



HAL
open science

Tight Hardness Results for Consensus Problems on Circular Strings and Time Series

Laurent Bulteau, Vincent Froese, Rolf Niedermeier

► **To cite this version:**

Laurent Bulteau, Vincent Froese, Rolf Niedermeier. Tight Hardness Results for Consensus Problems on Circular Strings and Time Series. *SIAM Journal on Discrete Mathematics*, 2020, 34 (3), pp.1854-1883. 10.1137/19M1255781 . hal-03031331

HAL Id: hal-03031331

<https://hal.science/hal-03031331>

Submitted on 30 Nov 2020

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.

TIGHT HARDNESS RESULTS FOR CONSENSUS PROBLEMS ON CIRCULAR STRINGS AND TIME SERIES*

LAURENT BULTEAU[†], VINCENT FROESE[‡], AND ROLF NIEDERMEIER[‡]

Abstract. Consensus problems for strings and sequences appear in numerous application contexts, ranging among bioinformatics, data mining, and machine learning. Closing some gaps in the literature, we show that several fundamental problems in this context are NP- and W[1]-hard, and that the known (including some brute-force) algorithms are close to optimality assuming the Exponential Time Hypothesis. Among our main contributions is to settle the complexity status of computing a mean in dynamic time warping spaces which, as pointed out by Brill et al. [DMKD 2019], suffered from many unproven or false assumptions in the literature. We prove this problem to be NP-hard and additionally show that a recent dynamic programming algorithm is essentially optimal. In this context, we study a broad family of circular string alignment problems. This family also serves as a key for our hardness reductions, and it is of independent (practical) interest in molecular biology. In particular, we show tight hardness and running time lower bounds for CIRCULAR CONSENSUS STRING; notably, the corresponding non-circular version is easily linear-time solvable.

Key words. Circular String Alignment, Time Series Averaging, Dynamic Time Warping, Fine-Grained Complexity and Reductions, Lower Bounds, Parameterized Complexity, Exponential Time Hypothesis

AMS subject classifications. 68Q17, 68T10, 92D20

1. Introduction. Consensus problems appear in many contexts of stringology and time series analysis, including applications in bioinformatics, data mining, machine learning, and speech recognition. Roughly speaking, given a set of input sequences, the goal is to find a consensus sequence that minimizes the “distance” (according to some specified distance measure) to the input sequences. Classic problems in this context are the NP-hard CLOSEST STRING [15, 25, 24, 18] (where the goal is to find a “closest string” that minimizes the maximum Hamming distance to a set of equal-length strings) or the more general CLOSEST SUBSTRING [13, 26]. Notably, the variant of CLOSEST STRING where one minimizes the sum of Hamming distances instead of the maximum distance is easily solvable in linear time.

In this work, we settle the computational complexity of prominent consensus problems on circular strings and time series. Despite their great importance in many applications, and a correspondingly rich set of heuristic solution strategies used in practice, to date, it has been unknown whether these problems are polynomial-time solvable or NP-hard. We prove their hardness, including also “tight” parameterized and fine-grained complexity results, thus justifying the massive use of heuristic solution strategies in real-world applications.

On the route to determining the complexity of exact mean computation in dynamic time warping spaces, a fundamental consensus problem in the context of time series analysis [33]¹, we first study a fairly general alignment problem² for circular

*Submitted to the editors DATE.

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

[‡]Technische Universität Berlin, Faculty IV, Algorithmics and Computational Complexity, Berlin, Germany (vincent.froese@tu-berlin.de, rolf.niedermeier@tu-berlin.de).

¹As of May 2020, according to Google Scholar the work by Petitjean et al. [33], who developed one of the most prominent heuristics for this problem, has already been cited around 460 times since 2011.

²Particularly from the viewpoint of applications in bioinformatics, consensus string problems can also be interpreted as alignment problems [23].

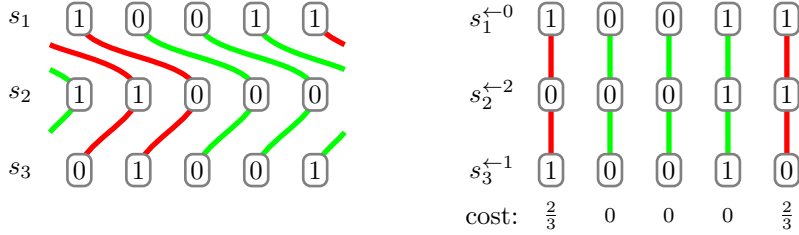


FIG. 1. An instance of σ -MSCS with three binary input strings, and an optimal multiple circular shift $\Delta = (0, 2, 1)$, using the sum of squared distances from the mean (σ) as a cost function. Columns of Δ are indicated with dark (red) or light (green) lines, depending on their cost. For example, column 1 with values $(1, 0, 1)$ has mean $\frac{2}{3}$ and cost $(\frac{1}{3})^2 + (\frac{2}{3})^2 + (\frac{1}{3})^2 = \frac{2}{3}$. The overall cost is $\frac{4}{3}$.

40 strings called MULTIPLE STRING CIRCULAR SHIFT (WITH COST f). Based on its
 41 analysis, we will also derive our results for two further, more specific problems. Given
 42 a set of input strings over a fixed alphabet Σ and a local cost function $f: \Sigma^* \rightarrow \mathbb{Q}$,
 43 the goal in MULTIPLE STRING CIRCULAR SHIFT (WITH COST f) (abbreviated by
 44 f -MSCS) is to find a cyclic shift of each input string such that the shifted strings
 45 “align well” in terms of the sum of local costs.³

f -MSCS

Input: A list of k strings $s_1, \dots, s_k \in \Sigma^n$ of length n and $c \in \mathbb{Q}$.

Question: Is there a multiple circular shift $\Delta = (\delta_1, \dots, \delta_k) \in \mathbb{N}^k$ with
 $\text{cost}_f(\Delta) := \sum_{i=1}^n f((s_1^{\leftarrow \delta_1}[i], \dots, s_k^{\leftarrow \delta_k}[i])) \leq c$?

47 Here, $s^{\leftarrow \delta}$ denotes a circular shift of s by δ (see Section 2 for details). See Fig-
 48 ure 1 for an example. We separately study the special case CIRCULAR CONSENSUS
 49 STRING for a binary alphabet, where the cost function $f: \{0, 1\}^* \rightarrow \mathbb{N}$ is defined as
 50 $f((x_1, \dots, x_k)) := \min\{\sum_{i=1}^k x_i, k - \sum_{i=1}^k x_i\}$. This corresponds to minimizing the
 51 sum of Hamming distances (not the maximum Hamming distance as in CLOSEST
 52 STRING). As we will show, allowing circular shifts makes consensus string problems
 53 much harder to solve.

54 Multiple circular string (sequence) alignment problems have been considered in
 55 different variations in bioinformatics, where circular strings naturally arise in several
 56 applications (for example, in multiple alignment of genomes, which often have a circu-
 57 lar molecular structure) [4, 5, 14, 19, 27, 37]. Depending on the application at hand,
 58 different cost functions are used. For example, non-trivial algorithms for computing
 59 a consensus string of three and four circular strings with respect to the Hamming
 60 distance have been developed [23]. However, most of the algorithmic work so far is
 61 heuristic in nature or only considers specific special cases. A thorough analysis of the
 62 computational complexity for these problems in general so far has been missing.

63 After having dealt with circular string alignment problems in a quite general
 64 fashion, we then study a fundamental (consensus) problem in time series analysis.
 65 *Dynamic time warping* (see Section 2 for details) defines a distance between two
 66 time series which is used in many applications in time series analysis [21, 28, 33, 36]
 67 (notably, dynamic time warping has also been considered in the context of circular

³We cast all problems in this work as decision problems for easier complexity-theoretic treatment. Our hardness results correspondingly hold for the associated optimization problems.

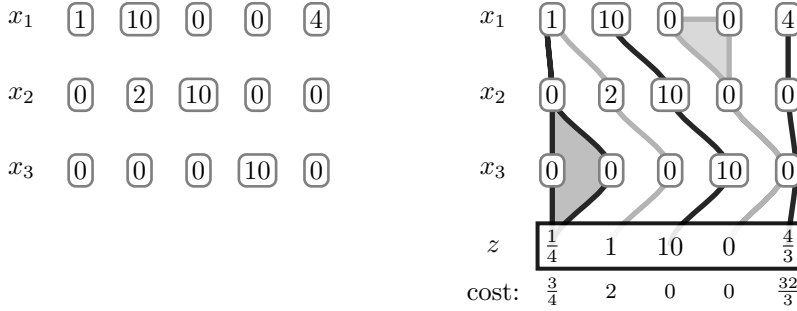


FIG. 2. A DTW-MEAN instance with three input sequences and an optimal length-5 mean (z). Alignments between the mean and input sequences can progress at different speeds. This is formalized using warping paths (see Section 2) represented by polygons (or lines in degenerate cases) with alternating shades. Every pair of aligned elements belongs to the same polygon. The cost of each mean element is the sum of squared differences over all aligned input elements, e.g. the cost of the first element is $(1 - \frac{1}{4})^2 + 3 \cdot (0 - \frac{1}{4})^2 = \frac{3}{4}$.

68 sequences [3, 29]). An important problem here is to compute an average of a given
 69 sample of time series under the dynamic time warping distance.

DTW-MEAN

Input: A list of k univariate rational time series x_1, \dots, x_k and $c \in \mathbb{Q}$.

Question: Is there a univariate rational time series z such that $\mathcal{F}(z) = \sum_{i=1}^k (\text{dtw}(z, x_i))^2 \leq c$?

71 Here, dtw denotes the dynamic time warping distance (see Section 2 for details).
 72 Intuitively, dynamic time warping allows for non-linear alignments between two series.
 73 Figure 2 depicts an example. The dtw -distance of two length- n time series can be
 74 computed via standard dynamic programming in $O(n^2)$ time. Some subquadratic
 75 algorithms are known [17, 22, 16]. For two binary time series, there exists an $O(n^{1.87})$ -
 76 time algorithm [1]. In general, however, a strongly subquadratic-time algorithm (that
 77 is, $O(n^{2-\varepsilon})$ time for some $\varepsilon > 0$) does not exist unless the Strong Exponential Time
 78 Hypothesis fails [1, 7, 22].

79 Regarding the computational complexity of DTW-MEAN, although more or less
 80 implicitly assumed in many publications presenting heuristic solution strategies⁴, NP-
 81 hardness still has been open (see Brill et al. [6, Section 3] for a discussion on some
 82 misconceptions and wrong statements in the literature). It is known to be solvable
 83 in $O(n^{2k+1}2^k k)$ time, where n is the maximum length of any input series [6]. Moreover,
 84 Brill et al. [6] presented a polynomial-time algorithm for the special case of binary
 85 time series which has been improved recently [34]. In practice, numerous heuristics
 86 are used [11, 31, 33, 35]. Note that DTW-MEAN is often described as closely related
 87 to multiple sequence alignment problems [2, 30, 32]. However, we are not aware
 88 of any formal proof regarding this connection. By giving a polynomial-time many-
 89 one reduction from MULTIPLE STRING CIRCULAR SHIFT (WITH COST f) to DTW-
 90 MEAN, we show that DTW-MEAN is actually connected to multiple circular sequence
 91 alignment problems. To the best of our knowledge, this is the first formally proven

⁴For instance, Petitjean et al. [31] write “Computational biologists have long known that averaging under time warping is a very complex problem, because it directly maps onto a multiple sequence alignment: the “Holy Grail” of computational biology.” Unfortunately, the term “directly maps” has not been formally defined and only sketchy explanations are given.

92 result regarding this connection.

93 *Our Results.* Using plausible complexity-theoretic assumptions, we provide a fine-
 94 grained picture of the exact computational complexity (including parameterized com-
 95 plexity) of the problems introduced above. We present two main results.

96 First, we show that, for a large class of natural cost functions f , f -MSCS on
 97 binary sequences is NP-hard, W[1]-hard with respect to the number k of inputs, and
 98 not solvable in $\rho(k) \cdot n^{o(k)}$ time for any computable function ρ (unless the Exponential
 99 Time Hypothesis fails). Note that f -MSCS is easily solvable in $\rho(k) \cdot n^{O(k)}$ time (for
 100 computable functions f) since there are at most n^{k-1} cyclic shifts to try out (without
 101 loss of generality, the first string is not shifted). Our running time lower bound thus
 102 implies that the brute-force approach can only be improved up to a constant factor in
 103 the exponent. Based on this, we can also prove the same hardness for the CIRCULAR
 104 CONSENSUS STRING problem. In fact, the general ideas of our reduction might also
 105 be used to develop hardness reductions for other circular string alignment problems.

106 As our second main contribution, we obtain the same list of hardness results as
 107 above for DTW-MEAN on binary time series. We achieve this by a polynomial-time
 108 reduction from a special case of f -MSCS. Our reduction implies that, unless the
 109 Exponential Time Hypothesis fails, the known $O(n^{2k+1}2^k k)$ -time algorithm [6] essen-
 110 tially can be improved only up to constants in the first exponent. Note that recently
 111 Buchin et al. [8] achieved the same hardness result for the problem of averaging time
 112 series under generalized (p, q) -DTW. Their reduction, however, does not yield binary
 113 input time series.

114 *Organization.* In Section 2 we fix notation and introduce basic concepts, also in-
 115 cluding the formal definition of dynamic time warping and the corresponding concept
 116 of warping paths. In Section 3, we identify a circular string problem (of independent
 117 interest in molecular biology) which forms the basis for the results in Section 5. More
 118 specifically, we prove the hardness results for MULTIPLE STRING CIRCULAR SHIFT
 119 (WITH COST f). The key ingredient here is a specially geared polynomial-time reduc-
 120 tion from the REGULAR MULTICOLORED CLIQUE problem. Moreover, we introduce
 121 the concept of polynomially bounded grouping functions f (for which our results hold).
 122 In Section 4, providing a reduction from MULTIPLE STRING CIRCULAR SHIFT (WITH
 123 COST f), we show analogous hardness results for CIRCULAR CONSENSUS STRING.
 124 Notably, the cost function corresponding to CIRCULAR CONSENSUS STRING is not a
 125 polynomially bounded grouping function, making the direct application of the result
 126 for MULTIPLE STRING CIRCULAR SHIFT (WITH COST f) impossible. In Section 5
 127 we prove analogous complexity results for DTW-MEAN, again devising a polynomial-
 128 time reduction from MULTIPLE STRING CIRCULAR SHIFT (WITH COST f). In Sec-
 129 tion 6, we conclude with some open questions and directions for future research.

130 **2. Preliminaries.** We briefly introduce our notation and formal definitions.

