



**HAL**  
open science

# Towards Efficient k-Mer Set Operations via Function-Assigned Masked Superstrings

Ondřej Sladký, Pavel Veselý, Karel Břinda

► **To cite this version:**

Ondřej Sladký, Pavel Veselý, Karel Břinda. Towards Efficient k-Mer Set Operations via Function-Assigned Masked Superstrings. 2024. hal-04573444v2

**HAL Id: hal-04573444**

**<https://hal.science/hal-04573444v2>**

Preprint submitted on 18 Nov 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

# 1 Towards Efficient $k$ -Mer Set Operations via 2 Function-Assigned Masked Superstrings

3 Ondřej Sladký ✉ 

4 Computer Science Institute of Charles University, Prague, Czechia  
5 ETH Zürich, Switzerland

6 Pavel Veselý ✉ 

7 Computer Science Institute of Charles University, Prague, Czechia

8 Karel Břinda ✉ 

9 Inria, Irista, Univ. Rennes, 35042 Rennes, France

## 10 — Abstract —

---

11 The design of efficient dynamic data structures for large  $k$ -mer sets belongs to central chal-  
12 lenges of sequence bioinformatics. Recent advances in compact  $k$ -mer set representations via  
13 simplitigs/Spectrum-Preserving String Sets, culminating with the masked superstring framework,  
14 have provided data structures of remarkable space efficiency for wide ranges of  $k$ -mer sets. However,  
15 the possibility to perform set operations remained limited due to the static nature of the underlying  
16 compact representations. Here, we develop  $f$ -masked superstrings, a concept combining masked  
17 superstrings with custom demasking functions  $f$  to enable efficient  $k$ -mer set operations via string  
18 concatenation. Combined with the FMSI index for masked superstrings, we obtain a memory-efficient  
19  $k$ -mer index supporting set operations via Burrows-Wheeler Transform merging. The framework  
20 provides a promising theoretical solution to a pressing bioinformatics problem and highlights the  
21 potential of  $f$ -masked superstrings to become an elementary data type for  $k$ -mer sets.

22 **2012 ACM Subject Classification** Applied computing → Bioinformatics; Theory of computation →  
23 Pattern matching

24 **Keywords and phrases**  $k$ -mer sets, data structures, set operations, masked superstrings

25 **Funding** *Ondřej Sladký*: Supported by GA ČR project 22-22997S and ERC-CZ project LL2406 of  
26 the Ministry of Education of Czech Republic.

27 *Pavel Veselý*: Supported by GA ČR project 22-22997S, ERC-CZ project LL2406 of the Ministry of  
28 Education of Czech Republic, and Center for Foundations of Modern Computer Science (Charles  
29 Univ. project UNCE 24/SCI/008).

30 *Karel Břinda*: Supported by French National Research Agency (ANR) under Grant ANR-24-CE45-  
31 1226 for the REALL project.

## 32 **1** Introduction

33 To store and analyze the vast volumes of DNA sequencing data [61], modern bioinformatics  
 34 methods increasingly rely on *k*-mers to bypass the computationally expensive sequence  
 35 alignment and mitigate data heterogeneity. *k*-mer-based methods are used in a range of  
 36 applications, including large-scale data search [6, 8, 14, 33], metagenomic classification [15, 62],  
 37 infectious disease diagnostics [9, 13], assembly evaluation [52], and transcript abundance  
 38 quantification [10, 48]. All such applications depend on the efficiency of the underlying *k*-mer  
 39 data structures.

40 The design of efficient *k*-mer data structures is a central challenge of contemporary  
 41 sequence bioinformatics [18, 43]. We summarize the desired properties of an ideal data  
 42 structure in Section 1.1, but the two main requirements are as follows. First, as even a  
 43 single *k*-mer set can be large, potentially containing up to hundreds of billions distinct  
 44 *k*-mers [33], the main challenge is to provide data structures that are efficient in space and  
 45 time simultaneously. Second, as modern genomic databases undergo rapid development,  
 46 thanks to their growing content and curation efforts, rapid and space-efficient updates across  
 47 *k*-mer indexes are increasingly needed to avoid repetitive and costly index recomputations.  
 48 This includes scenarios such as *k*-mer set operations across sets, and additions or removals of  
 49 individual *k*-mers.

50 The space efficiency has been particularly well addressed using the (repetitive) Spectrum-  
 51 Preserving String Sets, (r)SPSS, that exploit *k*-mer non-independence [17]. In particular,  
 52 the key observation is that genomic *k*-mer sets can be typically represented by *k*-long  
 53 substrings of a small number of (arbitrarily long) strings, which is the so-called *spectrum-like*  
 54 *property* (SLP) [18]. Building on this observation, several textual representations of *k*-mer  
 55 sets were proposed. The first ones were based on de Bruijn graphs, to which we jointly  
 56 refer as (repetitive) *Spectrum Preserving String Sets* or (r)SPSS [11, 12, 19, 20, 51, 57, 58].  
 57 (r)SPSS are currently highly standard and many data structures for *k*-mer sets base upon  
 58 them [2, 7, 45, 49, 55],

59 *Masked superstrings* (MS) have provided additional space gains [59], by the virtue of  
 60 a better *k*-mer set compaction. The core improvement over (r)SPSS lies in modeling the  
 61 structure of *k*-mer sets by overlap graphs instead of de Bruijn graphs, thus being able to  
 62 exploit overlaps of *any* length. MS represent *k*-mer sets using an approximately shortest  
 63 superstring of all *k*-mers and a binary mask to avoid false positive *k*-mers. MS generalize  
 64 any existing (r)SPSS representation as these can always be encoded as MS, but provide  
 65 further compression power, especially for sets without the SLP, such as those arising from  
 66 sketching or subsampling. The resulting representation is well indexable using a technique  
 67 called Masked Burrows Wheeler Transform [60], resulting in a *k*-mer data structure with  
 68  $2 + o(1)$  bits per *k*-mer under the SLP.

69 However, the lack of dynamicity of both (r)SPSS and MS – and of their derived data  
 70 structures – has been limiting their wider applicability. To the best of our knowledge, the  
 71 only supported operations were union, either via merging (r)SPSS of several *k*-mer sets  
 72 resulting in an rSPSS of their union, or by the Cdbg-Tricks [28] to calculate the union unitigs  
 73 from unitigs of multiple original *k*-mer sets. However, other operations than union, such as  
 74 intersection or symmetric difference have never been considered.

75 Here, we develop a dynamic variant of masked superstrings (and thus also of (r)SPSS)  
 76 called the *f-masked superstrings* (*f*-MS). The key idea is to equip masked superstrings  
 77 with so-called demasking functions *f* for more flexible mask interpretation (**Section 3**).  
 78 When complemented with the concatenation operation, this provides support for any set

79 operation, including union, intersection, and symmetric or asymmetric difference on the  
 80 represented  $k$ -mer sets, resulting in a complete algebraic type for  $k$ -mer sets (**Section 4**). We  
 81 implement the whole approach in the FMSI index for masked superstrings [60] (**Section 5**),  
 82 and demonstrate applicability of the concept on an example dataset (**Section 6**).

## 83 1.1 Problem formulation

84 Before delving into technical details of our contribution, we formulate properties of an ideal  
 85  $k$ -mer data structure. The desired data structure for representing a set of  $k$ -mers should  
 86 support the following operations in a space- and time-efficient way:

- 87 (i)  **$k$ -mer set index construction**, with time complexity linear with the number of  $k$ -mers.
- 88 (ii)  **$k$ -mer membership queries**, with low bits-per- $k$ -mer memory requirements, approach-  
 89 ing 2 bits per distinct  $k$ -mer for datasets following the Spectrum-Like Property [18] and  
 90 proportionally for more complex datasets.
- 91 (iii) **Set operations**, including union, intersection, difference, and symmetric difference.
- 92 (iv) **Single  $k$ -mer deletion and insertion**, or at least one of these.

93 In this paper, we make progress on efficient set operations while retaining high space  
 94 efficiency for  $k$ -mer membership queries and other desirable properties of (r)SPSS- or MS-  
 95 based data structures.

96 **Naïve approaches for performing set operations.** We note that one can add support  
 97 for set operations to a static data structure in a straightforward way: One option is extracting  
 98 the  $k$ -mer sets from the input indexes, performing the given operation with the sets, and  
 99 computing the new index for the resulting set; however, this process requires substantial time  
 100 and memory. Another option is to keep the indexes for input  $k$ -mer sets and process a  $k$ -mer  
 101 query on the set resulting from the operation by asking each index for the presence/absence  
 102 of the  $k$ -mer in each input set, which is however time and memory inefficient as all the  
 103 indexes need to be loaded into memory, as also noted in [30]. Therefore, we seek to perform  
 104 set operations without the costly operation of recomputing the index or making multiple  
 105 queries to original indexes.

## 106 1.2 Related Work

107 Many works have recently focused on data structures for single  $k$ -mer sets and their collections;  
 108 we refer to [18, 43] for recent surveys. Here, we primarily focus on those that offer some kind  
 109 of dynamicity, i.e., an efficient support for set operations or, at least, insertions/deletions of  
 110 individual  $k$ -mers. The recently introduced Conway-Bromage-Lyndon (CBL) structure [46]  
 111 builds on the work of Conway and Bromage [21] on sparse bit-vector encodings and combines  
 112 them with smallest cyclic rotations of  $k$ -mers (a.k.a. Lyndon words), which yields a dynamic  
 113 and exact  $k$ -mer index supporting set operations, such as union, intersection, and difference,  
 114 as well as insertions or deletions of  $k$ -mers. However, the memory requirements for processing  
 115 queries on dataset satisfying the spectrum-like property are substantially worse than for  
 116 other, static methods [46].

