

# String Sanitization: A Combinatorial Approach

Giulia Bernardini<sup>1</sup>, Huiping Chen<sup>2</sup>, Alessio Conte<sup>3,4</sup>, Roberto Grossi<sup>3,4</sup>,  
Grigorios Loukides<sup>2</sup> (✉), Nadia Pisanti<sup>3,4</sup>, Solon P. Pissis<sup>4,5</sup>, and Giovanna  
Rosone<sup>3\*</sup>

<sup>1</sup> Department of Informatics, Systems and Communication, University of  
Milano-Bicocca, Milan, Italy, [giulia.bernardini@unimib.it](mailto:giulia.bernardini@unimib.it)

<sup>2</sup> Department of Informatics, King's College London, London, UK  
[\[huiping.chen,grigorios.loukides\]@kcl.ac.uk](mailto:[huiping.chen,grigorios.loukides]@kcl.ac.uk)

<sup>3</sup> Department of Computer Science, University of Pisa, Pisa, Italy  
[\[conte,grossi,pisanti\]@di.unipi.it](mailto:[conte,grossi,pisanti]@di.unipi.it), [giovanna.rosone@unipi.it](mailto:giovanna.rosone@unipi.it)

<sup>4</sup> ERABLE Team, INRIA, Lyon, France

<sup>5</sup> CWI, Amsterdam, The Netherlands, [solon.pissis@cwi.nl](mailto:solon.pissis@cwi.nl)

**Abstract.** String data are often disseminated to support applications such as location-based service provision or DNA sequence analysis. This dissemination, however, may expose sensitive patterns that model confidential knowledge (*e.g.*, trips to mental health clinics from a string representing a user's location history). In this paper, we consider the problem of sanitizing a string by concealing the occurrences of sensitive patterns, while maintaining data utility. First, we propose a time-optimal algorithm, TFS-ALGO, to construct the shortest string preserving the order of appearance and the frequency of all non-sensitive patterns. Such a string allows accurately performing tasks based on the sequential nature and pattern frequencies of the string. Second, we propose a time-optimal algorithm, PFS-ALGO, which preserves a partial order of appearance of non-sensitive patterns but produces a much shorter string that can be analyzed more efficiently. The strings produced by either of these algorithms may reveal the location of sensitive patterns. In response, we propose a heuristic, MCSR-ALGO, which replaces letters in these strings with carefully selected letters, so that sensitive patterns are not reinstated and occurrences of spurious patterns are prevented. We implemented our sanitization approach that applies TFS-ALGO, PFS-ALGO and then MCSR-ALGO and experimentally show that it is effective and efficient.

---

\* The final authenticated publication is available online at [https://doi.org/10.1007/978-3-030-46150-8\\_37](https://doi.org/10.1007/978-3-030-46150-8_37). Thanks to Published source. Please, cite the publisher version: Giulia Bernardini, Huiping Chen, Alessio Conte, Roberto Grossi, Grigorios Loukides, Nadia Pisanti, Solon P. Pissis, Giovanna Rosone: String Sanitization: A Combinatorial Approach. Machine Learning and Knowledge Discovery in Databases. ECML PKDD 2019. Lecture Notes in Computer Science, vol 11906. Springer, Cham. [https://doi.org/10.1007/978-3-030-46150-8\\_37](https://doi.org/10.1007/978-3-030-46150-8_37).

## 1 Introduction

A large number of applications, in domains ranging from transportation to web analytics and bioinformatics feature data modeled as *strings*, *i.e.*, sequences of letters over some finite alphabet. For instance, a string may represent the history of visited locations of one or more individuals, with each letter corresponding to a location. Similarly, it may represent the history of search query terms of one or more web users, with letters corresponding to query terms, or a medically important part of the DNA sequence of a patient, with letters corresponding to DNA bases. Analyzing such strings is key in applications including location-based service provision, product recommendation, and DNA sequence analysis. Therefore, such strings are often disseminated beyond the party that has collected them. For example, location-based service providers often outsource their data to data analytics companies who perform tasks such as similarity evaluation between strings [?], and retailers outsource their data to marketing agencies who perform tasks such as mining frequent patterns from the strings [?].

However, disseminating a string intact may result in the exposure of confidential knowledge, such as trips to mental health clinics in transportation data [?], query terms revealing political beliefs or sexual orientation of individuals in web data [?], or diseases associated with certain parts of DNA data [?]. Thus, it may be necessary to sanitize a string prior to its dissemination, so that confidential knowledge is not exposed. At the same time, it is important to preserve the utility of the sanitized string, so that data protection does not outweigh the benefits of disseminating the string to the party that disseminates or analyzes the string, or to the society at large. For example, a retailer should still be able to obtain actionable knowledge in the form of frequent patterns from the marketing agency who analyzed their outsourced data; and researchers should still be able to perform analyses such as identifying significant patterns in DNA sequences.

**Our Model and Setting.** Motivated by the discussion above, we introduce the following model which we call *Combinatorial String Dissemination (CSD)*. In CSD, a party has a string  $W$  that it seeks to disseminate, while satisfying a set of *constraints* and a set of desirable *properties*. For instance, the constraints aim to capture privacy requirements and the properties aim to capture data utility considerations (*e.g.*, posed by some other party based on applications). To satisfy both,  $W$  must be transformed to a string  $X$  by applying a sequence of edit operations. The computational task is to determine this sequence of edit operations so that  $X$  satisfies the desirable properties subject to the constraints.

Under the CSD model, we consider a specific setting in which the sanitized string  $X$  must satisfy the following constraint **C1**: for an integer  $k > 0$ , no given length- $k$  substring (also called pattern) modeling confidential knowledge should occur in  $X$ . We call each such length- $k$  substring a *sensitive pattern*. We aim at finding the shortest possible string  $X$  satisfying the following desired properties: **(P1)** the order of appearance of all other length- $k$  substrings (*non-sensitive patterns*) is the same in  $W$  and in  $X$ ; and **(P2)** the frequency of these length- $k$  substrings is the same in  $W$  and in  $X$ . The problem of constructing  $X$  in this setting is referred to as TFS (Total order, Frequency, Sanitization). Clearly,

substrings of arbitrary lengths can be hidden from  $X$  by setting  $k$  equal to the length of the shortest substring we wish to hide, and then setting, for each of these substrings, any length- $k$  substring as sensitive.

Our setting is motivated by real-world applications involving string dissemination. In these applications, a *data custodian* disseminates the sanitized version  $X$  of a string  $W$  to a *data recipient*, for the purpose of analysis (*e.g.*, mining).  $W$  contains confidential information that the data custodian needs to hide, so that it does not occur in  $X$ . Such information is specified by the data custodian based on domain expertise, as in  $[?, ?, ?, ?]$ . At the same time, the data recipient specifies **P1** and **P2** that  $X$  must satisfy in order to be useful. These properties map directly to common data utility considerations in string analysis. By satisfying **P1**,  $X$  allows tasks based on the sequential nature of the string, such as blockwise  $q$ -gram distance computation [?], to be performed accurately. By satisfying **P2**,  $X$  allows computing the frequency of length- $k$  substrings [?] and hence mining frequent length- $k$  substrings with no utility loss. We require that  $X$  has minimal length so that it does not contain redundant information. For instance, the string which is constructed by concatenating all non-sensitive length- $k$  substrings in  $W$  and separating them with a special letter that does not occur in  $W$ , satisfies **P1** and **P2** but is not the shortest possible. Such a string  $X$  will have a negative impact on the efficiency of any subsequent analysis tasks to be performed on it.

Note, existing works for sequential data sanitization (*e.g.*, [?, ?, ?, ?, ?]) or anonymization (*e.g.*, [?, ?, ?]) cannot be applied to our setting (see Section 7).

