# Approximation Algorithms for LCS and LIS with Truly Improved Running Times[*]

Aviad Rubinstein[†]    Saeed Seddighin[‡]    Zhao Song[§]    Xiaorui Sun[¶]

November 23, 2021

## Abstract

Longest common subsequence (LCS) is a classic and central problem in combinatorial optimization. While LCS admits a quadratic time solution, recent evidence suggests that solving the problem may be impossible in truly subquadratic time. A special case of LCS wherein each character appears at most once in every string is equivalent to the longest increasing subsequence problem (LIS) which can be solved in quasilinear time. In this work, we present novel algorithms for approximating LCS in truly subquadratic time and LIS in truly sublinear time. Our approximation factors depend on the ratio of the optimal solution size over the input size. We denote this ratio by $\lambda$ and obtain the following results for LCS and LIS without any prior knowledge of $\lambda$.

- A truly subquadratic time algorithm for LCS with approximation factor $\Omega(\lambda^3)$.
- A truly sublinear time algorithm for LIS with approximation factor $\Omega(\lambda^3)$.

Triangle inequality was recently used by Boroujeni, Ehsani, Ghodsi, HajiAghayi and Seddighin [BEG+18] and Charkraborty, Das, Goldenberg, Koucky and Saks [CDG+18] to present new approximation algorithms for edit distance. Our techniques for LCS extend the notion of triangle inequality to non-metric settings.

---

# 1   Introduction

We consider three problems in combinatorial optimization: the longest common subsequence (LCS), the edit distance (ED), and the longest increasing subsequence (LIS). The LCS of two strings $A$ and $B$ is simply their longest (not necessarily contiguous) common subsequence. The edit distance is defined as the minimum number of character deletions, insertions, and substitutions required to transform $A$ into $B$. For the purpose of our discussion, we consider a more restricted definition of edit distance where substitutions are not allowed[1]. Longest increasing subsequence is equivalent to a special case of LCS where the input strings are permutations. All three problems are very fundamental and have been subject to a plethora of studies in the past few decades and specially in recent years  [LMS98, BYJKK04, BES06, AO09, AKO10, SS10, BI15, ABW15, AHWW16, AB17, AR18, CGL+19, BEG+18, CDG+18, HSSS19].

If the strings have length $n$, both LCS and ED can be solved in quadratic time ($O(n^2)$) with dynamic programming. These running times are slightly improved to $O(n^2/\log^2(n))$ by Masek and Paterson [MP80], however, efforts to improve the running time to $O(n^{2-\Omega(1)})$ for either edit distance or LCS were all futile.

In recent years, our understanding of the source of complexity for these problems improved thanks to a sequence of fine-grained complexity developments [ABW15, AHWW16]. We now know that assuming the strong exponential time hypothesis (SETH) [ABW15], or even weaker assumptions such as the orthogonal vectors conjecture (OVC) [ABW15] or branching-program-SETH [AHWW16], there are no truly sub-quadratic[2] time algorithms for LCS. Similar results also hold for edit distance [BI15].

The classic approach to break the quadratic barrier for these problems is approximation algorithms. Note that for (multiplicative) approximations, LCS and edit distance are no longer equivalent (much like we have a 2-approximation algorithm for Vertex Cover, but Independent Set is NP-hard to approximate within near-linear factors).

For edit distance, an $\widetilde{O}(n + \Delta^2)$-time algorithm of [LMS98] (where $\Delta$ is the true edit distance between the strings) implies a linear-time $\sqrt{n}$-approximation algorithm. The approximation factor has been significantly improved in a series of works to $O(n^{3/7})$ [BYJKK04], to $O(n^{0.34})$ [BES06], to $O(2^{\widetilde{O}(\sqrt{\log n})})$ [AO09][3], and finally to polylogarithmic [AKO10]. A recent work of Boroujeni *et al.* [BEG+18] obtains a constant factor approximation quantum algorithm for edit distance that runs in truly subquadratic time. Finally, the breakthrough of Chakraborty *et al.* [CDG+18] gave a classic (randomized) constant factor approximation for edit distance in truly subquadratic time. A key component in both of the latest constant factor approximation algorithms is the application of triangle inequality (for edit distance between certain substrings of the input). A particular challenge in extending these ideas to LCS is that LCS is not a metric and in particular does not satisfy the triangle inequality.

