



Elastic-Degenerate String Matching via Fast Matrix Multiplication

Giulia Bernardini, Pawel Gawrychowski, Nadia Pisanti, Solon Pissis,
Giovanna Rosone

► To cite this version:

Giulia Bernardini, Pawel Gawrychowski, Nadia Pisanti, Solon Pissis, Giovanna Rosone. Elastic-Degenerate String Matching via Fast Matrix Multiplication. SIAM Journal on Computing, 2022, 51 (3), pp.549-576. 10.1137/20M1368033 . hal-03676475

HAL Id: hal-03676475

<https://hal.inria.fr/hal-03676475>

Submitted on 24 May 2022

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.

Elastic-Degenerate String Matching via Fast Matrix Multiplication

Giulia Bernardini¹, Paweł Gawrychowski², Nadia Pisanti³, Solon P. Pissis⁴, and Giovanna Rosone⁵

¹Department of Informatics, Systems and Communication (DISCo), University of Milan - Bicocca, Italy, giulia.bernardini@unimib.it

²Institute of Computer Science, University of Wrocław, Poland, gawry@cs.uni.wroc.pl

³Department of Computer Science, University of Pisa, Italy and ERABLE Team, INRIA, France, pisanti@di.unipi.it

⁴CWI, Amsterdam, The Netherlands, solon.pissis@cwi.nl

⁵Department of Computer Science, University of Pisa, Italy, giovanna.rosone@unipi.it

Abstract

An elastic-degenerate (ED) string is a sequence of n sets of strings of total length N , which was recently proposed to model a set of similar sequences. The ED string matching (EDSM) problem is to find all occurrences of a pattern of length m in an ED text. The EDSM problem has recently received some attention in the combinatorial pattern matching community, and an $\mathcal{O}(nm^{1.5}\sqrt{\log m} + N)$ -time algorithm is known [Aoyama et al., CPM 2018]. The standard assumption in the prior work on this question is that N is substantially larger than both n and m , and thus we would like to have a linear dependency on the former. Under this assumption, the natural open problem is whether we can decrease the 1.5 exponent in the time complexity, similarly as in the related (but, to the best of our knowledge, not equivalent) *word break* problem [Backurs and Indyk, FOCS 2016].

Our starting point is a conditional lower bound for the EDSM problem. We use the popular combinatorial Boolean Matrix Multiplication (BMM) conjecture stating that there is no truly subcubic *combinatorial* algorithm for BMM [Abboud and Williams, FOCS 2014]. By designing an appropriate reduction we show that a combinatorial algorithm solving the EDSM problem in $\mathcal{O}(nm^{1.5-\epsilon} + N)$ time, for any $\epsilon > 0$, refutes this conjecture. Our reduction should be understood as an indication that decreasing the exponent requires fast matrix multiplication.

String periodicity and fast Fourier transform are two standard tools in string algorithms. Our main technical contribution¹ is that we successfully combine these tools with fast matrix multiplication to design a non-combinatorial $\tilde{\mathcal{O}}(nm^{\omega-1} + N)$ -time algorithm for EDSM, where ω denotes the matrix multiplication exponent and the $\tilde{\mathcal{O}}(\cdot)$ notation suppresses polylog factors. To the best of our knowledge, we are the first to combine these tools. In particular, using the fact that $\omega < 2.373$ [Le Gall, ISSAC 2014; Williams, STOC 2012], we obtain an $\mathcal{O}(nm^{1.373} + N)$ -time algorithm for EDSM. An important building block in our solution, that might find applications in other problems, is a method of selecting a small set of length- ℓ substrings of the pattern, called anchors, so that any occurrence of a string from an ED text set contains at least one but not too many such anchors inside.

¹A preliminary version of this work appeared in ICALP 2019 [12].

1 Introduction

Boolean matrix multiplication (BMM) is one of the most fundamental computational problems. Apart from its theoretical interest, it has a wide range of applications [30, 32, 40, 51, 58]. BMM is also the core combinatorial part of integer matrix multiplication. In both problems, we are given two $\mathcal{N} \times \mathcal{N}$ matrices and we are to compute \mathcal{N}^2 values. Integer matrix multiplication can be performed in *truly subcubic* time, i.e., in $\mathcal{O}(\mathcal{N}^{3-\epsilon})$ operations over the field, for some $\epsilon > 0$. The fastest known algorithms for this problem run in $\mathcal{O}(\mathcal{N}^{2.373})$ time [47, 60]. These algorithms are known as algebraic: they rely on the ring structure of matrices over the field.

There also exists a different family of algorithms for the BMM problem known as combinatorial. Their focus is on unveiling the combinatorial structure in the Boolean matrices to reduce redundant computations. A series of results [7, 9, 16] culminating in an $\hat{\mathcal{O}}(\mathcal{N}^3 / \log^4 \mathcal{N})$ -time algorithm [64, 65] (the $\hat{\mathcal{O}}(\cdot)$ notation suppresses polyloglog factors) has led to the popular combinatorial BMM conjecture stating that there is no combinatorial algorithm for BMM working in time $\mathcal{O}(\mathcal{N}^{3-\epsilon})$, for any $\epsilon > 0$ [2]. There has been ample work on applying this conjecture to obtain BMM hardness results: see, e.g., [2, 18, 36, 45, 46, 48, 54].

String matching is another fundamental problem, asking to find all fragments of a string text of length n that match a string pattern of length m . This problem has several linear-time solutions [24]. In many real-world applications, it is often the case that letters at some positions are either unknown or uncertain. A way of representing these positions is with a subset of the alphabet Σ . Such a representation is called *degenerate string*. A special case of a degenerate string is when at such unknown or uncertain positions the only subset of the alphabet allowed is the whole alphabet. These special degenerate strings are more commonly known as strings with wildcards. The first efficient algorithm for a text and a pattern, where both may contain wildcards, was published by Fischer and Paterson in 1974 [31]. It has undergone several improvements since then [21, 22, 39, 42]. The first efficient algorithm for a standard text and a degenerate pattern, which may contain any non-empty subset of the alphabet, was published by Abrahamson in 1987 [3], followed by several practically efficient algorithms [37, 52, 63].

Degenerate letters are used in the IUPAC notation [41] to represent a position in a DNA sequence that can have multiple possible alternatives. These are used to encode the consensus of a population of sequences [4, 33, 57] in a multiple sequence alignment (MSA). In the presence of insertions or deletions in the MSA, we may need to consider alternative representations. Consider the following MSA of three closely-related sequences (on the left):

GCAACGGGTA--TT

$$\begin{array}{l} \text{GCAACGGGTATATT} \\ \text{GCACCTGG----TT} \end{array} \quad \tilde{T} = \{ \text{GCA} \} \cdot \left\{ \begin{array}{c} \text{A} \\ \text{C} \end{array} \right\} \cdot \{ \text{C} \} \cdot \left\{ \begin{array}{c} \text{G} \\ \text{T} \end{array} \right\} \cdot \{ \text{GG} \} \cdot \left\{ \begin{array}{c} \text{TA} \\ \text{TATA} \\ \varepsilon \end{array} \right\} \cdot \{ \text{TT} \}$$

These sequences can be compacted into a single sequence \tilde{T} of sets of strings (on the right) containing some deterministic and some non-deterministic segments. A non-deterministic segment is a finite set of deterministic strings and may contain the empty string ε corresponding to a deletion. The total number of segments is the *length* of \tilde{T} and the total number of letters is the *size* of \tilde{T} . We denote the length by $n = |\tilde{T}|$ and the size by $N = ||\tilde{T}||$.

This representation has been defined in [38] by Iliopoulos et al. as an *elastic-degenerate* (ED) string. Being a sequence of subsets of Σ^* , it can be seen as a generalization of a degenerate string. The natural problem that arises is finding all matches of a deterministic pattern P in an ED text \tilde{T} . This is the *elastic-degenerate string matching* (EDSM) problem. Since its introduction in 2017 [38], it has attracted some attention in the combinatorial pattern matching community, and a series of results have been published. The simple algorithm by Iliopoulos et al. [38] for EDSM was first improved by Grossi et al. in the same year, who showed that, for a pattern of length m , the EDSM problem can be solved *on-line* in $\mathcal{O}(nm^2 + N)$ time [35]; on-line means that it reads the text segment-by-segment and reports an occurrence as soon as this is detected. This result

was improved by Aoyama et al. [6] who presented an $\mathcal{O}(nm^{1.5}\sqrt{\log m} + N)$ -time algorithm. An important feature of these bounds is their *linear dependency* on N . A different branch of on-line algorithms waiving the linear-dependency restriction exists [19, 20, 35, 53]. Moreover, the EDSM problem has been considered under Hamming and edit distance [13]. Recent results on founder block graphs [49] can also be casted on elastic-degenerate strings.

A question with a somewhat similar flavor is the *word break* problem. We are given a dictionary \mathcal{D} , $m = |\mathcal{D}|$, and a string S , $n = |S|$, and the question is whether we can split S into fragments that appear in \mathcal{D} (the same element of \mathcal{D} can be used multiple times). Backurs and Indyk [8] designed an $\tilde{\mathcal{O}}(nm^{1/2-1/18} + m)$ -time algorithm for this problem². Bringmann et al. [15] improved this to $\tilde{\mathcal{O}}(nm^{1/3} + m)$ and showed that this is optimal for combinatorial algorithms by a reduction from k -Clique. Their algorithm uses fast Fourier transform (FFT), and so it is not clear whether it should be considered combinatorial. While this problem seems similar to EDSM, there does not seem to be a direct reduction and so their lower bound does not immediately apply.

Our Results. It is known that BMM and triangle detection (TD) in graphs either both have truly subcubic combinatorial algorithms or none of them do [62]. Recall also that the currently fastest algorithm with linear dependency on N for the EDSM problem runs in $\mathcal{O}(nm^{1.5}\sqrt{\log m} + N)$ time [6]. In this paper we prove the following two theorems.

Theorem 1. *If the EDSM problem can be solved in $\mathcal{O}(nm^{1.5-\epsilon} + N)$ time, for any $\epsilon > 0$, with a combinatorial algorithm, then there exists a truly subcubic combinatorial algorithm for TD.*

Arguably, the notion of combinatorial algorithms is not clearly defined, and Theorem 1 should be understood as an indication that in order to achieve a better complexity one should use fast matrix multiplication. Indeed, there are examples where a lower bound conditioned on BMM was helpful in constructing efficient algorithms using fast matrix multiplication [1, 14, 17, 26, 50, 61, 66]. We successfully design such a non-combinatorial algorithm by combining three ingredients: a string periodicity argument, FFT, and fast matrix multiplication. While periodicity is the usual tool in combinatorial pattern matching [25, 43, 44] and using FFT is also not unusual (for example, it often shows up in approximate string matching [3, 5, 21, 34]), to the best of our knowledge, we are the first to combine these with fast matrix multiplication. Specifically, we show the following result for the EDSM problem, where ω denotes the matrix multiplication exponent.

Theorem 2. *The EDSM problem can be solved on-line in $\tilde{\mathcal{O}}(nm^{\omega-1} + N)$ time.*

In order to obtain a faster algorithm for the EDSM problem, we focus on the *active prefixes* (AP) problem that lies at the heart of all current solutions [6, 35]. In the AP problem, we are given a string P of length m and a set of arbitrary prefixes $P[1..i]$ of P , called *active prefixes*, stored in a bit vector U so that $U[i] = 1$ if $P[1..i]$ is active. We are further given a set \mathcal{S} of strings of total length N and we are asked to compute a bit vector V which stores the new set of active prefixes of P . A new active prefix of P is a concatenation of $P[1..i]$ (such that $U[i] = 1$) and some element of \mathcal{S} .

Using the algorithmic framework introduced in [35], EDSM is addressed by solving an instance of the AP problem per each segment i of the ED text corresponding to set \mathcal{S} of the AP problem. Hence, an $\mathcal{O}(f(m) + N_i)$ solution for the AP problem (with N_i being the size of a single segment of the ED text) implies an $\mathcal{O}(nf(m) + N)$ solution of EDSM, as $f(m)$ is repeated n times and $N = \sum_{i=1}^n N_i$. The algorithm of [6] solves the AP problem in $\mathcal{O}(m^{1.5}\sqrt{\log m} + N_i)$ time leading to $\mathcal{O}(nm^{1.5}\sqrt{\log m} + N)$ time for the EDSM problem. Our algorithm partitions the strings of each segment i of the ED text into three types according to a periodicity criterion, and then solves a restricted instance of the AP problem for each of the types. In particular, we solve the AP problem in $\tilde{\mathcal{O}}(m^{\omega-1} + N_i)$ time leading to $\tilde{\mathcal{O}}(nm^{\omega-1} + N)$ time for the EDSM problem.

²The $\tilde{\mathcal{O}}(\cdot)$ notation suppresses polylog factors.

Given this connection between the two problems and, in particular, between their size parameter N , in the rest of the paper we will denote with N also the parameter N_i of the AP problem.

An important building block in our solution that might find applications in other problems is a method of selecting a small set of length- ℓ substrings of the pattern, called *anchors*, so that any relevant occurrence of a string from an ED text set contains at least one but not too many such anchors inside. This is obtained by rephrasing the question in a graph-theoretical language and then generalizing the well-known fact that an instance of the hitting set problem with m sets over $[n]$, each of size at least k , has a solution of size $\mathcal{O}(n/k \cdot \log m)$. While the idea of carefully selecting some substrings of the same length is not new (for example Kociumaka et al. [44] used it to design a data structure for pattern matching queries on a string), our setting is different and hence so is the method of selecting these substrings.

In addition to the conditional lower bound for the EDSM problem (Theorem 1), which also appeared in [12], we also show here the following conditional lower bound for the AP problem.

Theorem 3. *If the AP problem can be solved in $\mathcal{O}(m^{1.5-\epsilon} + N)$ time, for any $\epsilon > 0$, with a combinatorial algorithm, then there exists a truly subcubic combinatorial algorithm for the BMM problem.*

Roadmap. Section 2 provides the necessary definitions and notation as well as the algorithmic toolbox used throughout the paper. In Section 3 we prove our lower bound result for the AP problem (Theorem 3). The lower bound result for the EDSM problem is proved in Section 4 (Theorem 1). In Section 5 we present our algorithm for EDSM (Theorem 2); this is the most technically involved part of the paper.