**Our Contributions.** We define the TFS problem for string sanitization and a variant of it, referred to as PFS (Partial order, Frequency, Sanitization), which aims at producing an even shorter string  $Y$  by relaxing **P1** of TFS. Our algorithms for TFS and PFS construct strings  $X$  and  $Y$  using a separator letter  $\#$ , which is not contained in the alphabet of  $W$ . This prevents occurrences of sensitive patterns in  $X$  or  $Y$ . The algorithms repeat proper substrings of sensitive patterns so that the frequency of non-sensitive patterns overlapping with sensitive ones does not change. For  $X$ , we give a deterministic construction which may be easily reversible (*i.e.*, it may enable a data recipient to construct  $W$  from  $X$ ), because the occurrences of  $\#$  reveal the exact location of sensitive patterns. For  $Y$ , we give a construction which breaks several ties arbitrarily, thus being less easily reversible. We further address the reversibility issue by defining the MCSR (Minimum-Cost Separators Replacement) problem and designing an algorithm for dealing with it. In MCSR, we seek to replace all separators, so that the location of sensitive patterns is not revealed, while preserving data utility. We make the following specific contributions:

1. We design an algorithm for solving the TFS problem in  $\mathcal{O}(kn)$  time, where  $n$  is the length of  $W$ . In fact we prove that  $\mathcal{O}(kn)$  time is worst-case optimal by showing that the length of  $X$  is in  $\Theta(kn)$  in the worst case. The output of the algorithm is a string  $X$  consisting of a sequence of substrings over the alphabet of  $W$  separated by  $\#$  (see Example 1 below). An important feature of our algorithm, which is useful in the efficient construction of  $Y$  discussed next, is that it can be

implemented to produce an  $\mathcal{O}(n)$ -sized representation of  $X$  with respect to  $W$  in  $\mathcal{O}(n)$  time. See Section 3.

*Example 1.* Let  $W = \text{aabaaaababbbaab}$ ,  $k = 4$ , and the set of sensitive patterns be  $\{\text{aaaa}, \text{baaa}, \text{bbaa}\}$ . The string  $X = \text{aabaa\#aaababbba\#baab}$  consists of three substrings over the alphabet  $\{\text{a}, \text{b}\}$  separated by  $\#$ . Note that no sensitive pattern occurs in  $X$ , while all non-sensitive substrings of length 4 have the same frequency in  $W$  and in  $X$  (e.g.,  $\text{aaba}$  appears twice) and appear in the same order in  $W$  and in  $X$  (e.g.,  $\text{babb}$  precedes  $\text{abbb}$ ). Also, note that any shorter string than  $X$  would either create sensitive patterns or change the frequencies (e.g., removing the last letter of  $X$  creates a string in which  $\text{baab}$  no longer appears).  $\square$

**2.** We define the PFS problem relaxing **P1** of TFS to produce shorter strings that are more efficient to analyze. Instead of a *total order* (**P1**), we require a *partial order* (**PI1**) that preserves the order of appearance only for sequences of successive non-sensitive length- $k$  substrings that overlap by  $k - 1$  letters. This makes sense because the order of two successive non-sensitive length- $k$  substrings with no length- $(k - 1)$  overlap has anyway been “interrupted” (by a sensitive pattern). We exploit this observation to shorten the string further. Specifically, we design an algorithm that solves PFS in the optimal  $\mathcal{O}(n + |Y|)$  time, where  $|Y|$  is the length of  $Y$ , using the  $\mathcal{O}(n)$ -sized representation of  $X$ . See Section 4.

*Example 2. (Cont’d from Example 1)* Recall that  $W = \text{aabaaaababbbaab}$ . A string  $Y$  is  $\text{aaababbba\#aabaab}$ . The order of  $\text{babb}$  and  $\text{abbb}$  is preserved in  $Y$  since they are successive, non-sensitive, and with an overlap of  $k - 1 = 3$  letters. The order of  $\text{abaa}$  and  $\text{aaab}$ , which are successive and non-sensitive, is not preserved since they do not have an overlap of  $k - 1 = 3$  letters.  $\square$

**3.** We define the MCSR problem, which seeks to produce a string  $Z$ , by deleting or replacing all separators in  $Y$  with letters from the alphabet of  $W$  so that: no sensitive patterns are reinstated in  $Z$ ; occurrences of spurious patterns that may not be mined from  $W$  but can be mined from  $Z$ , for a given support threshold, are prevented; the distortion incurred by the replacements in  $Z$  is bounded. The first requirement is to preserve privacy and the next two to preserve data utility. We show that MCSR is NP-hard and propose a heuristic to attack it. See Section 5.

**4.** We implemented our combinatorial approach for sanitizing a string  $W$  (i.e., all aforementioned algorithms implementing the pipeline  $W \rightarrow X \rightarrow Y \rightarrow Z$ ) and show its effectiveness and efficiency on real and synthetic data. See Section 6.

## 2 Preliminaries, Problem Statements, and Main Results

**Preliminaries.** Let  $T = T[0]T[1] \dots T[n - 1]$  be a *string* of length  $|T| = n$  over a finite ordered alphabet  $\Sigma$  of size  $|\Sigma| = \sigma$ . By  $\Sigma^*$  we denote the set of all strings over  $\Sigma$ . By  $\Sigma^k$  we denote the set of all length- $k$  strings over  $\Sigma$ . For two positions  $i$  and  $j$  on  $T$ , we denote by  $T[i..j] = T[i] \dots T[j]$  the *substring* of  $T$  that starts at position  $i$  and ends at position  $j$  of  $T$ . By  $\varepsilon$  we denote the *empty string* of length 0. A *prefix* of  $T$  is a substring of the form  $T[0..j]$ , and a *suffix* of  $T$  is a substring of the form  $T[i..n - 1]$ . A *proper prefix* (suffix) of a string is not

equal to the string itself. By  $\text{Freq}_V(U)$  we denote the number of occurrences of string  $U$  in string  $V$ . Given two strings  $U$  and  $V$  we say that  $U$  has a *suffix-prefix overlap* of length  $\ell > 0$  with  $V$  if and only if the length- $\ell$  suffix of  $U$  is equal to the length- $\ell$  prefix of  $V$ , i.e.,  $U[|U| - \ell .. |U| - 1] = V[0 .. \ell - 1]$ .

We fix a string  $W$  of length  $n$  over an alphabet  $\Sigma = \{1, \dots, n^{\mathcal{O}(1)}\}$  and an integer  $0 < k < n$ . We refer to a length- $k$  string or a *pattern* interchangeably. An occurrence of a pattern is uniquely represented by its starting position. Let  $\mathcal{S}$  be a set of positions over  $\{0, \dots, n - k\}$  with the following closure property: for every  $i \in \mathcal{S}$ , if there exists  $j$  such that  $W[j .. j + k - 1] = W[i .. i + k - 1]$ , then  $j \in \mathcal{S}$ . That is, if an occurrence of a pattern is in  $\mathcal{S}$  all its occurrences are in  $\mathcal{S}$ . A substring  $W[i .. i + k - 1]$  of  $W$  is called *sensitive* if and only if  $i \in \mathcal{S}$ .  $\mathcal{S}$  is thus the set of occurrences of sensitive patterns. The difference set  $\mathcal{I} = \{0, \dots, n - k\} \setminus \mathcal{S}$  is the set of occurrences of *non-sensitive* patterns.

