

Consequences of Faster Alignment of Sequences

Amir Abboud¹, Virginia Vassilevska Williams¹, and Oren Weimann²

¹ Stanford University, USA. {abboud,virgi}@cs.stanford.edu

² University of Haifa, Israel. oren@cs.haifa.ac.il

Abstract. The Local Alignment problem is a classical problem with applications in biology. Given two input strings and a scoring function on pairs of letters, one is asked to find the substrings of the two input strings that are most similar under the scoring function. The best algorithms for Local Alignment run in time that is roughly quadratic in the string length. It is a big open problem whether substantially subquadratic algorithms exist. In this paper we show that for all $\varepsilon > 0$, an $O(n^{2-\varepsilon})$ time algorithm for Local Alignment on strings of length n would imply breakthroughs on three longstanding open problems: it would imply that for some $\delta > 0$, 3SUM on n numbers is in $O(n^{2-\delta})$ time, CNF-SAT on n variables is in $O((2-\delta)^n)$ time, and Max Weight 4-Clique is in $O(n^{4-\delta})$ time. Our result for CNF-SAT also applies to the easier problem of finding the longest common substring of binary strings with don't cares. We also give strong conditional lower bounds for the more general Multiple Local Alignment problem on k strings, under both k -wise and SP scoring, and for other string similarity problems such as Global Alignment with gap penalties and normalized Longest Common Subsequence.

1 Introduction

Many basic string and pattern matching problems have overwhelming importance in current bioinformatics research. A well known such problem is the Local Alignment problem which asks to find the two substrings of two given strings that are most similar, under some given similarity measure. The fastest theoretical algorithm for this problem runs in time $O(n^2/\log n)$ [15,38,8] on n -length strings, and is not much faster than the classical dynamic programming algorithm of Smith and Waterman [47] which runs in $O(n^2)$ time. A faster algorithm for this problem, one that runs in, say $O(n^{1.5})$ time, would have tremendous impact, as witnessed by the 49,000 citations to the paper introducing the practical BLAST (Basic Local Alignment Search Tool) algorithm [3]. However, there seems to be very little optimism in the computational biology community that Local Alignment and other important string problems admit such “truly subquadratic” algorithms, i.e. running in time $O(n^{2-\varepsilon})$ for $\varepsilon > 0$. Yet, the theoretical computer science community has not provided any evidence for this impossibility, and in particular, it is yet to give an answer to this pressing question: can Local Alignment be solved in truly subquadratic time?

Perhaps the main reason for this lack of an answer, is that we do not have a clear technique for providing negative answers to such questions. The state of the art on unconditional lower bounds seems far from proving any significant superlinear lower bounds in the near future. The theory of NP-completeness

cannot distinguish between quadratic upper bounds and $n^{1.5}$ ones. The $W[t]$ -hardness approach of parameterized complexity requires parameterization, and like NP-hardness, does not distinguish between differing polynomial runtimes. Another approach is to prove lower bounds for a restricted family of algorithms. However, it is unclear what an appropriate candidate family would be, and either way, a restricted model lower bound only gives a partial answer to the question.

Our approach. We follow an approach that can be viewed as a refinement of NP-hardness. The importance of showing NP-hardness for a certain problem lies in the consequence that a polynomial time algorithm for this problem would also imply a polynomial time algorithm for many other problems that are widely believed to require superpolynomial solutions. The goal of this work and previous works that are mentioned below is to develop such theory that is able to prove that improving the exact running times of certain problems would also imply surprising algorithms for many other problems and is therefore unlikely.

Using this approach we are able to provide the following answer to our pressing question, which is stated more formally in our theorems: A truly subquadratic algorithm for Local Alignment is unlikely because it would also give truly faster algorithms for other famous problems like CNF-SAT, 3-SUM and Max-Weight-4-Clique, implying breakthroughs in three different areas of computer science: the satisfiability algorithms and circuit lower bounds, the computational geometry and the graph algorithms communities!

To provide such answers, to this and other important questions about the optimality of current upper bounds for string problems, we devise careful reductions to the string problems from famous problems that are widely believed to require certain running times, not necessarily quadratic.

1.1 3-SUM Hardness

The most prominent example of this approach is the theory of 3-SUM hardness which was introduced by Gajentaan and Overmars [25] and has been used to show that subquadratic upper bounds for many problems in Computational Geometry are unlikely.

In the 3-SUM problem we are given three lists of n numbers and are asked whether we can pick a number from each list so that the sum is 0. A simple algorithm solves the problem in $\tilde{O}(n^2)$ time, and Baran, Demaine and Pătraşcu [6] were able to get a $O(n^2/\log^2 n)$ solution, yet any improvement beyond this seems unlikely and it is a widely believed conjecture that a polynomial improvement on the upper bound is impossible. Support for this belief comes from the $\Omega(n^2)$ lower bound for the depth of an algebraic decision tree for the problem [23].

