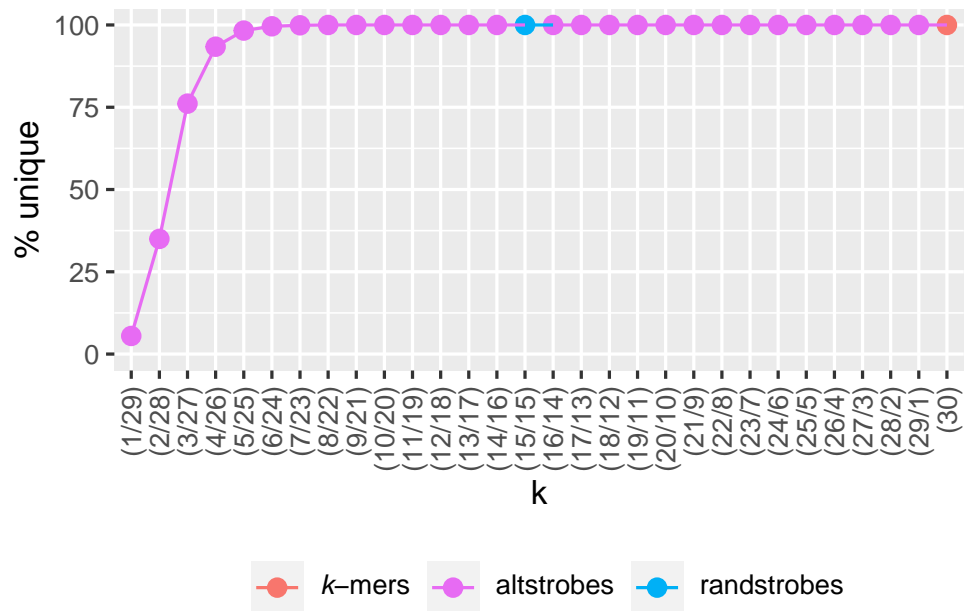


**Figure S2. Simulations showing how stochasticity in seed construct influence probability of  $w$  consecutive seeds producing at least one match in a region of length  $2w=128$  between sequences.**

Panel Entropy shows  $H(\mathbf{X})$  for  $k$ -mers, randstrokes, altstrokes (for different  $(k_s, k_l)$ ), mixedstrokes (for different  $q$ ) and multistrokes (for different  $(k_s, k_l)$ ). Panel Sensitivity shows  $P(N_m(30, 64) > 0)$  summed over  $m \in [1, 60]$  for various window sizes  $(w_{min}, 100)$ ,  $(w_{min} \in 16, 25, 50, 75, 95)$  and  $(w_{min}, 200)$ ,  $(w_{min} \in 16, 50, 150, 190, 195)$ .