For any string  $U$ , we denote by  $\mathcal{I}_U$  the set of occurrences of non-sensitive length- $k$  strings over  $\Sigma$ . (We have that  $\mathcal{I}_W = \mathcal{I}$ .) We call an occurrence  $i$  the *t-predecessor* of another occurrence  $j$  in  $\mathcal{I}_U$  if and only if  $i$  is the largest element in  $\mathcal{I}_U$  that is less than  $j$ . This relation induces a *strict total order* on the occurrences in  $\mathcal{I}_U$ . We call  $i$  the *p-predecessor* of  $j$  in  $\mathcal{I}_U$  if and only if  $i$  is the t-predecessor of  $j$  in  $\mathcal{I}_U$  and  $U[i .. i + k - 1]$  has a suffix-prefix overlap of length  $k - 1$  with  $U[j .. j + k - 1]$ . This relation induces a *strict partial order* on the occurrences in  $\mathcal{I}_U$ . We call a subset  $\mathcal{J}$  of  $\mathcal{I}_U$  a *t-chain* (resp., *p-chain*) if for all elements in  $\mathcal{J}$  except the minimum one, their t-predecessor (resp., p-predecessor) is also in  $\mathcal{J}$ . For two strings  $U$  and  $V$ , chains  $\mathcal{J}_U$  and  $\mathcal{J}_V$  are *equivalent*, denoted by  $\mathcal{J}_U \equiv \mathcal{J}_V$ , if and only if  $|\mathcal{J}_U| = |\mathcal{J}_V|$  and  $U[u .. u + k - 1] = V[v .. v + k - 1]$ , where  $u$  is the  $j$ th smallest element of  $\mathcal{J}_U$  and  $v$  is the  $j$ th smallest of  $\mathcal{J}_V$ , for all  $j \leq |\mathcal{J}_U|$ .

### Problem Statements and Main Results.

**Problem 1 (TFS).** *Given  $W$ ,  $k$ ,  $\mathcal{S}$ , and  $\mathcal{I}$  construct the shortest string  $X$ :*

**C1**  *$X$  does not contain any sensitive pattern.*

**P1**  *$\mathcal{I}_W \equiv \mathcal{I}_X$ , i.e., the t-chains  $\mathcal{I}_W$  and  $\mathcal{I}_X$  are equivalent.*

**P2**  *$\text{Freq}_X(U) = \text{Freq}_W(U)$ , for all  $U \in \Sigma^k \setminus \{W[i .. i + k - 1] : i \in \mathcal{S}\}$ .*

TFS requires constructing the shortest string  $X$  in which all sensitive patterns from  $W$  are concealed (**C1**), while preserving the order (**P1**) and the frequency (**P2**) of all non-sensitive patterns. Our first result is the following.

**Theorem 1.** *Let  $W$  be a string of length  $n$  over  $\Sigma = \{1, \dots, n^{\mathcal{O}(1)}\}$ . Given  $k < n$  and  $\mathcal{S}$ , TFS-ALGO solves Problem 1 in  $\mathcal{O}(kn)$  time, which is worst-case optimal. An  $\mathcal{O}(n)$ -sized representation of  $X$  can be built in  $\mathcal{O}(n)$  time.*

**P1** implies **P2**, but **P1** is a strong assumption that may result in long output strings that are inefficient to analyze. We thus relax **P1** to require that the order of appearance remains the same only for sequences of successive non-sensitive length- $k$  substrings that also overlap by  $k - 1$  letters (p-chains).

**Problem 2 (PFS).** *Given  $W$ ,  $k$ ,  $\mathcal{S}$ , and  $\mathcal{I}$  construct a shortest string  $Y$ :*

**C1**  *$Y$  does not contain any sensitive pattern.*

**II1** There exists an injective function  $f$  from the  $p$ -chains of  $\mathcal{I}_W$  to the  $p$ -chains of  $\mathcal{I}_Y$  such that  $f(\mathcal{J}_W) \equiv \mathcal{J}_Y$  for any  $p$ -chain  $\mathcal{J}_W$  of  $\mathcal{I}_W$ .

**P2**  $\text{Freq}_Y(U) = \text{Freq}_W(U)$ , for all  $U \in \Sigma^k \setminus \{W[i..i+k-1] : i \in \mathcal{S}\}$ .

Our second result, which builds on Theorem 1, is the following.

**Theorem 2.** Let  $W$  be a string of length  $n$  over  $\Sigma = \{1, \dots, n^{\mathcal{O}(1)}\}$ . Given  $k < n$  and  $\mathcal{S}$ , PFS-ALGO solves Problem 2 in the optimal  $\mathcal{O}(n + |Y|)$  time.

To arrive at Theorems 1 and 2, we use a special letter (separator)  $\# \notin \Sigma$  when required. However, the occurrences of  $\#$  may reveal the locations of sensitive patterns. We thus seek to delete or replace the occurrences of  $\#$  in  $Y$  with letters from  $\Sigma$ . The new string  $Z$  should not reinstate any sensitive pattern. Given an integer threshold  $\tau > 0$ , we call pattern  $U \in \Sigma^k$  a  $\tau$ -ghost in  $Z$  if and only if  $\text{Freq}_W(U) < \tau$  but  $\text{Freq}_Z(U) \geq \tau$ . Moreover, we seek to prevent  $\tau$ -ghost occurrences in  $Z$  by also bounding the total *weight* of the *letter choices* we make to replace the occurrences of  $\#$ . This is the MCSR problem. We show that already a restricted version of the MCSR problem, namely, the version when  $k = 1$ , is NP-hard via the *Multiple Choice Knapsack* (MCK) problem [?].

**Theorem 3.** The MCSR problem is NP-hard.

Based on this connection, we propose a non-trivial heuristic algorithm to attack the MCSR problem for the general case of an arbitrary  $k$ .

### 3 TFS-ALGO

We convert string  $W$  into a string  $X$  over alphabet  $\Sigma \cup \{\#\}$ ,  $\# \notin \Sigma$ , by reading the letters of  $W$ , from left to right, and appending them to  $X$  while enforcing the following two rules:

**R1:** When the last letter of a sensitive substring  $U$  is read from  $W$ , we append  $\#$  to  $X$  (essentially replacing this last letter of  $U$  with  $\#$ ). Then, we append the succeeding non-sensitive substring (in the t-predecessor order) after  $\#$ .

**R2:** When the  $k - 1$  letters before  $\#$  are the same as the  $k - 1$  letters after  $\#$ , we remove  $\#$  and the  $k - 1$  succeeding letters (inspect Fig. 1).

**Fig. 1:** Sensitive patterns are overlined in red; non-sensitive are under- or over-lined in blue;  $\tilde{X}$  is obtained by applying **R1**; and  $X$  by applying **R1** and **R2**. In green we highlight an overlap of  $k - 1 = 3$  letters. Note that substring aaaababbb, whose length is greater than  $k$ , is also not occurring in  $X$ .

$W = \underline{\text{aabaaa}}\underline{\text{ababbb}}\underline{\text{aab}}$

$\tilde{X} = \underline{\text{aabaaa}}\underline{\#}\underline{\text{aaaba}}\underline{\#}\underline{\text{babb}}\underline{\#}\underline{\text{bbbaab}}$

$X = \underline{\text{aabaaa}}\underline{\text{ba}}\underline{\#}\underline{\text{babb}}\underline{\#}\underline{\text{bbbaab}}$

**R1** prevents  $U$  from occurring in  $X$ , and **R2** reduces the length of  $X$  (*i.e.*, allows to protect sensitive patterns with fewer extra letters). Both rules leave unchanged the order and frequencies of non-sensitive patterns. It is crucial to observe that applying the idea behind **R2** on more than  $k - 1$  letters would decrease the frequency of some pattern, while applying it on fewer than  $k - 1$  letters would create new patterns. Thus, we need to consider just **R2** *as-is*.

Let  $C$  be an array of size  $n$  that stores the occurrences of sensitive and non-sensitive patterns:  $C[i] = 1$  if  $i \in \mathcal{S}$  and  $C[i] = 0$  if  $i \in \mathcal{I}$ . For technical reasons we set the last  $k-1$  values in  $C$  equal to  $C[n-k]$ ; i.e.,  $C[n-k+1] := \dots := C[n-1] := C[n-k]$ . Note that  $C$  is constructible from  $\mathcal{S}$  in  $\mathcal{O}(n)$  time. Given  $C$  and  $k < n$ , TFS-ALGO efficiently constructs  $X$  by implementing **R1** and **R2** concurrently as opposed to implementing **R1** and then **R2** (see the proof of Lemma 1 for details of the workings of TFS-ALGO and Fig. 1 for an example). We next show that string  $X$  enjoys several properties.