131 *Circular Shifts.* For a string $s = s[1] \dots s[n] \in \Sigma^n$, we denote its length n by $|s|$.
 132 For $0 \leq \delta < n$, we define the *circular (left) shift by δ* as the string

$$133 \quad s^{\leftarrow \delta} := s[\delta + 1] \dots s[n]s[1] \dots s[\delta] \quad (\text{note that } s^{\leftarrow \delta}[i] = s[(i + \delta - 1 \bmod n) + 1]),$$

135 that is, we circularly shift the string δ times to the left. Let s_1, \dots, s_k be strings
 136 of length n . A *multiple circular (left) shift* of s_1, \dots, s_k is defined by a k -tuple
 137 $\Delta = (\delta_1, \dots, \delta_k) \in \{0, \dots, n-1\}^k$ and yields the strings $s_1^{\leftarrow \delta_1}, \dots, s_k^{\leftarrow \delta_k}$. We define *col-*
 138 *umn* $i \in \{1, \dots, n\}$ of a multiple circular shift Δ as the k -tuple $(s_1^{\leftarrow \delta_1}[i], \dots, s_k^{\leftarrow \delta_k}[i])$.
 139 By *row* $j \in \{1, \dots, k\}$ of column i we denote the element $s_j^{\leftarrow \delta_j}[i]$.

140 *Cost Functions.* A *local cost function* is a function $f: \Sigma^* \rightarrow \mathbb{Q}$ assigning a cost
 141 to any tuple of values. Given such a function, the *overall cost* of a circular shift Δ
 142 for k length- n strings is defined as

$$143 \quad \text{cost}_f(\Delta) := \sum_{i=1}^n f((s_1^{\leftarrow \delta_1}[i], \dots, s_k^{\leftarrow \delta_k}[i])),$$

144 that is, we sum up the local costs of all columns of Δ .

145 An example for a local cost is the sum of squared distances from the arithmetic
 146 mean (i.e., k times the variance, here called σ), that is,

$$147 \quad \sigma((x_1, \dots, x_k)) = \sum_{i=1}^k \left(x_i - \frac{1}{k} \sum_{j=1}^k x_j \right)^2.$$

148
 149 Using a well-known formula for the variance, we get the following useful formula for σ :

$$150 \quad \sigma((x_1, \dots, x_k)) = \left(\sum_{j=1}^k x_j^2 \right) - \frac{1}{k} \left(\sum_{j=1}^k x_j \right)^2.$$

151 For binary strings (that is, $x_j \in \{0, 1\}$ for all $1 \leq j \leq k$), σ does only depend
 152 on the number $w := \sum_{j=1}^k x_j$ of 1's and the number $k - w$ of 0's and can be written
 153 (according to the formula above) as

$$154 \quad (2.1) \quad \sigma((x_1, \dots, x_k)) = w - \frac{w^2}{k} = \frac{w(k-w)}{k}.$$

155 We will repeatedly use this formula later on for cost calculations in the proof for
 156 DTW-MEAN ([Theorem 5.1](#)).

157 *Dynamic Time Warping.* A time series is a sequence $x = (x_1, \dots, x_n) \in \mathbb{Q}^n$. The
 158 dynamic time warping distance between two time series is based on the concept of a
 159 warping path.

160 **DEFINITION 2.1.** A warping path of order $m \times n$ is a sequence $p = (p_1, \dots, p_L)$,
 161 $L \in \mathbb{N}$, of index pairs $p_\ell = (i_\ell, j_\ell) \in \{1, \dots, m\} \times \{1, \dots, n\}$, $1 \leq \ell \leq L$, such that

- 162 (i) $p_1 = (1, 1)$,
- 163 (ii) $p_L = (m, n)$, and
- 164 (iii) $(i_{\ell+1} - i_\ell, j_{\ell+1} - j_\ell) \in \{(1, 0), (0, 1), (1, 1)\}$ for each $1 \leq \ell \leq L - 1$.

165 See [Figure 2](#) in [Section 1](#) for an example.

166 The set of all warping paths of order $m \times n$ is denoted by $\mathcal{P}_{m,n}$. A warping
 167 path $p \in \mathcal{P}_{m,n}$ defines an *alignment* between two time series $x = (x[1], \dots, x[m])$
 168 and $y = (y[1], \dots, y[n])$ in the following way: Every pair $(i, j) \in p$ *aligns* element x_i
 169 with y_j . Note that every element from x can be aligned with multiple elements
 170 from y , and vice versa. The *dtw-distance* (with squared cost function) between x
 171 and y is defined as

$$172 \quad \text{dtw}(x, y) := \min_{p \in \mathcal{P}_{m,n}} \left(\sum_{(i,j) \in p} (x[i] - y[j])^2 \right)^{1/2}.$$

173 Note that also other cost functions can be considered. In this work, we only consider
 174 the most common case of squared costs.

175 A *mean* of time series x_1, \dots, x_k is a time series that minimizes the *Fréchet*
176 function

$$177 \quad \mathcal{F}(z) := \sum_{j=1}^k (\text{dtw}(z, x_j))^2 = \sum_{j=1}^k \min_{p_j \in \mathcal{P}_{|x_j|, |z|}} \sum_{(u,v) \in p_j} (x_j[u] - z[v])^2.$$

178 Note that given, for each $j \in [k]$, a warping path p_j between z and x_j , the value of $z[i]$
179 that minimizes

$$180 \quad (2.2) \quad \sum_{j=1}^k \sum_{(u,v) \in p_j} (x_j[u] - z[v])^2$$

181 is the arithmetic mean of all values aligned to $z[i]$,

$$182 \quad z[i] = \frac{\sum_{j=1}^k \sum_{(i,u) \in p_j} x_j[u]}{\sum_{j=1}^k |\{(i,u) \in p_j\}|}.$$

183 That is, the length of a mean together with the optimal alignments to the input time
184 series determine the mean. The contribution of $z[i]$ to the sum (2.2) is the sum of
185 squared distances between $z[i]$ and all values aligned to $z[i]$,

$$186 \quad \sum_{j=1}^k \sum_{(u,i) \in p_j} (x_j[u] - z[i])^2.$$

187 Note that this corresponds to the cost function σ above.

188 We remark that for DTW-MEAN, often a normalized cost $F(z) := \frac{1}{k} \mathcal{F}(z)$ is
189 considered. Clearly, this does not affect the computational complexity of the problem,
190 so for simplification purposes we only consider the non-normalized cost $\mathcal{F}(z)$.

191 *Parameterized Complexity.* We assume familiarity with the basic concepts from
192 classic and parameterized complexity theory.

193 An instance of a parameterized problem is a pair (I, k) consisting of the classic
194 problem instance I and a natural number k (the *parameter*). A parameterized problem
195 is contained in the class XP if there is an algorithm solving an instance (I, k) in
196 polynomial time if k is a constant, that is, in time $O(|I|^{f(k)})$ for some computable
197 function f only depending on k (here $|I|$ is the size of I). A parameterized problem
198 is *fixed-parameter tractable* (contained in the class FPT) if it is solvable in time $f(k) \cdot$
199 $|I|^{O(1)}$ for some computable function f depending solely on k . The class W[1] contains
200 all problems which are parameterized reducible to CLIQUE parameterized by the clique
201 size. A parameterized reduction from a problem Q to a problem P is an algorithm
202 mapping an instance (I, k) of Q in time $f(k) \cdot |I|^{O(1)}$ to an equivalent instance (I', k')
203 of P such that $k' \leq g(k)$ (for some computable functions f and g). It holds that
204 $\text{FPT} \subseteq \text{W}[1] \subseteq \text{XP}$.

205 A parameterized problem that is W[1]-hard with respect to a parameter (such as
206 CLIQUE with parameter clique size) is widely believed not to be in FPT.

207 *Exponential Time Hypothesis.* Impagliazzo and Paturi [20] formulated the *Expo-*
208 *ponential Time Hypothesis* (ETH) which asserts that there exists a constant $c > 0$ such
209 that 3-SAT cannot be solved in $O(2^{cn})$ time, where n is the number n of variables in
210 the input formula. It is a stronger assumption than common complexity assumptions
211 such as $\text{P} \neq \text{NP}$ or $\text{FPT} \neq \text{W}[1]$.

212 Several conditional running time lower bounds have since been shown based on
213 the ETH, for example, CLIQUE cannot be solved in $\rho(k) \cdot n^{o(k)}$ time for any computable
214 function ρ unless the ETH fails [10].

215 **3. Hardness of f -MSCS on Binary Strings.** In this section, we focus on
 216 binary strings from $\{0, 1\}^*$. We prove hardness for a family of local cost functions
 217 that satisfy certain properties. The functions we consider have the common property
 218 that they only depend on the number of 0's and 1's in a column, and that they aim
 219 at grouping similar values together.

220 **DEFINITION 3.1.** *A function $f: \{0, 1\}^* \rightarrow \mathbb{Q}$ is called order-independent (or sym-*
 221 *metric) if, for each $k \in \mathbb{N}$, there exists a function $f_k: \{0, \dots, k\} \rightarrow \mathbb{Q}$ such that*
 222 *$f((x_1, \dots, x_k)) = f_k(\sum_{j=1}^k x_j)$ holds for all $(x_1, \dots, x_k) \in \{0, 1\}^k$.*

223 *For an order-independent function f , we define the function $f'_k: \{1, \dots, k\} \rightarrow \mathbb{Q}$*
 224 *as*

$$225 \quad f'_k(x) := \frac{f_k(x) - f_k(0)}{x}.$$

226 *An order-independent function f is grouping if $f'_k(k) < \min_{1 \leq x < k} f'_k(x)$ and $f'_k(2) <$*
 227 *$f'_k(1)$ holds for every $k \in \mathbb{N}$.*

228 For an order-independent function f , f'_k can be seen as the cost per 1-value (a
 229 column with x 1's and $k - x$ 0's has cost $f_k(x) = f_k(0) + x f'_k(x)$). It can also be
 230 seen as a discrete version of the derivative for f_k , so that if f_k is concave then f'_k is
 231 decreasing. The intuition behind a grouping function is that the cost per 1-value is
 232 minimal in columns containing only 1's, and that having two 1's in a column has less
 233 cost than having two columns with a single 1. In particular, any function f where
 234 all f_k are strictly concave is grouping. Finally, if f is grouping, then the function

$$235 \quad (x_1, \dots, x_k) \mapsto f_k \left(\sum_{j=1}^k x_j \right) + a \sum_{j=1}^k x_j + b$$

236 is also grouping for any $a, b \in \mathbb{Q}$.

237 The following definitions are required to ensure that our subsequent reduction
 238 ([Lemma 3.3](#)) is computable in polynomial time.

239 **DEFINITION 3.2.** *Let f be an order-independent function. The gap of f_k is*

$$240 \quad \varepsilon_k := \min\{f'_k(x) - f'_k(y) \mid 1 \leq x, y \leq k, f'_k(x) > f'_k(y)\}.$$

241 *The range of f_k is $\mu_k := \max_{1 \leq x \leq k} |f'_k(x)|$.*

242 *An order-independent function f is polynomially bounded if it is polynomial-time*
 243 *computable and if, for every $k \in \mathbb{N}$, μ_k and ε_k^{-1} are upper-bounded by a polynomial*
 244 *in k .*

245 For binary strings, the function σ (see [Section 2](#)) is a polynomially bounded
 246 grouping function. Indeed, it is order-independent since $\sigma((x_1, \dots, x_k)) = \frac{w(k-w)}{k}$,
 247 where $w = \sum_{j=1}^k x_j$. Thus, $\sigma_k(w) = \frac{w(k-w)}{k}$ and we have $\sigma_k(0) = 0$, and $\sigma'_k(w) =$
 248 $\frac{k-w}{k}$, so σ'_k is strictly decreasing, which is sufficient for σ to be grouping. Finally, it
 249 is polynomially bounded, with gap $\varepsilon_k = \frac{1}{k}$ and range $\mu_k = \frac{k-1}{k} \leq 1$.

250 We prove our hardness results with a polynomial-time reduction from a special
 251 version of the CLIQUE problem.

REGULAR MULTICOLORED CLIQUE (RMCC)

Input: A d -regular undirected graph $G = (V, E)$ where the vertices are colored with k colors such that each color class contains the same number of vertices.

Question: Does G have a size- k complete subgraph (containing $\binom{k}{2}$ edges, called a k -clique) with exactly one vertex from each color?

RMCC is known to be NP-hard, W[1]-hard with respect to k , and not solvable in $\rho(k) \cdot |V|^{o(k)}$ time for any computable function ρ unless the ETH fails [12].

The following lemma states the existence of a polynomial-time reduction from RMCC to f -MSCS which implies hardness of f -MSCS for polynomially bounded grouping functions.

LEMMA 3.3. *Let f be a polynomially bounded grouping function. Then there is a polynomial-time reduction that, given an RMCC instance $G = (V, E)$ with k colors, outputs binary strings s_0, \dots, s_k of equal length and $c \in \mathbb{Q}$ such that the following holds:*

- If G contains a properly colored k -clique, then there exists a multiple circular shift Δ of s_0, \dots, s_k with $\text{cost}_f(\Delta) = c$.
- If G does not contain a properly colored k -clique, then every multiple circular shift Δ of s_0, \dots, s_k has $\text{cost}_f(\Delta) \geq c + \varepsilon_{k+1}$.

To prove Lemma 3.3, we first describe the reduction and then prove several claims about the structure and the costs of multiple circular shifts in the resulting f -MSCS instance.

Reduction. Consider an instance of RMCC, that is, a graph $G = (V, E)$ with a partition of V into k subsets V_1, \dots, V_k of size $n := \frac{|V|}{k}$ each, such that each vertex has degree d . Let $V_j = \{v_{j,1}, \dots, v_{j,n}\}$, $m = |E|$, and $E = \{e_1, \dots, e_m\}$. We assume that $k \geq 3$ since the instance is trivially solvable otherwise.

We build an f -MSCS instance with $k+1$ binary strings. Hence, the local cost of a column of a multiple circular shift is given by the function f_{k+1} . For simplicity, we write f' , gap ε , and range μ for f'_{k+1} , ε_{k+1} , and μ_{k+1} .

For each $j \in \{1, \dots, k\}$, let p_j be the length- k string such that $p_j[h] = 1$ if $h = j$, and $p_j[h] = 0$ otherwise. For each vertex $v_{j,i}$, let $q_{j,i} \in \{0, 1\}^m$ be the string such that

$$q_{j,i}[h] := \begin{cases} 1, & \text{if } 1 \leq h \leq m \text{ and } v_{j,i} \in e_h \\ 0, & \text{otherwise} \end{cases}$$

and let $u_{j,i} := p_j q_{j,i}$ be the concatenation of p_j and $q_{j,i}$. Note that $u_{j,i}$ has length $m' := m + k$, contains $1 + d$ ones, and $m' - 1 - d$ zeros. Let $\mathbf{0} := 0^{m'}$ be the string containing m' zeros and define the numbers

$$\kappa := knd + kn + k,$$

$$\gamma := nk,$$

$$\lambda := \max \left\{ \left\lceil \kappa \left(\frac{2\mu}{\varepsilon} + 1 \right) \right\rceil, 2n(\gamma + k + 1) \right\} + 1.$$

For $1 \leq j \leq k$, we define the string

$$s_j := 1u_{j,1}(\mathbf{10})^{\gamma+j} 1u_{j,2}(\mathbf{10})^{\gamma+j} \dots 1u_{j,n}(\mathbf{10})^{\gamma+j} (\mathbf{10})^{\lambda-n(\gamma+j+1)}.$$

Note that $|1u_{j,i}| = |\mathbf{10}| = m' + 1$. Thus, each string s_j has length

$$n(m' + 1)(1 + \gamma + j) + (m' + 1)(\lambda - n(\gamma + j + 1)) = \lambda(m' + 1) =: \ell,$$

290 where $\ell \leq \text{poly}(nk)$. We further define the following length- ℓ *dummy* string

$$291 \quad s_0 = 11^k 0^m (\mathbf{10})^{\lambda-1}.$$

292 Finally, we define the target cost

$$\begin{aligned} 293 \quad c := & \ell f_{k+1}(0) \\ 294 \quad & + \lambda(k+1)f'(k+1) \\ 295 \quad & + 2 \left(k + \binom{k}{2} \right) (f'(2) - f'(1)) \\ 296 \quad & + \kappa f'(1). \end{aligned}$$

298 Clearly, the strings s_0, \dots, s_k and the value c can be computed in polynomial time.
299 Our construction is illustrated in [Figure 3](#).