2 Preliminaries

Let $T = T[1]T[2] \dots T[n]$ be a string of length $|T| = n$ over a finite ordered alphabet Σ of size $|\Sigma| = \sigma$. For two positions i and j on T , we denote by $T[i..j] = T[i] \dots T[j]$ the substring of T that starts at position i and ends at position j (it is of length 0 if $j < i$). By ε we denote the empty string of length 0. A prefix of T is a substring of the form $T[1..j]$, and a suffix of T is a substring of the form $T[i..n]$. T^r denotes the reverse of T , that is, $T[n]T[n-1] \dots T[1]$. We say that a string X is a power of a string Y if there exists an integer $k > 1$, such that X is expressed as k consecutive concatenations of Y , denoted by $X = Y^k$. A period of a string X is any integer $p \in [1, |X|]$ such that $X[i] = X[i+p]$ for every $i = 1, 2, \dots, |X| - p$, and *the period*, denoted by $\text{per}(X)$, is the smallest such p . We call a string X *strongly periodic* if $\text{per}(X) \leq |X|/4$.

Lemma 1 ([29]). *If p and q are both periods of the same string X , and additionally $p+q \leq |X|+1$, then $\gcd(p, q)$ is also a period of X .*

A *trie* is a tree in which every edge is labeled with a single letter, and every two edges outgoing from the same node have different labels. The label of a node u in such a tree T , denoted by $\mathcal{L}(u)$, is defined as the concatenation of the labels of all the edges on the path from the root of T to u . Thus, the label of the root of T is ε , and a trie is a representation of a set of strings consisting of the labels of all its leaves. By replacing each path p consisting of nodes with exactly one child by an edge labeled by the concatenation of the labels of the edges of p we obtain a *compact trie*. The nodes of the trie that are removed after this transformation are called *implicit*, while the remaining ones are referred to as *explicit*. The suffix tree of a string S is the compact trie representing all suffixes of $S\$$, $\$ \notin \Sigma$, where instead of explicitly storing the label $S[i..j]$ of an edge we represent it by the pair (i, j) .

A *heavy path decomposition* of a tree T is obtained by selecting, for every non-leaf node $u \in T$, its child v such that the subtree rooted at v is the largest. This decomposes the nodes of T into node-disjoint paths, with each such path p (called a heavy path) starting at some node, called

the *head* of p , and ending at a leaf. An important property of such a decomposition is that the number of distinct heavy paths above any leaf (that is, intersecting the path from a leaf to the root) is only logarithmic in the size of T [56].

Let $\tilde{\Sigma}$ denote the set of all finite non-empty subsets of Σ^* . Previous works (cf. [6, 13, 35, 38, 53]) define $\tilde{\Sigma}$ as the set of all finite non-empty subsets of Σ^* excluding $\{\varepsilon\}$ but we waive here the latter restriction as it has no algorithmic implications. An *elastic-degenerate string* $\tilde{T} = \tilde{T}[1] \dots \tilde{T}[n]$, or ED string, over alphabet Σ , is a string over $\tilde{\Sigma}$, i.e., an ED string is an element of $\tilde{\Sigma}^*$, and hence each $\tilde{T}[i]$ is a set of strings.

Let \tilde{T} denote an ED string of length n , i.e. $|\tilde{T}| = n$. We assume that for any $1 \leq i \leq n$, the set $\tilde{T}[i] \in \tilde{\Sigma}$ is implemented as an array and can be accessed by an index, i.e., $\tilde{T}[i] = \{\tilde{T}[i][k] \mid k = 1, \dots, |\tilde{T}[i]|\}$. For any $\tilde{\sigma} \in \tilde{\Sigma}$, $||\tilde{\sigma}||$ denotes the total length of all strings in $\tilde{\sigma}$, and for any ED string \tilde{T} , $||\tilde{T}||$ denotes the total length of all strings in all $\tilde{T}[i]$ s. We will denote $N_i = \sum_{k=1}^{|\tilde{T}[i]|} |\tilde{T}[i][k]|$ the total length of all strings in $\tilde{T}[i]$ and $N = \sum_{i=1}^n ||\tilde{T}[i]||$ the *size* of \tilde{T} . An ED string \tilde{T} can be thought of as a compact representation of the set of strings $\mathcal{A}(\tilde{T})$ which is the Cartesian product of all $\tilde{T}[i]$ s; that is, $\mathcal{A}(\tilde{T}) = \tilde{T}[1] \times \dots \times \tilde{T}[n]$ where $A \times B = \{xy \mid x \in A, y \in B\}$ for any sets of strings A and B .

For any ED string \tilde{X} and a pattern P , we say that P *matches* \tilde{X} if:

1. $|\tilde{X}| = 1$ and P is a substring of some string in $\tilde{X}[1]$, or,
2. $|\tilde{X}| > 1$ and $P = P_1 \dots P_{|\tilde{X}|}$, where P_1 is a suffix of some string in $\tilde{X}[1]$, $P_{|\tilde{X}|}$ is a prefix of some string in $\tilde{X}[|\tilde{X}|]$, and $P_i \in \tilde{X}[i]$, for all $1 < i < |\tilde{X}|$.

We say that an occurrence of a string P ends at position j of an ED string \tilde{T} if there exists $i \leq j$ such that P matches $\tilde{T}[i] \dots \tilde{T}[j]$. We will refer to string P as the *pattern* and to ED string \tilde{T} as the *text*. We define the main problem considered in this paper.

ELASTIC-DEGENERATE STRING MATCHING (EDSM)

INPUT: A string P of length m and an ED string \tilde{T} of length n and size $N \geq m$.

OUTPUT: All positions in \tilde{T} where at least one occurrence of P ends.

Example 1. Pattern $P = \text{GTAT}$ ends at positions 2, 6, and 7 of the following text \tilde{T} .

$$\tilde{T} = \left\{ \begin{array}{c} \text{AT} \end{array} \text{GTA} \right\} \cdot \left\{ \begin{array}{c} \text{A} \\ \text{T} \end{array} \right\} \cdot \left\{ \text{C} \right\} \cdot \left\{ \begin{array}{c} \text{G} \\ \text{T} \end{array} \right\} \cdot \left\{ \text{CG} \right\} \cdot \left\{ \begin{array}{c} \text{TA} \\ \text{TATA} \\ \varepsilon \end{array} \right\} \cdot \left\{ \begin{array}{c} \text{TATGC} \\ \text{TTTTA} \end{array} \right\}$$

Aoyama et al. [6] obtained an on-line $\mathcal{O}(nm^{1.5}\sqrt{\log m} + N)$ -time algorithm by designing an efficient solution for the following problem.

ACTIVE PREFIXES (AP)

INPUT: A string P of length m , a bit vector U of size m , a set \mathcal{S} of strings of total length N .

OUTPUT: A bit vector V of size m with $V[j] = 1$ if and only if there exists $S \in \mathcal{S}$ and $i \in [1, m]$, $U[i] = 1$, such that $P[1..i] \cdot S = P[1..i + |S|]$ and $j = i + |S|$.

In more detail, given an ED text $\tilde{T} = \tilde{T}[1] \dots \tilde{T}[n]$, one should consider an instance of the AP problem per each $\tilde{T}[i]$. Hence, an $\mathcal{O}(f(m) + N_i)$ solution for AP (N_i being the size of $\tilde{T}[i]$) implies an $\mathcal{O}(n \cdot f(m) + N)$ solution for EDSM, as $f(m)$ is repeated n times and $N = \sum_{i=1}^n N_i$. We provide an example of the AP problem.

Example 2. Let $P = \text{ababbababab}$ of length $m = 11$, $U = 01000100000$, and $\mathcal{S} = \{\varepsilon, \text{ab}, \text{abb}, \text{ba}, \text{baba}\}$. We have that $V = 01011101010$.

For our lower bound results we rely on BMM and the following closely related problem.

BOOLEAN MATRIX MULTIPLICATION (BMM)

INPUT: Two $\mathcal{N} \times \mathcal{N}$ Boolean matrices A and B .

OUTPUT: $\mathcal{N} \times \mathcal{N}$ Boolean matrix C , where $C[i, j] = \bigvee_k (A[i, k] \wedge B[k, j])$.

TRIANGLE DETECTION (TD)

INPUT: Three $\mathcal{N} \times \mathcal{N}$ Boolean matrices A, B and C .

OUTPUT: Are there i, j, k such that $A[i, j] = B[j, k] = C[k, i] = 1$?

An algorithm is called *truly subcubic* if it runs in $\mathcal{O}(\mathcal{N}^{3-\epsilon})$ time, for some $\epsilon > 0$. TD and BMM either both have truly subcubic combinatorial algorithms, or none of them do [62].

3 AP Conditional Lower Bound

To investigate the hardness of the EDSM problem, we first show that an $\mathcal{O}(m^{1.5-\epsilon} + N)$ -time solution to the active prefixes problem, that constitutes the core of the solutions proposed in [6, 35], would imply a truly subcubic combinatorial algorithm for Boolean matrix multiplication (BMM). We recall that in the AP problem, we are given a string P of length m and a set of prefixes $P[1..i]$ of P , called *active prefixes*, stored in a bit vector U ($U[i] = 1$ if and only if $P[1..i]$ is active). We are further given a set \mathcal{S} of strings of total length N and we are asked to compute a bit vector V storing the new set of active prefixes of P : a prefix of P that extends $P[1..i]$ (such that $U[i] = 1$) with some element of \mathcal{S} . Of course, we can solve BMM by working over integers and using one of the fast matrix multiplication algorithms; plugging in the best known bounds results in an $\mathcal{O}(\mathcal{N}^{2.373})$ -time algorithm [47, 60]. However, such an algorithm is not *combinatorial*, i.e., it uses *algebraic* methods. In comparison, the best known combinatorial algorithm for BMM works in $\hat{\mathcal{O}}(\mathcal{N}^3 / \log^4 \mathcal{N})$ time [65]. This leads to the following popular conjecture.

Conjecture 1 ([2]). There is no combinatorial algorithm for the BMM problem working in time $\mathcal{O}(\mathcal{N}^{3-\epsilon})$, for any $\epsilon > 0$.

Aoyama et al. [6] showed that the AP problem can be solved in $\mathcal{O}(m^{1.5} \sqrt{\log m} + N)$ time for constant-sized alphabets. Together with some standard string-processing techniques applied similarly as in [35], this is then used to solve the EDSM problem by creating an instance of the AP problem for every set $\tilde{T}[i]$ of \tilde{T} , i.e., with $\mathcal{S} = \tilde{T}[i]$.

We argue that, unless Conjecture 1 is false, the AP problem cannot be solved faster than $\mathcal{O}(m^{1.5-\epsilon} + N)$, for any $\epsilon > 0$, with a combinatorial algorithm (note that the algorithm of Aoyama et al. [6] uses FFT, and so it is not completely clear whether it should be considered to be combinatorial). We show this by a reduction from combinatorial BMM. Assume that, for the AP problem, we seek combinatorial algorithms with the running time $\mathcal{O}(m^{1.5-\epsilon} + N)$, i.e., with linear dependency on the total length of the strings. We need to show that such an algorithm implies that the BMM problem can be solved in $\mathcal{O}(\mathcal{N}^{3-\epsilon'})$ time, for some $\epsilon' > 0$, with a combinatorial algorithm, thus implying that Conjecture 1 is false.

Theorem 3. *If the AP problem can be solved in $\mathcal{O}(m^{1.5-\epsilon} + N)$ time, for any $\epsilon > 0$, with a combinatorial algorithm, then there exists a truly subcubic combinatorial algorithm for the BMM problem.*

Proof. Recall that in the BMM problem the matrices are denoted by A and B . In order to compute $C = A \times B$, we need to find, for every $i, j = 1, \dots, \mathcal{N}$, an index k such that $A[i, k] = 1$ and $B[k, j] = 1$. To this purpose, we split matrix A into blocks of size $\mathcal{N} \cdot L$ and B into blocks of size $L \cdot L$. This corresponds to considering values of j and k in intervals of size L , and clearly there are \mathcal{N}/L such intervals. Matrix B is thus split into $(\mathcal{N}/L)^2$ blocks, giving rise to an equal

number of instances of the AP problem, each one corresponding to an interval of j and an interval of k . This creates $(\mathcal{N}/L)^2$ blocks in matrix B ; we will thus create $(\mathcal{N}/L)^2$ separate instances of the AP problem corresponding to an interval of j and an interval of k . We will now describe the instance corresponding to the (K, J) -th block, where $1 \leq K, J \leq \mathcal{N}/L$.

We build the string P of the AP problem, for any block, as a concatenation of \mathcal{N} gadgets corresponding to $i = 1, \dots, \mathcal{N}$, and we construct the bit vector $U^{(K,J)}$ of the AP problem as a concatenation of \mathcal{N} bit vectors, one per gadget. Each gadget consists of the same string $\mathbf{a}^L \mathbf{b} \mathbf{a}^L$; we set to 1 the k' -th bit of the i -th gadget bit vector if $A[i, (K-1)L + k'] = 1$. The solution of the AP problem $V^{(K,J)}$ will allow us to recover the solution of BMM, as we will ensure that the bit corresponding to the j' -th \mathbf{a} in the second half of the gadget is set to 1 if and only if, for some $k' \in [L]$, $A[i, (K-1)L + k'] = 1$ and $B[(K-1)L + k', (J-1)L + j'] = 1$. In order to enforce this, we will include the following strings in set $\mathcal{S}^{(K,J)}$:

$$\mathbf{a}^{L-k'} \mathbf{b} \mathbf{a}^{j'}, \text{ for every } k', j' \in [L] \text{ such that } B[(K-1)L + k', (J-1)L + j'] = 1.$$

This guarantees that after solving the AP problem we have the required property, and thus, after solving all the instances, we have obtained matrix $C = A \times B$. Indeed, consider values j , i.e., the index that runs on the columns of C , in intervals of size L . By construction and by definition of BMM, the i -th line of the J -th column interval of C is obtained by taking the disjunction of the second half of the i -th interval of each (K, J) -th bit vector for every $K = 1, 2, \dots, \mathcal{N}/L$.