**Lemma 1.** *Let  $W$  be a string of length  $n$  over  $\Sigma$ . Given  $k < n$  and array  $C$ , TFS-ALGO constructs the shortest string  $X$  such that the following hold:*

1. *There exists no  $W[i..i+k-1]$  with  $C[i] = 1$  occurring in  $X$  (**C1**).*
2.  *$\mathcal{I}_W \equiv \mathcal{I}_X$ , i.e., the order of substrings  $W[i..i+k-1]$ , for all  $i$  such that  $C[i] = 0$ , is the same in  $W$  and in  $X$ ; conversely, the order of all substrings  $U \in \Sigma^k$  of  $X$  is the same in  $X$  and in  $W$  (**P1**).*
3.  *$\text{Freq}_X(U) = \text{Freq}_W(U)$ , for all  $U \in \Sigma^k \setminus \{W[i..i+k-1] : C[i] = 1\}$  (**P2**).*
4. *The occurrences of letter  $\#$  in  $X$  are at most  $\lfloor \frac{n-k+1}{2} \rfloor$  and they are at least  $k$  positions apart (**P3**).*
5.  *$0 \leq |X| \leq \lceil \frac{n-k+1}{2} \rceil \cdot k + \lfloor \frac{n-k+1}{2} \rfloor$  and these bounds are tight (**P4**).*

---

TFS-ALGO( $W \in \Sigma^n, C, k, \# \notin \Sigma$ )

---

```

1   $X \leftarrow \varepsilon; j \leftarrow |W|; \ell \leftarrow 0;$ 
2   $j \leftarrow \min\{i \mid C[i] = 0\};$  /*  $j$  is the leftmost pos of a non-sens. pattern */
3  if  $j + k - 1 < |W|$  then /* Append the first non-sens. pattern to  $X$  */
4       $X[0..k-1] \leftarrow W[j..j+k-1]; j \leftarrow j+k; \ell \leftarrow \ell+k;$ 
5  while  $j < |W|$  do /* Examine two consecutive patterns */
6       $p \leftarrow j-k; c \leftarrow p+1;$ 
7      if  $C[p] = C[c] = 0$  then /* If both are non-sens., append the last
8          letter of the rightmost one to  $X$  */
9           $X[\ell] \leftarrow W[j]; \ell \leftarrow \ell+1; j \leftarrow j+1;$ 
10         if  $C[p] = 0 \wedge C[c] = 1$  then /* If the rightmost is sens., mark it
11             and advance  $j$  */
12              $f \leftarrow c; j \leftarrow j+1;$ 
13         if  $C[p] = C[c] = 1$  then  $j \leftarrow j+1;$  /* If both are sens., advance  $j$  */
14         if  $C[p] = 1 \wedge C[c] = 0$  then /* If the leftmost is sens. and the
15             rightmost is not */
16             if  $W[c..c+k-2] = W[f..f+k-2]$  then /* If the last marked
17                 sens. pattern and the current non-sens. overlap by  $k-1$ ,
18                 append the last letter of the latter to  $X$  */
19                  $X[\ell] \leftarrow W[j]; \ell \leftarrow \ell+1; j \leftarrow j+1;$ 
20             else /* Else append  $\#$  and the current non-sensitive pattern
21                 to  $X$  */
22                  $X[\ell] \leftarrow \#; \ell \leftarrow \ell+1;$ 
23                  $X[\ell.. \ell+k-1] \leftarrow W[j-k+1..j]; \ell \leftarrow \ell+k; j \leftarrow j+1;$ 
24 report  $X$ 

```

---

*Proof.* Proofs of **C1** and **P1-P4** can be found in [?]. We prove here that  $X$  has minimal length. Let  $X_j$  be the prefix of string  $X$  obtained by processing  $W[0..j]$ . Let  $j_{\min} = \min\{i \mid C[i] = 0\} + k - 1$ . We will proceed by induction on  $j$ , claiming that  $X_j$  is the shortest string such that **C1** and **P1-P4** hold for  $W[0..j]$ ,  $\forall j_{\min} \leq j \leq |W| - 1$ . We call such a string *optimal*.

*Base case:*  $j = j_{\min}$ . By Lines 3-4 of TFS-ALGO,  $X_j$  is equal to the first non-sensitive length- $k$  substring of  $W$ , and it is clearly the shortest string such that **C1** and **P1-P4** hold for  $W[0..j]$ .

*Inductive hypothesis and step:*  $X_{j-1}$  is optimal for  $j > j_{\min}$ . If  $C[j-k] = C[j-k+1] = 0$ ,  $X_j = X_{j-1}W[j]$  and this is clearly optimal. If  $C[j-k+1] = 1$ ,  $X_j = X_{j-1}$  thus still optimal. Finally, if  $C[j-k] = 1$  and  $C[j-k+1] = 0$  we have two subcases: if  $W[f..f+k-2] = W[j-k+1..j-1]$  then  $X_j = X_{j-1}W[j]$ , and once again  $X_j$  is evidently optimal. Otherwise,  $X_j = X_{j-1}\#W[j-k+1..j]$ . Suppose by contradiction that there exists a shorter  $X'_j$  such that **C1** and **P1-P4** still hold: either drop  $\#$  or append less than  $k$  letters after  $\#$ . If we appended less than  $k$  letters after  $\#$ , since TFS-ALGO will not read  $W[j]$  ever again, **P2-P3** would be violated, as an occurrence of  $W[j-k+1..j]$  would be missed. Without  $\#$ , the last  $k$  letters of  $X_{j-1}W[j-k+1]$  would violate either **C1** or **P1** and **P2** (since we suppose  $W[f..f+k-2] \neq W[j-k+1..j-1]$ ). Then  $X_j$  is optimal.  $\square$

**Theorem 1.** *Let  $W$  be a string of length  $n$  over  $\Sigma = \{1, \dots, n^{\mathcal{O}(1)}\}$ . Given  $k < n$  and  $\mathcal{S}$ , TFS-ALGO solves Problem 1 in  $\mathcal{O}(kn)$  time, which is worst-case optimal. An  $\mathcal{O}(n)$ -sized representation of  $X$  can be built in  $\mathcal{O}(n)$  time.*

*Proof.* For the first part inspect TFS-ALGO. Lines 2-4 can be realized in  $\mathcal{O}(n)$  time. The *while* loop in Line 5 is executed no more than  $n$  times, and every operation inside the loop takes  $\mathcal{O}(1)$  time except for Line 13 and Line 17 which take  $\mathcal{O}(k)$  time. Correctness and optimality follow directly from Lemma 1 (**P4**).

For the second part, we assume that  $X$  is represented by  $W$  and a sequence of pointers  $[i, j]$  to  $W$  interleaved (if necessary) by occurrences of  $\#$ . In Line 17, we can use an interval  $[i, j]$  to represent the length- $k$  substring of  $W$  added to  $X$ . In all other lines (Lines 4, 8 and 14) we can use  $[i, i]$  as one letter is added to  $X$  per one letter of  $W$ . By Lemma 1 we can have at most  $\lfloor \frac{n-k+1}{2} \rfloor$  occurrences of letter  $\#$ . The check at Line 13 can be implemented in constant time after linear-time pre-processing of  $W$  for longest common extension queries [?]. All other operations take in total linear time in  $n$ . Thus there exists an  $\mathcal{O}(n)$ -sized representation of  $X$  and it is constructible in  $\mathcal{O}(n)$  time.  $\square$

## 4 PFS-ALGO

Lemma 1 tells us that  $X$  is the shortest string satisfying constraint **C1** and properties **P1-P4**. If we were to drop **P1** and employ the partial order  $\Pi 1$  (see Problem 2), the length of  $X = X_1\#\dots\#X_N$  would not always be minimal: if a *permutation* of the strings  $X_1, \dots, X_N$  contains pairs  $X_i, X_j$  with a suffix-prefix overlap of length  $\ell = k - 1$ , we may further apply **R2**, obtaining a shorter string while still satisfying  $\Pi 1$ .