300 In the strings s_0, \dots, s_k , any 1-value at a position i with $i \bmod (m' + 1) = 1$ is
301 called a *separator*, other 1-values are *coding* values. A coding value is either *vertex-*
302 *coding* if it belongs to some p_j (or to the k coding values of s_0), or *edge-coding*
303 otherwise (then it belongs to some $q_{i,j}$). There are $\lambda(k+1)$ separator values in total
304 and κ coding values.

305 Given a multiple circular shift Δ , we define the *weight* w of a column as the
306 number of 1-values it contains, that is, $0 \leq w \leq k+1$. The cost for such column is
307 $f_{k+1}(w) = f_{k+1}(0) + wf'(w)$. Each 1-value of this column is attributed a *local cost*
308 of $f'(w)$, so that the cost of any solution is composed of a *base cost* of $\ell f_{k+1}(0)$ and
309 of the sum of all local costs of all 1-values. In the following we mainly focus on local
310 costs.

311 It remains to be shown that there exists a multiple circular shift of s_0, \dots, s_k with
312 cost c if G contains a properly colored k -clique, and that otherwise every multiple
313 circular shift has cost at least $c + \varepsilon$. We proceed by analyzing the structure and costs
314 of optimal multiple circular shifts.

315 *Aligning Separators.* Let $\Delta = (\delta_0, \dots, \delta_k)$ be a multiple circular shift of s_0, \dots, s_k .
316 Without loss of generality, we can assume that $\delta_0 = 0$ since setting each δ_j to $(\delta_j -$
317 $\delta_0) \bmod \ell$ yields a shift with the same cost. First, we show that if $\delta_j \bmod (m' + 1) \neq 0$
318 holds for some $0 < j \leq k$, then Δ has large cost.

319 **CLAIM 3.4.** *For any multiple circular shift $\Delta = (\delta_0 = 0, \delta_1, \dots, \delta_k)$ with $\delta_j \bmod$*
320 *$(m' + 1) \neq 0$ for some $1 < j \leq k$, it holds that $\text{cost}_f(\Delta) \geq c + \varepsilon$.*

321 *Proof.* Assume that $\delta_j \bmod (m' + 1) = a \in \{1, \dots, m'\}$ for some $0 < j \leq k$.
322 We count the number of weight- $(k+1)$ columns: such a column cannot only contain
323 separator values since it cannot contain a separator value in both row 0 and row j .
324 Hence, it contains at least one coding value. Since there are κ coding values, there
325 are at most κ weight- $(k+1)$ columns, so at most $k\kappa$ separator values have local cost
326 $f'(k+1)$. All other separator values have local cost $f'(w)$ for some $w < k+1$, which
327 is at least $f'(k+1) + \varepsilon$. There are at least $\lambda(k+1) - k\kappa$ such separator values. Adding
328 the base cost of $\ell f_{k+1}(0)$, the cost of Δ is thus at least:

$$\begin{aligned} 329 \quad \text{cost}_f(\Delta) & \geq \ell f_{k+1}(0) + (\lambda(k+1) - k\kappa)(f'(k+1) + \varepsilon) \\ 330 \quad & \geq \ell f_{k+1}(0) + \lambda(k+1)f'(k+1) + \lambda k\varepsilon - k\kappa(\mu + \varepsilon). \end{aligned}$$

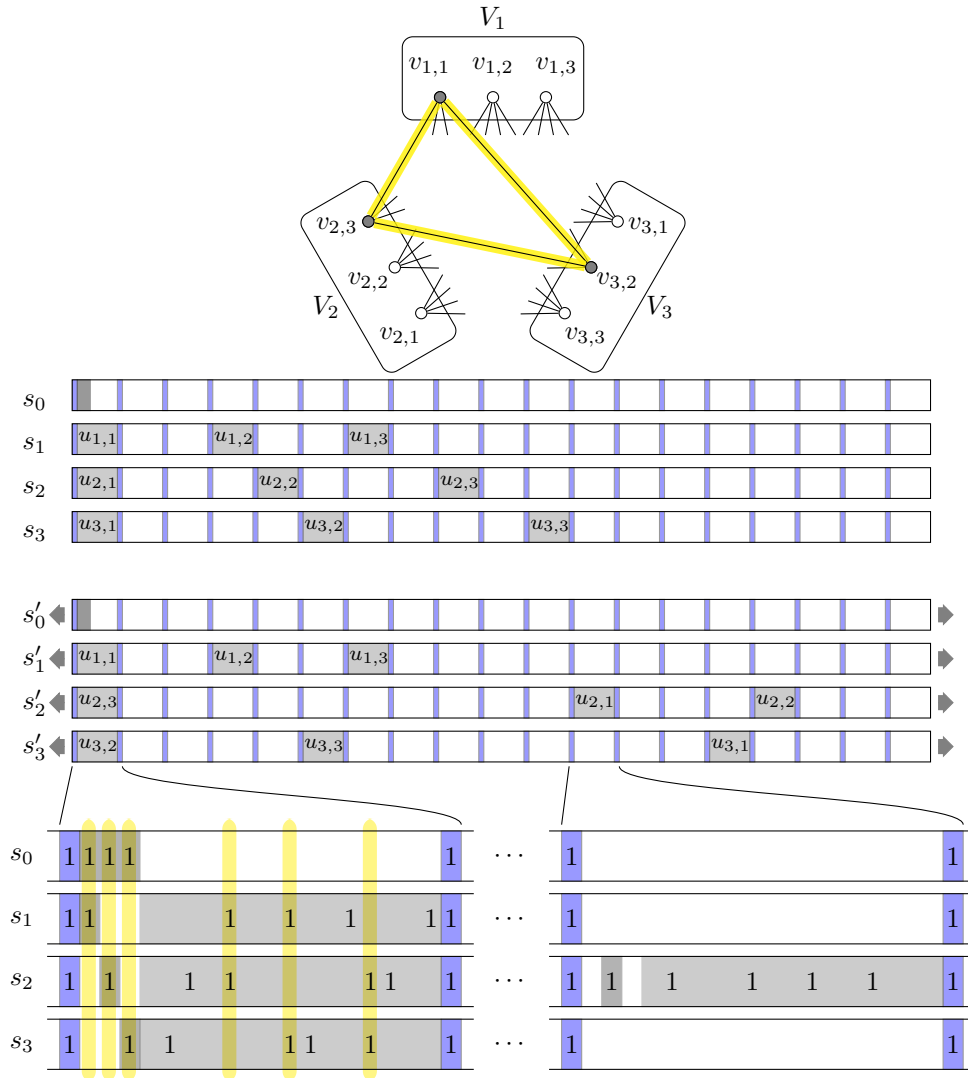


FIG. 3. Illustration of the reduction from an instance of RMCC (top) with $k = 3$. Middle: Sequences s_0 to s_3 , and their optimal circular shifts s'_0 to s'_3 . Blue stripes represent the regularly-spaced separator 1-values. The (light) gray intervals contain both 0's and 1's according to strings $u_{i,j}$, and white intervals contain only 0's. The spacing between consecutive $u_{i,j}$'s is defined using γ and the overall string length depends on λ , both values are chosen so as to restrict the possible alignments between different $u_{i,j}$'s; in this example we use $\gamma = 1$ and $\lambda = 19$. Bottom: a zoom-in on blocks 1 and 12 in the shifted strings (only non-0 values are indicated, weight-2 columns are highlighted). Through vertex columns, the dummy string s_0 ensures that one vertex occupies block 1 in each row, and weight-2 edge-columns ensure that $\binom{k}{2}$ edges (as highlighted in the graph) are induced by these vertices.

332 Recall that

$$\begin{aligned}
 333 \quad c &= \ell f_{k+1}(0) + \lambda(k+1)f'(k+1) + 2 \left(k + \binom{k}{2} \right) (f'(2) - f'(1)) + \kappa f'(1) \\
 334 \quad &\leq \ell f_{k+1}(0) + \lambda(k+1)f'(k+1) + \kappa \mu
 \end{aligned}$$

336 since $f'(2) - f'(1) < 0$. Combining the above bounds for c and $\text{cost}_f(\Delta)$ using
 337 $\lambda \geq \kappa \left(\frac{2\mu}{\varepsilon} + 1\right) + 1$ (by definition) yields

$$\begin{aligned}
 338 \quad \text{cost}_f(\Delta) - c &\geq \lambda k \varepsilon - k \kappa (\mu + \varepsilon) - \kappa \mu \\
 339 &\geq 2 \kappa k \mu + \kappa k \varepsilon + k \varepsilon - k \kappa (\mu + \varepsilon) - \kappa \mu \\
 340 &\geq \varepsilon. \quad \square
 \end{aligned}$$

342 *Cost of Circular Shifts.* We assume from now on that $\delta_j \bmod (m' + 1) = 0$ for all
 343 $j \in \{0, \dots, k\}$. We now provide a precise characterization of the cost of Δ .

344 For $l \in \{1, \dots, \lambda\}$, we define the l -th block consisting of the m' consecutive
 345 columns $(l-1)(m'+1)+2, \dots, l(m'+1)$. The block index of column i is $i-1 \bmod (m'+1)$.
 346 For $j \in \{1, \dots, k\}$, the substring $s_j^{\leftarrow \delta_j}[(l-1)(m'+1)+2] \dots s_j^{\leftarrow \delta_j}[l(m'+1)]$ cor-
 347 responding to the l -th block of $s_j^{\leftarrow \delta_j}$ either equals some $u_{j,i}$ or $\mathbf{0}$. We say that block l
 348 is occupied by vertex $v_{j,i} \in V_j$, if the corresponding substring of $s_j^{\leftarrow \delta_j}$ is $u_{j,i}$. Note
 349 that for each j there are n distinct blocks out of λ that are occupied by a vertex
 350 in V_j . Columns with block-index 1 to k are called *vertex-columns* and columns with
 351 block-index $k+1$ to $k+m = m'$ are *edge-columns* (they may only contain edge-coding
 352 values from some $q_{i,j}$). Let P denote the set of vertices occupying block 1.

353 **OBSERVATION 3.5.** *In block l , if the vertex-column with block-index h has weight 2,*
 354 *then $l = 1$, and $V_h \cap P \neq \emptyset$. No vertex-column can have weight 3 or more.*

355 *Proof.* Consider the vertex-column with block-index h . By construction, among
 356 s_1, \dots, s_k only s_h may have a 1 in this column (which is true if some vertex from V_h
 357 occupies this block). The string s_0 has a 1 in this column if it is a column in block 1.
 358 Thus, assuming that column h has weight 2 implies $l = 1$ and $V_h \cap P \neq \emptyset$. \square

359 **OBSERVATION 3.6.** *In block l , if the edge-column with block-index $k+h$, $1 \leq$
 360 $h \leq m$, has weight 2, then block l is occupied by both vertices of edge $e_h \in E$. No
 361 edge-column can have weight 3 or more.*

362 *Proof.* Consider an edge-column with block-index $k+h$, $1 \leq h \leq m$. Denote by
 363 v_{j_0, i_0} and v_{j_1, i_1} the endpoints of edge e_h . For any $1 \leq j \leq k$, s_j has a 1 in this column
 364 only if block l is occupied by some vertex $v_{j,i}$, and, moreover, only if $u_{j,i}$ has a 1 in
 365 column h , i.e., $v_{j,i} = v_{j_0, i_0}$ or $v_{j,i} = v_{j_1, i_1}$, hence $j = j_1$ or $j = j_2$. So this column
 366 may not have weight 3 or more, and if it has weight 2, then block l is occupied by
 367 both endpoints of e_h . \square

368 From Observations 3.5 and 3.6 it follows that no column (beside separators) can
 369 have weight 3 or more. Since the number of coding values is fixed, the cost is entirely
 370 determined by the number of weight-2 columns. In the following, we will first give an
 371 upper bound on the number of weight-2 columns and then analyze how it determines
 372 the cost.

373 **Observation 3.5** gives a direct upper bound of at most k weight-2 vertex columns
 374 (since they all are in block 1). We now focus on weight-2 edge-columns. The following
 375 claim will help us to show an upper bound on their number.

376 **CLAIM 3.7.** *For any two rows j, j' , there exists at most one block l that is occupied*
 377 *by vertices from both V_j and $V_{j'}$.*

378 *Proof.* First, note that if two distinct blocks l and l' are both occupied by a
 379 vertex from the same row j , then, by construction, there are two possible cases:
 380 either $|l-l'| = a(\gamma+j+1)$ or $|l-l'| = \lambda - a(\gamma+j+1)$, where $1 \leq a < n$ in both cases.

381 Indeed, there are n regularly-spaced substrings $u_{j,i}$ (having $\gamma + j$ blocks in between
382 them) in s_j (consisting of λ blocks in total).

383 Assume towards a contradiction that two distinct blocks l and l' are each occupied
384 by a vertex from V_j and $V_{j'}$. Then, there exists an $a \in \{1, \dots, n-1\}$ such that
385 $|l-l'| = a(\gamma+j+1)$ or $|l-l'| = \lambda - a(\gamma+j+1)$, and there exists an $a' \in \{1, \dots, n-1\}$
386 such that $|l-l'| = a'(\gamma+j'+1)$ or $|l-l'| = \lambda - a'(\gamma+j'+1)$. This gives four cases
387 to consider (in fact just three by symmetry of j and j').

388 If $|l-l'| = a(\gamma+j+1) = a'(\gamma+j'+1)$, then $(a-a')(\gamma+1) = a'j' - aj$. We
389 have $a \neq a'$, as otherwise this would imply $j = j'$. So $|a'j' - aj| \geq \gamma + 1$, but this
390 is impossible since $a < n$, $a' < n$, $j \leq k$, $j' \leq k$, that is, $|a'j' - aj| \leq kn$, whereas
391 $\gamma + 1 = kn + 1$ (by construction).

392 If $|l-l'| = a(\gamma+j+1) = \lambda - a'(\gamma+j'+1)$, then $\lambda = a(\gamma+j+1) + a'(\gamma+j'+1)$.
393 However, $\lambda > 2n(\gamma+k+1)$ by construction, so this case also leads to a contradiction.