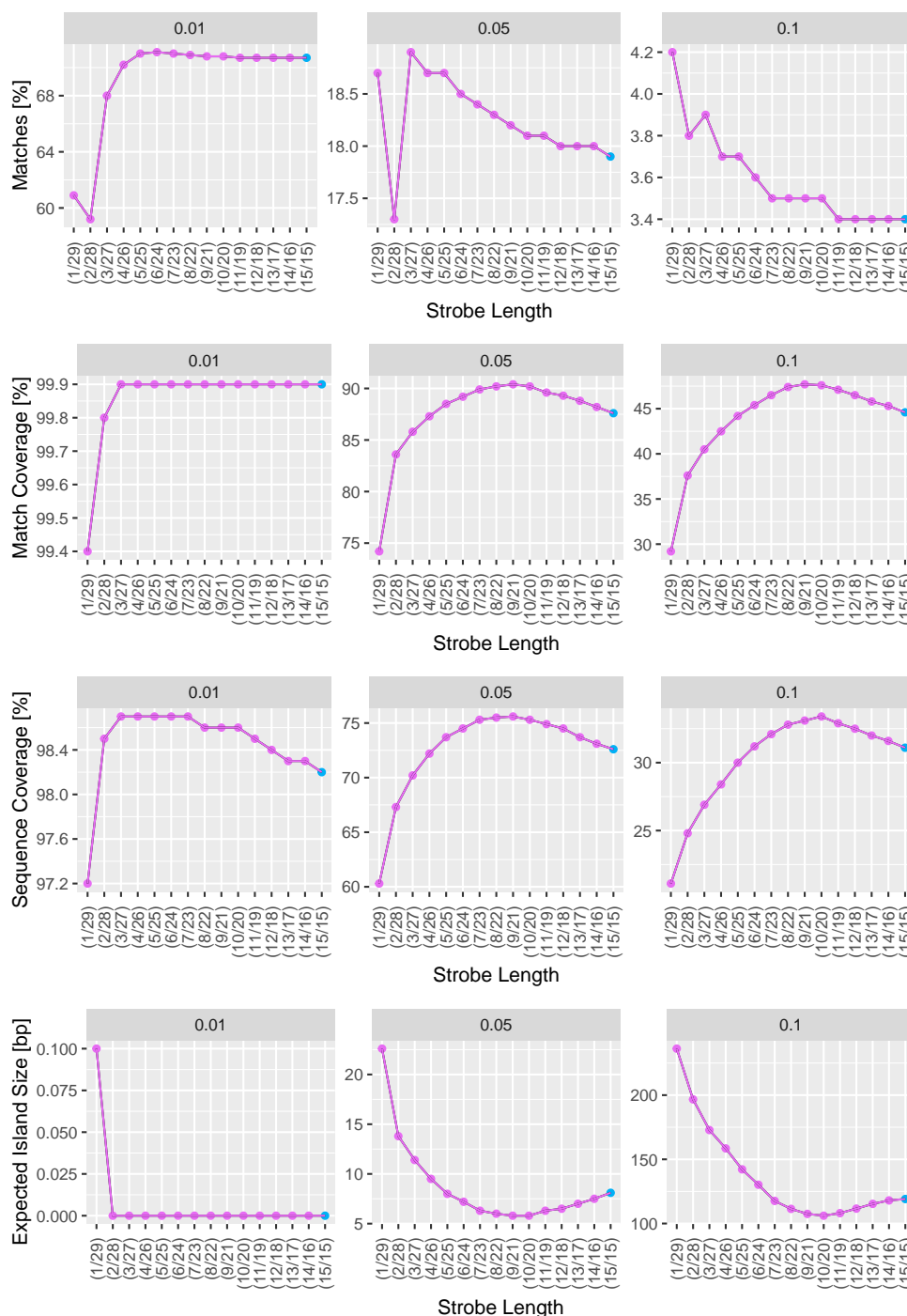
**Figure S3. Simulations showing entropies and sensitivity estimates for mixedstrokes, altstrokes, and multistrokes with very narrow window sizes ( $w_{min}, w_{max}$ ) of (49,50), (99,100), and (199,200)**  
 Panel Entropy shows  $H(\mathbf{X})$  for  $k$ -mers, randstrokes, altstrokes (for different  $(k_s, k_l)$ ), mixedstrokes (for different  $q$ ) and multistrokes (for different  $(k_s, k_l)$ ). Panel Sensitivity shows  $P(N_m(30, 64) > 0)$  summed over  $m \in [1, 60]$  for various small window sizes.