We propose PFS-ALGO to find such a permutation efficiently constructing a shorter string  $Y$  from  $W$ . The crux of our algorithm is an efficient method to



solve a variant of the classic NP-complete *Shortest Common Superstring* (SCS) problem [?]. Specifically our algorithm: (I) Computes string  $X$  using Theorem 1. (II) Constructs a collection  $\mathcal{B}'$  of strings, each of two symbols (two identifiers) and in a one-to-one correspondence with the elements of  $\mathcal{B} = \{X_1, \dots, X_N\}$ : the first (resp., second) symbol of the  $i$ th element of  $\mathcal{B}'$  is a unique identifier of the string corresponding to the length- $\ell$  prefix (resp., suffix) of the  $i$ th element of  $\mathcal{B}$ . (III) Computes a shortest string containing every element in  $\mathcal{B}'$  as a distinct substring. (IV) Constructs  $Y$  by mapping back each element to its distinct substring in  $\mathcal{B}$ . If there are multiple possible shortest strings, one is selected arbitrarily.

*Example 3 (Illustration of the workings of PFS-ALGO).* Let  $\ell = k - 1 = 3$  and

$$X = \text{aabc}\#\text{bccaab}\#\text{bbca}\#\text{aaabac}\#\text{aabcbbc}.$$

The collection  $\mathcal{B}$  is `aabc`, `bccaab`, `bbca`, `aaabac`, `aabcbbc`, and the collection  $\mathcal{B}'$  is `24`, `62`, `45`, `13`, `24` (id of prefix  $\cdot$  id of suffix). A shortest string containing all elements of  $\mathcal{B}'$  as distinct substrings is: `13 · 24 · 6245` (obtained by permuting the original string as `13, 24, 62, 24, 45` then applying **R2** twice). This shortest string is mapped back to the solution  $Y = \text{aaabac}\#\text{aabc}\#\text{bccaab}\#\text{bbca}$ . For example, `13` is mapped back to `aaabac`. Note,  $Y$  contains two occurrences of `#` and has length 24, while  $X$  contains 4 occurrences of `#` and has length 32.  $\square$

We now present the details of PFS-ALGO. We first introduce the *Fixed-Overlap Shortest String with Multiplicities* (FO-SSM) problem: Given a collection  $\mathcal{B}$  of strings  $B_1, \dots, B_{|\mathcal{B}|}$  and an integer  $\ell$ , with  $|B_i| > \ell$ , for all  $1 \leq i \leq |\mathcal{B}|$ , FO-SSM seeks to find a shortest string containing each element of  $\mathcal{B}$  as a distinct substring using the following operations on any pair of strings  $B_i, B_j$ :

1.  $\text{concat}(B_i, B_j) = B_i \cdot B_j$ ;
2.  $\ell\text{-merge}(B_i, B_j) = B_i[0..|B_i| - 1 - \ell]B_j[0..|B_j| - 1] = B_i[0..|B_i| - 1 - \ell] \cdot B_j$ .

Any solution to FO-SSM with  $\ell := k - 1$  and  $\mathcal{B} := X_1, \dots, X_N$  implies a solution to the PFS problem, because  $|X_i| > k - 1$  for all  $i$ 's (see Lemma 1, **P3**)

The FO-SSM problem is a variant of the SCS problem. In the SCS problem, we are given a *set* of strings and we are asked to compute the shortest common superstring of the elements of this set. The SCS problem is known to be NP-Complete, even for binary strings [?]. However, if all strings are of length two, the SCS problem admits a linear-time solution [?]. We exploit this crucial detail positively to show a linear-time solution to the FO-SSM problem in Lemma 3. In order to arrive to this result, we first adapt the SCS linear-time solution of [?] to our needs (see Lemma 2) and plug this solution to Lemma 3.

**Lemma 2.** *Let  $\mathcal{Q}$  be a collection of  $q$  strings, each of length two, over an alphabet  $\Sigma = \{1, \dots, (2q)^{\mathcal{O}(1)}\}$ . We can compute a shortest string containing every element of  $\mathcal{Q}$  as a distinct substring in  $\mathcal{O}(q)$  time.*

*Proof.* We sort the elements of  $\mathcal{Q}$  lexicographically in  $\mathcal{O}(q)$  time using radixsort. We also replace every letter in these strings with their *lexicographic rank* from  $\{1, \dots, 2q\}$  in  $\mathcal{O}(q)$  time using radixsort. In  $\mathcal{O}(q)$  time we construct the de Bruijn multigraph  $G$  of these strings [?]. Within the same time complexity, we find all

nodes  $v$  in  $G$  with in-degree, denoted by  $\text{IN}(v)$ , smaller than out-degree, denoted by  $\text{OUT}(v)$ . We perform the following two steps:

**Step 1:** While there exists a node  $v$  in  $G$  with  $\text{IN}(v) < \text{OUT}(v)$ , we start an arbitrary path (with possibly repeated nodes) from  $v$ , traverse consecutive edges and delete them. Each time we delete an edge, we update the in- and out-degree of the affected nodes. We stop traversing edges when a node  $v'$  with  $\text{OUT}(v') = 0$  is reached: whenever  $\text{IN}(v') = \text{OUT}(v') = 0$ , we also delete  $v'$  from  $G$ . Then, we add the traversed path  $p = v \dots v'$  to a set  $\mathcal{P}$  of paths. The path can contain the same node  $v$  more than once. If  $G$  is empty we halt. Proceeding this way, there are no two elements  $p_1$  and  $p_2$  in  $\mathcal{P}$  such that  $p_1$  starts with  $v$  and  $p_2$  ends with  $v$ ; thus this path decomposition is minimal. If  $G$  is not empty at the end, by construction, it consists of only cycles.

**Step 2:** While  $G$  is not empty, we perform the following. If there exists a cycle  $c$  that *intersects* with any path  $p$  in  $\mathcal{P}$  we splice  $c$  with  $p$ , update  $p$  with the result of splicing, and delete  $c$  from  $G$ . This operation can be efficiently implemented by maintaining an array  $A$  of size  $2q$  of linked lists over the paths in  $\mathcal{P}$ :  $A[\alpha]$  stores a list of pointers to all occurrences of letter  $\alpha$  in the elements of  $\mathcal{P}$ . Thus in constant time per node of  $c$  we check if any such path  $p$  exists in  $\mathcal{P}$  and splice the two in this case. If no such path exists in  $\mathcal{P}$ , we add to  $\mathcal{P}$  any of the path-linearizations of the cycle, and delete the cycle from  $G$ . After each change to  $\mathcal{P}$ , we update  $A$  and delete every node  $u$  with  $\text{IN}(u) = \text{OUT}(u) = 0$  from  $G$ .

The correctness of this algorithm follows from the fact that  $\mathcal{P}$  is a minimal path decomposition of  $G$ . Thus any concatenation of paths in  $\mathcal{P}$  represents a shortest string containing all elements in  $\mathcal{Q}$  as distinct substrings.  $\square$

Omitted proofs of Lemmas 3 and 4 can be found in [?].