Conjecture 1 (3-SUM Conjecture) *In the Word RAM model with words of $O(\log n)$ bits, any algorithm requires $n^{2-o(1)}$ time in expectation to determine whether three sets $A, B, C \subset \{-n^3, \dots, n^3\}$ with $|A| = |B| = |C| = n$ integers contain three elements $a \in A, b \in B, c \in C$ with $a + b + c = 0$.*

Since [25], there have been many papers proving the hardness of computational geometry problems, based on 3SUM, e.g. [19,37,22,13,7]. More recently, the 3-SUM Conjecture has been used in surprising ways to show polynomial

lower bounds for purely combinatorial problems in dynamic algorithms [40,2] and Graph algorithms [40,32,48]. The only previous work relating 3-SUM to a Stringology problem, to our knowledge, is the result of Chen et al. [12] showing that under the 3-SUM Conjecture, when the input strings are encodings of much longer strings, using Run-Length-Encoding, then the string matching with don't cares problem requires time that is quadratic in the lengths of the *compressions*. This string problem, however, is strongly related to geometric problems and is less “combinatorial” than the problems we consider here (e.g. it is solvable by a sweep-line algorithm). Thus our reductions require different techniques.

We expand the list of 3-SUM hard problems, showing a reduction from 3-SUM to the Local Alignment problem, proving that a truly subquadratic algorithm for Local Alignment is impossible under the 3-SUM Conjecture, provided the alphabet is large enough.

Theorem 1. *If for some $\varepsilon > 0$, $\delta \in (0,1)$, one can solve the Local Alignment problem on two strings of length n over an alphabet of size $n^{1-\delta}$ in time $O(n^{2-\delta-\varepsilon})$, then the 3-SUM Conjecture is false.*

To prove Theorem 1 we combine a recent reduction from k -SUM to k -Vector-SUM [1] with an efficient self reduction for 3-SUM using hashing [6,40], then we carefully construct a scoring scheme. The details are given in Section 3.

We note that it is not hard to argue that there is an unconditional lower bound of $\min\{|\Sigma|^2, n^2\}$ for Local Alignment where Σ is the alphabet. When $|\Sigma| = n^{1-\delta}$, this lower bound is $\Omega(n^{2-2\delta})$. Our lower bound for such alphabets is essentially $n^{2-\delta}$ which is a *polynomial* improvement over $n^{2-2\delta}$. Nevertheless, our reduction requires alphabet size at least n^ε for some arbitrarily small but constant $\varepsilon > 0$, and the really interesting case of Local Alignment is when the alphabet size is constant, e.g. 4 for DNA and RNA sequences and 20 for protein sequences. We are not able to use the 3-SUM Conjecture to conclude a lower bound for this case. To handle the constant alphabet case, we turn to the presumed hardness of CNF-SAT to prove a lower bound even for *binary* strings.

1.2 Strong ETH Hardness

Despite hundreds of papers on faster exponential algorithms for NP-Hard problems in recent years (see the surveys by Woeginger for an exposition [52]), and despite the remarkable effort put into obtaining faster satisfiability algorithms, the best upper bounds for CNF-SAT on n variables and m clauses remain of the form $2^{n-o(n)} \text{poly}(m)$ (e.g. [28,41,46]). The Strong Exponential Time Hypothesis (Strong ETH) of Impagliazzo, Paturi and Zane, which has received a lot of attention recently, states that better algorithms do not exist.

Conjecture 2 (Strong ETH) *For every $\varepsilon > 0$, there exists a k , such that SAT on k -CNF formulas on n variables cannot be solved in $O^*(2^{(1-\varepsilon)n})$ time.*

Strong ETH is an extremely popular conjecture in the exact exponential time algorithms community [11,18,35,17], and Cygan et al. [16] even showed it to be equivalent to assuming that several other NP-hard problems essentially require exhaustive search. Recently, many surprising lower bounds in several different areas were shown to hold under the SETH, including lower bounds for

approximating the diameter of a sparse graph [45], for maintaining the number of strongly connected components in a dynamic graph [2], and for the 3-party communication complexity of Set-Disjointness [43].

We show a reduction from CNF-SAT to the *longest common substring with don't cares* problem, which is one of the simplest string problems for which truly subquadratic algorithms are not known and is a very restricted version of the Local Alignment problem.

Definition 1 (The longest common substring with don't cares problem). *Given a string S over alphabet $\Sigma = \{0, 1\}$ and a string T over $\Sigma \cup \{\star\}$, find the length of the longest string that is a substring of both S and T , where a \star in T can be treated as either 0 or 1.*

If don't care letters are not allowed, the problem can be solved using a generalized suffix tree in $O(n)$ time [27]. If instead of looking for the longest common substring, one wants to check whether a binary string with don't cares appears as a substring in a length n binary string, then there are $O(n \log n)$ time algorithms based on the Fast Fourier Transform [24,31,33]. Thus only slight variations of the longest common substring with don't cares problem admit almost *linear* time solutions, yet our reduction implies that a truly *subquadratic* algorithm for it refutes Strong ETH!

Theorem 2. *If for some $\varepsilon > 0$ one can solve either the Local Alignment problem on two binary strings of length n , or the longest common substring with don't cares problem in time $O(n^{2-\varepsilon})$, then Strong ETH is false.*

To prove Theorem 2 we represent a partial assignment to the variables of our formula with a substring, and we make sure that two substrings will match if and only if the two partial assignments satisfy all the clauses of our formula. The details are given in Section 2.