394 Finally, if $|l-l'| = \lambda - a(\gamma+j+1) = \lambda - a'(\gamma+j'+1)$, then we have $a(\gamma+j+1) =$
395 $a'(\gamma+j'+1)$. This case yields, as the first case, a contradiction. \square

396 CLAIM 3.8. *There are at most $\binom{k}{2}$ weight-2 edge-columns.*

397 *Proof.* Consider any pair j, j' such that $1 \leq j < j' \leq k$. It suffices to show that
398 there exists at most one weight-2 edge-column with a 1 in rows j and j' . Aiming
399 at a contradiction, assume that two such columns exist. By Observation 3.6, they
400 must each belong to a block which is occupied by vertices both in V_j and $V_{j'}$. From
401 Claim 3.7 it follows that both columns belong to the same block. Let v and v' be the
402 vertices of V_j and $V_{j'}$, respectively, occupying this block. By Observation 3.6 again,
403 both edges are equal to $\{v, v'\}$, which contradicts the fact that they are distinct. \square

404 Having established an upper bound of $k + \binom{k}{2}$ for the number of weight-2 columns,
405 the following result describes the corresponding cost.

406 CLAIM 3.9. *Let W_2 be the number of weight-2 columns. If $W_2 = k + \binom{k}{2}$, then*
407 *$\text{cost}_f(\Delta) = c$. If $W_2 < k + \binom{k}{2}$, then $\text{cost}_f(\Delta) \geq c + \varepsilon$.*

408 *Proof.* The base cost $\ell f_{k+1}(0)$ of the solution only depends on the number ℓ of
409 columns. Separator values are in weight- $(k+1)$ columns. Since there are λ of them,
410 it follows that the total local cost of all separator values is $\lambda(k+1)f'(k+1)$.

411 The total number of coding values is κ , each coding value has a local weight
412 of $f'(1)$ if it belongs to a weight-1 column, and $f'(2)$ otherwise (since there is no
413 vertex- or edge-column with weight 3 or more). There are W_2 weight-2 columns, so
414 exactly $2W_2$ coding values within weight-2 columns. Summing the base cost with the
415 local costs of all separator and coding values, we get:

$$\begin{aligned} 416 \quad \text{cost}_f(\Delta) &= \ell f_{k+1}(0) \\ 417 &\quad + \lambda(k+1)f'(k+1) \\ 418 &\quad + 2W_2(f'(2) - f'(1)) \\ 419 &\quad + \kappa f'(1). \end{aligned}$$

421 Thus, by definition of c , we have $\text{cost}_f(\Delta) = c$ if $W_2 = k + \binom{k}{2}$. If $W_2 < k + \binom{k}{2}$, then
422 using the fact that, by assumption, $f'(2) - f'(1) \leq -\varepsilon$, we obtain

$$423 \quad \text{cost}_f(\Delta) = c + 2 \left(W_2 - k - \binom{k}{2} \right) (f'(2) - f'(1)) \geq c + \varepsilon. \quad \square$$

424 CLAIM 3.10. *If G does not contain a properly colored k -clique, then there are at*
 425 *most $k + \binom{k}{2} - 1$ weight-2 columns.*

426 *Proof.* We prove the contrapositive. Assume that there are at least $k + \binom{k}{2}$ weight-
 427 2 columns. By Claim 3.8, there are at least k weight-2 vertex-columns. By Obser-
 428 vation 3.5, only the k vertex-columns of block 1 may have weight 2, hence for each
 429 $1 \leq j \leq k$ the column of block 1 with block-index j has weight 2. Thus, for every j ,
 430 $P \cap V_j \neq \emptyset$.

431 By Claim 3.7, no other block than block 1 may be occupied by two vertices,
 432 hence any edge-column with weight 2 must be in block 1, and both endpoints are
 433 in P . There cannot be more than k weight-2 vertex-columns, hence there are $\binom{k}{2}$
 434 weight-2 edge-columns, and for each of these there exists a distinct edge with both
 435 endpoints in P . Thus, P is a properly colored k -clique. \square

436 *Cliques and Circular Shifts with Low Cost.* We are now ready to complete the
 437 proof of Lemma 3.3. First, assume that G contains a properly colored k -clique $P =$
 438 $\{v_{1,i_1}, \dots, v_{k,i_k}\}$. Consider the multiple circular shift $\Delta = (\delta_0, \dots, \delta_k)$, where $\delta_0 = 0$
 439 and

$$440 \quad \delta_j := (i_j - 1)(m' + 1)(\gamma + j + 1)$$

441 for $j \in \{1, \dots, k\}$. Note that $|P| = k$, and all edge-columns in block 1 corresponding
 442 to edges induced in P have weight 2. Hence there are $\binom{k}{2}$ weight-2 edge-columns and
 443 k weight-2 vertex-columns. By Claim 3.9, $\text{cost}_f(\Delta) = c$.

444 Now, assume that G does not contain a properly colored k -clique. Let $\Delta =$
 445 $(\delta_0, \dots, \delta_k)$ be a multiple circular shift with $\delta_0 = 0$ (recall that we can assume this
 446 without loss of generality). Clearly, if $\delta_j \bmod (m' + 1) \neq 0$ for some j , then $\text{cost}_f(\Delta) \geq$
 447 $c + \varepsilon$ (by Claim 3.4). Otherwise, by Claim 3.10 there are at most $k + \binom{k}{2} - 1$ weight-2
 448 columns. By Claim 3.9, $\text{cost}_f(\Delta) \geq c + \varepsilon$.

449 This completes the proof of Lemma 3.3 which directly leads to our main result of
 450 this section.

451 THEOREM 3.11. *Let f be a polynomially bounded grouping function. Then, f -*
 452 *MSCS on binary strings is*

- 453 (i) NP-hard,
- 454 (ii) W[1]-hard with respect to the number k of input strings, and
- 455 (iii) not solvable in $\rho(k) \cdot n^{o(k)}$ time for any computable function ρ unless the ETH
- 456 fails.

457 *Proof.* The polynomial-time reduction from Lemma 3.3 yields the NP-hardness.
 458 Moreover, the number of strings in the f -MSCS instance only depends on the size of
 459 the multicolored clique. Hence, it is a parameterized reduction from RMCC param-
 460 eterized by the size of the clique to f -MSCS parameterized by the number of input
 461 strings and thus yields W[1]-hardness. Lastly, the number $k' = k + 1$ of strings is lin-
 462 ear in the size k of the clique. Thus, any $\rho(k') \cdot n^{o(k')}$ -time algorithm for DTW-MEAN
 463 would imply a $\rho'(k) \cdot |V|^{o(k)}$ -time algorithm for RMCC, contradicting the ETH. \square

464 Note that Theorem 3.11 holds for the function σ since it is a polynomially bounded
 465 grouping function (as discussed earlier).

466 The assumption that f is polynomially bounded is only needed to obtain a
 467 polynomial-time reduction in Lemma 3.3. Without this assumption, we still obtain a
 468 parameterized reduction from RMCC parameterized by the clique size to f -MSCS
 469 parameterized by the number of input strings, which yields the following corollary for
 470 a larger class of functions.

471 COROLLARY 3.12. *Let f be a computable grouping function. Then, f -MSCS on*
 472 *binary strings is $W[1]$ -hard with respect to the number k of input strings and not*
 473 *solvable in $\rho(k) \cdot n^{o(k)}$ time for any computable function ρ unless the ETH fails.*

474 **4. Circular Consensus String.** In this section we briefly study the CIRCULAR
 475 CONSENSUS STRING problem:

CIRCULAR CONSENSUS STRING (SCC)

Input: A list of k strings $s_1, \dots, s_k \in \Sigma^n$ of length n and $c \in \mathbb{Q}$.

476 **Question:** Is there a string $s^* \in \Sigma^n$ and a multiple circular shift $(\delta_1, \dots, \delta_k)$
 such that $\sum_{j=1}^k d(s_j^{\leftarrow \delta_j}, s^*) \leq c$?

477 Here, d denotes the Hamming distance, that is, the number of mismatches between
 478 the positions of two strings. Although consensus string problems in general have been
 479 widely studied from a theoretical point of view [9], somewhat surprisingly this is not
 480 the case for the circular version(s). For CCS, only an $O(n^2 \log n)$ -time algorithm
 481 for $k = 3$ and an $O(n^3 \log n)$ -time algorithm for $k = 4$ is known [23]. However, for
 482 general k no hardness result is known. Note that without circular shifts the problem
 483 is solvable in linear time: It is optimal to set $s^*[i]$ to any element that appears a
 484 maximum number of times among the elements $s_1[i], \dots, s_k[i]$.

485 For binary strings, it can easily be seen that the cost induced by column i is the
 486 minimum of the number of 0's and the number of 1's. Let f^{CS} be the polynomially
 487 bounded order-independent function with $f_k^{\text{CS}}(w) = \min\{w, k - w\}$. It follows from
 488 the discussion above that CIRCULAR CONSENSUS STRING is exactly f^{CS} -MSCS. Note,
 489 however, that f^{CS} is not a grouping function since $f_k^{\text{CS}'}(2) = f_k^{\text{CS}'}(1) = 1$. That is, we
 490 do not immediately obtain hardness of CCS from Theorem 3.11. We can still prove
 491 hardness via a reduction using a properly chosen polynomially bounded grouping
 492 function.

493 THEOREM 4.1. CIRCULAR CONSENSUS STRING on binary strings is

- 494 (i) NP-hard,
 495 (ii) $W[1]$ -hard with respect to the number k of input strings, and
 496 (iii) not solvable in $\rho(k) \cdot n^{o(k)}$ time for any computable function ρ unless the ETH
 497 fails.

498 *Proof.* As discussed above, CCS is equivalent to f^{CS} -MSCS. To prove hardness,
 499 we define a local cost function g (similar to f^{CS}) and reduce from g -MSCS to f^{CS} -
 500 MSCS.

501 Let g be the order-independent local cost function such that

$$502 \quad g_k(w) := f_{2k-2}^{\text{CS}}(w + (k - 2)) = \min\{w + k - 2, k - w\}.$$

503 Note that the function g_k is linearly decreasing on $\{1, \dots, k\}$ and that $g'_k(w) = \frac{2-w}{w} =$
 504 $\frac{2}{w} - 1$. The range of g_k is $\mu_k = 1$ and its gap is $\varepsilon_k = \frac{2}{k-1} - \frac{2}{k} > \frac{2}{k^2}$. That is, g satisfies all
 505 conditions of Theorem 3.11 and the corresponding hardness results hold for g -MSCS.
 506 See Figure 4 for an illustration.

507 Given an instance $\mathcal{I} = (s_1, \dots, s_k, c)$ of g -MSCS, we define the strings $s_j := 1^{|s_1|}$
 508 for $j = k + 1, \dots, 2k - 2$. We show that \mathcal{I} is a yes-instance if and only if $\mathcal{I}' :=$
 509 $(s_1, \dots, s_{2k-2}, c)$ is a yes-instance for f^{CS} -MSCS.

510 For the forward direction, consider a multiple circular shift $\Delta = (\delta_1, \dots, \delta_k)$
 511 of s_1, \dots, s_k such that $\text{cost}_g(\Delta) \leq c$. We define the multiple circular shift $\Delta' :=$
 512 $(\delta_1, \dots, \delta_k, \delta_{k+1} = 0, \dots, \delta_{2k-2} = 0)$ of s_1, \dots, s_{2k-2} . Consider column i of Δ' and
 513 let w' be the number of 1's it contains. Then, $w' = w + k - 2$, where w is the number

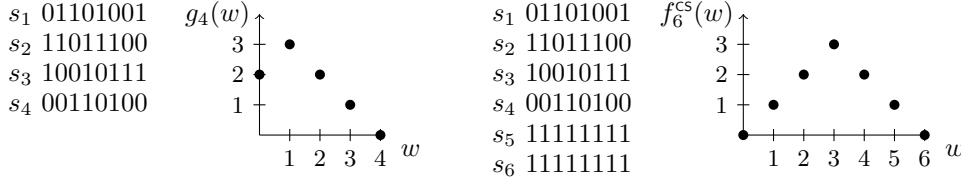


FIG. 4. Reduction from an instance of g -MSCS (left) to an instance of f^{CS} -MSCS, which is equivalent to the CIRCULAR CONSENSUS STRING problem. Plots of the (polynomially bounded and order-independent) local cost functions for $k = 4$ are shown. Note that g_4 is obtained from f_6^{CS} by cropping the first two values in order to become grouping.

514 of 1's in column i of Δ . The cost of column i is $f_{2k-2}^{\text{CS}}(w + k - 2) = g_k(w)$. Hence,
 515 column i has the same cost in both solutions. This implies $\text{cost}_g(\Delta) = \text{cost}_{f^{\text{CS}}}(\Delta')$.

516 The converse direction is similar. Any multiple circular shift Δ' of s_1, \dots, s_{2k-2}
 517 can be restricted to a multiple circular shift Δ of s_1, \dots, s_k with the same cost. \square

518 **5. Consensus for Time Series: DTW-Mean.** In this section we consider the
 519 DTW-MEAN problem.

DTW-MEAN

Input: A list of k univariate rational time series x_1, \dots, x_k and $c \in \mathbb{Q}$.

520 **Question:** Is there a univariate rational time series z such that $\mathcal{F}(z) =$
 $\sum_{i=1}^k (\text{dtw}(z, x_i))^2 \leq c$?

521 We prove the following theorem, settling the complexity status of this prominent
 522 consensus problem in time series analysis.

523 **THEOREM 5.1.** DTW-MEAN on binary time series is

- 524 (i) NP-hard,
- 525 (ii) W[1]-hard with respect to the number k of input series, and
- 526 (iii) not solvable in $\rho(k) \cdot n^{o(k)}$ time for any computable function ρ unless the ETH
 527 fails.

528 The proof is based on a polynomial-time reduction from a special variant of f -
 529 MSCS for which hardness holds via [Theorem 3.11](#) in [Section 3](#). At this point we
 530 make crucial use of the fact that the reduction described in the proof of [Lemma 3.3](#)
 531 actually shows that it is hard to decide whether there is a multiple circular shift of
 532 cost at most c or whether all multiple circular shifts have cost at least $c + \varepsilon$ for some ε
 533 polynomially bounded in the number of strings. This gap of ε guarantees that a no-
 534 instance of f -MSCS is reduced to a no-instance of DTW-MEAN. Being polynomially
 535 bounded is required for ε in order to obtain a polynomial-time reduction (otherwise
 536 our constructed time series are too long).

537 Before proving [Theorem 5.1](#), we introduce some definitions. A *position* i in a time
 538 series x is an integer $1 \leq i \leq |x|$, its *value* is $x[i]$. The *distance* between two positions i
 539 and i' is $|i' - i|$. A *block* in a binary time series is a maximal subseries of consecutive 0's
 540 (a *0-block*) or 1's (a *1-block*). Blocks are also represented by integers, indicating their
 541 rank in the series (a series with n blocks has blocks $1, 2, \dots, n$). The *distance* between
 542 two blocks of rank y and y' is $|y' - y|$. Note that the notion of distance is different
 543 in the context of positions and blocks (even between size-1 blocks, as larger blocks in
 544 between increase the position distance).