We have a total of $(\mathcal{N}/L)^2$ instances. In each of them, the total length of all strings is $\mathcal{O}(L^3)$, and the length of the input string P is $(2L+1)\mathcal{N} = \mathcal{O}(L \cdot \mathcal{N})$. Using our assumed algorithm for each instance, we obtain the following total time:

$$\mathcal{O}((\mathcal{N}/L)^2 \cdot (L^3 + (\mathcal{N} \cdot L)^{1.5-\epsilon})) = \mathcal{O}(\mathcal{N}^2 \cdot L + \mathcal{N}^{3.5-\epsilon}/L^{0.5+\epsilon}).$$

If we set $L = \mathcal{N}^{(1.5-\epsilon)/(1.5+\epsilon)}$, then the total time becomes:

$$\begin{aligned} & \mathcal{O}(\mathcal{N}^{2+(1.5-\epsilon)/(1.5+\epsilon)} + \mathcal{N}^{3.5-\epsilon-(0.5+\epsilon)(1.5-\epsilon)/(1.5+\epsilon)}) \\ &= \mathcal{O}(\mathcal{N}^{2+(1.5-\epsilon)/(1.5+\epsilon)} + \mathcal{N}^{2+(1.5-\epsilon)-(1.5-\epsilon)(0.5+\epsilon)/(1.5+\epsilon)}) \\ &= \mathcal{O}(\mathcal{N}^{2+(1.5-\epsilon)/(1.5+\epsilon)} + \mathcal{N}^{2+(1.5-\epsilon)(1.5+\epsilon-0.5-\epsilon)/(1.5+\epsilon)}) \\ &= \mathcal{O}(\mathcal{N}^{2+(1.5-\epsilon)/(1.5+\epsilon)}). \end{aligned}$$

Hence we obtain a combinatorial BMM algorithm with complexity $\mathcal{O}(\mathcal{N}^{3-\epsilon'})$, where $\epsilon' = 1 - (1.5 - \epsilon)/(1.5 + \epsilon) > 0$. \square

Example 3. Consider the following instance of the BMM problem with $\mathcal{N} = 6$ and $L = 3$.

$$\begin{array}{c} A \end{array} \quad \begin{array}{c} B \end{array} \quad \begin{array}{c} C \end{array}$$

$$\left[\begin{array}{ccc|ccc} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \hline 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{array} \right] \times \left[\begin{array}{ccc|ccc} 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \end{array} \right] = \left[\begin{array}{ccc|ccc} \mathbf{1} & \mathbf{0} & \mathbf{0} & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{array} \right]$$

From matrices A and B , we now show how the resulting matrix C can be found by building and solving 4 instances of the AP problem constructed as follows. The pattern is

$$P = \mathbf{aaabaaa} \cdot \mathbf{aaabaaa} \cdot \mathbf{aaabaaa} \cdot \mathbf{aaabaaa} \cdot \mathbf{aaabaaa} \cdot \mathbf{aaabaaa}$$

where the six gadgets are separated by a $'$ to be highlighted. For the AP instances, the vectors $U^{(K,J)}$ shown below are the input bit vectors, and the sets $\mathcal{S}^{(K,J)}$ are the input set of strings.

For each instance, the bit vector $V^{(K,J)}$ shown below is the output of the AP problem.

| i | 1 | 2 | 3 | 4 | 5 | 6 |
|---------------|-------------------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| $U^{(1,1)} :$ | [0 1 0 0 0 0 0] | [1 0 1 0 0 0 0] | [0 0 0 0 0 0 0] | [1 0 0 0 0 0 0] | [0 0 0 0 0 0 0] | [0 1 0 0 0 0 0] |
| $S^{(1,1)} :$ | {aba,baaa} | | | | | |
| $V^{(1,1)} :$ | [0 0 0 0 1 0 0] | [0 0 0 0 0 0 1] | [0 0 0 0 0 0 0] | [0 0 0 0 0 0 0] | [0 0 0 0 0 0 0] | [0 0 0 0 1 0 0] |
| $U^{(1,2)} :$ | [0 1 0 0 0 0 0] | [1 0 1 0 0 0 0] | [0 0 0 0 0 0 0] | [1 0 0 0 0 0 0] | [0 0 0 0 0 0 0] | [0 1 0 0 0 0 0] |
| $S^{(1,2)} :$ | {aabaaa,baa} | | | | | |
| $V^{(1,2)} :$ | [0 0 0 0 0 0 0] | [0 0 0 0 0 1 1] | [0 0 0 0 0 0 0] | [0 0 0 0 0 0 1] | [0 0 0 0 0 0 0] | [0 0 0 0 0 0 0] |
| $U^{(2,1)} :$ | [0 1 0 0 0 0 0] | [0 0 0 0 0 0 0] | [0 0 1 0 0 0 0] | [0 1 0 0 0 0 0] | [1 0 0 0 0 0 0] | [0 0 0 0 0 0 0] |
| $S^{(2,1)} :$ | {aabaa,ba} | | | | | |
| $V^{(2,1)} :$ | [0 0 0 0 0 0 0] | [0 0 0 0 0 0 0] | [0 0 0 0 1 0 0] | [0 0 0 0 0 0 0] | [0 0 0 0 0 1 0] | [0 0 0 0 0 0 0] |
| $U^{(2,2)} :$ | [0 1 0 0 0 0 0] | [0 0 0 0 0 0 0] | [0 0 1 0 0 0 0] | [0 1 0 0 0 0 0] | [1 0 0 0 0 0 0] | [0 0 0 0 0 0 0] |
| $S^{(2,2)} :$ | {aba,baa} | | | | | |
| $V^{(2,2)} :$ | [0 0 0 0 1 0 0] | [0 0 0 0 0 0 0] | [0 0 0 0 0 1 0] | [0 0 0 0 1 0 0] | [0 0 0 0 0 0 0] | [0 0 0 0 0 0 0] |

As an example on how to obtain matrix C , consider the bold part of C above (*i.e.*, the first line of block (1, 1) of C). This is obtained by taking the disjunction of the bold parts of $V^{(1,1)}$ and $V^{(2,1)}$.

4 EDSM Conditional Lower Bound

Since the lower bound for the AP problem does not imply *per se* a lower bound for the whole EDSM problem, in this section we show a conditional lower bound for the EDSM problem. Specifically, we perform a reduction from Triangle Detection to show that, if the EDSM problem could be solved in $\mathcal{O}(nm^{1.5-\epsilon} + N)$ time, this would imply the existence of a truly subcubic algorithm for TD. We show that TD can be reduced to the decision version of the EDSM problem: the goal is to detect whether there exists at least one occurrence of P in \tilde{T} . To this aim, given three matrices A, B, C , we first decompose matrix B into blocks of size $\mathcal{N}/s \times \mathcal{N}/s$, where s is a parameter to be determined later; the pattern P is obtained by concatenating a number (namely $z = \mathcal{N}s^2$) of constituent parts P_i of length $\mathcal{O}(\mathcal{N}/s)$, each one built with five letters from disjoint subalphabets. The ED text \tilde{T} is composed of three parts: the central part consists of three degenerate segments, the first one encoding the 1s of matrix A , the second one those of matrix B and the third one those of matrix C . These segments are built in such a way that the concatenation of strings of subsequent segments is of the same form as the pattern's building blocks. This central part is then padded to the left and to the right with sets containing appropriately chosen concatenations of substrings P_i of P , so that an occurrence of the pattern in the text implies that one of its building blocks matches the central part of the text, thus corresponding to a triangle. Formally:

Theorem 1. *If the EDSM problem can be solved in $\mathcal{O}(nm^{1.5-\epsilon} + N)$ time, for any $\epsilon > 0$, with a combinatorial algorithm, then there exists a truly subcubic combinatorial algorithm for TD.*

Proof. Consider an instance of TD, where we are given three $\mathcal{N} \times \mathcal{N}$ Boolean matrices A, B, C , and the question is to check if there exist i, j, k such that $A[i, j] = B[j, k] = C[k, i] = 1$. Let s

be a parameter, to be determined later, that corresponds to decomposing B into blocks of size $(\mathcal{N}/s) \times (\mathcal{N}/s)$. We reduce to an instance of EDSM over an alphabet Σ of size $\mathcal{O}(\mathcal{N})$.

Pattern P . We construct P by concatenating, in some fixed order, the following strings:

$$P(i, x, y) = v(i)xa^{\mathcal{N}/s}x\$y a^{\mathcal{N}/s}yv(i)$$

for every $i = 1, 2, \dots, \mathcal{N}$ and $x, y = 1, 2, \dots, s$, where $a \in \Sigma_1$, $\$ \in \Sigma_2$, $x \in \Sigma_3$, $y \in \Sigma_4$, $v(i) \in \Sigma_5$, and $\Sigma_1, \Sigma_2, \dots, \Sigma_5$ are disjoint subsets of Σ .

ED text \tilde{T} . The text \tilde{T} consists of three parts. Its middle part encodes all the entries equal to 1 in matrices A , B and C , and consists of three string sets $\mathcal{X} = \mathcal{X}_1 \cdot \mathcal{X}_2 \cdot \mathcal{X}_3$, where:

1. \mathcal{X}_1 contains all strings of the form $v(i)xa^j$, for some $i \in [\mathcal{N}]$, $x \in [s]$ and $j \in [\mathcal{N}/s]$ such that $A[i, (x-1) \cdot (\mathcal{N}/s) + j] = 1$;
2. \mathcal{X}_2 contains all strings of the form $a^{\mathcal{N}/s-j} x\$y a^{\mathcal{N}/s-k}$, for some $x, y \in [s]$ and $j, k \in [\mathcal{N}/s]$ such that $B[(x-1) \cdot (\mathcal{N}/s) + j, (y-1) \cdot (\mathcal{N}/s) + k] = 1$, i.e., if the corresponding entry of B is 1;
3. \mathcal{X}_3 contains all strings of the form $a^k yv(i)$, for some $i \in [\mathcal{N}]$, $y \in [s]$ and $k \in [\mathcal{N}/s]$ such that $C[(y-1) \cdot (\mathcal{N}/s) + k, i] = 1$.

It is easy to see that $|P(i, x, y)| = \mathcal{O}(\mathcal{N}/s)$. This implies the following:

1. The length of the pattern is $m = \mathcal{O}(\mathcal{N} \cdot s^2 \cdot \mathcal{N}/s) = \mathcal{O}(\mathcal{N}^2 \cdot s)$;
2. The total length of \mathcal{X} is $|\mathcal{X}| = \mathcal{O}(\mathcal{N} \cdot s \cdot \mathcal{N}/s \cdot \mathcal{N}/s + s^2 \cdot (\mathcal{N}/s)^2 \cdot \mathcal{N}/s + \mathcal{N} \cdot s \cdot \mathcal{N}/s \cdot \mathcal{N}/s) = \mathcal{O}(\mathcal{N}^3/s)$.

By the above construction, we obtain the following fact.

Fact 1. $P(i, x, y)$ matches \mathcal{X} if and only if the following holds for some $j, k = 1, 2, \dots, \mathcal{N}/s$:

$$A[i, (x-1) \cdot (\mathcal{N}/s) + j] = B[(x-1) \cdot (\mathcal{N}/s) + j, (y-1) \cdot (\mathcal{N}/s) + k] = C[(y-1) \cdot (\mathcal{N}/s) + k, i] = 1$$

Solving the TD problem thus reduces to taking the disjunction of all such conditions. Let us write down all strings $P(i, x, y)$ in some arbitrary but fixed order to obtain $P = P_1 P_2 \dots P_z$ with $z = \mathcal{N}s^2$, where every $P_t = P(i, x, y)$, for some i, x, y . We aim to construct a small number of sets of strings that, when considered as an ED text, match any prefix $P_1 P_2 \dots P_t$ of the pattern, $1 \leq t \leq z-1$; a similar construction can be carried on to obtain sets of strings that match any suffix $P_k \dots P_{z-1} P_z$, $2 \leq k \leq z$. These sets will then be added to the left and to the right of \mathcal{X} , respectively, to obtain the ED text \tilde{T} .

ED Prefix. We construct $\log z$ sets of strings as follows. The first one contains the empty string ε and $P_1 P_2 \dots P_{z/2}$. The second one contains ε , $P_1 P_2 \dots P_{z/4}$ and $P_{z/2+1} \dots P_{z/2+z/4}$. The third one contains ε , $P_1 P_2 \dots P_{z/8}$, $P_{z/4+1} \dots P_{z/4+z/8}$, $P_{z/2+1} \dots P_{z/2+z/8}$ and $P_{z/2+z/4+1} \dots P_{z/2+z/4+z/8}$. Formally, for every $i = 1, 2, \dots, \log z$, the i -th of such sets is:

$$\tilde{T}_i^p = \varepsilon \cup \{P_{j \frac{z}{2^{i-1}} + 1} \dots P_{j \frac{z}{2^{i-1}} + \frac{z}{2^i}} \mid j = 0, 1, \dots, 2^{i-1} - 1\}.$$

ED Suffix. We similarly construct $\log z$ sets to be appended to \mathcal{X} :

$$\tilde{T}_i^s = \varepsilon \cup \{P_{z-j \frac{z}{2^{i-1}} - \frac{z}{2^i} + 1} \dots P_{z-j \frac{z}{2^{i-1}}} \mid j = 0, 1, \dots, 2^{i-1} - 1\}.$$

The total length of all the ED prefix and ED suffix strings is $\mathcal{O}(\log z \cdot \mathcal{N}^2 \cdot s) = \mathcal{O}(\mathcal{N}^2 \cdot s \cdot \log \mathcal{N})$. The whole ED text \tilde{T} is thus: $\tilde{T} = \tilde{T}_1^p \dots \tilde{T}_{\log z}^p \cdot \mathcal{X} \cdot \tilde{T}_{\log z}^s \dots \tilde{T}_1^s$. We next show how a solution of such instance of EDSM corresponds to the solution of TD.

Lemma 2. *The pattern P occurs in the ED text \tilde{T} if and only if there exist i, j, k such that $A[i, j] = B[j, k] = C[k, i] = 1$.*

Proof. By Fact 1, if such i, j, k exist then P_t matches \mathcal{X} , for some $t \in \{1, \dots, z\}$. Then, by construction of the sets \tilde{T}_i^p and \tilde{T}_i^s , the prefix $P_1 \dots P_{t-1}$ matches the ED prefix (this can be proved by induction), and similarly the suffix $P_{t+1} \dots P_z$ matches the ED suffix, so the whole P matches \tilde{T} , and so P occurs therein. Because of the letters $\$$ appearing only in the center of P_i s and strings from \mathcal{X}_2 , every P_i s and a concatenation of $X_1 \in \mathcal{X}_1$, $X_2 \in \mathcal{X}_2$, $X_3 \in \mathcal{X}_3$ having the same length, and the P_i s being distinct, there is an occurrence of the pattern P in \tilde{T} if and only if $X_1 X_2 X_3 = P_t$ for some t and $X_1 \in \mathcal{X}_1$, $X_2 \in \mathcal{X}_2$, $X_3 \in \mathcal{X}_3$. But then, by Fact 1 there exists a triangle. \square

Note that for the EDSM problem we have $m = \mathcal{N}^2 \cdot s$, $n = 1 + 2\log z$ and $N = \|\mathcal{X}\| + \mathcal{O}(\mathcal{N}^2 \cdot s \cdot \log \mathcal{N})$. Thus if we had a solution running in $\mathcal{O}(\log z \cdot m^{1.5-\epsilon} + \|\mathcal{X}\| + \mathcal{N}^2 \cdot s \cdot \log \mathcal{N}) = \mathcal{O}(\log \mathcal{N} \cdot (\mathcal{N}^2 \cdot s)^{1.5-\epsilon} + \mathcal{N}^3/s)$ time, for some $\epsilon > 0$, by choosing a sufficiently small $\alpha > 0$ and setting $s = \mathcal{N}^\alpha$ we would obtain, for some $\delta > 0$, an $\mathcal{O}(\mathcal{N}^{3-\delta})$ -time algorithm for TD. \square

5 An $\tilde{\mathcal{O}}(nm^{\omega-1} + N)$ -time Algorithm for EDSM

Our goal is to design a non-combinatorial $\tilde{\mathcal{O}}(nm^{\omega-1} + N)$ -time algorithm for EDSM, which in turn can be achieved with a non-combinatorial $\tilde{\mathcal{O}}(m^{\omega-1} + N)$ -time algorithm for the AP problem, that is the bottleneck of EDSM (cf. [35]).