117 While, to the best of our knowledge, other  $k$ -mer indexes do not support efficient  
 118 set operations such as the intersection or difference, other tools, including BufBoss [1],  
 119 DynamicBoss [4], and FDBG [22] allow for efficient insertions and deletions of individual  
 120  $k$ -mers. Bifrost [29], VARI-merge [47], Metagraph [33], dynamic Mantis [5], or the very  
 121 recent Cdbg-Tricks [28] support insertions but not deletions. Other data structures, such as  
 122 COBS [6], RAMBO [27], or kmtricks [35] trade exactness for space compression, allowing a  
 123 certain false probability rate. These employ variants of the Bloom filter that can also process

124 insertions and compute unions efficiently but other set operations such as the intersection  
 125 are not directly possible. We note that there are many more highly efficient but static data  
 126 structures for individual or multiple *k*-mer sets, e.g., [2, 3, 7, 24, 25, 29, 40, 44, 45, 49].  
 127 Finally, one can also use the textual representations of (r)SPSS, with efficient string indexes  
 128 such as the FM index [26] or BWA [36–39], but this only yields static indexes.

129 Another line of work, e.g. [23, 32, 34, 41, 42, 52, 53], focused on *k*-mer counting, where  
 130 we additionally require to compute the *k*-mer frequencies. Out of the many *k*-mer counters  
 131 available, GenomeTester4 [32], KMC3 [34], or Meryl [52] support operations such as union,  
 132 intersection, or difference on multisets. However, *k*-mer counters require substantially larger  
 133 memory than the most efficient *k*-mer indexes or heavily utilize disk.

## 134 2 Preliminaries

135 **Strings and *k*-mers.** We use constant-size alphabets  $\Sigma$ , typically the nucleotide alphabet  
 136  $\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}$  (unless stated otherwise). The set  $\Sigma^*$  contains all finite strings over  $\Sigma$ , with  
 137  $\epsilon$  representing the empty string. For a given string  $S \in \Sigma^*$ ,  $|S|$  denotes its length, and  $|S|_c$   
 138 the number of occurrences of the letter  $c$  in  $S$ . For two strings  $S$  and  $T$ ,  $S + T$  denotes their  
 139 concatenation. A *k*-mer is a *k*-long string over  $\Sigma$ , and unless stated otherwise, we assume  
 140 *canonical k-mers*, i.e., a *k*-mer and its reverse complement are considered equal. For a string  
 141  $S$  and a fixed length  $k$ , the *k*-mers *generated* by  $S$  are all *k*-long substrings of  $S$ , and similarly  
 142 for a set of strings  $\mathcal{R}$ , they are those generated by the individual strings  $S \in \mathcal{R}$ .

143 **(r)SPSS representations of *k*-mer sets.** For a given *k*-mer set  $K$  (comprising *k*-mers of the  
 144 same size  $k$ ), a set of strings  $\mathcal{R}$  (called *simplitigs*) constitutes its *Spectrum Preserving String*  
 145 *Set* (SPSS) [11, 12, 50, 51] if (1) it generates  $K$  and (2) each *k*-mer occurs only once among  
 146 the simplitigs (uniqueness). Repetitive Spectrum Preserving String Sets (rSPSS) [56, 58],  
 147 including its members called *matchtigs*, are defined similarly, except that condition (2) is  
 148 omitted. To refer to either of these concepts, we use the abbreviation (r)SPSS .

149 **Masked superstrings.** Given a *k*-mer set  $K$ , a *masked superstring* (MS) [59] consists of a  
 150 pair  $(S, M)$ , where  $S$  is an arbitrary superstring of the *k*-mers in  $K$  and  $M$  is a binary mask  
 151 of the same length. An occurrence of a *k*-mer in an MS is said to be ON if there is 1 at the  
 152 corresponding position in the mask (i.e., the initial position of the occurrence), and OFF  
 153 otherwise. The set of *k*-mers generated by  $S$  are referred to as the *appearing k-mers*, and are  
 154 further classified based on their values in the mask into *represented k-mers* (at least one ON  
 155 occurrence) and *ghost k-mers* (all occurrences OFF). All masks  $M$  that represent a given  $K$   
 156 in a combination with a given superstring are called *compatible*.

157 **Encoding conventions.** To implicitly encode  $k$  in MS, the last  $k - 1$  positions of masks  
 158 are always set to 0 and preceded by 1 (always feasible). In this paper, MS are presented  
 159 using their mask-case encoding, i.e., the superstring with individual letters in either lower or  
 160 upper case indicating a 1 or 0 in the mask, respectively.

161 ► **Example 1.** Consider the *k*-mer set  $K = \{\mathbf{ACG}, \mathbf{GGG}\}$ . One possible superstring is  $\mathbf{ACGGGG}$ ,  
 162 with three compatible masks:  $\mathbf{101100}$ ,  $\mathbf{100100}$ ,  $\mathbf{101000}$ , resulting in the encodings  $\mathbf{AcGGgg}$ ,  
 163  $\mathbf{AcgGgg}$ ,  $\mathbf{AcGggg}$ , respectively. The latter mask would contravene our encoding assumptions,  
 164 and thus, would not be used for MS storage. When parsing  $\mathbf{AcgGgg}$ , the suffix implies  $k = 3$ ,  
 165 and the *k*-mers are decoded as  $\{\mathbf{ACG}, \mathbf{GGG}\}$ .

### 3 Function-Assigned Masked Superstrings

Let  $(M, S)$  be an MS (Masked Superstring) and suppose that our objective is to determine whether a given  $k$ -mer  $Q$  is among the MS-represented  $k$ -mers. Conceptually, this process consists of two steps: (1) identify the occurrences of  $Q$  in  $S$ , and (2) verify using the mask  $M$  whether at least one occurrence of  $Q$  is ON. We can formalize this process via a so-called *occurrence function*.

► **Definition 2.** For a superstring  $S$ , a mask  $M$ , and a  $k$ -mer  $Q$ , the occurrence function  $\lambda(S, M, Q) \rightarrow \{0, 1\}^*$  is a function returning a finite binary sequence with the mask symbols of the corresponding occurrences, i.e.,

$$\lambda(S, M, Q) := (M_i \mid S_i \cdots S_{i+k-1} = Q) . \quad (1)$$

In this notation, verifying  $k$ -mer presence corresponds to evaluating the composite function ‘ $\text{or} \circ \lambda$ ’; that is, a  $k$ -mer is present if  $\lambda(S, M, Q)$  is non-empty and the logical **or** operation on these values yields 1. For instance, for  $\text{AcgGgg}$  from Example 1, and query  $k$ -mer  $Q = \text{GGG}$ , it holds that  $\lambda(S, M, Q) = (0, 1)$ , as the first occurrence is OFF and the second is ON, with the **or** of these values being 1; therefore,  $\text{GGG}$  is represented.

The set of all MS-represented  $k$ -mers can thus be expressed as

$$K = \{Q \in \Sigma^k \mid \text{or}(\lambda(S, M, Q)) = 1\} . \quad (2)$$

The key observation of this work is that **or** is only one member of a large class of possible “demasking functions” (see examples in Table 1). For instance, MS could have been defined using the **xor** function, with a  $k$ -mer considered present if the number of its ON occurrences is odd. Indeed,  $k$ -mer demasking can use any symmetric Boolean function, which we further equip with a special return value, **invalid**, as a means to impose criteria on mask validity.

► **Definition 3.** We call a symmetric function  $f : \{0, 1\}^* \rightarrow \{0, 1, \text{invalid}\}$  a  $k$ -mer demasking function.

In addition, we will typically require our demasking functions to have several natural properties. First, the non-appearing  $k$ -mers should be treated as non-represented (P1) (unless we aim to compactly represent set complements). Second, all appearing  $k$ -mers should be encodable using a compatible mask as present or absent, irrespective of the number of their occurrences (P2). Third, a single  $k$ -mer occurrence that is also ON should be interpreted as  $k$ -mer presence (P3). Fourth, the function should be efficiently computable from the frequencies of 0s and 1s (P4).

► **Definition 4.** A demasking function  $f$  is comprehensive if it satisfies the following four properties:

(P1)  $f(\epsilon) = 0$ .

(P2) For every  $n > 0$ , exist  $x, y \in \{0, 1\}^n$  such that  $f(x) = 0$  and  $f(y) = 1$ .

(P3)  $f((1)) = 1$  and  $f((0)) = 0$ .

(P4) Given  $|x|_0$  and  $|x|_1$ , one can evaluate  $f(x)$  in constant time in the wordRAM model.

With the notation of demasking functions  $f$ , we can now generalize the concept of MS to so-called  $f$ -masked superstrings ( $f$ -MS).

► **Definition 5.** Given a demasking function  $f$ , a superstring  $S$ , and a binary mask  $M$ , such that  $|M| = |S|$ , we call a triplet  $\mathcal{S} = (f, S, M)$  a function-assigned masked superstring or  $f$ -masked superstrings, and abbreviate it as  $f$ -MS.

If  $\exists Q \in \Sigma^k$  such that  $f(\lambda(S, M, Q)) = \text{invalid}$ , we call the  $f$ -MS invalid.