**Figure S4. Fraction of unique seeds for different altstrobe combinations in a region of 10,000nt.**

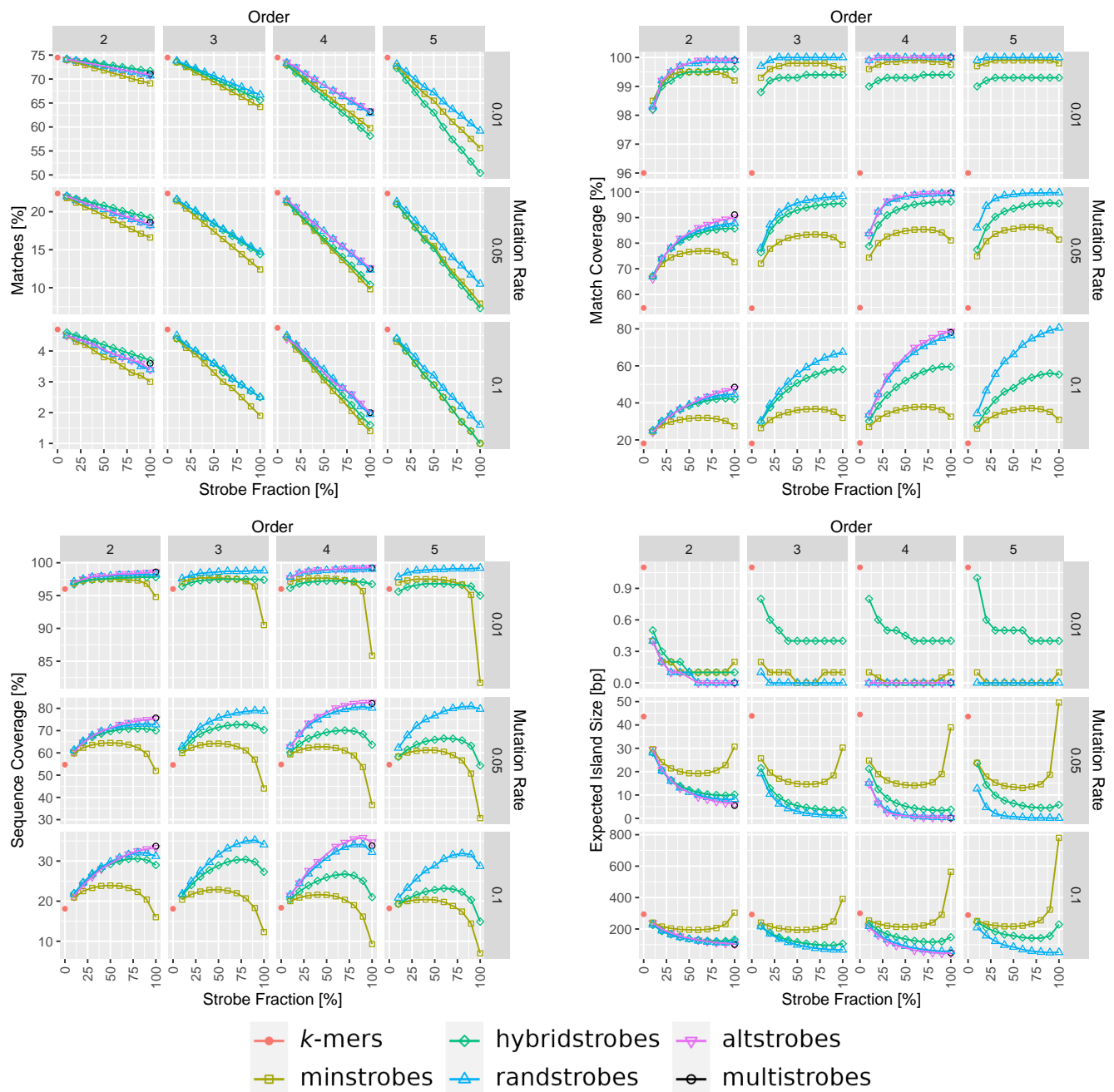
For a given altstrobe combination  $(k_1/k_2)$  in the plot, the uniqueness of seeds was computed with parameters  $(2, k_1/k_2, 25, 50)$  over simulated sequences of length 10,000nt. Note that all altstrobe seeds were sampled as  $k_1/k_2$ , which means that altstrokes(1/29) seeds only contain (1/29) seeds and no (29/1) seeds as opposed to all other altstrobe analysis in this manuscript.