We reduce AP to a logarithmic number of restricted instances of the same problem, based on the length of the strings in \mathcal{S} . We start by giving a lemma that we will use to process naively the strings of length up to a constant c , to be determined later, in $\mathcal{O}(m + N)$ time.

Lemma 3. *For any integer t , all strings in \mathcal{S} of length at most t can be processed in $\mathcal{O}(m \log m + mt + N)$ time.*

Proof. We first construct the suffix tree ST of P and store, for every node, the first letters on its outgoing edges in a static dictionary with constant access time. This can be done in $\mathcal{O}(m \log m)$ time [55]. For every $S \in \mathcal{S}$, find and mark its corresponding (implicit or explicit) node of ST . This takes $\mathcal{O}(N)$ time overall. For every possible length $t' \leq t$, scan P with a window of length t' while maintaining its corresponding node of ST . This takes $\mathcal{O}(m)$ time overall. If the current window $P[i \dots (i + t' - 1)]$ corresponds to a marked node of ST and additionally $U[i - 1] = 1$, we set $V[i + t' - 1] = 1$. \square

We build the rest of the restricted instances of the AP problem by restricting on strings in $\mathcal{S}_k \subseteq \mathcal{S}$ of length in $[(19/18)^k, (19/18)^{k+1})$ for each integer k ranging from $\left\lceil \frac{\log c}{\log(19/18)} \right\rceil$ to $\left\lfloor \frac{\log m}{\log(19/18)} \right\rfloor$. These intervals are a partition of the set of all strings in \mathcal{S} of length up to m ; longer strings are not addressed in EDSM by solving AP.

For each integer k from $\left\lceil \frac{\log c}{\log(19/18)} \right\rceil$ to $\left\lfloor \frac{\log m}{\log(19/18)} \right\rfloor$, let ℓ be an integer such that the length of every string in \mathcal{S}_k belongs to $[9/8 \cdot \ell, 5/4 \cdot \ell)$. Note that such an integer always exists for an appropriate choice of the integer constant c . In fact, it must hold that

$$\frac{9}{8} \cdot \ell \leq \left(\frac{19}{18}\right)^k < \left(\frac{19}{18}\right)^{k+1} \leq \frac{5}{4} \cdot \ell \iff \frac{4}{5} \cdot \left(\frac{19}{18}\right)^{k+1} \leq \ell \leq \frac{8}{9} \cdot \left(\frac{19}{18}\right)^k.$$

To ensure that there exists an integer ℓ satisfying such conditions, it must actually hold that

$$\frac{4}{5} \cdot \left(\frac{19}{18}\right)^{k+1} + 1 \leq \frac{8}{9} \cdot \left(\frac{19}{18}\right)^k \iff \frac{45}{2} \leq \left(\frac{19}{18}\right)^k.$$

The last equation holds for $k \geq 58$, implying that we will process naïvely the strings of length up to $c = 23$, and each \mathcal{S}_k , for k ranging from 58 to $\left\lfloor \frac{\log m}{\log(19/18)} \right\rfloor$, will be processed separately as described in the next paragraph.

Denoting by N_k the total size of strings in \mathcal{S}_k , we have that, if we solve every such instance of AP in $\mathcal{O}(N_k + f(m))$ time, then we can solve the original instance of AP in $\mathcal{O}(N + f(m) \log m)$ time by taking the results disjunction. Switching to $\tilde{\mathcal{O}}$ notation that disregards polylog factors, it thus suffices to solve each such instance of the AP problem in $\tilde{\mathcal{O}}(N + m^{\omega-1})$ time.

We further partition the strings in \mathcal{S}_k into three types, compute the corresponding bit vector V for each type separately and, finally, take the disjunction of the resulting bit vectors V to obtain the answer for each restricted instance.

Partitioning \mathcal{S}_k . Keeping in mind that from now on (until Section 5.4) we address the AP problem assuming that \mathcal{S} only contains strings of length in $[9/8 \cdot \ell, 5/4 \cdot \ell)$, and thus is in fact \mathcal{S}_k , to lighten the notation we now switch back to denote it simply with \mathcal{S} . The three types of strings are as follows:

Type 1: Strings $S \in \mathcal{S}$ such that every length- ℓ substring of S is not strongly periodic.

Type 2: Strings $S \in \mathcal{S}$ containing at least one length- ℓ substring that is not strongly periodic and at least one length- ℓ substring that is strongly periodic.

Type 3: Strings $S \in \mathcal{S}$ such that every length- ℓ substring of S is strongly periodic (in Lemma 4 we show that in this case $\text{per}(S) \leq \ell/4$).

These three types are evidently a partition of \mathcal{S} . We start with showing that, in fact, strings of type 3 are exactly strings with period at most $\ell/4$.

Lemma 4. *Let S be a string. If $\text{per}(S[j \dots j + \ell - 1]) \leq \ell/4$ for every j then $\text{per}(S) \leq \ell/4$.*

Proof. We first show that, for any string W and letters a, b , if $\text{per}(aW) \leq |aW|/4$ and $\text{per}(Wb) \leq |Wb|/4$ then $\text{per}(aW) = \text{per}(Wb)$. This follows from Lemma 1: since $\text{per}(aW)$ and $\text{per}(Wb)$ are both periods of W and $(1 + |W|)/4 \leq |W|/2$, then we have that $d = \gcd(\text{per}(aW), \text{per}(Wb))$ is a period of W . Assuming by contradiction that $\text{per}(aW) \neq \text{per}(Wb)$, then it must be that either $d < \text{per}(aW)$ or $d < \text{per}(Wb)$; by symmetry it is enough to consider the former possibility, and we claim that then d is a period of aW . Indeed, $a = W[\text{per}(aW) - 1]$ (observe that $\text{per}(aW) - 1 \leq |W|$) and $W[i] = W[i + d]$ for any $i = 1, 2, \dots, |W| - d$, so by $\text{per}(aW)$ being a multiple of d we obtain that $a = W[\text{per}(aW) - 1] = W[d - 1]$, which is a contradiction because by definition of $\text{per}(aW)$ we have that $d < \text{per}(aW)$ cannot be a period of aW .

If $\text{per}(S[j \dots j + \ell - 1]) \leq \ell/4$ for every j then by the above reasoning the periods of all substrings $S[j \dots j + \ell - 1]$ is the same and in fact equal to p . But then $S[i] = S[i + p]$ for every i , so $\text{per}(S) \leq \ell/4$. \square

Before proceeding with the algorithm, we show that, for each string $S \in \mathcal{S}$, we can determine its type in $\mathcal{O}(|S|)$ time.

Lemma 5. *Given a string S we can determine its type in $\mathcal{O}(|S|)$ time.*

Proof. It is well-known that $\text{per}(T)$ can be computed in $\mathcal{O}(|T|)$ time for any string T (cf. [24]). We partition S into blocks $T_\alpha = S[\alpha \lfloor \ell/2 \rfloor \dots (\alpha + 1) \lfloor \ell/2 \rfloor - 1]$ of size $\lfloor \ell/2 \rfloor$, and compute $\text{per}(T_\alpha)$ for every α in $\mathcal{O}(|S|)$ total time. Observe that every substring $S[i \dots i + \ell - 1]$ contains at least one whole block inside.

If $\text{per}(T_\alpha) > \ell/4$ then the period of any substring $S[i \dots i + \ell - 1]$ that contains T_α is also larger than $\ell/4$. Consequently, if $\text{per}(T_\alpha) > \ell/4$ for every α , then we declare S to be of type 1.

Consider any α such that $p = \text{per}(T_\alpha) \leq \ell/4$. If the period p' of a substring $S' = S[i \dots i + \ell - 1]$ that contains T_α is at most $\ell/4$, then in fact it must be equal to p , because $p' \geq p$ and so, by

Lemma 1 applied on T_α , p' must be a multiple of p and, by repeatedly applying $S'[j] = S'[j + p']$ and $T_\alpha[j] = T_\alpha[j + p]$ and using the fact that T_α occurs inside S' , we conclude that in fact $S'[j] = S'[j + p]$ for any j , and thus $p' = p$. This allows us to check whether there exists a substring $S' = S[i \dots i + \ell - 1]$ that contains T_α such that $\text{per}(S') \leq \ell/4$ by computing, in $\mathcal{O}(\ell)$ time, how far the period p extends to the left and to the right of T_α in $T_{\alpha-1}T_\alpha T_{\alpha+1}$ (if either $T_{\alpha-1}$ or $T_{\alpha+1}$ do not exist, then we do not extend the period in the corresponding direction). There exists such a substring S' if and only if the length of the extended substring with period p is at least ℓ . Therefore, for every α we can check in $\mathcal{O}(\ell)$ time if there exists a length- ℓ substring S' containing T_α with $\text{per}(S') \leq \ell/4$. By repeating this procedure for every α , we can distinguish between S of type 2 and S of type 3 in $\mathcal{O}(|S|)$ total time. \square

Since we have shown how to efficiently partition the strings of \mathcal{S} into the three types, in what follows we present our solution of the AP problem for each type of strings separately.

Remark 1. The length of every string in \mathcal{S} belonging to $[9/8 \cdot \ell, 5/4 \cdot \ell)$ implies that every string in \mathcal{S} contains at most $\ell/4$ length- ℓ substrings (and at least $1 + \ell/8$ of them).

5.1 Type 1 Strings

In this section we show how to solve a restricted instance of the AP problem where every string $S \in \mathcal{S}$ is of type 1, that is, each of its length- ℓ substrings is not strongly periodic, and furthermore $|S| \in [9/8 \cdot \ell, 5/4 \cdot \ell)$ for some $\ell \leq m$. Observe that all (and hence at most $\ell/4$ by Remark 1) length- ℓ substrings of any $S \in \mathcal{S}$ must be distinct, as otherwise we would be able to find two occurrences of a length- ℓ substring at distance at most $\ell/4$ in S , making the period of the substring at most $\ell/4$ and contradicting the assumption that S is of type 1.

We start with constructing the suffix tree ST of P (our pattern in the EDSM problem) in $\mathcal{O}(m \log m)$ time [59]. Let us remark that we are spending $\mathcal{O}(m \log m)$ time and not just $\mathcal{O}(m)$ so as to avoid any assumptions on the size of the alphabet. For every explicit node $u \in ST$, we construct a perfect hash function mapping the first letter on every edge outgoing from u to the corresponding edge. This takes $\mathcal{O}(m \log m)$ time [55] and allows us to navigate in ST in constant time per letter. Then, for every $S \in \mathcal{S}$, we check in $\mathcal{O}(|S|)$ time using ST if it occurs in P and, if not, we disregard it from further consideration. Therefore, from now on we assume that all strings S , and thus all their length- ℓ substrings, occur in P . We will select a set of length- ℓ substrings of P , called the *anchors*, each represented by one of its occurrences in P , such that:

1. The total number of occurrences of all anchors in P is $\mathcal{O}(m/\ell \cdot \log m)$.
2. For every $S \in \mathcal{S}$, at least one of its length- ℓ substrings is an anchor.
3. The total number of occurrences of all anchors in strings $S \in \mathcal{S}$ is $\mathcal{O}(|\mathcal{S}| \cdot \log m)$.

We formalize this using the following auxiliary problem, which is a strengthening of a well-known *Hitting Set* problem, which given a collection of m sets over $[n]$, each of size at least k , asks to choose a subset of $[n]$ of size $\mathcal{O}(n/k \cdot \log m)$ that nontrivially intersects every set.

NODE SELECTION (NS)

INPUT: A bipartite graph $G = (U, V, E)$ with $\deg(u) \in (d, 2d]$ for every $u \in U$ and weight $w(v)$ for every $v \in V$, where $W = \sum_{v \in V} w(v)$.

OUTPUT: A set $V' \subseteq V$ of total weight $\mathcal{O}(W/d \cdot \log |U|)$ such that $N[u] \cap V' \neq \emptyset$ for every node $u \in U$, and $\sum_{u \in U} |N[u] \cap V'| = \mathcal{O}(|U| \log |U|)$.

We reduce the problem of finding anchors to an instance of the NS problem, by building a bipartite graph G in which the nodes in U correspond to strings $S \in \mathcal{S}$, the nodes in V correspond to distinct length- ℓ substrings of P , and there is an edge (u, v) if the length- ℓ string corresponding to v occurs in the string S corresponding to u . Using suffix links, we can find the

node of the suffix tree corresponding to every length- ℓ substring of S in $\mathcal{O}(|S|)$ total time, so the whole construction takes $\mathcal{O}(m \log m + \sum_{S \in \mathcal{S}} |S|) = \mathcal{O}(m \log m + N)$ time. The size of G is $\mathcal{O}(m + N)$, and the degree of every node in U belongs to $(\ell/8, \ell/4]$. We set the weight of a node $v \in V$ to be its number of occurrences in P , and solve the obtained instance of the NS problem to obtain the set of anchors. It is not immediately clear that an instance of the NS problem always has a solution. We show that indeed it does, and that it can be found in linear time.

Lemma 6. *A solution to an instance of the NS problem always exists and can be found in linear time in the size of G .*

Proof. We first show a solution that uses the probabilistic method and leads us to an efficient Las Vegas algorithm; we will then derandomize the solution using the method of conditional expectations.

We independently choose each node of V with probability p to obtain the set V' of selected nodes. The expected total weight of V' is $\sum_{v \in V} p \cdot w(v) = p \cdot W$, so by Markov's inequality it exceeds $4p \cdot W$ with probability at most $1/4$. For every node $u \in U$, the probability that $N[u]$ does not intersect V' is at most $(1 - p)^d \leq e^{-pd}$. Finally, $\mathbb{E}[\sum_{u \in U} |N[u] \cap V'|] \leq |U| \cdot 2pd$, so by Markov's inequality $\sum_{u \in U} |N[u] \cap V'|$ exceeds $|U| \cdot 8pd$ with probability at most $1/4$. We set $p = \ln(4|U|)/d$ (observe that if $p > 1$ then we can select all nodes in V). By union bound, the probability that V' is not a valid solution is at most $3/4$, so indeed a valid solution exists. Furthermore, this reasoning gives us an efficient Las Vegas algorithm that chooses V' randomly as described above and then verifies if it constitutes a valid solution. Each iteration takes linear time in the size of G , and the expected number of required iterations is constant.

