



HAL
open science

Consensus Strings with Small Maximum Distance and Small Distance Sum

Laurent Bulteau, Markus L. Schmid

► **To cite this version:**

Laurent Bulteau, Markus L. Schmid. Consensus Strings with Small Maximum Distance and Small Distance Sum. 43rd International Symposium on Mathematical Foundations of Computer Science (MFCS 2018)., Aug 2018, Liverpool, United Kingdom. 10.4230/LIPIcs.MFCS.2018.1 . hal-01930623

HAL Id: hal-01930623

<https://hal.science/hal-01930623>

Submitted on 22 Nov 2018

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.

1 Consensus Strings with Small Maximum Distance 2 and Small Distance Sum

3 **Laurent Bulteau**

4 Université Paris-Est, LIGM (UMR 8049), CNRS, ENPC, ESIEE Paris, UPEM, F-77454,
5 Marne-la-Vallée, France
6 laurent.bulteau@u-pem.fr

7 **Markus L. Schmid**

8 Fachbereich 4 – Abteilung Informatikwissenschaften, Universität Trier, 54286 Trier, Germany
9 mlschmid@mlschmid.de

10 — Abstract —

11 The parameterised complexity of consensus string problems (CLOSEST STRING, CLOSEST SUB-
12 STRING, CLOSEST STRING WITH OUTLIERS) is investigated in a more general setting, i. e., with
13 a bound on the maximum Hamming distance *and* a bound on the sum of Hamming distances
14 between solution and input strings. We completely settle the parameterised complexity of these
15 generalised variants of CLOSEST STRING and CLOSEST SUBSTRING, and partly for CLOSEST
16 STRING WITH OUTLIERS; in addition, we answer some open questions from the literature re-
17 garding the classical problem variants with only one distance bound. Finally, we investigate the
18 question of polynomial kernels and respective lower bounds.

19 **2012 ACM Subject Classification** Theory of computation → Problems, reductions and com-
20 pleteness, Theory of computation → Fixed parameter tractability, Theory of computation → W
21 hierarchy

22 **Keywords and phrases** Consensus String Problems, Closest String, Closest Substring, Parame-
23 terised Complexity, Kernelisation

24 **Digital Object Identifier** 10.4230/LIPIcs.MFCS.2018.1

25 **1** Introduction

26 *Consensus string problems* have the following general form: given input strings $S =$
27 $\{s_1, \dots, s_k\}$ and a distance bound d , find a string s with distance at most d from the
28 input strings. With the Hamming distance as the central distance measure for strings,
29 there are two obvious types of distance between a single string and a set S of strings: the
30 maximum distance between s and any string from S (called *radius*) and the sum of all
31 distances between s and strings from S (called *distance sum*). The most basic consensus
32 string problem is CLOSEST STRING, where we get a set S of k length- ℓ strings and a bound
33 d , and ask whether there exists a length- ℓ *solution string* s with radius at most d . This
34 problem is NP-complete (see [?]), but fixed-parameter tractable for many variants (see [?]),
35 including the parameterisation by d , which in biological applications can often be assumed
36 to be small (see [?, ?]). A classical extension is CLOSEST SUBSTRING, where the strings of S
37 have length *at most* ℓ , the solution string must have a given length m and the radius bound d
38 is w. r. t. some length- m substrings of the input strings. A parameterised complexity analysis
39 (see [?, ?, ?]) has shown CLOSEST SUBSTRING to be harder than CLOSEST STRING. If we
40 bound the distance sum instead of the radius, then CLOSEST STRING collapses to a trivial
41 problem, while CLOSEST SUBSTRING, which is then called CONSENSUS PATTERNS, remains



© Laurent Bulteau and Markus L. Schmid;
licensed under Creative Commons License CC-BY

43rd International Symposium on Mathematical Foundations of Computer Science (MFCS 2018).

Editors: Igor Potapov, Paul Spirakis, and James Worrell; Article No. 1; pp. 1:1–1:13

Leibniz International Proceedings in Informatics



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

42 NP-complete. CLOSEST STRING WITH OUTLIERS is a recent extension, which is defined like
43 CLOSEST STRING, but with the possibility to ignore a given number of t input strings.

44 The main motivation for consensus string problems comes from the important task of
45 finding similar regions in DNA or other protein sequences, which arises in many different
46 contexts of computational biology, e. g., universal PCR primer design [?, ?, ?, ?], genetic
47 probe design [?], antisense drug design [?, ?], finding transcription factor binding sites in
48 genomic data [?], determining an unbiased consensus of a protein family [?], and motif-
49 recognition [?, ?, ?]. The consensus string problems are a formalisation of this computational
50 task and most variants of them are NP-hard. However, due to their high practical relevance,
51 it is necessary to solve them despite their intractability, which has motivated the study of
52 their approximability, on the one hand, but also their fixed-parameter tractability, on the
53 other (see the survey [?] for an overview of the parameterised complexity of consensus string
54 problems). This work is a contribution to the latter branch of research.

55 **Problem Definition.** Let Σ be a finite alphabet, Σ^* be the set of all strings over Σ ,
56 including the empty string ε and $\Sigma^+ = \Sigma^* \setminus \{\varepsilon\}$. For $w \in \Sigma^*$, $|w|$ is the length of w and,
57 for every i , $1 \leq i \leq |w|$, by $w[i]$, we refer to the symbol at position i of w . For every
58 $n \in \mathbb{N} \cup \{0\}$, let $\Sigma^n = \{w \in \Sigma^* \mid |w| = n\}$ and $\Sigma^{\leq n} = \bigcup_{i=0}^n \Sigma^i$. By \preceq , we denote the
59 *substring relation* over the set of strings, i. e., for $u, v \in \Sigma^*$, $u \preceq v$ if $v = xuy$, for some
60 $x, y \in \Sigma^*$. We use the concatenation of sets of strings as usually defined, i. e., for $A, B \subseteq \Sigma^*$,
61 $A \cdot B = \{uv \mid u \in A, v \in B\}$.

62 For strings $u, v \in \Sigma^*$ with $|u| = |v|$, $d_H(u, v)$ is the *Hamming distance* between u and v .
63 For a multi-set $S = \{u_i \mid 1 \leq i \leq n\} \subseteq \Sigma^\ell$ and a string $v \in \Sigma^\ell$, for some $\ell \in \mathbb{N}$, the *radius of S*
64 *(w. r. t. v)* is defined by $r_H(v, S) = \max\{d_H(v, u) \mid u \in S\}$ and the *distance sum of S* *(w. r. t. v)*
65 is defined by $s_H(v, S) = \sum_{u \in S} d_H(v, u)$.¹ Next, we state the problem (r, s)-CLOSEST STRING
66 in full detail, from which we then derive the other considered problems:

67 (r, s)-CLOSEST STRING

68 *Instance:* Multi-set $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$, $\ell \in \mathbb{N}$, and integers $d_r, d_s \in \mathbb{N}$.

69 *Question:* Is there an $s \in \Sigma^\ell$ with $r_H(s, S) \leq d_r$ and $s_H(s, S) \leq d_s$?