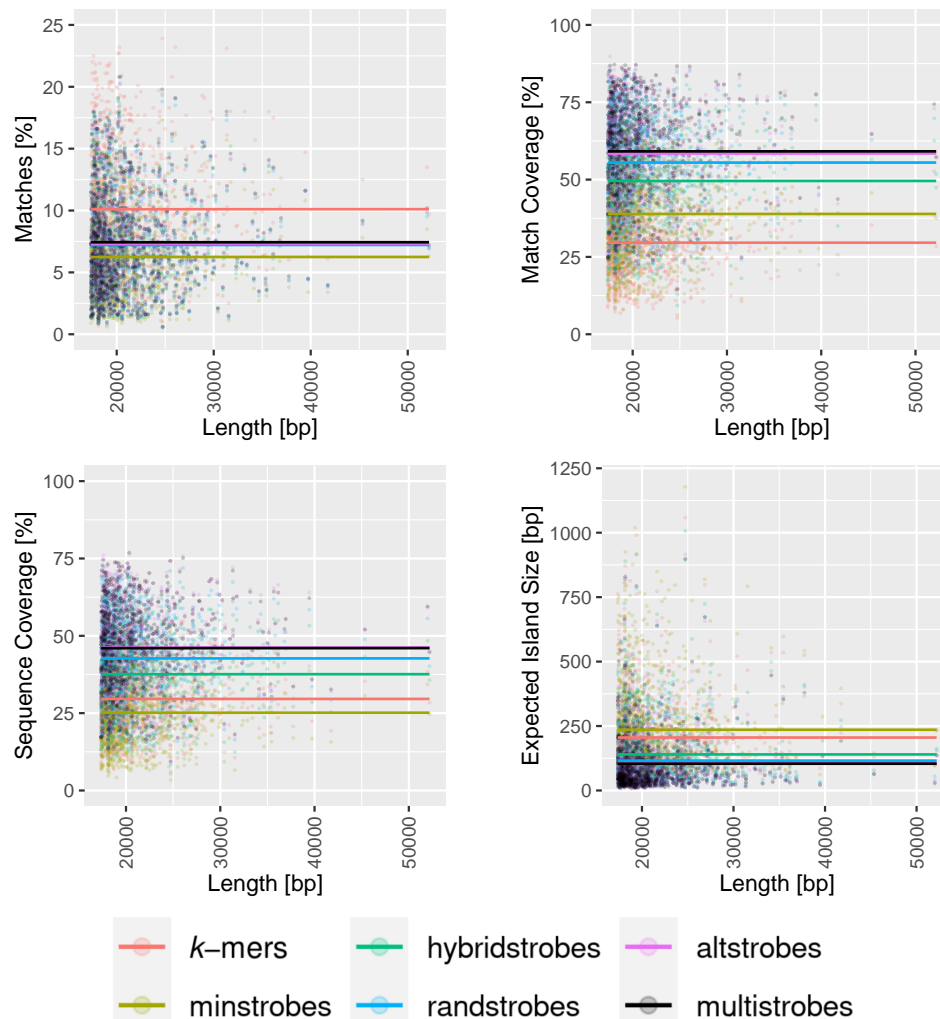
**Figure S5. Match Statistics for altstrokes with different strobe length combinations.**

For a given strobe length combination and mutation rate in the plot, 1,000 random sequences of length 10,000nt with randomly updated hash tables were generated. To guarantee an equal number of long-short and short-long altstrobe combinations even for very short strobos (see concerns prompted in section 1. Construction of Altstrobos), altstrobos were constructed with a modified altstrobe implementation. To this end, reference sequences were seeded with both altstrobos possibilities for each position and the decision about whether a short-long or a long-short query seed should be constructed was based on the longer strobe ensuring a large enough hash space to neglect systematic bias.



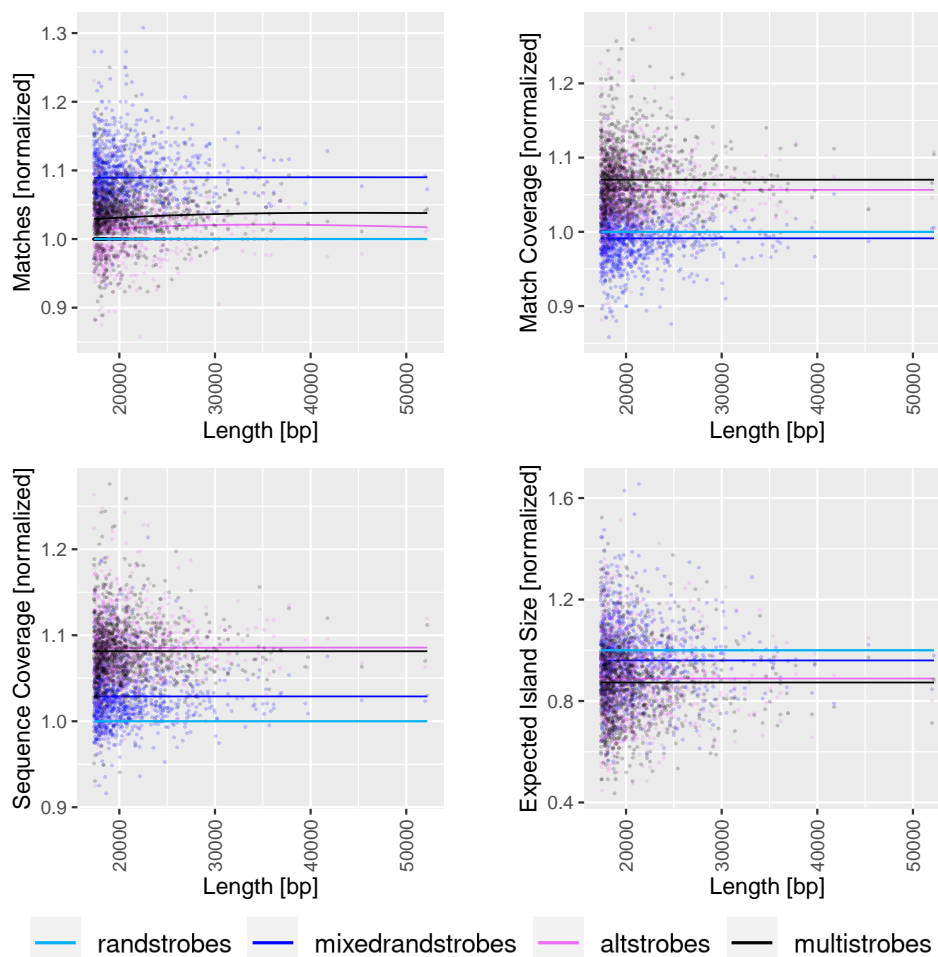
**Figure S6. Match Statistics for  $k$ -mers, strobemers, altstrobemers, multistrobemers and mixedstrobemers of orders 2 - 5.**

