

Fine-Grained Complexity of Analyzing Compressed Data: Quantifying Improvements over Decompress-And-Solve

Amir Abboud* Arturs Backurs† Karl Bringmann‡ Marvin Künnemann§

March 5, 2018

Abstract

Can we analyze data without decompressing it? As our data keeps growing, understanding the time complexity of problems on *compressed* inputs, rather than in convenient uncompressed forms, becomes more and more relevant. Suppose we are given a compression of size n of data that originally has size N , and we want to solve a problem with time complexity $T(\cdot)$. The naïve strategy of “decompress-and-solve” gives time $T(N)$, whereas “the gold standard” is time $T(n)$: to analyze the compression as efficiently as if the original data was small.

We restrict our attention to data in the form of a string (text, files, genomes, etc.) and study the most ubiquitous tasks. While the challenge might seem to depend heavily on the specific compression scheme, most methods of practical relevance (Lempel-Ziv-family, dictionary methods, and others) can be unified under the elegant notion of *Grammar-Compressions*. A vast literature, across many disciplines, established this as an influential notion for Algorithm design.

We introduce a framework for proving (conditional) lower bounds in this field, allowing us to assess whether decompress-and-solve can be improved, and by how much. Our main results are:

- The $O(nN\sqrt{\log N/n})$ bound for LCS and the $O(\min\{N \log N, nM\})$ bound for Pattern Matching with Wildcards are optimal up to $N^{o(1)}$ factors, under the Strong Exponential Time Hypothesis. (Here, M denotes the uncompressed length of the compressed pattern.)
- Decompress-and-solve is essentially optimal for Context-Free Grammar Parsing and RNA Folding, under the k -Clique conjecture.
- We give an algorithm showing that decompress-and-solve is *not* optimal for Disjointness.

*IBM Almaden Research Center, abboud@cs.stanford.edu. Work done while at Stanford University.

†MIT, backurs@mit.edu

‡Max Planck Institute for Informatics, Saarland Informatics Campus, Germany, kbringma@mpi-inf.mpg.de

§Max Planck Institute for Informatics, Saarland Informatics Campus, Germany, marvin@mpi-inf.mpg.de

Contents

1	Introduction	1
1.1	Previous Work	2
1.2	Our Work	3
1.3	Technical Overview	6
2	Preliminaries	8
2.1	Hardness Assumptions	10
3	Tight Bounds Assuming SETH	11
3.1	DFA Acceptance	12
3.2	Approximate Pattern Matching and Substring Hamming Distance	15
3.3	Longest Common Subsequence	21
3.3.1	Alignment Gadget Framework	22
3.3.2	General Lower Bound	23
3.3.3	Extended Alignment Gadget for LCS	27
4	Tight Bounds Assuming (Combinatorial) k-Clique	32
4.1	NFA Acceptance	33
4.2	Context-Free Grammar Parsing	37
4.3	RNA Folding	42
5	Disjointness, Hamming Distance, and Subsequence	47
5.1	Algorithms	48
5.2	Lower Bounds	50
6	Conclusion	55

1 Introduction

Computer Science is often called the science of processing digital data. A central goal of theoretical CS is to understand the time complexity of the tasks we want to perform on data. *Data compression* has been one of the most important notions in CS and Information Theory for decades, and it is increasingly relevant in our current age of “Big Data” where it is hard to think of reasons why *not* to compress our data: smaller data can be stored more efficiently, transmitting it takes less resources such as energy and bandwidth, and perhaps it can even be processed faster. Since nowadays and for years to come nearly all of our data comes in compressed form, a central question becomes:

What is the time complexity of analyzing compressed data?

Say we have a piece of data of size N given in a compressed form of size n . For a problem with time complexity $T(\cdot)$, the naïve strategy of “decompress and solve” takes $\Theta(T(N))$ time, while the “gold standard” is $O(T(n))$ time: we want to solve the problem on the compression as efficiently as if the original data was small. To provide meaningful statements we need to decide on three things: What type of *data* is it? What *problem* do we want to solve? Which *compression scheme* is being used?

For the first two questions, the focus of this paper will be on the most basic setting. We consider data that comes as strings, i.e. sequences of symbols such as text, computer code, genomes, and so on. And we study natural and basic questions one could ask about strings such as Pattern Matching, Language Membership, Longest Common Subsequence, Parsing, and Disjointness.

For the third question, we restrict our attention to *lossless* compression and, even then, there are multiple natural settings that we do not find to be the most relevant. We could consider Kolmogorov complexity, giving us the best possible compression of our data: assume that a string T is given by a short bitstring $K(T)$ which is a pair of Turing machine M and input x such that running M on x outputs T , i.e. $K(T) = \langle M, x \rangle$ such that $M(x) = T$. The issue with Kolmogorov-compressions is that none of our data comes in this form, for two good reasons: First, it is computationally intractable to compute $K(T)$ given T , not even approximately. And second, analyzing arbitrary Turing machines without just running them is an infamously hopeless task. Thus, while studying the time complexity of analyzing Kolmogorov-compressed strings is natural, it might not be the most relevant for computer science applications. Another option is to consider the mathematically simplest forms of compression such as *Run-Length Encoding* (RLE): we compress x consecutive letters σ into σ^x , so the compression has the form $0^{x_1}1^{x_2}0^{x_3} \dots 1^{x_\ell}$, and we only need $n = O(\ell \cdot \log N)$ bits to describe the potentially exponentially longer string of length N . This compression is at the other extreme of the spectrum: it is trivial to compute and easy to analyze, but it is far less “compressing” than popular schemes like Lempel-Ziv-compressions.

Instead, we consider what has proven to be one of the most influential kinds of compression for Algorithm design, namely *Grammar-Compressions*, a notion that has all the right properties. First, it is mathematically elegant and quite fun to reason about for theoreticians (as evidenced by the many pages of our paper). Second, it is equivalent [64] up to low order terms (moderate constants and log factors) to popular schemes like the Lempel-Ziv-family (LZ77, LZ78, LZW, etc.) [48, 81, 75], Byte-Pair Encoding [68], dictionary methods, and others [57, 50]. These compressions are used in ubiquitous applications such as the built-in Unix utility `compress`, `zip`, GIF, PNG, and even in PDF. Third, it is generic and likely to capture compression schemes that will be engineered in the future (after all, there is a whole industry on the topic and the quest might never

be over). Fourth, we can compute the optimal such compression (up to log factors) in linear time [64, 23, 44]. And last but not least, ingenious algorithmic techniques have shown that it is possible to computationally *analyze* grammar-compressed data, beating the “decompress and solve” bound for many important problems.

A grammar compression of a string X is simply a context-free grammar, whose language is exactly $\{X\}$, that is, the only string the grammar can produce is X . For the purposes of this paper, it is enough to focus on a restricted form of grammars, known as *Straight Line Programs* (SLP). An SLP is defined over some alphabet Σ , say $\{0, 1\}$, and it is a set of replacement rules (or productions) of a very simple form: a rule is either a symbol in Σ or it is the concatenation of two previous rules (under some fixed ordering of the rules). The last replacement rule is the sequence defined by the SLP. For example, we can compress the sequence 01011 with the rules $S_1 \rightarrow 0$; $S_2 \rightarrow 1$; $S_3 \rightarrow S_1 S_2$; $S_4 \rightarrow S_3 S_3$; $S_5 \rightarrow S_4 S_2$ and S_5 corresponds to the sequence 01011. For some strings this can give an exponential compression. A more formal definition and a figure are given in Section 2.

To learn more about the remarkable success of grammar-compressions, we refer the reader to the surveys [77, 47, 34, 67, 36, 63, 65, 53, 66]. As a side remark, one of the exciting developments in this context was the surprising observation that a “compress and solve” strategy could actually lead to theoretically new algorithms for some problems, e.g. [60, 45].

Thus, we focus on what we find the most important interpretation of the central question above:

What is the time complexity of basic problems on grammar-compressed strings?

1.1 Previous Work

As a motivating example, consider the Longest Common Subsequence (LCS) problem. Given two uncompressed strings of length N we can find the length of the longest common (not necessarily contiguous) subsequence in $O(N^2)$ time using dynamic programming, and there are almost-matching $N^{2-o(1)}$ conditional lower bounds [2, 17, 3]. Throughout the paper we mostly ignore log factors, and so we think of LCS as a problem with $\tilde{\Theta}(N^2)$ time complexity (on uncompressed data). Now, assume our sequences are given in compressed form of size n . A natural setting to keep in mind is where $n \approx N^{1/2}$. How much time do we need to solve LCS on these compressed strings? The naïve upper bound gives $O(N^2)$ and the gold standard is $O(n^2) \approx O(N)$, so which is it?

Besides being a very basic question, LCS and the closely related Edit Distance are a popular theoretical modeling of sequence alignment problems that are of great importance in Bioinformatics¹. Thus, this is a relatively faithful modeling of the question whether “compress-and-solve” can speed up genome analysis tasks, a question which has received extensive attention throughout the years [39, 57, 50, 38, 36].

A long line of work [18, 54, 7, 8, 27, 69, 70, 40] has shown that we can do *much* better than $O(N^2)$. The current best algorithm has the curious runtime $O(nN\sqrt{\log N/n})$ [35] which is tantalizingly close to a conjectured bound of $O(nN)$ from the seminal paper of Lifshits [49]. In our candidate setting of $n \approx N^{1/2}$, this is $\tilde{O}(N^{1.5})$. This is major speedup over the $\Omega(N^2)$ decompress-and-solve bound, but is still far away from the gold standard of $O(n^2)$ which in this case would be $O(N)$. Can we do better? For example, an $O(n^2 \cdot N^{0.1})$ bound could lead to major real-world improvements.

¹The *heuristic* algorithm BLAST for a generalized version of the problem has received sixty-thousand citations.

While there is a huge literature on the topic, both from the Algorithms community and from applied areas, in addition to the potential for real-world impact, studying these questions has not become a mainstream topic in the top algorithms conferences. In one of the only STOC/FOCS papers on the topic, Charikar et al. [23] write “*In short, the smallest grammar problem has been considered by many authors in many disciplines for many reasons over a span of decades. Given this level of interest, it is remarkable that the problem has not attracted greater attention in the general algorithms community.*”

We believe that one key reason for this is the lack of a relevant *complexity theory* and tools for proving *lower bounds*, leaving a confusing state of the art in which it is hard to distinguish algorithms providing fundamental new insights from *ad hoc* solutions. Most importantly, previous work has not given us the tools to know, when we encounter a data analysis problem in the real-world, what kind of upper bound we should expect. Instead, researchers have been proving P vs. NP-hard results, classifying problems into ones solvable in $\text{poly}(n, \log N)$ time and ones that probably require time $N^{\Omega(1)}$. In fact, even LCS is NP-hard [49]. This means that even if we have a compression of very small size $n = O(\log N)$ then we cannot solve LCS in $\text{poly}(n)$ time, unless $P = NP$. Dozens of such negative results have been proven (see [53]), and it has long been clear that almost any task of interest is “NP-hard”, including the basic $\text{poly}(N)$ time solvable problems we discuss in this paper. However, this is hardly relevant to the questions we ask in this paper since it does not address the possibility of highly desirable bounds such as $n^2 \cdot N^{0.1}$. What we would really like to know is whether the bound should be $\text{poly}(n) \cdot N^\epsilon$, or $\text{poly}(n) \cdot N$, or even higher: could it be that decompress-and-solve is impossible to beat for some problems?

1.2 Our Work

In this work, we introduce a framework for showing lower bounds on the time complexity of problems on grammar-compressed strings. Our lower bounds are based on popular conjectures from *Hardness in P* and Fine-Grained Complexity. This is perhaps surprising since the problems we consider are technically NP-hard. Our new complexity theoretic study of this field leads to three exciting developments: First, we resolve the exact time complexity up to $N^{o(1)}$ factors of some of the most classical problems such as LCS *on compressed data*. Second, we discover problems that *cannot be solved faster than the decompress-and-solve bound* by any N^ϵ factor. Third, we *fail* at proving tight lower bounds for some classical problems, which hints to us that known algorithms might be suboptimal. Indeed, in this paper we also find *new algorithms* for fundamental problems. We hope that our work will inspire increased interest in this important topic.

Longest Common Subsequence Our first result is a resolution of the time complexity of LCS on compressed data, up to $N^{o(1)}$ factors, under the Strong Exponential Time Hypothesis² (SETH). We complement the $O(nN\sqrt{\log N/n})$ upper bound of Gawrychowski [35] with an $(nN)^{1-o(1)}$ lower bound. Thus, in the natural setting $n \approx N^{1/2}$ from above, we should indeed be content with the $\tilde{O}(N^{1.5})$ upper bound since we will not be able to get much closer to the gold standard, unless SETH fails. Assuming SETH, our result confirms the conjecture of Lifshits, up to $N^{o(1)}$ factors. See Theorem 3.12 in Section 3.3 for the formal statement.

One way to view this result is as an *Instance Optimality* result for LCS. The exact complexity

²SETH is the pessimistic version of $P \neq NP$, stating that we cannot solve k -SAT in $O((2 - \epsilon)^n)$ time, for some $\epsilon > 0$ independent of and for all constant k [42, 19].

of LCS on two strings is precisely proportional to the product of the decompressed size N and the instance-inherent measure n of how compressible they are.

RNA Folding and CFG Parsing Next, we turn our attention to two other fundamental problems: Context-Free Grammar Recognition (aka Parsing) and RNA Folding. Parsing is the core computer science problem in which we want to decide whether a given string (e.g. computer code) can be derived from a given grammar (e.g. the grammar of a programming language). Having the ability to efficiently parse a *compressed* file is certainly desirable. In RNA Folding we are given a string over some alphabet (e.g. $\{A, C, G, T\}$) with a fixed pairing between its symbols (e.g. $A - T$ match and $C - G$ match), and the goal is to compute the maximum number of non-crossing arcs between matching letters that one can draw above the string (which corresponds to the minimum energy folding in two dimensions). RNA Folding is one of the most central problems in bioinformatics, and as we have discussed above, the ability to analyze compressed data is important in this field. How fast can we solve these problems?

Given an uncompressed string of size N , classical dynamic programming algorithms, such as the CYK parser [25, 80, 46], solve RNA Folding in $O(N^3)$ time and Parsing in $O(N^3 \cdot g)$ time if the grammar has size g . Wikipedia lists twenty-four parsing algorithms designed throughout the years, all of which take cubic time in the worst case. A theoretical breakthrough of Leslie Valiant [72] in 1975 showed that there are truly sub-cubic $O(gN^\omega)$ parsing algorithms, where $\omega < 2.38$ is the fast matrix multiplication (FMM) exponent. However, Valiant’s algorithm has not been used in practice due the inefficiency of FMM algorithms, and obtaining a *combinatorial*³ sub-cubic time algorithm would be of major interest. Alas, it was recently proved [1] that any improvement over these bounds implies breakthrough k -Clique algorithms: either finding such a combinatorial subcubic algorithm *or* getting any $O(N^{\omega-\varepsilon})$ time algorithm, for any $\varepsilon > 0$, would refute the k -Clique Conjecture⁴. The situation for RNA is even more interesting since Valiant’s sub-cubic algorithm does not generalize to this case. Under the k -Clique conjecture, the same lower bounds still apply [1, 22], implying that any improvement will have to use FMM. Indeed, an $O(N^{2.82})$ algorithm using FMM was recently achieved [15].

Cubic time is a real bottleneck when analyzing large genomic data. One would hope that if we are able to compress the data down to size n we could solve problems like RNA Folding and Parsing in time that is much faster than the N^3 lower bounds (to simplify the discussion we focus on combinatorial algorithms), such as $n^3 \cdot N^{o(1)}$ or at least $n^{1.5} N^{1.5}$, in certain analogy the LCS case. No such algorithms were found to date, and we provide an explanation: Decompress-and-solve *cannot be beaten* for Parsing and (essentially) for RNA Folding, under the k -Clique Conjecture. For both problems we prove a conditional lower bound of $N^{\omega-o(1)}$ for any kind of algorithm, and $N^{3-o(1)}$ for combinatorial algorithms, even restricted to $n = O(N^\varepsilon)$ for any $\varepsilon > 0$. See Theorem 4.4 in Section 4.2 for CFG Parsing and Theorem 4.10 in Section 4.3 for RNA Folding.

Approximate Pattern Matching We continue our quest for quantifying the possible improvements over decompress-and-solve for basic problems. Consider the following compressed versions

³For the purposes of this paper, “combinatorial” should be interpreted as any *practically efficient* algorithm that does not suffer from the issues of FMM such as large constants and inefficient memory usage.

⁴Given a graph on n nodes, the k -Clique conjecture [1] is in fact two independent conjectures: The first one states that we cannot solve k -clique in $O(n^{(1-\varepsilon) \cdot \omega k/3})$, for any $\varepsilon > 0$. The second one states that we cannot solve k -Clique combinatorially in $O(n^{(1-\varepsilon)k})$ time, for any $\varepsilon > 0$.

of important primitives in text analysis known as *Approximate Pattern Matching* problems. In all these problems we assume that we are given a compressed text T of size n (and decompressed size N), and a compressed pattern P of size m (and decompressed size M), both over some constant size alphabet.

- **Pattern Matching with Wildcards:** In this problem, the strings contain wildcard symbols that can be replaced by any letter, and our goal is to decide if P appears in T .
- **Substring Hamming Distance:** Compute the smallest Hamming distance of any substring of T to P .

And a problem that generalizes both is:

- **Generalized Pattern Matching:** Given some cost function on pairs of alphabet symbols, find the length- M substring T' of T minimizing the total cost of all pairs $(T'[i], P[i])$.

The above problems have been extensively studied both in the uncompressed (see [24]) and in the compressed [49, 13, 32] settings. All three problems can be solved in time $O(\min\{N \log N, nM\})$ (see Section 3.2). Note that this bound beats the decompress-and-solve bound when the pattern is small, but can we avoid decompressing the pattern? We show a completely tight SETH-based lower bound of $\min\{N, nM\}^{1-o(1)}$ for all three problems, even for constant size alphabets and in all settings where the parameters are polynomially related. See Theorems 3.9 and 3.10 in Section 3.2.

Language Membership Consider the compressed version of the most basic language membership problems. Assume we are given a compressed string T (again, from size N to n).

- **DFA Acceptance:** Given T and a DFA F with q states, decide whether F accepts T .
- **NFA Acceptance:** Given T and a NFA F with q states, decide whether F accepts T .

Classic algorithms solve the DFA Acceptance problem in time $O(\min\{nq, N + q\})$ [61, 41], and we prove a matching SETH-based lower bound of $\min\{nq, N + q\}^{1-o(1)}$. See Theorem 3.2 in Section 3.1.

For the NFA problem, the classic algorithms give $O(\min\{nq^\omega, Nq^2\})$ [55, 61, 41]. For combinatorial algorithms, we prove a matching lower bound of $\min\{nq^3, Nq^2\}^{1-o(1)}$, under the (combinatorial) k -Clique conjecture. See Theorem 4.2 in Section 4.1. Our lower bounds hold for constant size alphabets, and in all settings of n, N, q , even restricted to instances with $N = \Theta(n^{\alpha_N})$ and $q = \Theta(n^{\alpha_q})$ for any $\alpha_N > 1$ and $\alpha_q > 0$.

Disjointness, Hamming Distance, and Subsequence Could it be that for other, even more basic problems the decompress-and-solve bound cannot be beaten? One candidate might be Disjointness, the canonical hard problem in Communication Complexity.

- **Disjointness:** Given two equal-length bit-strings, is there a coordinate in which both are 1?

The following two natural problems are at least as hard as Disjointness (see Section 5).

- **Hamming Distance:** Compute the Hamming Distance of two strings.
- **Subsequence:** Decide if a pattern of length M is a subsequence of a text of length N .

Note that all these problems can be solved trivially in $O(N)$ time if our strings are uncompressed. Could it be that we cannot solve them without decompressing our data? We are not aware of any known algorithms solving any of these problems in $O(N^{1-\varepsilon})$ time, for any $\varepsilon > 0$, even when our strings are compressed into size $n = O(N^\alpha)$ for some small constant $\alpha > 0$. The only exceptions are the known $\tilde{O}(M)$ time algorithms [28, 20, 69, 79, 71, 12] for the Subsequence problem, which beat the decompress-and-solve bound when the pattern is significantly smaller than the text. However, in the case $M = \Theta(N)$ no improvements seem to be known.

In Section 5 we present our attempts at proving a matching lower bound. We prove the following: $N^{1-o(1)}$ for Subsequence in the setting $N = \Theta(M) = \Theta(n^2) = \Theta(m^2)$ and $|\Sigma| = O(N^\varepsilon)$ (Theorem 5.9). $N^{1/4-o(1)}$ for Disjointness (and thus also for the other two problems) in the setting $N = M$ and $n, m = O(N^\varepsilon)$ for any $\varepsilon > 0$, and constant alphabet size, assuming the k -SUM conjecture (Theorem 5.10). Similarly: $N^{1/3-o(1)}$ for Disjointness under Strong k -SUM conjecture (Theorem 5.11).

Motivated by our inability to prove tight lower bounds for these basic problems, despite seemingly having the right framework, we have turned our attention to upper bounds. In Section 5 we obtain the *first* improvement over the decompress-and-solve bound for Disjointness, Hamming Distance, and Subsequence. In particular, we obtain the first improvement over the decompress-and-solve bound for Disjointness, Hamming Distance, and Subsequence. Our algorithms solve all these problems in $O(n^{1.410} \cdot N^{0.593})$ time. As a side result, we also design a very simple algorithm for the Subsequence problem with $O((n|\Sigma| + M) \log N)$ runtime (Theorem 5.4), which is comparable to the known but more involved algorithms [12].

One of the biggest benefits of having complexity theoretic results is that algorithm designers know what to focus on. We believe that these upper bounds can be improved further and suggest it as an interesting open question: *What is the time complexity of computing Disjointness on two grammar-compressed strings?*

1.3 Technical Overview

From a technical perspective, our paper is most related to the conditional lower bounds for sequence similarity measures on strings and curves that have been shown in recent years, specifically, the SETH-based lower bounds for edit distance [10], longest common subsequence [2, 17], Fréchet distance [14], and others [4, 11, 16, 62].

These results all proceed as follows. Let ϕ be a given k -SAT instance on \tilde{n} variables and clauses $C_1, \dots, C_{\tilde{m}}$. We can assume that $\tilde{m} = O(\tilde{n})$ by the Sparsification Lemma [43]. Split the \tilde{n} variables into two halves X_1 and X_2 of size $\tilde{n}/2$. Enumerate all assignments $\alpha_1, \dots, \alpha_{2^{\tilde{n}/2}}$ of the variables in X_1 . For any assignment α_i and any clause C_ℓ , denote by $\text{sat}(\alpha_i, C_\ell)$ whether α_i satisfies C_ℓ , i.e., whether some variable in X_1 appears in C_ℓ (negated or unnegated) and is set by α_i so that C_ℓ is satisfied. Similarly, consider the assignments $\beta_1, \dots, \beta_{2^{\tilde{n}/2}}$ of X_2 . By construction, we can solve the k -SAT instance ϕ by testing whether there are α_i, β_j such that $\text{sat}(\alpha_i, C_\ell) \vee \text{sat}(\beta_j, C_\ell)$ holds for all $\ell \in [\tilde{m}]$. Making use of this fact, all previous conditional lower bounds for sequence similarity measures essentially construct the following natural sequence:

$$\begin{aligned} W &= \text{sat}(\alpha_1, C_1) \dots \text{sat}(\alpha_1, C_{\tilde{m}}) \dots \text{sat}(\alpha_{2^{\tilde{n}/2}}, C_1) \dots \text{sat}(\alpha_{2^{\tilde{n}/2}}, C_{\tilde{m}}) \\ &= \bigcirc_{i \in [2^{\tilde{n}/2}]} \bigcirc_{\ell \in [\tilde{m}]} \text{sat}(\alpha_i, C_\ell). \end{aligned}$$

One typical variation of this string is to replace the bits $\{0, 1\}$, indicating whether $\text{sat}(\alpha_i, C_\ell)$ holds, by two short strings $\{B(0), B(1)\}$. Other typical variations are to add appropriate padding strings around the substrings $\bigcirc_{\ell \in [\tilde{m}]} \text{sat}(\alpha_i, C_\ell)$ or around the whole sequence W . These paddings typically only depend on \tilde{n} and \tilde{m} . Constructing a second sequence W' with α_i replaced by β_i , one can then try to emulate the search for the half-assignments α_i, β_j by a similarity measure on W, W' . All previous reductions follow this recipe, and thus construct a sequence like W .

Is W compressible? For our purposes we need to construct compressible strings. Considering the entropy, the string W is very well compressible, since it only depends on the $\tilde{O}(\tilde{n})$ input bits of the sparse k -SAT instance ϕ . This entropy $\tilde{O}(\tilde{n})$ is extremely small compared to the length $O(\tilde{n}2^{\tilde{n}/2})$ of W . However, considering grammar-compression, the sequence W is a bad representation, since W is not generated by any SLP of size $o(2^{\tilde{n}/2}/\tilde{n})$ in general! To see this, first observe that all substrings $\bigcirc_{\ell \in [\tilde{m}]} \text{sat}(\alpha_i, C_\ell)$ of W can potentially be different, meaning that W can have $2^{\tilde{n}/2}$ different substrings of length \tilde{m} . This happens e.g. if for each variable $x_i \in X$ there is a clause C_i consisting only of x_i (which makes the k -SAT instance trivial, but shows that W may have many different substrings in general). Second, observe that for any SLP \mathcal{T} consisting of n non-terminals $S_1 \dots S_n$ and for any length $L \geq 1$ the generated string $\text{eval}(\mathcal{T})$ has at most $n \cdot L$ different substrings of length L . Indeed, a rule $S_i \rightarrow S_\ell S_r$ can only create a new substring, that is not already contained in $\text{eval}(S_\ell)$ or $\text{eval}(S_r)$, if this substring overlaps the boundary between $\text{eval}(S_\ell)$ and $\text{eval}(S_r)$ in $\text{eval}(S_i)$. Hence, the rule $S_i \rightarrow S_\ell S_r$ can contribute at most L new substrings of length L , amounting to at most nL different substrings overall. Combining these two facts, with $L = \tilde{m} = O(\tilde{n})$, we see that W in general has no SLP of size $o(2^{\tilde{n}/2}/\tilde{n})$.

Hence, the standard approach to conditional lower bounds for sequence similarity measures fails in the compressed setting, and it might seem like (SETH-based) conditional lower bounds are not applicable here.

A compressible sequence T On the contrary, we show that by simply inverting the ordering we obtain a very well compressible string:

$$\begin{aligned} T &= \text{sat}(\alpha_1, C_1) \dots \text{sat}(\alpha_{2^{\tilde{n}/2}}, C_1) \dots \text{sat}(\alpha_1, C_{\tilde{m}}) \dots \text{sat}(\alpha_{2^{\tilde{n}/2}}, C_{\tilde{m}}) \\ &= \bigcirc_{\ell \in [\tilde{m}]} \bigcirc_{i \in [2^{\tilde{n}/2}]} \text{sat}(\alpha_i, C_\ell). \end{aligned}$$

The difference between W and T might seem negligible, but it greatly changes the game of emulating k -SAT by a sequence similarity measure: In W we are looking for a local structure (a small substring) that “fits together” with a local structure in a different string W' . In T we have to ensure the choice of a consistent offset $\Delta \in [n]$ and “read” the symbols $T[\Delta], T[\Delta + 2^{\tilde{n}/2}], \dots, T[\Delta + (\tilde{m} - 1)2^{\tilde{n}/2}]$, which seems much more complicated.

T is compressible to an SLP \mathcal{T} of size $O(\tilde{n}^2)$, which is much smaller than the $\Omega(2^{\tilde{n}/2}/\tilde{n})$ bound for W . Indeed, consider a substring $\bigcirc_{i \in [2^{\tilde{n}/2}]} \text{sat}(\alpha_i, C_\ell)$. We may assume that no variable appears more than once in C_ℓ . Consider the following SLP rules, for $1 \leq i \leq \tilde{n}/2$,

$$\begin{aligned} A_0 &\rightarrow 1, \\ A_i &\rightarrow A_{i-1}A_{i-1}, \\ S_0 &\rightarrow 0, \end{aligned} \quad S_i \rightarrow \begin{cases} S_{i-1}A_{i-1} & \text{if } x_i \text{ appears in } C_\ell \\ A_{i-1}S_{i-1} & \text{if } \neg x_i \text{ appears in } C_\ell \\ S_{i-1}S_{i-1} & \text{otherwise} \end{cases}$$

We clearly have $\text{eval}(A_i) = 1^{2^i}$. Moreover, if $-x_i$ appears in C_ℓ , then for $x_i = 0$, no matter what we choose for x_1, \dots, x_{i-1} , we have $\text{sat}(\alpha_j, C_\ell) = 1$, and thus we may write A_{i-1} . For $x_i = 1$ we note that the value $\text{sat}(\alpha_j, C_\ell)$ only depends on the remaining variables x_1, \dots, x_{i-1} , and thus we may write S_{i-1} . Along these lines, one can check that $\text{eval}(S_{\tilde{n}/2}) = \bigcirc_{i \in [2^{\tilde{n}/2}]} \text{sat}(\alpha_i, C_\ell)$. Creating such an SLP for each $\ell \in [\tilde{m}]$ and constructing their concatenation, we obtain an SLP of size $O(\tilde{m}\tilde{n}) = O(\tilde{n}^2)$ generating T .

Example Lower Bound: Pattern Matching with Wildcards In the remainder of this section, we present an easy example for a conditional lower bound on compressed strings, namely for the problem Pattern Matching with Wildcards. Here we consider an alphabet Σ and we say that symbols $\sigma, \sigma' \in \Sigma \cup \{*\}$ *match* if $\sigma = *$ or $\sigma' = *$ or $\sigma = \sigma'$. We say that two equal-length strings X, Y (over alphabet $\Sigma \cup \{*\}$) *match* if $X[i]$ and $Y[i]$ match for all i . Given a text T of length N and a pattern P of length $M \leq N$, the task is to decide whether P matches some length- M substring of T .

Let ϕ be a k -SAT instance as above, but this time let $\alpha_1, \dots, \alpha_{2^{\tilde{n}}}$ be *all* the assignments of the \tilde{n} variables in ϕ . We define the text T and pattern P by

$$T = \bigcirc_{\ell \in [\tilde{m}]} \bigcirc_{i \in [2^{\tilde{n}}]} \text{sat}(\alpha_i, C_\ell) \quad P = 1(*^{2^{\tilde{n}}-1}1)^{\tilde{m}-1}.$$

Note that P matches some substring of T if and only if there is an offset $\Delta \in [2^{\tilde{n}}]$ such that $T[\Delta] = T[\Delta + 2^{\tilde{n}}] = \dots = T[\Delta + (\tilde{m} - 1)2^{\tilde{n}}] = 1$, which happens if and only if α_Δ is a satisfying assignment of ϕ . Hence, we constructed an equivalent instance of Pattern Matching with Wildcards.

Analogously to above, one can show that T is generated by an SLP \mathcal{T} of size $n = O(\tilde{n}^2)$ that can be computed in time $O(\tilde{n}^2)$. Similarly, it is easy to see that P is generated by an SLP \mathcal{P} of size $O(\tilde{n})$ that can be computed in time $O(\tilde{n})$. Hence, the reduction runs in time $O(\tilde{n}^2)$. We stress that we define strings T, P of exponential length in \tilde{n} , but in the reduction we never explicitly write down any such string, but we simply construct compressed representations. Since the resulting strings have length $O(2^{\tilde{n}}\tilde{n})$, any $O(N^{1-\varepsilon})$ time algorithm for Pattern Matching with Wildcards would imply an algorithm for k -SAT in time $O(2^{(1-\varepsilon)\tilde{n}}\text{poly}(\tilde{n}))$, contradicting the Strong Exponential Time Hypothesis (SETH). Note that this conditional lower bound of $N^{1-o(1)}$ holds even for strings compressible to size $\text{polylog}(N)$.

In Section 3.2 we analyze Pattern Matching with Wildcards in more detail and show that the optimal running time, conditional on SETH, is $\min\{N, nM\}^{1\pm o(1)}$, and this holds for all settings of the text length N , the compressed text size n , the pattern length M , and the compressed pattern size m .

In Pattern Matching with Wildcards, we got a consistent choice of an offset Δ for free. It is much more complicated to achieve this for other problems such as Longest Common Subsequence, CFG Parsing, or RNA Folding. This overview summarized the main technical contributions of this paper, but left out many problem-specific tricks that can be found in the subsequent proofs, and that we think will find more applications for analyzing problems on compressed strings.