Our understanding of the complexity of approximate solutions for LCS is embarrassingly limited. For general strings, there are several linear-time $1/\sqrt{n}$-approximation algorithms based on sampling techniques. For alphabet size $|\Sigma|$, there is a trivial $1/|\Sigma|$-approximation algorithm that runs in linear time. Whether or not these approximation factors can be improved by keeping the running time linear is one of the central problems in fine-grained complexity. Very recently, both the general $1/\sqrt{n}$-approximation factor, and, for binary strings, the $1/2$-approximation factor, have been slightly improved ([HSSS19] and [RS20], respectively). These works give improved algorithm for the two

---

[1] Alternatively, the cost of a substitution is doubled as it requires a deletion and an insertion.

[2] By *truly sub-quadratic* we mean $O(n^{2-\epsilon})$, for any constant $\epsilon > 0$

[3] We define $\widetilde{O}(f)$ to be $f \cdot \log^{O(1)}(f)$.

extreme cases where the size of the alphabet is very small or very large. In comparison, our approximation guarantee depends on the solution size rather than the size of the alphabet. Also, for the special case of balanced strings, we improve upon the result of [RS20] by obtaining an $o(|\Sigma|)$ approximate solution in subquadratic time. There are a few fine-grained complexity results for approximate LCS, but they only hold against deterministic algorithms, and rely on very strong assumptions [AB17, AR18, CGL$^+$19].

## 1.1 Our Results

For simplicity, we use $\mathsf{lcs}(A, B)$ to denote the size (not the whole sequence) of the longest common subsequence for two strings $A$ and $B$. Similarly, we use $\mathsf{ed}(A, B)$ to denote the edit distance and $\mathsf{lis}(A)$ for the size of the longest common subsequence. We sometimes normalize the solution by the length of the strings so that the size of the solution remains in the interval $[0, 1]$. We refer to the normalized solutions by $||\mathsf{lcs}(A, B)|| = \mathsf{lcs}(A, B)/n$ and $||\mathsf{ed}(A, B)|| = \mathsf{ed}(A, B)/2n$ (here both strings have equal length $n$), and $||\mathsf{lis}(A)|| = \mathsf{lis}(A)/n$. In this way, $||\mathsf{ed}(A, B)|| + ||\mathsf{lcs}(A, B)|| = 1$ (assuming both strings have equal length).

As mentioned earlier, recent developments for edit distance are based on a simple but rather useful observation. Edit distance satisfies triangle inequality, or in other words, given three strings $A_1, A_2, A_3$ of length $n$ such that $||\mathsf{ed}(A_1, A_2)|| \leq \delta$ and $||\mathsf{ed}(A_2, A_3)|| \leq \delta$ hold, we can easily imply that $||\mathsf{ed}(A_1, A_3)|| \leq 2\delta$. While $\mathsf{lcs}$ does not satisfy the triangle inequality in any meaningful way, it does, *on average*, satisfy the following birthday-paradox-like property that we call *birthday triangle inequality*.

**Property 1.1** (birthday triangle inequality). *Consider three equal-length strings $A_1$, $A_2$, and $A_3$ such that $||\mathsf{lcs}(A_1, A_2)|| \geq \lambda$ and $||\mathsf{lcs}(A_2, A_3)|| \geq \lambda$. If the common subsequences correspond to random indices of each string, we expect that $||\mathsf{lcs}(A_1, A_3)|| \geq \lambda^2$.*

Of course, this is not necessarily the case in general. More precisely, it is easy to construct examples[4] in which $||\mathsf{lcs}(A_1, A_2)|| = 1/2$ and $||\mathsf{lcs}(A_2, A_3)|| = 1/2$, but $||\mathsf{lcs}(A_1, A_3)|| = 0$. Our main result shows that while it only holds on average, we can algorithmically replace the triangle inequality for edit distance with the birthday triangle inequality *on worst case inputs*.

**Theorem 1.2** (Main Theorem, formally stated as Theorems 2.1 and 2.2). *Given strings $A, B$ both of length $n$ such that $||\mathsf{lcs}(A, B)|| = \lambda$, we can approximate the length of the LCS between the two strings within an $\Omega(\lambda^3)$ factor in subquadratic time. The approximation factor improves to $(1-\epsilon)\lambda^2$ when $1/\lambda$ is constant.*