1,000 random DNA sequences of 10,000 nucleotides were generated and subsequently mutated to obtain references and queries. Using these sequences, seeds were sampled with strobemer fractions from 0% ( $k$ -mers) to 100% (pure strobemers), downstream windows set to [25,50] and all strobemers combined adding up to equal length subsequences of size 30 for better comparison. Hence, the strobemer settings were as follows: (2,15,25,50, $q$ ), (3,10,25,50, $q$ ), and (5,6,25,50, $q$ ), whereby mixedstrobemers were sampled with a strobe fraction  $q$  ranging from 0% ( $k$ -mers) to 100% (strobemers) and a step size of 10% ( $f(0, 100, 10) = \{10k \mid k \in \{0, 10, \dots, 100\}\}$ ). For strobemers of order 4, seeds of lengths 28 and 32 were seeded and the average (mean) plotted for each metric to ensure similar sizes of subsequences between the protocols. Altstrobemers were seeded with strobemers of length  $k_s = 10/k_l = 20$  (2,10,20,25,50, $q$ ), and  $k_s = 5/k_l = 10$  (4,5,10,25,50, $q$ ), respectively. Multistrobemers were seeded as (2,5,25,25,50) and (4,5,10,25,50).



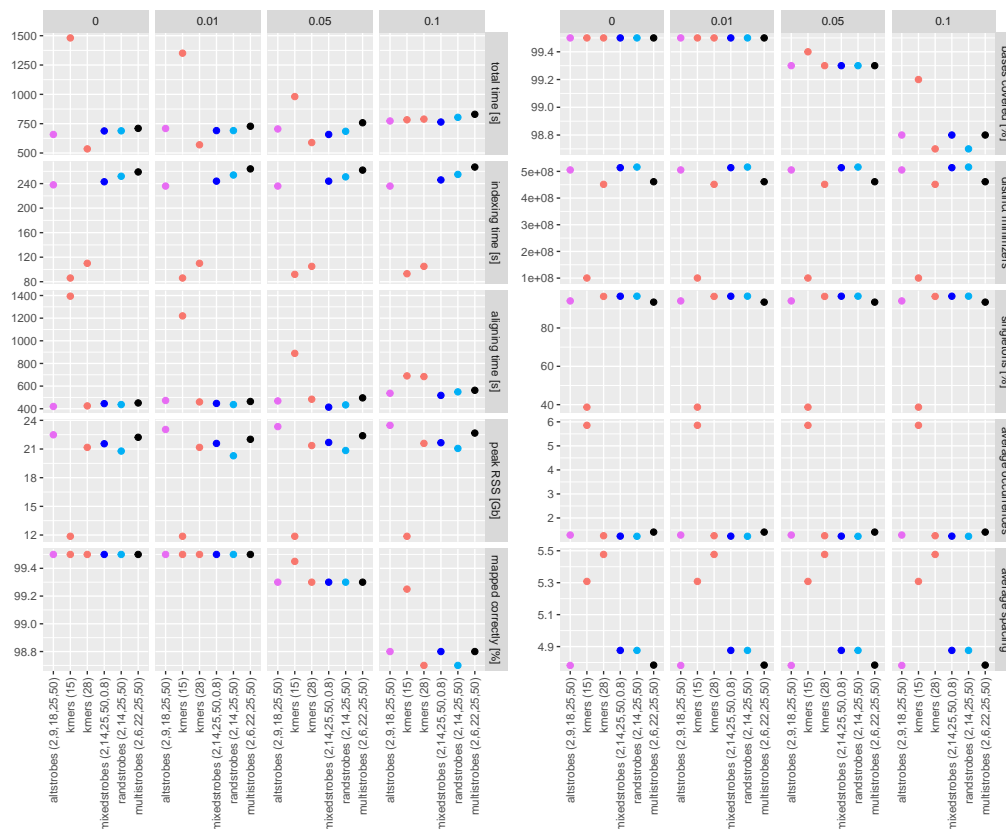
**Figure S7. Comparison between  $k$ -mers, strobemers, altstobes and multistobes when mapping genomic ONT reads for reads of different lengths (x-axis).**