**Lemma 3.** *Let  $\mathcal{B}$  be a collection of strings over an alphabet  $\Sigma = \{1, \dots, |\mathcal{B}|^{\mathcal{O}(1)}\}$ . Given an integer  $\ell$ , the FO-SSM problem for  $\mathcal{B}$  can be solved in  $\mathcal{O}(|\mathcal{B}|)$  time.*

Thus, PFS-ALGO applies Lemma 3 on  $\mathcal{B} := X_1, \dots, X_N$  with  $\ell := k - 1$  (recall that  $X_1 \# \dots \# X_N = X$ ). Note that each time the `concat` operation is performed, it also places the letter `#` in between the two strings.

**Lemma 4.** *Let  $W$  be a string of length  $n$  over an alphabet  $\Sigma$ . Given  $k < n$  and array  $C$ , PFS-ALGO constructs a shortest string  $Y$  with **C1**, **I1**, and **P2-P4**.*

**Theorem 2.** *Let  $W$  be a string of length  $n$  over  $\Sigma = \{1, \dots, n^{\mathcal{O}(1)}\}$ . Given  $k < n$  and  $S$ , PFS-ALGO solves Problem 2 in the optimal  $\mathcal{O}(n + |Y|)$  time.*

*Proof.* We compute the  $\mathcal{O}(n)$ -sized representation of string  $X$  with respect to  $W$  described in the proof of Theorem 1. This can be done in  $\mathcal{O}(n)$  time. If  $X \in \Sigma^*$ , then we construct and return  $Y := X$  in time  $\mathcal{O}(|Y|)$  from the representation. If  $X \in (\Sigma \cup \{\#\})^*$ , implying  $|Y| \leq |X|$ , we compute the LCP data structure of string  $W$  in  $\mathcal{O}(n)$  time [?]; and implement Lemma 3 in  $\mathcal{O}(n)$  time by avoiding to read string  $X$  explicitly: we rather rename  $X_1, \dots, X_N$  to a collection of two-letter strings by employing the LCP information of  $W$  directly. We then construct and report  $Y$  in time  $\mathcal{O}(|Y|)$ . Correctness follows directly from Lemma 4.  $\square$

## 5 The MCSR Problem and MCSR-ALGO

The strings  $X$  and  $Y$ , constructed by TFS-ALGO and PFS-ALGO, respectively, may contain the separator  $\#$ , which reveals information about the location of the sensitive patterns in  $W$ . Specifically, a malicious data recipient can go to the position of a  $\#$  in  $X$  and “undo” Rule **R1** that has been applied by TFS-ALGO, removing  $\#$  and the  $k-1$  letters after  $\#$  from  $X$ . The result will be an occurrence of the sensitive pattern. For example, applying this process to the first  $\#$  in  $X$  shown in Fig. 1, results in recovering the sensitive pattern **abab**. A similar attack is possible on the string  $Y$  produced by PFS-ALGO, although it is hampered by the fact that substrings within two consecutive  $\#$ s in  $X$  often swap places in  $Y$ .

To address this issue, we seek to construct a new string  $Z$ , in which  $\#$ s are either deleted or replaced by letters from  $\Sigma$ . To preserve privacy, we require separator replacements not to reinstate sensitive patterns. To preserve data utility, we favor separator replacements that have a small cost in terms of occurrences of  $\tau$ -ghosts (patterns with frequency less than  $\tau$  in  $W$  and at least  $\tau$  in  $Z$ ) and incur a bounded level of distortion in  $Z$ , as defined below. This is the MCSR problem, a restricted version of which is presented in Problem 3. The restricted version is referred to as  $\text{MCSR}_{k=1}$  and differs from MCSR in that it uses  $k=1$  for the pattern length instead of an arbitrary value  $k > 0$ .  $\text{MCSR}_{k=1}$  is presented next for simplicity and because it is used in the proof of Lemma 5 (see [?] for the proof). Lemma 5 implies Theorem 3.

**Problem 3 ( $\text{MCSR}_{k=1}$ ).** *Given a string  $Y$  over an alphabet  $\Sigma \cup \{\#\}$  with  $\delta > 0$  occurrences of letter  $\#$ , and parameters  $\tau$  and  $\theta$ , construct a new string  $Z$  by substituting the  $\delta$  occurrences of  $\#$  in  $Y$  with letters from  $\Sigma$ , such that:*

$$(I) \quad \sum_{\substack{i: Y[i]=\#, \text{Freq}_Y(Z[i]) < \tau \\ \text{Freq}_Z(Z[i]) \geq \tau}} \text{Ghost}(i, Z[i]) \text{ is minimum, and } (II) \quad \sum_{i: Y[i]=\#} \text{Sub}(i, Z[i]) \leq \theta.$$

The cost of  $\tau$ -ghosts is captured by a function  $\text{Ghost}$ . This function assigns a cost to an occurrence of a  $\tau$ -ghost, which is caused by a separator replacement at position  $i$ , and is specified based on domain knowledge. For example, with a cost equal to 1 for each gained occurrence of each  $\tau$ -ghost, we penalize more heavily a  $\tau$ -ghost with frequency much below  $\tau$  in  $Y$  and the penalty increases with the number of gained occurrences. Moreover, we may want to penalize positions towards the end of a temporally ordered string, to avoid spurious patterns that would be deemed important in applications based on time-decaying models [?].

The replacement distortion is captured by a function  $\text{Sub}$  which assigns a weight to a letter that could replace a  $\#$  and is specified based on domain knowledge. The maximum allowable replacement distortion is  $\theta$ . Small weights favor the replacement of separators with desirable letters (*e.g.*, letters that reinstate non-sensitive frequent patterns) and letters that reinstate sensitive patterns are assigned a weight larger than  $\theta$  that prohibits them from replacing a  $\#$ . Similarly, weights larger than  $\theta$  are assigned to letters which would lead to implausible patterns [?] if they replaced  $\#$ s.

**Lemma 5.** *The  $\text{MCSR}_{k=1}$  problem is NP-hard.*

**Theorem 3.** *The MCSR problem is NP-hard.*

**MCSR-ALGO.** Our MCSR-ALGO is a non-trivial heuristic that exploits the connection of the MCSR and MCK [?] problems and works by:

- (I) Constructing the set of all candidate  $\tau$ -ghost patterns (*i.e.*, length- $k$  strings over  $\Sigma$  with frequency below  $\tau$  in  $Y$  that can have frequency at least  $\tau$  in  $Z$ ).
- (II) Creating an instance of MCK from an instance of MCSR. For this, we map the  $i$ th occurrence of  $\#$  to a class  $C_i$  in MCK and each possible replacement of the occurrence with a letter  $j$  to a different item in  $C_i$ . Specifically, we consider all possible replacements with letters in  $\Sigma$  and also a replacement with the empty string, which models deleting (instead of replacing) the  $i$ th occurrence of  $\#$ . In addition, we set the costs and weights that are input to MCK as follows. The cost for replacing the  $i$ th occurrence of  $\#$  with the letter  $j$  is set to the sum of the Ghost function for all candidate  $\tau$ -ghost patterns when the  $i$ th occurrence of  $\#$  is replaced by  $j$ . That is, we make the worst-case assumption that the replacement forces all candidate  $\tau$ -ghosts to become  $\tau$ -ghosts in  $Z$ . The weight for replacing the  $i$ th occurrence of  $\#$  with letter  $j$  is set to  $\text{Sub}(i, j)$ .
- (III) Solving the instance of MCK and translating the solution back to a (possibly suboptimal) solution of the MCSR problem. For this, we replace the  $i$ th occurrence of  $\#$  with the letter corresponding to the element chosen by the MCK algorithm from class  $C_i$ , and similarly for each other occurrence of  $\#$ . If the instance has no solution (*i.e.*, no possible replacement can hide the sensitive patterns), MCSR-ALGO reports that  $Z$  cannot be constructed and terminates.

Lemma 6 below states the running time of MCSR-ALGO (see [?] for the proof on an efficient implementation of this algorithm).