We remark that our algorithm is actually able to output the whole sequence of the solution, but we only focus on estimating the size of the solution for simplicity. We begin by comparing our main theorem to previous work on edit distance. In this case, $1/\lambda$ is constant w.l.o.g.[5] and therefore the approximation factor of our algorithm is $(1 - \epsilon)\lambda^2$. If $\delta = ||\mathsf{ed}(A, B)||$, then our LCS algorithm outputs a transformation from $A$ to $B$ using at most $2n(1 - (1 - \epsilon)(1 - \delta)^3)$ operations. Observe that when the strings are not overly close and $\delta = \Omega(1)$ by scaling $\epsilon$, we already recover a $(3 + \epsilon')$-approximation for edit distance in truly subquadratic time. For mildly far strings, say $\delta = 0.1$, a more careful look at the expansion of $(1 - \delta)^3$ reveals that we save an additive $\Theta(\delta^2)$ in the approximation factor. For example, with $\delta = 0.1$ our approximation factor for edit distance is 2.71 instead of 3.

---

[4]For example, $A_1 = 0^{n/2}0^{n/2}, A_2 = 0^{n/2}1^{n/2}, A_3 = 1^{n/2}1^{n/2}$.

[5]When we use our solution to approximate edit distance, we can safely assume that $||\mathsf{lcs}(A, B)|| = \Omega(1)$ since otherwise the edit distance of the two strings is very close to $2n$.

An interesting implication of our main result is for LCS over a large alphabet $\Sigma$, where the optimum $||\mathsf{lcs}(A, B)||$ may be much smaller than 1. This is believed to be the hardest regime for approximation algorithms (and indeed the only one for which we have any conditional hardness of approximation results [AB17, AR18, CGL+19]). Here, we consider instances that satisfy a mild balance assumption: we assume that there is a character that appears with frequency at least $1/|\Sigma|$ in both strings[6]. Then, our main theorem implies an $O(1/|\Sigma|^{3/4})$-approximate solution in truly subquadratic time (the first improvement over the trivial $1/|\Sigma|$ approximation in this regime).

**Corollary 1.3** (LCS, formally stated as Corollary 2.3). *Given a pair of strings $(A, B)$ of length $n$ over alphabet $\Sigma$ that satisfy the balance condition, we can approximate their LCS within an $O(|\Sigma|^{3/4})$ factor in truly subquadratic time.*

Next, we show that a similar result can be obtained for LIS. Perhaps coincidentally, the approximation factor of our algorithm is also $\Omega(\lambda^3)$ which is same to LCS, but the technique is completely different. Although LIS can be solved exactly in time $O(n \log n)$, there have been several attempts to approximate the size of LIS and related problems in sublinear time [Sch61, Fre75, DGL+99, EKK+00, Fis04, ACCL07, SS10]. The best known solution is due to the work of Saks and Seshadhri [SS10] that obtains a $(1 + \epsilon)$-approximate algorithm for LIS in polylogarithmic time, when the solution size is at least a constant fraction of the input size [7]. In other words, if $||\mathsf{lis}(A)|| = \lambda$ and $1/\lambda$ is constant, their algorithm approximates $\mathsf{lis}(A)$ in polylogarithmic time. However, this only works if $1/\lambda$ is constant and even if $1/\lambda$ is logarithmically large, their method fails to run in sublinear time[8]. We complement the work of Saks and Seshadhri [SS10] by presenting a result for LIS similar to our result for LCS. More precisely, we show that when $||\mathsf{lis}(A)|| = \lambda$, an $\Omega(\lambda^3)$ approximation of LIS can be obtained in truly sublinear time. Although our approximation factor is worse than that of [SS10], our result works for any (not necessarily constant) $\lambda$.

**Theorem 1.4** (LIS, formally stated as Theorem 6.8). *Given an array $A$ of $n$ integer numbers such that $||\mathsf{lis}(A)|| = \lambda$. We can approximate the length of the LIS for $A$ in sublinear time within a factor $\Omega(\lambda^3)$.*

If one favors the running time over the approximation factor, it is possible to improve the exponent of $n$ in the running time down to any constant $\kappa > 0$ at the expense of incurring a larger multiplicative factor to the approximation.

## 1.2 Preliminaries

In LCS or edit distance, we are given two strings $A$ and $B$ as input. We assume for simplicity that the two strings have equal length and refer to that by $n$. In LCS, the goal is to find the largest subsequence of the characters which is shared between the two strings. In edit distance, the goal is to remove as few characters as possible from the two strings such that the remainders for the two strings are the same. We use $\mathsf{lcs}(A, B)$ and $\mathsf{ed}(A, B)$ to denote the size of the longest common subsequence and the edit distance of two strings $A$ and $B$.