1.3 Multiple Sequence Alignment

In Gusfield's book [27], algorithms for comparing multiple strings are called "the holy grail" of current research in computational biology. One of the important tasks is the Multiple Local Alignment (MLA) problem [34] defined as follows. Given k strings T_1, \dots, T_k over some alphabet, find substrings S_1, \dots, S_k (where S_i is a substring of T_i) whose alignment has maximum score. There are multiple ways to define the alignment score between k substrings but we focus on two options, the most general k -wise scoring scheme and the most popular Sum of Pairs (SP) scoring scheme [5].

The k -wise scoring function $s(\cdot, \dots, \cdot)$ which is given as a k dimensional matrix with $(|\Sigma|)^k$ entries. This case is called k -wise scoring, and the score of an alignment of k strings s_1, \dots, s_k is $\sum_i s(s_1[i], \dots, s_k[i])$ ³. The second option, which is called sum of pairs (SP) scoring, is to use a pairwise scoring function $s(\cdot, \cdot)$ and to define the score of an alignment to be $\sum_i \sum_{k < \ell} s(s_k[i], s_\ell[i])$.

³ In a more general alignment, one can align alphabet symbols with spaces, and the scoring function can take that into account. In this paper, we prove hardness even for the easier alignment problem where no spaces are allowed.

The MLA problem on $k \geq 3$ strings can be defined using either k -wise scoring or SP scoring. Local Alignment is the $k = 2$ case and both scoring rules coincide.

For both scoring schemes, unsurprisingly, the best upper bound for the problem is $O(n^k)$ using dynamic programming. The reduction of Wang and Jiang [49] from the shortest common superstring problem on k strings to a polynomial number of instances of the (global or local) alignment problem on $k + 2$ strings with SP scoring implies that our problem is NP-hard when the number of strings is unbounded. Moreover, the $W[1]$ -hardness results of Bodlaender et al. [10] for unbounded alphabets and of Pietrzak [42] for constant size alphabets also carry on to our problems, implying that upper bounds of the form $f(k) \cdot n^c$ for a constant c independent of k and n are unlikely. Huang [29] strengthens Pietrzak’s reduction from k -clique to the shortest common superstring problem and shows that $n^{o(k)}$ algorithms for our problems are not possible under the plausible Exponential Time Hypothesis [30].

These negative results deliver an important message to biologists, showing that an efficient algorithm for optimally aligning a hundred strings is not likely to exist, yet another pressing question remains widely open: is there an algorithm running in time $O(n^{k-1})$ or even $O(n^{k/5})$ for MLA?⁴ Such algorithms would imply a major advance in our ability to analyze biological data.

We extend our reduction from CNF-SAT to show that Strong ETH implies a negative answer to our question, when we are interested in k -wise scoring, even when the strings are binary.

Theorem 3. *If for some $\varepsilon > 0$ the MLA on k binary strings of length n with k -wise scoring can be solved in time $O(n^{k-\varepsilon})$, then Strong ETH is false.*

As is stressed in Gusfield’s book [27], the less general case of SP scoring has more applications in Bioinformatics. Our reduction from CNF-SAT requires the computation of an OR function, which is easy with k -wise scoring yet does not seem possible with SP scoring. We show, however, that the Weighted- k -Clique problem can explain the hardness of getting faster algorithms even for the SP case of MLA.

Weighted k -Clique. The Max-Weight k -Clique problem is as follows. Given a graph $G = (V, E)$ with integer edge weights, find a k -clique of maximum total weight, or determine that no k -clique exists. Since k -Clique is a special case of the problem, the Max-Weight k -Clique problem is NP-complete and $W[1]$ -hard. The unweighted k -Clique admits an $O(n^{0.792k})$ time algorithm using matrix multiplication [39]. When the edge weights are small, one can obtain $O(n^{k-\varepsilon})$ time algorithms for Max-Weight k -Clique as well, by reducing to a small number of instances of k -Clique. However, when the weights are larger than n^k , the trivial $O(n^k)$ algorithm is essentially the best known (ignoring $n^{o(1)}$ improvements).

We show a tight reduction from Max-Weight $2k$ -Clique to MLA on k strings with SP scoring, showing that improving on n^k time for MLA would also imply that Max-Weight $2k$ -Clique has faster than n^{2k} algorithms.

Theorem 4. *If for some $\varepsilon > 0$ the MLA on k strings of length n over an alphabet of size \sqrt{n} with SP scoring can be solved in time $O(n^{k-\varepsilon})$, then Max-Weight $2k$ -Clique on n nodes graphs can be solved in time $O(n^{2k-\varepsilon})$.*

⁴ In the reduction of Bodlaender et al. [10] k increases to k^2 and in Pietrzak’s [42] reduction n increases to n^7 .

An immediate corollary of the above theorem is a conditional lower bound for the Local Alignment problem for two strings, based on the assumption that Max-Weight 4-Clique does not have improved algorithms.