Function $f$	Definition	Compre- hensive	Use cases
<b>or</b>	1 if $ \lambda _1 > 0$ 0 if $ \lambda _1 = 0$	yes	<ul style="list-style-type: none"> <li>• MS (Sec. 3) and (r)SPSS (App. A)</li> <li>• <math>f</math> for union (Sec. 4)</li> </ul>
<b>xor</b>	1 if $ \lambda _1$ is odd 0 if $ \lambda _1$ is even	yes	<ul style="list-style-type: none"> <li>• <math>f</math> for sym. difference (Sec. 4)</li> </ul>
<b>and</b>	1 if $\lambda \neq \epsilon \wedge  \lambda _0 = 0$ 0 if $\lambda = \epsilon \vee  \lambda _0 > 0$	yes	<ul style="list-style-type: none"> <li>• allowing ON occurrences for ghost <math>k</math>-mers (App. C)</li> </ul>
<b><math>[a,b]</math>-threshold</b> ( $1 \leq a \leq b$ )	1 if $a \leq  \lambda _1 \leq b$ 0 otherwise	iff $a = 1$	<ul style="list-style-type: none"> <li>• <math>f_o</math> for intersection (Sec. 4)</li> <li>• <math>f_o</math> for set difference (Sec. 4)</li> </ul>
<b>one-or-nothing</b>	1 if $ \lambda _1 = 1$ 0 if $ \lambda _1 = 0$ <b>invalid</b> otherwise	yes	<ul style="list-style-type: none"> <li>• <math>f_1, f_2</math> for sym. difference (Sec. 4)</li> <li>• <math>f_1, f_2</math> for intersection (Sec. 4)</li> <li>• <math>f_1</math> for set difference (Sec. 4)</li> </ul>
<b>two-or-nothing</b>	1 if $ \lambda _1 = 2$ 0 if $ \lambda _1 = 0$ <b>invalid</b> otherwise	no	<ul style="list-style-type: none"> <li>• <math>f_2</math> for set difference (Sec. 4)</li> </ul>
<b>all-or-nothing</b>	1 if $\lambda \neq \epsilon \wedge  \lambda _0 = 0$ 0 if $ \lambda _1 = 0$ <b>invalid</b> otherwise	yes	<ul style="list-style-type: none"> <li>• omitting mask rank in membership queries (App. C)</li> </ul>

■ **Table 1 Common demasking functions used throughout the paper.**  $\lambda(f, S, M)$  is abbreviated as  $\lambda$ . All the mentioned non-comprehensive functions still satisfy properties (P1,4) from Definition 4.

209 Now, for a valid  $f$ -MS, we can generalize Equation (2) for  $k$ -mer decoding as

$$210 \quad K = \{Q \in \Sigma^k \mid f(\lambda(S, M, Q)) = 1\}. \quad (3)$$

211 As shown in [59], the most expensive part of MS computation (and thus also  $f$ -MS) is  
 212 finding a short superstring. It is therefore natural to ask whether we require a superstring  
 213 recomputation after performing a set operation. We now show that as long as function  $f$  in  
 214 the  $f$ -MS representation is comprehensive, the superstring may remain the same and only a  
 215 typically much simpler mask optimization may be needed.

216 ► **Lemma 6** (encoding). *Let  $f$  be a comprehensive demasking function,  $K$  be a  $k$ -mer set,*  
 217 *and  $S$  its arbitrary superstring. Then, there exists a mask  $M$  such that  $(f, S, M)$  is valid and*  
 218 *represents  $K$ .*

219 **Proof.** The lemma is a direct consequence of property (P2) in Definition 4. ◀

220 Similarly, we can change the demasking functions  $f$  to another comprehensive function  
 221 only by recoding the mask accordingly. For functions that are not comprehensive, such  
 222 recording may be impossible.

223 ► **Lemma 7** (recoding). *Let  $(f, S, M)$  be a valid  $f$ -MS representing a  $k$ -mer set  $K$ . Then,*  
 224 *for every comprehensive demasking function  $f'$ , there exists a valid mask  $M'$  such that*  
 225  *$(f', S, M')$  represents  $K$ .*

226 **Proof.** The existence of mask  $M'$  follows directly from Lemma 6 for the  $k$ -mer set  $K$ ,  
 227 superstring  $S$ , and function  $f'$ . ◀

228 ► **Example 8.** Consider Example 1 with the set of 3-mers  $K = \{\text{ACG}, \text{GGG}\}$  and the masked  
 229 superstring  $\text{AcGGgg}$ . For the query  $k$ -mer  $Q = \text{GGG}$ , the occurrence function for is  $\lambda(S, M, Q) =$   
 230  $(1, 1)$ . The result of demasking then  $f(\lambda(S, M, Q))$  then depends on the specific function  $f$ :  
 231 for **or** it evaluations to 1 and thus  $\text{GGG}$  is represented; however, for **xor** the result would be 0



232 and GGG would be a ghost  $k$ -mer, and thus not represented. As **xor** is comprehensive, we can  
 233 recode the  $f$ -MS from **or** to **xor** in a way it would still represent  $K$  by changing the mask,  
 234 for instance, to AcgGgg.

## 235 4 $f$ -MS as an Algebraic Framework

236 In this section, we describe on the conceptual level how to perform set operations on  $k$ -mer  
 237 sets by simply concatenating individual  $f$ -MS and choosing suitable demasking functions  
 238  $f$ . In Section 5, we deal with implementing this concept into data structures for  $k$ -mers,  
 239 specifically, into the FMSI index [60].

### 240 4.1 Concatenation as an Elementary Low-Level Operation

241 We define concatenation on  $f$ -MS as concatenating the underlying superstrings and masks  
 242 for all possible input and output functions  $f$ .

243 ► **Definition 9.** Given  $f$ -MS  $(f_1, S_1, M_1)$  and  $(f_2, S_2, M_2)$ , we define  $(f_1, f_2, f_o)$ -concatena-  
 244 tion as the operation taking these two  $f$ -MS and producing the result  $(f_o, S_1 + S_2, M_1 + M_2)$ .  
 245 We denote this operation by  $+_{f_1, f_2, f_o}$ .

246 Note that Definition 9 can be easily extended to more than two input  $f$ -MS. In the case  
 247 that all the functions are the same, i.e.  $f = f_1 = f_2 = f_o$ , we call it  $f$ -concatenation or just  
 248 concatenation if  $f$  is obvious from the context.

249 ► **Definition 10.** We call the set operations that can be performed with  $f_1 = f_2 = f_o$   
 250 function-preserving set operations. The operations that cannot be performed with a single  
 251 function are called function-transforming set operations.

252 Furthermore, note that while the set of appearing  $k$ -mers of  $S_1 + S_2$  clearly contains  
 253 the union of appearing  $k$ -mers of  $S_1$  and of  $S_2$ , additional new occurrences of  $k$ -mers may  
 254 appear at the boundary of the two superstrings. These newly appearing  $k$ -mers may not be  
 255 appearing in any of the superstrings  $S_1$  and  $S_2$ , and we refer to them as *boundary  $k$ -mers*.  
 256 The occurrences of appearing  $k$ -mers of  $S_1 + S_2$  that overlap both input superstrings are  
 257 called *boundary occurrences*.

### 258 4.2 Function-Preserving Set Operations

259 **Union.** As implicitly shown in [59], concatenating MS, which are **or**-MS in our notation,  
 260 acts as union on the represented sets, that is, the resulting represented set is the union  
 261 of the original represented sets. This allows **or**-MS to generalize (r)SPSS representations,  
 262 since any set of  $k$ -mers in the (r)SPSS representation can be directly viewed as an **or**-MS by  
 263 concatenating the individual simplitigs/matchtigs.

264 We show that **or** is the *only* comprehensive demasking function that acts as union on the  
 265 represented sets; see Appendix A for details. We further demonstrate this uniqueness even  
 266 on the level of matchtigs and therefore, **or**-MS are the only  $f$ -MS that generalize (r)SPSS  
 267 representations.

268 **Symmetric difference.** Next, we observe that **xor** naturally acts as the symmetric  
 269 difference set operation, i.e., concatenating two **xor**-MS results in a **xor**-MS representing  
 270 the symmetric difference of the original sets. Indeed, recall that using **xor** implies that a  
 271  $k$ -mer is represented if and only if there is a odd number of ON occurrences of that  $k$ -mer.



272 Observe that the boundary occurrences of *k*-mers do not affect the resulting represented  
 273 set as those have zeros in the mask. Thus, if a *k*-mer is present in both sets, it has an even  
 274 number of ON occurrences in total and hence, is not represented in the result. Likewise, if a  
 275 *k*-mer belongs to exactly one input set, it has an odd number of ON occurrences in this input  
 276 set and an even number (possibly zero) in the other; thus, it is represented in the result. As  
 277 any appearing *k*-mer is either boundary or appears in one of the MS, the result corresponds  
 278 to the symmetric difference.

### 279 4.3 Function-Transforming Set Operations

280 **Intersection.** After seeing functions for union and symmetric difference operations, it might  
 281 seem natural that there should be a function for intersection. This is however not the case  
 282 as there exists no comprehensive demasking function acting as intersection, which we show  
 283 in Appendix B.

284 We can circumvent the non-existence of a single demasking function acting as intersection  
 285 by using possibly non-comprehensive demasking functions that are different for the result  
 286 than for the input. We further show that such schemes have other applications beyond  
 287 intersection.

288 To this end, we will need two different types of demasking functions:

- 289 ■ **[*a*,*b*]-threshold** function (where  $0 < a \leq b$ ) is a demasking function that returns 1  
 290 whenever it receives an input of at least *a* ones and at most *b* ones and 0 otherwise. Note  
 291 that unless  $a = 1$ , **[*a*,*b*]-threshold** functions are not comprehensive as they do not satisfy  
 292 properties (P2) and (P3). The corresponding *f*-MS are denoted **[*a*,*b*]-threshold-MS**.
- 293 ■ The **one-or-nothing** function is a demasking function that returns 1 if there is exactly  
 294 one 1 in the input, 0 if there are no 1s, and **invalid** if there is more than a single ON  
 295 occurrence of the *k*-mer. Note that this function is comprehensive.

296 We now use these functions to perform any symmetric set operation on any number of  
 297 input *k*-mer sets. Given *N* sets of *k*-mers, we compute a **one-or-nothing-MS** for each. This  
 298 is always possible since **one-or-nothing** is a comprehensive demasking function and can be  
 299 done by directly using the superstrings and masks computed by KmerCamel [59].