2 Preliminaries

Here we give general preliminaries on strings, straight-line programs, and hardness assumptions. Problem definitions and additional problem specific preliminaries will be given in the corresponding

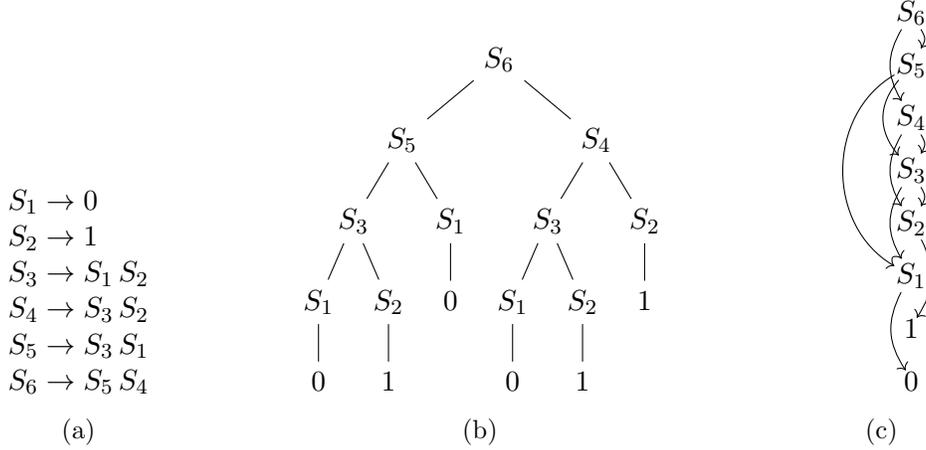


Figure 1: (a) An SLP generating the sequence 010011. (b) The corresponding parse tree. (c) The acyclic graph corresponding to the SLP.

sections. For a positive integer n we let $[n] = \{1, \dots, n\}$, while for a proposition A we let $[A]$ be 1 if A is true and 0 otherwise.

Strings Let Σ be a finite alphabet. In most parts of this paper we assume that $|\Sigma| = O(1)$, but in exceptional cases we allow the alphabet to grow with the input size. For a string T over alphabet Σ , we write $|T|$ for its length, $T[i]$ for its i -th symbol, and $T[i..j]$ for the substring from position i to position j . For two strings T, T' we write $T \circ T'$, or simply TT' , for their concatenation. For $k \geq 1$ we let $T^k := \bigcirc_{i=1}^k T$.

Straight-Line Programs (SLPs) An SLP \mathcal{T} is a set of non-terminals S_1, \dots, S_n , each equipped with a rule of the form (1) $S_i \rightarrow \sigma$ for some $\sigma \in \Sigma$ or (2) $S_i \rightarrow S_{\ell(i)}, S_{r(i)}$ with $\ell(i), r(i) < i$. The string T generated by SLP \mathcal{T} is recursively defined as follows. For a rule $S_i \rightarrow \sigma$ we let $\text{eval}(S_i) := \sigma$, and for a rule $S_i \rightarrow S_{\ell(i)}, S_{r(i)}$ we let $\text{eval}(S_i) := \text{eval}(S_{\ell(i)}) \circ \text{eval}(S_{r(i)})$. Then $T = \text{eval}(\mathcal{T}) := \text{eval}(S_n)$ is the string generated by SLP \mathcal{T} . Note that an SLP is a context-free grammar describing a unique string; so \mathcal{T} is a *grammar-compressed* representation of T . We call $|\mathcal{T}| = n$ the *size* of \mathcal{T} . See Figure 1 for the depiction of an SLP; in particular note the difference between the *directed acyclic graph* that is the compressed representation \mathcal{T} and the *parse tree* that we obtain by decompressing \mathcal{T} to a tree whose leaves spell the decompressed text T .

For an SLP \mathcal{T} with non-terminals S_1, \dots, S_n , we recursively define the depth $\text{depth}(S_i)$ as follows. For a rule $S_i \rightarrow \sigma$ we set $\text{depth}(S_i) := 0$. For a rule $S_i \rightarrow S_{\ell(i)}, S_{r(i)}$ we set $\text{depth}(S_i) = \max\{\text{depth}(S_{\ell(i)}), \text{depth}(S_{r(i)})\} + 1$. The depth of \mathcal{T} is defined as $\text{depth}(S_n)$. The SLP \mathcal{T} is called an *AVL-grammar* [64] if it is balanced: for any rule $S_i \rightarrow S_{\ell(i)}, S_{r(i)}$ in the SLP we have $|\text{depth}(S_{\ell(i)}) - \text{depth}(S_{r(i)})| \leq 1$. This implies that the depth of \mathcal{T} is $O(\log N)$, where $N = |\text{eval}(\mathcal{T})|$.

Theorem 2.1 ([64]). *Given a text T of length N by an SLP \mathcal{T} of size n , in $O(n \log N)$ time we can construct an AVL-grammar \mathcal{T}' for T with size $O(n \log N)$ and depth $O(\log N)$.*

Observation 2.2. *For any string T and $k \geq 1$, there is an SLP of size $O(|T| + \log k)$ generating the string T^k .*

In all problems considered in this paper, the input contains a text T given by a grammar-compressed representation \mathcal{T} , such that $T = \text{eval}(\mathcal{T})$. We always denote by $N = |T|$ the length of the text and by $n = |\mathcal{T}|$ the size of its representation. Sometimes we are additionally given a pattern P by a grammar-compressed representation \mathcal{P} , and we denote the pattern length by $M = |P|$ and its representation size by $m = |\mathcal{P}|$.

2.1 Hardness Assumptions

SETH and OV The Strong Exponential Time Hypothesis (SETH) was introduced by Impagliazzo, Paturi, and Zane [43] and asserts that the central NP-hard satisfiability problem has no algorithms that are much faster than exhaustive search.

Conjecture 2.3 (SETH). *There is no $\varepsilon > 0$ such that for all $k \geq 3$, k -SAT on n variables can be solved in time $O(2^{(1-\varepsilon)n})$.*

Effectively all known SETH-based lower bounds for polynomial-time problems use reductions via the *Orthogonal Vectors problem* (OV): Given sets $\mathcal{A}, \mathcal{B} \subseteq \{0, 1\}^d$ of size $|\mathcal{A}| = A, |\mathcal{B}| = B$, determine whether there exist vectors $a \in \mathcal{A}, b \in \mathcal{B}$ with $\sum_{i=1}^d a[i] \cdot b[i] = 0$. Simple algorithms solve OV in time $O(2^d(A+B))$ and $O(dAB)$. For $A = B$ and $d = c(A) \log A$ the fastest known algorithm runs in time $A^{2-1/O(\log c(A))}$ [5], which is only slightly subquadratic for $d \gg \log A$. This has led to the following conjecture, which follows from SETH [76].

Conjecture 2.4 (OV). *For any $\varepsilon > 0$ and $\beta > 0$, on instances with $B = \Theta(A^\beta)$ OV has no $O(A^{1+\beta-\varepsilon} \text{poly}(d))$ time algorithm.*

It is known that if this conjecture holds for some $\beta > 0$ then it holds for all $\beta > 0$, see e.g. [17].

More generally, for $k \geq 2$ we say that a tuple (a_1, \dots, a_k) with $a_i \in \{0, 1\}^d$ is *orthogonal* if for all $\ell \in [d]$ there exists an $i \in [k]$ such that $a_i[\ell] = 0$. In the k -OV problem we are given a set $\mathcal{A} \subseteq \{0, 1\}^d$ of size A and want to determine whether there is an orthogonal tuple (a_1, \dots, a_k) with $a_i \in \mathcal{A}$. The fastest known algorithm for k -OV is to run an easy reduction to OV and then solve OV. The following conjecture follows from SETH.

Conjecture 2.5 (k -OV). *For any $\varepsilon > 0$ and $k \geq 2$, k -OV is not in time $O(A^{k-\varepsilon} \text{poly}(d))$.*

k -Clique The fundamental k -Clique problem asks whether a given (undirected, unweighted) graph $G = (V, E)$ contains k nodes that are pairwise adjacent. k -Clique is among the most well-studied problems in theoretical computer science, and it is the canonical intractable (W[1]-complete) problem in parameterized complexity. With slight abuse of notation, we will denote the number of vertices and edges of G by V and E , respectively. The naive algorithm for k -Clique takes time $O(V^k)$. If k is divisible by 3, the fastest known algorithm runs in time $O(V^{\omega k/3})$, where $\omega < 2.373$ is the exponent of matrix multiplication [58]. See [31] for the case that k is not divisible by 3. To improve this bound is a longstanding open problem [78, 56]. Since fast matrix multiplication is considered impractical, researchers also studied *combinatorial* algorithms, that avoid fast matrix multiplication⁵. The fastest combinatorial algorithm runs in time $O(V^k / \log^k V)$ [73]. The following conjectures assert that these bounds are close to optimal, and have been used e.g. in [1, 16].

⁵Combinatorial algorithms are a notion without agreed upon definition; finding a formal definition is considered an open problem.

Conjecture 2.6 (*k*-Clique). *For any $\varepsilon > 0$ and $k \geq 3$, *k*-Clique has no $O(V^{(1-\varepsilon)\omega k/3})$ algorithm.*

Conjecture 2.7 (Combinatorial *k*-Clique). *For any $\varepsilon > 0$ and $k \geq 3$, *k*-Clique has no combinatorial $O(V^{(1-\varepsilon)k})$ algorithm.*

***k*-SUM** In the *k*-SUM problem, we are given integers $R, t \geq 0$ and a set $Z \subseteq \{0, 1, \dots, R\}$ of $|Z| = r$ integers, and the task is to decide whether there are *k* (not necessarily distinct) integers $z_1, \dots, z_k \in Z$ that sum to *t*, i.e., $z_1 + \dots + z_k = t$. This problem has well-known algorithms in time $O(r^{\lceil k/2 \rceil})$ and $O(r + R \log R)$, and it is conjectured that no much faster algorithms exist. The following conjectures, which generalize the more popular 3-SUM conjecture [33, 59] and Strong 3-SUM conjecture [6], remain believable despite recent algorithmic progress [9, 21, 37, 74].

Conjecture 2.8 (*k*-SUM). *For any $k \geq 3$ and $R = r^k$, the *k*-SUM problem is not in time $O(r^{\lceil k/2 \rceil - \varepsilon})$ for any $\varepsilon > 0$.*

Conjecture 2.9 (Strong *k*-SUM). *For any $k \geq 3$ and $R = r^{\lceil k/2 \rceil}$, the *k*-SUM problem is not in time $O(r^{\lceil k/2 \rceil - \varepsilon})$ for any $\varepsilon > 0$.*

3 Tight Bounds Assuming SETH

In this section we prove matching conditional lower bounds based on the Strong Exponential Time Hypothesis (SETH, see Conjecture 2.3) for the following problems:

- DFA Acceptance, i.e., deciding whether a given deterministic finite automaton accepts a given string,
- Substring Hamming Distance, i.e., determining the minimum Hamming distance that can be achieved by aligning a given pattern sequence with a substring of a given text sequence,
- Pattern Matching with Wildcards, i.e., deciding whether the given pattern sequence (containing wildcards that match any symbol) matches a substring of the given text,
- Longest Common Subsequence, i.e., computing the length of the longest common subsequence of two given strings.

See the respective subsections for precise problem definitions. In all our proofs, instead of using SETH directly, we use the more convenient OV conjecture (Conjecture 2.4) or *k*-OV conjecture (Conjecture 2.5), which are implied by SETH.

For DFA Acceptance, the compression used in our reduction from the given OV instance is extremely simple, in that we only rely on the fact that any repetition T^ℓ can be generated by an SLP of size $O(|T| + \log \ell)$ (Observation 2.2).

For Substring Hamming Distance, Pattern Matching with Wildcards and Longest Common Subsequence, however, our construction are more subtle. We crucially use the following idea: consider a *k*-OV instance \mathcal{A} on A vectors in d dimensions. There is a length- $O(dA^k)$ text T representing this instance so that (1) T is succinctly described by an SLP \mathcal{T} of size $O(dA)$ and (2) testing whether the *k*-OV instance has a solution corresponds to determining whether there is some $i = 1, \dots, A^k$ such that all bits $T[i], T[i + A^k], \dots, T[i + (d - 1)A^k]$ are equal to zero. Intuitively, $i \in \{1, \dots, A^k\}$ denotes the *i*-th *k*-tuple of vectors in \mathcal{A}^k , and $T[i] = 0$ holds if and only

if the k vectors in the i -th k -tuple are orthogonal in the 1st coordinate. In general, for $1 \leq \ell \leq d$, $T[i + (\ell - 1)A^k] = 0$ holds if and only if the k vectors are orthogonal in the ℓ -th coordinate. More formally, we set T to be

$$T = \bigcirc_{\ell=1}^d \bigcirc_{a_1 \in \mathcal{A}^{(1)}} \dots \bigcirc_{a_k \in \mathcal{A}^{(k)}} [a_1[\ell] = \dots = a_k[\ell] = 1],$$

where $[\cdot]$ is the Kronecker symbol, i.e., $[\text{true}] = 1$ and $[\text{false}] = 0$. For any ℓ , the sequence $\bigcirc_{a_1 \in \mathcal{A}^{(1)}} \dots \bigcirc_{a_k \in \mathcal{A}^{(k)}} [a_1[\ell] = \dots = a_k[\ell] = 1]$ is generated by an SLP of size $O(dA)$: if $a_1[\ell] = 0$, then for all a_2, \dots, a_k , the vectors will be orthogonal in this coordinate and we can write $0^{A^{k-1}}$, which is well compressible by Observation 2.2. Otherwise, if $a_1[\ell] = 1$, we recurse on a_2, \dots, a_k and the following A^{k-1} symbols do not depend on $a_1[\ell]$ anymore.

A modification of the above construction of T gives SETH hardness for Substring Hamming Distance and Pattern Matching with Wildcards. Showing hardness for Longest Common Subsequence requires more ideas. In particular, to be able to show tight hardness we extend the framework of [17].

We stress that if the sequence T would enumerate all k -tuples one after another (instead of iterating over the coordinates in the outer loop over ℓ), then it would not be compressible using SLPs, see Section 1.3. This makes our reductions quite different from all previously known hardness results where the sequences are concatenations of vector gadgets one after another.

Known Lower Bounds from Classic Complexity Theory We observe that the Substring Hamming Distance problem is a generalization of the Hamming Distance problem which asks to output the Hamming distance between a compressed text and a compressed pattern of equal length. The latter problem is known to be $\#P$ -complete and thus the Substring Hamming Distance problem is $\#P$ -hard (see the discussion at the beginning of Section 5.2). Similarly, Longest Common Subsequence is a generalization of the problem of deciding whether a given pattern is a subsequence of a given text. The latter problem is known to be PP -hard (see the aforementioned discussion) and this yields PP -hardness for Longest Common Subsequence.

The DFA Acceptance problem can be solved in polynomial time (see Section 3.1) and no conditional lower bounds were known for this problem. Finally, our reduction in Theorem 3.9 below shows that Pattern Matching with Wildcards is NP -hard.

3.1 DFA Acceptance

Recall that a finite-state automaton F over an alphabet Σ consists of a set of states Z of size q , a starting state $z_0 \in Z$, a set of accepting states $Z' \subseteq Z$, and a set of transitions $z \xrightarrow{\sigma} z'$ with $z, z' \in Z$ and $\sigma \in \Sigma$. We lift this notation to strings $T = T[1..N]$ by writing $z \xrightarrow{T} z'$ whenever there are states $z_1, \dots, z_{\ell-1}$ and transitions $z \xrightarrow{T[1]} z_1, z_1 \xrightarrow{T[2]} z_2, \dots, z_{\ell-1} \xrightarrow{T[N]} z'$. Furthermore, for a set $S \subseteq \Sigma$ we write $z \xrightarrow{S} z'$ whenever $z \xrightarrow{\sigma} z'$ for all $\sigma \in S$. The automaton F is deterministic if for any $z \in Z$ and $\sigma \in \Sigma$ there is at most one $z' \in Z$ with transition $z \xrightarrow{\sigma} z'$, and F is non-deterministic otherwise. The automaton F *accepts* a given string T if $z_0 \xrightarrow{T} z'$ holds for some accepting state z' .

Throughout this section, we assume the alphabet size to be constant. If F is a deterministic finite-state automaton (DFA), we may assume without loss of generality that for every state z and

symbol $\sigma \in \Sigma$, there always exists a (uniquely defined) state z' with $z \xrightarrow{\sigma} z'$.⁶ We fix the input description of F to a list of transitions of F as well as a list of accepting states. Observe that any DFA F on constant alphabet Σ has an input size of $O(q)$.

Consider the compressed variant of the acceptance problem of DFAs.

Problem 3.1 (DFA Acceptance). *Given a text T of length N by a grammar-compressed representation \mathcal{T} of size n as well as a DFA F with q states, decide whether T is accepted by F .*

The naive solution decompresses \mathcal{T} to obtain T and runs the obvious acceptance algorithm for DFAs, which takes time $O(|T| + q) = O(N + q)$. Exploiting the compressed setting, one can obtain an $O(nq)$ -time algorithm [61]: Recall that \mathcal{T} is a set of rules of the form $S_i \rightarrow S_{\ell(i)}S_{r(i)}$ or $S_i \rightarrow \sigma_i$, with $\ell(i), r(i) < i$ and $\sigma_i \in \Sigma$, for $1 \leq i \leq n$. We compute, for increasing i , the state transition function $f_i: [q] \rightarrow [q]$ (we denote states using integers $1, \dots, q$) that satisfies $z \xrightarrow{\text{eval}(S_i)} f_i(z)$. For $S_i \rightarrow S_{\ell(i)}S_{r(i)}$ we can compute f_i as $f_{r(i)} \circ f_{\ell(i)}$, where \circ is function composition. For $S_i \rightarrow \sigma_i$ we simply have $f_i(z) = z'$ for the unique transition $z \xrightarrow{\sigma_i} z'$. Hence, f_i can be computed in time $O(q)$ for every i . The text T is then accepted by F if and only if $f_n(z_0)$ is an accepting state, where z_0 is the starting state of F . Hence, the best-known algorithm takes time $O(\min\{nq, N + q\})$.

We prove that DFA Acceptance takes time $\min\{nq, N + q\}^{1-o(1)}$ assuming SETH, thus providing a conditional lower bound matching the known algorithmic results. It is straightforward to see that any algorithm must read the complete input description of F to always correctly decide the problem, yielding a lower bound of $\Omega(q)$. In the remainder, we provide the remaining conditional lower bound of $\min\{nq, N\}^{1-o(1)}$.

Theorem 3.2. *Assuming the OV conjecture, for no $\varepsilon > 0$ there is an $O(\min\{nq, N\}^{1-\varepsilon})$ -time algorithm for DFA Acceptance. This holds even restricted to instances with $N = \Theta(n^{\alpha_N})$ and $q = \Theta(n^{\alpha_q})$ for any $\alpha_N > 1$ and $\alpha_q > 0$.*

Proof. Let $\mathcal{A} = \{a_1, \dots, a_A\}, \mathcal{B} = \{b_1, \dots, b_B\}$ be a given OV instance in d dimensions. We construct a string T of length $N = O(dAB)$ with a representation \mathcal{T} of size $n = O(dA)$ and a DFA F with $q = O(dB)$. An $O(\min\{nq, N\}^{1-\varepsilon})$ -time algorithm for DFA Acceptance would then imply an algorithm for OV in time $O((d^2AB)^{1-\varepsilon}) = O((AB)^{1-\varepsilon} \text{poly}(d))$, contradicting the OV conjecture. At the end of this proof we show that this also holds for all restrictions $N = \Theta(n^{\alpha_N})$ and $q = \Theta(n^{\alpha_q})$ with $\alpha_N > 1$ and $\alpha_q > 0$.

Constructing the Text T We cast any vector $a \in \{0, 1\}^d$ to a string $T(a) := \bigcirc_{k=1}^d a[k]$ by simply concatenating its coordinates. We define the text T over the alphabet $\Sigma = \{0, 1, \#, !\}$ as

$$T := \left(! \circ \bigcirc_{i=1}^A \# T(a_i) \right)^B. \quad (1)$$

Here, we think of $!$ and $\#$ as “new group” and “new vector within group” indicators, respectively. Intuitively, the j -th repetition of $T(a_i)$ is supposed to lead to an accepting state of F if a_i and b_j are an orthogonal pair.

⁶Note that we can always define an absorbing non-accepting state z^{fail} with $z^{\text{fail}} \xrightarrow{\Sigma} z^{\text{fail}}$ and set, for any undefined transition from z under σ , $z \xrightarrow{\sigma} z^{\text{fail}}$, which increases the number of states only by one.

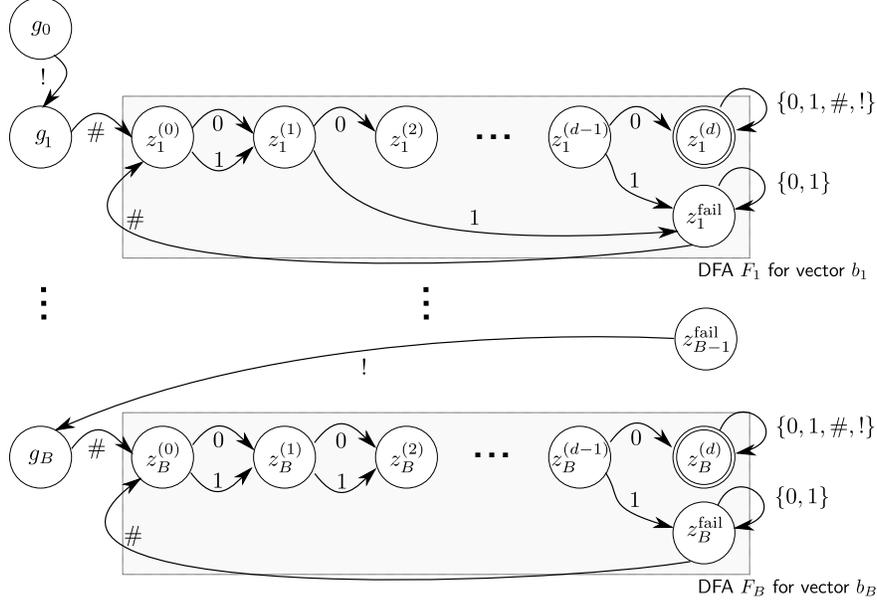


Figure 2: Illustration of the DFA F . Any transition not specified leads to a absorbing non-accepting state z^{fail} .

Constructing the DFA F For an illustration of the DFA construction see Figure 2. We start by defining “vector gadgets”: For any vector $b_j \in \mathcal{B}$ we construct a DFA F_j over alphabet $\{0, 1\}$ with states $z_j^{(0)}, z_j^{(1)}, \dots, z_j^{(d)}$ and z_j^{fail} . The initial state is $z_j^{(0)}$. For any $k \in [d]$ we have the transitions $z_j^{(k-1)} \xrightarrow{0} z_j^{(k)}$ and

$$z_j^{(k-1)} \xrightarrow{1} \begin{cases} z_j^{(k)} & \text{if } b_j[k] = 0, \\ z_j^{\text{fail}} & \text{otherwise.} \end{cases}$$

We let $z_j^{(d)}$ be an accepting state with transition $z_j^{(d)} \xrightarrow{\{0,1\}} z_j^{(d)}$. Furthermore, we have $z_j^{\text{fail}} \xrightarrow{\{0,1\}} z_j^{\text{fail}}$. It is easy to see that after reading a string $a \in \{0, 1\}^d$, F_j is in either $z_j^{(d)}$ or z_j^{fail} , and it is in $z_j^{(d)}$ if and only if a and b_j are orthogonal.

We combine these smaller DFAs to our final DFA F over the slightly larger alphabet $\Sigma = \{0, 1, \#, !\}$ as follows. We define additional states g_0, g_1, \dots, g_B and let g_0 be the initial state of F . We define the following additional transitions:

$$\begin{aligned} g_0 &\xrightarrow{!} g_1 \\ g_j &\xrightarrow{\#} z_j^{(0)} && \text{for } 1 \leq j \leq B, \\ z_{j-1}^{\text{fail}} &\xrightarrow{!} g_j && \text{for } 1 < j \leq B, \\ z_j^{\text{fail}} &\xrightarrow{\#} z_j^{(0)} && \text{for } 1 \leq j \leq B, \\ z_j^{(d)} &\xrightarrow{\{\#, !\}} z_j^{(d)} && \text{for } 1 \leq j \leq B. \end{aligned}$$

In this way, each $z_j^{(d)}$ is an absorbing accepting state, and the symbols $\#$ and $!$ satisfy the

semantics of jumping to the next vector in \mathcal{A} and \mathcal{B} , respectively. This finishes the definition of the reduction.

Correctness We claim that the constructed DFA F accepts T if and only if \mathcal{A}, \mathcal{B} contains an orthogonal pair. By structure of F and T , as well as the properties argued for $F_j, j \in [B]$, it is straightforward to show that after reading any prefix T' of T ending on $\#$, F is in the initial state of F_j , where j is the number of $!$'s in T' – this holds until F has encountered an accepting state for the first and final time. Thus, if T is accepted by F , then some prefix T' of T that ends on $\#a_i$ for some $i \in [A]$ has led an accepting state of F . This can only happen if a_i is orthogonal to b_j , where j is the number of $!$'s in T' , i.e., \mathcal{A}, \mathcal{B} contain an orthogonal pair. Conversely, if \mathcal{A}, \mathcal{B} contains an orthogonal pair, let a_i, b_j the smallest such pair in terms of the lexicographic order on (j, i) . Then the prefix T' that ends on $\#a_i$ and contains j $!$'s leads to the accepting state $z_j^{(d)}$.

Size Bounds We count that $|T| = B((d+1)A+1) = O(dAB)$. Since T consists of B repetitions of a string of length $(d+1)A+1$, we can compute an SLP \mathcal{T} of size $|\mathcal{T}| \leq O(\log(B) + dA) = O(dA)$ by Observation 2.2. The number of states of F is $O(dB)$. This satisfies the claimed size bounds. Note that the reduction can be implemented in linear time in the output size.

Strengthening the Statement In the remainder, we verify that our construction proves the desired lower bound even restricted to instances with $N = \Theta(n^{\alpha_N})$ and $q = \Theta(n^{\alpha_q})$ for any $\alpha_N > 1$ and $\alpha_q > 0$. Note that the number of states, the size of the SLP, and the text length can all three be increased by easy padding. E.g., to increase the text length we introduce a garbage symbol “ \natural ” that can be read at any state of the automaton, not changing the current state, and add a suitable number of copies of “ \natural ” to the text.

We now set $\beta := \min\{\alpha_q, \alpha_N - 1\}$ and only consider OV instances with $B = \Theta(A^\beta)$. Note that the OV conjecture asserts a lower bound of $A^{1+\beta-o(1)}$ in this setting. Note that the above construction yields $n = O(dA) = O(d^{\max\{1, 1/\beta\}}A)$, and we can pad to equality. Moreover, we have $q = O(dB) = O(d^{\max\{\beta, 1\}}B) = O(d^{\max\{\beta, 1\}}A^\beta) = O(n^\beta) = O(n^{\alpha_q})$, since $\beta \leq \alpha_q$, and we can pad to equality to obtain $q = \Theta(n^{\alpha_q})$. Similarly, we have $N = O(dAB) = O(d^{1+\max\{\beta, 1\}}A^{1+\beta}) = O(n^{1+\beta}) = O(n^{\alpha_N})$, since $\beta \leq \alpha_N - 1$, and we can pad to equality to obtain $N = \Theta(n^{\alpha_N})$. Finally, an $O(\min\{nq, N\}^{1-\varepsilon})$ -time algorithm for DFA Acceptance restricted to $N = \Theta(n^{\alpha_N})$ and $q = \Theta(n^{\alpha_q})$ would imply an algorithm for OV in time $O(\min\{n^{1+\alpha_q}, n^{\alpha_N}\}^{1-\varepsilon}) = O(n^{(1+\beta)(1-\varepsilon)}) = O((d^{\max\{1, 1/\beta\}}A)^{(1+\beta)(1-\varepsilon)}) = O(d^{\max\{1+\beta, 1+1/\beta\}}A^{(1+\beta)(1-\varepsilon)}) = O(A^{1+\beta-\varepsilon}\text{poly}(d))$, contradicting the OV conjecture. This finishes the proof. \square

3.2 Approximate Pattern Matching and Substring Hamming Distance

We study the following generalization of pattern matching.

Problem 3.3 (Generalized Pattern Matching). *Given a text T of length N by an SLP \mathcal{T} of size n , a pattern P of length M by an SLP \mathcal{P} of size m , both over some alphabet Σ , and given a cost function $\text{cost}: \Sigma \times \Sigma \rightarrow \mathbb{N}$, compute $\min_{0 \leq i \leq N-M} \sum_{j=1}^M \text{cost}(P[j], T[i+j])$, i.e., the minimum total cost of any alignment.*

In other words, we want to compute the length- M substring T' of T minimizing the total cost of aligned symbols in P and T' . This problem has two important special cases: (1) We obtain

Substring Hamming Distance when $\text{cost}(\sigma, \sigma') = [\sigma \neq \sigma']$ for any $\sigma, \sigma' \in \Sigma$. (2) We obtain Pattern Matching with Wildcards when T is over alphabet Σ and P is over alphabet $\Sigma \cup \{*\}$, we have $\text{cost}(*, \sigma) = 0$ for any $\sigma \in \Sigma$ and $\text{cost}(\sigma, \sigma') = [\sigma \neq \sigma']$ for any $\sigma, \sigma' \in \Sigma$, and the task is to decide whether the minimum total cost of any alignment is 0.

Problem 3.4 (Substring Hamming Distance). *Given a text T of length N by an SLP \mathcal{T} of size n and a pattern P of length M by an SLP \mathcal{P} of size m , both over some alphabet Σ , compute $\min_{0 \leq i \leq N-M} \sum_{j=1}^M [P[j] \neq T[i+j]]$, i.e., the minimum Hamming distance of any alignment.*

Problem 3.5 (Pattern Matching with Wildcards). *For some alphabet Σ , we are given a text T of length N by an SLP \mathcal{T} of size n over alphabet Σ and a pattern P of length M by an SLP \mathcal{P} of size m over alphabet $\Sigma \cup \{*\}$. We say that $\sigma' \in \Sigma \cup \{*\}$ and $\sigma \in \Sigma$ match if $\sigma' = *$ or $\sigma' = \sigma$. Decide whether for some offset $0 \leq i \leq N - M$ all pairs $P[j], T[i+j]$ match for $1 \leq j \leq M$.*

In this section, for all three problems we show an upper bound of $O(\min\{|\Sigma|N \log N, nM\})$ and a SETH-based lower bound of $\min\{N, nM\}^{1-o(1)}$. This yields a tight bound in case of constant alphabet size, as the lower bound constructs constant-alphabet strings. We leave it as an open problem to get tight bounds for larger alphabet size.

Note that it suffices to prove the upper bound for Generalized Pattern Matching and the lower bound for the special cases Substring Hamming Distance and Pattern Matching with Wildcards. We start with the following two upper bounds, which follow standard arguments.

Lemma 3.6. *Generalized Pattern Matching can be solved in time $O(|\Sigma|N \log N)$.*

Proof. Decompress both the text T and the pattern P . For each symbol $\sigma \in \Sigma$, build the vector $v^\sigma \in \mathbb{R}^N$ with $v_i^\sigma := \text{cost}(\sigma, T[i])$ and the vector $u^\sigma \in \{0, 1\}^M$ with $u_j^\sigma := [P[j] = \sigma]$. Compute their convolution $w^\sigma \in \mathbb{R}^{N-M+1}$ with $w_i^\sigma = \sum_{j=1}^M u_j^\sigma v_{i+j}^\sigma$. Using FFT, w^σ can be computed in time $O(N \log N)$. Finally, compute the vector $r \in \mathbb{R}^{N-M+1}$ with $r_i = \sum_{\sigma \in \Sigma} w_i^\sigma$ and return the minimal entry of r . Note that

$$r_i = \sum_{\sigma \in \Sigma} w_i^\sigma = \sum_{\sigma \in \Sigma} \sum_{j=1}^M u_j^\sigma v_{i+j}^\sigma = \sum_{\sigma \in \Sigma} \sum_{j=1}^M [P[j] = \sigma] \cdot \text{cost}(\sigma, T[i+j]) = \sum_{j=1}^M \text{cost}(P[j], T[i+j]),$$

which proves correctness. □

Lemma 3.7. *Generalized Pattern Matching can be solved in time $O(nM)$.*

Proof Sketch. Let S_1, \dots, S_n be the non-terminals of the SLP \mathcal{T} that generates the text T . In this proof, for simplicity we write $T_i := \text{eval}(S_i)$. We decompress the pattern P . For $1 \leq i \leq n$ we define

$$\text{Match}(i) := \min_{0 \leq d \leq |T_i| - M} \sum_{j=1}^M \text{cost}(P[j], T_i[j+d]),$$

or ∞ , if $|T_i| < M$. This solves the Generalized Pattern Matching problem restricted to the substring T_i of T . Clearly, we can solve the given Generalized Pattern Matching instance (T, P) by calling $\text{Match}(n)$. Moreover, for any offset d and any $i \in [n]$ we define

$$\text{FixMatch}(i, d) := \sum_{\substack{j: 1 \leq j+d \leq |T_i| \\ 1 \leq j \leq M}} \text{cost}(P[j], T_i[j+d]).$$

In other words, $\text{FixMatch}(i, d)$ is equal to the total cost between T_i and a shifted pattern P (by d symbols to the right, or $-d$ symbols to the left), where we consider only the symbols that have an aligned counterpart.