In LIS, the input contains an array $A$ of $n$ integer numbers and the goal is to find a sequence of elements of $A$ whose values (strictly) increase as their indices increase. For LIS, we denote the solution size for an array $A$ by $\mathsf{lis}(A)$. We also use $\mathsf{lis}^{[\alpha,\beta]}(A)$ to denote the size of the longest

---

[6]Note that in every instance in each string there is a character that appears with frequency at least $1/|\Sigma|$, but in general that may not be the same character.

[7]Their algorithm obtains an additive error of $\delta n$ in time $2^{\tilde{O}(1/\delta)}$. When the solution size is bounded by $\lambda n$, one needs to set $\delta < \lambda$ in order to guarantee a multiplicative factor approximation.

[8]There is a term $(1/\lambda)^{1/\lambda}$ in the running time.

increasing subsequence subject to elements whose values lie in range $[\alpha, \beta]$. Longest increasing subsequence is equivalent to LCS when the inputs are two permutations of $n$ distinct characters.

Finally, we define a notation to denote the normalized solution sizes. For LCS, we denote the normalized solution size by $||\mathsf{lcs}(A, B)|| = \mathsf{lcs}(A, B)/\sqrt{|A||B|}$ for $A$ and $B$ and we use $||\,\mathsf{ed}(A, B)|| = \mathsf{ed}(A, B)/(2\sqrt{|A||B|})$ for edit distance. Note that, when the two strings have equal length we have $||\,\mathsf{ed}(A, B)|| + ||\mathsf{lcs}(A, B)|| = 1$. Similarly, for longest increasing subsequence, we denote by $||\mathsf{lis}(A)|| = \mathsf{lis}(A)/|A|$ the normalized solution size. We usually refer to the size of the input array by $n$.

Throughout this paper, we call an algorithm $f(\lambda)$-approximation for LCS if it is able to distinguish the following two cases: i) $||\mathsf{lcs}(A, B)|| \geq \lambda$ or ii) $||\mathsf{lcs}(A, B)|| < \lambda f(\lambda)$. A similar definition carries over to LIS. Once an $f(\lambda)$-approximation algorithm is provided for either LCS or LIS, one can turn it into an algorithm that outputs a solution with size $f(\lambda)(1-\epsilon)\lambda n$ provided that the optimal solution has a size $\lambda n$. The algorithm is not aware of the value of $\lambda$ but will start with $\lambda_0 = 1$ and iteratively multiply $\lambda_0$ by $1 - \epsilon$ until a solution is found.

## 1.3 Techniques Overview

Our algorithm for LCS is inspired by the recent developments for edit distance [BEG$^+$18, CDG$^+$18]. We begin by briefly explaining the previous ideas for approximating edit distance and then we show how we use these techniques to obtain a solution for LCS. Finally, in Section 1.3.2 we outline our algorithm for LIS.

### 1.3.1 Summary of Previous ED Techniques

Indeed, edit distance is hard only if the two strings are far ($||\,\mathsf{ed}(A, B)|| = \delta$ and $\delta = n^{-o(1)}$) otherwise the $O(n + (n\delta)^2)$ algorithm of Landau *et al.* [LMS98] computes the solution in truly subquadratic time. The algorithm of Chakraborty *et al.* [CDG$^+$18] for edit distance has three main steps that we briefly discuss in the following.

**Step 0 (window-compatible solutions):** In the first step, they construct a set of *windows*, or (contiguous) substrings, $W_A$ for string $A$, and $W_B$ for string $B$. Let $k$ denote the total number of windows of $W_A \cup W_B$. For simplicity, let all the windows have the same size $d$ and $n \simeq O(kd)$[9]. The construction features two key properties: 1) provided that the edit distances of the windows between $W_A$ and $W_B$ are available, one can recover a $1 + \epsilon$ approximation of edit distance in time $\widetilde{O}(n + k^2)$ via dynamic programming. 2) $k^2 \times d^2 \simeq O(n^2)$. That is, if we naively compute the edit distance of every pair of windows, the overall running time would still asymptotically be the same as that of the classic algorithm.