Corollary 1. *If for some $\varepsilon > 0$ the Local Alignment on two strings of length n over an alphabet of size \sqrt{n} can be solved in time $O(n^{2-\varepsilon})$, then Max-Weight 4-Clique on n nodes graphs can be solved in time $O(n^{4-\varepsilon})$.*

We note that the special case of Max-Weight k -clique for $k = 3$ is especially interesting. The Max-Weight 3-Clique problem on n -node graphs is known to be essentially equivalent to the All Pairs Shortest Paths (APSP) problem, in the sense that if Max-Weight 3-Clique has a truly subcubic algorithm, so does APSP, and vice versa. It is a longstanding open problem whether APSP on n -node graphs can be solved in $O(n^{3-\varepsilon})$ time [51]. It is thus a major open problem whether Max-Weight 3-Clique is in $O(n^{3-\varepsilon})$ time for some $\varepsilon > 0$. Our current reductions show that a $O(n^{1.5-\varepsilon})$ time algorithm for Local Alignment on two strings would give $O(n^{3-\varepsilon})$ time for APSP, but we suspect that they can be strengthened to show a tighter relationship.

Extensions. In the full version of the paper we also show quadratic lower bounds for well-known generalizations of the *Global Alignment* problem like Alignment with Gap penalties [26] and Alignment with moves [36], and for other string problems like Normalized LCS [4,21] and Partial Match [14,44,9].

2 From CNF-SAT to Alignment

In this section we give a reduction from CNF-SAT on n variables and m clauses to the longest common substring with don't cares problem on binary strings of length $N = O(2^{n/2} \cdot m)$ in $O^*(2^{n/2} \cdot m)$ time. Thus, given an algorithm for this problem that runs in time $O(N^{2-\varepsilon})$ for some $\varepsilon > 0$, we can solve CNF-SAT in time $O^*((2^{n/2} \cdot m)^{2-\varepsilon}) = O^*(2^{(1-\varepsilon/2)n} \cdot \text{poly}(m))$, refuting Strong ETH. In the full version of this paper we explain how to get a reduction to Local Alignment on binary strings, proving Theorem 2, and give the extension to MLA on k strings to prove Theorem 3.

Our reduction follows the split and list technique introduced by Williams [50]. In particular, our reduction from CNF-SAT to the longest common substring problem can be obtained by combining his reduction from CNF-SAT to the *orthogonal vectors* problem on binary vectors and a simple reduction from the latter problem to the longest common substring problem. Below we present a direct reduction from CNF-SAT that makes our extension to MLA on k strings follow more clearly.

Lemma 1. *CNF-SAT over formulas on n variables and m clauses can be reduced to the longest common substring with don't cares problem over strings of length $O(2^{n/2} \cdot m)$ and constant-size alphabet in $O^*(2^{n/2} \cdot m)$ time.*

Proof. Let $V = \{x_1, \dots, x_n\}$ be the n variables of the input CNF formula φ . We split V into two sets of $n/2$ variables, $U = \{x_1, \dots, x_{n/2}\}$ and $V \setminus U$. Let $A = \{\alpha_1, \dots, \alpha_N\}$ and $B = \{\beta_1, \dots, \beta_N\}$ be the sets of all $N = 2^{n/2}$ partial

assignments of boolean values to the variables in U and $V \setminus U$, respectively. Combining two partial assignments $\alpha \in A$ and $\beta \in B$ gives an assignment $(\alpha \cdot \beta)$ to all the variables in the formula. We say that a partial assignment α satisfies a clause C if α either assigns TRUE to a variable that appears in C as a positive literal or it assigns FALSE to a variable that appears in C as a negative literal. The key idea of the reduction is the following simple observation, which gives a way of checking the satisfiability of ϕ : *The formula ϕ is satisfiable if and only if there are two partial assignments $\alpha \in A, \beta \in B$ such that for every clause C in the formula at least one of α and β satisfies the clause C .*

The reduction will generate two strings S and T . The string S will have N segments of length $5m$ corresponding to the N partial assignments in A and each segment will encode which clauses are satisfied by the partial assignment. Similarly, T will encode the partial assignments in B . A common substring of S and T will have to be entirely contained in these segments, and it will be able to be the whole segment only if the two corresponding assignments satisfy all the clauses of the formula. Therefore, the longest common substring will be of length $5m$ if and only if ϕ is satisfiable.

Let C_1, \dots, C_m be the clauses of our CNF formula ϕ . For $\alpha \in A$ we define the segment string S_α to contain a different symbol in the $(5j - 2)^{th}$ position $S_\alpha[j]$ according to whether α satisfies C_j . Then, $\forall j \in [m]$:

$$S_\alpha[(5j - 4) \dots (5j)] = [01x_j01], \text{ where } x_j = 1 \text{ if } \alpha \text{ satisfies } C_j, \text{ and } 0 \text{ otherwise.}$$

Similarly, for $\beta \in B$ we define the segment T_β as follows. $\forall j \in [m]$:

$$T_\beta[(5j - 4) \dots (5j)] = [\star 1y_j \star 1], \text{ where } y_j = \star \text{ if } \beta \text{ satisfies } C_j, \text{ and } 1 \text{ otherwise.}$$

Note that we can construct these strings for every $\alpha \in A$ and $\beta \in B$ given ϕ in $O^*(2^{n/2}m)$ time. Finally, we create S by concatenating all the segment strings S_{α_i} for all $\alpha_i \in A$ and placing “unmatchable” [000] segments between them, and we create T similarly by concatenating the T_{β_i} segment strings and placing [111] segments between them.