300 We then concatenate the individual **one-or-nothing-MS**. The result is not a valid  
 301 **one-or-nothing-MS** in general, but it has the special property that each *k*-mer has as  
 302 many ON occurrences as the number of sets in which it appears. We can therefore change  
 303 the demasking function of the resulting *f*-MS from **one-or-nothing** to an **[*a*,*b*]-threshold**  
 304 function. This will result in an **[*a*,*b*]-threshold-MS** that is always valid and the represented  
 305 set will be exactly the *k*-mers that appear in at least *a* sets and at most *b* sets. Important  
 306 **[*a*,*b*]-threshold-MS** in this setting include the following:

- 307 ■ The **[*N*,*N*]-threshold-MS** corresponds to taking the intersection of the represented sets.
- 308 ■ The **[1,*N*]-threshold-MS** is the **or-MS** and corresponds to taking the union.
- 309 ■ The **[1,1]-threshold-MS** corresponds to taking those *k*-mers that appear in exactly one  
 310 of the original sets. In case of  $N = 2$ , this corresponds to the symmetric difference.

311 It is important to emphasize that we can use different **[*a*,*b*]-threshold** functions to alter  
 312 the resulting *k*-mer set without changing the superstring or the mask. For instance, we can  
 313 use the same superstring and mask to consider intersection and union simply by changing  
 314 the function from **[*N*,*N*]-threshold** to **[1,*N*]-threshold**.

315 **Arbitrary symmetric set operations.** The same scheme, with more general demasking  
 316 functions, can be used to implement *any symmetric set operation op* on any number of sets.

317 Indeed, given  $N$ , we again concatenate their **one-or-nothing**-MS in an arbitrary order. The  
 318 symmetry of **op** implies that there is a set  $S_N \subseteq \{0, 1, \dots, N\}$  such that a  $k$ -mer belongs  
 319 to the set resulting from applying **op** if and only if it is in  $a$  input sets for some  $a \in S_N$ .  
 320 The sets  $S_N$  for  $N = 1, 2, \dots$  can be directly transferred into a demasking function  $f_{\text{op}}$  that  
 321 models **op**; however,  $f_{\text{op}}$  may not satisfy the property (P4) from Definition 4, i.e., that we  
 322 can compute it in  $O(1)$  time.

323 **Set difference.** Having seen how to perform symmetric set operations, we deal with  
 324 asymmetric ones, focusing on the set difference of  $k$ -mer sets  $A \setminus B$ . Clearly, we cannot  
 325 use the same demasking function  $f$  to represent both  $A$  and  $B$  as it would be impossible to  
 326 distinguish the sets after concatenation. Hence, we use different functions to represent  $A$   
 327 and  $B$ , namely,

- 328 ■ represent  $A$  using a [1,1]-**threshold**-MS,
- 329 ■ represent  $B$  using a [2,2]-**threshold**-MS, and
- 330 ■ interpret the result as a [1,1]-**threshold**-MS.

331 This computes the difference correctly as all  $k$ -mers represented in  $B$  are treated as ghosts  
 332 in the result, the  $k$ -mers from  $A$  but not from  $B$  still have a single ON occurrence and thus  
 333 are correctly considered represented, and finally, the ghost  $k$ -mers in either of the initial sets  
 334 or the boundary  $k$ -mers have no influence on the result. The same functions can be used if  
 335 we subtract more than a single set. Furthermore, this scheme can be generalized to *any set*  
 336 *operations on any number of sets*, by representing the  $i$ -th input set with [i,i]-**threshold**-MS  
 337 and using a suitable demasking function for the result of the concatenation (constructed  
 338 similarly as  $f_{\text{op}}$  for symmetric operation **op** above).

339 The downside to this approach is that the [2,2]-**threshold** function is not comprehensive  
 340 and we cannot simply use any superstring of  $k$ -mers in  $B$ , but we need a superstring such  
 341 that every  $k$ -mer of  $B$  appears at least twice, which can for instance be achieved by doubling  
 342 the computed superstring of  $B$ . We remark that this is the best we can do as set difference  
 343 cannot be achieved with comprehensive functions solely as we show in Appendix B.

344 **Other applications.** Furthermore, there are many more demasking functions that can be  
 345 used with  $f$ -MS, although they may not correspond to set operations; we mention the **and**  
 346 and **all-or-nothing** demasking functions in Appendix C (see also Table 1).

## 347 5 Indexing $f$ -MS with the FMSI index

348 To support set operations on  $k$ -mer sets, while allowing fast  $k$ -mer queries, we utilize the  
 349 FMSI index [60] introduced in our concurrent work. We first give an overview of the FMSI  
 350 index [60] and then describe how to generalize the query algorithm in FMSI to support  
 351 arbitrary demasking functions  $f$ . Next, we show that performing set operations narrows  
 352 down to merging the Burrows-Wheeler Transform and provide algorithms for changing the  
 353 demasking function  $f$ . Last, if after several set operations the size of the index gets too large,  
 354 one can improve space efficiency by a *compaction*, i.e., recomputing the representation.

355 **FMSI index.** The FMSI index [60] for a  $k$ -mer set  $K$  is constructed from its masked  
 356 superstring  $(M, S)$  maximizing the number of ones, which can also be seen as **all-or-nothing**-  
 357 masked superstring (Appendix C). The FMSI index consists of the Burrows-Wheeler transform  
 358 (BWT) [16] of  $S$  with an associated rank data structure [31], and the *SA-transformed mask*  
 359  $M'$  which is a bit-vector of the same length as  $S$ , where  $M'[i] = M[j_i - 1 \bmod |M|]$  where  
 360  $j_i$  is the starting position of the lexicographically  $i$ -th suffix of  $S$ . Optionally, FMSI index  
 361 can also construct the *kLCP array* [54], where  $kLCP[i] = 1$  if the  $i$ -th and  $(i + 1)$ -th suffix

362 share a prefix of  $k - 1$  characters. FMSI can be constructed in linear time through Masked  
 363 BWT [60], a tailored variant of the classical BWT [16]. FMSI can index a  $k$ -mer  $Q$  in  $O(k)$   
 364 time by first computing the range of occurrences of  $Q$  in the suffix-array coordinates. Then,  
 365 since the bits of the SA-transformed mask in this range correspond to the mask symbols  
 366 of occurrences of  $Q$  and since the mask has the maximum number of ones, the presence or  
 367 absence of  $Q$  in the represented  $k$ -mer set can be determined from any of those bits [60].  
 368 If  $k$ LCP is used, streaming queries can be answered in  $O(1)$  time per  $k$ -mer [60]. Here, we  
 369 only consider the FMSI index without  $k$ LCP, which requires  $2 + o(1)$  bits of memory per  
 370 distinct  $k$ -mer under the spectrum-like property and at most  $3 + o(1)$  bits per superstring  
 371 character in the general case [60]. In addition, we consider the rank data structure also for  
 372 the SA-transformed mask, which does not asymptotically increase complexity, i.e., costs only  
 373 another  $o(1)$  bits per superstring character [31].

374 **Efficient queries with arbitrary demasking functions.** In [60, Algorithm 2] we describe  
 375 the  $O(k)$ -time query algorithm for **or**-MS with masks having the maximum number of ones,  
 376 which can be viewed as **all-or-nothing**-MS (Appendix C). In Lemma 11 we generalize the  
 377 result of [60, Lemma 4] to  $f$ -MS with comprehensive  $f$ , which can be directly translated to  
 378 an algorithm for querying general  $f$ -MS with the same time guarantees.

379 **► Lemma 11.** *Consider a query for  $k$ -mer  $Q$  on an  $f$ -MS  $(f, S, M)$  representing a  $k$ -mer  
 380 set  $K$  such that  $f$  is comprehensive. Let  $M'$  be the corresponding SA-transformed mask [60]  
 381 and let  $[i, j)$  be the range of sorted rotations of  $S$  starting with a  $k$ -mer  $Q$ . Then the presence  
 382 or absence of  $Q$  in  $K$  can be determined in  $O(1)$  time.*

383 **Proof.** From [60, Lemma 1],  $M'[x]$  for  $x \in [i, j)$  corresponds to the mask symbol of a  
 384 particular occurrence of  $Q$ . Therefore  $|\lambda(S, M, Q)|_1 = \text{rank}_1(M', j) - \text{rank}_1(M', i)$ , which  
 385 can be computed in  $O(1)$  time using two rank queries on the mask; here,  $\text{rank}_1(M', i) =$   
 386  $\sum_{a=0}^{i-1} M'[a]$  is the number of ones on coordinates  $0, \dots, i-1$  in  $M'$ , computed by the rank data  
 387 structure. Furthermore,  $|\lambda(S, M, Q)|_0 = |\lambda(S, M, Q)| - |\lambda(S, M, Q)|_1 = j - i - |\lambda(S, M, Q)|_1$ .  
 388 Since  $f$  is comprehensive and in particular, commutative,  $f(\lambda(S, M, Q))$  can be evaluated  
 389 from the two quantities in constant time. ◀

390 To query a  $k$ -mer we can then simply use backwards search to get the range of occurrences  
 391 of  $Q$  and then apply Lemma 11. The same adjustment works also for positive streaming  
 392 queries which only require  $O(1)$  time per query  $k$ -mer [60].

393 **Set operations in linear time in the index sizes.** In Section 4, we describe how to  
 394 implement set operations with  $f$ -MS via masked superstring concatenation. Performing  
 395 the operation on indexes boils down to merging two BWTs using any algorithm for BWT  
 396 merging, for example [30] which runs in linear time. To merge the SA-transformed masks,  
 397 we attach the mask symbols to the corresponding characters of BWT, as described in [60];  
 398 hence, the existing algorithms for BWT merging can be adjusted in a straightforward way to  
 399 merge the SA-transformed masks.