545 Recall from [Section 2](#) that once the length of a mean z and the alignments to

546 the input time series are known, then its values are determined. In the binary case, a
 547 position i that is aligned to a 0's and b 1's contributes a cost of $az[i]^2 + b(1 - z[i])^2$
 548 to $\mathcal{F}(z)$. This cost is minimal for $z[i]$ being the arithmetic mean $b/(a + b)$ of all values
 549 aligned to position i . Thus, the *cost* of position i is

$$550 \quad C(i) = \min_{x \in \mathbb{Q}} (ax^2 + b(1 - x)^2) = \frac{ab}{a + b},$$

551 where the second equality follows from [Equation \(2.1\)](#). The reduction in the following
 552 proof of [Theorem 5.1](#) is from ϕ -MSCS for a polynomially bounded grouping function ϕ
 553 specifically chosen in relation to the above cost of a mean position.

554 *Proof.* We will reduce from ϕ -MSCS for a specially chosen cost function ϕ which
 555 we explain first. To this end, we briefly sketch the idea of the reduction. Given k bi-
 556 nary strings for ϕ -MSCS, we construct $k + 1$ binary time series: k time series encoding
 557 the original strings and a dummy time series containing a single 1. The construction
 558 is such that a column of a multiple circular shift with x 1's and $k - x$ 0's will corre-
 559 spond to a position in a time series that is aligned to $k + x$ 0's and one 1. Now, the
 560 cost of that column should be equal to the cost $(k + x)/(k + x + 1)$ of that position
 561 minus $k/(k + 1)$ (for technical reasons the cost should be 0 if $x = 0$). That is, we define
 562 $\phi: \{0, 1\}^* \rightarrow [0, 1]$ with $\phi((x_1, \dots, x_k)) = \phi_k(\sum_{j=1}^k x_j)$, where $\phi_k: \{0, \dots, k\} \rightarrow [0, 1]$
 563 with

$$564 \quad \phi_k(x) = \frac{k + x}{k + x + 1} - \frac{k}{k + 1} = \frac{x}{(k + x + 1)(k + 1)}.$$

565 Note that ϕ is a polynomially bounded grouping function since

$$566 \quad \phi'_k(x) = \frac{\phi_k(x) - \phi_k(0)}{x} = \frac{1}{(k + x + 1)(k + 1)}$$

567 is strictly decreasing (ϕ_k is strictly concave) with gap

$$568 \quad \varepsilon_k = \phi'_k(k - 1) - \phi'_k(k) = \frac{1}{(k + 1)(2k)(2k + 1)}$$

569 and range

$$570 \quad \mu_k = \phi'_k(1) = \frac{1}{(k + 2)(k + 1)}.$$

571 See [Figure 5](#) for an example of the functions ϕ_k and ϕ'_k . Hence, by [Theorem 3.11](#)
 572 hardness holds for ϕ -MSCS. We now give the polynomial-time reduction from ϕ -
 573 MSCS, see [Figure 6](#) for an illustration of the reduction.

574 *Reduction.* In the following, we assume to have an instance of ϕ -MSCS with
 575 $k \geq 15$ length- n binary strings $s_1, \dots, s_k \in \{A, B\}^n$ (where $A := 0$ and $B := 1$) and
 576 a target cost $0 \leq c \leq n$. We write ε for the gap ε_k of ϕ_k (note that $\varepsilon^{-1} \in O(k^3)$).
 577 The task is to decide whether there exists a multiple circular shift of cost at most c or
 578 whether all multiple circular shifts have cost at least $c + \varepsilon$ (the reduction in [Lemma 3.3](#)
 579 implies the hardness of this decision problem).

580 First, we encode characters A and B via certain binary strings. To this end,
 581 we define the number $m := 1600k \lceil c + \varepsilon \rceil$ and the binary strings $t_A := (10)^m$ and
 582 $t_B := 100(10)^{m-1}$, each string having m 0-blocks (all of length one, except for the
 583 first 0-block of t_B which has length two). The first 0-block of t_A and t_B is called a
 584 *coding* block (respectively, an A -coding or a B -coding block).

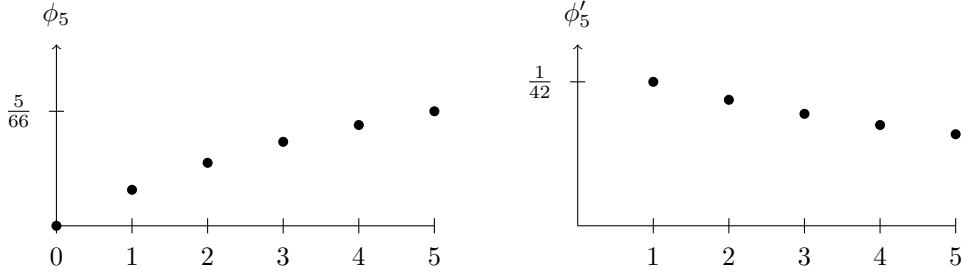


FIG. 5. Left: The function $\phi_5(x) = \frac{x}{6(6+x)}$. Right: The function $\phi'_5(x) = \frac{1}{6(6+x)}$.

585 Now, for each $1 \leq j \leq k$, let s'_j be the string obtained by concatenating the strings
 586 $t_{s_j[i]}$ for $1 \leq i \leq |s_j|$. The final time series x_j is obtained by concatenating r copies
 587 of s'_j , where

588
$$r := \left\lceil \frac{1}{\varepsilon} (3mnk + 2(c + \varepsilon)) \right\rceil + 1.$$

589 Note that each x_j contains $2mnr$ blocks and $|x_j| \leq \text{poly}(nk)$. We also define the *extra*
 590 series $x_{k+1} = (1)$ and set the target cost to

591
$$c' := r \left(c + \frac{mnk}{k+1} \right) + 3mnk.$$

592 For the correctness of this reduction we need to prove the following:

- 593 (i) If (s_1, \dots, s_k, c) is a yes-instance, that is, there exists a multiple circular shift Δ
 594 with $\text{cost}_\phi(\Delta) \leq c$, then there exists a time series z with $\mathcal{F}(z) \leq c'$.
 595 (ii) If (s_1, \dots, s_k, c) is a no-instance, that is, $\text{cost}_\phi(\Delta) \geq c + \varepsilon$ holds for every multiple
 596 circular shift Δ , then $\mathcal{F}(z) > c'$ holds for every time series z .

597 (i) *Yes-instance of ϕ -MSCS.* Consider a multiple circular shift $\Delta = (\delta_1, \dots, \delta_k)$
 598 of s_1, \dots, s_k with $\text{cost}_\phi(\Delta) \leq c$. Without loss of generality, we assume that $0 \leq \delta_j < n$
 599 holds for every $1 \leq j \leq k$.

600 We construct a time series z of length $2mn(r-1) + 2$ (also see Figure 6) such
 601 that $\mathcal{F}(z) \leq c'$. To this end, we describe the alignments between z and x_1, \dots, x_{k+1}
 602 (recall that this determines the values and costs of positions in z). For each $j =$
 603 $1, \dots, k$, the first position is aligned with the leftmost $2\delta_j m$ blocks of each x_j (or
 604 with the first block if $\delta_j = 0$) and the last position is aligned with the rightmost
 605 $2(n - \delta_j)m$ blocks of x_j . For each $1 < i < 2mn(r-1) + 2$, position i is aligned
 606 with the $(i-1 + 2\delta_j m)$ -th block of x_j . These positions are called *regular* positions,
 607 whereas the first and last position are called *extreme*. Clearly, all positions of z are
 608 also aligned with the single 1 in x_{k+1} .

609 Given the above alignments, the sum of costs of all positions in z is clearly an
 610 upper bound for $\mathcal{F}(z)$. The following two claims give an upper bound for this sum.

611 CLAIM 5.2. *The total cost of regular positions is at most $(r-1)(c + \frac{mnk}{k+1})$.*

612 *Proof.* Due to the alternation of 1- and 0-blocks in each x_j and the fact that
 613 $i + (2\delta_j m) \equiv i \pmod{2}$, it follows that the i -th regular position (which is $z[i+1]$) is
 614 mapped only to 1's if i is odd (*odd* position) or only to 0's (and the single 1 in x_{k+1})
 615 if i is even (*even* positions). Thus, odd positions have cost $\frac{(k+1) \cdot 0}{k+1} = 0$, and even
 616 positions have a cost depending on the size of the 0-blocks to which they are mapped.

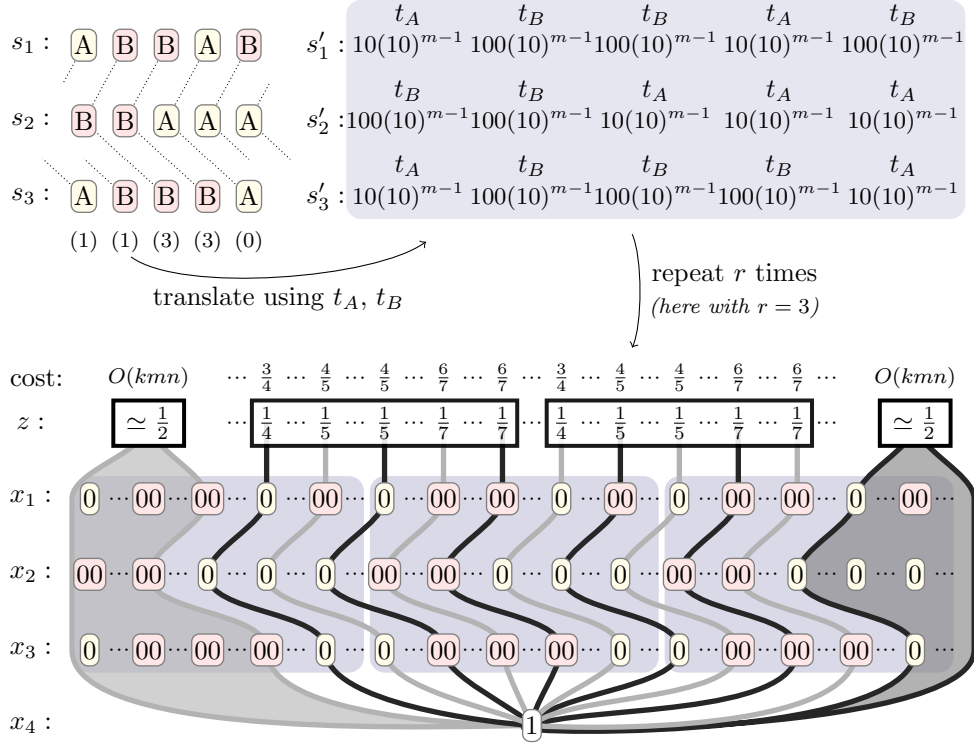


FIG. 6. *Top left:* Illustration of the reduction to DTW-MEAN from an instance of ϕ -MSCS with $k = 3$ and $n = 5$. An optimal circular shift $\Delta = (3, 2, 4)$ is indicated by dotted lines, and the number of B's in a shifted column is below each column. The total cost is $\text{cost}_\phi(\Delta) = 2\phi_3(1) + 2\phi_3(3) + \phi_3(0)$. *Top right:* The intermediate strings s'_1, s'_2, s'_3 encoding the original strings s_1, s_2, s_3 . *Bottom:* The resulting instance x_1, \dots, x_4 of DTW-MEAN (only coding blocks are shown) and an alignment to a time series z mimicking the circular shift Δ . The values of z are shown along with the cost of each position (positions that are only aligned to non-coding blocks are ignored, and contribute a background cost of either 0 or $\frac{3}{4}$). Note that the cost function ϕ is chosen so that the cost of a position aligned to k coding blocks equals the cost of the corresponding column of the original circular shift (plus the background cost $\frac{3}{4}$). For example, a position aligned to two A-coding blocks and one B-coding block has cost $\frac{4}{5} = \frac{3}{4} + \phi_3(1)$, where $\phi_3(1)$ is the cost of a column with two A's and one B in ϕ -MSCS. The value of m is chosen large enough to yield a large cost for misalignments, such as two consecutive coding blocks of the same series aligned together. The value of r is chosen such that only a periodic pattern ensures low cost of z , even though it requires to pay a high (but bounded) misalignment cost for the first and last positions.

617 Consider an even position i such that $i \bmod 2m \neq 2$. The i -th regular position is
 618 not aligned with any coding block in any x_j . Thus, it is aligned to k 0's and a single 1,
 619 and has cost $C(i+1) = \frac{k}{k+1}$. There are $(m-1)n(r-1)$ such positions, which thus
 620 contribute a total cost of

$$621 \quad (r-1) \frac{(m-1)nk}{k+1}.$$

622 For an even position i with $i \bmod 2m = 2$, the i -th regular position is aligned
 623 with a coding block in each x_j (except for the single 1 in x_{k+1}). Let $i = 2mi' + 2$.
 624 Then, $z[i+1]$ is aligned to coding blocks corresponding to column $i' \bmod n$ of Δ . If
 625 this column contains a A's and $k-a$ B's, then $z[i+1]$ is aligned to $a + 2(k-a)$ 0's

626 and a single 1 and has cost

$$627 \quad C(i+1) = \frac{a+2(k-a)}{a+2(k-a)+1} = \frac{2k-a}{2k-a+1} = \phi_k(k-a) + \frac{k}{k+1}.$$

628 Note that $\phi_k(k-a)$ is the cost of column $i' \bmod n$ of Δ . Hence, the overall cost of
629 the $(r-1)n$ regular positions i with $i \bmod 2m = 2$ is

$$630 \quad (r-1) \text{cost}_\phi(\Delta) + (r-1) \frac{nk}{k+1} \leq (r-1) \left(c + \frac{nk}{k+1} \right).$$

631 Overall, the regular positions have a total cost of at most

$$632 \quad (r-1) \frac{(m-1)nk}{k+1} + (r-1) \left(c + \frac{nk}{k+1} \right) = (r-1) \left(c + \frac{mnk}{k+1} \right). \quad \square$$

633 CLAIM 5.3. *The total cost of extreme positions is at most $2knm + 2$.*

634 *Proof.* Since $0 \leq \delta_j < n$ for $j \in \{1, \dots, k\}$, an extreme position $i \in \{1, |z|\}$ is
635 aligned to at most $2nm$ consecutive blocks in each x_j , thus accounting for at most
636 nm 1's in x_j . Moreover, position i is aligned with the additional 1 in x_{k+1} . Thus,
637 position i is aligned with at most $nmk + 1$ 1's, which implies $C(i) \leq nmk + 1$ (this
638 bound is achieved if $z[i] = 0$). \square

639 Combining Claims 5.2 and 5.3, we obtain

$$640 \quad \mathcal{F}(z) \leq \sum_{i=1}^{|z|} C(i) \leq (r-1) \left(c + \frac{mnk}{k+1} \right) + 2knm + 2 \leq c'.$$

641 Hence, $(x_1, \dots, x_{k+1}, c')$ is a yes-instance of DTW-MEAN.

642 (ii) *No-instance of ϕ -MSCS.* We assume that every multiple circular shift Δ
643 has $\text{cost}_\phi(\Delta) \geq c + \varepsilon$. Consider a fixed mean time series z (minimizing \mathcal{F}) together
644 with optimal warping paths between z and x_1, \dots, x_{k+1} . We show that $\mathcal{F}(z) > c'$.
645 We will do this hierarchically, starting with a lower bound on the cost of an individual
646 position of z . We then proceed to derive lower bounds for the cost of certain intervals
647 of positions until finally obtaining the desired lower bound on $\mathcal{F}(z)$. Before doing so,
648 we start with some preliminary observations about the structure of a mean.