In the remainder we show how to compute these functions by simple recursive algorithms. We precompute all lengths $|T_i|$ in time $O(n)$. For $\text{FixMatch}(\cdot, \cdot)$, observe that for a rule $S_i \rightarrow S_\ell S_r$ we have

$$\text{FixMatch}(i, d) = \text{FixMatch}(\ell, d) + \text{FixMatch}(r, d - |T_\ell|),$$

since the offset with respect to the first symbol of T_r differs to the offset with respect to the first symbol of T_i by $|T_\ell|$. Moreover, for a rule $S_i \rightarrow \sigma \in \Sigma$ we can compute $\text{FixMatch}(i, d)$ in constant time. Note that whenever the offset d is such that no symbols get aligned, we can immediately return 0. This completes our algorithm for $\text{FixMatch}(\cdot, \cdot)$.

Now consider $\text{Match}(i)$. For a rule $S_i \rightarrow S_\ell S_r$, the optimal alignment of the pattern in T_i is either completely contained in T_ℓ or completely contained in T_r or it has a non-empty intersection with both of them, in which case it has an offset $-M < d < 0$ with respect to the starting symbol of T_r , or equivalently an offset $|T_\ell| + d$ with respect to the starting symbol of T_ℓ . Hence, we have

$$\text{Match}(i) = \min \left\{ \text{Match}(\ell), \text{Match}(r), \min_{-M < d < 0} \text{FixMatch}(r, d) + \text{FixMatch}(\ell, |T_\ell| + d) \right\}.$$

Again, for a rule $S_i \rightarrow \sigma \in \Sigma$ we can compute $\text{Match}(i)$ in constant time. This completes the algorithm for $\text{Match}(\cdot)$.

To obtain the claimed running time, we use memoization to ensure that each argument is called at most once. Clearly, there are n possible arguments for $\text{Match}(\cdot)$, and each call takes time $O(M)$, resulting in time $O(nM)$. Note that $\text{Match}(\cdot)$ only calls $\text{FixMatch}(i, d)$ for offsets d such that the pattern crosses the left or right boundary of T_i . This property also holds as an invariant in the recursive subproblems of $\text{FixMatch}(i, d)$. Hence, there are less than $2M$ possible offsets d (i.e., less than M offsets for the left and right boundary). As there are n possible values for i , and each call to $\text{FixMatch}(\cdot, \cdot)$ takes time $O(1)$, we obtain the claimed total running time of $O(nM)$. \square

This completes the upper bound $O(\min\{|\Sigma|N \log N, nM\})$ for Generalized Pattern Matching. It remains to prove the SETH-based lower bound of $\min\{N, nM\}^{1-o(1)}$ for Substring Hamming Distance and Pattern Matching with Wildcards.

We now make the intuition given at the beginning of Section 3 formal, by designing a text T that enumerates all combinations of k vectors in a given k -OV instance, while still being well compressible. We give a slightly more general construction that will also be useful later for our SETH-based lower bounds for LCS, see Section 3.3. As usual, we consider k as a constant.

Lemma 3.8. *Consider a k -OV instance $\mathcal{A} = \{a_1, \dots, a_A\} \subseteq \{0, 1\}^d$. Let $b \in \{0, 1\}^d$ be an additional vector, and let $S(0), S(1)$ be strings of length γ ($S(i)$ is a sequence that represents an entry that is equal to i). We define the tuplified representation as follows:*

$$\begin{aligned} V &= \text{tuplify}(\mathcal{A}, k, b, S(0), S(1)) \\ &:= \bigcirc_{\ell=1}^d \bigcirc_{i_1, \dots, i_k \in [A]} S\left(b[\ell] \cdot a_{i_1}[\ell] \cdots a_{i_k}[\ell]\right), \end{aligned}$$

where the second \bigcirc goes over all tuples $(i_1, \dots, i_k) \in [A]^k$ in lexicographic order. This representation satisfies the following properties.

1. We can compute, in linear time in the output size, an SLP \mathcal{V} generating V of size $O(dA + \gamma)$ or, when given SLPs $\mathcal{S}(0), \mathcal{S}(1)$ generating $S(0), S(1)$, of size $O(dA + |\mathcal{S}(0)| + |\mathcal{S}(1)|)$.
2. Write $V = \bigcirc_{i=1}^{dA^k} V_i$ with $V_i \in \{S(0), S(1)\}$. Then there exist $i_1, \dots, i_k \in [A]$ such that $(b, a_{i_1}, \dots, a_{i_k})$ is orthogonal if and only if there is an offset $1 \leq \Delta \leq A^k$ such that

$$V_\Delta = V_{\Delta+A^k} = \dots = V_{\Delta+(d-1)A^k} = S(0).$$

Proof. For the second property, note that by definition $V_\Delta, V_{\Delta+A^k}, \dots, V_{\Delta+(d-1)A^k}$ are all equal to $S(0)$ for $\Delta \in [A^k]$ if and only if the Δ -th tuple $(i_1, \dots, i_k) \in [A]^k$ in the lexicographic ordering of $[A]^k$ satisfies

$$b[\ell] \cdot a_{i_1}[\ell] \cdots a_{i_k}[\ell] = 0 \quad \text{for all } \ell \in [d].$$

This condition is equivalent to $(b, a_{i_1}, \dots, a_{i_k})$ being an orthogonal pair, so the claim follows.

It remains to construct a short SLP \mathcal{V} generating V . We construct non-terminals $P_{S(0)}, P_{S(1)}$ with $\text{eval}(P_{S(i)}) = S(i)$ by an SLP of size $\gamma_S = O(\gamma)$ as in Observation 2.2, or of size $\gamma_S = O(|\mathcal{S}(0)| + |\mathcal{S}(1)|)$ by using given SLPs $\mathcal{S}(0), \mathcal{S}(1)$. We can extend this, using Observation 2.2, to a slightly larger SLP of size $O(\log A + \gamma_S)$ that includes, for every $1 \leq j \leq k$, a non-terminal $P_{S(0)}^j$ with $\text{eval}(P_{S(0)}^j) = S(0)^{A^j}$.

The crucial observation is the following: for any tuple $(i_1, \dots, i_k) \in [A]^k$, let $p_\ell(i_1, \dots, i_k) = a_{i_1}[\ell] \cdots a_{i_k}[\ell]$. Then for any $\ell \in [d], j \in [k]$ and $(i_1, \dots, i_j) \in [A]^j$, we have that $a_{i_j}[\ell] = 0$ implies $p_\ell(i_1, \dots, i_j, i'_{j+1}, \dots, i'_k) = 0$ for all $(i'_{j+1}, \dots, i'_k) \in [A]^{k-j}$. We now define the final SLP using the starting non-terminal S_0 and the following productions

$$\begin{aligned} S_0 &\rightarrow \text{Test}_1 \dots \text{Test}_d \\ \text{Test}_\ell &\rightarrow \begin{cases} P_{S(0)}^k & \text{if } b[\ell] = 0 \\ \text{List}_\ell^{(1)} & \text{otherwise} \end{cases} & \ell \in [d], \\ \text{List}_\ell^{(j)} &\rightarrow \bigcirc_{i \in [A]} \begin{cases} P_{S(0)}^{k-j} & \text{if } a_i[\ell] = 0, \\ \text{List}_\ell^{(j+1)} & \text{otherwise} \end{cases} & \ell \in [d], j \in [k], \\ \text{List}_\ell^{(k+1)} &\rightarrow P_{S(1)}. \end{aligned}$$

It is straight-forward to verify that $\text{eval}(S_0) = V$. Note that the size of this SLP, i.e., the total number of non-terminals on the right hand side of the above rules, is bounded by $O(\gamma_S + dA)$. Moreover, the SLP can be constructed in linear time in its size. \square

After this preparation, we can prove our conditional lower bounds.

Theorem 3.9. *Assuming the k -OV conjecture, Pattern Matching with Wildcards over alphabet $\{0, 1\}$ (plus wildcards $*$) takes time $\min\{N, nM\}^{1-o(1)}$. This holds even restricted to instances with $n = \Theta(N^{\alpha_n})$, $M = \Theta(N^{\alpha_M})$ and $m = \Theta(N^{\alpha_m})$ for any $0 < \alpha_n < 1$ and $0 < \alpha_m \leq \alpha_M \leq 1$.*

Before we prove Theorem 3.9, let us sketch the main idea by providing a simple $N^{1-o(1)}$ -time conditional lower bound in the setting $n, m = O(N^\epsilon)$ and $N = \Theta(M)$. Let $\mathcal{A} \subseteq \{0, 1\}^d$ of size A be an arbitrary k -OV instance with $k > 1/\epsilon$, and assume for simplicity $d \leq A^{o(1)}$. Using Lemma 3.8 on \mathcal{A} , k , $S(0) = 0, S(1) = 1$ and $b = (1, \dots, 1) \in \{0, 1\}^d$, we compute an SLP \mathcal{T} for

$$T = \text{tuplify}(\mathcal{A}, k, b, S(0), S(1)).$$

We define the pattern P as

$$P = 0(*^{A^k-1}0)^{d-1}.$$

Note that Pattern Matching with Wildcards on instance T, P checks whether for some offset Δ we have $T[\Delta] = T[\Delta + A^k] = \dots = T[\Delta + (d-1)A^k] = 0$. Hence, by Lemma 3.8, pattern P matches T if and only if there is an orthogonal tuple $(a_1, \dots, a_k) \in \mathcal{A}^k$, showing correctness of the reduction.

Note that we have $N = \Theta(M) = \Theta(dA^k)$. By Lemma 3.8, T has an SLP of size $O(dA)$, and by Observation 2.2, P has an SLP of size $O(d \log A)$. By $d \leq A^{o(1)}$ and $k > 1/\varepsilon$, we are indeed in the setting $n, m = O(N^\varepsilon)$ and $N = \Theta(M)$. An $O(N^{1-\varepsilon})$ algorithm for Pattern Matching with Wildcards would now imply an $O(A^{k(1-\varepsilon)})\text{poly}(d)$ for k -OV, contradicting the k -OV conjecture.

We now give the slightly more involved general construction.

Proof of Theorem 3.9. For $k \geq 2$, let $\mathcal{A} = \{a_1, \dots, a_A\}$ be a k -OV instance in d dimensions, and let $k_1, k_2 \geq 1$ with $k_1 + k_2 = k$. We will construct an equivalent instance of Pattern Matching with Wildcards with $N = O(dA^k)$, $M = O(dA^{k_1})$, $n = O(dA^{k_2+1})$, and $m = O(d \log A)$. Any $O(\min\{N, nM\}^{1-\varepsilon})$ algorithm for Pattern Matching with Wildcards would then imply an algorithm for k -OV in time $O(A^{(k+1)(1-\varepsilon)})\text{poly}(d) = O(A^{k(1-\varepsilon/2)})\text{poly}(d)$ for $k \geq 2/\varepsilon$, contradicting the k -OV conjecture. Below we strengthen this statement to hold restricted to instances with $n = \Theta(N^{\alpha_n})$, $M = \Theta(N^{\alpha_M})$ and $m = \Theta(N^{\alpha_m})$ for any $0 < \alpha_n < 1$ and $0 < \alpha_m \leq \alpha_M \leq 1$.

To give such a reduction, we define the text as

$$T = \bigcirc_{(j_1, \dots, j_{k_2}) \in [A]^{k_2}} 1^{A^{k_1}} \circ \text{tuplify}(\mathcal{A}, k_1, \min(a_{j_1}, \dots, a_{j_{k_2}}), 0, 1),$$

where $\min(b_1, \dots, b_\ell)$ denotes the component-wise minimum of b_1, \dots, b_ℓ .

We define the pattern P as

$$P = 0(*^{A^{k_1}-1}0)^{d-1}.$$

Correctness Observe that P cannot overlap any $1^{A^{k_1}}$ -block, since never more than $A^{k_1} - 1$ wildcards are followed by a 0 in P . Thus, P matches T if and only if there is a tuple $(j_1, \dots, j_{k_2}) \in [A]^{k_2}$ such that P matches $T((j_1, \dots, j_{k_2})) := \text{tuplify}(\mathcal{A}, k_1, \min(a_{j_1}, \dots, a_{j_{k_2}}), 0, 1)$. By the structure of the pattern, P matches any string S if and only if there is an offset Δ such that $S[\Delta] = S[\Delta + A^{k_1}] = \dots = S[\Delta + (d-1)A^{k_1}] = 0$. Thus, by Lemma 3.8, P matches $T((j_1, \dots, j_{k_2}))$ if and only if there are vectors $a_1, \dots, a_{k_1} \in \mathcal{A}$ for which $(a_1, \dots, a_{k_1}, \min(a_{j_1}, \dots, a_{j_{k_2}}))$ is an orthogonal tuple. The latter condition is equivalent to $(a_1, \dots, a_{k_1}, a_{j_1}, \dots, a_{j_{k_2}})$ being an orthogonal tuple. Since $k_1 + k_2 = k$ and T contains $T((j_1, \dots, j_{k_2}))$ for all $(j_1, \dots, j_{k_2}) \in [A]^{k_2}$, this proves that P matches T if and only if there is an orthogonal k -tuple in the instance \mathcal{A} .

Size Bounds Note that $N = |T| = O(dA^k)$. By Lemma 3.8 and Observation 2.2, we can compute an SLP \mathcal{T} of size $n = O(dA^{k_2+1})$ generating T , in linear time. Similarly, note that $M = |P| = O(dA^{k_1})$. By Observation 2.2, we can compute an SLP \mathcal{P} of length $m = O(d \log A)$ generating P , in linear time. This proves the claimed bounds.

Strengthening the Statement We now prove the lower bound restricted to instances with $n = \Theta(N^{\alpha_n})$, $M = \Theta(N^{\alpha_M})$ and $m = \Theta(N^{\alpha_m})$ for any $0 < \alpha_n < 1$ and $0 < \alpha_m \leq \alpha_M \leq 1$. Let $\varepsilon > 0$ and set $\beta := \min\{1, \alpha_M + \alpha_n\}$. We choose $k_1, k_2 \geq 1$ such that $k_1 + k_2 = k$ and

$k_1 \approx \min\{\alpha_M, 1 - \alpha_n\}k/\beta$ and $k_2 \approx \alpha_n k/\beta$. Note that k_1, k_2 are restricted to be integers, however, for sufficiently large k depending only on $\varepsilon, \alpha_M, \alpha_n$, we can ensure $k_1 \leq (1 + \varepsilon/4) \min\{\alpha_M, 1 - \alpha_n\}k/\beta$ and $k_2 + 1 \leq (1 + \varepsilon/4)\alpha_n k/\beta$. Note that for the dimension d we can assume $d \leq A$, since otherwise an $O(A^{k-\varepsilon} \text{poly}(d))$ algorithm clearly exists. In particular, for sufficiently large k we have $d \leq A^{(\varepsilon/4) \cdot \min\{\alpha_M, \alpha_n, 1 - \alpha_n\}k/\beta}$. This yields

$$\begin{aligned} N &= O(dA^k) = O(A^{(1+\varepsilon/2)k/\beta}), \\ M &= O(dA^{k_1}) = O(A^{(1+\varepsilon/2) \min\{\alpha_M, 1 - \alpha_n\}k/\beta}) = O(A^{(1+\varepsilon/2)\alpha_M k/\beta}), \\ n &= O(dA^{k_2+1}) = O(A^{(1+\varepsilon/2)\alpha_n k/\beta}), \\ m &= O(d \log A) = O(A^{(1+\varepsilon/2)\alpha_m k/\beta}). \end{aligned}$$

Standard padding⁷ of these four parameters allows us to achieve equality, up to constant factors, in the above inequalities, which yields the desired $n = \Theta(N^{\alpha_n})$, $M = \Theta(N^{\alpha_M})$ and $m = \Theta(N^{\alpha_m})$. Any $O(\min\{N, nM\}^{1-\varepsilon})$ algorithm for Pattern Matching with Wildcards in this setting would now imply an algorithm for k -OV in time $O(\min\{A^{(1+\varepsilon/2)k/\beta}, A^{(1+\varepsilon/2)(\alpha_M + \alpha_n)k/\beta}\}^{1-\varepsilon}) = O(A^{(1+\varepsilon/2)(1-\varepsilon) \min\{1, \alpha_M + \alpha_n\}k/\beta}) = O(A^{(1-\varepsilon/2)k})$, where we used the definition of β and $(1 + \varepsilon/2)(1 - \varepsilon) \leq 1 - \varepsilon/2$. This contradicts the k -OV conjecture, finishing the proof. \square

We next prove a lower bound similar to Theorem 3.9 for another special case of generalized pattern matching, namely Substring Hamming Distance. Instead of a direct reduction from k -OV, we present a linear-time reduction from Pattern Matching with Wildcards over alphabet $\{0, 1\}$ to Substring Hamming Distance.

Theorem 3.10. *Assuming the k -OV conjecture, Substring Hamming Distance on constant-size alphabet takes time $\min\{N, nM\}^{1-o(1)}$. This holds even restricted to instances with $n = \Theta(N^{\alpha_n})$, $M = \Theta(N^{\alpha_M})$ and $m = \Theta(N^{\alpha_m})$ for any $0 < \alpha_n < 1$ and $0 < \alpha_m \leq \alpha_M \leq 1$.*

Proof. For short, we write $d_H(X, Y)$ for the Hamming distance of strings X, Y . We prove the result by reducing any Pattern Matching with Wildcards instance $T_{\text{PM}}, P_{\text{PM}}$ over alphabet $\Sigma = \{0, 1\}$ to an instance $T_{\text{HD}}, P_{\text{HD}}$ of Substring Hamming Distance. We first define coordinate strings

$$\begin{aligned} s_T(0) &:= 100, & s_T(1) &:= 010, \\ s_P(0) &:= 101, & s_P(1) &:= 011, & s_P(*) &:= 000. \end{aligned}$$

Observe that these strings are defined in such a way that $d_H(s_P(*), s_T(y)) = 1$ for $y \in \{0, 1\}$, $d_H(s_P(x), s_T(y)) = 1$ for $x = y \in \{0, 1\}$, and $d_H(s_P(x), s_T(y)) = 3$ if $x \neq y, x, y \in \{0, 1\}$.

We introduce the *guarding* $G(s) := s \circ 234$ for length-3 strings $s \in \{0, 1\}^3$. This allows us to reduce $T_{\text{PM}}, P_{\text{PM}}$ to the following instance, using alphabet $\Sigma = \{0, 1, 2, 3, 4\}$,

$$\begin{aligned} T_{\text{HD}} &:= G(s_T(T_{\text{PM}}[1])) \dots G(s_T(T_{\text{PM}}[N])), \\ P_{\text{HD}} &:= G(s_P(P_{\text{PM}}[1])) \dots G(s_P(P_{\text{PM}}[M])). \end{aligned}$$

Note that for any $0 \leq i \leq N - M$,

$$\begin{aligned} d_H(T_{\text{HD}}[6i + 1..6i + 6M], P_{\text{HD}}) &= \sum_{j=1}^M d_H(s_P(P_{\text{PM}}[j]), s_T(T_{\text{PM}}[i + j])) \\ &= M + 2 \cdot \text{mismatch}(T_{\text{PM}}[i + 1..i + M], P_{\text{PM}}), \end{aligned}$$

⁷Add a prefix of wildcards to the pattern and a prefix of 1's to the text, and partially decompress the SLPs.

where $\text{mismatch}(z, z') = \#\{i \mid z'[i] \neq *, z[i] \neq z'[i]\}$ is the number of mismatches of z and z' .

We now observe that for all i with $i \bmod 6 \neq 0$, we have $d_H(T_{\text{HD}}[i + 1..i + 6M], P_{\text{HD}}) \geq 3M$, as no two symbols 2, 3, 4 in P_{HD} are aligned, so that each $G(s_P(P_{\text{PM}}[j]))$ contributes at least 3 to the Hamming distance. Since $d_H(T_{\text{HD}}[6i + 1..6i + 6M], P_{\text{HD}}) \leq 3M$ for all i , the substring with smallest Hamming distance has thus a Hamming distance of $M + 2 \cdot \min_{0 \leq i \leq N-M} \text{mismatch}(T_{\text{PM}}[i + 1..i + M], P_{\text{PM}})$. This value is equal to M if and only if P_{PM} matches T_{PM} , proving correctness.

The corresponding reduction of the compressed problems is straightforward: We can augment the SLP \mathcal{T}_{PM} for T_{PM} by $O(1)$ -sized productions to obtain an SLP \mathcal{T}_{HD} for T_{HD} , by replacing each terminal $\sigma \in \{0, 1, *\}$ by a non-terminal evaluating to $G(s_T(\sigma))$. Analogously, we can compute an SLP for P_{HD} of size $|\mathcal{P}_{\text{HD}}| = |\mathcal{P}_{\text{PM}}| + O(1)$ in linear time. Overall, since also $|T_{\text{HD}}| = O(|T_{\text{PM}}|)$, $|P_{\text{HD}}| = O(|P_{\text{PM}}|)$, all parameters are preserved up to constant factors. By this linear-time parameter-preserving reduction, the lower bound of Theorem 3.9 translates to Substring Hamming Distance, yielding the claim. \square

3.3 Longest Common Subsequence

In this section, we study the Longest Common Subsequence (LCS) problem. Recall that a string S of length ℓ is a substring of a string X if there are $1 \leq i_1 < \dots < i_\ell \leq |X|$ with $S[j] = X[i_j]$ for any $j \in [\ell]$. In the LCS problem, given two strings X, Y , the task is to determine the longest string S that is a subsequence of both X and Y . We denote the length of the LCS by $L(X, Y) = |S|$, and more precisely consider the problem of computing $L(X, Y)$. In the whole section, the alphabet Σ has constant size.

Problem 3.11 (LCS). *Given strings X, Y of length at most N by grammar-compressed representations \mathcal{X}, \mathcal{Y} of size at most n , compute the length of the LCS of X and Y .*

As discussed in the introduction, the $O(nN\sqrt{\log N/n})$ time algorithm by Gawrychowski [35] is the fastest known. Here we prove a matching lower bound of $(Nn)^{1-o(1)}$, assuming the k -OV conjecture.

Theorem 3.12. *Assuming the k -OV conjecture, there is no $(nN)^{1-\varepsilon}$ -time algorithm for LCS for any $\varepsilon > 0$. This even holds restricted to instances with $n = \Theta(N^{\alpha_n})$ for any $0 < \alpha_n < 1$, and an alphabet of constant size.*

The general approach is very similar to the lower bound for Pattern Matching with Wildcards given in Section 3.2. In particular, we again use the tuplified representation $T = \bigcirc_{i=1}^{dA^k} T[i]$ of Lemma 3.8 for a k -OV instance \mathcal{A} . Recall that this allows us to decide the k -OV instance by testing whether there is a subsequence of d substrings $T[\Delta], T[\Delta + A^k], \dots, T[\Delta + (d-1)A^k]$ all equal to a certain 0-coordinate string. Finding a pattern to test this was quite simple for Pattern Matching with Wildcards, yielding an $N^{1-o(1)}$ lower bound. For LCS, enforcing a coherent offset is much more complicated, since the “pattern” is matched as a subsequence not as a substring. Furthermore, the extension to a $(nN)^{1-o(1)}$ lower bound is more involved and relies on the quadratic-time nature of LCS. Fortunately, we can overcome the technical obstacles for LCS using (an extension of) alignment gadgets developed in [17]. We first redevelop and extend the corresponding alignment gadget tools in Section 3.3.1, then give the lower bound for compressed instances for general distance measures in Section 3.3.2 and then finish our LCS lower bound by designing an alignment gadget for LCS in Section 3.3.3.

3.3.1 Alignment Gadget Framework

We start by reviewing and adapting the definitions of [17]. In particular, we extend the alignment gadget definition for our purposes.

More generally than LCS, we consider an arbitrary *similarity measure* $\delta : \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{N}$. For LCS, the set of inputs \mathcal{I} is the set of all strings over some sufficiently large constant-sized alphabet Σ , and $\delta(X, Y) := |X| + |Y| - 2L(X, Y)$, where $L(X, Y)$ is the length of the LCS of X and Y .

Any sequence $X \in \mathcal{I}$ is assigned an (abstract) type $\text{type}(X)$. For LCS, we use $\text{type}(X) := (|X|, \Sigma)$, where $|X|$ is the length of X and Σ the alphabet over which X is defined. We define $\mathcal{I}_t := \{X \in \mathcal{I} \mid \text{type}(X) = t\}$ as the set of all inputs of type t .

Alignments Let $n \geq m$. An *alignment* is a set $\Lambda = \{(i_1, j_1), \dots, (i_k, j_k)\}$ with $0 \leq k \leq m$ such that $1 \leq i_1 < \dots < i_k \leq n$ and $1 \leq j_1 < \dots < j_k \leq m$. We say that $(i, j) \in \Lambda$ are *aligned*. Any $i \in [n]$ or $j \in [m]$ that is not contained in any pair in Λ is called *unaligned*. We denote the set of all alignments (with respect to n, m) by $\mathbf{\Lambda}_{n,m}$.

We call the alignment $\{(\Delta + 1, 1), \dots, (\Delta + m, m)\}$, with $0 \leq \Delta \leq n - m$, a *structured alignment*. We denote the set of all structured alignments by $\mathcal{S}_{n,m}$.

Defining the *cost* of an alignment $\Lambda \in \mathbf{\Lambda}_{n,m}$, we deviate from [17]: for any $X_1, \dots, X_n \in \mathcal{I}$ and $Y_1, \dots, Y_m \in \mathcal{I}$, we define the cost of $\Lambda = \{(i_1, j_1), \dots, (i_{|\Lambda|}, j_{|\Lambda|})\}$ as

$$\text{cost}(\Lambda) = \text{cost}_{Y_1, \dots, Y_m}^{X_1, \dots, X_n}(\Lambda) := \sum_{k=1}^{|\Lambda|} \delta(X_{i_k}, Y_{j_k}) + \begin{cases} (m - |\Lambda|)\gamma, & \text{if } |\Lambda| < m \\ (i_m - i_1 - m + 1)\gamma & \text{if } |\Lambda| = m, \end{cases}$$

where we set $\gamma := \max_{i,j} \delta(X_i, Y_j)$. In other words, (1) for any $j \in [m]$ which is aligned to some i , we “pay” the distance $\delta(X_i, Y_j)$, (2) if Λ is unstructured because it contains an unaligned j , we “pay” a penalty of γ for each such unaligned j (note that there are $m - |\Lambda|$ unaligned $j \in [m]$) and (3) if Λ is unstructured because it aligns all j but leaves out some i between the first and last aligned i , then for any unaligned i that is between the first aligned i_1 and last aligned $i_{|\Lambda|}$, we also “pay” a penalty of γ (note that $\sum_{k=1}^{|\Lambda|-1} (i_{k+1} - i_k - 1) = i_{|\Lambda|} - i_1 - |\Lambda| + 1$). This means that we incur punishment for *any* deviation from a structured alignment.

In [17], the cost of an alignment was defined to be the smaller quantity $\sum_{k=1}^{|\Lambda|} \delta(X_{i_k}, Y_{j_k}) + (m - |\Lambda|)\gamma$, i.e., unstructured alignments (that still align all $j \in [m]$) were punished less. For structured alignments both definitions coincide. Hence, the following extended alignment gadget is more powerful than the alignment gadget defined in [17].

Definition 3.13 (Extended alignment gadget). *The similarity measure δ admits an extended alignment gadget, if the following conditions hold: given instances $X_1, \dots, X_n \in \mathcal{I}_{t_x}$, $Y_1, \dots, Y_m \in \mathcal{I}_{t_y}$ with $m \leq n$ and types $t_x = (\ell_x, \Sigma), t_y = (\ell_y, \Sigma)$, we can construct new instances $X = \text{GA}_X^{m, t_y}(X_1, \dots, X_n)$ and $Y = \text{GA}_Y^{n, t_x}(Y_1, \dots, Y_m)$ and $C \in \mathbb{Z}$ such that*

$$\min_{\Lambda \in \mathbf{\Lambda}_{n,m}} \text{cost}(\Lambda) \leq \delta(X, Y) - C \leq \min_{\Lambda \in \mathcal{S}_{n,m}} \text{cost}(\Lambda). \quad (2)$$

Moreover, $\text{type}(X)$, $\text{type}(Y)$ and C only depend on n, m, t_x, t_y . Finally, $|X|, |Y| = \Theta((n + m)(\ell_x + \ell_y))$.

Definition 3.14 (Compressible alignment gadget). *We call an extended alignment gadget compressible, if X and Y are of the form $X = X_L (\bigcirc_{i=1}^n \text{pad}_X(X_i)) X_R$ and $Y = Y_L (\bigcirc_{j=1}^m \text{pad}_Y(Y_j)) Y_R$ for some strings X_L, X_R, Y_L, Y_R and functions $\text{pad}_X : \mathcal{I}_{\ell_X} \rightarrow \mathcal{I}$ and $\text{pad}_Y : \mathcal{I}_{\ell_Y} \rightarrow \mathcal{I}$ that satisfy the following properties:*

1. X_L, X_R, Y_L, Y_R have SLPs of size $O(\log n + \log(\ell_X + \ell_Y))$, computable in linear time in the output.
2. Given SLPs $\mathcal{X}_i, \mathcal{Y}_j$ for X_i, Y_j , we can compute SLPs for $\text{pad}_X(X_i), \text{pad}_Y(Y_j)$ of size $O(|\mathcal{X}_i| + \log(\ell_X + \ell_Y)), O(|\mathcal{Y}_j| + \log(\ell_X + \ell_Y))$ in linear time in the output.

In Section 3.3.3, we provide a compressible extended alignment gadget for LCS.

At the lowest level of our construction, we need the following notion.

Definition 3.15. *The similarity measure δ admits coordinate values, if there exist $\mathbf{0}_X, \mathbf{0}_Y, \mathbf{1}_X, \mathbf{1}_Y \in \mathcal{I}$ satisfying*

$$\delta(\mathbf{1}_X, \mathbf{1}_Y) > \delta(\mathbf{0}_X, \mathbf{1}_Y) = \delta(\mathbf{0}_X, \mathbf{0}_Y) = \delta(\mathbf{1}_X, \mathbf{0}_Y),$$

and, moreover, $\text{type}(\mathbf{0}_X) = \text{type}(\mathbf{1}_X)$ and $\text{type}(\mathbf{0}_Y) = \text{type}(\mathbf{1}_Y)$.

3.3.2 General Lower Bound

The following theorem proves a conditional lower bound of $(Nn)^{1-o(1)}$ for any similarity measure admitting a compressible extended alignment gadget and coordinate values.

Theorem 3.16. *Let δ be a similarity measure admitting a compressible extended alignment gadget and coordinate values. Then unless the k -OV conjecture fails, there is no $(nN)^{1-o(1)}$ -time algorithm for computing the value $\delta(X, Y)$, given SLPs \mathcal{X}, \mathcal{Y} of size at most n generating strings X, Y of length at most N . This even holds restricted to instances with $n = \Theta(N^{\alpha_n})$ for any $0 < \alpha_n < 1$, and constant alphabet size.*

We prove this theorem in the remainder of this section.

Let $\mathcal{A} = \{a_1, \dots, a_A\}$ be a k -OV instance in $d - 1$ dimensions. We augment all vectors in \mathcal{A} by another dimension where all vectors are 0 to obtain \mathcal{A}_0 , or where all vectors are 1 to obtain \mathcal{A}_1 . For any $k' \geq 1$ we let $\mathcal{A}^{(k')} := \{\min(a_{i_1}, \dots, a_{i_{k'}}) \mid i_1, \dots, i_{k'} \in [A]\}$, i.e., for each k' -tuple of vectors in \mathcal{A} the set $\mathcal{A}^{(k')}$ contains the pointwise minimum of this k' -tuple. Note that $\mathcal{A}^{(k')}$ is in general a multiset, it has size $|\mathcal{A}^{(k')}| = A^{k'}$, and is naturally ordered by the lexicographic ordering on k' -tuples $(i_1, \dots, i_{k'}) \in [A]^{k'}$. Similarly, we define $\mathcal{A}_0^{(k')}$ and $\mathcal{A}_1^{(k')}$ for the augmented vectors. We split $k = k_1 + 2k_2$ for some $k_1, k_2 \geq 1$ and set

$$\mathbf{A} := \mathcal{A}_0^{(k_1)}, \quad \mathbf{B} = \mathcal{A}_0^{(k_2)}, \quad \mathbf{C} := \mathcal{A}_1^{(k_2)}.$$

Observe that deciding the given k -OV instance is equivalent to testing whether there are orthogonal vectors $(\mathbf{a}, \mathbf{b}, \mathbf{c})$ with $\mathbf{a} \in \mathbf{A}$, $\mathbf{b} \in \mathbf{B}$ and $\mathbf{c} \in \mathbf{C}$. In particular, the additional dimension is irrelevant for orthogonality, since we choose at least one vector in \mathbf{A} and any such vector has the last coordinate equal to 0. For any $\ell \in [A^{k_1}]$, we denote by $\mathbf{a}(\ell)$ the ℓ -th vector in \mathbf{A} .