70 For (r, s)-CLOSEST SUBSTRING, we have $S \subseteq \Sigma^{\leq \ell}$ and an additional input integer $m \in \mathbb{N}$, and
71 we ask whether there is a multi-set $S' = \{s'_i \mid s'_i \preceq s_i, 1 \leq i \leq k\} \subseteq \Sigma^m$ with $r_H(s, S') \leq d_r$ and
72 $s_H(s, S') \leq d_s$. For (r, s)-CLOSEST STRING WITH OUTLIERS (or (r, s)-CLOSEST STRING-WO
73 for short) we have an additional input integer $t \in \mathbb{N}$, and we ask whether there is a multi-
74 set $S' \subseteq S$ with $|S'| = k - t$ such that $r_H(s, S') \leq d_r$ and $s_H(s, S') \leq d_s$. We also call
75 (r, s)-CLOSEST STRING the *general variant* of CLOSEST STRING, while (r)-CLOSEST STRING
76 and (s)-CLOSEST STRING denote the variants, where the only distance bound is d_r or d_s ,
77 respectively; we shall also call them the (r)- and (s)-*variant* of CLOSEST STRING. Analogous
78 notation apply to the other consensus string problems. The problem names that are also com-
79 monly used in the literature translate into our terminology as follows: CLOSEST STRING = (r)-
80 CLOSEST STRING, CLOSEST SUBSTRING = (r)-CLOSEST SUBSTRING, CONSENSUS PATTERNS
81 = (s)-CLOSEST SUBSTRING and CLOSEST STRING-WO = (r)-CLOSEST STRING-WO.

82 The motivation for our more general setting with respect to the bounds d_r and d_s is the
83 following. While the distance measures of radius and distance sum are well-motivated, they
84 have, if considered individually, also obvious deficiencies. In the distance sum variant, we

¹ Note that we slightly abuse notation with respect to the subset relation: for a multi-set A and a set B ,
 $A \subseteq B$ means that $A' \subseteq B$, where A' is the set obtained from A by deleting duplicates; for multi-sets
 A, B , $A \subseteq B$ is defined as usual. Moreover, whenever it is clear from the context that we talk about
multi-sets, we also simply use the term *set*.

85 may consider strings as close enough that are very close to some, but totally different to the
 86 other input strings. In the radius variant, on the other hand, we may consider strings as too
 87 different, even though they are very similar to all input strings except one, for which the
 88 bound is exceeded by only a small amount. Using an upper bound on the distance per each
 89 input string and an upper bound on the total sum of distances caters for these cases.²

90 For any problem K , by $K(p_1, p_2, \dots)$, we denote the variant of K parameterised by the
 91 parameters p_1, p_2, \dots . For unexplained concepts of parameterised complexity, we refer to the
 92 textbooks [?, ?, ?].

93 **Known Results.** In contrast to graph problems, where interesting parameters are often
 94 hidden in the graph structure, string problems typically contain a variety of obvious, but
 95 nevertheless interesting parameters that could be exploited in terms of fixed-parameter
 96 tractability. For the consensus string problems these are the number of input strings k ,
 97 their length ℓ , the radius bound d_r , the distance sum bound d_s , the alphabet size $|\Sigma|$, the
 98 substring length m (in case of (r, s) -CLOSEST SUBSTRING), the number of *outliers* t and
 99 *inliers* $k - t$ (in case of (r, s) -CLOSEST STRING-WO). This leads to a large number of different
 100 parameterisations, which justifies the hope for fixed-parameter tractable variants.

101 The parameterised complexity (w. r. t. the above mentioned parameters) of the radius
 102 as well as the distance sum variant of CLOSEST STRING and CLOSEST SUBSTRING has
 103 been settled by a sequence of papers (see [?, ?, ?, ?, ?] and, for a survey, [?]), except
 104 (s) -CLOSEST SUBSTRING with respect to parameter ℓ , which has been neglected in these
 105 papers and mentioned as an open problem in [?], in which it is shown that the fixed-parameter
 106 tractability results from (r) -CLOSEST STRING carry over to (r) -CLOSEST SUBSTRING, if we
 107 additionally parameterise by $(\ell - m)$. The parameterised complexity analysis of the radius
 108 variant of CLOSEST STRING WITH OUTLIERS has been started more recently in [?] and, to
 109 the knowledge of the authors, the distance sum variant has not yet been considered.

110 The parameterised complexity of the general variants, where we have a bound on both the
 111 radius and the distance sum, has not yet been considered in the literature. While there are
 112 obvious reductions from the (r) - and (s) -variants to the general variant, these three variants
 113 describe, especially in the parameterised setting, rather different problems.

114 **Our Contribution.** In this work, we answer some open questions from the literature
 115 regarding the (r) - and (s) -variants of the consensus string problems, and we initiate the
 116 parameterised complexity analysis of the general variants.

117 We extend all the FPT-results from (r) -CLOSEST STRING to the general variant; thus, we
 118 completely settle the fixed-parameter tractability of (r, s) -CLOSEST STRING. While for some
 119 parameterisations, this is straightforward, the case of parameter d_r follows from a non-trivial
 120 extension of the known branching algorithm for (r) -CLOSEST STRING(d_r) (see [?]).

121 For (r, s) -CLOSEST SUBSTRING, we classify all parameterised variants as being in FPT or
 122 W[1]-hard, which is done by answering the open question whether (s) -CLOSEST SUBSTRING(ℓ)
 123 is in FPT (see [?]) in the negative (which also settles the parameterised complexity of
 124 (s) -CLOSEST SUBSTRING) and by slightly adapting the existing FPT-algorithms.

125 Regarding (r, s) -CLOSEST STRING-WO, we solve an open question from [?] w. r. t. the
 126 radius variant, we show W[1]-hardness for a strong parameterisation of the (s) -variant, we
 127 show fixed-parameter tractability for some parameter combinations of the general variant and,
 128 as our main result, we present an FPT-algorithm (for the general variant) for parameters d_r ,

² To the knowledge of the authors, optimising both the radius and the distance sum has been considered first in [?], where algorithms for the special case $k = 3$ are considered.

129 and t (which is the same algorithm that shows (r, s) -CLOSEST STRING(d_r) \in FPT mentioned
130 above). Many other cases are left open for further research.

131 Finally, we investigate the question whether the fixed-parameter tractable variants of the
132 considered consensus string problems allow polynomial kernels; thus, continuing a line of work
133 initiated by Basavaraju et al. [?], in which kernelisation lower bounds for (r) -CLOSEST STRING
134 and (r) -CLOSEST SUBSTRING are proved. Our respective main result is a cross-composition
135 from (r) -CLOSEST STRING into (r) -CLOSEST STRING-WO.

136 Due to space constraints, proofs for results marked with $(*)$ are omitted.

137 2 Closest String and Closest String-wo

138 In this section, we investigate (r, s) -CLOSEST STRING and (r, s) -CLOSEST STRING-WO (and
139 their (r) - and (s) -variants) and we shall first give some useful definitions.

140 It will be convenient to treat a set $S = \{s_i \mid 1 \leq i \leq k\} \subseteq \Sigma^\ell$ as a $k \times \ell$ matrix with
141 entries from Σ . By the term *column of S* , we refer to the transpose of a column of the matrix
142 S , which is an element from Σ^k ; thus, the introduced string notations apply, e. g., if c is the
143 i^{th} column of S , then $c[j]$ corresponds to $s_j[i]$. A string $s \in \Sigma^\ell$ is a *majority string* (for a
144 set $S \subseteq \Sigma^\ell$) if, for every i , $1 \leq i \leq \ell$, $s[i]$ is a symbol with majority in the i^{th} column of S .
145 Obviously, $s_{\text{H}}(s, S) = \min\{s_{\text{H}}(s', S) \mid s' \in \Sigma^\ell\}$ if and only if s is a majority string for S . We
146 call a string $s \in \Sigma^\ell$ *radius optimal* or *distance sum optimal* (with respect to a set $S \subseteq \Sigma^\ell$) if
147 $r_{\text{H}}(s, S) = \min\{r_{\text{H}}(s', S) \mid s' \in \Sigma^\ell\}$ or $s_{\text{H}}(s, S) = \min\{s_{\text{H}}(s', S) \mid s' \in \Sigma^\ell\}$, respectively.