400  **$f$ -MS mask recasting for  $f$  transformation.** To change the demasking function  $f$  to  
 401 a different one without altering the represented  $k$ -mer set and the underlying superstring,  
 402 we may need to *recast* the mask. Although the recasting procedure depends on the specific  
 403 function  $f$  used, for all comprehensive functions mentioned in Table 1, recasting can be  
 404 done either by maximizing the number of 1s in the mask (**and** and **all-or-nothing**), or by  
 405 minimizing the number of 1s (all other functions in Table 1). If an  $f$ -MS is represented in the  
 406 original string form, this can be achieved in linear time using a single/two-pass algorithms,  
 407 respectively, as described in [59]. For  $f$ -MS in the suffix-array coordinates, we can export  
 408 the  $f$ -MS, then recast the mask, and index the result.

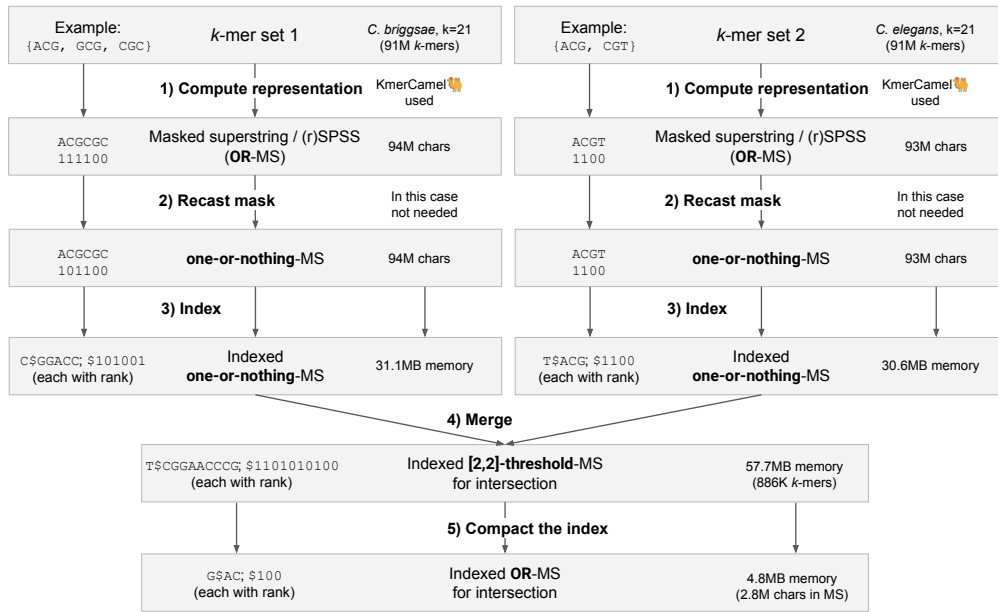
409 ***f*-MS compaction.** If an *f*-MS contains too many redundant copies of individual *k*-mers,  
410 e.g., if an *f*-MS is obtained by concatenating multiple input *f*-MS, it might be desirable to  
411 *compact* it, i.e., reoptimize its support superstring. This can be performed in time  $O(k|K|)$ ,  
412 by exporting the *f*-MS, counting the number of ON and OFF occurrences of each *k*-mer,  
413 constructing the represented *k*-mer set, and computing its **or**-MS [59]. We leave it to future  
414 work to design algorithms for compaction and for mask recasting directly in the FMSI index  
415 without the need to export the *f*-MS.

## 416 **6 Implementation and a Proof-of-Concept Experiment**

417 **Implementation in the FMSI tool.** We implemented this functionality as a proof-of-  
418 concept in the FMSI program [60], which is available under the MIT license on GitHub (<https://github.com/OndrejSladky/fmsi>) and distributed via BioConda. Merging of several  
419 indexes is currently implemented only by exporting individual superstrings, concatenating  
420 them as masked superstrings, and reindexing the result. Querying is a slight modification  
421 of the original implementation as described in [60], according to Lemma 11. We emphasize  
422 that our implementation is prototype, to demonstrate feasibility of performing set operations  
423 without the necessity to recompute a representation, and can be substantially optimized.  
424

425 **Proof-of-concept experiment.** We conducted a simple experiment with genomes of *C.*  
426 *elegans* (NC\_003279.8, 100M base pairs) and *C. briggsae* (NC\_013489.2, 108M base pairs),  
427 performing the union, intersection, and symmetric difference operations on their sets of  
428 *k*-mers. Our proposed pipeline for these operations, as depicted in Figure 1, consists of five  
429 steps: First, we compute a textual representation of the *k*-mer sets interpreted as **or**-MS. In  
430 our experiments, this was done using KmerCamel’s global greedy algorithm [59]. Second, we  
431 recast the mask to the desired demasking function, specifically we keep **or** for union and  
432 change to **one-or-nothing** for intersection and to **xor** for symmetric difference. In the case  
433 of **or**-MS computed by KmerCamel, mask recasting is actually not needed as the output  
434 already minimizes the number of 1s in the mask. Then we index the *f*-MS using FMSI [60]  
435 (this can be done even before mask recasting). The last two steps are concatenating the two  
436 indexes by index merging and compacting the resulting index if needed. Note that once  
437 indexed, we can ask membership queries on the resulting *k*-mer sets.

438 We evaluated the superstring length of each computed *f*-MS, and the memory requirements  
439 to perform queries on the indexed individual and concatenated *f*-MS, both before and after  
440 compaction. The results for intersection of the two roundworm genomes for  $k = 21$  are  
441 depicted in Figure 1. On the *k*-mer sets of both source genomes, which satisfy the spectrum-  
442 like property (SLP), the indexed *f*-MS required around 2.7 bits per distinct *k*-mer to perform  
443 queries. After taking the intersection, however, only about 1% of distinct *k*-mers remained,  
444 and the resulting *k*-mer set is far from SLP as the new MS representation after compaction  
445 requires three times more characters than the number of distinct *k*-mers. The memory  
446 requirements of FMSI per distinct *k*-mer are high as the resulting set has relatively few  
447 *k*-mers, less than one million, so latent memory to run FMSI was relatively significant. Note  
448 also that the fact that compaction significantly reduces the MS length highly depends on  
449 the particular use case, namely on the proportion of represented *k*-mers in the result. For  
450 union and symmetric difference for the same data, the compaction lead to negligible length  
451 reduction.



■ **Figure 1** Illustration of the *k*-mer set operations workflow using *f*-masked superstrings: an example of intersection of two sets using FMSI [60]. The workflow also contains an illustrative example on a set of 3-mers as well as experimental results on *C. briggsae* and *C. elegans* genomes with *k* = 21; namely, we show the number of masked superstring characters after each change, the number of distinct *k*-mers in each set, and for each indexed representation, the memory required to perform queries on the underlying set.

452 In contrast, CBL [46]<sup>1</sup> required about 950 MB of memory for processing queries for both  
 453 original genomes, which translates to about 83 bits per distinct *k*-mer. After computing the  
 454 index for the intersection (from the indexes of the source genomes), the memory dropped to  
 455 46.5 MB, still significantly larger than the memory of FMSI.

456 For results with symmetric difference, union, and other values of *k*, see the supplementary  
 457 repository (<https://github.com/OndrejSladky/f-masked-superstrings-supplement>).

## 458 7 Conclusion and Outlook

459 We have proposed *f*-masked superstrings (*f*-MS) as an algebraic data type for *k*-mer sets  
 460 that allows for seamless execution of set operations. It is primarily based on equipping  
 461 masked superstrings (MS) from [59] with a demasking function *f*, and we have thoroughly  
 462 investigated several natural demasking functions, demonstrating that set operations on  
 463 *k*-mer sets can be carried out simply by masked superstring concatenation or, if indexed, by  
 464 merging their masked Burrows-Wheeler transform from [60]. This leads to a simple data  
 465 structure that simultaneously allows for beyond worst-case compressibility, answering exact  
 466 membership queries, and efficiently performing set operations on the *k*-mer sets, without the

<sup>1</sup> Obtained from <https://github.com/imartayan/CBL>, version at commit '328bcc6'. The index was computed on canonical *k*-mers, to handle reverse complements, that is, we ran 'cbl build -c'. CBL was compiled using 'RUSTFLAGS="-C target-cpu=native" K={kmer-size} PREFIX\_BITS=28 \\  
 cargo +nightly build ---release ---examples ---target-dir target.k\_{kmer-size}'.

467 costly operation of recomputing the underlying representation. Another major advantage  
468 is the versatility of our concept as it can in fact be combined with (repetitive) Spectrum  
469 Preserving String Sets [12, 51, 58] instead of (more general) masked superstrings.

470 The main practical limitation of our work is the current implementation of the index merg-  
471 ing, which is very slow and not using the state-of-the-art algorithms for BWT merging [30].  
472 Furthermore, our proof-of-concept experiment is only meant to demonstrate feasibility, and  
473 we leave a more thorough evaluation, using various datasets and including a comparison to  
474 other tools for set operations, such as CBL [46], to future work.

475 Our work opens up several research directions for future theoretical investigation. Cur-  
476 rently, to deal with non-comprehensive demasking functions, which are necessary for asym-  
477 metric set operations, multiple copies of input masked superstrings are needed. This leads to  
478 a natural question whether this framework can be further generalized, possibly by extending  
479 the mask alphabet, to allow for more efficient asymmetric set operations. On the algorithmic  
480 level, our work relies on efficient algorithms for merging the Burrows-Wheeler transform, as  
481 well as merging the supporting data structures such as the  $k$ LCP array [54]. Moreover, it is  
482 open how to directly perform certain operations with  $f$ -masked superstrings indexed with the  
483 masked Burrows-Wheeler transform [60], namely, we seek algorithms for mask recasting and  
484 index compaction that do not necessitate to export the masked superstring, but rather work  
485 locally with the BWT of the superstring and the SA-transformed mask. Furthermore, as our  
486 work deals only with set operations, we open the question of performing single insertions  
487 and deletions in a more efficient way than performing these through set operations; note  
488 that for comprehensive demasking functions, deletions in the representation can be handled  
489 efficiently by changing the corresponding mask bits.