Tuple gadgets. For any $\mathbf{b} \in \mathbf{B}$, $\mathbf{c} \in \mathbf{C}$, we define vectors $u_{\mathbf{b}} \in \{0, 1\}^{dA^{k_1}}$ and $v_{\mathbf{c}} \in \{0, 1\}^{(d-1)A^{k_1}+1}$:

$$\begin{aligned} u_{\mathbf{b}} &:= (\mathbf{a}(1)[1] \cdot \mathbf{b}[1], \dots, \mathbf{a}(A^{k_1})[1] \cdot \mathbf{b}[1], \dots, \mathbf{a}(1)[d] \cdot \mathbf{b}[d], \dots, \mathbf{a}(A^{k_1})[d] \cdot \mathbf{b}[d]) \\ v_{\mathbf{c}} &:= (\mathbf{c}[1], \underbrace{0, \dots, 0}_{A^{k_1}-1 \text{ times}}, \mathbf{c}[2], \dots, \underbrace{0, \dots, 0}_{A^{k_1}-1 \text{ times}}, \mathbf{c}[d]). \end{aligned}$$

In other words, for $j \in [d]$ and $\ell \in [A^{k_1}]$ we have $(u_{\mathbf{b}})_{j \cdot d + \ell} = \mathbf{a}(\ell)[j + 1] \cdot \mathbf{b}[j + 1]$ as well as $(v_{\mathbf{c}})_{j \cdot d + \ell} = \mathbf{c}[j + 1]$ if $\ell = 1$ and $(v_{\mathbf{c}})_{j \cdot d + \ell} = 0$ otherwise.

The key idea is as follows. Consider a structured alignment $\Lambda = \{(\Delta + 1, 1), \dots, (\Delta + m, m)\} \in \mathcal{S}_{n,m}$ for the above vectors, where $n = dA^{k_1}$ and $m = (d-1)A^{k_1} + 1$. This chooses some tuple $\mathbf{a}(\Delta + 1) \in \mathbf{A}$ and aligns the pairs $(\mathbf{a}(\Delta + 1)[\ell] \cdot \mathbf{b}[\ell], \mathbf{c}[\ell])$ for all $\ell \in [d]$, additional to some trivial pairs where the coordinate of $v_{\mathbf{c}}$ is 0. This allows us to determine whether $(\mathbf{a}(\Delta + 1), \mathbf{b}, \mathbf{c})$ is orthogonal.

To formalize this, create $\tilde{u}_{\mathbf{b}}$ by replacing each 0- and 1-entry in $u_{\mathbf{b}}$ by $\mathbf{0}_X$ and $\mathbf{1}_X$ (from Definition 3.15), and create $\tilde{v}_{\mathbf{c}}$ by replacing each 0- and 1-entry in $v_{\mathbf{c}}$ by $\mathbf{0}_Y$ and $\mathbf{1}_Y$, respectively. Let t_X and t_Y be the types of $\mathbf{0}_X, \mathbf{1}_X$ and $\mathbf{0}_Y, \mathbf{1}_Y$, respectively. Set $\delta_0 := \delta(\mathbf{0}_X, \mathbf{0}_Y) = \delta(\mathbf{0}_X, \mathbf{1}_Y) = \delta(\mathbf{1}_X, \mathbf{0}_Y)$ and $\delta_1 := \delta(\mathbf{1}_X, \mathbf{1}_Y)$. We define the *tuple gadgets*

$$\begin{aligned} \text{TG}_X(\mathbf{b}) &:= \text{GA}^{(d-1)A^{k_1}+1, t_Y}(\tilde{u}_{\mathbf{b}}), \\ \text{TG}_Y(\mathbf{c}) &:= \text{GA}^{dA^{k_1}, t_X}(\tilde{v}_{\mathbf{c}}). \end{aligned}$$

Let t'_X, t'_Y denote the types of $\text{TG}_X(\mathbf{b})$, $\text{TG}_Y(\mathbf{c})$, and let C be the number obtained from Definition 3.13 when creating $\text{TG}_X(\mathbf{b})$, $\text{TG}_Y(\mathbf{c})$. Note that t'_X, t'_Y , and C do not depend on the choice of $\mathbf{b} \in \mathbf{B}, \mathbf{c} \in \mathbf{C}$.

Claim 3.17. *Let $\mathbf{b} \in \mathbf{B}, \mathbf{c} \in \mathbf{C}$ and set $n := dA^{k_1}$ and $m := (d-1)A^{k_1} + 1$. If there exists $\mathbf{a} \in \mathbf{A}$ such that $(\mathbf{a}, \mathbf{b}, \mathbf{c})$ are orthogonal, then $\delta(\text{TG}_X(\mathbf{b}), \text{TG}_Y(\mathbf{c})) = C + m \cdot \delta_0$. Otherwise $\delta(\text{TG}_X(\mathbf{b}), \text{TG}_Y(\mathbf{c})) \geq C + (m-1)\delta_0 + \delta_1$.*

Proof. If there is an $\mathbf{a} \in \mathbf{A}$ for which $(\mathbf{a}, \mathbf{b}, \mathbf{c})$ are orthogonal, let Δ be such that $\mathbf{a} = \mathbf{a}(\Delta + 1)$, where $\mathbf{a}(\ell)$ is the ℓ -th tuple in the lexicographic ordering of \mathbf{A} . The structured alignment $\Lambda = \{(\Delta + 1, 1), \dots, (\Delta + m, m)\}$ satisfies

$$\text{cost}(\Lambda) = \left(\sum_{\ell=1}^d \delta_{\mathbf{a}(\Delta+1)[\ell] \cdot \mathbf{b}[\ell] \cdot \mathbf{c}[\ell]} \right) + (A^{k_1} - 1)(d-1)\delta_0 = m \cdot \delta_0.$$

Furthermore, for any $\Lambda \in \mathbf{\Lambda}_{n,m}$, we have $\text{cost}(\Lambda) \geq m \cdot \delta_0$, since $\text{cost}(\Lambda)$ contains at least m summands of value at least $\min\{\gamma, \min_{i,j} \delta(X_i, Y_j)\} = \min_{i,j} \delta(X_i, Y_j) \geq \delta_0$. Thus $\delta(\text{TG}_X(\mathbf{b}), \text{TG}_Y(\mathbf{c})) = C + m \cdot \delta_0$ by Definition 3.13.

Otherwise, if no such \mathbf{a} exists, let $\Lambda \in \mathbf{\Lambda}_{n,m}$ be arbitrary. If $\Lambda = \{(\Delta + 1, 1), \dots, (\Delta + m, m)\}$ is a structured alignment, then

$$\text{cost}(\Lambda) = \left(\sum_{\ell=1}^d \delta_{\mathbf{a}(\Delta+1)[\ell] \cdot \mathbf{b}[\ell] \cdot \mathbf{c}[\ell]} \right) + (A^{k_1} - 1)(d-1)\delta_0 \geq (m-1) \cdot \delta_0 + \delta_1,$$

since there exists some $\ell \in [d]$ with $\mathbf{a}(\Delta + 1)[\ell] = \mathbf{b}[\ell] = \mathbf{c}[\ell] = 1$ which contributes a value of δ_1 .

If $\Lambda = \{(i_1, j_1), \dots, (i_{|\Lambda|}, j_{|\Lambda|})\}$ is unstructured, then either $|\Lambda| < m$, in which case we have

$$\text{cost}(\Lambda) \geq |\Lambda| \cdot \delta_0 + (m - |\Lambda|)\gamma \geq (m - 1)\delta_0 + \delta_1,$$

or $|\Lambda| = m$ and $i_m - i_1 > m - 1$, and thus

$$\text{cost}(\Lambda) \geq m\delta_0 + (i_m - i_1 - (m - 1))\gamma \geq (m - 1)\delta_0 + \delta_1.$$

Thus by Definition 3.13, $\delta(\text{TG}_X(\mathbf{b}), \text{TG}_Y(\mathbf{c})) \geq C + (m - 1)\delta_0 + \delta_1$. \square

Normalization. As usual in these kinds of reductions, we need a normalization step. We define a normalization sequence as

$$\text{TG}_{\text{norm}} = \text{GA}^{(d-1)A^{k_1}+1, t_Y}(\underbrace{\mathbf{0}_X, \dots, \mathbf{0}_X}_{(d-1)A^{k_1} \text{ times}}, \underbrace{\mathbf{1}_X, \dots, \mathbf{1}_X}_{A^{k_1} \text{ times}}).$$

Claim 3.18. For any $\mathbf{c} \in \mathbf{C}$, we have $\delta(\text{TG}_{\text{norm}}, \text{TG}_Y(\mathbf{c})) = C + (m - 1)\delta_0 + \delta_1$.

Proof. Let $n = dA^{k_1}$ and $m = (d - 1)A^{k_1} + 1$. Let $\Lambda = \{(\Delta + 1, 1), \dots, (\Delta + m, m)\} \in \mathcal{S}_{n,m}$ be a structured alignment. Then by construction of TG_{norm} and $\text{TG}_Y(\mathbf{c})$, the only pair corresponding to $\mathbf{1}_X$ and possibly $\mathbf{1}_Y$ entries is the pair $(\Delta + m, m)$, since only the last A^{k_1} entries of TG_{norm} are $\mathbf{1}_X$, and the only possible $\mathbf{1}_Y$ -entry of $\text{TG}_Y(\mathbf{c})$ that could be aligned with one of them is its final entry. Now we use that we constructed the vectors \mathbf{C} as $\mathcal{A}_1^{(k_2)}$, i.e., we augmented all vectors by a d -th coordinate 1, which implies that the m -th entry of $\text{TG}_Y(\mathbf{c})$ is indeed $\mathbf{1}_Y$. Hence, the pair $(\Delta + m, m)$ contributes a distance of δ_1 while all others contribute δ_0 . This yields $\text{cost}(\Lambda) = (m - 1)\delta_0 + \delta_1$.

Let $\Lambda \in \mathbf{\Lambda}_{n,m} \setminus \mathcal{S}_{n,m}$ be an unstructured alignment. Then its cost is at least $\text{cost}(\Lambda) \geq (m - 1)\delta_0 + \gamma \geq (m - 1)\delta_0 + \delta_1$, since it contains at least $m - 1$ summands of value $\min\{\gamma, \min_{i,j} \delta(X_i, Y_j)\} \geq \delta_0$ and at least one punishment term $\gamma \geq \delta_1$ for a deviation from a structured assignment. Thus by Definition 3.13, we have $\delta(\text{TG}_{\text{norm}}, \text{TG}_Y(\mathbf{c})) = C + (m - 1)\delta_0 + \delta_1$. \square

We now define for any $\mathbf{b} \in \mathbf{B}$, $\mathbf{c} \in \mathbf{C}$ the *normalized tuple gadgets*

$$\begin{aligned} \text{NTG}_X(\mathbf{b}) &:= \text{GA}^{1, t'_Y}(\text{TG}_X(\mathbf{b}), \text{TG}_{\text{norm}}), \\ \text{NTG}_Y(\mathbf{c}) &:= \text{GA}^{2, t'_X}(\text{TG}_Y(\mathbf{c})). \end{aligned}$$

We let t''_X, t''_Y denote the resulting types of $\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})$, and C' be the number obtained from Definition 3.13 when creating $\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})$. Note that t''_X, t''_Y , and C' do not depend on the choice of $\mathbf{b} \in \mathbf{B}, \mathbf{c} \in \mathbf{C}$. This definitions satisfies the following properties.

Claim 3.19. Let $\mathbf{b} \in \mathbf{B}, \mathbf{c} \in \mathbf{C}$. If there exists $\mathbf{a} \in \mathbf{A}$ such that $(\mathbf{a}, \mathbf{b}, \mathbf{c})$ are orthogonal, then $\delta(\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})) = \delta_{\text{orth}}$, otherwise we have $\delta(\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})) = \delta_{\text{non}}$, where

$$\begin{aligned} \delta_{\text{orth}} &:= C' + C + ((d - 1)A^k + 1) \cdot \delta_0, \\ \delta_{\text{non}} &:= C' + C + (d - 1)A^k \cdot \delta_0 + \delta_1. \end{aligned}$$

Proof. We check all possible alignments $\Lambda \in \mathbf{\Lambda}_{2,1}$: If $\Lambda = \{(1, 1)\}$, then $\text{cost}(\Lambda) = \delta(\text{TG}_X(\mathbf{b}), \text{TG}_Y(\mathbf{c}))$. If $\Lambda = \{(2, 1)\}$, we have $\text{cost}(\Lambda) = \delta(\text{TG}_{\text{norm}}, \text{TG}_Y(\mathbf{c}))$. For the only unstructured alignment $\Lambda = \emptyset$, we have $\text{cost}(\Lambda) = \gamma \geq \max\{\delta(\text{TG}_X(\mathbf{b}), \text{TG}_Y(\mathbf{c})), \delta(\text{TG}_{\text{norm}}, \text{TG}_Y(\mathbf{c}))\}$. Thus by Definition 3.13, we have $\delta(\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})) = C' + \min\{\delta(\text{TG}_X(\mathbf{b}), \text{TG}_Y(\mathbf{c})), \delta(\text{TG}_X(\mathbf{b}), \text{TG}_{\text{norm}})\}$. The claim now follows from Claims 3.17 and 3.18. \square

Final construction. To obtain our final instance, we enumerate all $\mathbf{b}(1), \dots, \mathbf{b}(A^{k_2}) \in \mathbf{B}$ and $\mathbf{c}(1), \dots, \mathbf{c}(A^{k_2}) \in \mathbf{C}$ in an arbitrary fashion. We finally combine their corresponding normalized tuple gadgets by defining

$$\begin{aligned} X &:= \text{GA}^{A^{k_2}, t''_Y}(\text{NTG}_X(\mathbf{b}(1)), \dots, \text{NTG}_X(\mathbf{b}(A^{k_2})), \text{NTG}_Y(\mathbf{b}(1)), \dots, \text{NTG}_Y(\mathbf{b}(A^{k_2}))), \\ Y &:= \text{GA}^{2A^{k_2}, t''_X}(\text{NTG}_Y(\mathbf{c}(1)), \dots, \text{NTG}_Y(\mathbf{c}(A^{k_2}))). \end{aligned}$$

Let C'' be the number obtained from Definition 3.13 when creating X, Y .

Claim 3.20. *We have $\delta(X, Y) \leq C'' + (A^{k_2} - 1)\delta_{\text{non}} + \delta_{\text{orth}}$ if and only if there are $\mathbf{a} \in \mathbf{A}, \mathbf{b} \in \mathbf{B}, \mathbf{c} \in \mathbf{C}$ such that $(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is orthogonal.*

Proof. Assume that there exists an orthogonal set of vectors and let $\mathbf{a} \in \mathbf{A}, \mathbf{b}(i) \in \mathbf{B}, \mathbf{c}(j) \in \mathbf{C}$ be the vectors representing them. Let $n = 2A^{k_2}$ and $m = A^{k_2}$. If $i \geq j$, we consider the structured alignment $\Lambda = \{(i-j+1, 1), \dots, (i-j+m, m)\}$. Then Λ aligns $\text{NTG}_X(\mathbf{b}(i)), \text{NTG}_Y(\mathbf{c}(j))$, yielding cost δ_{orth} by Claim 3.19. Since $\delta(\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})) \leq \delta_{\text{non}}$ for any \mathbf{b}, \mathbf{c} , we conclude that $\text{cost}(\Lambda) \leq (m-1)\delta_{\text{non}} + \delta_{\text{orth}}$. Similarly, if $i < j$, we define the structured alignment $\Lambda = \{(n+i-j+1, 1), \dots, (n+i-j+m, m)\}$. Then, again, Λ aligns $\text{NTG}_X(\mathbf{b}(i)), \text{NTG}_Y(\mathbf{c}(j))$. As before, we obtain $\text{cost}(\Lambda) \leq (m-1)\delta_{\text{non}} + \delta_{\text{orth}}$. Thus, in both cases Definition 3.13 yields $\delta(X, Y) \leq C'' + (m-1)\delta_{\text{non}} + \delta_{\text{orth}}$.

To prove the claim, it remains to prove that $\delta(X, Y) \geq C'' + m\delta_{\text{non}}$ if all choices of vectors are non-orthogonal. Note that for any $\Lambda \in \mathbf{\Lambda}_{n,m}$, $\text{cost}(\Lambda)$ consists of m summands with a value of at least $\min_{i,j} \delta(\text{NTG}_X(\mathbf{b}(i)), \text{NTG}_Y(\mathbf{c}(j))) \geq \delta_{\text{non}}$. This concludes the claim by Definition 3.13. \square

It remains to prove bounds on the lengths and compressed sizes of the constructed strings.

Claim 3.21. *The strings X, Y have length $O(dA^{k_1+k_2})$. We can, in linear time in the output size, compute SLPs \mathcal{X}, \mathcal{Y} for X, Y of size $O(dA^{k_2+1})$.*

Proof. We will frequently make use of the compressibility of the alignment gadget (Definition 3.14). We start by constructing an SLP $\mathcal{TG}_X(\mathbf{b})$ for $\text{TG}_X(\mathbf{b})$ for any $\mathbf{b} \in [A]^{k_2}$. Note that we can split $\text{TG}_X(\mathbf{b})$ into $\text{TG}_X(\mathbf{b})_L \circ (\bigcirc_{i=1}^{dA^{k_1}} \text{pad}_X(\tilde{u}_i)) \circ \text{TG}_X(\mathbf{b})_R$. We can apply Lemma 3.8 by observing that

$$\bigcirc_{i=1}^{dA^{k_1}} \text{pad}_X(\tilde{u}_i) = \text{tuplify}(\mathcal{A}_0, k_1, \mathbf{b}, \text{pad}_X(\mathbf{0}_X), \text{pad}_X(\mathbf{1}_X)), \quad (3)$$

where $S^X(0) := \text{pad}_X(\mathbf{0}_X)$ and $S^X(1) := \text{pad}_X(\mathbf{1}_X)$. Since $\mathbf{0}_X, \mathbf{1}_X$ are of constant size, we can compute SLPs $\mathcal{S}(0), \mathcal{S}(1)$ for $S^X(0), S^X(1)$ of size $O(1)$ by the compressibility assumption. Thus we can compute an SLP for (3) of size $O(dA)$. Since the left and right bounding string of $\text{TG}_X(\mathbf{b})$ have SLPs of size $O(\log A)$, we obtain an SLP $\mathcal{TG}_X(\mathbf{b})$ for $\text{TG}_X(\mathbf{b})$ of size $O(dA)$, while $|\text{TG}_X(\mathbf{b})| = \Theta(dA^{k_1})$.

To compute an SLP $\mathcal{TG}_Y(\mathbf{c})$ for $\text{TG}_Y(\mathbf{c})$ for any $\mathbf{c} \in \mathbf{C}$, we note that

$$\text{TG}_Y(\mathbf{c}) = \text{TG}_Y(\mathbf{c})_L \circ \left(\bigcirc_{i=1}^{(d-1)A^{k_1+1}} \text{pad}_Y(\tilde{v}_i) \right) \circ \text{TG}_Y(\mathbf{c})_R,$$

where

$$\bigcirc_{i=1}^{(d-1)A^{k_1+1}} \text{pad}_Y(\tilde{v}_i) = \left(\bigcirc_{\ell=1}^{d-1} S^Y(\mathbf{c}[\ell]) \circ S^Y(0)^{A^{k_1-1}} \right) \circ S^Y(\mathbf{c}[d]),$$

where $S^Y(0) := \text{pad}_Y(\mathbf{0}_Y)$ and $S^Y(1) := \text{pad}_Y(\mathbf{1}_Y)$. This immediately admits an SLP of size $O(d + \log A)$ by Observation 2.2. Again, using SLPs of size $O(\log A)$ for $\text{TG}_Y(\mathbf{c})_L, \text{TG}_Y(\mathbf{c})_R$, we obtain an SLP $\mathcal{TG}_Y(\mathbf{c})$ for $\text{TG}_Y(\mathbf{c})$ of size $O(d + \log A)$, while $|\text{TG}_Y(\mathbf{c})| = O(dA^{k_1})$.

In the construction of $\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})$ we use constant n, m . Together with the compressibility of the alignment gadget, we obtain SLPs $\mathcal{NTG}_X(\mathbf{b}), \mathcal{NTG}_Y(\mathbf{c})$ for $\text{NTG}_X(\mathbf{b}), \text{NTG}_Y(\mathbf{c})$ of size $O(|\mathcal{TG}_X(\mathbf{b})|), O(|\mathcal{TG}_Y(\mathbf{c})|)$. Furthermore, $|\text{NTG}_X(\mathbf{b})| = \Theta(|\text{TG}_X(\mathbf{b})|), |\text{NTG}_Y(\mathbf{c})| = \Theta(|\text{TG}_Y(\mathbf{c})|)$.

Finally, to obtain SLPs \mathcal{X}, \mathcal{Y} for X, Y , we use a final application of the compressibility of the alignment gadget. This yields $|\mathcal{X}| = O(\log A + \sum_{\mathbf{b} \in \mathcal{B}} |\mathcal{NTG}_X(\mathbf{b})|) = O(dA^{k_2+1})$ and $|\mathcal{Y}| = O(\log A + \sum_{\mathbf{c} \in \mathcal{C}} |\mathcal{NTG}_Y(\mathbf{c})|) = O(A^{k_2}(d + \log A))$. Note that $|X|, |Y| = \Theta(dA^{k_1+k_2})$. It is easy to verify that constructing \mathcal{X}, \mathcal{Y} takes time $O(dA^{k_1+k_2})$. \square

We are now ready to prove the theorem.

Proof of Theorem 3.16. Let $0 < \alpha_n < 1$ and set $\beta := \frac{\alpha_n}{1+\alpha_n}$. Let $k \geq 2$ and let \mathcal{A} be a k -OV instance with A vectors in dimension d . We split $k = k_1 + 2k_2$ with $k_1, k_2 \geq 1$ and $k_2 \approx \beta k$ and $k_1 \approx (1-2\beta)k$. Note that k_1, k_2 are restricted to be integers, however, for any $\varepsilon > 0$ and sufficiently large k depending only on ε and α_n we can ensure $k_2 + 1 \leq (1-\varepsilon/8)\beta k$ and $k_1 \leq (1+\varepsilon/4)(1-2\beta)k$. Since $k = k_1 + 2k_2$, it follows that $k_1 + k_2 \geq (1-\beta)k$. Note that for the dimension d we can assume $d \leq A$, since otherwise an $O(A^{k-\varepsilon} \text{poly}(d))$ algorithm clearly exists. In particular, for sufficiently large k we have $d \leq A^{(\varepsilon/8) \cdot \min\{\beta, 1-\beta\}k}$. By Claim 3.21, the constructed strings X, Y have length N bounded from above by $\Theta(dA^{k_1+k_2}) = O(dA^{(1+\varepsilon/4)(1-\beta)k}) = O(A^{(1+\varepsilon/2)(1-\beta)k})$ and bounded from below by $\Theta(dA^{k_1+k_2}) = \Omega(A^{(1-\beta)k})$. The constructed SLPs have size $n = O(dA^{k_2+1}) = O(dA^{(1-\varepsilon/8)\beta k}) = O(A^{\beta k})$. Since $\beta/(1-\beta) = \alpha_n$, it follows that $n = O(N^{\alpha_n})$, and by partially decompressing the SLPs we can ensure the desired $n = \Theta(N^{\alpha_n})$, while keeping $n = O(A^{(1+\varepsilon/2)\beta k})$. By Claim 3.20, computing $\delta(X, Y)$ allows us to decide feasibility of the given k -OV instance. Hence, any $O((nN)^{1-\varepsilon})$ time algorithm for $\delta(\cdot, \cdot)$ in the setting $n = \Theta(N^{\alpha_n})$ would yield an algorithm for k -OV in time $O((A^{(1+\varepsilon/2)k})^{1-\varepsilon}) = O(A^{(1-\varepsilon/2)k})$, contradicting the k -OV conjecture. \square

3.3.3 Extended Alignment Gadget for LCS

In this section, we fix the distance measure to be the LCS distance $\delta(X, Y) = |X| + |Y| - 2 \cdot L(X, Y)$, where $L(X, Y)$ denotes the length of an LCS S of X and Y . Note that $\delta(X, Y)$ counts the number of symbols to be deleted in X to obtain S plus the number of symbols to be deleted in Y to obtain S . We show that δ admits coordinate values and a compressible extended alignment gadget. Together with Theorem 3.16, this will yield our conditional lower bound for LCS.

We make use of the same coordinate values as in [17].

Lemma 3.22 ([17, Lemma V.2]). *LCS admits coordinate values by setting*

$$\mathbf{1}_X := 11100, \mathbf{0}_X := 10011, \mathbf{1}_Y := 00111, \mathbf{0}_Y := 11001.$$

These strings have type $(5, \{0, 1\})$.

It remains to implement a compressible extended alignment gadget. Let us first disregard compressibility.

Lemma 3.23. *The following construction implements an extended alignment gadget: Let X_1, \dots, X_n of length ℓ_X and Y_1, \dots, Y_m of length ℓ_Y be strings over Σ . We introduce new symbols $\sigma, \rho, \mu \notin \Sigma$, define $\kappa_1 := 4(\ell_X + \ell_Y)$ and $\kappa_2 := 2\kappa_1 + \ell_X$, and set*

$$G(S) := \sigma^{\kappa_1} S \rho^{\kappa_1},$$

The alignment gadget is now defined as

$$\begin{aligned} X &= G(X_1) Z_1^X G(X_2) \dots Z_{n-1}^X G(X_n), \\ Y &= L^Y G(Y_1) Z_1^Y G(Y_2) \dots Z_{m-1}^Y G(Y_m) R^Y, \end{aligned}$$

where $Z_i^X = Z_j^Y = \mu^{\kappa_2}$ for $i \in [n-1], j \in [m-1]$ and $L^Y = R^Y = \mu^{n\kappa_2}$. This satisfies property (2) of Definition 3.13 with $C := 2n\kappa_2$.

Proof. To analyze our alignment gadget construction (adapting the proof of the LCS gadget of the full version of [17]), we prepare some useful facts.

Claim 3.24 ([17, Fact V.7]). *Let X and Z_1, \dots, Z_k be strings. Set $Z := Z_1 \circ \dots \circ Z_k$. We have*

$$\delta(X, Z) = \min_{X(Z_1), \dots, X(Z_k)} \sum_{j=1}^k \delta(X(Z_j), Z_j),$$

where $X(Z_1), \dots, X(Z_k)$ range over all ordered partitions of X into k substrings, i.e., $X(Z_1) = x[i_0 + 1..i_1], X(Z_2) = x[i_1 + 1..i_2], \dots, X(Z_k) = x[i_{k-1} + 1..i_k]$ for any $0 = i_0 \leq i_1 \leq \dots \leq i_k = |X|$.

Claim 3.25. *Let U, V be strings over Σ , $\alpha \in \Sigma$ and $k \in \mathbb{N}_0$. Then we have*

$$(i) \delta(U, V) \geq ||U| - |V||,$$

$$(ii) \delta(\alpha^k U, \alpha^k V) = \delta(U, V),$$

$$(iii) \text{ Let } W \text{ be a string not containing } \alpha. \text{ Then } \delta(W\alpha U, \alpha^k V) \geq \min\{k, \delta(\alpha U, \alpha^k V)\}.$$

We obtain symmetric statements by reversing all involved strings.

Proof. (i) Suppose $|U| \geq |V|$, then at least $|U| - |V|$ many symbols must be deleted in U . The claim follows by symmetry.

(ii) It suffices to show the claim for $k = 1$, then the general statement follows by induction. Consider a LCS S of $(\alpha U, \alpha V)$. At least one α is matched in S , as otherwise we can extend S by matching both α 's. If exactly one α is matched in S , then the other α is free, so we may instead match the two α 's. Thus, without loss of generality a LCS of $(\alpha U, \alpha V)$ matches the two α 's. This yields $L(\alpha U, \alpha V) = 1 + L(U, V)$. Hence, $\delta(\alpha U, \alpha V) = |\alpha U| + |\alpha V| - 2L(\alpha U, \alpha V) = |U| + |V| - 2L(U, V) = \delta(U, V)$.

(iii) Fix an LCS S of $W\alpha U$ and $\alpha^k V$. If S starts with a symbol other than α , then S cannot use any symbol from the α^k -prefix of $\alpha^k V$, i.e., the α^k -prefix has to be deleted and thus $\delta(W\alpha U, \alpha^k V) \geq k$. Otherwise, if S starts with an α , then S cannot use any symbol from W (which is a string over $\Sigma \setminus \{\alpha\}$), i.e., S is an LCS of αU and $\alpha^k V$. Thus $\delta(W\alpha U, \alpha^k V) = |W\alpha U| + |\alpha^k V| - 2L(\alpha U, \alpha^k V) = |W| + \delta(\alpha U, \alpha^k V)$ and the claim follows. \square

Claim 3.26. *Let $\ell \geq 0$. For any prefix X' of X we have $\delta(X', \mu^\ell) \geq \ell$. Moreover, if X' is of the form $G(X_1)Z_1^X \dots G(X_i)Z_i^X$ for some $0 \leq i < n$ and $\ell \geq i \cdot \kappa_2$, then $\delta(X', \mu^\ell) = \ell$. Symmetric statements hold for any suffix of X .*

Proof. Note that for any $i \in [n]$ the string $G(X_i)Z_i^X$ contains $|Z_i^X| = \kappa_2$ many μ 's and $|G(X_i)| = 2\kappa_1 + \ell_X = \kappa_2$ many non- μ 's. Furthermore, any prefix of $G(X_i)Z_i^X$ contains at least as many non- μ 's as μ 's. Hence, the LCS of X' and μ^ℓ has a length of at most $|X'|/2$. This yields $\delta(X', \mu^\ell) = |X'| + |\mu^\ell| - 2L(X', \mu^\ell) \geq \ell$. If X' is of the form $G(X_1)Z_1^X \dots G(X_i)Z_i^X$ and μ^ℓ has at least $|X'|/2 = i\kappa_2$ many μ 's, we have equality. \square

We now prove that our construction yields an extended alignment gadget. We start with the upper bound of property (2), i.e., $\delta(X, Y) \leq 2n\kappa_2 + \min_{\Lambda \in \mathcal{S}_{n,m}} \text{cost}(\Lambda)$.

Let $\Lambda = \{(\Delta+1, 1), \dots, (\Delta+m, m)\}$ be a structured alignment and consider an ordered partition of X as in Claim 3.24 defined as follows:

$$\begin{aligned} X(G(Y_j)) &:= G(X_{\Delta+j}) \quad \text{for } j \in [m], \\ X(Z_j^Y) &:= Z_{\Delta+j}^X \quad \text{for } j \in [m-1], \\ X(L^Y) &:= G(X_1)Z_1^X \dots G(X_\Delta)Z_\Delta^X, \\ X(R^Y) &:= Z_{\Delta+m}^X G(X_{\Delta+m+1}) \dots Z_{n-1}^X G(X_n). \end{aligned}$$

Claim 3.24 thus yields

$$\delta(X, Y) \leq \delta(X(L^Y), L^Y) + \delta(X(R^Y), R^Y) + \sum_{j=1}^m \delta(G(X_{\Delta+j}), G(Y_j)) + \sum_{j=1}^{m-1} \delta(Z_{\Delta+j}^X, Z_j^Y).$$

By Claim 3.26, we obtain $\delta(X(L^Y), L^Y) = n\kappa_2$ and symmetrically, $\delta(X(R^Y), R^Y) = n\kappa_2$. Trivially, $\delta(Z_{\Delta+j}^X, Z_j^Y) = 0$. Finally, by matching the padding around X_i, Y_j in $G(X_i), G(Y_j)$, we obtain $\delta(G(X_{\Delta+j}), G(Y_j)) = \delta(X_{\Delta+j}, Y_j)$ by Claim 3.25(ii). Summing up all contributions, we obtain

$$\delta(X, Y) \leq 2n\kappa_2 + \sum_{(i,j) \in \Lambda} \delta(X_i, Y_j),$$

which holds for an arbitrary $\Lambda \in \mathcal{S}_{n,m}$, thus concluding the upper bound.