To derandomize the above procedure we apply the method of conditional expectations. Let X_1, X_2, \dots be the binary random variables corresponding to the nodes of V . Recall that in the above proof we set $X_i = 1$ with probability p . Now we will choose the values of X_1, X_2, \dots one-by-one. Define a function $f(X_1, X_2, \dots)$ that bounds the probability that X_1, X_2, \dots corresponds to a valid solution as follows:

$$f(X_1, X_2, \dots) = \frac{\sum_v X_v \cdot w(v)}{4W/d \cdot \ln(4|U|)} + \sum_{u \in U} \prod_{v \in N[u]} (1 - X_v) + \frac{\sum_{u \in U} \sum_{v \in N[u]} X_v}{8|U| \ln(4|U|)}.$$

As explained above, we have $\mathbb{E}[f(X_1, X_2, \dots)] = 3/4$. Assume that we have already fixed the values $X_1 = x_1, \dots, X_i = x_i$. Then there must be a choice of $X_{i+1} = x_{i+1}$ that does not increase the expected value of $f(X_1, X_2, \dots)$ conditioned on the already chosen values. We want to compare the following two quantities:

$$\begin{aligned} &\mathbb{E}[f(X_1, X_2, \dots) \mid X_1 = x_1, \dots, X_i = x_i, X_{i+1} = 0] \\ &\mathbb{E}[f(X_1, X_2, \dots) \mid X_1 = x_1, \dots, X_i = x_i, X_{i+1} = 1] \end{aligned}$$

and choose x_{i+1} corresponding to the smaller one. Cancelling out the shared terms, we need to compare the expected values of:

$$\begin{aligned} &0 + \sum_{u \in N[i+1]} \prod_{v \in N[u]} (1 - X_v) + 0 \quad \text{and} \\ &\frac{w(i+1)}{4W/d \cdot \ln(4|U|)} + 0 + \frac{\deg(i+1)}{8|U| \ln(4|U|)}. \end{aligned}$$

The second quantity can be computed in constant time. We claim that (ignoring the issue of numerical precision) the first quantity can be computed in time $\mathcal{O}(\deg(i+1))$ after a linear-time preprocessing as follows. In the preprocessing we compute and store $E[i] = \mathbb{E}[\prod_{j=1}^i (1 - Y_j)]$, where the Y_j 's are independent indicator variables with $\Pr[Y_j = 1] = p$, for every $i = 0, 1, \dots, |V|$. It is straightforward to compute $E[i+1]$ from $E[i]$ in constant time. Then, during the computation

we maintain, for every $u \in U$, the number $c[u]$ of $v \in N[u]$ for which we still need to choose the value X_v , and a single bit $b[u]$ denoting whether for some $v \in N[u] \cap \{1, \dots, i\}$ we already have $x_v = 1$. This information can be updated in $\mathcal{O}(\deg(i+1))$ time after selecting x_{i+1} . Now to compute the first quantity, we iterate over $u \in N[i+1]$ and, if $b[u] = 0$ then we add $E[c[u]]$ to the result. Finally, we claim that it is enough to implement all calculations with precision $\Theta(\log |V|)$ bits. This is because such precision allows us to calculate both quantities with relative accuracy $1/(8|V|)$, so the expected value of $f(X_1, X_2, \dots)$ might increase by a factor of $(1 + 1/(4|V|))$ in every step, which is at most $(1 + 1/(4|V|))^{|V|} \leq e^{1/4}$ overall. This still guarantees that the final value is at most $3/4 \cdot e^{1/4} < 1$, so we obtain a valid solution. \square

In the rest of this section we explain how to compute the bit vector V from the bit vector U , and thus solve the AP problem, after having obtained a set \mathcal{A} of anchors. For any $S \in \mathcal{S}$, since S contains an occurrence of at least one anchor $H \in \mathcal{A}$, say $S[j..(j+|H|-1)] = H$, so any occurrence of S in P can be generated by choosing some occurrence of H in P , say $P[i..(i+|H|-1)] = H$, and then checking that $S[1..(j-1)] = P[(i-j+1)..(i-1)]$ and $S[(j+|H|)..|S|] = P[(i+|H|)..(i+|S|-j)]$. In other words, $S[1..(j-1)]$ should be a suffix of $P[1..(i-1)]$ and $S[(j+|H|)..|S|]$ should be a prefix of $P[(i+|H|)..|P|]$. In such case, we say that the occurrence of S in P is generated by H . By the properties of \mathcal{A} , any occurrence of $S \in \mathcal{S}$ is generated by $\text{occ}_S \geq 1$ occurrences of anchors, where $\sum_{S \in \mathcal{S}} \text{occ}_S = \mathcal{O}(|\mathcal{S}| \log m)$. For every $H \in \mathcal{A}$ we create a separate data structure $D(H)$ responsible for setting $V[i+|S|-1] = 1$, when $U[i-1] = 1$ and $P[i..(i+|S|-1)] = S$ is generated by H . We now first describe what information is used to initialize each $D(H)$, and how this is later processed to update V .

Initialization. $D(H)$ consists of two compact tries $T(H)$ and $T^r(H)$. For every occurrence of H in P , denoted by $P[i..(i+|H|-1)] = H$, $T(H)$ should contain a leaf corresponding to $P[(i+|H|)..|P|]$ and $T^r(H)$ should contain a leaf corresponding to $(P[1..(i-1)])^r$, both decorated with position i . Additionally, $D(H)$ stores a list $L(H)$ of pairs of nodes (u, v) , where $u \in T^r(H)$ and $v \in T(H)$. Each such pair corresponds to an occurrence of H in a string $S \in \mathcal{S}_h$, $S[j..(j+|H|-1)] = H$, where u is the node of $T^r(H)$ corresponding to $(S[1..(j-1)])^r$ and v is the node of $T(H)$ corresponding to $S[(j+|H|+1)..|S|]$. We claim that $D(H)$, for all H , can be constructed in $\mathcal{O}(m \log m + N)$ total time.

We first construct the suffix tree ST of $P\$$ and the suffix tree ST^r of $P^r\$$ (again in $\mathcal{O}(m \log m)$ time not to make assumptions on the alphabet). We augment both trees with data for answering both *weighted ancestor* (WA) and *lowest common ancestor* (LCA) queries, that are defined as follows. For a rooted tree T on n nodes with an integer weight $\mathcal{D}(v)$ assigned to every node u , such that the weight of the root is zero and $\mathcal{D}(u) < \mathcal{D}(v)$ if u is the parent of v , we say that a node v is a weighted ancestor of a node u at depth ℓ , denoted by $\text{WA}_T(u, \ell)$, if v is the highest ancestor of u with weight at least ℓ . Such queries can be answered in $\mathcal{O}(\log n)$ time after an $\mathcal{O}(n)$ preprocessing [28]. For a rooted tree T , $\text{LCA}_T(u, v)$ is the lowest node that is an ancestor of both u and v . Such queries can be answered in $\mathcal{O}(1)$ time after an $\mathcal{O}(n)$ preprocessing [10]. Recall that every anchor H is represented by one of its occurrences in P . Using WA queries, we can access in $\mathcal{O}(\log m)$ time the nodes corresponding to H and H^r , respectively, and extract a lexicographically sorted list of suffixes following an occurrence of H in $P\$$ and a lexicographically sorted list of reversed prefixes preceding an occurrence of H in $P^r\$$ in time proportional to the number of such occurrences. Then, by iterating over the lexicographically sorted list of suffixes and using LCA queries on ST we can build $T(H)$ in time proportional to the length of the list, and similarly we can build $T^r(H)$. To construct $L(H)$ we start by computing, for every $S \in \mathcal{S}$ and $j = 1, \dots, |S|$, the node of ST^r corresponding to $(S[1..j])^r$ and the node of ST corresponding to $S[(j+1)..|S|]$ (the nodes might possibly be implicit). This takes only $\mathcal{O}(|\mathcal{S}|)$ time, by using suffix links. We also find, for every length- ℓ substring $S[j..(j+\ell-1)]$ of S , an anchor $H \in \mathcal{A}$ such that $S[j..(j+\ell-1)] = H$, if any exists. This can be done by finding the nodes (implicit or explicit) of ST that correspond to the anchors, and then scanning over all

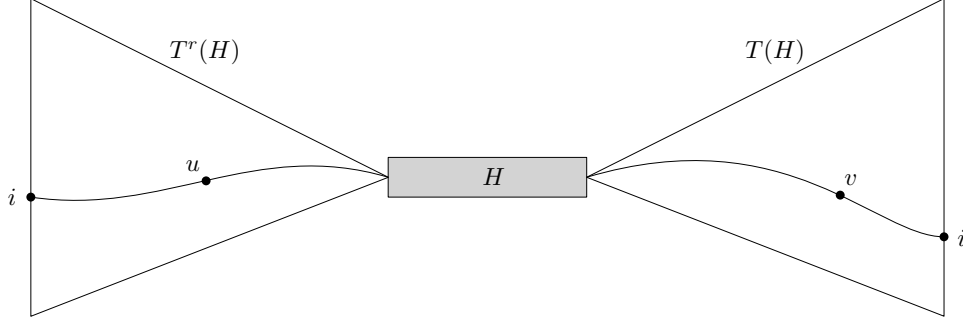


Figure 1: An occurrence of S starting at position i in P is generated by H : (u, v) corresponds to $S[j \dots (j + |H| - 1)] = H$ and i appears in the subtree of $T^r(H)$ rooted at u , as well as in the subtree of $T(H)$ rooted at v .

length- ℓ substrings while maintaining the node of ST corresponding to the current substring using suffix links in $\mathcal{O}(|S|)$ total time. After having determined that $S[j \dots (j + \ell - 1)] = H$ we add (u, v) to $L(H)$, where u and v are the previously found nodes of ST^r and ST corresponding to $(S[1 \dots (j - 1)])^r$ and $S[(j + \ell) \dots |S|]$, respectively. By construction, we have the following property, also illustrated in Figure 1.

Fact 2. *A string $S \in \mathcal{S}$ starts at position $i - j + 1$ in P if and only if, for some anchor $H \in \mathcal{A}$, $L(H)$ contains a pair (u, v) corresponding to $S[j \dots (j + |H| - 1)] = H$, such that the subtree of $T^r(H)$ rooted at u and that of $T(H)$ rooted at v contain a leaf decorated with i .*

Note that the overall size of all lists $L(H)$, when summed up over all $H \in \mathcal{A}$, is $\sum_{S \in \mathcal{S}} \text{occ}_S = \mathcal{O}(|S| \log m)$, and since each S is of length at least ℓ this is $\mathcal{O}(N/\ell \cdot \log m)$.

Processing. The goal of processing $D(H)$ is to efficiently process all occurrences generated by H . As a preliminary step, we decompose $T^r(H)$ and $T(H)$ into heavy paths. Then, for every pair of leaves $u \in T^r(H)$ and $v \in T(H)$ decorated by the same i , we consider all heavy paths above u and v . Let $p = u_1 - u_2 - \dots$ be a heavy path above u in $T^r(H)$ and $q = v_1 - v_2 - \dots$ be a heavy path above v in $T(H)$, where u_1 is the head of p and v_1 is the head of q , respectively. Further, choose the largest x such that u is in the subtree rooted at u_x , and the largest y such that v is in the subtree rooted at v_y (this is well-defined by the choice of p and q , as u is in the subtree rooted at u_1 and v is in the subtree rooted at v_1). We add $(i, |\mathcal{L}(u_x)|, |\mathcal{L}(v_y)|)$ to an auxiliary list associated with the pair of heavy paths (p, q) . In the rest of the processing we work with each such lists separately. Notice that the overall size of all auxiliary lists, when summed up over all $H \in \mathcal{A}$, is $\mathcal{O}(m/\ell \cdot \log^3 m)$, because there are at most $\log^2 m$ pairs of heavy paths above u and v decorated by the same i , and the total number of leaves in all trees $T^r(H)$ and $T(H)$ is bounded by the total number of occurrences of all anchors in P , which is $\mathcal{O}(m/\ell \cdot \log m)$. By Fact 2, there is an occurrence of a string \mathcal{S} generated by H and starting at position $i - j + 1$ in P if and only if $L(H)$ contains a pair (u, v) corresponding to $S[j \dots (j + |H| - 1)] = H$ such that, denoting by p the heavy path containing u in $T^r(H)$ and by q the heavy path containing v in $T(H)$, the auxiliary list associated with (p, q) contains a triple (i, x, y) such that $x \geq |\mathcal{L}(u)|$ and $y \geq |\mathcal{L}(v)|$. This is illustrated in Figure 2. Henceforth, we focus on the problem of processing a single auxiliary list associated with (p, q) , together with a list of pairs (u, v) , such that u belongs to p and v belongs to q .

An auxiliary list can be interpreted geometrically as follows: for every (i, x, y) we create a red point (x, y) , and for every (u, v) we create a blue point $(|\mathcal{L}(u)|, |\mathcal{L}(v)|)$. Then, each occurrence of $S \in \mathcal{S}$ generated by H corresponds to a pair of points (p_1, p_2) such that p_1 is red, p_2 is blue, and p_1 dominates p_2 . We further reduce this to a collection of simpler instances in which all red points already dominate all blue points. This can be done with a divide-and-conquer procedure

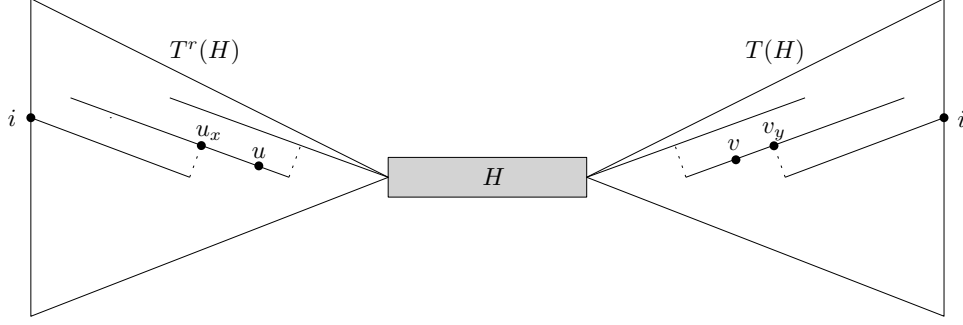


Figure 2: An occurrence of S starting at position i in P corresponds to a triple $(i, \mathcal{L}(u_x), \mathcal{L}(v_y))$ on some auxiliary list.

which is essentially equivalent to constructing a 2D range tree [11]. The total number of points in all obtained instances increases by a factor of $\mathcal{O}(\log^2 m)$, making the total number of red points in all instances $\mathcal{O}(m/\ell \cdot \log^5 m)$, while the total number of blue points is $\mathcal{O}(N/\ell \cdot \log^3 m)$. There is an occurrence of a string $S \in \mathcal{S}$ generated by H and starting at position $i - j + 1$ in P if and only if some simpler instance contains a red point created for some (i, x, y) and a blue point created for some (u, v) corresponding to $S[j \dots (j + |H| - 1)] = H$. In the following we focus on processing a single simpler instance.