649 *Structural Observations.* We say that position i of x_j is *matched* to position i'
650 of z if (i, i') is in the warping path between x_j and z . Clearly, the single position
651 in x_{k+1} is matched to every position of z . We write $\#_1(i)$ and $\#_0(i)$ respectively for
652 the number of positions with value 1 (resp. 0) among x_1, \dots, x_k (ignoring the extra
653 sequence x_{k+1}) which are matched to position i of z . Clearly, the cost of i is

$$654 \quad C(i) = \frac{\#_0(i)(\#_1(i) + 1)}{\#_0(i) + \#_1(i) + 1}.$$

655 We will use the following monotonicity property of the cost.

656 LEMMA 5.4. *For any $a \geq a' \geq 0$ and $b \geq b' \geq 1$, it holds $\frac{ab}{a+b} \geq \frac{a'b'}{a'+b'}$.*

657 *Proof.* It suffices to see that the partial derivatives

$$658 \quad \frac{\partial}{\partial a} \frac{ab}{a+b} = \frac{b^2}{(a+b)^2} \quad \text{and} \quad \frac{\partial}{\partial b} \frac{ab}{a+b} = \frac{a^2}{(a+b)^2}$$

659 are non-negative for $a \geq 0$ and $b \geq 1$. \square

660 We define some further notation for the remaining part of the proof. The *range*
 661 of a position i of z in x_j is the set of positions of x_j that are matched to i . The
 662 range is a subinterval of $\{1, 2, \dots, |x_j|\}$ and by construction of x_j , the corresponding
 663 subseries cannot have three consecutive 0's or two consecutive 1's. More precisely, its
 664 values alternate between 0 and 1, except for the (rare) occasions where it includes
 665 a B -coding block. The number of blocks of x_j intersecting the range of i in x_j is
 666 denoted $r_j(i)$. The number of B -coding blocks contained in the range of i in x_j is
 667 denoted $r_j^B(i)$ and we define $r^B(i) := \sum_{j=1}^k r_j^B(i)$.

668 A position i of z is called *0-simple* (resp. *1-simple*) if $\#_1(i) = 0$ (resp. $\#_0(i) = 0$).
 669 It is *simple* if it is 0- or 1-simple, and it is *bad* otherwise. Clearly, the cost of a 1-simple
 670 position is 0. For a 0-simple position i , we have $\#_1(i) = 0$ and $k \leq \#_0(i) \leq 2k$ (more
 671 precisely, $\#_0(i) = k + r^B(i)$ and $r^B(i) \leq k$). Thus, $\frac{k}{k+1} \leq C(i) \leq \frac{2k}{2k+1}$. Since we
 672 assumed $k \geq 15$, the cost of a 0-simple position is always contained in $[0.9, 1]$.

673 We continue with several structural observations regarding a mean.

674 **OBSERVATION 5.5.** *There exists a mean without consecutive 1-simple positions or*
 675 *consecutive 0-simple positions. Such a mean is called non-redundant.*

676 *Proof.* Any two consecutive 1-simple (or 0-simple) positions of a non-redundant
 677 mean z have consecutive or intersecting ranges in each x_j with the same value (1
 678 or 0). Hence, they can be merged into one single 1-simple (or 0-simple) position.

679 Since the warping of the other positions in z remains unchanged, we focus on the
 680 cost of the merged position. For 1-simple positions, the cost remains unchanged (both
 681 solutions yield a cost of 0 for the 1-simple positions). For 0-simple positions, the cost
 682 of the two 0-simple positions in the original solution is at least 0.9 each. However, the
 683 cost of the merged 0-simple position is at most 1, which is a contradiction to z being
 684 a mean. \square

685 In the following, we assume z to be a non-redundant mean.

686 We say a block b of some input x_j is *matched* (*fully matched*) to a position i
 687 in z if some position (all positions) in b is (are) matched to i . That is, a matched
 688 block intersects the range of i and a fully matched block is included in the range of i .
 689 Note that the distinction is only relevant for B -coding blocks, as all other blocks have
 690 size 1. Moreover, the number of B -coding blocks that are fully matched to a position i
 691 equals $r^B(i)$.

692 **OBSERVATION 5.6.** *For a non-redundant mean z , any B -coding block of some x_j*
 693 *that is not fully matched to a position in z is matched to at least one bad position*
 694 *in z .*

695 *If any block is matched to two consecutive positions in z , then at least one of them*
 696 *is bad.*

697 *Proof.* Consider a B -coding block b and all positions of z to which it is matched.
 698 There are at least two of them, which cannot all be 0-simple (since z is non-redundant).
 699 Also, none of them can be 1-simple (since b is a 0-block). Thus, at least one of them
 700 is bad.

701 We prove the contrapositive of the second statement: If two simple positions have
 702 a common block matched to them, then they are both a -simple, $a \in \{0, 1\}$, and cannot
 703 be consecutive in a non-redundant mean. \square

704 We now introduce an *assignment* relation between a block b of some input series
 705 and a position i of z . We say that b is *assigned* to the position i if i is the leftmost
 706 simple position to which b is fully matched (if any), or (if no such simple position

exists) if i is the leftmost bad position to which b is matched. Note that any size-1 block is fully matched to at least one position (simple or bad), and every size-2 block is either fully matched to a simple position or matched to a bad position (by [Observation 5.6](#)). Thus, every block is assigned to exactly one position.

For a position i of z , we introduce the following quantities:

$$\begin{aligned} \diamond_0(i) &:= \text{number of 0-blocks in } x_1, \dots, x_k \text{ matched to } i, \\ \diamond_1(i) &:= \text{number of 1-blocks in } x_1, \dots, x_k \text{ matched to } i, \\ q(i) &:= \text{number of } B\text{-coding blocks assigned to } i, \\ g(i) &:= \begin{cases} 0, & \text{if } i \text{ is simple} \\ \max\{1, r_1(i) - 1, \dots, r_k(i) - 1\}, & \text{if } i \text{ is bad} \end{cases}. \end{aligned}$$

We quickly observe the following:

(I) If i is simple, then $q(i) = r^B(i)$.

Proof. By definition, every B -coding block assigned to i is fully matched to i , that is, $q(i) \leq r^B(i)$. Furthermore, a B -coding block can be fully matched to only one position. Thus, if i is simple, then all B -coding blocks fully matched to it are also assigned to it, that is, $q(i) \geq r^B(i)$.

(II) If i is 0-simple, then $\diamond_0(i) = k$, $\diamond_1(i) = 0$, and $C(i) = \frac{k+q(i)}{k+q(i)+1}$.

Proof. A 0-simple position has exactly one 0-block in each x_1, \dots, x_k matched to it. The cost follows from (I).

(III) If i is 1-simple, then $\diamond_0(i) = 0$, $\diamond_1(i) = k$, and $C(i) = 0$.

Proof. A 1-simple position has exactly one 1-block in each x_1, \dots, x_k matched to it.

(IV) If i is bad, then $g(i) \leq 2 \min\{\diamond_1(i), \diamond_0(i)\}$.

Proof. For a bad position i , we have

$$\min\{\diamond_1(i), \diamond_0(i)\} \geq \max\left\{1, \frac{1}{2} \max_{j=1, \dots, k} (r_j(i) - 1)\right\}.$$

(V) If i is bad, then $C(i) \geq \frac{1}{2} \min\{\diamond_1(i) + 1, \diamond_0(i)\}$.

Proof. Let $\eta := \min\{\diamond_1(i) + 1, \diamond_0(i)\} \geq 1$. Then i is aligned to at least η 0's and at least η 1's. Thus, by [Lemma 5.4](#), $C(i) \geq \frac{\eta^2}{2\eta} = \frac{\eta}{2}$.

(VI) For every position i , it holds that $|\diamond_0(i) - \diamond_1(i)| \leq k$ and $\diamond_0(i) + \diamond_1(i) \geq k$.

Proof. For each x_j , $1 \leq j \leq k$, the difference between the number of 0-blocks matched to i and the number of 1-blocks matched to i is at most one. Clearly, there is at least one 0- or 1-block matched to i in each x_j .

Cost of a Single Position. We consider a fixed position i of z . For simplification, we write $\diamond_0 := \diamond_0(i)$, $\diamond_1 := \diamond_1(i)$, $q := q(i)$, $g := g(i)$, and $C := C(i)$. The goal is to provide a lower bound for C that can be decomposed into the following elements:

- a *background cost* $C_{\text{back}}(i) := \frac{\diamond_0}{k+1}$,
- a *coding cost* $C_{\text{code}}(i) := \phi_k(q)$ reflecting the extra cost induced by a matched coding block,
- a *gap cost* of $C_{\text{gap}}(i) := 0.01g$, which is 0 if i is simple, and which increases with the number of blocks matched to i if i is bad.

CLAIM 5.7. *The cost C of position i is at least $LB(\diamond_0, q, g, k)$, defined as follows:*

$$\begin{aligned} LB(\diamond_0, q, g, k) &:= C_{\text{back}}(i) + C_{\text{code}}(i) + C_{\text{gap}}(i) \\ &= \frac{\diamond_0}{k+1} + \phi_k(q) + 0.01g. \end{aligned}$$

751 *Proof.* In the following, we write LB for $LB(\diamond_0, q, g, k)$. We prove $C \geq LB$ by
752 case distinction.

753 For a 0-simple position, by (II), we have $\diamond_0 = k$ and $g = 0$. Thus,

$$754 \quad LB = \frac{k}{k+1} + \phi_k(q) + 0 = \frac{k+q}{k+q+1} = C.$$

755 For a 1-simple position, by (III), we have $\diamond_0 = 0$, $q = 0$, and $g = 0$. Thus,

$$756 \quad LB = 0 + \phi_k(0) + 0 = 0 = C.$$

757 For a bad position, we have $\diamond_0 \geq 1$ and $\diamond_1 \geq 1$. First, note that

$$760 \quad LB - 0.01g = \frac{\diamond_0}{k+1} + \frac{q}{(k+q+1)} \cdot \frac{1}{(k+1)} \leq \frac{\diamond_0}{k} + \frac{1}{k+1} \leq 2\frac{\diamond_0}{k},$$

761 that is, $LB \leq 2\frac{\diamond_0}{k} + 0.01g$.

762 We now use this upper bound for the following three sub-cases. First, if $\diamond_0 \leq \diamond_1$,
763 then we have

$$764 \quad \begin{aligned} 765 \quad C &\geq \frac{1}{2}\diamond_0 \quad (\text{by (V)}) \\ 766 \quad &\geq \frac{2}{k}\diamond_0 + 0.02\diamond_0 \quad (\text{since } k \geq 15) \\ 767 \quad &\geq LB \quad (\text{by (IV)}). \end{aligned}$$

768 Second, if $6 \leq \diamond_1 < \diamond_0$, we have

$$770 \quad \begin{aligned} 771 \quad C &\geq \frac{1}{2}\diamond_1 \quad (\text{by (V)}) \\ &\geq 2 + \left(\frac{2}{k} + 0.02\right)\diamond_1 \quad (\text{since } k \geq 15 \text{ and } \diamond_1 \geq 6) \\ 772 \quad &= 2 + \frac{2\diamond_1}{k} + 0.02\diamond_1 \\ 773 \quad &\geq \frac{2\diamond_0}{k} + 0.02\diamond_1 \quad (\text{since } \diamond_0 \leq \diamond_1 + k \text{ by (VI)}) \\ 774 \quad &\geq LB \quad (\text{by (IV)}). \end{aligned}$$

775 Finally, if $\diamond_1 \leq 5 < \diamond_0$ (note that $\diamond_1 < \diamond_0 \leq 5$ is not possible since, by (VI), we
776 have $\diamond_1 + \diamond_0 \geq k$ and we assumed $k \geq 15$), then $k - 5 \leq \diamond_0 \leq k + 5$ (by (VI)), and
777 $g \leq 10$ (by (IV)). We have

$$779 \quad C = \frac{(\diamond_1 + 1)\#_0(i)}{\diamond_1 + 1 + \#_0(i)} \geq \frac{2(k-5)}{k-3} \geq 1.66$$

780 using Lemma 5.4, with $\diamond_1 \geq 1$, $\#_0(i) \geq \diamond_0 \geq k - 5$, and $k \geq 15$. On the other side,
781 we have

$$782 \quad \begin{aligned} 783 \quad LB &= \frac{\diamond_0}{k+1} + \phi_k(q) + 0.01g \\ 784 \quad &\leq \frac{5}{k+1} + \frac{k}{k+1} + \phi_k(q) + 0.1 \quad (\text{since } \diamond_0 \leq k+5 \text{ and } g \leq 10) \\ 785 \quad &\leq \frac{5}{k} + 1.1 \quad (\text{since } \frac{k}{k+1} + \phi_k(q) \leq 1) \\ 786 \quad &\leq 1.44 \quad (\text{since } k \geq 15) \\ 787 \quad &\leq C. \end{aligned}$$

□

789 *Cost of (Ir)regular Intervals.* We aim at computing lower bounds on the cost of
 790 intervals of positions. Two positions i, i' of the mean z at distance $|i - i'| = \ell$ form
 791 an *irregular pair* if for some x_j a block b is matched to i and a block b' is matched
 792 to i' such that either $|b - b'| \leq \ell - \frac{m}{2k}$ or $|b - b'| \geq \ell + \frac{m}{2k}$. An interval I of positions
 793 in z is called *regular* if it does not contain any irregular pair (otherwise it is called
 794 *irregular*). The background, coding, and gap cost of I is the sum of the respective
 795 costs of its positions. The structure of regular intervals allows us to bound the coding
 796 cost from below using the minimum cost of the original ϕ -MSCS instance. Irregular
 797 intervals contain bad positions, which allow us to derive a lower bound on their gap
 798 cost.

799 We first introduce some notation: a position i of z is *j-coding* if there is a coding
 800 block b in x_j such that b is assigned to i ; it is *coding* if it is *j-coding* for some j ,
 801 otherwise it is *non-coding*. A non-coding position i is *free* if all positions at distance
 802 at most $\frac{m}{2k} + 2$ from i are non-coding (see Figure 7). We first make the following
 803 technical claim before proving the main bound on the coding cost of regular intervals
 804 (Claim 5.9).

805 CLAIM 5.8. *In a regular interval, if two positions $i < i'$ are at distance at most*
 806 *$2\alpha m - \frac{m}{2k}$ for some α , then, for any j , there are at most α j -coding positions in $[i, i']$.*

807 *Conversely, if i and i' are at distance at least $2\alpha m + \frac{m}{2k} + 1$, then, for any j , there*
 808 *are at least α j -coding positions in $[i, i']$.*

809 *Proof.* Fix $j \in \{1, \dots, k\}$ and consider the first block b in x_j matched to i and
 810 the last block b' in x_j matched to i' (then $b' \geq b$). Note that all *j-coding* positions
 811 in $[i, i']$ have been assigned a distinct coding block in $[b, b']$. Since i and i' are not an
 812 irregular pair, it holds that

$$813 \quad b' - b < i' - i + \frac{m}{2k} \leq 2\alpha m.$$

814 That is, $b' < b + 2\alpha m$, and thus x_j contains at most α coding blocks in $[b, b']$. These
 815 coding blocks are assigned to positions in $[i, i']$. Hence, $[i, i']$ contains at most α
 816 *j-coding* positions.