It remains to prove the lower bound of property (2), i.e., $\delta(X, Y) \geq 2n\kappa_2 + \min_{\Lambda \in \mathbf{\Lambda}_{n,m}} \text{cost}(\Lambda)$. Set $M^Y = G(Y_1)Z_1^Y \dots Z_{m-1}^Y G(Y_m)$. Using Claim 3.24, we let $X(L^Y)$, $X(M^Y)$ and $X(R^Y)$ be an ordered partition of X such that

$$\delta(X, Y) = \delta(X(L^Y), L^Y) + \delta(X(M^Y), M^Y) + \delta(X(R^Y), R^Y).$$

Since $L^Y = \mu^{n\kappa_2}$ and $X(L^Y)$ is a prefix of X , by Claim 3.26 we have $\delta(X(L^Y), L^Y) \geq n\kappa_2$, and similarly we get $\delta(X(R^Y), R^Y) \geq n\kappa_2$. It remains to construct an alignment $\Lambda \in \mathbf{\Lambda}_{n,m}$ satisfying

$$\text{cost}(\Lambda) \leq \delta(X(M^Y), M^Y), \tag{4}$$

then together we have shown the desired inequality $\delta(X, Y) \geq 2n\kappa_2 + \min_{\Lambda \in \mathbf{\Lambda}_{n,m}} \text{cost}(\Lambda)$.

As in Claim 3.24, we let $X(G(Y_j))$ for $j \in [m]$ and $X(Z_j^Y)$ for $j \in [m-1]$ be an ordered partition of $X(M^Y)$ such that

$$\delta(X(M^Y), M^Y) = \sum_{j=1}^m \delta(X(G(Y_j)), G(Y_j)) + \sum_{j=1}^{m-1} \delta(X(Z_j^Y), Z_j^Y).$$

Let $\mu(U)$ be the number of μ 's in a string U and let $\delta_{\text{del}-\mu}(U, V)$ denote the LCS distance of U and V after deleting all μ 's in U and V . Clearly, since $|\mu(U) - \mu(V)|$ μ 's have to be deleted in any LCS, we have

$$\delta(X(M^Y), M^Y) \geq \left(\sum_{j=1}^m \delta_{\text{del}-\mu}(X(G(Y_j)), G(Y_j)) \right) + |\mu(U) - \mu(V)|. \quad (5)$$

Let us construct an alignment Λ satisfying (4). For any $j \in [m]$, if $X(G(Y_j))$ contains more than half of some $X_{i'}$ (which is part of $G(X_{i'})$), then let i be the leftmost such index and align i and j . Note that the set Λ of all these aligned pairs (i, j) is a valid alignment in $\mathbf{\Lambda}_{n,m}$, since no X_i or Y_j can be aligned more than once.

We prove the following claims:

Claim 3.27. *For any aligned pair $(i, j) \in \Lambda$, we have $\delta_{\text{del}-\mu}(X(G(Y_j)), G(Y_j)) \geq \delta(X_i, Y_j)$.*

Proof. Let U be $X(G(Y_j))$ with all μ 's deleted (note that $G(Y_j)$ contains no μ 's). We will prove $\delta(U, G(Y_j)) \geq \delta(X_i, Y_j)$. Recall that $X(G(Y_j))$ contains more than half of X_i , thus so does U . If $||U| - |G(Y_j)|| \geq \ell_x + \ell_y$, then we have $\delta(U, G(Y_j)) \geq \ell_x + \ell_y \geq \delta(X_i, Y_j)$ by Claim 3.25(i). Since $|G(Y_j)| = 2\kappa_1 + \ell_y$, we may hence assume $2\kappa_1 - \ell_x \leq |U| \leq 2\kappa_1 + \ell_x + 2\ell_y$.

We distinguish three cases: Either U contains X_i fully (C1), or at least its right half but not fully (C2), or at least its left half but not fully (C3).

In case (C2), U is of the form $X'_i \rho^{\kappa_1} \sigma^a X'_{i+1} \rho^b$ where X'_i is a suffix of X_i , $a \leq \kappa_1$, X'_{i+1} is a prefix of X_{i+1} and $b \leq 2\ell_y$. In this case, by Claim 3.25(iii) with $\alpha = \sigma$ and $W = X'_i \rho^{\kappa_1}$, we have $\delta(U, G(Y_j)) \geq \min\{\kappa_1, \delta(\sigma^a X'_{i+1} \rho^b, \sigma^{\kappa_1} Y_j \rho^{\kappa_1})\}$. Note that since the second string contains κ_1 ρ 's and the first string contains less than $2\ell_y$ ρ 's, we have $\delta(\sigma^a X'_{i+1} \rho^b, \sigma^{\kappa_1} Y_j \rho^{\kappa_1}) \geq \kappa_1 - 2\ell_y$. Thus $\delta(U, G(Y_j)) \geq \kappa_1 - 2\ell_y \geq \ell_x + \ell_y \geq \delta(X_i, Y_j)$.

The case (C3) is symmetric to (C2).

Finally, in case (C1), U takes one of three forms: either (F1) $\sigma^a X_i \rho^{\kappa_1} \sigma^b X'_{i+1} \rho^c$, where $a \geq 0$, $b \leq \kappa_1$, X'_{i+1} is a (possibly empty) prefix of X_{i+1} and $c \leq 2\ell_y$, or the symmetric version (F2) $X'_{i-1} \rho^b \sigma^{\kappa_1} X_i \rho^a$ with X'_{i-1} a suffix of X_{i-1} and all other parameters as before, or finally (F3) $\rho^a \sigma^{\kappa_1} X_i \rho^{\kappa_1} \sigma^b$ with $a, b \leq 2\ell_y$.

For form (F1), we compute

$$\begin{aligned} \delta(U, G(Y_j)) &= \delta(\sigma^a X_i \rho^{\kappa_1} \sigma^b X'_{i+1}, \sigma^{\kappa_1} Y_j \rho^{\kappa_1 - c}) \\ &\geq \min\{\kappa_1 - c, \delta(\sigma^a X_i \rho^{\kappa_1}, \sigma^{\kappa_1} Y_j \rho^{\kappa_1})\}, \end{aligned}$$

where we used Claim 3.25(ii) in the first line and Claim 3.25(iii) with $\alpha = \rho$ and $W = \sigma^b X'_{i+1}$ in the second line. Note that by Claim 3.25(ii) and by deleting all σ 's only occurring in one string, $\delta(\sigma^a X_i \rho^{\kappa_1}, \sigma^{\kappa_1} Y_j \rho^{\kappa_1}) = (\kappa_1 - a) + \delta(X_i, Y_j) \geq \delta(X_i, Y_j)$. Since $\kappa_1 - c \geq \ell_x + \ell_y \geq \delta(X_i, Y_j)$, the claim follows for (F1). Symmetrically, we can do the same for (F2).

For the final form (F3), we compute, using Claim 3.25(iii) from the left with $\alpha = \sigma$ and $W = \rho^a$ and from the right with $\alpha = \rho$ and $W = \sigma^b$, $\delta(U, G(Y_j)) \geq \min\{\kappa_1, \delta(\sigma^{\kappa_1} X_i \rho^{\kappa_1}, \sigma^{\kappa_1} Y_j \rho^{\kappa_1})\} = \delta(X_i, Y_j)$, where the last equality follows from Claim 3.25(ii) and $\kappa_1 \geq \delta(X_i, Y_j)$. \square

Claim 3.28. *If j is unaligned in Λ , then $\delta_{\text{del}-\mu}(X(G(Y_j)), G(Y_j)) \geq \ell_X + \ell_Y$.*

Proof. Let U be $X(G(Y_j))$ with all μ 's deleted (note that $G(Y_j)$ contains no μ 's). We will prove $\delta(U, G(Y_j)) \geq \ell_X + \ell_Y$. Since $X(G(Y_j))$ contains less than half of any $G(X_i)$, U is of the form $X'_i \rho^a \sigma^b X'_{i+1}$ for some i , a suffix X'_i of X_i , some $a, b \leq \kappa_1$ and a prefix X'_{i+1} of X_{i+1} .

Furthermore using Claim 3.25(iii) with $\alpha = \sigma$ and $W = X'_i \rho^a$, we obtain that $\delta(U, G(Y_j)) \geq \min\{\kappa_1, \delta(\sigma^b X'_{i+1}, \sigma^{\kappa_1} Y_j \rho^{\kappa_1})\}$. Since $\sigma^{\kappa_1} Y_j \rho^{\kappa_1}$ contains κ_1 ρ 's, while $\sigma^b X'_{i+1}$ contains none, we conclude $\delta(U, G(Y_j)) \geq \kappa_1 \geq \ell_X + \ell_Y$. \square

Let us prove (4). If $|\Lambda| < m$, that is, there is an unaligned j , combining the two previous claims with (5) results in

$$\delta(X(M^Y), M^Y) \geq \left(\sum_{(i,j) \in \Lambda} \delta(X_i, Y_j) \right) + (m - |\Lambda|)(\ell_X + \ell_Y) + |\mu(X(M^Y)) - \mu(M^Y)|\kappa_2 \geq \text{cost}(\Lambda),$$

since $\ell_X + \ell_Y \geq \max_{i,j} \delta(X_i, Y_j)$ and $|\mu(X(M^Y)) - \mu(M^Y)| \geq 0$.

Otherwise, if $|\Lambda| = m$, we have $\Lambda = \{(i_1, 1), \dots, (i_m, m)\}$ with $i_1 < i_2 < \dots < i_m$. Note that $X(M^Y)$ is a substring that contains at least half of all X_{i_1}, \dots, X_{i_m} by definition of the alignment Λ . Thus, $\mu(X(M^Y)) \geq (i_m - i_1)\kappa_2$, since it contains all $Z_{i_1}^X, \dots, Z_{i_m-1}^X$. Since $\mu(M^Y) = (m-1)\kappa_2$, we obtain by (5) and Claim 3.28,

$$\delta(X(M^Y), M^Y) \geq \left(\sum_{(i,j) \in \Lambda} \delta(X_i, Y_j) \right) + (i_m - i_1 - m + 1)\kappa_2 \geq \text{cost}(\Lambda),$$

where we used that $\kappa_2 \geq \ell_X + \ell_Y \geq \max_{i,j} \delta(X_i, Y_j)$.

This concludes the proof of Lemma 3.23, showing that our construction yields an extended alignment gadget. \square

It remains to argue that a slight adaption of this gadget is compressible.

Lemma 3.29. *Consider the setting of Lemma 3.23. Adapt the definition of the extended alignment gadget slightly by defining*

$$\begin{aligned} X' &= Z_0^X X Z_n^X = Z_0^X G(X_1) Z_1^X \dots Z_{n-1}^X G(X_n) Z_n^X, \\ Y' &= Z_0^Y Y Z_m^Y = L^Y Z_0^Y G(Y_1) Z_1^Y \dots Z_{m-1}^Y G(Y_m) Z_m^Y R^Y, \end{aligned}$$

where we define the additional blocks $Z_i^X, Z_j^Y = \mu^{\kappa_2}$ with $i \in \{0, n\}, j \in \{0, m\}$. This construction (X', Y') yields a compressible extended alignment gadget.

Proof. By Claim 3.25(ii), we see that $\delta(X', Y') = \delta(X, Y)$, and thus X', Y' satisfies the extended alignment gadget condition (2) of Definition 3.13 by Lemma 3.23.

We define $\text{pad}_X(S) = \text{pad}_Y(S) = \mu^{\kappa_2/2} G(S) \mu^{\kappa_2/2}$ and $X_L = X_R = \mu^{\kappa_2/2}$ and $Y_L = Y_R = \mu^{n\kappa_2 + \kappa_2/2}$. Then we have $X' = X_L (\bigcirc_{i=1}^n \text{pad}_X(X_i)) X_R$ and $Y' = Y_L (\bigcirc_{j=1}^m \text{pad}_X(Y_j)) Y_R$. By

Observation 2.2, we can construct SLPs $\mathcal{X}_L, \mathcal{X}_R, \mathcal{Y}_L, \mathcal{Y}_R$ for X_L, X_R, Y_L, Y_R of size $O(\log n \kappa_2) = O(\log n + \log(\ell_x + \ell_y))$. Likewise, given SLPs $\mathcal{X}_i, \mathcal{Y}_j$ for X_i, Y_j , we can construct SLPs for $\text{pad}_x(X_i)$, $\text{pad}_y(Y_j)$ of size $O(|\mathcal{X}_i| + \log(\ell_x + \ell_y))$, $O(|\mathcal{Y}_j| + \log(\ell_x + \ell_y))$, respectively, as we can generate the paddings $\mu^{\kappa_2/2} \sigma^{\kappa_1}$ and $\rho^{\kappa_1} \mu^{\kappa_2/2}$ around X_i and Y_j using Observation 2.2. This concludes the proof. \square

Our LCS lower bound now follows.

Proof of Theorem 3.12. Since δ admits coordinate values and a compressible extended alignment gadget by Lemmas 3.22 and 3.29, we obtain the claim by the general lower bound of Theorem 3.16, as computing the length of the LCS of X and Y is equivalent to computing $\delta(X, Y)$. \square

4 Tight Bounds Assuming (Combinatorial) k -Clique

In this section we prove matching conditional lower bounds based on the k -Clique conjecture or combinatorial k -Clique conjecture for the following problems:

- NFA Acceptance, i.e., deciding whether a given non-deterministic finite automaton accepts a given string,
- CFG Parsing, i.e., deciding whether a given context-free grammar generates a given string,
- RNA Folding, i.e., computing the maximum number of non-crossing matching pairs of indices in a given string.

See the respective subsections for precise problem definitions.

For NFA Acceptance, the compression used in our proof is extremely simple, in that we only rely on the fact that any repetition T^ℓ can be generated by an SLP of size $O(|T| + \log \ell)$ (Observation 2.2). For CFG Parsing and RNA Folding, our construction is much more subtle. For both problems, we use that the following string and some variants thereof are compressible:

$$S_v := \bigcirc_{u_1, \dots, u_k \in V} [v \text{ is adjacent to every } u_i]$$

That is, we enumerate all k -tuples $(u_1, \dots, u_k) \in V^k$ and for each one check whether all u_i 's are adjacent to a fixed vertex v , writing 1 or 0 depending on this check. This string is generated by an SLP of size $O(V)$: Enumerate all $u_1 \in V$. If u_1 is not adjacent to v , then for all u_2, \dots, u_k the check results in 0, so we can simply write $0^{V^{k-1}}$, which is well compressible by Observation 2.2. Otherwise, if u_1 is adjacent to v , then we can recurse to u_2 , and the following V^{k-1} symbols do not depend on u_1 anymore. More formally, denote by $\text{Repeat}_0^{(d)}$ an SLP generating the string 0^{V^d} . Then with the following SLP rules, for $1 \leq d \leq k$, we have $S_v = \text{eval}(\text{Adj}_v^{(k)})$.

$$\begin{aligned} \text{Adj}_v^{(0)} &\rightarrow 1, \\ \text{Adj}_v^{(d)} &\rightarrow \bigcirc_{u \in V} \begin{cases} \text{Adj}_v^{(d-1)}, & \text{if } \{u, v\} \in E \\ \text{Repeat}_0^{(d-1)}, & \text{otherwise} \end{cases} \end{aligned}$$

Here we use the ‘‘syntactic sugar’’ of having more than two SLP symbols on the right hand side, but clearly this can be converted to a proper SLP of size $O(V)$.

We stress that if in the string S_v we would enumerate only the k -cliques instead of all k -tuples, then S would no longer be easily compressible, since then even the length of a substring depends on the “history” of choosing u_1, \dots, u_{k-d} , and thus the above recursive way of writing S would fail. This demonstrates how subtle our argument is.

Known Lower Bounds from Classic Complexity Theory Plandowski and Rytter [61] showed that deciding whether a given compressed text can be generated by a given CFG is PSPACE-complete. Later, Lohrey [51] showed that this holds even if we restrict the CFG to be fixed (i.e., not part of the input) and deterministic. We observe that the RNA Folding problem is at least as hard as Longest Common Subsequence (see, e.g. [1]). This implies that RNA Folding is PP-hard (see the discussion at the beginning of Section 5.2). Finally, the NFA Acceptance problem can be solved in polynomial $O(nq^\omega)$ time (see below) and previously no conditional lower bounds were known.

4.1 NFA Acceptance

For general notation regarding finite automata, see Section 3.1. Consider the compressed variant of the acceptance problem of nondeterministic finite automata (NFAs).

Problem 4.1 (NFA Acceptance). *We are given a text T of length N by a grammar-compressed representation \mathcal{T} of size n as well as a NFA F with q states, i.e., for any two states z, z' and any symbol $\sigma \in \Sigma$ we are given whether $z \xrightarrow{\sigma} z'$. Decide whether T is accepted by F .*

Note that the input size is $\tilde{O}(n+q^2)$, since we again assume the alphabet size $|\Sigma|$ to be constant.

The naive solution is to decompress \mathcal{T} to obtain T and run the standard acceptance algorithm for NFAs, which takes time $O(|T|q^2) = O(Nq^2)$. Exploiting the compressed setting, one can obtain an $O(nq^\omega)$ -time algorithm [61]: Recall that \mathcal{T} is a set of rules of the form $S_i \rightarrow S_{\ell(i)}S_{r(i)}$ or $S_i \rightarrow \sigma_i$, with $\ell(i), r(i) < i$ and $\sigma_i \in \Sigma$, for $1 \leq i \leq n$. We compute, for increasing i , the state transition matrix A_i , where $(A_i)_{z,z'} = 1$ if we can start in state z , read the string $\text{eval}(S_i)$, and end in state z' , and $(A_i)_{z,z'} = 0$ otherwise. For $S_i \rightarrow S_{\ell(i)}S_{r(i)}$ we can compute A_i as $A_{\ell(i)} \cdot A_{r(i)}$, where \cdot is Boolean matrix multiplication. For $S_i \rightarrow \sigma_i$ we simply have $(A_i)_{z,z'} = 1$ if $z \xrightarrow{\sigma_i} z'$, and 0 otherwise. Hence, A_i can be computed in time $O(q^\omega)$ for every i . The text T is then accepted by F if there is an accepting state z such that $(A_n)_{z_0,z} = 1$, where z_0 is the starting state of F .

Note that this best-known upper bound $O(\min\{nq^\omega, Nq^2\})$ contains “mixed terms” with some factors having exponent ω but not all. Since no standard conjecture contains such mixed terms, we cannot hope to prove a matching lower bound of $\min\{nq^\omega, Nq^2\}^{1-o(1)}$. However, restricting our attention to combinatorial algorithms the best-known running time simplifies to $O(\min\{nq^3, Nq^2\})$, and we can hope to prove a matching lower bound under some assumption on combinatorial algorithms, say for matrix multiplication or k -Clique. For matrix multiplication, the typical issue that we would need to considerably compress the input graph [1] is a barrier for a reduction. Hence, we can only hope to prove a matching lower bound for combinatorial algorithms assuming the k -Clique conjecture. We prove such a result in the following.

Theorem 4.2. *Assuming the combinatorial k -Clique conjecture, there is no combinatorial algorithm for NFA Acceptance in time $O(\min\{nq^3, Nq^2\}^{1-\varepsilon})$ for any $\varepsilon > 0$. This holds even restricted to instances with $n = \Theta(q^{\alpha_n})$ and $N = \Theta(q^{\alpha_N})$ for any $\alpha_N \geq \alpha_n > 0$.*

Proof. Let $k \geq 3$ and let $G = (V, E)$ be a k -Clique instance. In the following, for any $\kappa, \kappa' \geq 1$ with $3\kappa + \kappa' = k$ we will construct an equivalent NFA Acceptance instance with $q = O(V^{\kappa+1} \log V)$, $N = |T| = O(V^{\kappa+\kappa'} \log V)$, and $n = |\mathcal{T}| = O(V^{\kappa'} \log V)$. Note that a combinatorial $O(\min\{nq^3, Nq^2\}^{1-\varepsilon})$ time algorithm for NFA Acceptance then yields a combinatorial algorithm for k -Clique in time $O(V^{(3\kappa+\kappa'+3)(1-\varepsilon)} \log^4 V) = O(V^{(k+4)(1-\varepsilon)})$, which for $k \geq 8/\varepsilon$ is $O(V^{k(1+\varepsilon/2)(1-\varepsilon)}) = O(V^{k(1-\varepsilon/2)})$, contradicting the combinatorial k -Clique conjecture. This yields the desired conditional lower bound. At the end of this proof we will strengthen this statement to even hold for all restrictions $n = \Theta(q^{\alpha n})$ and $N = \Theta(q^{\alpha N})$.

Our construction uses the following gadgets.

Neighborhood Gadgets Let $V = \{v_1, \dots, v_n\}$ and denote by $NG_T(v_i)$ the binary encoding of the number i using $\lceil \log V \rceil$ bits. For any $v \in V$, let $NG_F(v)$ be the NFA that has start state s and target state t , and $|N(v)|$ disjoint directed paths from s to t such that the path corresponding to neighbor $u \in N(v)$ spells $NG_T(u)$. Clearly, we can walk from s to t in $NG_F(v)$ parsing the string $NG_T(u)$ if and only if u is a neighbor of v .

Clique Gadgets For two neighborhood gadgets $NG_F(u), NG_F(v)$ as above, we define their *concatenation* $NG_F(u) \circ NG_F(v)$ as the NFA where we identify the target state t of $NG_F(u)$ with the starting state s of $NG_F(v)$. The start state of the concatenation is the start state of $NG_F(u)$, and the target state is the target state of $NG_F(v)$. We combine neighborhood gadgets to clique gadgets as follows. Let $\kappa, \kappa' \geq 1$. Let $C = \{u_1, \dots, u_\kappa\}$ be a κ -clique and $C' = \{u'_1, \dots, u'_{\kappa'}\}$ be a κ' -clique in G . We define the following concatenation of NFAs and strings, respectively:

$$CG_F(C, \kappa, \kappa') := \bigcirc_{i=1}^{\kappa} \bigcirc_{j=1}^{\kappa'} NG_F(u_i),$$

$$CG_T(C', \kappa, \kappa') := \bigcirc_{i=1}^{\kappa} \bigcirc_{j=1}^{\kappa'} NG_T(u'_j).$$

Observe that we can walk from start to target state of $CG_F(C, \kappa, \kappa')$ parsing $CG_T(C', \kappa, \kappa')$ if and only if $C \cup C'$ forms a $(\kappa + \kappa')$ -clique, since the neighborhood gadgets check adjacency for each pair of nodes $u_i \in C$ and $u'_j \in C'$.

Complete Construction For $\kappa \geq 1$, let $\mathcal{C}(\kappa)$ be the set of κ -cliques in G , and set $m(\kappa) := |\mathcal{C}(\kappa)|$. Let $\kappa, \kappa' \geq 1$ such that $3\kappa + \kappa' = k$. The final text is defined as

$$T := \$ \circ \bigcirc_{C' \in \mathcal{C}(\kappa')} \left((\# \circ CG_T(C', \kappa, \kappa'))^{m(\kappa)+4} \circ \$ \right),$$

using alphabet $\{0, 1, \#, \$\}$.

The NFA F consists of four copies of the clique gadgets $CG_F(C, \kappa, \kappa')$ for any κ -clique C , denoted by $CG^r(i)$ for $1 \leq i \leq m(\kappa)$ and $1 \leq r \leq 4$. Additionally, we have states $s, s_1, s_2, \dots, s_{m(\kappa)}$ and $t, t_1, t_2, \dots, t_{m(\kappa)}$. These states are connected as follows. In the starting state s we can stay as long as we want, reading any symbol in the alphabet $\{0, 1, \#, \$\}$. When reading $\$$ we can alternatively go to state s_1 . In any state s_i when reading 0 or 1 we stay in s_i , while when reading $\#$ we either go to the starting state of $CG^1(i)$ or to s_{i+1} (the latter is only possible if $i < m(\kappa)$). For any $1 \leq r < 4$ and i, j , from the ending state of $CG^r(i)$ when reading $\#$ we can go to the

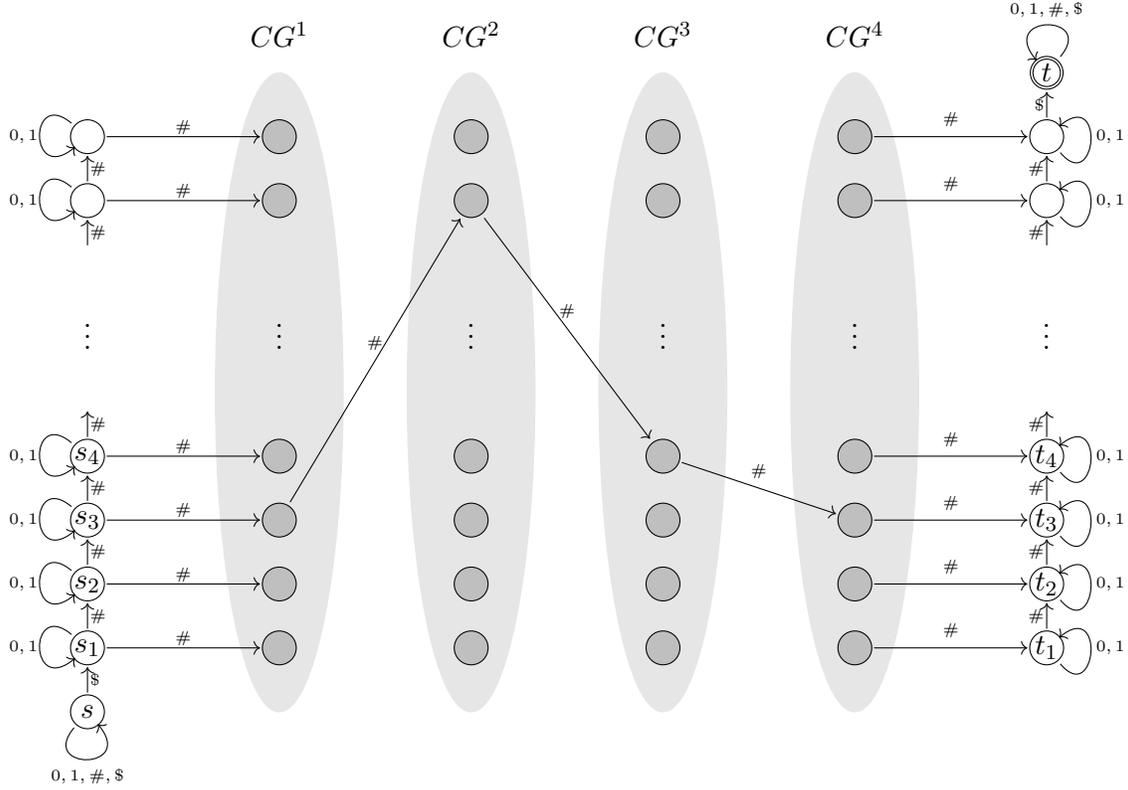


Figure 3: An illustration of the NFA constructed in the proof of Theorem 4.2. s is the starting state and t is the only accepting state. The first column consists of states $s, s_1, s_2, \dots, s_{m(\kappa)}$, the last column consists of $t_1, t_2, \dots, t_{m(\kappa)}, t$. The second, third, fourth and fifth columns contain the gadgets $CG^1(i), CG^2(i), CG^3(i), CG^4(i)$, respectively. We do not show all transitions between gadgets $CG^r(i)$ and $CG^{r+1}(j)$ for $r = 1, 2, 3$. As shown in the picture and as we prove, an accepting execution must visit $CG^1(i)$ and $CG^4(j)$ for $i = j$.

starting state of $CG^{r+1}(j)$ if the corresponding cliques together form a 2κ -clique. From the ending state of $CG^4(i)$ when reading $\#$ we go to t_i . In any state t_i when reading 0 or 1 we stay in t_i , while when reading $\#$ we go to t_{i+1} , or to t if $i = m(\kappa)$. Finally, t is the only accepting state and we stay in t reading any symbol in the alphabet. This finishes the construction of the NFA Acceptance instance. See Figure 3 for the illustration of the NFA.

Correctness Let us first show that if G contains a $(3\kappa + \kappa')$ -clique C then F accepts T . Write $C = C_1 + C_2 + C_3 + C'$, where C' is a κ' -clique and C_1, C_2, C_3 are κ -cliques (with indices i_1, i_2, i_3 in $\mathcal{C}(\kappa)$). We can stay in s until the beginning of the substring $T' := \$ \circ (\# \circ CG_T(C', \kappa, \kappa'))^{m(\kappa)+4} \circ \$$. With the first symbol $\$$ in T' we go to s_1 . We then walk to s_{i_1} reading $(\# \circ CG_T(C', \kappa, \kappa'))^{i_1-1}$. With $\#$ we then step to the starting state of $CG^1(i_1)$, corresponding to clique C_1 . Since $C_1 \cup C'$ forms a $(\kappa + \kappa')$ -clique, we can walk to the ending state of $CG^1(i_1)$ reading $CG_T(C', \kappa, \kappa')$. Since $C_1 \cup C_2$ forms a 2κ -clique, we can next step to the starting state of $CG^2(i_2)$ (corresponding to C_2). Similarly, we can then walk through $CG^2(i_2), CG^3(i_3)$ (corresponding to C_3), and $CG^4(i_1)$

(corresponding to C_1 again). Next we step to t_{i_1} reading $\#$, and then we simply walk to $t_{m(\kappa)}$ reading $(\# \circ CG_T(C', \kappa, \kappa'))^{m(\kappa)-i_1}$. Note that the number of times we read a symbol $\#$ is $i_1 - 1$ (for walking to s_{i_1}) plus 5 (for walking from s_{i_1} to t_{i_1}) plus $m(\kappa) - i_1$ (for walking from t_{i_1} to $t_{m(\kappa)}$), summing to $m(\kappa) + 4$. Hence, indeed we parse all symbols $\#$ in T' . Thus, we can next step to t reading the final symbol $\$$ of T' . We then stay in t reading the remainder of T . Since t is accepting, we are done.

For the other direction, note that if F accepts T then it also accepts some substring $T' := \$ \circ (\# \circ CG_T(C', \kappa, \kappa'))^{m(\kappa)+4} \circ \$$. Moreover, when reading T' we must walk through some clique gadgets $CG^1(i_1), CG^2(i_2), CG^3(i_3)$, and $CG^4(i_4)$, corresponding to κ -cliques C_1, C_2, C_3 , and C_4 . Note that the number of symbols $\#$ on such a walk is $i_1 - 1$ (for walking to s_{i_1}) plus 5 (for walking from s_{i_1} to t_{i_4}) plus $m(\kappa) - i_4$ (for walking from t_{i_4} to $t_{m(\kappa)}$), summing to $m(\kappa) + 4 + i_1 - i_4$. Since T' contains exactly $m(\kappa) + 4$ symbols $\#$, we obtain $i_1 = i_4$ and thus $C_1 = C_4$. By the restrictions on the edges from $CG^r(i)$ to $CG^{r+1}(j)$ we see that $C_1 \cup C_2, C_2 \cup C_3$, and $C_3 \cup C_4 = C_3 \cup C_1$ form 2κ -cliques. Moreover, since we walked through the clique gadgets we see that $C' \cup C_1, C' \cup C_2$, and $C' \cup C_3$ form $(\kappa + \kappa')$ -cliques. In total, we obtain that $C_1 \cup C_2 \cup C_3 \cup C'$ forms a $(3\kappa + \kappa' = k)$ -clique, finishing the correctness argument.

Size Bounds Note that clique gadgets CG_T in the text have length $O(\log V)$, while the clique gadgets CG_F in the automaton have $O(V \log V)$ states. We can thus read off a text length of $N = O(m(\kappa)m(\kappa') \log V) = O(V^{\kappa+\kappa'} \log V)$. Since the repetition $(\# \circ CG_T(C', \kappa, \kappa'))^{m(\kappa)+4}$ can be easily compressed to size $O(\log V)$ by Observation 2.2, we obtain a compressed size of $n = O(m(\kappa') \log V) = O(V^{\kappa'} \log V)$. Finally, the number of states is $q = O(m(\kappa)V \log V) = O(V^{\kappa+1} \log V)$. Note also that the output of this reduction can be computed in time $O(n + q^2)$, i.e., in linear time in the output description. We thus obtain the desired reduction which, as argued in the beginning of this proof, rules out a combinatorial $O(\min\{nq^3, Nq^2\}^{1-\varepsilon})$ algorithm for NFA Acceptance, assuming the combinatorial k -Clique conjecture.

Strengthening the Statement In the remainder, we verify that our construction proves the desired lower bound even restricted to instances with $n = \Theta(q^{\alpha_n})$ and $N = \Theta(q^{\alpha_N})$ for any $\alpha_N \geq \alpha_n > 0$. Note that the number of states, the size of the SLP, and the text length can all three be increased by easy padding. E.g., to increase the text length we introduce a garbage symbol “!” that can be read at any state of the automaton, not changing the current state, and add a suitable number of copies of “!” to the text. We now consider two cases.