To process a simpler instance we need to check if $U[i - j] = 1$, for a red point created for some (i, x, y) and a blue point created for some (u, v) corresponding to $S[j \dots (j + |H| - 1)] = H$, and if so set $V[i - j + |S|] = 1$. This has a natural interpretation as an instance of BMM: we create a $\lceil 5/4 \cdot \ell \rceil \times \lceil 5/4 \cdot \ell \rceil$ matrix M such that $M[\lceil 5/4 \cdot \ell \rceil - j, \lceil 5/4 \cdot \ell \rceil + 1 - j] = 1$ if and only if there is a blue point created for some (u, v) corresponding to $S[j \dots (j + |H| - 1)] = H$; then for every red point created for some (i, x, y) we construct a bit vector $U_i = U[(i - \lceil 5/4 \cdot \ell \rceil) \dots (i - 1)]$ (if $i < \lceil 5/4 \cdot \ell \rceil$, we pad U_i with 0s to make its length always equal to $\lceil 5/4 \cdot \ell \rceil$); calculate $V_i = M \times U_i$; and finally set $V[i + j] = 1$ whenever $V_i[j] = 1$ (and $i + j \leq m$).

Lemma 7. $V_i[k] = 1$ if and only if there is a blue point created for some (u, v) corresponding to $S[j \dots (j + |H| - 1)] = H$ such that $U[i - j] = 1$ and $k = |S| - j$.

Proof. By definition of $V_i = M \times U_i$, we have that $V_i[k] = 1$ if and only if $M[k, t] = 1$ for some t such that $U_i[t] = 1$. By definition of U_i , we have that $U_i[t] = 1$ if and only if $U[i - \lceil 5/4 \cdot \ell \rceil + t - 1] = 1$, and hence the previous condition can be rewritten as $M[k, t] = 1$ and $U[i - \lceil 5/4 \cdot \ell \rceil + t - 1] = 1$, or equivalently, by substituting $j = \lceil 5/4 \cdot \ell \rceil + 1 - t$, $M[k, \lceil 5/4 \cdot \ell \rceil + 1 - j] = 1$ and $U[i - j] = 1$. By definition of M , we have that $M[k, \lceil 5/4 \cdot \ell \rceil + 1 - j] = 1$ if and only if there is a blue point created for some (u, v) corresponding to $S[j \dots (j + |H| - 1)] = H$ with $k = |S| - j$, which proves the lemma. \square

The total length of all vectors U_i and V_i is $\mathcal{O}(m \log^5 m)$, so we can afford to extract the appropriate fragment of U and then update the corresponding fragment of V . The bottleneck is computing the matrix-vector product $V_i = M \times U_i$. Since the total number of 1s in all matrices M is bounded by the total number of blue points, a naïve method would take $\mathcal{O}(N/\ell \cdot \log^3 m)$ time; we overcome this by processing together all multiplications concerning the same matrix M , thus amortizing the costs. Let $U_{i_1}, U_{i_2}, \dots, U_{i_s}$ be all bit vectors that need to be multiplied with M , and let z a parameter to be determined later. We distinguish between two cases: (i) if $s < z$, then we compute the products naïvely by iterating over all 1s in M , and the total computation time, when summed up over all such matrices M , is $\mathcal{O}(N/\ell \cdot \log^3 m \cdot z)$; (ii) if $s \geq z$, then we partition the bit vectors into $\lceil s/z \rceil \leq s/z + 1$ groups of z (padding the last group with bit vectors containing all 0s) and, for every group, we create a single matrix whose columns contain all the bit vectors belonging to the group. Thus, we reduce the problem of computing all matrix-vector

products $M \times U_i$ to that of computing $\mathcal{O}(s/z)$ matrix-matrix products of the form $M \times M'$, where M' is an $\lceil 5/4 \cdot \ell \rceil \times z$ matrix. Even if M' is not necessarily a square matrix, we can still apply the fast matrix multiplication algorithm to compute $M \times M'$ using the standard trick of decomposing the matrices into square blocks.

Lemma 8. *If two $\mathcal{N} \times \mathcal{N}$ matrices can be multiplied in $\mathcal{O}(\mathcal{N}^\omega)$ time, then, for any $\mathcal{N} \geq \mathcal{N}'$, an $\mathcal{N} \times \mathcal{N}$ and an $\mathcal{N} \times \mathcal{N}'$ matrix can be multiplied in $\mathcal{O}((\mathcal{N}/\mathcal{N}')^2 \mathcal{N}'^\omega)$ time.*

Proof. We partition both matrices into blocks of size $\mathcal{N}' \times \mathcal{N}'$. There are $(\mathcal{N}/\mathcal{N}')^2$ such blocks in the first matrix and \mathcal{N}/\mathcal{N}' in the second matrix. Then, to compute the product we multiply each block from the first matrix by the appropriate block in the second matrix in $\mathcal{O}(\mathcal{N}'^\omega)$ time, resulting in the claimed complexity. \square

By applying Lemma 8, we can compute $M \times M'$ in $\mathcal{O}(\ell^2 z^{\omega-2})$ time (as long as we later verify that $5/4 \cdot \ell \geq z$), so all products $M \times U_i$ can be computed in $\mathcal{O}(\ell^2 z^{\omega-2} \cdot (s/z + 1))$ time. Note that this case can occur only $\mathcal{O}(m/(\ell \cdot z) \cdot \log^5 m)$ times, because all values of s sum up to $\mathcal{O}(m/\ell \cdot \log^5 m)$. This makes the total computation time, when summed up over all such matrices M , $\mathcal{O}(\ell^2 z^{\omega-2} \cdot m/(\ell \cdot z) \cdot \log^5 m) = \mathcal{O}(\ell z^{\omega-3} \cdot m \log^5 m)$. We can now prove our final result for strings of type 1.

Theorem 4. *An instance of the AP problem where all strings are of type 1 can be solved in $\tilde{\mathcal{O}}(m^{\omega-1} + N)$ time.*

Proof. The total time complexity is first $\mathcal{O}(m+N)$ to construct the graph G , then $\mathcal{O}(m \log m + N)$ to solve its corresponding instances of the NODESELECTION problem and obtain the set of anchors H . The time to initialize all structures $D(H)$ is $\mathcal{O}(m \log m + N)$. For every $D(H)$, we obtain in $\mathcal{O}(m/\ell \cdot \log^5 m + N/\ell \cdot \log^3 m)$ time a number of simpler instances, and then construct the corresponding Boolean matrices M and bit vectors U_i in additional $\mathcal{O}(m \log^5 m)$ time. Note that some M might be sparse, so we need to represent them as a list of 1s. Then, summing up over all matrices M and both cases, we spend $\mathcal{O}(N/\ell \cdot \log^3 m \cdot z + \ell z^{\omega-3} \cdot m \log^5 m)$ time. We would like to assume that $\ell \geq \log^3 m$ so that we can set $z = \ell/\log^3 m$. This is indeed possible, because for any t we can switch to a more naïve approach to process all strings of length at most t in $\mathcal{O}(mt^2 + N)$ time as described in 3. After applying it with $t = \log^3 m$ in $\mathcal{O}(m \log^6 m + N)$ time, we can set $z = \ell/\log^3 m$ (so that indeed $5/4 \cdot \ell \geq z$ as required in case $s \geq z$) and the overall time complexity for all matrices M and both cases becomes $\mathcal{O}(N + \ell^{\omega-2} \cdot m \log^{5+3(3-\omega)} m)$. Summing up over all values of ℓ and taking the initialization into account we obtain $\mathcal{O}(m \log^7 m + m^{\omega-1} \log^{5+3(3-\omega)} m + N) = \tilde{\mathcal{O}}(m^{\omega-1} + N)$ total time. \square

5.2 Type 2 Strings

In this section we show how to solve a restricted instance of the AP problem where every string $S \in \mathcal{S}$ is of type 2, that is, S contains a length- ℓ substring that is not strongly periodic as well as a length- ℓ substring that is strongly periodic, and furthermore $|S| \in [9/8 \cdot \ell, 5/4 \cdot \ell)$ for some $\ell \leq m$.

Similarly as in Section 5.1, we select a set of anchors. In this case, instead of the NODESELECTION problem we need to exploit periodicity. We call a string T ℓ -periodic if $|T| \geq \ell$ and $\text{per}(T) \leq \ell/4$. We consider all maximal ℓ -periodic substrings of S , that is, ℓ -periodic substrings $S[i..j]$ such that either $i = 1$ or $\text{per}(S[(i-1)..j]) > \ell/4$, and $j = |S|$ or $\text{per}(S[i..(j+1)]) > \ell/4$. We know that S contains at least one such substring (because there exists a length- ℓ substring that is strongly periodic), and that the whole S is not such a substring (because otherwise S would be of type 3). Further, two maximal ℓ -periodic substrings cannot overlap too much, as formalized in the following lemma.

Lemma 9. *Any two distinct maximal ℓ -periodic substrings of the same string S overlap by less than $\ell/2$ letters.*

Proof. Assume (by contradiction) the opposite; then we have two distinct ℓ -periodic substrings $S[i..j]$ and $S[i'..j']$ such that $i < i' \leq j < j'$ and $j - i' + 1 \geq \ell/2$. Then, both $p = \text{per}(S[i..j])$ and $p' = \text{per}(S[i'..j'])$ are periods of $S[i'..j]$, and hence by Lemma 1 we have that $\gcd(p, p')$ is a period of $S[i'..j]$. If $p \neq p'$ then, because $S[i'..j]$ contains an occurrence of both $S[i..(i+p-1)]$ and $S[i'..(i'+p'-1)]$, we obtain that one of these two substrings is a power of a shorter string, thus contradicting the definition of p or p' . So $p = p'$, but then $p \leq \ell/4$ is actually a period of the whole $S[i..j']$, meaning that $S[i..j]$ and $S[i'..j']$ are not maximal, a contradiction. \square

By Lemma 9, every $S \in \mathcal{S}$ contains exactly one maximal ℓ -periodic substring, and by the same argument P contains $\mathcal{O}(m/\ell)$ such substrings. The set of anchors will be generated by considering the unique maximal ℓ -periodic substring of every $S \in \mathcal{S}$, so we first need to show how to efficiently generate such substrings.

Lemma 10. *Given a string S of length at most $5/4 \cdot \ell$, we can generate its (unique) maximal ℓ -periodic substring in $\mathcal{O}(|S|)$ time.*

Proof. We start with observing that any length- ℓ substring of S must contain $S[(\lfloor \ell/2 \rfloor + 1)..\ell]$ inside. Consequently, we can proceed similarly as in the proof of Lemma 5. We compute $p = \text{per}(S[(\lfloor \ell/2 \rfloor + 1)..\ell])$ in $\mathcal{O}(|S|)$ time. If $p > \ell/4$ then S does not contain any ℓ -periodic substrings. Otherwise, we compute in $\mathcal{O}(|S|)$ time how far the period p extends to the left and to the right; that is, we compute the smallest $i \leq \lfloor \ell/2 \rfloor + 1$ such that $S[k] = S[k+p]$ for every $k = i, i+1, \dots, \lfloor \ell/2 \rfloor$ and the largest $j \geq \ell$ such that $S[k] = S[k-p]$ for every $k = \ell+1, \ell+2, \dots, j$. If $j - i + 1 \geq \ell$ then $S[i..j]$ is a maximal ℓ -periodic substring of S , and, as shown earlier by Lemma 9, S cannot contain any other maximal ℓ -periodic substrings. We return $S[i..j]$ as the (unique) maximal ℓ -periodic substring of S . \square

For every $S \in \mathcal{S}$, we apply Lemma 10 on S to find its (unique) maximal ℓ -periodic substring $S[i..j]$ in $\mathcal{O}(|S|)$ time. If $i > 1$ then we designate $S[(i-1)..(i-1+\ell)]$ as an anchor, and similarly if $j < |S|$ we designate $S[(j+1-\ell)..(j+1)]$ as an anchor. Observe that because S is of type 2 (and not of type 3) either $i > 1$ or $j < |S|$, so for every $S \in \mathcal{S}$ we designate at least one of its length- $(\ell+1)$ substrings as an anchor. As in Section 5.1, we represent each anchor by one of its occurrences in P , and so need to find its corresponding node in the suffix tree of P (if any). This can be done in $\mathcal{O}(|S|)$ time, so $\mathcal{O}(N)$ overall. During this process we might designate the same string as an anchor multiple times, but we can easily remove the possible duplicates to obtain the set \mathcal{A} of anchors in the end. Then, we generate the occurrences of all anchors in P by accessing their corresponding nodes in the suffix tree of P and iterating over all leaves in their subtrees. We claim that the total number of all these occurrences is only $\mathcal{O}(m/\ell)$. This follows from the following characterization.

Lemma 11. *If $P[x..(x+\ell)]$ is an occurrence of an anchor then either $P[(x+1)..y]$ is a maximal ℓ -periodic substring of P , for some $y \geq x+\ell$, or $P[x'..(x+\ell-1)]$ is a maximal ℓ -periodic substring of P , for some $x' \leq x$.*

Proof. By symmetry, it is enough to consider an anchor H created because of a maximal ℓ -periodic substring $S[i..j]$ such that $i > 1$, when we add $S[(i-1)..(i-1+\ell)]$ to \mathcal{A} . Thus, $\text{per}(H[2..|H|]) \leq \ell/4$ and if $P[x..(x+\ell)] = H$ then $\text{per}(P[(x+1)..(x+\ell)]) \leq \ell/4$, making $P[(x+1)..(x+\ell)]$ a substring of some maximal ℓ -periodic substring of $P[(x'+1)..y]$, where $x' \leq x$ and $y \geq x+\ell$. If $x' < x$ then $\text{per}(H) \leq \ell/4$. But then $H = S[(i-1)..(i-1+\ell)]$ can be extended to some maximal ℓ -periodic substring $S[i'..j']$ such that $i' \leq i-1$ and $j' \geq i-1+\ell$. The overlap between $S[i..j]$ and $S[i'..j']$ is at least ℓ , so by Lemma 9 $i = i'$ and $j = j'$, which is a contradiction. Consequently, $x' = x$ and we obtain the lemma. \square

By Lemma 11, the number of occurrences of all anchors in P is at most two per each maximal ℓ -periodic substrings, so $\mathcal{O}(m/\ell)$ in total. We thus obtain a set of length- $(\ell+1)$ anchors with the following properties:

1. The total number of occurrences of all anchors in P is $\mathcal{O}(m/\ell)$.
2. For every $S \in \mathcal{S}$, at least one of its length- $(\ell + 1)$ substrings is an anchor.
3. For every $S \in \mathcal{S}$, at most two of its length- $(\ell + 1)$ substrings are anchors.