148 It is a well-known fact that (r) -CLOSEST STRING allows FPT-algorithms for any of the
149 single parameters k , d_r or ℓ , and it is still NP-hard for $|\Sigma| = 2$ (see [?]). While the latter
150 hardness result trivially carries over to (r, s) -CLOSEST STRING (by setting $d_s = k d_r$), we
151 have to modify the FPT-algorithms for extending the fixed-parameter tractability results
152 to (r, s) -CLOSEST STRING. We start with parameter k , for which we can extend the ILP-
153 approach that is used in [?] to show (r) -CLOSEST STRING(k) \in FPT.

154 ► **Theorem 1** (*). (r, s) -CLOSEST STRING(k) \in FPT.

155 Next, we consider the parameter d_r . For the (r) -variant of (r, s) -CLOSEST STRING,
156 the fixed-parameter tractability with respect to d_r is shown in [?] by a branching algo-
157 rithm, which proved itself as rather versatile: it has successfully been extended in [?] to
158 (r) -CLOSEST STRING-WO(d_r, t) and in [?] to (r) -CLOSEST SUBSTRING($d_r, (\ell - m)$).

159 We propose an extension of the same branching algorithm, that allows for a bound d_s on the
160 distance sum; thus, it works for (r, s) -CLOSEST STRING(d_r). In fact, we prove in Theorem 7
161 an even stronger result, where we also extend the algorithm to exclude up to t outlier strings
162 from the input set S , i. e., we extend it to the problem (r, s) -CLOSEST STRING-WO(d_r, t).
163 Since Theorem 3 can therefore be seen as a corollary of this result by taking $t = 0$, we only
164 give an informal description of a direct approach that solves (r, s) -CLOSEST STRING(d_r) (and
165 refer to Theorem 7 for a formal proof).

166 The core idea is to apply the branching algorithm starting with the majority string for
167 the input set S , instead of any random string from S . Then, as in [?], the algorithm would
168 replace some characters of the current string with characters of the solution string. This way,
169 it can be shown that the distance sum of the current string is always a lower bound of the
170 distance sum of the optimal string, which allows to cut any branch where the distance sum
171 goes beyond the threshold d_s . We prove the following lemma, which allows to bound the
172 depth of the search tree (and shall also be used in the proof of Theorem 7 later on):

k	d_r	d_s	$ \Sigma $	ℓ	Result	Note/Ref.
p	–	–	–	–	FPT	Thm. 1
–	p	–	–	–	FPT	Thm. 3
–	–	p	–	–	FPT	Cor. 4
–	–	–	2	–	NP-hard	from (r)-variant [?]
–	–	–	–	p	FPT	Cor. 4

■ **Table 1** Results for (r, s)-CLOSEST STRING.

173 ► **Lemma 2 (*)**. Let $S \subseteq \Sigma^\ell$, $s \in \Sigma^\ell$ such that $r_H(s, S) \leq d_r$, and let s_m be a majority string
174 for S . Then $d_H(s_m, s) \leq 2d_r$.

175 ► **Theorem 3**. (r, s)-CLOSEST STRING(d_r) \in FPT.

176 Obviously, we can assume $d_r \leq \ell$ and we can further assume that every column of S
177 contains at least two different symbols (all columns without this property could be removed),
178 which implies $s_H(s_i, S) \geq \ell$ for every $s \in \Sigma^\ell$; thus, we can assume $\ell \leq d_s$. Consequently, we
179 obtain the following corollary:

180 ► **Corollary 4**. (r, s)-CLOSEST STRING(ℓ) \in FPT, (r, s)-CLOSEST STRING(d_s) \in FPT.

181 This completely settles the parameterised complexity of (r, s)-CLOSEST STRING with
182 respect to parameters k , d_r , d_s , $|\Sigma|$ and ℓ ; recall that the (r)-variant is already settled, while
183 the (s)-variant is trivial.

184 2.1 (r, s)-Closest String-wo

185 We now turn to the problem (r, s)-CLOSEST STRING-WO and we first prove several fixed-
186 parameter tractability results for the general variant; in Sec. 2.2, we consider the (r)- and
187 (s)-variants separately.

188 First, we note that solving an instance of (r, s)-CLOSEST STRING-WO(k) can be reduced
189 to solving $f(k)$ many (r, s)-CLOSEST STRING(k)-instances, which, due to the fixed-parameter
190 tractability of the latter problem, yields the fixed-parameter tractability of the former.

191 ► **Theorem 5 (*)**. (r, s)-CLOSEST STRING-WO(k) \in FPT.

192 If the number $k - t$ of inliers exceeds d_s , then an (r, s)-CLOSEST STRING-WO-instance
193 becomes easily solvable; thus, $k - t$ can be bounded by d_s , which yields the following result:

194 ► **Theorem 6 (*)**. (r, s)-CLOSEST STRING-WO(d_s, t) \in FPT.

195 The algorithm introduced in [?] to prove (r)-CLOSEST STRING(d_r) \in FPT has been
196 extended in [?] with an additional branching that guesses whether a string s_j should be consid-
197 ered an outlier or not; thus, yielding fixed-parameter tractability of (r)-CLOSEST STRING-WO(d_r, t).
198 We present a non-trivial extension of this algorithm, with a carefully selected starting string,
199 to obtain the fixed-parameter tractability of (r, s)-CLOSEST STRING-WO(d_r, t) (and, as ex-
200 plained in Section 2, also of (r, s)-CLOSEST STRING(d_r)):

201 ► **Theorem 7**. (r, s)-CLOSEST STRING-WO(d_r, t) \in FPT.

202 **Proof.** Let (S, d_s, d_r, t) be a positive instance of (r, s)-CLOSEST STRING-WO(d_r, t) with $k \geq$
203 $5t$ (otherwise k can be considered as a parameter). A character x is *frequent in column* i if it

	Input:		$s_1 =$	dbaddcbcdbbdbb		
	$d_r = 5$		$s_2 =$	daaaacbcdccdbd		
	$d_s = 14$		$s_3 =$	daaddabcaccdbd		
	$t = 2$		$s_4 =$	aacdaccdccabd		
			$s_5 =$	aacdaabdaccadd		
	$D = 10$		$s_6 =$	acaaaabccddbadd		
Step	S'	t	s'	d	$r_H(s', S')$	action
1	$\{s_1, s_2, \dots, s_6\}$	2	$\diamond a \diamond \diamond \diamond b \diamond \diamond c \diamond \diamond d$	20	13	$s[3] \leftarrow s_1[3]$
2	$\{s_1, s_2, \dots, s_6\}$	2	$\diamond a a \diamond \diamond b \diamond \diamond c \diamond \diamond d$	19	12	$s[12] \leftarrow s_1[12]$
3	$\{s_1, s_2, \dots, s_6\}$	2	$\diamond a a \diamond \diamond b \diamond \diamond c \diamond d \diamond d$	18	11	remove s_6
4	$\{s_1, s_2, \dots, s_5\}$	1	$\diamond a a \diamond \diamond b \diamond \diamond c \diamond d \diamond d$	18	11	$s[6] \leftarrow s_1[6]$
5	$\{s_1, s_2, \dots, s_5\}$	1	$\diamond a a \diamond \diamond c b \diamond \diamond c \diamond d \diamond d$	17	10	remove s_5
6	$\{s_1, \dots, s_4\}$	0	$\diamond a a \diamond \diamond c b \diamond \diamond c \diamond d \diamond d$	17	10	
			$s'' =$	daadacbcdccdbd		$s[7] \leftarrow s_4[7]$
7	$\{s_1, \dots, s_4\}$	0	$\diamond a a \diamond \diamond c c \diamond \diamond c \diamond d \diamond d$	16	10	
			$s'' =$	daadaccdccdbd		return S', s''

■ **Figure 1** Example for Algorithm 1 on an instance of (r, s) -CLOSEST STRING-WO. The shown steps correspond to one branch that yields a correct solution. The algorithm starts with the majority string where disputed characters are replaced by \diamond . At each step, the algorithm either inserts a character from an input string at maximal distance from s' (note that even non-disputed characters may be replaced), or removes one such string. When $t = 0$, it is checked whether the completion s'' of s' is a correct solution. At step 7, we return a solution with $r_H(s'', S') = 5$ and $s_H(s'', S') = 14$.