490 In conclusion, while the primary contributions of this paper are conceptual, they pave  
491 the path towards a space- and time-efficient library for  $k$ -mer sets that would include all of  
492 these features. In the light of advances in efficient superstring approximation algorithms  
493 and BWT-based indexing, we believe that the  $f$ -masked superstring framework is a useful  
494 step towards designing appropriate data structures for this library, which is now mainly an  
495 engineering challenge.

## 496 References

- 497 1. Jarno Alanko, Bahar Alipanahi, Jonathen Settle, Christina Boucher, and Travis Gagie.  
498 Buffering updates enables efficient dynamic de bruijn graphs. *Computational and*  
499 *structural biotechnology journal*, 19:4067–4078, 2021.
- 500 2. Jarno N Alanko, Simon J Puglisi, and Jaakko Vuhtoniemi. Small searchable  $\kappa$ -spectra  
501 via subset rank queries on the spectral Burrows-Wheeler transform. In *SIAM Conference*  
502 *on Applied and Computational Discrete Algorithms (ACDA23)*, pages 225–236. SIAM,  
503 2023. doi:10.1137/1.9781611977714.20.
- 504 3. Jarno N Alanko, Jaakko Vuhtoniemi, Tommi Mäklin, and Simon J Puglisi. Themisto: a  
505 scalable colored  $k$ -mer index for sensitive pseudoalignment against hundreds of thousands  
506 of bacterial genomes. *Bioinformatics*, 39(Supplement\_1):i260–i269, 2023. doi:10.1093/  
507 *bioinformatics/btad233*.
- 508 4. Bahar Alipanahi, Alan Kuhnle, Simon J Puglisi, Leena Salmela, and Christina Boucher.  
509 Succinct dynamic de bruijn graphs. *Bioinformatics*, 37(14):1946–1952, 2021.
- 510 5. Fatemeh Almodaresi, Jamshed Khan, Sergey Madaminov, Michael Ferdman, Rob John-  
511 son, Prashant Pandey, and Rob Patro. An incrementally updatable and scalable system



- 512 for large-scale sequence search using the Bentley-Saxe transformation. *Bioinformatics*,  
513 38(12):3155–3163, 2022. doi:10.1093/bioinformatics/btac142.
- 514 6. Timo Bingmann, Phelim Bradley, Florian Gauger, and Zamin Iqbal. COBS: a compact bit-  
515 sliced signature index. In *String Processing and Information Retrieval: 26th International*  
516 *Symposium, SPIRE 2019, Segovia, Spain, October 7–9, 2019, Proceedings 26*, pages  
517 285–303. Springer, 2019. doi:10.1007/978-3-030-32686-9\_21.
- 518 7. Alexander Bowe, Taku Onodera, Kunihiko Sadakane, and Tetsuo Shibuya. Succinct de  
519 Bruijn graphs. In Benjamin J. Raphael and Jijun Tang, editors, *Algorithms in Bioin-*  
520 *formatics - 12th International Workshop, WABI 2012, Ljubljana, Slovenia, September*  
521 *10-12, 2012. Proceedings*, volume 7534 of *Lecture Notes in Computer Science*, pages  
522 225–235. Springer, 2012. doi:10.1007/978-3-642-33122-0\_18.
- 523 8. Phelim Bradley, Henk C Den Bakker, Eduardo PC Rocha, Gil McVean, and Zamin Iqbal.  
524 Ultrafast search of all deposited bacterial and viral genomic data. *Nature Biotechnology*,  
525 37(2):152–159, 2019.
- 526 9. Phelim Bradley, N Claire Gordon, Timothy M Walker, Laura Dunn, Simon Heys, Bill  
527 Huang, Sarah Earle, Louise J Pankhurst, Luke Anson, Mariateresa De Cesare, et al.  
528 Rapid antibiotic-resistance predictions from genome sequence data for staphylococcus  
529 aureus and mycobacterium tuberculosis. *Nature Communications*, 6(1):10063, 2015.  
530 doi:10.1038/ncomms10063.
- 531 10. Nicolas L Bray, Harold Pimentel, Páll Melsted, and Lior Pachter. Near-optimal  
532 probabilistic RNA-seq quantification. *Nature Biotechnology*, 34(5):525–527, 2016.  
533 doi:10.1038/nbt.3519.
- 534 11. Karel Břinda. Novel computational techniques for mapping and classification of Next-  
535 Generation Sequencing data. PhD thesis, Université Paris-Est, 2016. doi:10.5281/  
536 zenodo.1045317.
- 537 12. Karel Břinda, Michael Baym, and Gregory Kucherov. Simplitigs as an efficient and  
538 scalable representation of de Bruijn graphs. *Genome Biology*, 22(96), 2021. doi:  
539 10.1186/s13059-021-02297-z.
- 540 13. Karel Břinda, Alanna Callendrello, Kevin C Ma, Derek R MacFadden, Themoula  
541 Charalampous, Robyn S Lee, Lauren Cowley, Crista B Wadsworth, Yonatan H Grad,  
542 Gregory Kucherov, et al. Rapid inference of antibiotic resistance and susceptibility by  
543 genomic neighbour typing. *Nature Microbiology*, 5(3):455–464, 2020. doi:10.1038/  
544 s41564-019-0656-6.
- 545 14. Karel Břinda, Leandro Lima, Simone Pignotti, Natalia Quinones-Olvera, Kamil Salikhov,  
546 Rayan Chikhi, Gregory Kucherov, Zamin Iqbal, and Michael Baym. Efficient and robust  
547 search of microbial genomes via phylogenetic compression. *bioRxiv*, 2023.04.15.536996,  
548 2023. doi:10.1101/2023.04.15.536996.
- 549 15. Karel Břinda, Kamil Salikhov, Simone Pignotti, and Gregory Kucherov. Prophyle  
550 0.3.1.0. *Zenodo*, 5281, 2017. URL: <https://prophyle.github.io>, doi:10.5281/  
551 zenodo.5237391.
- 552 16. Michael Burrows and David Wheeler. A block-sorting lossless data compression algorithm.  
553 Technical Report 124, Digital Equipment Corporation, 1994.
- 554 17. Rayan Chikhi. K-mer data structures in sequence bioinformatics. HDR thesis, Insti-  
555 tute Pasteur Ecole Doctorale “EDITE”, 2021. URL: [http://rayan.chikhi.name/pdf/  
556 RChikhi-HDR.pdf](http://rayan.chikhi.name/pdf/RChikhi-HDR.pdf).
- 557 18. Rayan Chikhi, Jan Holub, and Paul Medvedev. Data structures to represent a set  
558 of  $k$ -long DNA sequences. *ACM Computing Surveys*, 54(1):17:1–17:22, 2022. doi:  
559 10.1145/3445967.



- 560 19. Rayan Chikhi, Antoine Limasset, Shaun Jackman, Jared T. Simpson, and Paul Medvedev.  
561 On the representation of de Bruijn graphs. In Roded Sharan, editor, *Research in*  
562 *Computational Molecular Biology*, pages 35–55, Cham, 2014. Springer International  
563 Publishing. doi:10.1007/978-3-319-05269-4\_4.
- 564 20. Rayan Chikhi, Antoine Limasset, and Paul Medvedev. Compacting de Bruijn graphs  
565 from sequencing data quickly and in low memory. *Bioinformatics*, 32(12):i201–i208, 2016.  
566 doi:10.1093/bioinformatics/btw279.
- 567 21. Thomas C. Conway and Andrew J. Bromage. Succinct data structures for assembling large  
568 genomes. *Bioinformatics*, 27(4):479–486, 2011. doi:10.1093/bioinformatics/btq697.
- 569 22. Victoria G. Crawford, Alan Kuhnle, Christina Boucher, Rayan Chikhi, and Travis  
570 Gagie. Practical dynamic de Bruijn graphs. *Bioinformatics*, 34(24):4189–4195, 2018.  
571 doi:10.1093/bioinformatics/bty500.
- 572 23. Sebastian Deorowicz, Agnieszka Debudaj-Grabysz, and Szymon Grabowski. Disk-based  
573 k-mer counting on a pc. *BMC bioinformatics*, 14:1–12, 2013.
- 574 24. Jason Fan, Jamshed Khan, Giulio Ermanno Pibiri, and Rob Patro. Spectrum preserving  
575 tilings enable sparse and modular reference indexing. In Haixu Tang, editor, *Research in*  
576 *Computational Molecular Biology - 27th Annual International Conference, RECOMB*  
577 *2023, Istanbul, Turkey, April 16-19, 2023, Proceedings*, volume 13976 of *Lecture Notes in*  
578 *Computer Science*, pages 21–40. Springer, 2023. doi:10.1007/978-3-031-29119-7\\_2.
- 579 25. Jason Fan, Jamshed Khan, Noor Pratap Singh, Giulio Ermanno Pibiri, and Rob Patro.  
580 Fulgor: a fast and compact k-mer index for large-scale matching and color queries.  
581 *Algorithms for Molecular Biology*, 19(1):3, 2024. doi:10.1186/S13015-024-00251-9.
- 582 26. Paolo Ferragina and Giovanni Manzini. Opportunistic data structures with applications.  
583 In *Proceedings 41st Annual Symposium on Foundations of Computer Science, SFCS-00*.  
584 IEEE Comput. Soc, 2000. doi:10.1109/sfcs.2000.892127.
- 585 27. Gaurav Gupta, Minghao Yan, Benjamin Coleman, Bryce Kille, R. A. Leo Elworth,  
586 Tharun Medini, Todd Treangen, and Anshumali Shrivastava. Fast processing and  
587 querying of 170TB of genomics data via a Repeated And Merged BloOm filter (RAMBO).  
588 In *Proceedings of the 2021 International Conference on Management of Data, SIGMOD*  
589 *'21*, page 2226–2234, New York, NY, USA, 2021. Association for Computing Machinery.  
590 doi:10.1145/3448016.3457333.
- 591 28. Khodor Hannoush, Camille Marchet, and Pierre Peterlongo. Cdbgtricks: Strategies  
592 to update a compacted de bruijn graph. *bioRxiv*, 2024. URL: <https://www.biorxiv.org/content/early/2024/05/28/2024.05.24.595676>, arXiv:<https://www.biorxiv.org/content/early/2024/05/28/2024.05.24.595676.full.pdf>, doi:10.1101/2024.  
593 05.24.595676, doi:10.1101/2024.  
594 05.24.595676.  
595
- 596 29. Guillaume Holley and Páll Melsted. Bifrost: highly parallel construction and indexing  
597 of colored and compacted de Bruijn graphs. *Genome Biology*, 21(1):1–20, 2020. doi:  
598 10.1186/s13059-020-02135-8.
- 599 30. James Holt and Leonard McMillan. Merging of multi-string bwts with applications. *Bioin-*  
600 *formatics*, 30(24):3524–3531, August 2014. doi:10.1093/bioinformatics/btu584.
- 601 31. Guy Joseph Jacobson. *Succinct static data structures*. PhD thesis, Carnegie Mellon  
602 University, USA, 1988. AAI8918056.
- 603 32. Lauris Kaplinski, Maarja Lepamets, and Maido Remm. GenomeTester4: a toolkit for  
604 performing basic set operations - union, intersection and complement on k-mer lists.  
605 *GigaScience*, 4(1):s13742–015–0097–y, 12 2015. doi:10.1186/s13742-015-0097-y.