These properties are even stronger than what we had used in Section 5.1 (except that now we are working with length- $(\ell + 1)$ substrings, which is irrelevant), we can now prove our final result also for strings of type 2.

Theorem 5. *An instance of the AP problem where all strings are of type 2 can be solved in $\tilde{\mathcal{O}}(m^{\omega-1} + N)$ time.*

5.3 Type 3 Strings

In this section we show how to solve a restricted instance of the AP problem where every string $S \in \mathcal{S}$ is of type 3, that is, $\text{per}(S) \leq \ell/4$. Furthermore $|S| \in [9/8 \cdot \ell, 5/4 \cdot \ell]$ for some $\ell \leq m$. Recall that strings $S \in \mathcal{S}$ are such that every length- ℓ substring of S is strongly periodic and, by Lemma 4, in this case, $\text{per}(S) \leq \ell/4$. An occurrence of such S in P must be contained in a maximal ℓ -periodic substring of P . Recall that a string T is called ℓ -periodic if $|T| \geq \ell$ and $\text{per}(T) \leq \ell/4$. For an ℓ -periodic string T , let its *root*, denoted by $\text{root}(T)$, be the lexicographically smallest cyclic shift of $T[1 \dots \text{per}(T)]$. Because $\text{per}(T) \leq \ell/4$ and $|T| \geq \ell$ by definition, there are at least four repetitions of the period in T , so we can write $T = R[i \dots |R|]R^\alpha R[1 \dots j]$, where $R = \text{root}(T)$, for some $i, j \in [1, |R|]$ and $\alpha \geq 2$. It is well known that $\text{root}(T)$ can be computed in $\mathcal{O}(|T|)$ time [27].

Example 4. Let $T = \text{babababab}$ and $\ell = 8$. We have $|T| = 9 \geq \ell = 8$ and $\text{per}(T) = 2 \leq \ell/4 = 2$, so T is ℓ -periodic. We have $\text{root}(T) = R = \text{ab}$, and T can be written as $T = \text{b} \cdot (\text{ab})^3 \cdot \text{ab}$, for $i = 2$ and $j = 2$.

We will now make a partition of type 3 strings based on their root. We start with extracting all maximal ℓ -periodic substrings of P using Lemma 10 and compute the root of every such substring. This can be done in $\mathcal{O}(m)$ total time because two maximal ℓ -periodic substrings cannot overlap by more than $\ell/2$ letters, and hence their total length is at most $3/2 \cdot \ell$. We also extract the root of every $S \in \mathcal{S}$ in $\mathcal{O}(N)$ total time. We then partition maximal ℓ -periodic substrings of P and strings $S \in \mathcal{S}$ into groups that have the same root. In the remaining part we describe how to process each such group corresponding to root R in which all maximal ℓ -periodic substrings of P have total length m' , and the strings $S \in \mathcal{S}$ have total length N' .

Recall that the bit vector U stores the active prefixes input to the AP problem, and the bit vector V encodes the new active prefixes we aim to compute. For every maximal ℓ -periodic substring of P with root R we extract the corresponding fragment of the bit vector U and need to update the corresponding fragment of the bit vector V . To make the description less cluttered, we assume that each such substring of P is a power of R , that is, R^α for some $\alpha \geq 4$. This can be assumed without loss of generality as it can be ensured by appropriately padding the extracted fragment of U and then truncating the results, while increasing the total length of all considered substrings of P by at most half of their length. In the description below, for simplicity of presentation, U and V denote these padded fragments of the original U and V . When computing V from U we use two different methods for processing the elements $S = R[i \dots |R|]R^\beta R[1 \dots j]$ of \mathcal{S} depending on their length: either $\beta > \alpha/|R|$ (large β) or $\beta \leq \alpha/|R|$ (small β).

Large β . We proceed in phases corresponding to $\beta = \alpha/|R| + 1, \dots, \alpha$. In each single phase, we consider all strings $S \in \mathcal{S}$ with $S = R[i \dots |R|]R^\beta R[1 \dots j]$, for some i and j . Let $C(\beta)$ be the set of the corresponding pairs (i, j) , and observe that $\sum_\beta |C(\beta)| \cdot |R^\beta| \leq N'$. This is because the length of R^β is not greater than that of $S = R[i \dots |R|]R^\beta R[1 \dots j]$, there are $|C(\beta)|$ distinct strings of

the latter form, and the total length of all $S \in \mathcal{S}$ is N' . The total number of occurrences of a string $S = R[i \dots |R|]R^\beta R[1 \dots j]$ in R^α is bounded by $\mathcal{O}(\alpha)$, and all such occurrences can be generated in time proportional to their number. Thus, for every $(i, j) \in C(\beta)$, we can generate all occurrences of the corresponding string and appropriately update V in $\mathcal{O}(\alpha \cdot |C(\beta)|)$ total time.

Small β . We start by giving a technical lemma on the complexity of multiplying two $r \times r$ matrices whose cells are polynomials of degree up to d .

Lemma 12. *If two $r \times r$ matrices over \mathbb{Z} can be multiplied in $\mathcal{O}(r^\omega)$ time, then two $r \times r$ matrices over $\mathbb{Z}[X]$ with degrees up to d can be multiplied in $\tilde{\mathcal{O}}(dr^\omega)$ time.*

Proof. Let A and B be two $r \times r$ matrices over $\mathbb{Z}[X]$ with degrees up to d . We reduce the product $A \cdot B = C$ to $2d$ products of $r \times r$ matrices over \mathbb{Z} as follows. We evaluate the polynomials of each matrix in the complex $(2d)$ th roots of unity: let A_i and B_i be the matrices obtained by evaluating the polynomials of A and B in the i -th such root, respectively. We then perform the $2d$ products $A_1 \cdot B_1, \dots, A_{2d} \cdot B_{2d}$ to obtain matrices C_1, \dots, C_{2d} : the $2d$ values $C_1[i, j], \dots, C_{2d}[i, j]$ are finally interpolated to obtain the coefficient representation of $C[i, j]$, for each $i, j = 1, \dots, r$, in $\mathcal{O}(d \log d)$ time for each polynomial [23]. Since we perform $2d$ products of matrices in $\mathbb{Z}^{r \times r}$, and we evaluate and interpolate r^2 polynomials of degree up to $2d$, the overall time complexity is $2d\mathcal{O}(r^\omega) + r^2\mathcal{O}(d \log d) = \tilde{\mathcal{O}}(dr^\omega)$. \square

Unlike in the large β case, we process $\beta = 2, \dots, \alpha/|R|$ simultaneously as follows. For each β we construct an $|R| \times |R|$ matrix M_β , with $M_\beta[i, j] = 1$ if and only if $(i, j) \in C(\beta)$ (and $M_\beta[i, j] = 0$ otherwise), and collect them in a single 3D matrix $M \in \{0, 1\}^{|R| \times |R| \times (\alpha/|R|-1)}$ with the third dimension corresponding to the value of β . We then create another $\alpha \times |R|$ matrix, denoted by M' , by setting $M'[\gamma, i] = 1$ if and only if $U[\gamma \cdot |R| + i - 1] = 1$ (and $M_\beta[i, j] = 0$ otherwise). Observe that M' can be interpreted as a vector of length $|R|$ over $\mathbb{Z}[X]$ with degrees up to α , and M as an $|R| \times |R|$ matrix over $\mathbb{Z}[X]$ with degrees up to $\alpha/|R|$: in this way, x^γ appears with non-zero coefficient in the polynomial at $M'[i]$ if and only if $U[\gamma \cdot |R| + i - 1] = 1$, and x^β appears with non-zero coefficient in the polynomial at $M[i, j]$ if and only if $(i, j) \in C(\beta)$. M can be constructed in total $\mathcal{O}(N')$ time by first iterating over all $S \in \mathcal{S}$ and adding x^β to the polynomial at $M[i, j]$, where $S = R[i \dots |R|]R^\beta R[1 \dots j]$, and then extracting a prefix of each polynomial consisting of monomials of degree less than $\alpha/|R|$.

The product $M' \cdot M = M''$ allows us to recover the updates to V by observing that $V[(q+1) \cdot |R| + j] = 1$ if and only if x^q appears with non-zero coefficient in the polynomial at $M''[j]$. We aim at reducing this product to a matrix-matrix product over $\mathbb{Z}[X]$ with degrees up to $\alpha/|R|$, so as to compute it efficiently by applying Lemma 12.

The idea now is to decompose the columns of M' into $|R|$ chunks of size $\alpha/|R|$ in order to transform it into another 3D matrix. To this end, we transform M' into an $|R| \times |R| \times (\alpha/|R|)$ matrix A by setting

$$A[k, i, \gamma] = 1 \Leftrightarrow M'[(k-1)\alpha/|R| + \gamma, i] = 1 \Leftrightarrow U[(k-1)\alpha/|R| + \gamma + i - 1] = 1.$$

By interpreting A as an $|R| \times |R|$ matrix over $\mathbb{Z}[X]$ with degrees up to $\alpha/|R|$, and interpreting M' as a vector of length $|R|$ over $\mathbb{Z}[X]$ with degrees up to α , we have that the first row of A consists of the coefficients of $x^1, \dots, x^{\alpha/|R|}$ of each of the $|R|$ polynomials of M' , the second row consists of the coefficients of $x^{\alpha/|R|+1}, \dots, x^{2\alpha/|R|}$ of each of the $|R|$ polynomials of M' , and so on. In general, $A[k, i]$ consists of the coefficients of $x^{(k-1)\alpha/|R|+1}, \dots, x^{k\alpha/|R|}$ of polynomial $M'[i]$.

The product $A \cdot M = C$ still allows us to recover the updates of V , by observing that $V[((k-1)\alpha/|R| + 1 + q + 1)|R| + j] = 1$ if x^q appears with non-zero coefficient in the polynomial at $C[k, j]$. This is because at row $A[k, \cdot]$ there are the coefficients that correspond to $U[\gamma \cdot |R| + i - 1]$ for $\gamma = (k-1)\alpha/|R| + 1, \dots, k\alpha/|R|$ and $i = 1, \dots, |R|$, and hence a x^q appearing at $C[k, j]$ is equivalent to a $x^{q+(k-1)\alpha/|R|+1}$ at $M''[j]$.

We are now in a position to prove the following result for type 3 strings.

Theorem 6. *An instance of the AP problem where all strings are of type 3 can be solved in $\tilde{O}(m^{\omega-1} + N)$ time.*

Proof. Recall that we consider strings S of type 3 with root R and substrings of P with root R together. We first analyze the time to process a single group containing a number of substrings of P of total length m' and a number of strings $S \in \mathcal{S}$ of total length N' . Let us denote by R^{α_i} the i -th considered substring of P and further define $\alpha = \sum_i \alpha_i = m'/|R|$.

If $\beta > \alpha/|R|$ we use the first method and spend $\mathcal{O}(\alpha_i \cdot |C(\beta)|)$ time, where $C(\beta)$ is the set of (i, j) for this specific β . The overall time used for all applications of the first method is:

$$\begin{aligned} \sum_i \mathcal{O}(\alpha_i \cdot \sum_{\beta > \alpha/|R|} |C(\beta)|) &= \mathcal{O}(\alpha/|R^{\alpha/|R|}| \cdot \sum_{\beta > \alpha/|R|} |C(\beta)| \cdot |R^{\alpha/|R|}|) \\ &= \mathcal{O}(\sum_{\beta > \alpha/|R|} |C(\beta)| \cdot |R^\beta|) = \mathcal{O}(N'), \end{aligned}$$

using the fact that $\sum_\beta |C(\beta)| \cdot |R^\beta| \leq N'$ and $\alpha/|R^{\alpha/|R|}| = \alpha/(\alpha/|R|)|R| = \mathcal{O}(1)$.

For each α_i , we process together all $\beta \leq \alpha_i/|R|$ using the second method, and we need to multiply two $|R| \times |R|$ matrices of polynomials of degree up to $\alpha_i/|R|$, that we can build in time $\mathcal{O}(N')$ and multiply in time $\mathcal{O}(|R|^\omega \cdot \alpha_i/|R| + |R|^2(\alpha_i/|R|) \log(\alpha_i/|R|))$ by Lemma 12. The overall time used for all applications of the second method is:

$$\mathcal{O}(N') + \sum_i \mathcal{O}(|R|^\omega \cdot \alpha_i/|R| + |R|^2(\alpha_i/|R|) \log(\alpha_i/|R|)) = \mathcal{O}(|R|^{\omega-2}m' + m' \log m' + N'),$$

using the fact that $\alpha = m'/|R|$. Since $|R| \leq m'$, this is in fact $\mathcal{O}((m')^{\omega-1} + m' \log m' + N')$.

Because all values of N' sum up to N and all values of m' sum up to $\mathcal{O}(m)$, by convexity of $x^{\omega-1}$ we obtain that the overall time complexity is $\tilde{O}(m^{\omega-1} + N)$. \square

5.4 Wrapping Up

In Sections 5.1, 5.2 and 5.3 we design three $\tilde{O}(m^{\omega-1} + N)$ -time algorithms for an instance of the AP problem where all strings are of type 1, 2 and 3, respectively. We thus obtain Theorem 2, and using the fact that $\omega < 2.373$ [47, 60] we obtain the following corollary.

Corollary 7. *The EDSM problem can be solved on-line in $\mathcal{O}(nm^{1.373} + N)$ time.*

Note that the polylog factors are shaved from $\tilde{O}(nm^{\omega-1} + N)$ by using the fact that the inequality of $\omega < 2.373$ is strict.