204 has at least as many occurrences as the majority character minus t (thus, for any $S' \subseteq S$,
 205 $|S'| \geq |S| - t$, all majority characters for S' are frequent characters). A column i is *disputed*
 206 if it contains at least two frequent characters. Let D be the number of disputed columns.

207 Let (S^*, s^*) be a solution for this instance. In a disputed column i , no character
 208 occurs more than $\frac{k+t}{2}$ times, hence, among the $k - t$ strings of S^* , there are at least
 209 $(k - t) - \frac{k+t}{2} = \frac{k-3t}{2}$ mismatches at position i . The disputed columns thus introduce at least
 210 $D \frac{k-3t}{2}$ mismatches. Since the overall number of mismatches is upper-bounded by $d_r(k - t)$,
 211 we have $D \leq \frac{2d_r(k-t)}{k-3t} = 2d_r \left(1 + \frac{2t}{k-3t}\right)$, and, with $k \geq 5t$, the upper-bound $D \leq 4d_r$ follows.

212 We introduce a new character $\diamond \notin \Sigma$. A string $s' \in (\Sigma \cup \{\diamond\})^\ell$ is a *lower bound* for a
 213 solution s^* , if, for every i such that $s'[i] \neq s^*[i]$, either i is a disputed column and $s'[i] = \diamond$, or
 214 i is not disputed and $s'[i]$ is the majority character for column i of S^* (which is equal to the
 215 majority character for column i of S). Intuitively speaking, whenever a character $s'[i]$ differs
 216 from $s^*[i]$, it is the majority character of its column (except for disputed columns in which
 217 we use an “undecided” character \diamond). Finally, the *completion for S'* of a string $s' \in (\Sigma \cup \{\diamond\})^*$
 218 is the string obtained by replacing each occurrence of \diamond by a majority character of the
 219 corresponding column in S' .

220 We now prove that Algorithm 1 solves (r, s) -CLOSEST STRING-WO in time at most
 221 $O^*((d_r + 1)^{6d_r} 2^{6d_r+t})$, using the following three claims (see Figure 1 for an example).

222 *Claim 1:* Any call to SOLVE CLOSEST STRING-WO(S', t, s', d) always returns after a time
 223 $O^*((d_r + 1)^d 2^{d+t})$

224 *Proof of Claim 1:* We prove this running time by induction: if $d = t = 0$, then the function
 225 returns in Line 3 or 4; thus, it returns after polynomial time. Otherwise, it performs at most

ALGORITHM 1: SOLVE CLOSEST STRING-WO**Input** : $S' \subseteq S$, $t \in \mathbb{N}$, $s' \in (\Sigma \cup \{\diamond\})^\ell$, $d \in \mathbb{N}$ **Output**: a pair (S^*, s^*) or the symbol ∇

```

1 if  $t = 0$  then
2    $s'' =$  completion of  $s'$  in  $S'$ ;
3   if  $\mathfrak{s}_H(s'', S') \leq d_s$ , and  $\mathfrak{r}_H(s'', S') \leq d_r$  then return  $(S', s'')$ ;
4   if  $d = 0$  then return  $\nabla$ ;
5 Let  $s_j \in S'$  be such that  $\mathfrak{d}_H(s', s_j)$  is maximal;
6 if  $t > 0$  then
7    $sol =$  SOLVE CLOSEST STRING-WO( $S' \setminus \{s_j\}, t - 1, s', d$ );
8   if  $sol \neq \nabla$  then return  $sol$ ;
9 if  $d > 0$  then
10  Let  $I \subseteq \{1, \dots, \ell\}$  contain  $d_r + 1$  indices  $i$  s. t.  $s'[i] \neq s_j[i]$  (or all indices if  $\mathfrak{d}_H(s_j, s') \leq d_r$ );
11  for  $i \in I$  do
12     $s'' = s'$ ,  $s''[i] = s_j[i]$ ;
13     $sol =$  SOLVE CLOSEST STRING-WO( $S', t, s'', d - 1$ );
14    if  $sol \neq \nabla$  then return  $sol$ ;
15 return  $\nabla$ ;
```

226 $d_r + 1$ recursive calls with parameters $(d - 1, t)$, and one recursive call with parameters $(d, t - 1)$.
 227 By induction, the complexity of this step is $O^*((d_r + 1)(d_r + 1)^{d-1}2^{d+t-1} + (d_r + 1)d2^{d+t-1}) =$
 228 $O^*((d_r + 1)^d 2^{d+t})$. ◀ (Claim 1)

229 A tuple (S', t', s', d) is *valid* if $|S'| - t' = |S| - t$, there exists an optimal solution (S^*, s^*) for
 230 which $S^* \subseteq S'$, $|S^*| = |S'| - t'$, $\mathfrak{d}_H(s', s^*) \leq d$, and s' is a lower bound for s^* . A call of the
 231 algorithm is *valid* if its parameters form a valid tuple, its *witness* is the pair (S^*, s^*) .

232 *Claim 2:* Any valid call to SOLVE CLOSEST STRING-WO either directly returns a solution or
 233 performs at least one recursive valid call.

234 *Proof of Claim 2:* Let $S' \subseteq \Sigma^\ell$, $t' \geq 0$, $s' \in (\Sigma \cup \{\diamond\})^\ell$, and $d \geq 0$. Consider the call to
 235 SOLVE CLOSEST STRING-WO(S', t', s', d). Assume it is valid, with witness (S^*, s^*) .

236 **Case 1:** If $d = t' = 0$, then $s^* = s'$ and $S^* = S'$. The completion s'' of s' is exactly s' , and
 237 since (S', s') is a solution, it satisfies the conditions of Line 3 and is returned on Line 3.

238 **Case 2:** If $t' = 0$ and $\forall s \in S' : \mathfrak{d}_H(s, s') \leq d_r$. Then $S^* = S'$ and s' is a lower bound for s^* .
 239 Let s'' be the completion of s' . We show that $\mathfrak{s}_H(s'', S') \leq \mathfrak{s}_H(s^*, S') \leq d_s$. Indeed, consider
 240 any column i with $s''[i] \neq s^*[i]$. Either $s'[i] = \diamond$, in which case $s''[i]$ is the majority character
 241 for column i of S' , or $s'[i] \neq \diamond$, in which case by the definition of lower bound, i is not a
 242 disputed column and $s'[i] = s''[i]$ contains the only frequent character of this column, which
 243 is the majority character for S' . In both cases, $s''[i]$ is a majority character for S' in any
 244 column where it differs from s^* ; thus, it satisfies the upper-bound on the distance sum. Since
 245 it also satisfies the distance radius (by the case hypothesis: $\mathfrak{d}_H(s, s'') \leq \mathfrak{d}_H(s, s') \leq d_r$ for all
 246 $s \in S'$), it satisfies the conditions of Line 3; thus, solution (S', s'') is returned on Line 3.

247 In the following cases, we can thus assume that the algorithm reaches Line 5. Indeed,
 248 if it returns on Line 3 then it returns a solution, and if it returns on Line 4 then we have
 249 $d = t' = 0$, which is dealt in Case 1 above (the algorithm may not return on this line when it
 250 has a valid input). We can thus define s_j to be the string selected in Line 5.

251 **Case 3:** $s_j \in S' \setminus S^*$. Then in particular $t' > 0$; and since $S^* \subseteq S' \setminus \{s_j\}$, the recursive call
 252 in Line 7 is valid, with the same witness (S^*, s^*) .