- 606 33. Mikhail Karasikov, Harun Mustafa, Daniel Danciu, Marc Zimmermann, Christopher  
607 Barber, Gunnar Rätsch, and André Kahles. Indexing all life’s known biological sequences.  
608 *bioRxiv*, 2020.10.01.322164, 2024. doi:10.1101/2020.10.01.322164.
- 609 34. Marek Kokot, Maciej Długosz, and Sebastian Deorowicz. KMC 3: counting and  
610 manipulating k-mer statistics. *Bioinformatics*, 33(17):2759–2761, 05 2017. doi:  
611 10.1093/bioinformatics/btx304.
- 612 35. Téo Lemane, Paul Medvedev, Rayan Chikhi, and Pierre Peterlongo. kmtricks: efficient and  
613 flexible construction of Bloom filters for large sequencing data collections. *Bioinformatics  
614 Advances*, 2(1):vbac029, 04 2022. doi:10.1093/bioadv/vbac029.
- 615 36. Heng Li. Exploring single-sample SNP and INDEL calling with whole-genome de novo  
616 assembly. *Bioinformatics*, 28(14):1838–1844, 2012. doi:10.1093/bioinformatics/  
617 bts280.
- 618 37. Heng Li. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
619 *arXiv preprint arXiv:1303.3997*, 2013. doi:10.48550/arXiv.1303.3997.
- 620 38. Heng Li and Richard Durbin. Fast and accurate short read alignment with  
621 Burrows-Wheeler transform. *Bioinformatics*, 25(14):1754–1760, 2009. doi:10.1093/  
622 bioinformatics/btp324.
- 623 39. Heng Li and Richard Durbin. Fast and accurate long-read alignment with  
624 Burrows-Wheeler transform. *Bioinformatics*, 26(5):589–595, 2010. doi:10.1093/  
625 bioinformatics/btp698.
- 626 40. Antoine Limasset, Guillaume Rizk, Rayan Chikhi, and Pierre Peterlongo. Fast and  
627 scalable minimal perfect hashing for massive key sets. In Costas S. Iliopoulos, Solon P.  
628 Pissis, Simon J. Puglisi, and Rajeev Raman, editors, *16th International Symposium  
629 on Experimental Algorithms, SEA 2017, June 21-23, 2017, London, UK*, volume 75  
630 of *LIPICs*, pages 25:1–25:16. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2017.  
631 doi:10.4230/lipics.sea.2017.25.
- 632 41. Abdullah-Al Mamun, Soumitra Pal, and Sanguthevar Rajasekaran. Kcmbt: ak-mer  
633 counter based on multiple burst trees. *Bioinformatics*, 32(18):2783–2790, 2016.
- 634 42. Guillaume Marçais and Carl Kingsford. A fast, lock-free approach for efficient parallel  
635 counting of occurrences of k-mers. *Bioinformatics*, 27(6):764–770, 2011.
- 636 43. Camille Marchet, Christina Boucher, Simon J Puglisi, Paul Medvedev, Mikaël Salson,  
637 and Rayan Chikhi. Data structures based on k-mers for querying large collections of  
638 sequencing data sets. *Genome Research*, 31(1):1–12, 2021. doi:10.1101/gr.260604.119.
- 639 44. Camille Marchet, Zamin Iqbal, Daniel Gautheret, Mikaël Salson, and Rayan Chikhi.  
640 REINDEER: efficient indexing of k-mer presence and abundance in sequencing datasets.  
641 *Bioinformatics*, 36(Supplement-1):i177–i185, 2020. doi:10.1093/bioinformatics/  
642 btaa487.
- 643 45. Camille Marchet, Maël Kerbirou, and Antoine Limasset. Blight: efficient exact asso-  
644 ciative structure for k-mers. *Bioinformatics*, 37(18):2858–2865, 2021. doi:10.1093/  
645 bioinformatics/btab217.
- 646 46. Igor Martayan, Bastien Cazaux, Antoine Limasset, and Camille Marchet. Con-  
647 way–bromage–lyndon (cbl): an exact, dynamic representation of k-mer sets.  
648 *Bioinformatics*, 40(Supplement\_1):i48–i57, 06 2024. arXiv:[https://academic.oup.com/bioinformatics/article-pdf/40/Supplement\\_1/i48/58354678/btae217.pdf](https://academic.oup.com/bioinformatics/article-pdf/40/Supplement_1/i48/58354678/btae217.pdf),  
649 doi:10.1093/bioinformatics/btae217.
- 650 47. Martin D. Muggli, Bahar Alipanahi, and Christina Boucher. Building large updatable  
651 colored de Bruijn graphs via merging. *Bioinformatics*, 35(14):i51–i60, 2019. doi:  
652 10.1093/bioinformatics/btz350.  
653

- 654 48. Rob Patro, Geet Duggal, Michael I Love, Rafael A Irizarry, and Carl Kingsford. Salmon  
655 provides fast and bias-aware quantification of transcript expression. *Nature Methods*,  
656 14(4):417–419, 2017. doi:10.1038/nmeth.4197.
- 657 49. Giulio Ermanno Pibiri. Sparse and skew hashing of K-mers. *Bioinformatics*, 38(Supple-  
658 ment\_1):i185–i194, 2022. doi:10.1093/bioinformatics/btac245.
- 659 50. Amatur Rahman. Compression algorithms for de Bruijn graphs and uncovering hidden  
660 assembly artifacts. PhD thesis, The Pennsylvania State University, 2023.
- 661 51. Amatur Rahman and Paul Medvedev. Representation of k-mer sets using spectrum-  
662 preserving string sets. *Journal of Computational Biology*, 28(4):381–394, 2021. PMID:  
663 33290137. doi:10.1089/cmb.2020.0431.
- 664 52. Arang Rhie, Brian P Walenz, Sergey Koren, and Adam M Phillippy. Merqury: reference-  
665 free quality, completeness, and phasing assessment for genome assemblies. *Genome*  
666 *biology*, 21:1–27, 2020.
- 667 53. Guillaume Rizk, Dominique Lavenier, and Rayan Chikhi. Dsk: k-mer counting with very  
668 low memory usage. *Bioinformatics*, 29(5):652–653, 2013.
- 669 54. Kamil Salikhov. Efficient algorithms and data structures for indexing dna sequence data.  
670 PhD thesis, Université Paris-Est, 2017.
- 671 55. Kamil Salikhov, Karel Břinda, Simone Pignotti, and Gregory Kucherov. ProPhex.  
672 <https://github.com/prophyle/prophex>, 2018. doi:10.5281/zenodo.1247432.
- 673 56. Sebastian Schmidt. Unitigs are not enough: the advantages of superunitig-based algo-  
674 rithms in bioinformatics. PhD thesis, University of Helsinki, 2023.
- 675 57. Sebastian Schmidt and Jarno N. Alanko. Eulertigs: minimum plain text representation of  
676 k-mer sets without repetitions in linear time. *Algorithms for Molecular Biology*, 18(1):5,  
677 2023. doi:10.1186/s13015-023-00227-1.
- 678 58. Sebastian Schmidt, Shahbaz Khan, Jarno N. Alanko, Giulio E. Pibiri, and Alexandru I.  
679 Tomescu. Matchtigs: minimum plain text representation of k-mer sets. *Genome Biology*,  
680 24(1):136, 2023. doi:10.1186/s13059-023-02968-z.
- 681 59. Ondřej Sladký, Pavel Veselý, and Karel Břinda. Masked superstrings as a unified  
682 framework for textual k-mer set representations. *bioRxiv*, 2023.02.01.526717, 2023.  
683 doi:10.1101/2023.02.01.526717.
- 684 60. Ondřej Sladký, Pavel Veselý, and Karel Břinda. FroM Superstring to Indexing: a space-  
685 efficient index for unconstrained k-mer sets using the masked burrows-wheeler transform  
686 (MBWT). *bioRxiv*, November 2024. URL: [http://dx.doi.org/10.1101/2024.10.30.](http://dx.doi.org/10.1101/2024.10.30.621029)  
687 621029, doi:10.1101/2024.10.30.621029.
- 688 61. Zachary D Stephens, Skylar Y Lee, Faraz Faghri, Roy H Campbell, Chengxiang Zhai,  
689 Miles J Efron, Ravishankar Iyer, Michael C Schatz, Saurabh Sinha, and Gene E Robinson.  
690 Big data: Astronomical or genetical? *PLoS Biology*, 13(7):e1002195, 2015. doi:  
691 10.1371/journal.pbio.1002195.
- 692 62. Derrick E Wood and Steven L Salzberg. Kraken: ultrafast metagenomic sequence  
693 classification using exact alignments. *Genome Biology*, 15(3):1–12, 2014. doi:10.1186/  
694 gb-2014-15-3-r46.