$$S = S_{\alpha_1} \circ 0^3 \circ S_{\alpha_2} \circ \dots \circ 0^3 \circ S_{\alpha_N}, \quad T = T_{\beta_1} \circ 1^3 \circ T_{\beta_2} \circ \dots \circ 1^3 \circ T_{\beta_N}$$

Claim 1 *The longest common substring of S and T is of length $5m$ iff there are $\alpha \in A, \beta \in B$ such that every clause C_j in ϕ is satisfied by at least one of α, β .*

To prove Claim 1 we first observe that two segments S_α and T_β match, giving a common substring of length $5m$, if and only if every clause in ϕ is satisfied by at least one of α and β , which implies that the formula is satisfiable. Then, we need to argue that any common substring of length $5m$ cannot contain any of our “unmatchable” parts and must therefore correspond to two segments S_α and T_β that match. The details are given in the full version of this paper.

3 From 3-SUM to Alignment

In this section we prove Theorem 1 by proving the following lemma.

Lemma 2. *For any $0 < \delta < 1$, the 3-SUM problem on n numbers can be reduced in $n^{2-\delta+o(1)}$ time to $n^{\delta+o(1)}$ instances of the local alignment problem on two strings of length $\tilde{O}(n)$ over an alphabet Σ of size $\tilde{O}(n^{1-\delta})$.*

Proof. Given three lists $A, B, C \subseteq \{-n^3, \dots, n^3\}$ of n numbers each, we want to find a triple of numbers $a \in A, b \in B, c \in C$ that sum to 0. We start by applying a hashing scheme that is due to Dietzfelbinger [20] and that has been used in recent works on 3-SUM [6,40,32]. Let $R = n^{1-\delta}$. There is a simple family \mathcal{H} of hash functions $h : \{-n^3, \dots, n^3\} \rightarrow [R]$ such that if we pick a function $h \in \mathcal{H}$ from the family at random, and hash each element x in our input sets $A \cup B \cup C$ to the bucket $B(h(x))$, we get the following properties:⁵

- *Almost linearity.* For any three numbers a, b, c , if $a + b + c = 0$ then $(h(a) + h(b) + h(c))$ modulo R is either 0 or 1.
- *Good load balancing.* The average number of elements x hashed into a bucket $B(x)$ is $3n/R$ and, in expectation, at most $O(R)$ elements are hashed to buckets with load exceeding $9n/R$.

The reduction picks a random hash function $h \in \mathcal{H}$ and hashes each element x in our lists to a bucket $B(h(x))$. For the $O(R)$ elements that fall in overloaded buckets we can run a brute force check to see if they participate in a 3-sum in $O(nR)$ time. Therefore, we can assume that we have at most R buckets $B(1), \dots, B(R)$, each containing at most $t = 9n/R = O(n^\delta)$ elements. We order the elements in these buckets $B(i) = \{x_1, \dots, x_t\}$, and for each index j from 1 to t , we will have a separate stage. In stage j we check whether there is any element c in C such that c is the j^{th} element of its bucket $B(h(c))$ and $a + b + c = 0$ for some $a \in A, b \in B$. By the “almost linearity” property, it is enough to search for the pair $a \in A, b \in B$ among the elements a, b for which either $h(a) + h(b) = -h(c)$ or $h(a) + h(b) = -h(c) + 1$ (modulo R). To do these checks, in every stage j we create $n^{o(1)}$ instances of the local alignment problem, as described below. The total number of instances is therefore $t \cdot n^{o(1)} = n^{\delta+o(1)}$.

In a recent result by Abboud, Lewi and Williams (Lemma 3.1 in [1]), the authors show that we can construct a set of simple $N = n^{O(1/\log \log n)} = n^{o(1)}$ mappings f_1, \dots, f_N from numbers in $\{-n^3, \dots, n^3\}$ to vectors in $\{-p, \dots, p\}^d$ where $p = \log n$ and $d = O(\log n / \log \log n)$ with the following useful property⁶. If three elements $a, b, c \in \{-n^3, \dots, n^3\}$ sum to 0, then for some $i \in [N]$, the vectors $f_i(a), f_i(b), f_i(c)$ will sum to the all-zero vector $\mathbf{0}$, while if three numbers a, b, c do not sum to 0, then for every $i \in [N]$, the vectors $f_i(a), f_i(b), f_i(c)$ will not sum to the all-zero vector. Note that the entries in each coordinate of our vectors are very small since $p = \log n$.

For every triple of numbers (i, j, z) where $i \in [N], j \in [t], z \in \{0, 1\}$ we create an instance of the Local Alignment problem, i.e. two strings S, T over alphabet Σ and a scoring function $s(\cdot, \cdot)$. The optimal solution to the (i, j, z) local alignment instance will determine whether there are three numbers $a \in A, b \in B, c \in C$ such that

1. c is the j^{th} element in its bucket $B(h(c))$,
2. $h(a) + h(b) + h(c) = z \pmod{R}$, and
3. $f_i(a) + f_i(b) + f_i(c) = \mathbf{0}$.