Case 1: If $\alpha_N \geq \alpha_n + 1$, then set $\kappa, \kappa' \geq 1$ such that $3\kappa + \kappa' = k$ and $\kappa \approx k/(\alpha_n + 3)$ (recall that κ, κ' are restricted to be integers). We can ensure that $\kappa < k/(\alpha_n + 3) + 2$ and $\kappa' < \alpha_n k/(\alpha_n + 3) + 3$. Note that for any $\varepsilon > 0$, for sufficiently large $k = k(\varepsilon, \alpha_n)$ we have $\kappa + 1 < (1 + \varepsilon/2)k/(\alpha_n + 3)$ and $\kappa' < (1 + \varepsilon/2)\alpha_n k/(\alpha_n + 3)$. We can thus pad the number of states from $O(V^{\kappa+1} \log V)$ to $q = \Theta(V^{(1+\varepsilon/2)k/(\alpha_n+3)})$ and the compressed size from $O(V^{\kappa'} \log V)$ to $n = \Theta(V^{(1+\varepsilon/2)\alpha_n k/(\alpha_n+3)}) = \Theta(q^{\alpha_n})$. Similarly, for the decompressed text length, using $\alpha_N \geq \alpha_n + 1$, we have $N = O(V^{\kappa+\kappa'} \log V) = O(V^{(1+\varepsilon/2)(\alpha_n+1)k/(\alpha_n+3)}) = O(V^{(1+\varepsilon/2)\alpha_N k/(\alpha_n+3)}) = O(q^{\alpha_N})$, which we can pad to equality. Then we indeed end up with an instance with $N = \Theta(q^{\alpha_N})$ and $n = \Theta(q^{\alpha_n})$. Hence, if NFA Acceptance can be solved in combinatorial time $O(\min\{nq^3, Nq^2\}^{1-\varepsilon})$ restricted to such instances, then we obtain a combinatorial algorithm for k -Clique in time $O((nq^3)^{1-\varepsilon}) = O(V^{(\alpha_n+3) \cdot (1-\varepsilon)(1+\varepsilon/2)k/(\alpha_n+3)}) = O(V^{k(1-\varepsilon/2)})$, contradicting the combinatorial k -Clique conjecture.

Case 2: If $\alpha_N < \alpha_n + 1$, then we have to slightly adapt the above construction. We introduce a third parameter $\hat{\kappa} \leq \kappa$ and let the first and fourth column of clique gadgets $CG^1(i)$ and $CG^4(i)$ range over $\hat{\kappa}$ -cliques. At the same time, we change the number of repetitions of each part in the text from $m(\kappa) + 4$ to $m(\hat{\kappa}) + 4$. We are now detecting $(2\kappa + \hat{\kappa} + \kappa')$ -cliques in G . It can be checked that this does not violate the correctness of the construction. The new size bounds are $N = O(V^{\hat{\kappa} + \kappa'} \log V)$, $n = O(V^{\kappa'} \log V)$, and $q = O(V^\kappa \log V)$. Furthermore, we now allow to set $\kappa' = 0$, in which case the text is not responsible for choosing any part of the clique. Since in this case we do not need any clique gadgets, we define $CG_F(C, \kappa, 0)$ to consist of a single state $s = t$ and $CG_T(C', \kappa, 0)$ to be the empty string. In this case we set the final string to be $T := \$ \#^{m(\kappa)+4} \$$. The same correctness proof goes through.

We now choose integers $\kappa, \hat{\kappa} \geq 1$ and $\kappa' \geq 0$ with $2\kappa + \hat{\kappa} + \kappa' = k$ and $\hat{\kappa} \leq \kappa$ such that $\kappa \approx k/(\alpha_N + 2)$, $\kappa' \approx \max\{0, (\alpha_N - 1)k/(\alpha_N + 2)\}$, and $\hat{\kappa} \approx \min\{\alpha_N, 1\} \cdot k/(\alpha_N + 2)$. Similarly to case 1, we can ensure for any $\varepsilon > 0$ and sufficiently large $k = k(\varepsilon, \alpha_N)$ that $\kappa < (1 + \varepsilon/2)k/(\alpha_N + 2)$, $\kappa' \leq \max\{0, (1 + \varepsilon/2)(\alpha_N - 1)k/(\alpha_N + 2)\}$, and $\hat{\kappa} < (1 + \varepsilon/2) \min\{\alpha_N, 1\} \cdot k/(\alpha_N + 2)$. We can thus pad the number of states to $q = \Theta(V^{(1+\varepsilon/2)k/(\alpha_N+2)})$ and since $\hat{\kappa} + \kappa' < (1 + \varepsilon/2)\alpha_N k/(\alpha_N + 2)$ we can pad the decompressed text length to $N = \Theta(V^{(1+\varepsilon/2)\alpha_N k/(\alpha_N+2)}) = \Theta(q^{\alpha_N})$. For the compressed size, note that by the assumptions $\alpha_N < \alpha_n + 1$ and $\alpha_n > 0$ we have $\kappa' \leq \max\{0, (1 + \varepsilon/2)(\alpha_N - 1)k/(\alpha_N + 2)\} < (1 + \varepsilon/2)\alpha_n k/(\alpha_N + 2)$, and thus $n = O(V^{(1+\varepsilon/2)\alpha_n k/(\alpha_N+2)}) = O(q^{\alpha_n})$, which we can pad to equality. Then we indeed end up with an instance with $N = \Theta(q^{\alpha_N})$ and $n = \Theta(q^{\alpha_n})$. Hence, if NFA Acceptance can be solved in combinatorial time $O(\min\{nq^3, Nq^2\}^{1-\varepsilon})$ restricted to such instances, then we obtain a combinatorial algorithm for k -Clique in time $O((Nq^2)^{1-\varepsilon}) = O(V^{(\alpha_N+2) \cdot (1-\varepsilon)(1+\varepsilon/2)k/(\alpha_N+2)}) = O(V^{k(1-\varepsilon/2)})$, contradicting the combinatorial k -Clique conjecture. \square

4.2 Context-Free Grammar Parsing

We again assume that the alphabet size $|\Sigma|$ is constant throughout this section.

In this section we show a strong conditional lower bound for context-free grammar parsing. Recall that a context-free grammar (CFG) Γ consists of a set of terminals Σ , a set of non-terminals Ω , a starting non-terminal $S \in \Omega$, and a set of productions Φ , each of the form $A \rightarrow \alpha$, where $A \in \Omega$ and $\alpha \in (\Sigma \cup \Omega)^*$. The size $|\Gamma|$ is the total length of all α over all productions. Applying a production $A \rightarrow \alpha$ to a string $\beta = \beta_1 A \beta_2 \in (\Sigma \cup \Omega)^*$ means to generate the string $\beta_1 \alpha \beta_2$. The language $L(\Gamma)$ is the set of strings in Σ^* that can be generated by starting with S and repeatedly applying productions. More generally, for any non-terminal A the language $L(A)$ is the set of strings in Σ^* that can be generated by starting with A .

Problem 4.3 (CFG Recognition). *Given a text T of length N by a grammar-compressed representation \mathcal{T} of size n as well as a CFG Γ , decide whether $T \in L(\Gamma)$.*

(CFG parsing is an augmentation of this decision problem where in case $T \in L(\Gamma)$ we also need to return a sequence of productions as a certificate.)

As discussed in the introduction, after decompressing the text T we can use classic parsers to solve CFG recognition in time $O(N^3 \text{poly}(|\Gamma|))$ [25, 46, 80, 29], while Valiant's parser uses fast matrix multiplication to obtain an improved running time of $O(N^\omega \text{poly}(|\Gamma|))$ [72].⁸ In the uncompressed setting, matching lower bounds based on the k -Clique conjecture were shown by Abboud et al. [1].

⁸We ignore the specific polynomial dependence on $|\Gamma|$, since we are more interested in the dependence on N .

In the compressed setting no improved algorithms are known, even for, say, $n = N^{0.01}$. Below we prove a matching lower bound for both running times $O(N^3)$ and $O(N^\omega)$, even restricted to very small grammars and quite compressible strings. Our proof differs considerably from the conditional lower bound in the uncompressed setting by Abboud et al. [1], as their strings are not compressible in a strong sense. On a high level, their construction implements adjacency tests locally, around three chosen positions that encode three k -cliques. In our construction, we instead implement adjacency tests on a more global level, by choosing three offsets and reading all text positions that adhere to these offsets. This global view makes it possible to construct a compressible text.

Theorem 4.4. *Assuming the k -Clique conjecture, there is no $O(N^{\omega-\varepsilon})$ time algorithm for CFG recognition for any $\varepsilon > 0$. Assuming the combinatorial k -Clique conjecture, there is no combinatorial $O(N^{3-\varepsilon})$ time algorithm for CFG recognition for any $\varepsilon > 0$. Both results hold even restricted to instances with $|\Gamma| = O(\log N)$ and $n = O(N^\varepsilon)$.*

Proof. Let $k \geq 1$ and let $G = (V, E)$ be a k -Clique instance. We will construct a CFG Γ of size $O(\log V)$ and a text T of length $N = O(V^{k+2})$ generated by an SLP \mathcal{T} of size $n = O(V^3)$ such that $T \in L(\Gamma)$ holds if and only if G contains a $3k$ -clique. Note that an $O(N^{\omega-\varepsilon}) = O(N^{\omega(1-\varepsilon/3)})$ algorithm for CFG recognition would then imply an algorithm for $3k$ -Clique in time $O(V^{(k+2)\omega(1-\varepsilon/3)})$, which for $k \geq 12/\varepsilon$ is bounded by $O(V^{k(1+\varepsilon/6)\omega(1-\varepsilon/3)}) = O(V^{\omega k(1-\varepsilon/3)})$, contradicting the $3k$ -Clique conjecture. The argument for combinatorial algorithms is analogous. Moreover, we have $|\Gamma| = O(\log V) = O(\log N)$ and $n = O(V^3) = O(N^{3/(k+2)}) = O(N^\varepsilon)$ for $k \geq 3/\varepsilon$.⁹

In our construction we enumerate all k -tuples of vertices $U = (u_1, \dots, u_k)$. Choosing three such k -tuples U_1, U_2, U_3 we then need to check that (1) each k -tuple U_i forms a k -clique and (2) each pair U_i, U_j forms a biclique for $i \neq j$. We remark that it is indeed necessary to enumerate all k -tuples and not just, say, all k -cliques, as the k -tuples are much more structured, leading to compressible strings. In the following we construct gadgets that perform these tests. We will use alphabet $\Sigma = \{0, 1, \#, \$, x, y, z\}$.

Offsets Let $U(i)$ be the i -th k -tuple $(u_1, \dots, u_k) \in V^k$ in lexicographic order. Choosing a k -tuple thus correspond to choosing a number $1 \leq i \leq V^k$, which we will interpret as an offset in the text T , resulting in *relevant* positions of the form $i + V^k \cdot \mathbb{N}$. In order to only read the relevant positions, we need to implement jumping over $V^k - 1$ symbols, so that after reading one relevant symbol we can jump to the next one. To this end, we construct a non-terminal X of Γ with $L(X) = \Sigma^{V^k-1}$. This can be build by constructing non-terminals X_d with $L(X_d) = \Sigma^{2^d}$ by the productions

$$\begin{aligned} X_0 &\rightarrow \sigma && \text{for any } \sigma \in \Sigma, \\ X_d &\rightarrow X_{d-1}X_{d-1} && \text{for } 1 \leq d \leq \log(V^k - 1). \end{aligned}$$

Then the production $X \rightarrow X_{i_1} \dots X_{i_\ell}$, where i_1, \dots, i_ℓ are the 1-bits in the binary encoding of $V^k - 1$, yields the desired non-terminal X . Note that this yields a grammar of size $O(\log V)$.

Clique Test We now design gadgets that allow to test for any offset i whether $U(i)$ forms a k -clique. Let $\bar{E} = \binom{V}{2} \setminus E$ be the non-edges of G . Let $[\cdot]$ be the Kronecker symbol, i.e., $[\text{true}] = 1$

⁹Strictly speaking, we need to pad the text length to $\Theta(V^{k+2})$ first. This can easily be accomplished by adding garbage to the text and garbage handling rules to the grammar.

and $[\text{false}] = 0$. We use the following text:

$$T_C := \$^{V^k} \circ \left(\bigcirc_{\{u,v\} \in \bar{E}} \bigcirc_{1 \leq i \leq V^k} [u \text{ and } v \text{ appear in } U(i)] \right) \circ \$^{V^k}.$$

For any offset $1 \leq i \leq V^k$, if $U(i) = (u_1, \dots, u_k)$ forms a k -clique then no non-edge appears among $\{u_1, \dots, u_k\}$, and thus $T_C[i + j \cdot V^k] = 0$ for all $1 \leq j \leq |\bar{E}|$. The opposite implication holds as well. This leads us to testing for a k -clique via the following CFG rules:

$$C \rightarrow \$X\tilde{C}, \quad \tilde{C} \rightarrow 0X\tilde{C} \mid \$.$$

Lemma 4.5. *We call $T_C(i) := T_C[i..i + 1 + (|\bar{E}| + 1) \cdot V^k]$ for $1 \leq i \leq V^k$ the valid substrings of T_C . Any substring of T_C that is parsable by C is valid. Moreover, substring $T_C(i)$ is parsable by C if and only if the k -tuple $U(i)$ forms a k -clique in G .*

Proof. The first statement follows by C starting and ending with a $\$$ symbol and advancing by $V^k - 1$ steps via X . The second statement follows from the argument above this lemma. \square

Lemma 4.6. *The string T_C has an SLP of size $O(V^3)$.*

Proof. For any $1 \leq d \leq k$, $\sigma \in \Sigma = \{0, 1, \$, \#, x, y, z\}$, and $S \subseteq V$ with $|S| \leq 2$ we define the following SLP rules:

$$\begin{aligned} \text{Repeat}_\sigma^{(0)} &\rightarrow \sigma, \\ \text{Repeat}_\sigma^{(d)} &\rightarrow \bigcirc_{v \in V} \text{Repeat}_\sigma^{(d-1)}, \\ \text{Incl}_S^{(0)} &\rightarrow \begin{cases} 1, & \text{if } S = \emptyset, \\ 0, & \text{otherwise} \end{cases} \\ \text{Incl}_S^{(d)} &\rightarrow \bigcirc_{v \in V} \text{Incl}_{S \setminus \{v\}}^{(d-1)}, \\ \text{C-Test} &\rightarrow \text{Repeat}_\$^{(k)} \circ \left(\bigcirc_{\{u,v\} \in \bar{E}} \text{Incl}_{\{u,v\}}^{(k)} \right) \circ \text{Repeat}_\$^{(k)}. \end{aligned}$$

We claim that $\text{eval}(\text{C-Test}) = T_C$. Note that $\text{Repeat}_\sigma^{(d)}$ generates the string σ^{V^d} , and thus the prefix and suffix $\$^{V^k}$ is correct. Further, it can be checked that $\text{Incl}_S^{(d)}$ generates a string of length V^d where the i -th position, corresponding to a d -tuple $(u_1, \dots, u_d) \in V^d$, is 1 if $S \subseteq \{u_1, \dots, u_d\}$ and 0 otherwise. Hence, writing the string $\text{Incl}_{\{u,v\}}^{(k)}$ for all $\{u,v\} \in \bar{E}$ yields the middle part of the string T_C . This proves the claim.

Note that the total size of the above SLP for T_C , i.e., the total number of symbols on the right hand sides of the above rules, is indeed $O(V^3)$. \square

Biclique Test We next design gadgets that allow us to test for two offsets i, j whether $u \sim v$ for all $u \in U(i), v \in U(j)$, i.e., whether $U(i), U(j)$ form a biclique. To this end, we let V^{rev} be the reverse ordering of the vertices in V and define the texts

$$\begin{aligned} T_B &:= \#^{V^k} \circ \left(\bigcirc_{u \in V} \bigcirc_{1 \leq i \leq V^k} [u \text{ appears in } U(i)] \right) \circ \#^{V^k}, \\ T'_B &:= \#^{V^k} \circ \left(\bigcirc_{u \in V^{\text{rev}}} \bigcirc_{1 \leq i \leq V^k} [u \text{ is adjacent to every vertex in } U(i)] \right) \circ \#^{V^k}. \end{aligned}$$

Note that $U(i), U(j)$ form a biclique if every vertex that appears in $U(i)$ is adjacent to every vertex in $U(j)$. Thus, for every $1 \leq \ell \leq V$ we want that if $T_B[i + \ell \cdot V^k] = 1$ then also $T'_B[j + (V + 1 - \ell) \cdot V^k] = 1$. This leads us to testing for a biclique via the following CFG rules:

$$\begin{aligned} B_{\text{in}} &\rightarrow \# X B X \# \\ B &\rightarrow 1 X B X 1 \quad | \quad 0 X B X 1 \quad | \quad 0 X B X 0 \quad | \quad \# B_{\text{out}} \# \end{aligned}$$

We view this part of the grammar as a subroutine that is started by invoking B_{in} and that can be followed by further operations by adding productions starting from B_{out} . Note that each call of a rule of B_{in} or B reads V^k symbols from the left and from the right, except for the last one, which reads 1 symbol from the left and from the right. That is, the offsets are never changed throughout the parsing process. The parsing rules check that a 1 at a certain position in T_B implies a 1 at the corresponding position in T'_B . Hence, when starting with offsets i in T_B and j in T'_B , this process checks that $U(i), U(j)$ form a biclique. It stops when we reach the $\#$ -blocks at the end of T_B and at the beginning of T'_B , where we exit to B_{out} . Then it depends on the (not yet defined) productions involving B_{out} whether the remainder of the string can be parsed. In summary, we obtain the following.

Lemma 4.7. *We call $T_B(i) := T_B[i..i + 1 + (V + 1) \cdot V^k]$ and $T'_B(j) := T'_B[j..j + 1 + (V + 1) \cdot V^k]$ for $1 \leq i, j \leq V^k$ the valid substrings of T_B and T'_B , respectively. Let R be any string. Then B_{in} can parse $T_B(i) R T'_B(j)$ if and only if $U(i), U(j)$ form a biclique and B_{out} can parse R . Moreover, if \tilde{T}_B and \tilde{T}'_B are substrings of T_B and T'_B , respectively, and B_{in} can parse $\tilde{T}_B R \tilde{T}'_B$ such that B_{out} parses R , then \tilde{T}_B and \tilde{T}'_B are valid.*

Lemma 4.8. *The strings T_B and T'_B have SLPs of size $O(V^2)$.*

Proof. Note that T_B is the string generated by the following SLP, where we use notation as in Lemma 4.6:

$$\text{B-Test} \rightarrow \text{Repeat}_{\#}^{(k)} \circ \left(\bigcirc_{v \in V} \text{Incl}_{\{v\}}^{(k)} \right) \circ \text{Repeat}_{\#}^{(k)}.$$

This has size $O(V^2)$ as shown in the proof of Lemma 4.6.

For T'_B we use the following SLP rules for $1 \leq d \leq k$ and $v \in V$:

$$\begin{aligned} \text{Adj}_v^{(0)} &\rightarrow 1, \\ \text{Adj}_v^{(d)} &\rightarrow \bigcirc_{u \in V} \begin{cases} \text{Adj}_v^{(d-1)}, & \text{if } \{u, v\} \in E \\ \text{Repeat}_0^{(d-1)}, & \text{otherwise} \end{cases} \\ \text{B}'\text{-Test} &\rightarrow \text{Repeat}_{\#}^{(k)} \circ \left(\bigcirc_{v \in V^{\text{rev}}} \text{Adj}_v^{(k)} \right) \circ \text{Repeat}_{\#}^{(k)} \end{aligned}$$

An easy inductive proof shows that $\text{Adj}_v^{(d)}$ generates a string of length V^d where the i -th position, corresponding to a d -tuple $(u_1, \dots, u_d) \in V^d$, is 1 if v is adjacent to every u_i , and 0 otherwise. Hence, writing $\text{Adj}_v^{(k)}$ for all $v \in V$ (in reverse order) yields the middle part of T'_B , and thus $\text{B}'\text{-Test}$ generates T'_B . Again, the total size of the right hand sides is $O(V^2)$, so the SLP has size $O(V^2)$. \square

Complete Construction The final string is

$$T := x^{V^k} T_C T_B T_B T'_B y^{V^k} T_C T_B T'_B T'_B T_C z^{V^k}.$$

Here, the parts x^{V^k} , y^{V^k} , and z^{V^k} are used to choose three offsets i_1, i_2, i_3 , corresponding to three k -tuples $U(i_1), U(i_2), U(i_3)$. The three copies of T_C are used to check that each $U(i_j)$ forms a k -clique. The left copy of $T_B T'_B$ is used for checking that $U(i_1), U(i_2)$ forms a biclique, similarly for the right copy and $U(i_2), U(i_3)$. Finally, the leftmost T_B and rightmost T'_B are used to check that $U(i_1), U(i_3)$ form a biclique. Note that T uses alphabet $\Sigma = \{0, 1, \#, \$, x, y, z\}$.

We now describe the final grammar Γ . We copy the non-terminals $B_{\text{in}}, B, B_{\text{out}}$ to $\tilde{B}_{\text{in}}, \tilde{B}, \tilde{B}_{\text{out}}$, since we need this subroutine twice with different productions starting from B_{out} . We let S be a new starting symbol and define the following productions, additional to the ones defined above:

$$\begin{aligned} S &\rightarrow xS \mid Sz \mid XCXB_{\text{in}}XCX \\ B_{\text{out}} &\rightarrow X\tilde{B}_{\text{in}}XyXCX\tilde{B}_{\text{in}}X \\ \tilde{B}_{\text{out}} &\rightarrow \#\tilde{B}_{\text{out}} \mid \epsilon, \end{aligned}$$

where ϵ denotes the empty string. This finishes the construction of the CFG recognition instance.

Correctness We show that $T \in L(\Gamma)$ holds if and only if there is a $3k$ -clique in G . Assume that G contains a $3k$ -clique and let $1 \leq i_1, i_2, i_3 \leq V^k$ be such that $U(i_1) \cup U(i_2) \cup U(i_3)$ forms a $3k$ -clique. Remove i_1 symbols x from the left end of T and $V^k - i_3 + 1$ symbols z from the right, leaving offsets i_1 and i_3 , respectively. Then apply the rule $S \rightarrow XCXB_{\text{in}}XCX$. The outer calls to X keep the offsets i_1 and i_3 by advancing to the next relevant positions w.r.t. offsets i_1 and i_3 , respectively. By Lemma 4.5, the calls of C parse valid substrings of T_C starting and ending with offset i_1 and i_3 , respectively. The lemma is applicable since $U(i_1)$ and $U(i_3)$ form k -cliques. The further calls to X again advance to the next relevant positions w.r.t. offsets i_1 and i_3 , now lying in the outer $\#$ -blocks in the leftmost T_B and rightmost T'_B , respectively. Finally, by Lemma 4.7 the call to B_{in} reads valid substrings of the leftmost T_B and rightmost T'_B and ends with B_{out} . The lemma is applicable since $U(i_1), U(i_3)$ forms a biclique. The outer calls to X in the rule $B_{\text{out}} \rightarrow X\tilde{B}_{\text{in}}XyXCX\tilde{B}_{\text{in}}X$ then advances the left and right end to the first relevant position w.r.t. offset i_1 in the second copy of T_B and the last relevant position w.r.t. offset i_3 in the second-to-last copy of T'_B . We match the y appearing in this rule to the i_2 -th y in the y^{V^k} part of T . To the right of y , C parses a valid substring of T_C , which works since $U(i_2)$ forms a k -clique. The remaining \tilde{B}_{in} then has to parse valid substrings of the right copy of $T_B T'_B$, starting with offset i_2 and ending with offset i_3 . Similarly, to the left of y , \tilde{B}_{in} has to parse valid substrings of the left copy of $T_B T'_B$, starting with offset i_1 and ending with offset i_2 . This works as $U(i_1), U(i_2)$ and $U(i_2), U(i_3)$ form bicliques. Note that after reaching \tilde{B}_{out} we are left with some symbols of the last $\#$ -block of T_B and some symbols of the first $\#$ -block of T'_B . Both can be parsed completely using the rules involving \tilde{B}_{out} . Thus, we have $T \in L(\Gamma)$.

For the other direction, we follow the same line of arguments, observing that there was no choice except for the offsets i_1, i_2, i_3 . The core of the argument is that $U(i_1) \cup U(i_2) \cup U(i_3)$ forms a $3k$ -clique if and only if each $U(i_j)$ forms a k -clique and each pair $U(i_j), U(i_{j'})$ forms a biclique.

Size Bounds Since T consists of $O(|V| + |\bar{E}|)$ parts of length V^k , the text length is $O(V^{k+2})$. By Lemmas 4.6 and 4.8 and since x^{V^k} has an SLP of size $O(\log V)$, T has an SLP of size $O(V^3)$.

Finally, the size of the grammar Γ is $O(\log V)$, the bottleneck being the non-terminal X that ensures offset consistency. Hence, all claimed size bounds are met. Note also that the constructed instance can be computed in time linear in the output size. This finishes the proof of Theorem 4.4. \square

4.3 RNA Folding

We now give a variant of the construction for CFG recognition, proving a matching conditional lower bound for RNA folding.

Again we consider a constant-size alphabet Σ , however, now each symbol $\sigma \in \Sigma$ has a unique counterpart $\bar{\sigma} \in \Sigma$ such that $\bar{\bar{\sigma}} = \sigma$. We say that $\sigma \in \Sigma$ and its counterpart $\bar{\sigma}$ *match*.

Two pairs of indices $(i, j), (i', j')$ with $i < j$ and $i' < j'$ are said to *cross* if at least one of the following conditions holds: (1) $i = i'$ or $i = j'$ or $j = i'$ or $j = j'$, (2) $i < i' < j < j'$, or (3) $i' < i < j' < j$. In other words, $(i, j), (i', j')$ with $i < j$ and $i' < j'$ are non-crossing if they are disjoint, i.e., $i < j < i' < j'$ or $i' < j' < i < j$, or they are nesting, i.e., $i < i' < j' < j$ or $i' < i < j < j'$.

Problem 4.9 (RNA Folding). *Given a text T of length N by a grammar-compressed representation \mathcal{T} of size n , compute the maximum number of pairs $R \subseteq \{(i, j) \mid 1 \leq i < j \leq N\}$ such that for every $(i, j) \in R$ the symbols $T[i]$ and $T[j]$ match and there are no crossing pairs in R . We denote this maximum number by $\text{RNA}(T)$.*

We refer to the set R as a *matching* of T .

In the uncompressed setting, RNA Folding has an easy dynamic programming solution in time $O(N^3)$ [30]. Using fast matrix multiplication, this was recently improved to $O(N^{2.82})$ [15]. For combinatorial algorithms, a matching lower bound of $N^{3-o(1)}$ assuming the combinatorial k -Clique conjecture was recently shown by Abboud et al. [1]. They also prove a conditional lower bound of $N^{\omega-o(1)}$ assuming the k -Clique conjecture, however, this leaves a gap to the current upper bound.

As for CFG parsing, no improved algorithms are known in the compressed setting, even for, say, $n = N^{0.01}$. Here we prove lower bounds of $N^{3-o(1)}$ for combinatorial algorithms and $N^{\omega-o(1)}$ in general, assuming the (combinatorial) k -Clique conjecture.

Theorem 4.10. *Assuming the k -Clique conjecture, there is no $O(N^{\omega-\varepsilon})$ time algorithm for RNA Folding for any $\varepsilon > 0$. Assuming the combinatorial k -Clique conjecture, there is no combinatorial $O(N^{3-\varepsilon})$ time algorithm for RNA Folding for any $\varepsilon > 0$. Both results hold even restricted to instances with $n = O(N^\varepsilon)$.*

Abboud et al. [1] showed that RNA Folding is equivalent to the following weighted variant.

Problem 4.11 (Weighted RNA Folding). *We are given a text T of length N by a grammar-compressed representation \mathcal{T} of size n as well as a weight function $w: \Sigma \rightarrow [M]$ with $w(\sigma) = w(\bar{\sigma})$ for all $\sigma \in \Sigma$. For any set $R \subseteq \{(i, j) \mid 1 \leq i < j \leq N\}$ define its weight as $\sum_{(i, j) \in R} w(T[i])$. Compute the maximum weight of any set R such that for every $(i, j) \in R$ the symbols $T[i]$ and $T[j]$ match and there are no crossing pairs in R . We denote this maximum weight by $\text{WRNA}(T)$.*

Lemma 4.12 (Lemma 2 in [1]). *For an instance T of Weighted RNA Folding, consider the string $\tilde{T} := T[1]^{w(T[1])} \dots T[N]^{w(T[N])}$, i.e., each symbol $T[i]$ is repeated $w(T[i])$ times. Then we have $\text{WRNA}(T) = \text{RNA}(\tilde{T})$.*

Proof of Theorem 4.10. Let $k \geq 1$ and let $G = (V, E)$ be a k -Clique instance. We will construct a Weighted RNA Folding instance T of length $O(V^{k+2})$ (and $\Omega(V^k)$) generated by an SLP \mathcal{T} of size $O(V^3)$ and a number λ such that $\text{WRNA}(T) \geq \lambda$ holds if and only if G contains a $3k$ -clique. The alphabet size will be $|\Sigma| = 48$ and the weights are bounded by $O(V^2)$. By Lemma 4.12, the corresponding unweighted text \tilde{T} has $\text{RNA}(\tilde{T}) = \text{WRNA}(T)$ and thus $\text{RNA}(\tilde{T}) \geq \lambda$ holds if and only if G contains a $3k$ -clique. Moreover, since the weights in T are bounded by $O(V^2)$ we have $N = |\tilde{T}| = O(V^2|T|) = O(V^{k+4})$. Finally, by compressing $O(V^2)$ repetitions to $O(\log V)$ SLP rules, \tilde{T} has an SLP $\tilde{\mathcal{T}}$ of size $n = O(|\mathcal{T}| \log V) = O(V^3 \log V)$.

Hence, an $O(N^{\omega-\varepsilon}) = O(N^{\omega(1-\varepsilon/3)})$ algorithm for RNA Folding would imply an algorithm for $3k$ -Clique in time $O(V^{(k+4)\omega(1-\varepsilon/3)})$, which for $k \geq 24/\varepsilon$ is bounded by $O(V^{k(1+\varepsilon/6)\omega(1-\varepsilon/3)}) = O(V^{\omega k(1-\varepsilon/6)})$, contradicting the $3k$ -Clique conjecture. The argument for combinatorial algorithms is analogous. Moreover, we have $n = O(V^3 \log V) = O(N^{3/k} \log V) = O(N^\varepsilon)$ for $k \geq 3/\varepsilon$.

To construct the desired instance T, \mathcal{T} of Weighted RNA Folding, we again enumerate all k -tuples $U(i)$ for $1 \leq i \leq V^k$, as in the proof for CFG parsing. We again choose three such k -tuples $U(i_1), U(i_2), U(i_3)$ and check that each $U(i_j)$ forms a k -clique and all pairs $U(i_j), U(i_{j'})$ form a biclique for $j \neq j'$.