695  
696

## 697 **A** Proof of Uniqueness of Union Function

698 Masked superstrings [59], which we call **or-MS**, have the important property that concatenat-  
699 ing them results in the union of represented  $k$ -mers. This property makes it possible for MS

700 to generalize (r)SPSS representations [59] as unifying individual simplitigs/matchtigs results  
 701 in correctly representing the union of respective represented  $k$ -mer sets. In this section,  
 702 we show that **or**-MS are the only  $f$ -MS with any of these properties (acting as union and  
 703 generalizing (r)SPSS).

704 ► **Theorem 12.** *or is the only comprehensive demasking function  $f$  such that for any two*  
 705  *$k$ -mer sets  $K$  and  $K'$  and any of their valid  $f$ -MS  $(f, S, M)$  and  $(f, S', M')$ , respectively,*  
 706 *their concatenation  $(f, S + S', M + M')$  is a valid  $f$ -MS representing the set  $K \cup K'$ .*

707 **Proof.** For a contradiction assume there is a comprehensive demasking function  $f$  different  
 708 than **or** that satisfies the above. Consider the smallest  $n$  such that  $f$  behaves differently  
 709 than **or** for a length- $n$  input, meaning that there is  $x \in \{0, 1\}^n$  not equal to the all-zeros  
 710 vector such that  $f(x) \neq 1$ . As  $f$  is comprehensive, it cannot happen that  $f(x) = 1$  for all  
 711  $x \in \{0, 1\}^n$  by Definition 4, and moreover,  $f((1)) = 1$ , implying  $n > 1$ . Fix a  $k$ -mer, for  
 712 simplicity  $A^k$  (although similar approach works for all  $k$ -mers). We take the first  $f$ -MS to  
 713 be the  $k$ -mer with mask being  $M_0 = x_0$  and  $M_i = 0$  for the remaining  $k - 1$  positions  $i > 0$ .  
 714 The second  $f$ -MS is  $CA \dots A$  where  $A$  appears  $n + k - 2$  times with the mask being:  $M_0 = 0$ ,  
 715  $M_i = x_i$  for  $i = 1, \dots, k$ , and  $M_i = 0$  for  $i > k$ . At least one of the represented sets contains  
 716 the  $k$ -mer as  $x \neq 0$  and  $n$  is the but the resulting  $f$ -MS is either invalid or does not contain  
 717 the  $k$ -mer in the represented set as  $f(x) \neq 1$ , a contradiction. ◀

718 ► **Theorem 13.** *or is the only comprehensive demasking function  $f$  such that for any*  
 719 *sequence of  $f$ -MS, where individual superstrings are matchtigs, the concatenation of all the*  
 720  *$f$ -MS represents the union of represented  $k$ -mers.*

721 **Proof.** It is sufficient to find a construction of matchtigs such that we can construct an  
 722 arbitrary sequence of ones and zeros at the occurrences of a given  $k$ -mer and the rest follows  
 723 similarly as in Theorem 12.

724 We do this with  $k$ -mer  $CG$  and matchtigs  $Cg$  and  $Gc$ . Consider the counterexample sequence  
 725 of occurring ones and zeros from Theorem 12. For every one in the sequence, we add the  
 726 matchtig  $Cg$  and for each  $m$  consecutive zeros, we add  $m + 1$  times the matchtig  $Gc$ , since  
 727 at the boundary of two  $Gc$  matchtigs an OFF occurrence of  $k$ -mer  $CG$  appears. At any other  
 728 boundary, the  $k$ -mer  $CG$  does not appear, therefore the construction is correct. The rest of  
 729 the proof follows a similar argument as in Theorem 12. ◀

730 Note, however, that the same does not hold if we want to represent simplitigs/SPSS  
 731 solely. As an individual  $k$ -mer cannot appear more than once with an ON occurrence, any  
 732 comprehensive function generalizes SPSS representations if it satisfies that if there is one ON  
 733 occurrence of a  $k$ -mer, it returns 1, and if there is none, it returns 0.

## 734 **B Limits of Performing Set Operations using Comprehensive Functions**

735 In this section, we prove theoretical limitations of performing set operations by concatenation  
 736 of  $f$ -MS with comprehensive demasking functions. First, we show that intersection cannot  
 737 be a function-preserving set operation and second, we prove that it is impossible to use  
 738 only comprehensive demasking functions for input sets for the set difference operation. In  
 739 Section 4, we show how to overcome these limits via careful choice of demasking functions  
 740 that are not comprehensive.

## 741 B.1 Non-Existence of a Preserved Comprehensive Function for 742 Intersection.

743 We show that there is no comprehensive demasking function that acts as the intersection  
744 when concatenating  $f$ -MS. In a nutshell, this impossibility is caused by the fact that if  
745 there is a  $k$ -mer  $Q$  that occurs exactly once in the input MS with 1 in the mask, then after  
746 concatenation, it will still occur once with 1 in the mask, so under any comprehensive  $f$  the  
747  $k$ -mer would appear as if it was in the intersection.

748 ► **Theorem 14.** *There is no comprehensive demasking function  $f$  with the property that the*  
749 *result of  $f$ -concatenation of two  $f$ -MS always represents the intersection of the originally*  
750 *represented  $k$ -mer sets.*

751 **Proof.** Let  $f$  be any comprehensive demasking function. Consider MS A and C, each repre-  
752 senting a single 1-mer. Their concatenation is AC. Since  $f((1)) = 1$  by the comprehensiveness  
753 of  $f$ , the concatenation represents both 1-mer A and C. However, the intersection is empty  
754 and thus,  $f$  cannot be used to compute the intersection from the concatenation. ◀

755 Note that the proof cannot be generally extended to all demasking function as there exist  
756 non-comprehensive demasking functions acting as the intersection on the represented sets  
757 upon concatenation, for instance the constant zero function. However, since the constant  
758 zero function always represents the empty set, it is of no use in practice.

759 We further remark that although we have for convenience used the property (P3) from  
760 the definition of comprehensive functions, the proof in fact relies only on the property (P2)  
761 and holds even if we consider not only 1-mers, that is, a similar example can be constructed  
762 for any  $k$ .

## 763 B.2 Non-Existence of Comprehensive Input Functions for Set Difference.

764 We show that it is impossible to perform set difference of  $k$ -mer sets using  $f$ -MS with  
765 comprehensive functions solely.

766 ► **Theorem 15.** *There is no demasking function  $f_o$  and no comprehensive demasking*  
767 *functions  $f_1$  and  $f_2$ , such that the result of  $(f_1, f_2, f_o)$ -concatenation would always represent*  
768 *the set difference of the originally represented  $k$ -mer sets.*

769 **Proof.** Consider two masked superstrings and a  $k$ -mer  $Q$  which appears only once in each of  
770 the superstring. As an example, take the set of 2-mers  $K_1 = \{CA, CG\}$  and  $K_2 = \{AC\}$  and  
771 their (shortest) superstrings  $S_1 = CACG$  and  $S_2 = AC$ . Since the  $k$ -mer  $Q$  appears only once,  
772 from comprehensiveness, it must be represented as ON if it is present and as OFF otherwise,  
773 resulting that it is represented only based on its presence, the same in both MS. By symmetry  
774 of the output function, the result is the same if we concatenate  $S_1$  and  $S_2$  or the other way  
775 around. But it is not true in general that  $K_1 - K_2 = K_2 - K_1$  (consider the example above).  
776 This gives a contradiction and concludes the proof.

777 ◀

## 778 C Alternative Demasking Functions

779 In this section, we provide two other demasking functions that can be useful for some  
780 applications.

781 **C.1 The all-or-nothing-masked superstrings**

782 Perhaps the simplest approach to representing a set of  $k$ -mers is to mark all occurrences of  
783 represented  $k$ -mers with one, all ghost  $k$ -mers with zero, and treat all other masks as invalid.  
784 This corresponds to a function that returns 1 if it receives a list of ones, 0 if a list of zeros  
785 (or an empty list), and `invalid` otherwise. Alternatively, **all-or-nothing**-MS can be viewed  
786 as **or**-MS that maximize the number of ones in the mask.

787 This representation has been use in the FMSI index [60] as it allows to determine the  
788 presence of a  $k$ -mer from the mask symbol at an arbitrary occurrence of a particular  $k$ -mer,  
789 which enables the FMSI index to infer  $k$ -mer presence without any rank queries on the  
790 mask [60].

791 As another possible use case of **all-or-nothing**-MS, we could potentially achieve higher  
792 compressibility of the mask by realizing that we can infer the presence or absence of a  $k$ -mer  
793 from its first occurrence, which comes from the fact that a mask for a given set is unique.  
794 Thus, we can omit all symbols in the mask corresponding to any further occurrences of the  
795  $k$ -mer, making the mask shorter and easier to store, while it can be easily reconstructed  
796 afterwards.

797 **C.2 The and-MS**

798 We could easily replace the **or** function with **and**. That is, we could consider a  $k$ -mer present  
799 if it is marked as present at *all* its occurrences, with the small difference that we consider a  
800  $k$ -mer not represented if it does not appear, i.e., we consider the **and** of an empty binary  
801 string equal to 0, unlike in typical definitions of **and**. This ensures that the **and** function is  
802 comprehensive.

803 The potential advantage of **and**-MS over **or**-MS is that we can mark ghost  $k$ -mers with  
804 ones at some occurrences and therefore obtain masks with more ones in them, which could  
805 be beneficial, for instance, for additional improvements in mask compressibility.