⁵ The value $h(x)$ is computed by taking a random odd integer a and returning the $\log R$ most significant bits from the number $a \cdot x$.

⁶ The idea is simple: each number is mapped to a vector containing the base- p representation of the number, then enumerate over all guesses for the carries when summing k numbers in their base- p representation.

If we find a triple satisfying these conditions then we have found a 3-sum, since by the property of the mappings from numbers to vectors described above, condition 3 can only happen if $a+b+c=0$. On the other hand, if there is a 3-sum $a \in A, b \in B, c \in C, a+b+c=0$ in our input set, then for some $(i, j, z) \in [N] \times [t] \times \{0, 1\}$, the above three conditions will hold and we will find this 3-sum. To see this, note that by the property of the mappings there must exist an $i \in [N]$ for which condition 3 holds, and by choosing $j \in [t]$ to be such that c is the j^{th} element in its bucket $B(h(c))$ we satisfy condition 2, and finally, by the ‘‘almost linearity’’ property of our hash function, $z = h(a) + h(b) + h(c) \pmod R$ is in the set $\{0, 1\}$ and condition 2 holds as well.

We now describe the Local Alignment instances that we generate for each triple $(i, j, z) \in [N] \times [t] \times \{0, 1\}$. Our alphabet Σ will contain a letter (h, y) for every pair of integers $h \in [R]$ (which will be used to indicate the value of a hash of a number) and $y \in \{-p, \dots, p\}$ (which will represent the value in a coordinate of our vectors). We also add two symbols $\$1$ and $\$2$ to the alphabet. Note that by our choices of $p = \log n$ and $R = n^{1-\delta}$, we get that $|\Sigma| = \tilde{O}(n^{1-\delta})$. As in the reduction from CNF-SAT, the strings will be composed of segments. For every number a in A we create the segment S_a which will have length d and in its ℓ^{th} coordinate it will contain the letter $(h(a), f_i(a)[\ell])$. This letter encodes both the hash of a and the value in the ℓ^{th} entry of the vector $f_i(a)$ (corresponding to a in our current mapping i). Similarly, for every number b in B we define the segment T_b so that $T_b[\ell] = (h(b), f_i(b)[\ell])$ for every $\ell \in [d]$. The strings S, T of our instance (i, j, z) are constructed by concatenating the segments with $\$$ signs between them. Let $A = \{a_1, \dots, a_n\}$ and $B = \{b_1, \dots, b_n\}$, then $S = S_{a_1} \circ \$1 \circ S_{a_2} \circ \$1 \circ \dots \circ \$1 \circ S_{a_n}$ and $T = T_{b_1} \circ \$2 \circ T_{b_2} \circ \$2 \circ \dots \circ \$2 \circ T_{b_n}$.

The scoring function $s(\cdot, \cdot)$ is defined as follows. Given two letters (h_1, y_1) and (h_2, y_2) , the scoring function will lookup $c \in C$, where c is the j^{th} element in bucket number $-(h_1 + h_2) + z \pmod R$, and return a score of 1 if $f_i(c) = -(y_1 + y_2)$. In any other case, the returned score is $-\infty$. Formally, for any pair of letters $(h_1, y_1), (h_2, y_2) \in [R] \times \{-p, \dots, p\}$,

$$s((h_1, y_1), (h_2, y_2)) = \begin{cases} 1 & \text{if } c \in C \text{ is the } j^{\text{th}} \text{ element in } B((- (h_1 + h_2) + z) \pmod R) \\ & \text{and } f_i(c) = -(y_1 + y_2). \\ -\infty & \text{otherwise} \end{cases}$$

We also disallow $\$$ symbols and gaps in the optimal alignment by giving a score of $-\infty$ to any pair containing them.

The proof is completed with the following claim that shows that our construction will find a 3-sum if it exists. The details are given in the full version of this paper.

Claim 2 *There are two substrings of S, T in the (i, j, z) instance achieving a score of d if and only if there is a triple $a \in A, b \in B, c \in C$ satisfying conditions 1 to 3 above.*

4 From Weighted Clique to Alignment

In this section we obtain efficient reductions for all $k \geq 2$ from the Max-Weight $2k$ -Clique problem to the MLA on k strings *with SP scoring*. We can assume that the input graph is complete, by making each nonedge have weight $-\infty$.

Lemma 3. *Max-Weight $2k$ -Clique on a weighted graph with n nodes and m edges with weights in $\{-M, \dots, M\}$ can be reduced in $O(mk)$ time to MLA on k strings of length $O(m(k + \log M))$ and alphabet of size $O(n + \log M)$.*

An interesting observation about our reduction is that the length of the substrings in the optimal alignment is only $O(k + \log M)$, which is quite short, while one could have hoped for faster algorithms for the restricted problem in which we are only looking for short substrings. With more work, one can strengthen the reduction to make the length of the optimal substrings only $O(\log k + \log M)$. Theorem 4 in the introduction is an immediate consequence of Lemma 3.