The *E.coli* reads were split up in long disjoint segments of 2,000nt. Next, the segments were seeded with  $k$ -mers, strobemers and altstobes, downstream windows set to [25,50] and all stobes combined adding up to equal length subsequences of size 30 for better comparison. Then for each segment, the collinear solution of raw hits was computed to subsequently quantify number of matches, match coverage, sequence coverage and expected island size for each read. Each dot represents one read while the line displays a smoothed conditional mean (GAM curve with cubic spline:  $y \sim s(x, bs = "cs")$ ).



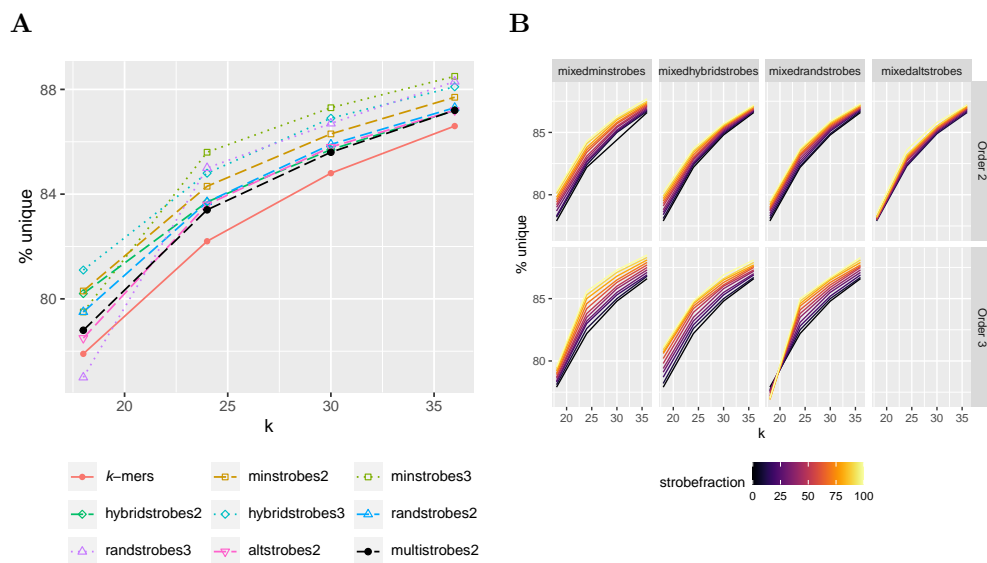
**Figure S8. Comparison between randstrobos, mixedrandstrobos, altstrobos and multistrobos when mapping genomic ONT reads for reads of different lengths (x-axis).**

The *E.coli* reads were split up in long disjoint segments of 2,000nt. Next, the segments were seeded with randstrobos (2,15,25,50), mixedrandstrobos (2,15,25,50,0.8), altstrobos (2,10,20,25,50) and multistrobos (2,5,25,25,50). For each segment, the collinear solution of raw hits was computed to subsequently quantify number of matches, match coverage, sequence coverage and expected island size for each read. For each read, the matching metrics from mixedrandstrobos (blue), altstrobos (pink) and multistrobos (black) were subsequently normalized by randstrobos (turquoise) for better visualisation. Multistrobos and altstrobos perform best indicated by similar number of matches, higher sequence and match coverage as well as lower gap size. Mixedrandstrobos perform better than randstrobos for all metrics besides match coverage where mixedrandstrobos perform roughly 1% worse. Each dot represents one read while the line displays a smoothed conditional mean (GAM curve with cubic spline:  $y \sim s(x, bs = "cs")$ ).