253 **Case 4:** $s_j \in S^*$, $d = 0$ and $t' > 0$. Then $s' = s^*$, let s'_j be any string of $S' \setminus S^*$, and
 254 $S^+ = S^* \setminus \{s'_j\} \cup \{s_j\}$. Then the pair (S^+, s^*) is a solution, since $d_H(s^*, s'_j) \leq d_H(s^*, s_j)$ by
 255 definition of s_j . Thus the recursive call on Line 7 is valid, with witness (S^+, s^*) .

256 **Case 5:** $s_j \in S^*$, $d > 0$ and $d_H(s_j, s') > d_r$. Consider the set I defined in Line 10. I has
 257 size $d_r + 1$, hence there exists $i_0 \in I$ such that $s_j[i_0] = s^*[i_0]$. Then the recursive call with
 258 parameters $(S', t, s'', d - 1)$ in Line 13 with $i = i_0$ is valid with the same witness (S^*, s^*) .
 259 Indeed, s'' is obtained from s' by setting $s''[i_0] = s^*[i_0] \neq s'[i_0]$, hence, all mismatches
 260 between s'' and s^* already exist between s' and s^* , which implies that s'' is still a lower
 261 bound for s^* . Moreover, $d_H(s'', s^*) = d_H(s', s^*) - 1 \leq d - 1$.

262 From now on, we can assume that $d > 0$ and $t' > 0$. Indeed, $d = 0$ is dealt with in cases
 263 1, 3 and 4, and $t' = 0, d > 0$ is dealt with in cases 2 and 5. Moreover, with cases 3 and 5, we
 264 can assume that $s_j \in S^*$ and $d_H(s_j, s') \leq d_r$ (i.e. $d_H(s, s') \leq d_r$ for all $s \in S^*$).

265 **Case 6:** There exists i_0 such that $s_j[i_0] = s^*[i_0] \neq s'[i_0]$. Then again consider the set I
 266 defined in Line 10. Since $d_H(s_j, s') \leq d_r$, we have $i_0 \in I$, and, with the same argument as in
 267 Case 5, there is a valid recursive call in Line 13 when $i = i_0$.

268 **Case 7:** For all i , $s_j[i] \neq s'[i] \Rightarrow s_j[i] \neq s^*[i]$. In this case no character from s_j can be
 269 used to improve our current solution, so the character switching procedure Line 13 will not
 270 improve the solution, but still s_j is part of our witness set S^* , so it is not clear a priori that
 271 we can remove s_j from our current solution, i.e. that the recursive call on Line 7 is valid.

272 We handle this situation as follows. Let s^+ be obtained from s' by filling the \diamond -positions
 273 of s' with the corresponding symbols of s^* . We now show that (S^*, s^+) is a solution. To this
 274 end, let $s \in S^*$. For every i , $1 \leq i \leq \ell$, if $s[i] \neq s^+[i]$, then either $s'[i] = \diamond$ or $s'[i] \in \Sigma$ with
 275 $s'[i] = s^+[i]$. In both cases, we have $s[i] \neq s'[i]$, which implies $d_H(s, s^+) \leq d_H(s, s') \leq d_r$, i.e.,
 276 the radius is satisfied. Regarding the distance sum, we note that if $s^+[i] \neq s^*[i]$, then, since
 277 occurrences of \diamond of s' have been replaced by the corresponding symbol from s^* , $s'[i] \neq \diamond$,
 278 which, by the definition of lower bound, implies that $s^+[i] = s'[i]$ is the majority character
 279 for column i of S^* . Consequently, $\sum_{s \in S^*} d_H(s^+[i], s[i]) \leq \sum_{s \in S^*} d_H(s^*[i], s[i])$, which implies
 280 $s_H(s^+, S^*) \leq s_H(s^*, S^*) \leq d_s$.

281 Having defined a new solution string s^+ (with respect to S^*), we now prove that s^+ is also a
 282 solution string with respect to $S^+ = (S^* \setminus \{s_j\}) \cup \{s'_j\}$, where s'_j is any string of $S' \setminus S^*$. To this
 283 end, we prove that $d_H(s'_j, s^+) \leq d_H(s_j, s^+)$; together with the fact that $d_H(s'_j, s') \leq d_r$, this
 284 implies that (S^+, s^+) is a solution. For two strings $s_1, s_2 \in \Sigma^\ell$, let $d_\diamond(s_1, s_2)$ be the number of
 285 mismatches between s_1 and s_2 at positions i such that $s'[i] = \diamond$, and $d_\Sigma(s_1, s_2)$ be the number of
 286 mismatches at other positions. Clearly $d_H(s_1, s_2) = d_\diamond(s_1, s_2) + d_\Sigma(s_1, s_2)$. Comparing
 287 strings s_j and s'_j to s' , we have $d_\diamond(s_j, s') = d_\diamond(s'_j, s')$ (both distances are equal to the number
 288 of occurrences of \diamond in s'). Since $d_H(s_j, s')$ is maximal, we have $d_\Sigma(s'_j, s') \leq d_\Sigma(s_j, s')$. Consider
 289 now s^+ . Since s^+ is equal to s' in every non- \diamond characters, we have $d_\Sigma(s'_j, s^+) \leq d_\Sigma(s_j, s^+)$.
 290 Finally, for any i such that $s'[i] = \diamond$, by hypothesis of this case we have $s_j[i] \neq s^*[i] = s^+[i]$,
 291 hence $d_\diamond(s_j, s^+)$ is equal to the number of occurrences of \diamond in s' , which is an upper bound
 292 for $d_\diamond(s'_j, s^+)$. Overall, $d(s'_j, s^+) \leq d(s_j, s^+)$, and (S^+, s^+) is a solution.

293 Thus, (S^+, s^+) is a solution such that $S^+ \subseteq S' \setminus \{s_j\}$, s' is a lower bound for s^+ , and
 294 $d_H(s', s^+) \leq d$, hence the recursive call in Line 7 is valid. ◀ (Claim 2)

295 It follows from the claim above that any valid call to SOLVE CLOSEST STRING-WO returns
 296 a solution. Indeed, if it does not directly return a solution, then it receives a solution of a
 297 more constrained instance from a valid recursive call, which is returned on Line 8 or 14.

298 *Claim 3:* Let s' be the majority string for S where for every disputed column i , $s'[i] = \diamond$.
 299 Then SOLVE CLOSEST STRING-WO($S, t, s', 2d_r + D$) is a valid call.

300 *Proof of Claim 3:* Consider a solution (S^*, s^*) . We need to check whether $d_H(s^*, s') \leq 2d_r + D$,
 301 and whether s' is a lower bound of s^* . The fact that s' is a lower bound follows from the
 302 definition, since \diamond is selected in every disputed column, and the majority character is selected
 303 in the other columns. String s^* can be seen as a solution of (r, s)-CLOSEST STRING over
 304 S^*, d_r, d_s , thus, we can use Lemma 2: the distance between s^* and the majority string of S^*
 305 is at most $2d_r$. Hence there are at most $2d_r$ mismatches between s' and s^* in non-disputed
 306 columns (since in those columns, the majority characters are identical in S and S^*). Adding
 307 the D mismatches from disputed columns, we get the $2d_r + D$ upper bound. ◀ (Claim
 308 3) ▶

309 2.2 The (r)- and (s)-Variants of Closest String-wo

310 In [?], the fixed-parameter tractability of (r)-CLOSEST STRING-wo w. r. t. parameter k and
 311 w. r. t. parameters $(|\Sigma|, d_r, k - t)$ are reported as open problems. Since Theorem 5 also applies
 312 to (r)-CLOSEST STRING-wo, the only open cases left for the (r)-variant are the following:

313 ▶ **Open Problem 8.** *What is the fixed-parameter tractability of (r)-CLOSEST STRING-wo*
 314 *with respect to $(|\Sigma|, k - t)$, $(|\Sigma|, d_r)$ and $(|\Sigma|, d_r, k - t)$?*

315 Next, we consider the (s)-variant of CLOSEST STRING-wo. We recall that replacing
 316 the radius bound by a bound on the distance sum turns (r)-CLOSEST STRING into a triv-
 317 ial problem, while (s)-CLOSEST SUBSTRING remains hard. The next result shows that
 318 CLOSEST STRING-wo behaves like CLOSEST SUBSTRING in this regard. For the proof, we
 319 use MULTI-COLOURED CLIQUE (which is W[1]-hard, see [?]), which is identical to the
 320 standard parameterisation of CLIQUE, but the input graph $G = (V, E)$ has a partition
 321 $V = V_1 \cup \dots \cup V_{k_c}$, such that every V_i , $1 \leq i \leq k_c$, is an independent set (we denote the
 322 parameter by k_c to avoid confusion with the number of input strings k).