**Lemma 6.** *MCSR-ALGO runs in  $\mathcal{O}(|Y| + k\delta\sigma + \mathcal{T}(\delta, \sigma))$  time, where  $\mathcal{T}(\delta, \sigma)$  is the running time of the MCK algorithm for  $\delta$  classes with  $\sigma + 1$  elements each.*

## 6 Experimental Evaluation

We evaluate our approach, referred to as TPM, in terms of *data utility* and *efficiency*. Given a string  $W$  over  $\Sigma$ , TPM sanitizes  $W$  by applying TFS-ALGO, PFS-ALGO, and then MCSR-ALGO, which uses the  $\mathcal{O}(\delta\sigma\theta)$ -time algorithm of [?] for solving the MCK instances. The final output is a string  $Z$  over  $\Sigma$ .

**Experimental Setup and Data.** We do not compare TPM against existing methods, because they are not alternatives to our approach (see Section 7). Instead, we compared against a greedy baseline referred to as BA.

BA initializes its output string  $Z_{\text{BA}}$  to  $W$  and then considers each sensitive pattern  $R$  in  $Z_{\text{BA}}$ , from left to right. For each  $R$ , it replaces the letter  $r$  of  $R$  that has the largest frequency in  $Z_{\text{BA}}$  with another letter  $r'$  that is not contained in  $R$  and has the smallest frequency in  $Z_{\text{BA}}$ , breaking all ties arbitrarily. If no such  $r'$  exists,  $r$  is replaced by  $\#$  to ensure that a solution is produced (even if it may reveal the location of a sensitive pattern). Each replacement removes the occurrence of  $R$  and aims to prevent  $\tau$ -ghost occurrences by selecting an  $r'$  that will not substantially increase the frequency of patterns overlapping with  $R$ .

We considered the following publicly available datasets used in [?, ?, ?, ?]: Oldenburg (OLD), Trucks (TRU), MSNBC (MSN), the complete genome of

*Escherichia coli* (DNA), and synthetic data (uniformly random strings, the largest of which is referred to as SYN). See Table 1 for the characteristics of these datasets and the parameter values used in experiments, unless stated otherwise.

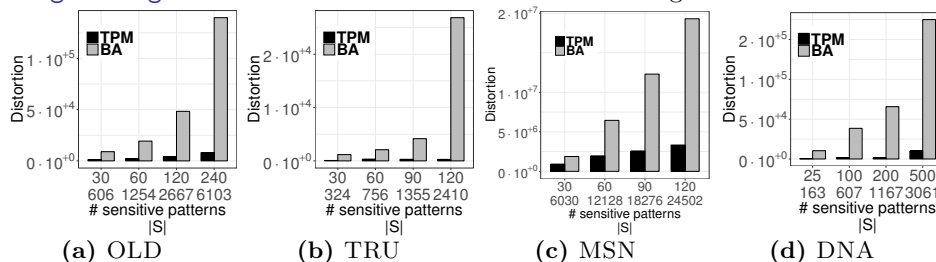
Dataset	Data domain	Length $n$	Alphabet size $ \Sigma $	# sensitive patterns	# sensitive positions $ \mathcal{S} $	Pattern length $k$
OLD	Movement	85,563	100	[30, 240] <b>(60)</b>	[600, 6103]	[3, 7] <b>(4)</b>
TRU	Transportation	5,763	100	[30, 120] <b>(10)</b>	[324, 2410]	[2, 5] <b>(4)</b>
MSN	Web	4,698,764	17	[30, 120] <b>(60)</b>	[6030, 320480]	[3, 8] <b>(4)</b>
DNA	Genomic	4,641,652	4	[25, 500] <b>(100)</b>	[163, 3488]	[5, 15] <b>(13)</b>
SYN	Synthetic	20,000,000	10	[10, 1000] <b>(1000)</b>	[10724, 20171]	[3, 6] <b>(6)</b>

**Table 1:** Characteristics of datasets and values used (default values are in bold).

The sensitive patterns were selected randomly among the frequent length- $k$  substrings at minimum support  $\tau$  following [?, ?, ?]. We used the fairly low values  $\tau = 10$ ,  $\tau = 20$ ,  $\tau = 200$ , and  $\tau = 20$  for TRU, OLD, MSN, and DNA, respectively, to have a wider selection of sensitive patterns. We also used a uniform cost of 1 for every occurrence of each  $\tau$ -ghost, a weight of 1 (resp.,  $\infty$ ) for each letter replacement that does not (resp., does) create a sensitive pattern, and we further set  $\theta = \delta$ . This setup treats all candidate  $\tau$ -ghost patterns and all candidate letters for replacement uniformly, to facilitate a fair comparison with BA which cannot distinguish between  $\tau$ -ghost candidates or favor specific letters.

To capture the utility of sanitized data, we used the (*frequency*) *distortion* measure  $\sum_U (\text{Freq}_W(U) - \text{Freq}_Z(U))^2$ , where  $U \in \Sigma^k$  is a non-sensitive pattern. The distortion measure quantifies changes in the frequency of non-sensitive patterns with low values suggesting that  $Z$  remains useful for tasks based on pattern frequency (*e.g.*, identifying motifs corresponding to functional or conserved DNA [?]). We also measured the number of  $\tau$ -ghost and  $\tau$ -lost patterns in  $Z$  following [?, ?, ?], where a pattern  $U$  is  $\tau$ -lost in  $Z$  if and only if  $\text{Freq}_W(U) \geq \tau$  but  $\text{Freq}_Z(U) < \tau$ . That is,  $\tau$ -lost patterns model knowledge that can no longer be mined from  $Z$  but could be mined from  $W$ , whereas  $\tau$ -ghost patterns model knowledge that can be mined from  $Z$  but not from  $W$ . A small number of  $\tau$ -lost/ghost patterns suggests that frequent pattern mining can be accurately performed on  $Z$  [?, ?, ?]. Unlike BA, by design TPM *does not* incur any  $\tau$ -lost pattern, as TFS-ALGO and PFS-ALGO preserve frequencies of nonsensitive patterns, and MCSR-ALGO can only increase pattern frequencies.

All experiments ran on an Intel Xeon E5-2640 at 2.66GHz with 16GB RAM. Our source code, written in C++, is available at <https://bitbucket.org/stringssanitization>. The results have been averaged over 10 runs.



**Fig. 2:** Distortion vs. number of sensitive patterns and their total number  $|\mathcal{S}|$  of occurrences in  $W$  (first two lines on the  $X$  axis).

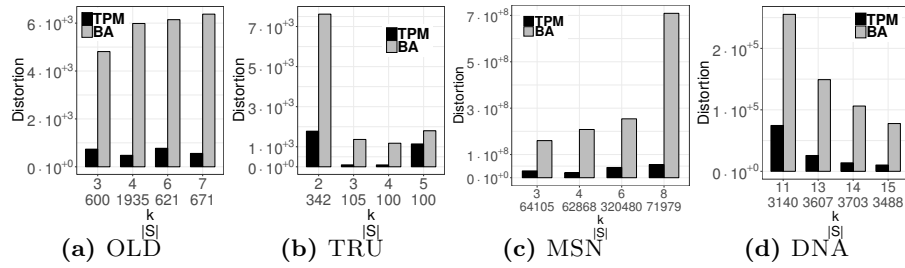


Fig. 3: Distortion vs. length of sensitive patterns  $k$  (and  $|S|$ ).

**Data Utility.** We first demonstrate that TPM incurs *very low distortion*, which implies high utility for tasks based on the frequency of patterns (e.g., [?]). Fig. 2 shows that, for varying number of sensitive patterns, TPM incurred on average 18.4 (and up to 95) times lower distortion than BA over all experiments. Also, Fig. 2 shows that TPM remains effective even in challenging settings, with many sensitive patterns (e.g., the last point in Fig. 2b where about 42% of the positions in  $W$  are sensitive). Fig. 3 shows that, for varying  $k$ , TPM caused on average 7.6 (and up to 14) times lower distortion than BA over all experiments.