Clique Test Consider alphabet $\{0, \bar{0}, 1, \bar{1}\}$ (with weights 1) and set for $e \in \bar{E}$ and $1 \leq i \leq V^k$

$$r_{e,i} := \begin{cases} \bar{1}, & \text{if some node in } e \text{ does not appear in } U(i) \\ \bar{0}, & \text{otherwise} \end{cases}$$

Since $U(i)$ forms a k -clique iff for every non-edge at least one of the endpoints does not appear in $U(i)$, we obtain:

Lemma 4.13. *Set $r_i := \bigcirc_{e \in \bar{E}} r_{e,i}$. We have $\text{WRNA}(1^{\bar{E}} r_i) \leq \bar{E}$, with equality if and only if $U(i)$ forms a k -clique.*

Biclique Test Consider alphabet $\{2, \bar{2}, 3, \bar{3}, 4, \bar{4}\}$ (with weights 1) and set for $v \in V$ and $i \in [V^k]$

$$p_{v,i} := \begin{cases} 24, & \text{if } v \text{ appears in } U(i) \\ 2\bar{3}4, & \text{otherwise} \end{cases} \quad q_{v,i} := \begin{cases} \bar{2}\bar{4}, & \text{if } v \text{ is adjacent to every node in } U(i) \\ \bar{3}\bar{4}, & \text{otherwise} \end{cases}$$

Lemma 4.14. *Set $p_i := \bigcirc_{v \in V} p_{v,i}$ and $q_i := \bigcirc_{v \in V} q_{v,i}$. For any i, j , we have $\text{WRNA}(p_i q_j) \leq 2V$, with equality if and only if $U(i), U(j)$ form a biclique.*

Proof. Note that the total weight of q_j is $2V$, which shows the upper bound $\text{WRNA}(p_i q_j) \leq 2V$. To obtain equality, all symbols in q_j must be matched. In particular, the $\bar{4}$ in $q_{v,j}$ must be matched to the 4 in $p_{v,i}$. It follows that the $\bar{2}$ or $\bar{3}$ in $q_{v,j}$ can only be matched to a 2 or 3 in $p_{v,i}$. Hence, we have $\text{WRNA}(p_i q_j) = 2V$ if and only if there is no $v \in V$ such that v appears in $U(i)$ but v is not adjacent to every node in $U(j)$, which happens if and only if $U(i), U(j)$ form a biclique. \square

Complete Construction For any symbol σ used so far, we introduce two copies σ' and σ'' . For the strings $r_{e,i}, p_{v,i}, q_{v,i}$ defined above, we write $r'_{e,i}, p'_{v,i}, q'_{v,i}$ and $r''_{e,i}, p''_{v,i}, q''_{v,i}$ to denote that we replace all symbols by their primed copies. For $i_1, i_2, i_3 \in [V^k]$ consider the string

$$T(i_1, i_2, i_3) := 1^{\bar{E}} r_{i_1} p_{i_1} p'_{i_1} 1^{\bar{E}} r'_{i_2} q'_{i_2} p''_{i_2} 1^{\bar{E}} r''_{i_3} q''_{i_3} q_{i_3}.$$

Note that the alphabet is partitioned such that the only possible matchings are among $1^{\bar{E}} r_{i_1}$, $1'^{\bar{E}} r'_{i_2}$, $1''^{\bar{E}} r''_{i_3}$ as well as $p_{i_1} q_{i_3}$, $p'_{i_1} q'_{i_2}$, $p''_{i_2} q''_{i_3}$. Also note that these pairs are non-crossing. Hence, by Lemmas 4.13 and 4.14, we have $\text{WRNA}(T(i_1, i_2, i_3)) \leq 6V + 3\bar{E}$, with equality if and only if $U(i_j)$ forms a k -clique and $U(i_j), U(i_{j'})$ form a biclique for any $j \neq j'$, which happens if and only if $U(i_1) \cup U(i_2) \cup U(i_3)$ forms a $3k$ -clique.

This is close to a complete reduction. It remains to force the choice of consistent offsets i_1, i_2, i_3 , which we accomplish with the following lemma. Its proof is technical and deferred to the end of this section.

Lemma 4.15. *Let $A, B, W \geq 1$. Let $x_{a,b}$ for $a \in [A]$, $b \in [B]$ be strings over alphabet Σ , each with total weight $\sum_i w(x_{a,b}[i]) \leq W$. Assume that no two symbols in $\bigcirc_{a,b} x_{a,b}$ match. Let $5, \bar{5}, 6, \bar{6}, 7, \bar{7}$ be new symbols not appearing in Σ , with weights $w(5) = w(\bar{5}) = w(7) = w(\bar{7}) = 4AW$ and $w(6) = w(\bar{6}) = 8AW$. Set $\rho := (8A + 12)ABW$ and*

$$G(\{x_{a,b}\}) := 5^B (6 \bar{5})^B \circ \left(\bigcirc_{a \in [A]} \left(\bigcirc_{b \in [B]} \bar{6} x_{a,b} \right) \circ 6^B \right) \circ \bar{6}^B (7 \bar{6})^B \bar{7}^B.$$

Then for any strings y_1, y_2 over alphabet Σ we have

$$\text{WRNA}(y_1 G(\{x_{a,b}\}) y_2) = \rho + \max_{b \in [B]} \text{WRNA}\left(y_1 \circ \left(\bigcirc_{a \in [A]} x_{a,b} \right) \circ y_2\right).$$

We apply the above lemma as follows. Let $B = V^k$ and $A = V + 2\bar{E}$, and for $b \in [B]$ set $x_{a,b} := r_{a,b}$ for $a \in [\bar{E}]$, $x_{\bar{E}+a,b} := p_{a,b}$ for $a \in [V]$, and $x_{\bar{E}+V+a,b} := p'_{a,b}$ for $a \in [V]$. Note that $\bigcirc_{a \in [A]} x_{a,i_1} = r_{i_1} p_{i_1} p'_{i_1}$, which is a substring of $T(i_1, i_2, i_3)$. Construct $G(\{x_{a,b}\})$. Similarly define $y_{a,b}$ so that $\bigcirc_{a \in [A]} y_{a,i_2} = r'_{i_2} q'_{i_2} p''_{i_2}$, and construct $G'(\{y_{a,b}\})$, where the new symbols are now $5', \bar{5}', 6', \bar{6}', 7', \bar{7}'$. Similarly define $z_{a,b}$ so that $\bigcirc_{a \in [A]} z_{a,i_3} = r''_{i_3} q''_{i_3} q_{i_3}$, and construct $G''(\{z_{a,b}\})$, where the new symbols are now $5'', \bar{5}'', 6'', \bar{6}'', 7'', \bar{7}''$.

The final text is

$$T := 1^{\bar{E}} G(\{x_{a,b}\}) 1'^{\bar{E}} G'(\{y_{a,b}\}) 1''^{\bar{E}} G''(\{z_{a,b}\}).$$

Applying Lemma 4.15 three times, we see that

$$\text{WRNA}(T) = 3\rho + \max_{i_1, i_2, i_3 \in [V^k]} \text{WRNA}(T(i_1, i_2, i_3)).$$

Since $\text{WRNA}(T(i_1, i_2, i_3)) \leq 6V + 3\bar{E}$ with equality if and only if $U(i_1) \cup U(i_2) \cup U(i_3)$ forms a $3k$ -clique, we obtain that $\text{WRNA}(T) \geq 3\rho + 6V + 3\bar{E}$ if and only if G contains a $3k$ -clique. This finishes the construction and proves the correctness.

Size Bounds Note that for each symbol $\sigma \in \{0, 1, \dots, 7\}$ we have a counterpart $\bar{\sigma}$, and both have three primed variants. Thus, the alphabet size is $|\Sigma| = 8 \cdot 2 \cdot 3 = 48$. Since $A = O(V + \bar{E}) = O(V^2)$ and $B = V^k$, the text length is $N = O(V^{k+2})$. Note that each $x_{a,b}, y_{a,b}$, and $z_{a,b}$ has total weight $W \leq 3$. Hence, the weight of the symbols introduced by the guarding $G(\cdot)$ is $8AW = O(A) = O(V^2)$. The following lemma analyzes the compressibility of the constructed text. We thus obtain all size bounds as claimed in the beginning of this proof.

Lemma 4.16. *The text T has an SLP \mathcal{T} of size $O(V^3)$.*

Proof. As in Lemmas 4.6 and 4.8, for any $a \in A$ there are SLPs for the strings $\bigcirc_{b \in [B]} \bar{6} x_{a,b}$, $\bigcirc_{b \in [B]} \bar{6} y_{a,b}$, and $\bigcirc_{b \in [B]} \bar{6} z_{a,b}$ of size $O(V)$. Indeed, any such string is equal to $\bigcirc_{i \in [V^k]} \bar{6} r_{e,i}$, $\bigcirc_{i \in [V^k]} \bar{6} p_{v,i}$, or $\bigcirc_{i \in [V^k]} \bar{6} q_{v,i}$, or their primed variants, for some $v \in V, e \in \bar{E}$. By definition of $r_{e,i}, p_{v,i}, q_{v,i}$, these strings are generated by $\text{Incl}_e^{(k)}$, $\text{Incl}_v^{(k)}$, and $\text{Adj}_v^{(k)}$, respectively, except that the terminals 0, 1 are replaced by some constant-length strings over $\{0, \bar{0}, \dots, 4, \bar{4}\}$. The final text T consists of $O(A) = O(V^2)$ strings of the form $\bigcirc_{b \in [B]} \bar{6} x_{a,b}$, $\bigcirc_{b \in [B]} \bar{6} y_{a,b}$, or $\bigcirc_{b \in [B]} \bar{6} z_{a,b}$, plus some very repetitive padding strings that can be compressed to length $O(\log V)$ by Observation 2.2. The bound follows. \square

It remains to prove Lemma 4.15 to finish the proof of Theorem 4.10.

Proof of Lemma 4.15. Let $x := G(\{x_{a,b}\})$ and fix $b \in [B]$. In every block $\bigcirc_{b \in [B]} \bar{6} x_{a,b}$ or $\bar{6}^B$ of x , we match the first b $\bar{6}$'s to the directly preceding 6's, and match the last $B - b$ $\bar{6}$'s to the directly succeeding 6's. At the beginning, this leaves $B - b$ $\bar{5}$'s to be matched to the first $\bar{5}$'s, and at the end this leaves b $\bar{7}$'s to be matched to the last $\bar{7}$'s. Since we match all $(A + 1)B$ $\bar{6}$'s and $B - b$ $\bar{5}$'s and b $\bar{7}$'s, the total weight of this matching is $(A + 1)B \cdot 8AW + (b + (B - b)) \cdot 4AW = \rho$. Note that this matching leaves all $x_{a,b}$ for $a \in A$ unmatched and uncovered, i.e., for no two matched symbols $x[i], x[j]$ we have that $x[i]$ is to the left of $x_{a,b}$ and $x[j]$ is to the right of $x_{a,b}$ in x . Hence, any solution to $\text{WRNA}(y_1 \circ (\bigcirc_{a \in [A]} x_{a,b}) \circ y_2)$ can be added to the pairs matched so far. This yields

$$\text{WRNA}(y_1 x y_2) \geq \rho + \max_{b \in [B]} \text{WRNA}\left(y_1 \circ \left(\bigcirc_{a \in [A]} x_{a,b}\right) \circ y_2\right).$$

For the other direction, consider an optimal matching R of $y_1 x y_2$, realizing $\text{WRNA}(y_1 x y_2)$. Write w_x for the total weight of pairs in R with both indices in x , and let $w_{x,y}$ be the total weight of pairs in R with one end in x and the other in y_1 or y_2 . Note that $w_x + w_{x,y} \geq \rho$, since otherwise, as shown above, we could replace the pairs of R incident with x to obtain $w_x = \rho$ and $w_{x,y} = 0$, yielding a higher total weight, which contradicts optimality of R .

Note that symbols in $x_{a,b}$ can only be matched to symbols in y_1 or y_2 , and the only possible matchings between x and y_1 or y_2 happen in the strings $x_{a,b}$. Let $Z \subseteq [A] \times [B]$ be the set of all pairs (a, b) such that $x_{a,b}$ contains at least one position matched by R . Consider first the case $Z = \emptyset$, so that $w_{x,y} = 0$. Denote by m_5, m_6, m_7 the number of matched symbols 5, 6, 7 in x . Note that each matched $\bar{5}$ and each matched $\bar{7}$ covers one 6. Hence, at most $B - m_5 + B - m_7 + AB$ 6's can be matched. Since the number of $\bar{6}$'s is $(A + 1)B$, we have $m_6 \leq \min\{B - m_5 + B - m_7 + AB, (A + 1)B\}$. We thus obtain an upper bound on w_x of

$$\begin{aligned} (m_5 + 2m_6 + m_7) \cdot 4AW &\leq (m_5 + m_7) \cdot 4AW + \min\{B - m_5 + B - m_7 + AB, (A + 1)B\} \cdot 8AW \\ &= \min\{2(A + 2)B - m_5 - m_7, 2(A + 1)B + m_5 + m_7\} \cdot 4AW. \end{aligned}$$

Optimizing over m_5, m_7 yields

$$w_x \leq (2A + 3)B = \rho.$$

Hence, in the current case $Z = \emptyset$ we have $w_{x,y} = 0$ and $w_x = \rho$, which yields

$$\text{WRNA}(y_1 x y_2) \leq \rho + \text{WRNA}(y_1 y_2) \leq \rho + \max_{b \in [B]} \text{WRNA}\left(y_1 \circ \left(\bigcirc_{a \in [A]} x_{a,b}\right) \circ y_2\right).$$

Now consider the remaining case $|Z| \geq 1$. Write $Z = \{x_{a_1, b_1}, \dots, x_{a_\ell, b_\ell}\}$, lexicographically sorted by (a, b) . Then we can bound $w_{x,y} \leq \ell \cdot W$, since the total weight of each $x_{a,b}$ is bounded from above by W .

In the following we bound w_x . Note that between x_{a_i, b_i} and $x_{a_{i+1}, b_{i+1}}$ the only symbols contributing to w_x are $\bar{6}$ and $\bar{6}$. We count $(a_{i+1} - a_i)B$ $\bar{6}$'s and $(a_{i+1} - a_i)B + b_{i+1} - b_i$ $\bar{6}$'s in this substring. Hence, this contribution is bounded from above by

$$\min\{(a_{i+1} - a_i)B, (a_{i+1} - a_i)B + b_{i+1} - b_i\} \cdot 8AW = (2(a_{i+1} - a_i)B + \min\{0, 2b_{i+1} - 2b_i\}) \cdot 4AW.$$

Using the identity $\min\{0, 2z\} = z - |z|$, we can rewrite this bound as

$$(2(a_{i+1} - a_i)B + b_{i+1} - b_i - |b_{i+1} - b_i|) \cdot 4AW.$$

We next analyze the contribution to w_x before x_{a_1, b_1} . We count $(a_1 - 1)B + b_1$ $\bar{6}$'s and $a_1 B$ $\bar{6}$'s as well as B $\bar{5}$'s and $\bar{5}$'s in this substring of x . Denote by m_5 the number of matched $\bar{5}$'s, and note that this covers m_5 $\bar{6}$'s from matching with $\bar{6}$'s. Hence, we can match at most $\min\{(a_1 - 1)B + b_1, a_1 B - m_5\}$ $\bar{6}$'s. Summing up the weights, we obtain an upper bound on the contribution to w_x before x_{a_1, b_1} of

$$m_5 \cdot 4AW + \min\{(a_1 - 1)B + b_1, a_1 B - m_5\} \cdot 8AW = \min\{2(a_1 - 1)B + 2b_1 + m_5, 2a_1 B - m_5\} \cdot 4AW.$$

Optimizing over m_5 , we obtain an upper bound of $((2a_1 - 1)B + b_1) \cdot 4AW$.

Lastly, we analyze the contribution to w_x after x_{a_ℓ, b_ℓ} . We count $(A - a_\ell + 2)B - b_\ell$ $\bar{6}$'s and $(A - a_\ell + 2)B$ $\bar{6}$'s as well as B $\bar{7}$'s and $\bar{7}$'s. Similarly to the last paragraph, when matching m_7 $\bar{7}$'s we obtain an upper bound on the contribution of

$$\begin{aligned} & m_7 \cdot 4AW + \min\{(A - a_\ell + 2)B - b_\ell, (A - a_\ell + 2)B - m_7\} \cdot 8AW \\ & = \min\{2(A - a_\ell + 2)B - 2b_\ell + m_7, 2(A - a_\ell + 2)B - m_7\} \cdot 4AW. \end{aligned}$$

Optimizing over m_7 yields an upper bound of $(2(A - a_\ell + 2)B - b_\ell) \cdot 4AW$.

Summing over all three cases, we obtain an upper bound on w_x of

$$\left((2a_1 - 1)B + b_1 + 2(A - a_\ell + 2)B - b_\ell + \sum_{i=1}^{\ell-1} (2(a_{i+1} - a_i)B + b_{i+1} - b_i - |b_{i+1} - b_i|) \right) \cdot 4AW.$$

Note that all a_i 's and almost all b_i 's cancel as they form telescoping sums. What remains is

$$w_x \leq \left(-B + 2(A + 2)B - \sum_{i=1}^{\ell} |b_{i+1} - b_i| \right) \cdot 4AW = \rho - 4AW \sum_{i=1}^{\ell-1} |b_{i+1} - b_i|.$$

In combination with the inequalities $w_{x,y} \leq \ell W$ and $w_x + w_{x,y} \geq \rho$ shown above, we obtain

$$\sum_{i=1}^{\ell-1} |b_{i+1} - b_i| \leq \frac{\ell}{4A}.$$

Note that we have $|b_{i+1} - b_i| = 0$ for at most $A - 1$ i 's, since $b_{i+1} = b_i$ implies $a_{i+1} > a_i$. This yields

$$\sum_{i=1}^{\ell-1} |b_{i+1} - b_i| \geq \ell - 1 - (A - 1) = \ell - A.$$

Together with the upper bound, we obtain $\ell - A \leq \ell/(4A) \leq \ell/2$, which yields $\ell \leq 2A$. Hence, we have

$$\sum_{i=1}^{\ell-1} |b_{i+1} - b_i| \leq \frac{\ell}{4A} \leq 1/2 < 1,$$

which implies that $b_{i+1} = b_i$ for all i . Let $b := b_1 = \dots = b_\ell$. Then R matches only the strings $x_{a,b}$ for $a \in A$, among all strings in X . Since we showed $w_x \leq \rho$, we indeed obtain

$$\text{WRNA}(T) \leq \rho + \max_{b \in [B]} \text{WRNA}\left(y_1 \circ \left(\bigcirc_{a \in [A]} x_{a,b}\right) \circ y_2\right). \quad \square$$

\square

5 Disjointness, Hamming Distance, and Subsequence

In this section we consider the following three problems on compressed sequences. In all problems we are given SLPs \mathcal{T} and \mathcal{P} of size n and m , representing a text $T = \text{eval}(\mathcal{T})$ of length N and a pattern $P = \text{eval}(\mathcal{P})$ of length M .

Problem 5.1 (Disjointness). *Given two compressed sequences \mathcal{T} and \mathcal{P} of equal decompressed lengths $N = M$ over alphabet $\{0, 1\}$, decide whether there is a position such that both sequences have symbol 1 at that position, i.e., whether $T[i] = P[i] = 1$ holds for some i .*

Problem 5.2 (Hamming Distance). *Given two compressed sequences \mathcal{T} and \mathcal{P} of equal decompressed lengths $N = M$, output $\text{Hamming}(P, T) = |\{i | P[i] \neq T[i]\}|$. That is, output the number of positions where the decompressed sequences differ.*

Problem 5.3 (Subsequence). *Given two compressed sequences \mathcal{T} and \mathcal{P} of decompressed length $N \geq M$, decide whether the pattern sequence P is a subsequence of the text sequence T .*

We note that in the uncompressed setting all three problems have linear time trivial algorithms. This immediately implies that all three problems can be solved in time $O(N)$ by decompressing the sequences and running the trivial algorithms. Below we show that this running time is not optimal and can be improved for all three problems for sufficiently compressible strings. Furthermore, we show conditional lower bounds for the three problems assuming the Combinatorial k -Clique conjecture, k -SUM conjecture, and Strong k -SUM conjecture (see Section 2.1 for definitions). We were, however, not able to establish matching upper and lower bounds and we leave it as an open problem to close the gap.

Known Lower Bounds from Classic Complexity Theory In [49] it was shown that the Hamming Distance problem is $\#\text{P}$ -complete and thus a polynomial time $(nm)^{O(1)}$ algorithm for it is unlikely to exist. Lohrey [52] showed that the Subsequence problem is at least as hard as PP and is contained in PSPACE . It is conjectured that the subsequence problem is PSPACE -complete [53]. Note that the class PP contains computationally very difficult problems. In particular, Toda's theorem states that the entire polynomial hierarchy PH is contained in P^{PP} .

We can easily check that the Disjointness problem is in NP . A variant of our Theorem 5.10 below implies that the Subset Sum problem can be reduced to the Disjointness problem and thus Disjointness is in fact NP -complete.

5.1 Algorithms

We start this section by showing a simple algorithm for the Subsequence problem that runs in time $O((n|\Sigma| + M) \log N)$ (see Theorem 5.4). An algorithm with very similar guarantees was obtained in [12]. Note that in a natural setting, namely when $|\Sigma| \leq O(1)$, $n \leq M$ and $N \leq M^{O(1)}$, the algorithm runs in time $\tilde{O}(M)$. That is, we do not need to decompress the text sequence to be able to solve the Subsequence problem.

In Theorems 5.5 and 5.6 below we show $O(\max(m, n)^{1.5} \cdot N^{0.6})$ time algorithms for the Hamming Distance and Subsequence problems, respectively. We observe that both running times that we obtain for the Subsequence problem are incomparable. Finally, by Theorem 5.7 from Section 5.2, the Disjointness problem can be reduced to the Subsequence problem. This implies an $O(\max(m, n)^{1.5} \cdot N^{0.6})$ time algorithm for the Disjointness problem. To the best of our knowledge these upper bounds are new.

Theorem 5.4. *The Subsequence problem can be solved in time $O((n|\Sigma| + M) \log N)$.*

Proof. We start by decompressing the pattern sequence \mathcal{P} in $O(M)$ time. To decide whether P is a subsequence of the text sequence T , for $i = 1, \dots, M$ (in this order) we will find the smallest $j \in \{1, \dots, N\}$ such that $P[1..i]$ (the prefix of the decompressed pattern of length i) is a subsequence of $T[1..j]$. In the rest of the proof we will describe how to do this efficiently.

We start by transforming the compressed text \mathcal{T} into an AVL-grammar of size $O(n \log N)$ and depth $O(\log N)$ according to Theorem 2.1. This takes $O(n \log N)$ time. Additionally, for every alphabet symbol $\sigma \in \Sigma$ and every non-terminal T_i of the AVL-grammar, we decide whether the sequence produced by the non-terminal T_i contains the symbol σ . For every symbol, this can be done in $O(n \log N)$ time. Since the size of the alphabet is $|\Sigma|$, this takes $O(n|\Sigma| \log N)$ total time.

Given an index $i = 1, \dots, M$, suppose that we know the smallest index $j \in \{1, \dots, N\}$ such that $P[1..i]$ is a subsequence of $T[1..j]$. We will show how to find the smallest $j' > j$ such that $P[1..i + 1]$ is a subsequence of $T[1..j']$. The required running time will follow since we will be able to do this in $O(\log N)$ time for every index i . We find the smallest $j' > j$ in two steps. In the first step we traverse the parse tree bottom-up from the symbol $T[j]$ until the current node has $T[j]$ in the left subtree and the right subtree contains symbol $P[i + 1]$. In the second step we go to the right subtree and then keep going to the left-most child that contains the symbol $P[i + 1]$. Since the height of the parse tree is $O(\log N)$, this takes $O(\log N)$ time. This finishes the description of the algorithm. Note that we did not decompress the text sequence T in this process. \square

Theorem 5.5. *The Hamming Distance problem can be solved in time*

$$\tilde{O}\left(\max(m, n)^{2-1/\log_2(2\varphi)} \cdot N^{1/\log_2(2\varphi)}\right) = \tilde{O}\left(\max(m, n)^{1.409\dots} \cdot N^{0.592\dots}\right),$$

where $\varphi = \frac{1+\sqrt{5}}{2}$ is the golden ratio.

Proof. Let P_1, P_2, \dots, P_m be the SLP \mathcal{P} corresponding to the decompressed pattern sequence P and let T_1, T_2, \dots, T_n be the SLP \mathcal{T} corresponding to the decompressed text sequence T . We assume that the decompressed length of the sequences P and T is $|P| = |T| = N$.

By Theorem 2.1 we can assume that P_1, P_2, \dots, P_m and T_1, T_2, \dots, T_n are AVL-grammars. This increases the running time by a factor of at most polylog N , which is hidden in the $\tilde{O}(\cdot)$ notation. Fix an $i = 1, \dots, m$ and consider the sequence $\text{eval}(P_i)$ with the corresponding parse tree

of height $\text{depth}(P_i)$. Then one can verify that the length of the sequence is bounded from above by $|\text{eval}(P_i)| \leq 2^{\text{depth}(P_i)}$ and from below by

$$|\text{eval}(P_i)| \geq F_{\text{depth}(P_i)} \geq \Omega\left(\varphi^{\text{depth}(P_i)}\right), \quad (6)$$

where $F_{\text{depth}(P_i)}$ is the $\text{depth}(P_i)$ -th Fibonacci number and φ is the golden ratio [26]. Analogous properties hold for T_j for $j = 1, \dots, n$.

For every P_i and T_j we precompute the length of $\text{eval}(P_i)$ and $\text{eval}(T_j)$, respectively. We define the function

$$\text{Ham}(P_i, T_j, d) := \sum_{\substack{r \in \{1, \dots, |\text{eval}(P_i)|\}, \\ r+d \in \{1, \dots, |\text{eval}(T_j)|\}}} [\text{eval}(P_i)_r \neq \text{eval}(T_j)_{r+d}],$$

where d is a negative or a non-negative integer. In other words, $\text{Ham}(P_i, T_j, d)$ is equal to the Hamming distance between T_j and a shifted P_i (by d symbols to the right if $d > 0$ and by $|d|$ symbols to the left otherwise), where we consider only the symbols that have aligned counterparts. Clearly, we can solve the Hamming Distance problem by outputting $\text{Hamming}(P, T) = \text{Ham}(P_m, T_n, 0)$.

A simple algorithm for computing the Hamming distance is the following recursive method. Assume that the sequence $\text{eval}(P_i)$ is longer than the sequence $\text{eval}(T_j)$, and $P_i \rightarrow P_{\ell(i)}, P_{r(i)}$. Then

$$\text{Ham}(P_i, T_j, d) = \text{Ham}(P_{\ell(i)}, T_j, d) + \text{Ham}(P_{r(i)}, T_j, d + |\text{eval}(P_{\ell(i)})|).$$

Otherwise, if the sequence $\text{eval}(T_j)$ is longer and $T_j \rightarrow T_{\ell'(j)}, T_{r'(j)}$, then

$$\text{Ham}(P_i, T_j, d) = \text{Ham}(P_i, T_{\ell'(j)}, d) + \text{Ham}(P_i, T_{r'(j)}, d - |\text{eval}(T_{\ell'(j)})|).$$

Clearly, for any recursive subproblem where the argument d is such that no symbols get aligned, we can immediately return 0. When P_i or T_j encode a single symbol, we compute their Hamming distance in a constant time.

We use this recursive algorithm with memoization, i.e., if we call the same inputs twice, then we return the stored answer.

Running Time We crucially use the fact that we split the longer text in each step, and property 6. Both together imply that

$$|\text{eval}(T_j)| \geq \Omega\left(|\text{eval}(P_i)|^{\log_2 \varphi}\right) = \Omega\left(|\text{eval}(P_i)|^{0.694\dots}\right)$$

for each call $\text{Ham}(P_i, T_j, d)$. We bound the running time by counting for each P_i how many different calls there are of the form $\text{Ham}(P_i, T_j, d)$ with $|\text{eval}(T_j)| \leq |\text{eval}(P_i)|$. The running time corresponding to the calls with $|\text{eval}(T_j)| > |\text{eval}(P_i)|$ can be analyzed analogously. Note that $|\text{eval}(T_j)| \leq |\text{eval}(P_i)|$ implies $|d| \leq O(|\text{eval}(P_i)|)$, as larger shifts immediately give answer 0. Let $0 < \alpha < 1$ to be fixed later. If $|\text{eval}(P_i)| < N^\alpha$ we can thus bound the contribution of P_i to the running time by nN^α (there are n T_j 's and N^α possible offsets d). Otherwise, if $|\text{eval}(P_i)| \geq N^\alpha$, then $|\text{eval}(T_j)| \geq N^{\alpha \log_2 \varphi}$, and thus there are at most $N^{1-\alpha \log_2 \varphi}$ calls to such T_j in the parse tree for $T = \text{eval}(T_n)$. Thus, there are at most this many calls $\text{Ham}(P_i, T_j, d)$, so the contribution of P_i to the running time is at most $N^{1-\alpha \log_2 \varphi}$. Summed over all m different P_i 's the total running time is bounded by $O(m(nN^\alpha + N^{1-\alpha \log_2 \varphi}))$. Minimizing over α gives the running time $O(m \cdot n^{1-1/\log_2(2\varphi)} \cdot N^{1/\log_2(2\varphi)})$. The running time corresponding to the calls with $|\text{eval}(T_j)| > |\text{eval}(P_i)|$ can be similarly bounded by $O(n \cdot m^{1-1/\log_2(2\varphi)} \cdot N^{1/\log_2(2\varphi)})$. It remains to observe that the total running time is bounded by $O(\max(m, n)^{2-1/\log_2(2\varphi)} \cdot N^{1/\log_2(2\varphi)})$ as required. \square

Theorem 5.6. *The Subsequence problem can be solved in time*

$$\tilde{O}\left(\max(m, n)^{2-1/\log_2(2\varphi)} \cdot N^{1/\log_2(2\varphi)}\right) = \tilde{O}\left(\max(m, n)^{1.409\dots} \cdot N^{0.592\dots}\right),$$

where $\varphi = \frac{1+\sqrt{5}}{2}$ is the golden ratio.

Proof. The algorithm follows a similar recursive method as in Theorem 5.5. As above, we assume that the SLPs \mathcal{P} and \mathcal{T} are AVL-grammars.

For non-terminals P_i and T_j and an integer d we define the function $\text{Subseq}(P_i, T_j, d)$. If $d \geq 0$, then we assume that we already matched a prefix of $\text{eval}(P_i)$ of length d (the prefix is a subsequence of an earlier part of the text) and our goal is to match the rest of $\text{eval}(P_i)$ with $\text{eval}(T_j)$. On the other hand, if $d < 0$, then we assume that we already matched a prefix of $\text{eval}(T_j)$ of length $|d|$ (a previous part of the pattern is a subsequence of the prefix) and our goal is to match $\text{eval}(P_i)$ to the rest of $\text{eval}(T_j)$. The function returns an integer as follows. Let d' be the length of the longest prefix of $\text{eval}(P_i)$ that can be matched to $\text{eval}(T_j)$. (If $d \geq 0$, then we match only the remainder of $\text{eval}(P_i)$ to $\text{eval}(T_j)$. If $d < 0$, then we match $\text{eval}(P_i)$ to the remainder of $\text{eval}(T_j)$.) If $d' < |\text{eval}(P_i)|$, that is, we cannot match entire $\text{eval}(P_i)$ to $\text{eval}(T_j)$, then the function $\text{Subseq}(P_i, T_j, d)$ returns d' . Otherwise, if $d' = |\text{eval}(P_i)|$, the function returns the length of the shortest prefix of (the remainder of) $\text{eval}(T_j)$ that can be matched to (the remainder of) $\text{eval}(P_i)$.

Given the description of the function, the recursive implementation of it is straightforward and is described below. To evaluate $\text{Subseq}(P_i, T_j, d)$, we consider three cases.

Case 1 P_i or T_j represents a single symbol. The problem is trivial to solve in this case.

Case 2 $|\text{eval}(P_i)| \geq |\text{eval}(T_j)|$. Let $P_i \rightarrow P_{\ell(i)}, P_{r(i)}$ be the SLP rule corresponding to P_i . If $d \geq |\text{eval}(P_{\ell(i)})|$, then the function returns $\text{Subseq}(P_{r(i)}, T_j, d - |\text{eval}(P_{\ell(i)})|)$, which we compute recursively. If, on the other hand, $d < |\text{eval}(P_{\ell(i)})|$, we recursively compute $d' := \text{Subseq}(P_{\ell(i)}, T_j, d)$ and return d' if $d' \geq 0$ or return $\text{Subseq}(P_{r(i)}, T_j, d')$ if $d' < 0$.

Case 3 $|\text{eval}(P_i)| < |\text{eval}(T_j)|$. This case is similar to the previous one. Let $T_j \rightarrow T_{\ell'(j)}, T_{r'(j)}$ be the SLP rule corresponding to T_j . If $-d > |\text{eval}(T_{\ell'(j)})|$, we return $\text{Subseq}(P_i, T_{r'(j)}, -d - |\text{eval}(T_{\ell'(j)})|)$, which we compute recursively. Otherwise, we define $d' := \text{Subseq}(P_i, T_{\ell'(j)}, d)$ and return d' if $d' < 0$ or return $\text{Subseq}(P_i, T_{r'(j)}, d')$ if $d' \geq 0$.

The correctness of the algorithm follows from the description and the definition of the function $\text{Subseq}(P_i, T_j, d)$. The running time analysis is similar to Theorem 5.5 and we omit it. \square

5.2 Lower Bounds

In this section we show conditional lower bounds for the Disjointness, Hamming Distance and Subsequence problems. First, we show that the Disjointness problem can be reduced to the Subsequence problem (Theorem 5.7) and to the Hamming Distance problem (Theorem 5.8). Thus, any algorithmic improvement for the latter two problems implies a faster algorithm for the Disjointness problem. Alternatively, we can think about the Disjointness problem as the core hard problem explaining hardness for the two other problems. Second, we show a matching $N^{1-o(1)}$ lower bound for combinatorial algorithms for the Subsequence problem in the setting where $N \approx M \approx n^2 \approx m^2$. We use the combinatorial k -Clique conjecture to establish this hardness. Finally, we use the k -SUM

conjecture (Conjecture 2.8) for all three aforementioned problems. The lower bounds that we show are not tight. We show that assuming a stronger version of the k -SUM conjecture (Conjecture 2.9) allows us to get higher lower bounds, but still not matching.