323 ▶ **Theorem 9.** (s)-CLOSEST STRING-wo($d_s, \ell, k - t$) is W[1]-hard.

324 **Proof.** Let $G = (V_1 \cup \dots \cup V_{k_c}, E)$ be a MULTI-COLOURED CLIQUE-instance. We assume
 325 that, for some $q \in \mathbb{N}$, $V_i = \{v_{i,1}, v_{i,2}, \dots, v_{i,q}\}$, $1 \leq i \leq k_c$, i. e., each vertex has an index
 326 depending on its colour-class and its rank within its colour-class. Let $\Sigma = V \cup \Gamma$, where
 327 Γ is some alphabet with $|\Gamma| = |E|(k_c - 2)$. For every $e = (v_{i,j}, v_{i',j'}) \in E$, let $s_e \in \Sigma^{k_c}$
 328 with $s_e[i] = v_{i,j}$, $s_e[i'] = v_{i',j'}$ and all other non-defined positions are filled with symbols
 329 from Γ such that each $x \in \Gamma$ has exactly one occurrence in the strings s_e , $e \in E$. We set
 330 $S = \{s_e \mid e \in E\}$, $t = |E| - \binom{k_c}{2}$ (i. e., the number of inliers is $\binom{k_c}{2}$) and $d_s = \binom{k_c}{2}(k_c - 2)$.

331 Let K be a clique of G of size k_c , let $s \in \Sigma^{k_c}$ be defined by $\{s[i]\} = K \cap V_i$, $1 \leq i \leq k_c$, and
 332 let $S' = \{s_e \mid e \subseteq K\}$. Since $d_H(s, s') = k_c - 2$, for every $s' \in S'$, $s_H(s, S') = d_s$. Consequently,
 333 S' and s is a solution for the (s)-CLOSEST STRING-wo-instance S, t, d_s .

334 Now let $s \in \Sigma^{k_c}$ and $S' \subseteq S$ with $|S'| = \binom{k_c}{2}$ be a solution for the (s)-CLOSEST STRING-wo-
 335 instance S, t, d_s . If, for some $s'_1 \in S'$, $d_H(s'_1, s) \geq k_c - 1$, then there is an $s'_2 \in S'$ with
 336 $d_H(s'_2, s) \leq k_c - 3$. Thus, for some i , $1 \leq i \leq k_c$, $s[i] = s'_2[i]$ and $s'_2[i] \in \Gamma$, which implies that
 337 replacing $s[i]$ by $s'_1[i]$ does not increase $s_H(s, S')$. Moreover, after this modification, $d_H(s'_1, s)$
 338 has decreased by 1, while $d_H(s'_2, s) \leq k_c - 2$. By repeating such operations, we can transform
 339 s such that $d_H(s', s) \leq k_c - 2$, $s' \in S'$. Next, assume that, for some i , $1 \leq i \leq k_c$, there is an
 340 $S'' \subseteq S'$ with $|S''| = k_c$ and, for every $s' \in S''$, $s[i] = s'[i]$. Since $d_H(s', s) \leq k_c - 2$ for every
 341 $s' \in S''$, pigeon-hole principle implies that there are $s'_1, s'_2 \in S''$ with $s'_1[i'] = s'_2[i'] = s[i']$, for
 342 some i' , $1 \leq i' \leq k_c$, and $i' \neq i$, which, by the structure of the strings of S , is a contradiction.
 343 Consequently, for every i , $1 \leq i \leq k_c$, s matches with at most $k_c - 1$ strings from S' at
 344 position i . Since there are at least $2\binom{k_c}{2} = k_c(k_c - 1)$ matches, we conclude that, for every

k	t	$ \Sigma $	ℓ	d_r	d_s	$k-t$	Result	Note/Ref.
p	–	–	–	–	–	–	FPT	Thm. 5, Open Prob. in [?]
–	0	2	–	–	–	–	NP-hard	even for d_r -var., but P for d_s -var.
–	p	–	p	–	–	–	FPT	$d_r \leq \ell$
–	p	–	–	p	–	–	FPT	Thm. 7, and [?] for d_r -var.
–	p	–	–	–	p	–	FPT	Thm. 6
–	p	–	–	–	–	p	FPT	$k = t + (k - t)$
–	–	p	p	–	–	–	FPT	trivial
–	–	p	–	*	*	*	Open	param. $ \Sigma $ and some of $d_r, d_s, k-t$
–	–	–	p	p	p	p	W[1]-hard	even for d_r -var. [?] and d_s -var. (Thm. 9)

■ **Table 2** Results for (r, s) -CLOSEST STRING-WO, including (r) - and (s) -variants.

345 i , $1 \leq i \leq k_c$, $s[i]$ matches exactly $k_c - 1$ times with the i^{th} position of a string from
346 S' . Hence, $s[i] \in V_i$, $1 \leq i \leq k_c$, i. e., $s = v_{1,r_1}v_{2,r_2} \dots v_{k_c,r_{k_c}}$, for some $r_j \in \{1, 2, \dots, q\}$,
347 $1 \leq j \leq k_c$. Let $K = \{v_{1,r_1}, v_{2,r_2}, \dots, v_{k_c,r_{k_c}}\}$. For every $s' \in S'$, by definition of the strings
348 s_e , we have $d_H(s, s') \geq k_c - 2$, combining with the lower-bound proved earlier, we conclude
349 $d_H(s, s') = k_c - 2$, for every $s' \in S$. Now let $e = (v_{i,j}, v_{i',j'}) \in E$ be such that $s_e \in S'$. From
350 $d_H(s, s_e) = k_c - 2$ it follows that $s[i] = v_{i,j}$ and $s[i'] = v_{i',j'}$, which implies $e \subseteq K$. Since
351 $|S| = \binom{k_c}{2}$, there are $\binom{k_c}{2}$ edges connecting vertices from K ; thus, K is a clique. ◀

352 Setting $d_r = k_c - 2$ instead of $d_s = \binom{k_c}{2}(k_c - 2)$ in the reduction of Theorem 9 leads to a
353 simpler proof for the W[1]-hardness of (r) -CLOSEST STRING-WO($d_r, \ell, k - t$) shown in [?] (on
354 the other hand, the reduction of [?] does not work for (s) -CLOSEST STRING-WO($d_s, \ell, k - t$)).
355 The results obtained in this section are summarized in Table 2.