**Figure S9. Minimap2 Benchmarking.**

100,000 sequences of lengths 10,000nt without any incompletely specified bases (e.g. "n") were selected from random positions of the CHM13v2.0 human assembly and mutated (insertions, deletions and base substitutions) with mutation frequencies of 0%, 1%, 5% and 10%. The selected reads were mapped and aligned back to the reference using  $k$ -mers ( $k = 15$  and  $28$ ), altstrobobes (2, 9, 18, 25, 50), randstrobobes (2, 14, 25, 50), multistrobobes (2, 5, 23, 25, 50) and mixedstrobobes (2, 14, 25, 40, 0.8). All experiments were repeated 5 times and the average (mean) taken to account for variance in computer processing speed.



**Figure S10. Fraction of unique seeds for different strobemer seeding techniques and fractions on Human chromosome 21.**

Chromosome 21 of the human GRCh38 assembly was seeded with *k*-mers, strobemers of orders 2 ( $2, k/2, 25, 50, q$ ) and 3 ( $3, k/3, 25, 50, q$ ) as well as altstrobes of order 2 ( $2, (k/3), (2k/3), 25, 50, q$ ). Multistrobes was seeded with  $(2, 5, k-5, 25, 50)$ . In panel B, the numbers of extracted nucleotides ( $k=30$ ) is the same for all seeding techniques for fair comparison.

w = 1

Mutation Rate		0.01				0.05				0.1			
Seeding	Settings	m	mc	sc	e	m	mc	sc	e	m	mc	sc	e
<i>k</i> -mers	30	<b>74.5</b>	96.0	96.0	1.1	<b>22.4</b>	54.7	54.7	43.6	<b>4.7</b>	18.1	18.1	293.8
spaced <i>k</i> -mers	dense	67.6	95.6	96.2	1.5	13.8	53.9	50.9	65.7	1.8	16.3	14.2	481.0
spaced <i>k</i> -mers	sparse	50.5	89.8	87.8	11.1	3.5	26.8	21.4	493.5	0.1	3.6	2.1	4120.8
altstobes	(2,10,20,25,50)	70.9	<b>99.9</b>	<b>98.6</b>	<b>0.0</b>	18.3	90.0	75.3	6.1	3.4	47.7	33.3	104.4
randstobes	(2,15,25,50)	70.7	<b>99.9</b>	98.2	<b>0.0</b>	18.2	87.8	72.7	8.2	3.4	44.6	21.2	118.1
hybridstobes	(2,15,25,50)	71.7	99.6	97.8	0.1	19.2	85.7	70.0	10.2	3.7	42.0	29.0	132.7
minstobes	(2,15,25,50)	69.1	99.2	94.8	0.2	16.6	72.6	51.9	30.7	3.0	27.4	16.0	304.0
mixedstobes	(2,15,25,50,0.8)	71.5	<b>99.9</b>	98.2	0.1	19.0	86.8	72.9	8.4	3.7	43.8	32.1	113.4
multistobes	(2,5,25,25,50)	71.1	<b>99.9</b>	<b>98.6</b>	<b>0.0</b>	18.6	<b>91.1</b>	<b>75.7</b>	<b>5.5</b>	3.6	<b>48.5</b>	<b>33.7</b>	<b>101.2</b>