Theorem 5.7. *The Disjointness problem can be reduced to the Subsequence problem. The reduction loses at most constant factors in the length of compressed and decompressed sequences.*

Proof. Let P and T be two binary sequences, forming an instance of the Disjointness problem. We construct a sequence P' from P by replacing every symbol 0 with symbol “0” and every symbol 1 with two symbols “10”. Similarly, we construct a sequence T' from T by replacing every symbol 0 with two symbols “10” and every symbol 1 with “0”.

The resulting sequences P' and T' are compressible similarly as P and T . We can check that P' is a subsequence of T' if and only if we have $P[i] = 0$ or $T[i] = 0$ for all i . This completes the reduction. \square

Theorem 5.8. *The Disjointness problem can be reduced to the Hamming Distance problem. The reduction loses at most constant factors in the length of compressed and decompressed sequences.*

Proof. Let P and T be two binary sequences, forming an instance of the Disjointness problem. We construct a sequence P' from P by replacing every symbol 0 with three symbols “011” and every symbol 1 with three symbols “000”. Similarly, we construct a sequence T' from T by replacing every symbol 0 with “001” and every symbol 1 with “111”.

These four gadget sequences have Hamming distance 1 for all pairs except when both original symbols are 1. In this case the Hamming distance between the two gadgets is 3. We conclude that $\text{Hamming}(P', T') > N = |P| = |T|$ if and only if there exists i with $P[i] = T[i] = 1$. This concludes the reduction. \square

Theorem 5.9. *The Subsequence problem has no combinatorial $O(N^{1-\varepsilon})$ time algorithm for any $\varepsilon > 0$ in the setting $N = \Theta(M) = \Theta(n^2) = \Theta(m^2)$ and $|\Sigma| = O(N^\varepsilon)$, assuming the combinatorial k -Clique conjecture.*

Proof. The reduction will rule out combinatorial algorithms with running time $N^{1-\varepsilon}$ by using the Combinatorial k -Clique conjecture 2.7 with $k = O(1/\varepsilon)$. Let $k \geq 4$ be even, and let $G = (V, E)$ be an instance of k -Clique. In the following we will construct an equivalent instance of the Subsequence problem, i.e., a text $T = \text{eval}(\mathcal{T})$ and a pattern $P = \text{eval}(\mathcal{P})$, satisfying $N = |T| = O(V^{k+1})$, $n = |\mathcal{T}| = O(V^{(k/2)+1})$, $M = |P| = O(V^k)$, and $m = |\mathcal{P}| = O(V^{k/2})$. The alphabet size will be $|\Sigma| = O(V)$. By a simple padding¹⁰, we can then ensure that $N, M = \Theta(V^{k+2})$ and $n, m = \Theta(V^{k/2+1})$, so that indeed $N = \Theta(M) = \Theta(n^2) = \Theta(m^2)$, and we have $|\Sigma| = O(N^\varepsilon)$ for any $k \geq 1/\varepsilon$. Finally, a combinatorial $O(N^{1-\varepsilon})$ algorithm for the Subsequence problem in this setting would yield a combinatorial algorithm for k -Clique in time $O(V^{(k+2)(1-\varepsilon)}) = O(V^{k(1-\varepsilon/2)})$ for any $k \geq 4/\varepsilon$, contradicting the combinatorial k -Clique conjecture.

We first construct clique gadgets and then the pattern and the text. The alphabet will be $\Sigma = V \cup \{\#, \$\}$.

¹⁰Specifically, let \dagger be a fresh symbol and add $\dagger^{V^{k+2}}$ as a prefix to T and P . Compress this string $\dagger^{V^{k+2}}$ to length $V^{k/2+1}$ by writing it as $(\dagger^{V^{k/2+1}})^{V^{k/2+1}}$ and using Observation 2.2.

Construction of the clique gadgets CG Given a $(k/2)$ -clique $C = \{v_1, \dots, v_{k/2}\}$, we construct the clique gadget $CG(C)$ as:

$$CG(C) = (v_1 v_2 \dots v_{k/2} \#)^{k/2}.$$

That is, we write down the labels of the vertices (in increasing order), put “#” at the end and repeat the resulting sequence $k/2$ times.

Construction of the clique gadgets CG' Given a $(k/2)$ -clique $C' = \{u_1, \dots, u_{k/2}\}$, we construct $CG'(C')$ as:

$$CG'(C') = \text{Neighbors}(u_1) \# \text{Neighbors}(u_2) \# \dots \# \text{Neighbors}(u_{k/2}) \#$$

where $\text{Neighbors}(u)$ lists all neighbors of vertex u in increasing order.

We can check that for any $(k/2)$ -cliques C, C' , $CG(C)$ is a subsequence of $CG'(C')$ if and only if $C \cup C'$ forms a k -clique.

Construction of the sequence Z We construct Z as:

$$Z := (L \#)^{k/2},$$

where L is the sequence containing all V vertices in the graph in increasing order. We can verify the any clique gadget $CG(C)$ is a subsequence of Z .

Construction of the Pattern The pattern consists of clique gadgets as follows. Enumerate all $(k/2)$ -cliques C_1, \dots, C_Q with $Q \leq V^{k/2}$ in G . The pattern sequence P is constructed as:

$$P := (CG(C_1) \$ CG(C_2) \$ \dots \$ CG(C_Q) \$)^Q.$$

That is, we concatenate the Q clique gadgets $CG(C_1), \dots, CG(C_Q)$ in one sequence and put “\$” after every gadget, and repeat the resulting sequence Q times. Note that the symbol “\$” does not appear in any clique gadget.

Construction of the Text The text is somewhat similar to the pattern, defined by:

$$T := (CG'(C_1) \$ Z \$)^Q (CG'(C_2) \$ Z \$)^Q \dots (CG'(C_{Q-1}) \$ Z \$)^Q (CG'(C_Q) \$ Z \$)^{Q-1} CG'(C_Q) \$.$$

Correctness The pattern consists of Q^2 clique gadgets with the symbol \$ in between any two of them. The text consist of Q^2 cliques gadgets with the sequence $\$Z\$$ in between any two of them. Since there are only $Q^2 - 1$ Z 's in the text, we cannot match all clique gadgets of the pattern to Z 's in the text. Hence, if P is a subsequence of T , then at least one clique gadget $CG(C_i)$ is a subsequence of $CG'(C_j)$ for some i, j . This happens only if $C_i \cup C_j$ form as k -clique in G .

For the other direction, we show that if G contains a k -clique, so that there are i, j with $C_i \cup C_j$ forming a k -clique, implying that $CG(C_i)$ is a subsequence of $CG'(C_j)$, then the pattern is a subsequence of the text. Indeed, let $q = j \cdot Q + i$. The q -th clique gadget in the pattern is $CG(C_i)$ and the q -th clique gadget in the text is $CG'(C_j)$. We match all clique gadgets before the q -th one as well as after the q -th one to Z 's, and we match $CG(C_i)$ to $CG'(C_j)$. This shows that P is a subsequence of T .

Since $|L| = V$, $Q \leq V^{k/2}$, and k is a constant, the length bounds $N = O(V^{k+1})$ and $M = O(V^k)$ are immediate. Using Observation 2.2 to compress strings of the form X^Q to size $O(|X| + \log Q)$, we also immediately obtain $n = O(V^{k/2+1})$ and $m = O(V^{k/2})$. This finishes the proof. \square

Theorem 5.10. *Let $k \geq 1$ be an integer. Consider the Disjointness problem with $N = M = \Theta(n^{4k+1}) = \Theta(m^{4k+1})$. Solving the Disjointness problem in this setting requires $N^{\frac{1}{4} + \frac{3}{16k+4} - o(1)}$ time assuming the $(2k+1)$ -SUM conjecture.*

Theorem 5.11. *Let $k \geq 1$ be an integer. Consider the Disjointness problem with $N = M = \Theta(n^{3k+1}) = \Theta(m^{3k+1})$. Solving the Disjointness problem in this setting requires $N^{\frac{1}{3} + \frac{2}{9k+3} - o(1)}$ time assuming the Strong $(2k+1)$ -SUM conjecture.*

By Theorems 5.7 and 5.8, the same kind of hardness holds for the Subsequence and Hamming Distance problems.

Proof of Theorems 5.10 and 5.11. Let $k \geq 1$ be an integer and let $A \subseteq \{0, 1, \dots, R-1, R\}$ be an instance of the $(2k+1)$ -SUM problem with $|A| = r$ and target sum t . Without loss of generality, R is divisible by $k+1$ and t is divisible by k . We define the set $B := \{\frac{t}{k} + R - a \mid a \in A\}$ and the set $C := \{\frac{Rk}{k+1} + a \mid a \in A\}$. We can verify that there exist $b_1, \dots, b_k \in B$ and $c_1, \dots, c_{k+1} \in C$ with $b_1 + \dots + b_k = c_1 + \dots + c_{k+1}$ if and only if there exist $a_1, \dots, a_{2k+1} \in A$ with $a_1 + \dots + a_{2k+1} = t$. We note that $B, C \subseteq \{1, 2, \dots, R'\}$ for $R' := 2R$.

In $O(r \log r)$ time we will construct an instance to the Disjointness problem with the following properties.

- Pattern $P = \text{eval}(\mathcal{P})$ is constructed from the set B and has length $M = R' \cdot r^{2k}$ and compressed size $m = O(r \log r)$,
- Text $T = \text{eval}(\mathcal{T})$ is constructed from the set C and has length $N = R' \cdot r^{2k}$ and compressed size $n = O(r \log r)$,
- There exists i such that $P[i] = T[i] = 1$ if and only if there exist $b_1, \dots, b_k \in B$ and $c_1, \dots, c_{k+1} \in C$ with $b_1 + \dots + b_k = c_1 + \dots + c_{k+1}$.

Simply padding allows us to increase the text length and pattern length to $R' r^{2k} \log^{k'} r$ for any $k' \geq 0$, and to achieve $n, m = \Theta(r \log r)$. Setting $R = r^{2k+1}$, we thus have $N = M = 2r^{4k+1} \log^{4k+1} r = \Theta(n^{4k+1}) = \Theta(m^{4k+1})$. Any $O(N^{1/4+3/(16k+4)-\varepsilon}) = O(N^{(k+1-\varepsilon)/(4k+1)})$ time algorithm for Disjointness would now imply an algorithm for $(2k+1)$ -SUM in time $O((r \log r)^{k+1-\varepsilon}) = O(r^{k+1-\varepsilon/2})$, contradicting the $(2k+1)$ -SUM conjecture (Conjecture 2.8). This proves Theorem 5.10. Similarly, setting $R = r^{k+1}$ and using the Strong $(2k+1)$ -SUM conjecture (Conjecture 2.9) we obtain Theorem 5.11.

In the remainder of the proof we present the promised construction.

Without loss of generality, we have $R' > 10k \cdot \max(B \cup C)$.

Construction of the Pattern We define the pattern as

$$P := \left(\bigcirc_{b_1, \dots, b_k \in B} 0^{b_1 + \dots + b_k} 1 0^{R' - (b_1 + \dots + b_k) - 1} \right)^{r^k},$$

where the \bigcirc goes over all tuples $(b_1, \dots, b_k) \in B^k$ in lexicographic order. That is, P consists of r^k repetitions of a sequence Z of length $R' \cdot r^k$. The sequence Z consists of sequences Z_1, \dots, Z_{r^k} , corresponding to k -tuples $(b_1, \dots, b_k) \in B^k$. Each sequence Z_i has length R' , and the sequence Z_i corresponding to tuples (b_1, \dots, b_k) has 0's everywhere except at position $b_1 + \dots + b_k + 1$.

Construction of the Text We define the text as

$$T := \bigcirc_{c_1, \dots, c_k \in C} \left(Y(c_1, \dots, c_k) \right)^{r^k}, \quad (7)$$

where $Y(c_1, \dots, c_k)$ is a string of length R' with $Y(c_1, \dots, c_k)[j+1] = 1$ if $j \in \{c + c_1 + \dots + c_k \mid c \in C\}$, and $Y(c_1, \dots, c_k)[j+1] = 0$ otherwise.

Analysis Note that there is an index i with $P[i] = T[i] = 1$ if and only if there exist $b_1, \dots, b_k \in B$ and $c_1, \dots, c_{k+1} \in C$ with $b_1 + \dots + b_k = c_1 + \dots + c_{k+1}$. Hence, correctness of the reduction can be easily verified. The length $N = M = R' r^{2k}$ is immediate. It remains to show that the pattern and the text are compressible.

Compressing the Pattern Since $P = Z^{r^k}$, by Observation 2.2 it suffices to compress Z . We construct the sequence Z inductively. We write $B = \{B_1, \dots, B_r\}$. We define $S_0 \rightarrow 1$ to be a non-terminal generating a sequence of length 1 containing a single symbol 1. For $i \in [k]$ we define the non-terminal S_i as follows:

$$S_i \rightarrow \left(\bigcirc_{w=1}^{r-1} 0^{B_w} S_{i-1} 0^{R' r^{i-1} - B_w - |\text{eval}(S_{i-1})|} \right) \circ 0^{B_r} S_{i-1}. \quad (8)$$

Finally, we set $S \rightarrow S_k \circ 0^{R' r^k - |S_k|}$. Here the right hand side contains more than two SLP non-terminals, but using Observation 2.2 it is easy to convert this into a proper SLP of size $O(r \log r)$ as required. It remains to check that $Z = \text{eval}(S)$, i.e., $\text{eval}(S) = \bigcirc_{b_1, \dots, b_k \in B} 0^{b_1 + \dots + b_k} 1 0^{R' - (b_1 + \dots + b_k) - 1}$. Indeed, a straightforward induction shows that we constructed S_i , $i \in [k]$ such that

$$\text{eval}(S_i) \circ 0^{R' r^i - |\text{eval}(S_i)|} = \bigcirc_{b_1, \dots, b_i \in B} 0^{b_1 + \dots + b_i} 1 0^{R' - (b_1 + \dots + b_i) - 1}.$$

The induction step is performed by using the derivation rule (8).

Compressing the Text Let W be a string of length R' consisting only of 0's except $W[j+1] = 1$ for any $j \in C$. We define an SLP non-terminal Y' that generates the shortest prefix of W containing all 1's of W . We set

$$Y_0 \rightarrow \left(Y' 0^{R' - |\text{eval}(Y')|} \right)^{r^k - 1} Y'.$$

Note that $\text{eval}(Y') 0^{R' - |\text{eval}(Y')|} = W$. Hence, Y_0 generates the string W^{r^k} where we removed the longest suffix of 0's. We write $C = \{C_1, \dots, C_r\}$.

For $i = 1, \dots, k$ we define sequence Y_i as follows:

$$Y_i \rightarrow \left(\bigcirc_{w=1}^{r-1} 0^{C_w} Y_{i-1} 0^{R' r^{k+i-1} - C_w - |\text{eval}(Y_{i-1})|} \right) \circ 0^{C_r} Y_{i-1}. \quad (9)$$

Finally, we set $\mathcal{T} \rightarrow Y_k \circ 0^{R'r^{2k}-|\text{eval}(Y_k)|}$. It is easy to verify that the size of the above SLP \mathcal{T} is $O(r \log r)$. It remains to show that $\text{eval}(\mathcal{T}) = T$ as in (7). That is, we want to show that $\text{eval}(\mathcal{T}) = \bigcirc_{c_1, \dots, c_k \in C} \left(Y(c_1, \dots, c_k) \right)^{r^k}$. This follows by a straightforward induction. We can check that for $i = 0, 1, \dots, k$ we have

$$\text{eval}(Y_i) \circ 0^{R'r^{k+i}-|\text{eval}(Y_i)|} = \bigcirc_{c_1, \dots, c_i \in C} \left(Y(c_1, \dots, c_i) \right)^{r^k}.$$

The induction step is performed by using the derivation rule (9). □

6 Conclusion

With this paper we started the fine-grained complexity of analyzing compressed data, thus providing lower bound tools for a practically highly relevant area. We focused on the most basic problems on strings, leaving many other stringology problems for future work. Besides strings, there is a large literature on grammar-compressed other forms of data, e.g. graphs. It would be interesting to apply our framework and classify the important problems in these contexts as well.

Specifically, we leave the following open problems.

- Determine the optimal running time for the Disjointness, Hamming Distance, and Subsequence problems.
- Generalize our lower bound for LCS to Edit Distance.
- For NFA Acceptance we obtained tight bounds in case of a potentially dense automaton with q states and up to $O(q^2)$ transitions. Prove tight bounds for the case of sparse automata with $O(q)$ transitions.
- For large (i.e. superconstant) alphabet size, some bounds given in this paper are not tight, most prominently for Generalized Pattern Matching, Substring Hamming Distance, and Pattern Matching with Wildcards. Determine the optimal running time in this case.
- For all lower bounds presented in this paper, check whether they can be improved to work for binary strings.

Acknowledgements

This paper would not have been possible without Oren Weimann and *Schloss Dagstuhl*. Inspired by a Dagstuhl seminar on Compressed Pattern Matching in October, and while attending a Dagstuhl seminar on Fine-Grained Complexity in November, Oren asked in the open problems session whether SETH can explain the lack of $O((nN)^{1-\varepsilon})$ algorithms for problems like LCS on compressed strings. Later, in January, three of the authors of this paper attended a Dagstuhl seminar on Parameterized Complexity and made key progress towards the results of this work. Part of the work was also performed while visiting the Simons Institute for the Theory of Computing, Berkeley, CA. We thank Paweł Gawrychowski for helpful comments.

A.A. was supported by Virginia Vassilevska Williams' NSF Grants CCF-1417238 and CCF-1514339, and BSF Grant BSF:2012338. Arturs Backurs was supported by an IBM PhD Fellowship, the NSF and the Simons Foundation. While performing part of this work, M. Künnemann was affiliated with University of California, San Diego.

References

- [1] A. Abboud, A. Backurs, and V. Vassilevska Williams. If the current clique algorithms are optimal, so is Valiant's parser. In *Proc. 56th IEEE Annual Symposium on Foundations of Computer Science (FOCS'15)*, pages 98–117. IEEE, 2015.
- [2] A. Abboud, A. Backurs, and V. Vassilevska Williams. Tight Hardness Results for LCS and other Sequence Similarity Measures. In *Proc. 56th IEEE Annual Symposium on Foundations of Computer Science (FOCS'15)*, pages 59–78, 2015.
- [3] A. Abboud, T. D. Hansen, V. Vassilevska Williams, and R. Williams. Simulating branching programs with edit distance and friends: or: a polylog shaved is a lower bound made. In *Proc. 48th Annual ACM Symposium on Theory of Computing (STOC'16)*, pages 375–388, 2016.
- [4] A. Abboud, V. Vassilevska Williams, and O. Weimann. Consequences of faster sequence alignment. In *Proc. 41st International Colloquium on Automata, Languages, and Programming (ICALP'14)*, pages 39–51, 2014.
- [5] A. Abboud, R. Williams, and H. Yu. More applications of the polynomial method to algorithm design. In *Proc. 26th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'15)*, pages 218–230, 2015.
- [6] A. Amir, T. M. Chan, M. Lewenstein, and N. Lewenstein. On hardness of jumbled indexing. In *Proc. 41st International Colloquium on Automata, Languages, and Programming (ICALP'14)*, pages 114–125. Springer, 2014.
- [7] A. Apostolico, G. M. Landau, and S. Skiena. Matching for run-length encoded strings. In *Proc. 1997 International Conference on Compression and Complexity of Sequences (SEQUENCES'97)*, pages 348–356. IEEE, 1997.
- [8] O. Arbell, G. M. Landau, and J. S. Mitchell. Edit distance of run-length encoded strings. *Information Processing Letters*, 83(6):307–314, 2002.
- [9] P. Austrin, P. Kaski, M. Koivisto, and J. Määttä. Space–time tradeoffs for subset sum: An improved worst case algorithm. In *Proc. 40th International Colloquium on Automata, Languages, and Programming (ICALP'13)*, pages 45–56, 2013.
- [10] A. Backurs and P. Indyk. Edit Distance Cannot Be Computed in Strongly Subquadratic Time (unless SETH is false). In *Proc. 47th Annual ACM Symposium on Theory of Computing (STOC'15)*, pages 51–58, 2015.
- [11] A. Backurs and P. Indyk. Which regular expression patterns are hard to match? In *Proc. 57th IEEE Annual Symposium on Foundations of Computer Science (FOCS'16)*, 2016.

- [12] P. Bille, P. H. Cording, and I. L. Gørtz. Compressed subsequence matching and packed tree coloring. In *Proc. Annual Symposium on Combinatorial Pattern Matching (CPM'14)*, pages 40–49, 2014.
- [13] P. Bille, G. M. Landau, R. Raman, K. Sadakane, S. R. Satti, and O. Weimann. Random access to grammar-compressed strings and trees. *SIAM Journal on Computing*, 44(3):513–539, 2015.
- [14] K. Bringmann. Why walking the dog takes time: Frechet distance has no strongly subquadratic algorithms unless seth fails. In *Proc. of 55th IEEE Annual Symposium on Foundations of Computer Science (FOCS'14)*, pages 661–670, 2014.
- [15] K. Bringmann, F. Grandoni, B. Saha, and V. Vassilevska Williams. Truly sub-cubic algorithms for language edit distance and rna-folding via fast bounded-difference min-plus product. In *Proc. 57th IEEE Annual Symposium on Foundations of Computer Science (FOCS'16)*, pages 375–384. IEEE, 2016.
- [16] K. Bringmann, A. Grønlund, and K. G. Larsen. A dichotomy for regular expression membership testing. In *Proc. 58th IEEE Annual Symposium on Foundations of Computer Science (FOCS'17)*, 2017.
- [17] K. Bringmann and M. Künnemann. Quadratic Conditional Lower Bounds for String Problems and Dynamic Time Warping. In *Proc. 56th IEEE Annual Symposium on Foundations of Computer Science (FOCS'15)*, pages 79–97, 2015.
- [18] H. Bunke and J. Csirik. An improved algorithm for computing the edit distance of run-length coded strings. *Information Processing Letters*, 54(2):93–96, 1995.
- [19] C. Calabro, R. Impagliazzo, and R. Paturi. A duality between clause width and clause density for SAT. In *Proc. 21st IEEE Conference on Computational Complexity (CCC'06)*, pages 252–260, 2006.
- [20] P. Cégielski, I. Guessarian, Y. Lifshits, and Y. Matiyasevich. Window subsequence problems for compressed texts. In *Proc. 1st International Computer Science Symposium in Russia (CSR'06)*, pages 127–136. Springer, 2006.
- [21] T. M. Chan and M. Lewenstein. Clustered Integer 3SUM via Additive Combinatorics. In *Proc. 47th Annual ACM Symposium on Theory of Computing (STOC'15)*, 2015.
- [22] Y. Chang. Conditional lower bound for RNA folding problem. *CoRR*, abs/1511.04731, 2015.
- [23] M. Charikar, E. Lehman, D. Liu, R. Panigrahy, M. Prabhakaran, A. Sahai, and A. Shelat. The smallest grammar problem. *STOC'02 and IEEE Transactions on Information Theory*, 51(7):2554–2576, 2005.
- [24] R. Clifford, A. Fontaine, E. Porat, B. Sach, and T. Starikovskaya. The k-mismatch problem revisited. In *Proc. 27th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'17)*, pages 2039–2052, 2016.
- [25] J. Cocke. Programming languages and their compilers. 1970.

- [26] T. H. Cormen, C. Stein, R. L. Rivest, and C. E. Leiserson. *Introduction to Algorithms*. McGraw-Hill Higher Education, 2nd edition, 2001.
- [27] M. Crochemore, G. M. Landau, and M. Ziv-Ukelson. A subquadratic sequence alignment algorithm for unrestricted scoring matrices. *SIAM Journal on Computing*, 32(6):1654–1673, 2003.
- [28] G. Das, R. Fleischer, L. Gasieniec, D. Gunopulos, and J. Kärkkäinen. Episode matching. In *Proc. Annual Symposium on Combinatorial Pattern Matching (CPM’97)*, pages 12–27. Springer, 1997.
- [29] J. Earley. An efficient context-free parsing algorithm. *Communications of the ACM*, 13(2):94–102, 1970.
- [30] S. R. Eddy. How do rna folding algorithms work? *Nature biotechnology*, 22(11):1457–1458, 2004.
- [31] F. Eisenbrand and F. Grandoni. On the complexity of fixed parameter clique and dominating set. *Theoretical Computer Science*, 326(1-3):57–67, 2004.
- [32] T. Gagie, P. Gawrychowski, and S. J. Puglisi. Faster approximate pattern matching in compressed repetitive texts. In *International Symposium on Algorithms and Computation*, pages 653–662. Springer, 2011.
- [33] A. Gajentaan and M. H. Overmars. On a class of $O(N^2)$ problems in computational geometry. *Comput. Geom. Theory Appl.*, 45(4):140–152, 2012.
- [34] L. Gasieniec, M. Karpinski, W. Plandowski, and W. Rytter. Efficient algorithms for Lempel-Ziv encoding. *Proc. 5th Scandinavian Workshop on Algorithm Theory (SWAT’96)*, pages 392–403, 1996.
- [35] P. Gawrychowski. Faster algorithm for computing the edit distance between slp-compressed strings. In *International Symposium on String Processing and Information Retrieval*, pages 229–236. Springer, 2012.
- [36] R. Giancarlo, D. Scaturro, and F. Utro. Textual data compression in computational biology: a synopsis. *Bioinformatics*, 25(13):1575–1586, 2009.
- [37] A. Grønlund and S. Pettie. Threesomes, degenerates, and love triangles. In *Proc. 55th IEEE Annual Symposium on Foundations of Computer Science (FOCS’14)*, pages 621–630, 2014.
- [38] S. Grumbach and F. Tahi. Compression of DNA sequences. In *Proc. Data Compression Conference (DCC’93)*, pages 340–350, 1993.
- [39] S. Grumbach and F. Tahi. A new challenge for compression algorithms: genetic sequences. *Information Processing & Management*, 30(6):875–886, 1994.
- [40] D. Hermelin, G. M. Landau, S. Landau, and O. Weimann. Unified compression-based acceleration of edit-distance computation. *Algorithmica*, 65(2):339–353, 2013.

- [41] J. E. Hopcroft, R. Motwani, and J. D. Ullman. Automata theory, languages, and computation. *International Edition*, 24, 2006.
- [42] R. Impagliazzo and R. Paturi. On the complexity of k-sat. *Journal of Computer and System Sciences*, 62(2):367–375, 2001.
- [43] R. Impagliazzo, R. Paturi, and F. Zane. Which problems have strongly exponential complexity? *Journal of Computer and System Sciences*, 63:512–530, 2001.
- [44] A. Jez. A really simple approximation of smallest grammar. *Theoretical Computer Science*, 616:141–150, 2016.
- [45] A. Jez. Recompression: a simple and powerful technique for word equations. *Journal of the ACM (JACM)*, 63(1):4, 2016.
- [46] T. Kasami. An efficient recognition and syntax algorithm for context-free algorithms. In *Technical Report AFCRL-65-758 Air Force Cambridge Research Lab Bedford, Mass.* 1965.
- [47] N. J. Larsson. *Structures of string matching and data compression*. Department of Computer Science, Lund University, 1999.
- [48] A. Lempel and J. Ziv. On the complexity of finite sequences. *IEEE Transactions on Information Theory*, 22(1):75–81, 1976.
- [49] Y. Lifshits. Processing compressed texts: A tractability border. In *Proc. Annual Symposium on Combinatorial Pattern Matching (CPM'07)*, pages 228–240. Springer, 2007.
- [50] Q. Liu, Y. Yang, C. Chen, J. Bu, Y. Zhang, and X. Ye. RNACompress: Grammar-based compression and informational complexity measurement of RNA secondary structure. *BMC bioinformatics*, 9(1):176, 2008.
- [51] M. Lohrey. Word problems and membership problems on compressed words. *SIAM Journal on Computing*, 35(5):1210–1240, 2006.
- [52] M. Lohrey. Leaf languages and string compression. *Information and Computation*, 209(6):951–965, 2011.
- [53] M. Lohrey. Algorithmics on SLP-compressed strings: A survey. *Groups Complexity Cryptology*, 4(2):241–299, 2012.
- [54] U. Manber. A text compression scheme that allows fast searching directly in the compressed file. *ACM Transactions on Information Systems (TOIS)*, 15(2):124–136, 1997.
- [55] N. Markey and P. Schnoebelen. A ptime-complete matching problem for slp-compressed words. *Information Processing Letters*, 90(1):3–6, 2004.
- [56] Miscellaneous Authors. Queries and problems. *SIGACT News*, 16(3):38–47, 1984.
- [57] C. G. Nevill-Manning and I. H. Witten. Compression and explanation using hierarchical grammars. *The Computer Journal*, 40(2 and 3):103–116, 1997.

- [58] J. Nešetřil and S. Poljak. On the complexity of the subgraph problem. *Commentationes Math. Universitatis Carolinae*, 026(2):415–419, 1985.
- [59] M. Patrascu. Towards polynomial lower bounds for dynamic problems. In *Proc. 42nd ACM Symposium on Theory of Computing (STOC'10)*, pages 603–610, 2010.
- [60] W. Plandowski and W. Rytter. Application of Lempel-Ziv encodings to the solution of word equations. *Automata, Languages and Programming*, pages 731–742, 1998.
- [61] W. Plandowski and W. Rytter. Complexity of language recognition problems for compressed words. In *Jewels are forever*, pages 262–272. Springer, 1999.
- [62] A. Polak. Why is it hard to beat $O(n^2)$ for longest common weakly increasing subsequence? *Information Processing Letters*, 132:1–5, 2018.
- [63] R. Radicioni and A. Bertoni. Grammatical compression: compressed equivalence and other problems. *Discrete Mathematics and Theoretical Computer Science*, 12(4):109, 2010.
- [64] W. Rytter. Application of Lempel–Ziv factorization to the approximation of grammar-based compression. *Theoretical Computer Science*, 302(1-3):211–222, 2003.
- [65] W. Rytter. Grammar compression, LZ-encodings, and string algorithms with implicit input. In *Proc. 31st International Colloquium on Automata, Languages, and Programming (ICALP'04)*, pages 15–27. Springer, 2004.
- [66] H. Sakamoto. Grammar compression: Grammatical inference by compression and its application to real data. In *ICGI*, pages 3–20, 2014.
- [67] D. Sculley and C. E. Brodley. Compression and machine learning: A new perspective on feature space vectors. In *Proc. Data Compression Conference (DCC'06)*, pages 332–341, 2006.
- [68] Y. Shibata, T. Kida, S. Fukamachi, M. Takeda, A. Shinohara, T. Shinohara, and S. Arikawa. Byte pair encoding: A text compression scheme that accelerates pattern matching. Technical report, Technical Report DOI-TR-161, Department of Informatics, Kyushu University, 1999.
- [69] A. Tiskin. Faster subsequence recognition in compressed strings. *Journal of Mathematical Sciences*, 158(5):759–769, 2009.
- [70] A. Tiskin. Fast distance multiplication of unit-Monge matrices. In *Proc. 21st Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'10)*, pages 1287–1296. SIAM, 2010.
- [71] A. Tiskin. Towards approximate matching in compressed strings: Local subsequence recognition. In *Proc. International Computer Science Symposium in Russia (CSR'11)*, pages 401–414. Springer, 2011.
- [72] L. G. Valiant. General context-free recognition in less than cubic time. *Journal of Computer and System Sciences*, 10(2):308–315, 1975.
- [73] V. Vassilevska. Efficient algorithms for clique problems. *Inf. Process. Lett.*, 109(4):254–257, 2009.

- [74] J. Wang. Space-efficient randomized algorithms for k-sum. In *Proc. 22nd Annual European Symposium on Algorithms (ESA'14)*, pages 810–829, 2014.
- [75] T. A. Welch. A technique for high-performance data compression. *Computer*, 6(17):8–19, 1984.
- [76] R. Williams. A new algorithm for optimal 2-constraint satisfaction and its implications. *Theoretical Computer Science*, 348(2):357–365, 2005.
- [77] I. H. Witten, A. Moffat, and T. C. Bell. *Managing gigabytes: compressing and indexing documents and images*. Morgan Kaufmann, 1999.
- [78] G. J. Woeginger. Space and time complexity of exact algorithms: Some open problems. In *Proc. 1st International Workshop on Parameterized and Exact Computation (IWPEC'04)*, pages 281–290, 2004.
- [79] T. Yamamoto, H. Bannai, S. Inenaga, and M. Takeda. Faster subsequence and don't-care pattern matching on compressed texts. In *Proc. Annual Symposium on Combinatorial Pattern Matching (CPM'11)*, pages 309–322. Springer, 2011.
- [80] D. H. Younger. Recognition and parsing of context-free languages in time n^3 . *Information and Control*, 10(2):189–208, 1967.
- [81] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3):337–343, 1977.