Next, we demonstrate that TPM permits *accurate frequent pattern mining*: Fig. 4 shows that TPM led to no  $\tau$ -lost or  $\tau$ -ghost patterns for the TRU and MSN datasets. This implies no utility loss for mining frequent length- $k$  substrings with threshold  $\tau$ . In all other cases, the number of  $\tau$ -ghosts was on average 6 (and up to 12) times smaller than the total number of  $\tau$ -lost and  $\tau$ -ghost patterns for BA. BA performed poorly (e.g., up to 44% of frequent patterns became  $\tau$ -lost for TRU and 27% for DNA). Fig. 5 shows that, for varying  $k$ , TPM led to on average 5.8 (and up to 19) times fewer  $\tau$ -lost/ghost patterns than BA. BA performed poorly (e.g., up to 98% of frequent patterns became  $\tau$ -lost for DNA).

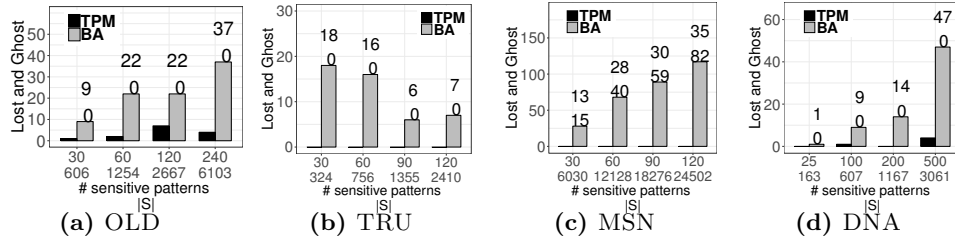


Fig. 4: Total number of  $\tau$ -lost and  $\tau$ -ghost patterns vs. number of sensitive patterns (and  $|S|$ ).  $x$ / $y$  on the top of each bar for BA denotes  $x$   $\tau$ -lost and  $y$   $\tau$ -ghost patterns.

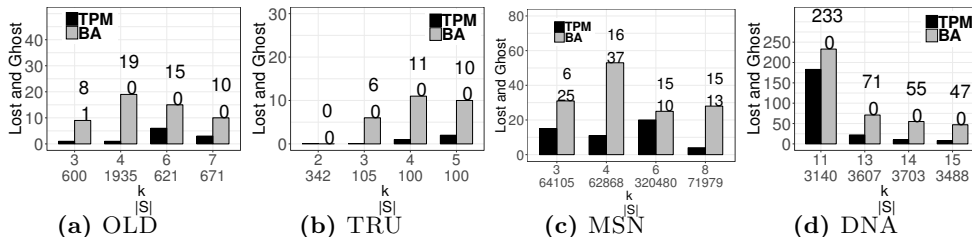
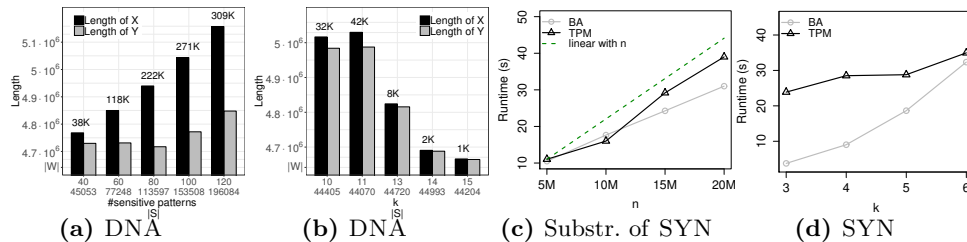


Fig. 5: Total number of  $\tau$ -lost and  $\tau$ -ghost patterns vs. length of sensitive patterns  $k$  (and  $|S|$ ).  $x$ / $y$  on the top of each bar for BA denotes  $x$   $\tau$ -lost and  $y$   $\tau$ -ghost patterns.

We also demonstrate that PFS-ALGO reduces the length of the output string  $X$  of TFS-ALGO substantially, creating a string  $Y$  that contains *less redundant information* and allows for more efficient analysis. Fig. 6a shows the length of  $X$  and of  $Y$  and their difference for  $k = 5$ .  $Y$  was much shorter than  $X$  and its length decreased with the number of sensitive patterns, since more substrings had a suffix-prefix overlap of length  $k - 1 = 4$  and were removed (see Section 4). Interestingly, the length of  $Y$  was close to that of  $W$  (the string before sanitization). A larger  $k$  led to less substantial length reduction as shown in Fig. 6b (but still few thousand letters were removed), since it is less likely for long substrings to have an overlap and be removed.



**Fig. 6:** Length of  $X$  and  $Y$  (output of TFS-ALGO and PFS-ALGO, resp.) for varying: (a) number of sensitive patterns (and  $|\mathcal{S}|$ ), (b) length of sensitive patterns  $k$  (and  $|\mathcal{S}|$ ). On the top of each pair of bars we plot  $|X| - |Y|$ . Runtime on synthetic data for varying: (c) length  $n$  of string and (d) length  $k$  of sensitive patterns. Note that  $|Y| = |Z|$ .

**Efficiency.** We finally measured the runtime of TPM using prefixes of the synthetic string SYN whose length  $n$  is 20 million letters. Fig. 6c (resp., Fig. 6d) shows that TPM scaled linearly with  $n$  (resp.,  $k$ ), as predicted by our analysis in Section 5 (TPM takes  $\mathcal{O}(n + |Y| + kd\sigma + \delta\sigma\theta) = \mathcal{O}(kn + kd\sigma + \delta\sigma\theta)$  time, since the algorithm of [?] was used for MCK instances). In addition, TPM is efficient, with a runtime similar to that of BA and less than 40 seconds for SYN.

## 7 Related Work

Data sanitization (*a.k.a.* knowledge hiding) aims at concealing patterns modeling confidential knowledge by limiting their frequency, so that they are not easily mined from the data. Existing methods are applied to: (I) a *collection* of set-valued data (transactions) [?] or spatiotemporal data (trajectories) [?]; (II) a *collection* of sequences [?,?]; or (III) a *single* sequence [?,?,?]. Yet, none of these methods follows our CSD setting: Methods in category I are not applicable to string data, and those in categories II and III do not have guarantees on privacy-related constraints [?] or on utility-related properties [?,?,?,?]. Specifically, unlike our approach, [?] cannot guarantee that all sensitive patterns are concealed (constraint C1), while [?,?,?,?] do not guarantee the satisfaction of utility properties (*e.g.*, I1 and P2).

Anonymization aims to prevent the disclosure of individuals' identity and/or information that individuals are not willing to be associated with [?]. Anonymization works (*e.g.*, [?,?,?]) are thus not alternatives to our work (see [?] for details).

## 8 Conclusion

In this paper, we introduced the Combinatorial String Dissemination model. The focus of this model is on *guaranteeing* privacy-utility trade-offs (*e.g.*, C1 vs. I1

and **P2**). We defined a problem (TFS) which seeks to produce the shortest string that preserves the order of appearance and the frequency of all non-sensitive patterns; and a variant (PFS) that preserves a partial order and the frequency of the non-sensitive patterns but produces a shorter string. We developed two time-optimal algorithms, TFS-ALGO and PFS-ALGO, for the problem and its variant, respectively. We also developed MCSR-ALGO, a heuristic that prevents the disclosure of the location of sensitive patterns from the outputs of TFS-ALGO and PFS-ALGO. Our experiments show that sanitizing a string by TFS-ALGO, PFS-ALGO and then MCSR-ALGO is effective and efficient.

**Acknowledgments.** HC is supported by a CSC scholarship. GR and NP are partially supported by MIUR-SIR project CMACBioSeq grant n. RBSI146R5L. We acknowledge the use of the Rosalind HPC cluster hosted by King's College London.