817 For the other direction, consider again blocks b and b' as above. In this case
 818 there is a slight difference: If block b or b' is coding, then it might be assigned to a
 819 coding position outside of the interval $[i, i']$, which then would not count in the lower
 820 bound. Thus, we consider only blocks strictly between b and b' , among which all
 821 coding blocks are assigned to a coding position in $[i, i']$. Since $i' - i \geq 2\alpha m + \frac{m}{2k} + 1$,
 822 we have $b' - b > i' - i + \frac{m}{2k} \geq 2\alpha m + 1$, so there are at least $2\alpha m$ blocks strictly
 823 between b and b' , including at least α coding blocks. These are assigned to at least α
 824 *j-coding* positions in $[i, i']$. \square

825 CLAIM 5.9. *The coding cost of a length- ℓ regular interval I is at least*

$$826 \quad C_{\text{code}}(I) \geq \left(\frac{\ell + 1}{2mn} - 2 \right) (c + \varepsilon).$$

827 *Proof.* Let I be a regular length- ℓ interval. Assume that $\ell \geq 3mn$ (otherwise the
 828 stated lower bound is negative which is trivial).

829 The first part of the proof consists of splitting interval I into consecutive length-
 830 $2m$ segments, each one containing exactly one *j-coding* position for each j . To this
 831 end, a few positions need to be cropped from both ends of I . In other words, we need
 832 to find a good starting point (a free position) close to the left end of I .

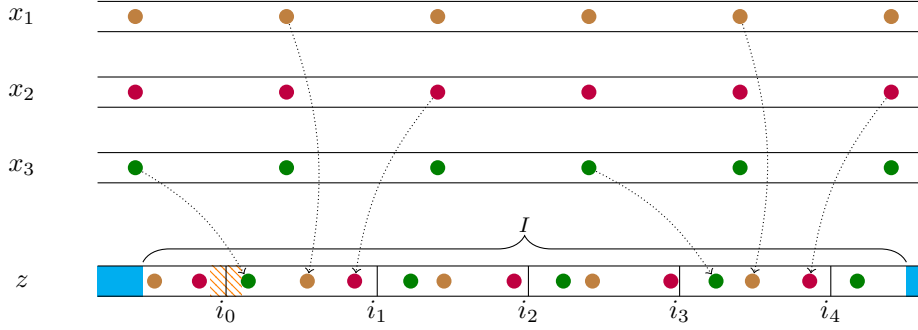


FIG. 7. Decomposition of a regular interval (I) into segments delimited at positions i_l . Coding blocks in x_j are indicated with colored bullets, as well as their assigned positions in z (dotted arcs). Position i_0 is free (no coding position within the striped area), and each segment $[i_l, i_{l+1}]$ contains exactly one j -coding position from each x_j . These coding positions correspond to columns of a multiple circular shift of the ϕ -MSCS input strings. The coding cost of n consecutive segments can be bounded from below by the minimum cost of the ϕ -MSCS instance.

833 Consider the first position i of I and position $i' := i + 2m - \frac{m}{2k}$ (in I). By
 834 [Claim 5.8](#), for any j , there is at most one j -coding position in $[i, i']$ (and so at most k
 835 such coding positions in total). Accordingly, the coding positions split the interval
 836 $[i, i']$ into at most $(k + 1)$ disjoint intervals of non-coding positions, with total size at
 837 least $(i' - i + 1) - k = 2m - \frac{m}{2k} - k + 1$. Hence, there is one interval of non-coding
 838 positions in $[i, i']$ with size at least

$$839 \quad \frac{2m - \frac{m}{2k} - k + 1}{k + 1} \geq \frac{2m - \frac{m}{2k}}{k + 1} - 1 \geq \frac{m}{k} + 5.$$

840 Note that the second inequality above is equivalent to

$$841 \quad \frac{m(2 - \frac{1}{2k})}{k + 1} \geq \frac{m}{k} + 6 \quad \Leftrightarrow \quad m \geq \frac{6k(k + 1)}{k - \frac{3}{2}},$$

842 which is true since $m \geq 1600k$ and $k \geq 15$. Hence, this interval contains a free
 843 position, denoted i_0 , with $i + \frac{m}{2k} + 2 \leq i_0 \leq i + 2m - \frac{m}{2k} - 2$.

844 Let $\lambda := \lfloor \frac{\ell}{2m} \rfloor - 2$ and $i_l := i_0 + 2lm$ for $0 < l \leq \lambda$. Note that $\lambda \geq n$ and that
 845 every i_l is in I (following the assumption on ℓ).

846 We fix some input series x_j , $1 \leq j \leq k$. Intuitively, positions i_l are the cutting
 847 points of our segments within interval I (see [Figure 7](#)). We now aim at showing that
 848 there is exactly one j -coding position in each segment $[i_{l-1}, i_l - 1]$. First, consider
 849 positions $h = i_0 - \frac{m}{2k} - 2$ and $i_l - 1$. By [Claim 5.8](#), since $(i_l - 1) - h = 2lm + \frac{m}{2k} + 1$,
 850 interval $[h, i_l - 1]$ contains at least l j -coding positions. Since i_0 is free, these coding
 851 positions cannot be before i_0 (as $i_0 - h = \frac{m}{2k} + 2$), so they are in $[i_0, i_l - 1]$. Consider
 852 now positions $h' = i_0 + \frac{m}{2k} + 1$ and $i_l - 1$. By [Claim 5.8](#), since $(i_l - 1) - h' = 2lm - \frac{m}{2k}$,
 853 there are at most l j -coding positions in $[h', i_l - 1]$, and therefore at most l j -coding
 854 positions in $[i_0, i_l - 1]$. That is, there are exactly l j -coding positions in $[i_0, i_l - 1]$.

855 Overall, there is exactly one j -coding position in $[i_{l-1}, i_l - 1]$ for every $0 < l \leq \lambda$
 856 and every j . We write $C_{l,j}$ for the corresponding coding block in x_j .

857 Let q_l be the number of B -coding blocks among $C_{l,1}, \dots, C_{l,k}$. Then,

$$858 \quad \sum_{h=i_{l-1}}^{i_l-1} q(h) = q_l,$$

859 and thus

$$\begin{aligned}
 860 \quad \sum_{h=i_{l-1}}^{i_l-1} \phi_k(q(h)) &= \sum_{h=i_{l-1}}^{i_l-1} \frac{q(h)}{(k+q(h)+1)(k+1)} \\
 861 \quad &\geq \sum_{h=i_{l-1}}^{i_l-1} \frac{q(h)}{(k+q_l+1)(k+1)} = \phi_k(q_l). \\
 862
 \end{aligned}$$

863 Let δ_j be such that $C_{1,j}$ is the δ_j -th coding block of x_j . Then, $C_{l,j}$ is the (δ_j+l-1) -
 864 th coding block of x_j , for every $0 < l \leq \lambda$. Note that the coding blocks $C_{l,1}, \dots, C_{l,k}$
 865 correspond to the $((l-2) \bmod n + 1)$ -th column in the multiple circular shift $\Delta =$
 866 $(\delta_1 \bmod n, \dots, \delta_k \bmod n)$ of s_1, \dots, s_k . Thus, q_l is the number of B 's in this column
 867 and $\phi_k(q_l)$ is the corresponding cost of this column.

868 That is, for any integer a with $0 < a \leq \lambda - n$, the sum $\sum_{l=a}^{a+n-1} \phi_k(q_l)$ corresponds
 869 to the cost of some multiple circular shift of s_1, \dots, s_k . Since, by assumption, every
 870 multiple circular shift of s_1, \dots, s_k has cost at least $c + \varepsilon$, we have

$$871 \quad \sum_{l=a}^{a+n-1} \phi_k(q_l) \geq c + \varepsilon.$$

872 We can now compute the lower bound on the coding cost of interval I . To this end,
 873 we first extract $\lfloor \frac{\lambda}{n} \rfloor \geq 1$ length- $2mn$ subintervals of I , each consisting of n segments
 874 of the form $[i_{l-1}, i_l - 1]$. It follows

$$\begin{aligned}
 875 \quad C_{\text{code}}(I) &= \sum_{h \in I} \phi_k(q(h)) \geq \sum_{h=i_0}^{i_\lambda-1} \phi_k(q(h)) \\
 876 \quad &\geq \sum_{a=0}^{\lfloor \frac{\lambda}{n} \rfloor - 1} \sum_{l=an}^{an+n-1} \sum_{h=i_{l-1}}^{i_l-1} \phi_k(q(h)) \\
 877 \quad &\geq \sum_{a=0}^{\lfloor \frac{\lambda}{n} \rfloor - 1} \sum_{l=an}^{an+n-1} \phi_k(q_l) \\
 878 \quad &\geq \sum_{a=0}^{\lfloor \frac{\lambda}{n} \rfloor - 1} c + \varepsilon \\
 879 \quad &= \left\lfloor \frac{\lambda}{n} \right\rfloor (c + \varepsilon) \\
 880 \quad &\geq \left(\frac{\lambda}{n} - 1 \right) (c + \varepsilon) \\
 881 \quad &= \left(\frac{\lfloor \frac{\ell}{2m} \rfloor - 2}{n} - 1 \right) (c + \varepsilon) \\
 882 \quad &\geq \left(\frac{\frac{\ell}{2m} - 3}{n} - 1 \right) (c + \varepsilon) \\
 883 \quad &\geq \left(\frac{\ell + 1}{2mn} - 2 \right) (c + \varepsilon), \\
 884
 \end{aligned}$$

885 where we assume $n \geq 4$ for the last inequality. \square

886 Next, we prove a lower bound on the gap cost of an interval I . In the follow-
887 ing, the total number of matches between blocks of x_1, \dots, x_k and positions in I is
888 denoted $W(I) = \sum_{i \in I} (\diamond_0(i) + \diamond_1(i))$.

889 CLAIM 5.10. *For any interval I of length ℓ , the gap cost $C_{\text{gap}}(I)$ fulfills*

$$890 \quad C_{\text{gap}}(I) \geq \frac{1}{100} \left(\frac{W(I)}{k} - \ell \right).$$

891 Moreover, if I is irregular, then

$$892 \quad C_{\text{gap}}(I) \geq \frac{m}{400k}.$$

893 *Proof.* For the first lower bound, it suffices to note that for any position i (simple
894 or bad), it holds

$$895 \quad g(i) \geq \max_{1 \leq j \leq k} r_j(i) - 1 \geq \frac{W(i)}{k} - 1,$$

896 where $W(i) := \diamond_0(i) + \diamond_1(i)$.

897 For the second lower bound, consider an irregular pair $i < i'$ in I and an integer j
898 such that a block b in x_j is matched to i and a block $b' > b$ in x_j is matched to i'
899 where $|(b' - b) - (i' - i)| > \frac{m}{2k}$.

900 If $b' - b > i' - i + \frac{m}{2k}$, then

$$901 \quad \sum_{h=i}^{i'} g(h) \geq \sum_{h=i}^{i'} (r_j(h) - 1) = \sum_{h=i}^{i'} r_j(h) - (i' - i + 1) \geq b' - b - (i' - i) > \frac{m}{2k} > \frac{m}{4k}.$$

902 If $b' - b < i' - i - \frac{m}{2k}$, then there are at least $\frac{m}{2k}$ pairs of consecutive positions having
903 the same block in x_j matched to them, and for every such pair at least one of the two
904 positions is bad (by [Observation 5.6](#)). Since any bad position may be counted in at
905 most two such pairs, the interval has at least $\frac{m}{4k}$ bad positions. Hence, using $g(h) \geq 1$
906 for bad positions, we obtain $\sum_{h=i}^{i'} g(h) \geq \frac{m}{4k}$. \square

907 *Cost of a Mean.* To obtain a lower bound for $\mathcal{F}(z)$, we now partition the posi-
908 tions $[1, |z|]$ into minimal irregular intervals (from left). To this end, let $\alpha_1 := 1$ and
909 let β_1 be the position such that the interval $[\alpha_1, \beta_1]$ is irregular (if such a position
910 does not exist, then $\beta_1 := |z|$) and $[\alpha_1, \beta_1 - 1]$ is regular. If $\beta_1 < |z|$, then we continue
911 analogously and define $\alpha_2 := \beta_1 + 1$ and β_2 to be the position such that $[\alpha_2, \beta_2]$ is
912 irregular and $[\alpha_2, \beta_2 - 1]$ is regular. This procedure is repeated until we obtain a
913 partition

$$914 \quad [\alpha_1 := 1, \beta_1], [\alpha_2 := \beta_1 + 1, \beta_2], \dots, [\alpha_L := \beta_{L-1} + 1, \beta_L := |z|]$$

915 of $[1, |z|]$ into $L \geq 1$ intervals of which the first $L - 1$ are irregular and the last is
916 possibly regular. The following lower bounds hold.

917 CLAIM 5.11. *For $1 \leq l < L$, it holds that*

$$918 \quad C_{\text{gap}}([\alpha_l, \beta_l]) + C_{\text{code}}([\alpha_l, \beta_l]) \geq (c + \varepsilon) \frac{W([\alpha_l, \beta_l])}{2knm}.$$

919 *For the coding and gap costs of $[\alpha_L, \beta_L]$, it holds that*

$$920 \quad C_{\text{gap}}([\alpha_L, \beta_L]) + C_{\text{code}}([\alpha_L, \beta_L]) \geq (c + \varepsilon) \frac{W([\alpha_L, \beta_L])}{2knm} - 2(c + \varepsilon).$$

921 *Proof.* Consider an interval $[\alpha_l, \beta_l]$, $1 \leq l \leq L$, and let ℓ be its length. Let
 922 $W := W([\alpha_l, \beta_l])$. Since $[\alpha_l, \beta_l - 1]$ is a regular interval of length $\ell - 1$, by [Claim 5.9](#),
 923 we have the following lower bound on the coding cost:

$$924 \quad (5.1) \quad C_{\text{code}}([\alpha_l, \beta_l]) \geq C_{\text{code}}([\alpha_l, \beta_l - 1]) \geq \left(\frac{\ell}{2mn} - 2 \right) (c + \varepsilon).$$

925 For $l < L$, we combine both bounds on the gap cost of [Claim 5.10](#) (by averaging their
 926 values):

$$927 \quad C_{\text{gap}}([\alpha_l, \beta_l]) \geq \frac{1}{200} \left(\frac{W}{k} - \ell \right) + \frac{m}{800k}.$$

928 Using $m \geq 1600k(c + \varepsilon)$ (by definition) and $m \geq 100 \frac{c + \varepsilon}{n}$, we obtain

$$929 \quad (5.2) \quad C_{\text{gap}}([\alpha_l, \beta_l]) \geq \frac{c + \varepsilon}{2nm} \left(\frac{W}{k} - \ell \right) + 2(c + \varepsilon).$$

930 The sum of Inequations [5.1](#) and [5.2](#) yields the claimed lower bound.

931 For interval $[\alpha_L, \beta_L]$, we use the general lower bound from [Claim 5.10](#), which
 932 yields

$$933 \quad (5.3) \quad C_{\text{gap}}([\alpha_L, \beta_L]) \geq \frac{1}{100} \left(\frac{W}{k} - \ell \right) \geq \frac{c + \varepsilon}{2mn} \left(\frac{W}{k} - \ell \right).$$

934 The sum of Inequations [5.1](#) and [5.3](#) yields the claimed lower bound. \square

935 Finally, to finish the proof of [Theorem 5.1](#), we show that the mean z has high
 936 cost, that is, $(x_1, \dots, x_{k+1}, c')$ is a no-instance of DTW-MEAN.