We explain the idea of the proof next and give the formal details in the full version of this paper. Each of the k input strings will contain m segments, one for each edge in the graph. These segments will be separated by “unmatchable” $\$$ symbols so that the scoring function $s(\cdot, \cdot)$ will have $s(\$, a) = -\infty$ for any symbol a in the alphabet. This will enforce that any solution of positive value must pick a substring from a single segment from each input string. If a solution picks an entire segment corresponding to an edge (u_i, v_i) from each input string i , then our construction will guarantee that the score of the segment alignment is exactly $\sum_{i,j} w(u_i, v_j)$. Notice that if the vertices in $\{u_i, v_i\}_i$ are distinct, i.e. the edges (u_i, v_i) form a matching, then $\sum_{i,j} w(u_i, v_j)$ is exactly the weight of the $2k$ -clique formed by these vertices. We make sure that when we align two segments (u_i, v_i) and (u_j, v_j) of two different strings i, j , for each of the weights $w \in \{w(u_i, v_j), w(u_j, v_i), w(u_i, u_j), w(v_i, v_j), w(u_i, v_i) + w(u_j, v_j)\}$ there is some coordinate of the segment alignment in which the pairwise score contributes to w . We carefully devise these contributions to obtain that the sum of the scores is exactly the weight of the clique.

Finally, our construction will guarantee that the best solution always picks an entire segment for each input string. We ensure this by beginning and ending each segment with a symbol Λ , where $s(\Lambda, \Lambda) = k^2 M$ and $s(\Lambda, a) = -\infty$ if $a \neq \Lambda$. The choice of $s(\Lambda, \Lambda)$ guarantees that the optimum solution would always align full segments on top of each other in order to include the Λ, Λ alignments, and hence the optimum solution corresponds to the maximum weight of a $2k$ -clique.

Acknowledgements. We would like to thank Kevin Lewi, Ryan Williams and the anonymous reviewers for helpful discussions and comments. The first and second authors were supported by a Stanford School of Engineering Hoover Fellowship, NSF Grant CCF-1417238 and BSF Grant BSF:2012338. The third author was supported by Israel Science Foundation grant 794/13.

References

1. A. Abboud, K. Lewi, and R. Williams. On the parameterized complexity of k -sum. *CoRR*, abs/1311.3054, 2013.
2. A. Abboud and V. Vassilevska Williams. Popular conjectures imply strong lower bounds for dynamic problems. *arXiv*, arXiv:1402.0054, 2014.
3. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, 1990.
4. A.N. Arslan, Ö. Egecioglu, and P.A. Pevzner. A new approach to sequence comparison: Normalized sequence alignment. *Bioinformatics*, 17(4):327–337, 2001.

5. David J. Bacon and Wayne F. Anderson. Multiple sequence alignment. *Journal of Molecular Biology*, 191(2):153 – 161, 1986.
6. Ilya Baran, Erik D. Demaine, and Mihai Pătraşcu. Subquadratic algorithms for 3SUM. *Algorithmica*, 50(4):584–596, 2008. See also WADS’05.
7. Gill Barequet and Sarel Har-Peled. Some variants of polygonal containment and minimum hausdorff distance undertranslation are 3SUM-hard. In *SODA*, pages 862–863, 1999.
8. P. Bille and M. Farach-Colton. Fast and compact regular expression matching. *Theoretical Computer Science*, 409(3):486–496, 2008.
9. Philip Bille, Inge Li Gørtz, Hjalte Wedel Vildhøj, and Søren Vind. Less space: Indexing for queries with wildcards. In *SWAT*, pages 283–294, 2012.
10. Hans L Bodlaender, Rodney G Downey, Michael R Fellows, and Harold T Wareham. The parameterized complexity of sequence alignment and consensus. *Theoretical Computer Science*, 147(1):31–54, 1995.
11. C. Calabro, R. Impagliazzo, and R. Paturi. The complexity of satisfiability of small depth circuits. In *Parameterized and Exact Computation*, pages 75–85. Springer, 2009.
12. K. Chen, P. Hsu, and K. Chao. Approximate matching for run-length encoded strings is 3sum-hard. In *CPM*, pages 168–179. Springer, 2009.
13. Otfried Cheong, Alon Efrat, and Sarel Har-Peled. On finding a guard that sees most and a shop that sells most. In *SODA*, pages 1098–1107, 2004.
14. Richard Cole, Lee-Ad Gottlieb, and Moshe Lewenstein. Dictionary matching and indexing with errors and don’t cares. In *STOC*, pages 91–100, 2004.
15. M. Crochemore, G.M. Landau, and M. Ziv-Ukelson. A subquadratic sequence alignment algorithm for unrestricted scoring matrices. *SIAM Journal on Computing*, 32:1654–1673, 2003.
16. M. Cygan, H. Dell, D. Lokshtanov, D. Marx, J. Nederlof, Y. Okamoto, R. Paturi, S. Saurabh, and M. Wahlstrom. On problems as hard as CNFSAT. In *CCC*, pages 74–84, 2012.
17. Marek Cygan, Stefan Kratsch, and Jesper Nederlof. Fast Hamiltonicity checking via bases of perfect matchings. In *STOC*, pages 301–310, 2013.
18. E. Dantsin and A. Wolpert. On moderately exponential time for SAT. In *13th International Conference on Theory and Applications of Satisfiability Testing*, pages 313–325, 2010.
19. Mark de Berg, Marko de Groot, and Mark H. Overmars. Perfect binary space partitions. *Computational Geometry: Theory and Applications*, 7(81):81–91, 1997.
20. Martin Dietzfelbinger. Universal hashing and k-wise independent random variables via integer arithmetic without primes. In *STACS*, pages 569–580, 1996.
21. N. Efraty and G.M. Landau. Sparse normalized local alignment. In *CPM*, pages 333–346, 2004.
22. J. Erickson. New lower bounds for convex hull problems in odd dimensions. *SIAM Journal on Computing*, 28(4):1198–1214, 1999.
23. Jeff Erickson. Bounds for linear satisfiability problems. *Chicago J. Theor. Comput. Sci.*, 1999, 1999.
24. M.J. Fischer and M.S. Paterson. String matching and other products. *SIAM-AMS Proc.*, 7:113–125, 1973.
25. A. Gajentaan and M. Overmars. On a class of $o(n^2)$ problems in computational geometry. *Computational Geometry*, 5(3):165–185, 1995.
26. O. Gotoh. An improved algorithm for matching biological sequences. *J. Mol. Biol.*, 162:705–708, 1982.
27. Dan Gusfield. *Algorithms on strings, trees and sequences: computer science and computational biology*. Cambridge University Press, 1997.
28. E. A. Hirsch. Two new upper bounds for SAT. In *SODA*, pages 521–530, 1998.
29. Xiuzhen Huang. *Parameterized complexity and polynomial-time approximation schemes*. PhD thesis, Citeseer, 2004.