w = 10

Mutation Rate		0.01				0.05				0.1			
Seeding	Settings	m	mc	sc	e	m	mc	sc	e	m	mc	sc	e
<i>k</i> -mers	30	<b>73.2</b>	90.3	90.3	2.6	<b>20.7</b>	42.8	42.8	73.2	<b>3.9</b>	11.4	11.4	501.8
spaced <i>k</i> -mers	dense	65.5	90.9	87.3	3.9	12.1	36.5	30.9	147.8	1.4	6.9	5.3	1265.4
spaced <i>k</i> -mers	sparse	47.9	84.9	74.4	17.0	2.7	16.7	9.5	945.9	0.1	1.2	0.5	7140.2
altstobes	(2,10,20,25,50)	69.7	98.5	92.3	<b>0.4</b>	17.1	62.7	46.0	46.9	3.0	19.2	12.1	423.3
randstobes	(2,15,25,50)	69.6	98.4	92.3	0.5	16.8	62.3	45.9	48.7	2.9	18.5	11.8	451.6
hybridstobes	(2,15,25,50)	68.4	97.3	90.2	1.2	15.8	58.6	42.5	58.0	2.6	16.8	10.5	506.6
minstobes	(2,15,25,50)	68.0	98.1	87.8	0.7	15.2	58.8	37.1	67.4	2.5	16.5	8.8	611.2
mixedstobes	(2,15,25,50,0.8)	70.5	98.0	92.2	0.7	17.8	60.4	46.8	49.2	3.2	17.8	12.4	425.6
multistobes	(2,5,25,25,50)	70.3	<b>98.6</b>	<b>92.4</b>	<b>0.4</b>	17.2	<b>64.2</b>	<b>47.0</b>	<b>44.1</b>	3.3	<b>20.2</b>	<b>12.7</b>	<b>404.2</b>

w = 20

Mutation Rate		0.01				0.05				0.1			
Seeding	Settings	m	mc	sc	e	m	mc	sc	e	m	mc	sc	e
<i>k</i> -mers	30	<b>71.8</b>	84.2	84.2	5.8	<b>19.3</b>	33.2	33.2	111.1	<b>3.6</b>	7.8	7.8	737.4
spaced <i>k</i> -mers	dense	64.4	85.0	77.9	8.0	11.4	27.9	21.9	221.3	1.4	4.5	3.3	1931.3
spaced <i>k</i> -mers	sparse	47.3	79.9	62.2	26.6	2.6	12.2	5.9	1256.4	0.1	0.7	0.3	8242.6
altstobes	(2,10,20,25,50)	68.9	94.8	83.6	2.3	16.4	47.0	32.1	98.5	2.9	11.6	7.1	766.8
randstobes	(2,15,25,50)	68.3	94.8	<b>83.9</b>	2.4	15.6	46.0	31.2	101.3	2.7	11.0	6.7	804.9
hybridstobes	(2,15,25,50)	66.2	92.6	81.8	3.8	13.8	41.1	28.1	124.8	2.2	9.2	5.6	993.7
minstobes	(2,15,25,50)	66.7	95.3	75.8	2.6	14.0	45.4	25.6	129.3	2.2	10.4	5.2	1020.3
mixedstobes	(2,15,25,50,0.8)	69.3	93.3	83.8	2.9	16.6	44.3	32.4	99.0	2.9	10.6	7.2	768.5
multistobes	(2,5,25,25,50)	69.8	<b>95.0</b>	83.6	<b>2.2</b>	17.2	<b>48.3</b>	<b>32.8</b>	<b>93.9</b>	3.1	<b>12.2</b>	<b>7.4</b>	<b>722.8</b>

**Table S1: Match statistics for simulated sequences ( $L = 10000$ ) under different sampling protocols under mutations rates of 0.01, 0.05, 0.1 using minimizer thinning with  $w = 1$  (no thinning),  $w = 10$ , and  $w = 20$ .**

Here, m denotes the number of matches as a percentage of the total number of extracted subsequences for the protocol, sc (sequence coverage) and mc (match coverage) is shown as the percentage of the total sequence length, and E is the expected island size. Boldfaced values indicate the most desirable result across protocols for each of the match statistics.

305 **Supplementary References**

- 306 [SR1] Ilie, L., Ilie, S.: Multiple spaced seeds for homology search. *Bioinformatics* **23**(22),  
307 2969–2977 (09 2007). <https://doi.org/10.1093/bioinformatics/btm422>, [https://](https://doi.org/10.1093/bioinformatics/btm422)  
308 [doi.org/10.1093/bioinformatics/btm422](https://doi.org/10.1093/bioinformatics/btm422)
- 309 [SR2] Sahlin, K.: Effective sequence similarity detection with strobemers. *Genome*  
310 *research* **31**(11), 2080–2094 (Nov 2021). <https://doi.org/10.1101/gr.275648.121>,  
311 <https://pubmed.ncbi.nlm.nih.gov/34667119>, 34667119[pmid]