937 **CLAIM 5.12.** $\mathcal{F}(z) > c'$.

938 *Proof.* Using [Claim 5.7](#) on each position of $I := [1, |z|]$, we obtain the following
 939 lower bound

$$940 \quad \mathcal{F}(z) = \sum_{i=1}^{|z|} C(i) \geq C_{\text{code}}(I) + C_{\text{gap}}(I) + C_{\text{back}}(I).$$

941 For the coding and gap cost of I , we use [Claim 5.11](#) together with the fact that
 942 all $2knmr$ blocks of x_1, \dots, x_k are involved in at least one match with a position of z ,
 943 which yields $W(I) = \sum_{l=1}^L W([\alpha_l, \beta_l]) \geq 2knmr$. Thus,

$$944 \quad C_{\text{code}}(I) + C_{\text{gap}}(I) \geq (c + \varepsilon)r - 2(c + \varepsilon).$$

945 The overall background cost is $C_{\text{back}}(I) = \sum_{i=1}^{|z|} \frac{\diamond_0(i)}{k+1}$. Since overall there are
 946 $knmr$ 0-blocks in x_1, \dots, x_k , and each of those is matched to at least one position
 947 of z , we have $\sum_{i=1}^{|z|} \diamond_0(i) \geq knmr$ and thus

$$948 \quad C_{\text{back}}(I) \geq \frac{nmrk}{k+1}.$$

949 Combining the two bounds above yields

$$950 \quad C_{\text{code}}(I) + C_{\text{gap}}(I) + C_{\text{back}}(I) \geq \frac{nmrk}{k+1} + (c + \varepsilon)r - 2(c + \varepsilon).$$

951 Since $\varepsilon r > 3mnk + 2(c + \varepsilon)$, we get

$$952 \quad \mathcal{F}(z) > \frac{nmrk}{k+1} + rc + 3mnk = c'. \quad \square$$

953 Since the above reduction is a polynomial-time reduction from ϕ -MSCS where
 954 the resulting number of time series is linear in the number of strings in the ϕ -MSCS
 955 instance, [Theorem 5.1](#) now follows from [Theorem 3.11](#). \square

956 Closing this section, we remark that Buchin et al. [[8](#), Theorem 7] recently obtained
 957 the same hardness results as in [Theorem 5.1](#) for the problem of computing an average
 958 series z that minimizes

$$959 \quad \mathcal{F}_p^q(z) := \sum_{j=1}^k \left(\min_{p_j \in \mathcal{P}_{|x_j|, |z|}} \sum_{(u,v) \in p_j} |x_j[u] - z[v]|^p \right)^{q/p}$$

960

961 for all integers $p, q \geq 1$. Their reduction, however, builds time series containing three
 962 different values. Hence, [Theorem 5.1](#) yields a stronger hardness on binary inputs for
 963 $p = q = 2$. Note that if also the mean is restricted to be a binary time series, then
 964 the problem is solvable in polynomial time [[6](#), [34](#)].

965 **6. Conclusion.** Shedding light on the computational complexity of prominent
 966 consensus problems in stringology and time series analysis, we proved several tight
 967 computational hardness results for circular string alignment problems and time series
 968 averaging in dynamic time warping spaces. Notably, we have shown that the compu-
 969 tational complexity of consensus string problems can drastically change (that is, they
 970 become hard) when considering *circular* strings instead of classic strings. Our results
 971 imply that these problems with a rich set of applications are intractable in the worst
 972 case (even on binary data). Hence, it is unlikely to find algorithms which significantly
 973 improve the worst-case running times of the best known algorithms. This now partly
 974 justifies the use of heuristics as has been done for a long time in many real-world
 975 applications.

976 We conclude with some open questions and directions for future work.

- 977 • We conjecture that the idea of the reduction for f -MSCS can be used to
 978 prove the same hardness result for most non-linear (polynomially bounded)
 979 order-independent cost functions (note that f -MSCS is trivially solvable if f_k
 980 is linear since every shift has the same cost). Proving a complexity dichotomy
 981 with respect to the cost function is a worthwhile goal.
- 982 • From an algorithmic point of view, it would be nice to improve the constant in
 983 the exponent of the running time for DTW-MEAN, that is, to find algorithms
 984 running in $O(n^{\alpha k})$ time for small α . In particular, we ask to find an $O(n^k)$ -
 985 time algorithm for DTW-MEAN.
- 986 • What about the parameter maximum sequence length n ? Are the considered
 987 problems polynomial-time solvable if n is a constant? Are they even fixed-
 988 parameter tractable with respect to n ?
- 989 • Finally, can the hardness result for averaging time series with respect to (p, q) -
 990 DTW by Buchin et al. [[8](#), Theorem 7] be strengthened to binary inputs?

991 **Acknowledgments.** We are very grateful to two anonymous reviewers of *SIAM*
 992 *Journal on Discrete Mathematics* whose very detailed and constructive feedback
 993 helped to improve the presentation of the paper significantly.

994 **References.**

- 995 [1] A. ABBOUD, A. BACKURS, AND V. V. WILLIAMS, *Tight hardness results for*
 996 *LCS and other sequence similarity measures*, in Proceedings of the 56th Annual
 997 IEEE Symposium on Foundations of Computer Science (FOCS '15), IEEE, 2015,
 998 pp. 59–78.

999 [2] S. AGHABOZORGI, A. S. SHIRKHORSHIDI, AND T. Y. WAH, *Time-series clustering – A decade review*, Information Systems, 53 (2015), pp. 16–38.
 1000
 1001 [3] N. ARICA, *Cyclic sequence comparison using dynamic warping*, in Proceedings of
 1002 the International Conference on Image and Video Retrieval (CIVR '05), vol. 3568
 1003 of LNCS, Springer, 2005, pp. 328–335.
 1004 [4] L. A. K. AYAD AND S. P. PISSIS, *MARS: improving multiple circular sequence*
 1005 *alignment using refined sequences*, BMC Genomics, 18 (2017), p. 86.
 1006 [5] C. BARTON, C. S. ILIOPOULOS, R. KUNDU, S. P. PISSIS, A. RETHA, AND
 1007 F. VAYANI, *Accurate and efficient methods to improve multiple circular sequence*
 1008 *alignment*, in Proceedings of the 14th International Symposium on Experimental
 1009 Algorithms (SEA '15), vol. 9125 of LNCS, Springer, 2015, pp. 247–258.
 1010 [6] M. BRILL, T. FLUSCHNIK, V. FROESE, B. JAIN, R. NIEDERMEIER, AND
 1011 D. SCHULTZ, *Exact mean computation in dynamic time warping spaces*, Data
 1012 Mining and Knowledge Discovery, 33 (2019), pp. 252–291.
 1013 [7] K. BRINGMANN AND M. KÜNNEMANN, *Quadratic conditional lower bounds for*
 1014 *string problems and dynamic time warping*, in Proceedings of the 56th Annual
 1015 IEEE Symposium on Foundations of Computer Science (FOCS '15), IEEE, 2015,
 1016 pp. 79–97.
 1017 [8] K. BUCHIN, A. DRIEMEL, AND M. STRUIJS, *On the hardness of comput-*
 1018 *ing an average curve*, CoRR, abs/1902.08053 (2019). Accepted for publica-
 1019 tion at the 17th Scandinavian Symposium and Workshops on Algorithm Theory
 1020 (SWAT '20).
 1021 [9] L. BULTEAU, F. HÜFFNER, C. KOMUSIEWICZ, AND R. NIEDERMEIER, *Multi-*
 1022 *variate algorithmics for NP-hard string problems*, Bulletin of the EATCS, 114
 1023 (2014).
 1024 [10] J. CHEN, B. CHOR, M. FELLOWS, X. HUANG, D. W. JUEDES, I. A. KANJ,
 1025 AND G. XIA, *Tight lower bounds for certain parameterized NP-hard problems*,
 1026 Information and Computation, 201 (2005), pp. 216–231.
 1027 [11] M. CUTURI AND M. BLONDEL, *Soft-DTW: a differentiable loss function for time-*
 1028 *series*, in Proceedings of the 34th International Conference on Machine Learning
 1029 (ICML '17), vol. 70 of Proceedings of Machine Learning Research, PMLR, 2017,
 1030 pp. 894–903.
 1031 [12] M. CYGAN, F. V. FOMIN, Ł. KOWALIK, D. LOKSHTANOV, D. MARX,
 1032 M. PILIPCZUK, M. PILIPCZUK, AND S. SAURABH, *Parameterized Algorithms*,
 1033 Springer, 2015.
 1034 [13] M. R. FELLOWS, J. GRAMM, AND R. NIEDERMEIER, *On the parameterized*
 1035 *intractability of motif search problems*, Combinatorica, 26 (2006), pp. 141–167.
 1036 [14] F. FERNANDES, L. PEREIRA, AND A. T. FREITAS, *CSA: An efficient algorithm*
 1037 *to improve circular DNA multiple alignment*, BMC Bioinformatics, 10 (2009),
 1038 p. 230.
 1039 [15] M. FRANCES AND A. LITMAN, *On covering problems of codes*, Theory of Com-
 1040 puting Systems, 30 (1997), pp. 113–119.
 1041 [16] V. FROESE, B. JAIN, M. RYMAR, AND M. WELLER, *Fast exact dynamic time*
 1042 *warping on run-length encoded time series*, CoRR, abs/1903.03003 (2020).
 1043 [17] O. GOLD AND M. SHARIR, *Dynamic time warping and geometric edit distance:*
 1044 *Breaking the quadratic barrier*, ACM Transactions on Algorithms, 14 (2018),
 1045 pp. 50:1–50:17.
 1046 [18] J. GRAMM, R. NIEDERMEIER, AND P. ROSSMANITH, *Fixed-parameter algo-*
 1047 *rithms for Closest String and related problems*, Algorithmica, 37 (2003), pp. 25–
 1048 42.

- 1049 [19] R. GROSSI, C. S. ILIOPOULOS, R. MERCAS, N. PISANTI, S. P. PISSIS,
1050 A. RETHA, AND F. VAYANI, *Circular sequence comparison: algorithms and ap-*
1051 *plications*, Algorithms for Molecular Biology, 11 (2016).
- 1052 [20] R. IMPAGLIAZZO AND R. PATURI, *On the complexity of k -SAT*, Journal of Com-
1053 *puter and System Sciences*, 62 (2001), pp. 367–375.
- 1054 [21] E. KEOGH AND C. A. RATANAMAHATANA, *Exact indexing of dynamic time*
1055 *warping*, Knowledge and Information Systems, 7 (2005), pp. 358–386.
- 1056 [22] W. KUSZMAUL, *Dynamic time warping in strongly subquadratic time: Algo-*
1057 *rithms for the low-distance regime and approximate evaluation*, in Proceedings of
1058 *the 46th International Colloquium on Automata, Languages, and Programming*
1059 *(ICALP '19)*, 2019, pp. 80:1–80:15.
- 1060 [23] T. LEE, J. C. NA, H. PARK, K. PARK, AND J. S. SIM, *Finding consensus and*
1061 *optimal alignment of circular strings*, Theoretical Computer Science, 468 (2013),
1062 pp. 92–101.
- 1063 [24] M. LI, B. MA, AND L. WANG, *Finding similar regions in many sequences*,
1064 *Journal of Computer and System Sciences*, 65 (2002), pp. 73–96.
- 1065 [25] M. LI, B. MA, AND L. WANG, *On the closest string and substring problems*,
1066 *Journal of the ACM*, 49 (2002), pp. 157–171.
- 1067 [26] D. MARX, *Closest substring problems with small distances*, SIAM Journal on
1068 *Computing*, 38 (2008), pp. 1382–1410.
- 1069 [27] R. A. MOLLINEDA, E. VIDAL, AND F. CASACUBERTA, *Cyclic sequence align-*
1070 *ments: Approximate versus optimal techniques*, International Journal of Pattern
1071 *Recognition and Artificial Intelligence*, 16 (2002), pp. 291–299.
- 1072 [28] M. MOREL, C. ACHARD, R. KULPA, AND S. DUBUISSON, *Time-series averaging*
1073 *using constrained dynamic time warping with tolerance*, Pattern Recognition, 74
1074 (2018), pp. 77–89.
- 1075 [29] V. PALAZÓN-GONZÁLEZ AND A. MARZAL, *On the dynamic time warping of*
1076 *cyclic sequences for shape retrieval*, Image and Vision Computing, 30 (2012),
1077 pp. 978–990.
- 1078 [30] J. PAPARRIZOS AND L. GRAVANO, *Fast and accurate time-series clustering*,
1079 *ACM Transactions on Database Systems*, 42 (2017), pp. 8:1–8:49.
- 1080 [31] F. PETITJEAN, G. FORESTIER, G. I. WEBB, A. E. NICHOLSON, Y. CHEN, AND
1081 E. KEOGH, *Faster and more accurate classification of time series by exploiting*
1082 *a novel dynamic time warping averaging algorithm*, Knowledge and Information
1083 *Systems*, 47 (2016), pp. 1–26.
- 1084 [32] F. PETITJEAN AND P. GAŃCARSKI, *Summarizing a set of time series by averag-*
1085 *ing: From Steiner sequence to compact multiple alignment*, Theoretical Computer
1086 *Science*, 414 (2012), pp. 76–91.
- 1087 [33] F. PETITJEAN, A. KETTERLIN, AND P. GAŃCARSKI, *A global averaging method*
1088 *for dynamic time warping, with applications to clustering*, Pattern Recognition,
1089 44 (2011), pp. 678–693.
- 1090 [34] N. SCHAAR, V. FROESE, AND R. NIEDERMEIER, *Faster binary mean compu-*
1091 *tation under dynamic time warping*, in Proceedings of the 31th Annual Symposi-
1092 *um on Combinatorial Pattern Matching (CPM '20)*, LIPIcs, Schloss Dagstuhl
1093 *- Leibniz-Zentrum für Informatik*, 2020. Accepted for publication (arXiv version
1094 available).
- 1095 [35] D. SCHULTZ AND B. JAIN, *Nonsmooth analysis and subgradient methods for av-*
1096 *eraging in dynamic time warping spaces*, Pattern Recognition, 74 (2018), pp. 340–
1097 358.
- 1098 [36] S. SOHEILY-KHAH, A. DOUZAL-CHOUAKRIA, AND E. GAUSSIÉ, *Generalized*

- 1099 *k*-means-based clustering for temporal data under weighted and kernel time warp,
1100 Pattern Recognition Letters, 75 (2016), pp. 63–69.
- 1101 [37] S. WILL AND P. F. STADLER, *A common framework for linear and cyclic multi-*
1102 *ple sequence alignment problems*, in Proceedings of the 14th International Work-
1103 shop on Algorithms in Bioinformatics (WABI '14), vol. 8701 of LNCS, Springer,
1104 2014, pp. 135–147.