6 Final Remarks

Our contribution in this paper is twofold. First, we designed an appropriate reduction showing that a combinatorial algorithm solving the EDSM problem in $\mathcal{O}(nm^{1.5-\epsilon} + N)$ time, for any $\epsilon > 0$, refutes the well-known BMM conjecture. Second, we designed a non-combinatorial $\tilde{O}(nm^{\omega-1} + N)$ -time algorithm to attack the same problem. By using the fact that $\omega < 2.373$, our algorithm runs in $\mathcal{O}(nm^{1.373} + N)$ time thus breaking the combinatorial conditional lower bound for the EDSM problem. Let us point out that if $\omega = 2$ then our algorithm for the AP problem is time-optimal up to polylog factors.

Acknowledgments

GR and NP are partially supported by MIUR-SIR project CMACBioSeq “Combinatorial methods for analysis and compression of biological sequences” grant n. RBSI146R5L.



This project has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 872539.

References

- [1] A. ABBOUD, A. BACKURS, AND V. WILLIAMS, *If the current clique algorithms are optimal, so is Valiant’s parser*, in 56th IEEE Symposium on Foundations Of Computer Science (FOCS), 2015, pp. 98–117.
- [2] A. ABBOUD AND V. WILLIAMS, *Popular conjectures imply strong lower bounds for dynamic problems*, in 55th IEEE Symposium on Foundations Of Computer Science (FOCS), 2014, pp. 434–443.
- [3] K. ABRAHAMSON, *Generalized string matching*, SIAM J. Comput., 16 (1987), pp. 1039–1051.
- [4] M. ALZAMEL, L. AYAD, G. BERNARDINI, R. GROSSI, C. ILIOPOULOS, N. PISANTI, S. PISSIS, AND G. ROSONE, *Degenerate string comparison and applications*, in 18th Workshop on Algorithms in Bioinformatics (WABI), vol. 113 of LIPIcs, 2018, pp. 21:1–21:14.
- [5] A. AMIR, M. LEWENSTEIN, AND E. PORAT, *Faster algorithms for string matching with k mismatches*, J. Algorithms, 50 (2004), pp. 257–275.
- [6] K. AOYAMA, Y. NAKASHIMA, T. I, S. INENAGA, H. BANNAI, AND M. TAKEDA, *Faster online elastic degenerate string matching*, in 29th Symposium on Combinatorial Pattern Matching (CPM), vol. 105 of LIPIcs, 2018, pp. 9:1–9:10.
- [7] V. ARLAZAROV, E. DINIC, M. KRONROD, AND I. FARADŽEV, *On economical construction of the transitive closure of a directed graph*, Soviet Mathematics Doklady, 11 (1970), pp. 1209–1210.
- [8] A. BACKURS AND P. INDYK, *Which regular expression patterns are hard to match?*, in 57th IEEE Symposium on Foundations Of Computer Science (FOCS), 2016, pp. 457–466.
- [9] N. BANSAL AND R. WILLIAMS, *Regularity lemmas and combinatorial algorithms*, in 50th IEEE Symposium on Foundations Of Computer Science (FOCS), 2009, pp. 745–754.
- [10] M. BENDER AND M. FARACH-COLTON, *The LCA problem revisited*, in 4th Latin American symposium on Theoretical INformatics (LATIN), vol. 1776 of Springer LNCS, 2000, pp. 88–94.
- [11] J. BENTLEY, *Multidimensional binary search trees used for associative searching*, Commun. ACM, 18 (1975), pp. 509–517.
- [12] G. BERNARDINI, P. GAWRYCHOWSKI, N. PISANTI, S. P. PISSIS, AND G. ROSONE, *Even Faster Elastic-Degenerate String Matching via Fast Matrix Multiplication*, in 46th International Colloquium on Automata, Languages, and Programming (ICALP), vol. 132 of LIPIcs, 2019, pp. 21:1–21:15.
- [13] G. BERNARDINI, N. PISANTI, S. P. PISSIS, AND G. ROSONE, *Approximate pattern matching on elastic-degenerate text*, Theor. Comput. Sci., 812 (2020), pp. 109–122.

- [14] K. BRINGMANN, F. GRANDONI, B. SAHA, AND V. WILLIAMS, *Truly sub-cubic algorithms for language edit distance and RNA-folding via fast bounded-difference min-plus product*, in 56th IEEE Symposium on Foundations Of Computer Science (FOCS), 2016, pp. 375–384.
- [15] K. BRINGMANN, A. GRØNLUND, AND K. LARSEN, *A dichotomy for regular expression membership testing*, in 58th IEEE Symposium on Foundations Of Computer Science (FOCS), 2017, pp. 307–318.
- [16] T. CHAN, *Speeding up the four Russians algorithm by about one more logarithmic factor*, in 26th ACM-SIAM Symposium On Discrete Algorithms (SODA), 2015, pp. 212–217.
- [17] Y.-J. CHANG, *Hardness of RNA folding problem with four symbols*, in 27th Symposium on Combinatorial Pattern Matching (CPM), vol. 54 of LIPIcs, 2016, pp. 13:1–13:12.
- [18] K. CHATTERJEE, B. CHOUDHARY, AND A. PAVLOGIANNIS, *Optimal Dyck reachability for data-dependence and alias analysis*, in 45th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL), 2018, pp. 30:1–30:30.
- [19] A. CISLAK AND S. GRABOWSKI, *Sopang 2: online searching over a pan-genome without false positives*, CoRR, abs/2004.03033 (2020), <https://arxiv.org/abs/2004.03033>.
- [20] A. CISLAK, S. GRABOWSKI, AND J. HOLUB, *SOPanG: online text searching over a pan-genome*, Bioinformatics, 34 (2018), pp. 4290–4292.
- [21] P. CLIFFORD AND R. CLIFFORD, *Simple deterministic wildcard matching*, Inf. Process. Lett., 101 (2007), pp. 53–54.
- [22] R. COLE AND R. HARIHARAN, *Verifying candidate matches in sparse and wildcard matching*, in 34th ACM Symposium on Theory Of Computing (STOC), 2002, pp. 592–601.
- [23] J. COOLEY AND J. TUKEY, *An algorithm for the machine calculation of complex Fourier series*, Mathematics of Computation, 19 (1965), pp. 297–301.
- [24] M. CROCHEMORE, C. HANCART, AND T. LECROQ, *Algorithms on strings*, Cambridge University Press, 2007.
- [25] M. CROCHEMORE AND D. PERRIN, *Two-way string matching*, J. ACM, 38 (1991), pp. 651–675.
- [26] A. CZUMAJ AND A. LINGAS, *Finding a heaviest vertex-weighted triangle is not harder than matrix multiplication*, SIAM J. Comput., 39 (2009), pp. 431–444.
- [27] J. DUVAL, *Factorizing words over an ordered alphabet*, J. Algorithms, 4 (1983), pp. 363–381.
- [28] M. FARACH-COLTON AND S. MUTHUKRISHNAN, *Perfect hashing for strings: formalization and algorithms*, in 7th Annual Symposium on Combinatorial Pattern Matching (CPM), vol. 1075 of Springer LNCS, 1996, pp. 130–140.
- [29] N. FINE AND H. WILF, *Uniqueness theorems for periodic functions*, Proceedings of the American Mathematical Society, 16 (1965), pp. 109–114.
- [30] M. FISCHER AND A. MEYER, *Boolean matrix multiplication and transitive closure*, in 12th IEEE Symposium on Switching and Automata Theory (SWAT/FOCS), 1971, pp. 129–131.
- [31] M. FISCHER AND M. PATERSON, *String matching and other products*, in 7th SIAM-AMS Complexity of Computation, 1974, pp. 113–125.

- [32] M. FURMAN, *Application of a method of fast multiplication of matrices in the problem of finding the transitive closure of a graph*, Soviet Mathematics Doklady, 11 (1970), p. 1252.
- [33] P. GAWRYCHOWSKI, S. GHAZAWI, AND G. M. LANDAU, *On indeterminate strings matching*, in 31st Symposium on Combinatorial Pattern Matching (CPM), vol. 161 of LIPIcs, 2020, pp. 23:1–23:12.
- [34] P. GAWRYCHOWSKI AND P. UZNANSKI, *Towards unified approximate pattern matching for Hamming and L_1 distance*, in 45th International Colloquium on Automata, Languages and Programming (ICALP), vol. 107 of LIPIcs, 2018, pp. 62:1–62:13.
- [35] R. GROSSI, C. ILIOPOULOS, C. LIU, N. PISANTI, S. PISSIS, A. RETHA, G. ROSONE, F. VAYANI, AND L. VERSARI, *On-line pattern matching on similar texts*, in 28th Symposium on Combinatorial Pattern Matching (CPM), vol. 78 of LIPIcs, 2017, pp. 9:1–9:14.
- [36] M. HENZINGER, S. KRINNINGER, D. NANONGKAI, AND T. SARANURAK, *Unifying and strengthening hardness for dynamic problems via the online matrix-vector multiplication conjecture*, in 47th ACM Symposium on Theory Of Computing (STOC), 2015, pp. 21–30.
- [37] J. HOLUB, W. SMYTH, AND S. WANG, *Fast pattern-matching on indeterminate strings*, J. Discrete Algorithms, 6 (2008), pp. 37–50.
- [38] C. ILIOPOULOS, R. KUNDU, AND S. PISSIS, *Efficient pattern matching in elastic-degenerate texts*, in 11th International Conference on Language and Automata Theory and Applications (LATA), vol. 10168 of Springer LNCS, 2017, pp. 131–142.
- [39] P. INDYK, *Faster algorithms for string matching problems: Matching the convolution bound*, in 39th Symposium on Foundations Of Computer Science (FOCS), 1998, pp. 166–173.
- [40] A. ITAI AND M. RODEH, *Finding a minimum circuit in a graph*, in 9th ACM Symposium on Theory Of Computing (STOC), 1977, pp. 1–10.
- [41] IUPAC-IUB COMMISSION ON BIOCHEMICAL NOMENCLATURE, *Abbreviations and symbols for nucleic acids, polynucleotides, and their constituents*, Biochemistry, 9 (1970), pp. 4022–4027.
- [42] A. KALAI, *Efficient pattern-matching with don't cares*, in 13th ACM-SIAM Symposium on Discrete Algorithms (SODA), 2002, pp. 655–656.
- [43] D. KNUTH, J. M. JR., AND V. PRATT, *Fast pattern matching in strings*, SIAM J. Comput., 6 (1977), pp. 323–350.
- [44] T. KOCIUMAKA, J. RADOSZEWSKI, W. RYTTER, AND T. WALEŃ, *Internal pattern matching queries in a text and applications*, in 26th ACM-SIAM Symposium on Discrete Algorithms (SODA), 2015, pp. 532–551.
- [45] T. KOPELOWITZ AND R. KRAUTHGAMER, *Color-distance oracles and snippets*, in 27th Symposium on Combinatorial Pattern Matching (CPM), vol. 54 of LIPIcs, 2016, pp. 24:1–24:10.
- [46] K. LARSEN, I. MUNRO, J. NIELSEN, AND S. THANKACHAN, *On hardness of several string indexing problems*, Theor. Comput. Sci., 582 (2015), pp. 74–82.
- [47] F. LE GALL, *Powers of tensors and fast matrix multiplication*, in 39th International Symposium on Symbolic and Algebraic Computation (ISSAC), 2014, pp. 296–303.
- [48] L. LEE, *Fast context-free grammar parsing requires fast boolean matrix multiplication*, J. ACM, 49 (2002), pp. 1–15.

- [49] V. MÄKINEN, B. CAZAUX, M. EQUI, T. NORRI, AND A. I. TOMESCU, *Linear time construction of indexable founder block graphs*, in 20th Workshop on Algorithms in Bioinformatics (WABI), vol. 172 of LIPIcs, 2020, pp. 7:1–7:18.
- [50] J. MATOUŠEK, *Computing dominances in E^n* , Inf. Process. Lett., 38 (1991), pp. 277–278.
- [51] I. MUNRO, *Efficient determination of the transitive closure of a directed graph*, Inf. Process. Lett., 1 (1971), pp. 56–58.
- [52] G. NAVARRO, *Nr-grep: a fast and flexible pattern-matching tool*, Softw., Pract. Exper., 31 (2001), pp. 1265–1312.
- [53] S. PISSIS AND A. RETHA, *Dictionary matching in elastic-degenerate texts with applications in searching VCF files on-line*, in 17th International Symposium on Experimental Algorithms (SEA), vol. 103 of LIPIcs, 2018, pp. 16:1–16:14.
- [54] L. RODITTY AND U. ZWICK, *On dynamic shortest paths problems*, in 12th European Symposium on Algorithms (ESA), vol. 3221 of Springer LNCS, 2004, pp. 580–591.
- [55] M. RUŽIĆ, *Constructing efficient dictionaries in close to sorting time*, in 35th International Colloquium on Automata, Languages and Programming (ICALP), vol. 5125 of Springer LNCS, 2008, pp. 84–95.
- [56] D. SLEATOR AND R. TARJAN, *A data structure for dynamic trees*, J. Comput. Syst. Sci., 26 (1983), pp. 362–391.
- [57] THE COMPUTATIONAL PAN-GENOMICS CONSORTIUM, *Computational pan-genomics: status, promises and challenges*, Briefings in Bioinformatics, 19 (2018), pp. 118–135.
- [58] L. VALIANT, *General context-free recognition in less than cubic time*, J. Comput. Syst. Sci., 10 (1975), pp. 308–315.
- [59] P. WEINER, *Linear pattern matching algorithms*, in 14th IEEE Annual Symposium on Switching and Automata Theory (SWAT/FOCS), 1973, pp. 1–11.
- [60] V. WILLIAMS, *Multiplying matrices faster than Coppersmith-Winograd*, in 44th ACM Symposium on Theory Of Computing Conference (STOC), 2012, pp. 887–898.
- [61] V. WILLIAMS AND R. WILLIAMS, *Finding a maximum weight triangle in $n^{3-\delta}$ time, with applications*, in 38th ACM Symposium on Theory Of Computing Conference (STOC), 2006, pp. 225–231.
- [62] V. WILLIAMS AND R. WILLIAMS, *Subcubic equivalences between path, matrix and triangle problems*, in 51st IEEE Symposium on Foundations Of Computer Science (FOCS), 2010, pp. 645–654.
- [63] S. WU AND U. MANBER, *Agrep – a fast approximate pattern-matching tool*, in USENIX Technical Conference, 1992, pp. 153–162.
- [64] H. YU, *An improved combinatorial algorithm for boolean matrix multiplication*, in 42nd International Colloquium on Automata, Languages, and Programming (ICALP), vol. 9134 of Springer LNCS, 2015, pp. 1094–1105.
- [65] H. YU, *An improved combinatorial algorithm for boolean matrix multiplication*, Inf. Comput., 261 (2018), pp. 240–247.
- [66] U. ZWICK, *All pairs shortest paths using bridging sets and rectangular matrix multiplication*, J. ACM, 49 (2002), pp. 289–317.