356 **3 Closest Substring**

357 In this section, we consider the problem (r, s) -CLOSEST SUBSTRING and, as done in Section 2
358 for (r, s) -CLOSEST STRING, we classify all parameterisations of (r, s) -CLOSEST SUBSTRING
359 (and its (r) - and (s) -variants) with respect to the parameters ℓ, k, m, d_r, d_s and $|\Sigma|$ into
360 either fixed-parameter tractable or W[1]-hard. Of course, many of those questions are already
361 solved in the literature, but, unlike for (r, s) -CLOSEST STRING, not all cases of the (r) - and
362 (s) -variants are settled, i. e., the status of (s) -CLOSEST SUBSTRING(ℓ) is unknown, which is
363 mentioned as open problem in [?]. We shall first close this gap by defining a reduction from
364 MULTI-COLOURED CLIQUE to (s) -CLOSEST SUBSTRING.

365 Let $G = (V_1 \cup \dots \cup V_{k_c}, E)$ be a MULTI-COLOURED CLIQUE-instance. We assume that, for
366 some $q \in \mathbb{N}$, $V_i = \{v_{i,1}, v_{i,2}, \dots, v_{i,q}\}$, $1 \leq i \leq k_c$, i. e., each vertex has an index depending on
367 its colour-class and its rank within its colour-class. Let $\Sigma = V \cup \{\$, \diamond\}$. For every j , $1 \leq j \leq q$,
368 we list all j^{th} elements of the colour-classes as a string $\mathcal{V}_j = \$v_{1,j}v_{2,j} \dots v_{k_c,j}$. For every edge
369 $e = (v_{i,j}, v_{i',j'})$ with $i < i'$, we define a string $\mathcal{E}_e = \$\diamond^i v_{i,j} \diamond^{i'-i-1} v_{i',j'} \diamond^{k_c-i'-1}$. Note that
370 $\mathcal{E}_e = \$\diamond \mathcal{E}'_e$, where $|\mathcal{E}'_e| = k_c$, the positions i and i' of \mathcal{E}'_e are $v_{i,j}$ and $v_{i',j'}$, respectively, and all
371 remaining positions are \diamond . The (s) -CLOSEST SUBSTRING-instance is now defined as follows.
372 Let S contain $N = |E|(k_c + 2) + 1$ occurrences of each \mathcal{V}_j , $1 \leq j \leq q$, and one occurrence of
373 each \mathcal{E}_e , $e \in E$, and let $m = k_c + 1$. We note that $\ell = k_c + 2$. See Figure 2 for an example.

374 In the following, we extend the notation of *radius optimal* and *distance sum optimal*
375 to sets $S \subseteq \Sigma^{\leq \ell}$ and strings $s \in \Sigma^m$ in the natural way by taking all sets S' of length- m
376 substrings of the string in S into account. The next lemma shows that distance sum optimal
377 strings (with respect to S and m) are basically lists of vertices from each colour-class.

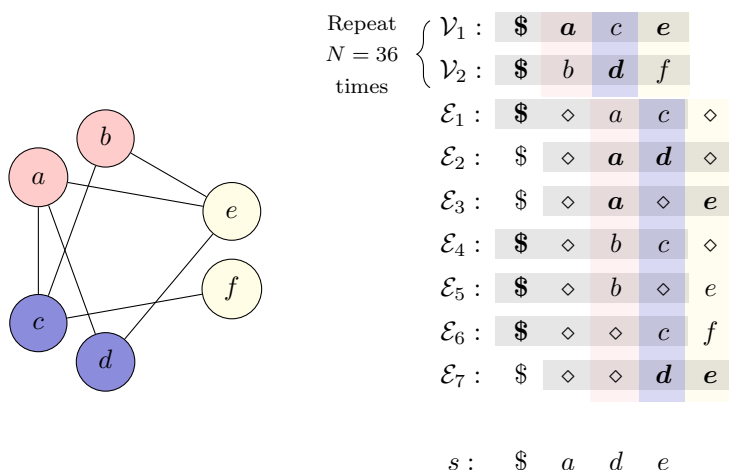


Figure 2 Illustration of the parameterized reduction from a MULTI-COLOURED CLIQUE-instance to (s)-CLOSEST SUBSTRING. The colour-classes of the graph are $V_1 = \{a, b\}$ (red), $V_2 = \{c, d\}$ (blue) and $V_3 = \{e, f\}$ (yellow), the occurrences of symbols from V in the strings \mathcal{V}_j and \mathcal{E}_i are coloured according to their colour-classes. The string $s = \$ade$ is an optimal solution with respect to the substrings emphasised with grey background (positions producing a match are in bold). Note that vertices $\{a, d, e\}$ form a clique in G .

378 **► Lemma 10 (*)**. *If $s \in \Sigma^{k+1}$ is distance sum optimal w. r. t. S , then $s \in \{\$\} \cdot V_1 \cdot V_2 \cdot \dots \cdot V_k$.*

379 Now let s be distance sum optimal with respect to S and m . From Lemma 10, we can
 380 conclude that $s = \$v_{1,r_1}v_{2,r_2} \dots v_{k_c,r_{k_c}}$, for some $r_j \in \{1, 2, \dots, q\}$, $1 \leq j \leq k_c$. Let K be the
 381 corresponding set of vertices, i. e., $K = \{v_{1,r_1}, v_{2,r_2}, \dots, v_{k_c,r_{k_c}}\}$.

382 **► Lemma 11 (*)**. *Let $e \in E$. The optimal distance between s and a length- $(k_c + 1)$ substring
 383 of \mathcal{E}_e is $k_c - 1$ if $e \subseteq K$, and k_c otherwise.*

384 Using the lemmas from above, we can now show the correctness of the reduction.

385 **► Theorem 12**. *(s)-CLOSEST SUBSTRING(ℓ, m) is $W[1]$ -hard.*

386 **Proof.** Let $s \in \Sigma^{k_c+1}$ be distance sum optimal with respect to S and m , and let K
 387 be the corresponding set of vertices. We first note that the total distance from s to
 388 the N copies of the strings \mathcal{V}_j , $1 \leq j \leq q$, is exactly Nqk_c . According to Lemma 11,
 389 for every $e \in E$, the optimal distance sum between s and the respective substring of
 390 \mathcal{E}_e is $k_c - 1$ if $e \subseteq K$, and k_c otherwise. Hence, the total distance sum from s to the
 391 respective substrings of \mathcal{E}_e , $e \in E$, is $|E|k_c - r$, where $r = |\{e \in E \mid e \subseteq K\}|$, and the
 392 total distance sum between s and S is therefore $Nqk_c + |E|k_c - r$. This implies that the
 393 distance sum between s and S is $Nqk_c + |E|k_c - \frac{k_c(k_c-1)}{2}$ if and only if $r = \frac{k_c(k_c-1)}{2}$ if and
 394 only if K is a clique of size k_c . Consequently, the above reduction, with the addition of
 395 $d_s = Nqk_c + |E|k_c - \frac{k_c(k_c-1)}{2}$, is a parameterised reduction from MULTI-COLOURED CLIQUE
 396 to (s)-CLOSEST SUBSTRING(ℓ, m). ◀

397 As illustrated by Table 3, Theorem 12 together with known results from the literature
 398 completely settle the parameterised complexity of (s)-CLOSEST SUBSTRING.

399 Moving on to the problem (r,s)-CLOSEST SUBSTRING, we first observe that reducing
 400 (s)-CLOSEST SUBSTRING to (r,s)-CLOSEST SUBSTRING by setting $d_r = m$ is a parameterised

ℓ	k	m	d_s	$ \Sigma $	Result	Reference
–	–	p	–	p	FPT	trivial
p	–	–	–	p	FPT	[?]
p	p	–	–	–	FPT	[?]
p	–	–	p	–	FPT	[?]
–	–	–	p	p	FPT	[?]
–	p	–	–	2	W[1]-hard	[?]
–	p	p	p	–	W[1]-hard	[?]
p	–	p	–	–	W[1]-hard	Thm. 12

■ **Table 3** Results for (s)-CLOSEST SUBSTRING.