30. R. Impagliazzo and R. Paturi. On the complexity of k -sat. *J. Comput. Syst. Sci.*, 62(2):367–375, 2001.
31. P. Indyk. Faster algorithms for string matching problems: Matching the convolution bound. In *FOCS*, pages 166–, 1998.
32. Z. Jafarholi and E. Viola. 3sum, 3xor, triangles. *Electronic Colloquium on Computational Complexity (ECCC)*, 20:9, 2013.
33. A. Kalai. Efficient pattern-matching with don't cares. In *SODA*, pages 655–656, 2002.
34. C. E. Lawrence, S. F. Altschul, M. S. Boguski, J. S. Liu, A. F. Neuwald, and J. C. Wootton. Detecting subtle sequence signals: a gibbs sampling strategy for multiple alignment. *science*, 262(5131):208–214, 1993.
35. Daniel Lokshtanov, Dániel Marx, and Saket Saurabh. Known algorithms on graphs on bounded treewidth are probably optimal. In *SODA*, pages 777–789, 2011.
36. Daniel Lopresti and Andrew Tomkins. Block edit models for approximate string matching. *Theoretical Computer Science*, 181:159–179, 1997.
37. J. Erickson M. Soss and M. H. Overmars. Preprocessing chains for fast dihedral rotations is hard or even impossible. *Computational Geometry: Theory and Applications*, 26(3):235–246, 2002.
38. W.J. Masek and M.S. Paterson. A faster algorithm computing string edit distances. *Journal of Computer and System Sciences*, 20, 1980.
39. J. Nešetřil and S. Poljak. On the complexity of the subgraph problem. *Comment. Math. Univ. Carolin.*, 26(2):415–419, 1985.
40. Mihai Pătraşcu. Towards polynomial lower bounds for dynamic problems. In *STOC*, pages 603–610, 2010.
41. R. Paturi, P. Pudlák, M. E. Saks, and F. Zane. An improved exponential-time algorithm for k -SAT. *J. ACM*, 52(3):337–364, 2005.
42. Krzysztof Pietrzak. On the parameterized complexity of the fixed alphabet shortest common supersequence and longest common subsequence problems. *Journal of Computer and System Sciences*, 67(4):757–771, 2003.
43. M. Pătraşcu and R. Williams. On the possibility of faster SAT algorithms. In *SODA*, pages 1065–1075, 2010.
44. M.S. Rahman, C. Iliopoulos, I. Lee, M. Mohamed, and W.F. Smyth. Finding patterns with variable length gaps or don't cares. In *COCOON*, pages 146–155, 2006.
45. L. Roditty and V. Vassilevska Williams. Fast approximation algorithms for the diameter and radius of sparse graphs. In *STOC*, pages 515–524, 2013.
46. U. Schöning. A probabilistic algorithm for k -SAT and constraint satisfaction problems. In *FOCS*, pages 410–414, 1999.
47. T.F. Smith and M.S. Waterman. The identification of common molecular subsequences. *Journal of Molecular Biology*, 147:195–197, 1981.
48. V. Vassilevska and R. Williams. Finding, minimizing, and counting weighted subgraphs. In *STOC*, pages 455–464, 2009.
49. Lusheng Wang and Tao Jiang. On the complexity of multiple sequence alignment. *Journal of computational biology*, 1(4):337–348, 1994.
50. Ryan Williams. A new algorithm for optimal constraint satisfaction and its implications. In *ICALP*, pages 1227–1237, 2004.
51. V. Vassilevska Williams and R. Williams. Subcubic equivalences between path, matrix and triangle problems. In *FOCS*, pages 645–654, 2010.
52. G. J. Woeginger. Space and time complexity of exact algorithms: Some open problems. In *IWPEC*, pages 281–290, 2004.