ℓ	k	m	d_r	d_s	$ \Sigma $	Result	Reference
–	–	p	–	–	p	FPT	Thm. 14
p	p	–	–	–	–	FPT	Thm. 14
p	–	–	–	p	–	FPT	Thm. 14
p	–	–	–	–	p	FPT	Thm. 14
p	–	p	p	–	–	W[1]-hard	Cor. 13, Open Prob. in [?]
–	p	–	p	p	p	W[1]-hard	[?]
–	p	p	p	p	–	W[1]-hard	[?]

■ **Table 4** Results for (r, s)-CLOSEST SUBSTRING.

401 reduction from (s)-CLOSEST SUBSTRING(ℓ, m) to (r, s)-CLOSEST SUBSTRING(ℓ, m, d_r), which
 402 implies the following corollary:

403 ► **Corollary 13.** (r, s)-CLOSEST SUBSTRING(ℓ, m, d_r) is W[1]-hard.

404 Next, we consider several fixed-parameter tractable variants of (r, s)-CLOSEST SUBSTRING.

405 ► **Theorem 14 (*).** (r, s)-CLOSEST SUBSTRING(x) \in FPT, for every $x \in \{(m, |\Sigma|), (\ell, k),$
 406 $(\ell, |\Sigma|), (\ell, d_s)\}$.

407 It remains to observe that all remaining parameterisations of (r, s)-CLOSEST SUBSTRING
 408 are W[1]-hard. More precisely, it is known that (r)-CLOSEST SUBSTRING is W[1]-hard for
 409 parameterisations ($k, d_r, |\Sigma|$) (see [?]) and (k, m, d_r) (see [?]). Hence, the obvious reduc-
 410 tion from (r)-CLOSEST SUBSTRING to (r, s)-CLOSEST SUBSTRING, i. e., setting $d_s = k d_r$,
 411 shows that (r, s)-CLOSEST SUBSTRING is W[1]-hard for parameterisations ($k, d_r, d_s, |\Sigma|$) and
 412 (k, m, d_r, d_s). As can be checked with the help of Table 4, this now classifies all parameterised
 413 variants of (r, s)-CLOSEST SUBSTRING.

414 **4 Kernelisation**

415 Neither (r)-CLOSEST STRING($d_r, \ell, |\Sigma|$) nor (r)-CLOSEST SUBSTRING(k, m, d_r) admit poly-
 416 nomial kernels unless $\text{coNP} \subseteq \text{NP/Poly}$ (see [?]), and (r)-CLOSEST STRING(k, d_r) has a kernel
 417 of size $O(k^2 d_r \log k)$ (see [?]). From these results, we can conclude the following:

418 ► **Proposition 15 (*).**

- 419 ■ (r, s) -CLOSEST STRING($d_r, \ell, |\Sigma|$) has no polynomial kernel unless $\text{coNP} \subseteq \text{NP/Poly}$.
 420 ■ (r, s) -CLOSEST STRING(k, d_r) has a kernel of size $O(k^2 d_r \log k)$.
 421 ■ (r, s) -CLOSEST STRING(d_s) has a kernel of size $O((d_s)^3 \log d_s)$.

422 This only leaves the case open, where only k (or k and $|\Sigma|$, which, due to the dependency
 423 $|\Sigma| \leq k$ (see [?]), is the same question) is a parameter (regarding this case, note that for
 424 (r) -CLOSEST STRING(k) no combinatorial kernel or combinatorial FPT-algorithm is known).

425 ► **Proposition 16 (*)**.

- 426 ■ (r, s) -CLOSEST SUBSTRING($k, m, d_r, d_s, |\Sigma|$) has no polynomial kernel unless $\text{coNP} \subseteq$
 427 NP/Poly .
 428 ■ (r, s) -CLOSEST SUBSTRING(ℓ, k) and (r, s) -CLOSEST SUBSTRING(ℓ, d_s) have kernels of
 429 size $O(\ell k)$ and $O(\ell d_s)$, respectively.

430 This almost settles the (r, s) -variant, for which only the parameterisation $(\ell, |\Sigma|)$ is open.
 431 For the (r) -variant, the parameterisations ℓ , (ℓ, d_r) and $(\ell, |\Sigma|)$, and for the (s) -variant, the
 432 parameterisations $(m, |\Sigma|)$ and $(d_s, |\Sigma|)$ are open.

433 For (r) -CLOSEST STRING-WO no kernelisation lower bounds are known so far. However,
 434 the following can be concluded from [?]:

435 ► **Proposition 17 (*)**. (r) -CLOSEST STRING-WO($d_r, \ell, t, |\Sigma|$) has no polynomial kernel unless
 436 $\text{coNP} \subseteq \text{NP/Poly}$.

437 By a cross-composition³ from (r) -CLOSEST STRING into (r) -CLOSEST STRING-WO, we
 438 can rule out a polynomial kernel for the parameterisation $(d_r, d_s, \ell, (k - t), |\Sigma|)$.

439 To this end, we define a polynomial equivalence relation \sim over the (r) -CLOSEST STRING-
 440 instances as follows. For $j \in \{1, 2\}$, let $S_j = \{s_{j,i} \mid 1 \leq i \leq k_j\} \subseteq \Sigma^{\ell_j}$ and $d_{r,j} \in \mathbb{N}$. Then
 441 $(S_1, d_{r,1}) \sim (S_2, d_{r,2})$ if $k_1 = k_2$, $\ell_1 = \ell_2$ and $d_{r,1} = d_{r,2}$. Now let $(S_1, d_r), (S_2, d_r), \dots, (S_q, d_r)$
 442 be \sim -equivalent (r) -CLOSEST STRING-instances, where, for the sake of convenience, $S_i =$
 443 $\{s_{i,1}, s_{i,2}, \dots, s_{i,k}\} \subseteq \Sigma^\ell$, $1 \leq i \leq q$. For every i , $1 \leq i \leq q$, let B_i denote the binary
 444 representation of i with exactly $\lceil \log(q) \rceil$ bits, and let $C_i = (B_i)^{2^{d_r+1}}$. Moreover, for every i ,
 445 $1 \leq i \leq q$, let $S'_i = \{s'_{i,1}, s'_{i,2}, \dots, s'_{i,k}\}$, where, for every j , $1 \leq j \leq k$, $s'_{i,j} = s_{i,j} C_i$. Finally,
 446 let the (r, s) -CLOSEST STRING-WO-instance be (S', d'_r, d'_s, t) with $S' = \bigcup_{i=1}^q S'_i$, $d'_r = d_r$,
 447 $d'_s = k d_r$ and $t = (q - 1)k$.

448 ► **Theorem 18 (*)**. (r, s) -CLOSEST STRING-WO($d_r, d_s, \ell, (k - t), |\Sigma|$) does not admit a poly-
 449 nomial kernel unless $\text{coNP} \subseteq \text{NP/Poly}$.

450 5 Conclusions

451 The parameterised complexity of the (r) -, (s) - and general variant of CLOSEST STRING
 452 and CLOSEST SUBSTRING with respect to $\ell, k, m, d_r, d_s, |\Sigma|$ is now completely settled. For
 453 (r, s) -CLOSEST SUBSTRING, where positive results are less abundant, it might be worthwhile to
 454 identify other parameters that yield fixed-parameter tractability. For (r, s) -CLOSEST STRING,
 455 it should be pointed out that the FPT-algorithms with respect to k are based on ILP and are
 456 most likely practically not relevant; direct combinatorial FPT-algorithms are still unknown.
 457 For the outlier variant of (r, s) -CLOSEST STRING, many cases are left open, most prominently,
 458 the ones with $|\Sigma|$ as parameter, and we expect those to be challenging. Moreover, for several
 459 FPT-variants, the existence of polynomial kernels is not yet answered.

³ For the technique of cross-composition, see Bodlaender et al. [?].