In order to obtain a solution for edit distance, it suffices to know the distances between the windows. However, Chakraborty *et al.* [CDG$^+$18] show that knowing the distances between some of the window pairs is enough to obtain an approximately optimal solution for edit distance. Step 1 provides estimates for the distances of the windows which is approximately correct except for $O(k^{2-\Omega(1)})$ many pairs and Step 2 shows how this can be used to obtain a solution for edit distance. Discretization simplifies the problem substantially. For a fixed $0 \leq \delta \leq 1$, they introduce a graph $\mathsf{G}_\delta$ where the nodes correspond to the windows and an edge between window $w_i \in W_A$ and window $w_j \in W_B$ means that $||\,\mathsf{ed}(w_i, w_j)|| \leq \delta$. If we are able to construct $\mathsf{G}_\delta$ for logarithmically different choices of $\delta$, we can as well estimate the distances within a $1 + \epsilon$ factor for the windows. Therefore

---

[9]The equality holds if we assume $\delta = \Omega(1)$.

the problem boils down to constructing $\mathsf{G}_\delta$ for a fixed given $\delta$ without computing the edit distance between all pairs of windows.

**Step 1 (sparsification via triangle inequality):** This step is the heart of the algorithm. Suppose we choose a high-degree vertex $v$ from $\mathsf{G}_\delta$ and discover all its incident edges by computing its edit distance to the rest of the windows. Triangle inequality implies that every pair of windows in $N(v)$ has a distance bounded by $2\delta$. Therefore by losing a factor 2 in the approximation, one can put all these edges in $\mathsf{G}_\delta$ and not compute the edit distances explicitly. Although this does save some running time, in order to make sure the running time is truly subquadratic, we need to make a similar argument for paths of length 3 and thereby lose a factor 3 in the approximation. This method sparsifies the graph and what remains is to discover the edges of a sparse graph.

**Step 2 (discovering the edges of the sparse graph):** Step 1 uses triangle inequality and discovers many edges between the vertices of $\mathsf{G}_\delta$. However, it may not discover all the edges completely. When in the remainder graph, the degrees are small (and hence the graph is sparse) triangle inequality does not offer an improvement and thus a different approach is required. Roughly speaking, Chakraborty *et al.* [CDG$^+$18] subsample the windows of $W_A$ into a smaller set $S$ and discover all pairs of windows $w_i \in S$ and $w_j \in W_B$ such that edge $(i, j)$ is not discovered in Step 1. Next, they compute the edit distance of each pair of windows $(w_i, w_j), w_i \in W_A, w_j \in W_B$ such that there exist two nearby windows $(w_a, w_b)$ satisfying $w_a \in S, w_b \in W_B$ and the edge between $w_a$ and $w_b$ was missed in Step 1. The key observation is that even though this procedure does not discover all the edges, the approximated distances lead to an approximate solution for edit distance.

### 1.3.2 LCS

Our algorithm for LCS mimics the same guideline. In addition to this, Steps 0 and 2 of our algorithm are LCS analogues of the ones used by Chakraborty *et al.* [CDG$^+$18]. The main novelty of our algorithm is Step 1 which is a replacement for triangle inequality. Recall that unlike edit distance, triangle inequality does not hold for LCS.

**Challenge 1.5.** *How can we introduce a notion similar to triangle inequality to a non-metric setting such as LCS?*

We introduce the notion of birthday triangle inequality to overcome the above difficulty. Given windows $w_1$, $w_2$, and $w_3$ of size $d$ such that $||\mathsf{lcs}(w_1, w_2)|| \geq \lambda$ and $||\mathsf{lcs}(w_2, w_3)|| \geq \lambda$ hold, what can we say about the LCS of $w_1$ and $w_3$? In general, nothing! $||\mathsf{lcs}(w_1, w_3)||$ could be as small as 0. However, let us add some randomness to the setting. Think of the LCS of $w_1$ and $w_2$ as a matching from the characters of $w_1$ to $w_2$ and similarly the LCS for $w_2$ and $w_3$ as another matching between characters of $w_2$ and $w_3$. Assume (for the sake of the thought experiment) that the characters of $w_2$ appear randomly in each matching. Since $||\mathsf{lcs}(w_1, w_2)|| \geq \lambda$, each character of $w_2$ appears with probability at least $\lambda$ in the matching between $w_1$ and $w_2$. A similar argument implies that each character of $w_2$ appears with probability $\lambda$ in the matching of $w_2$ and $w_3$. Thus, (assuming independence), each character of $w_2$ appears in both matchings with probability $\lambda^2$. This means that in expectation, there are $\lambda^2 d$ paths of length 2 between $w_1$ and $w_3$ which suggests $||\mathsf{lcs}(w_1, w_3)|| \geq \lambda^2$ as shown in Figure 1. This is basically birthday paradox used for the sake of triangle inequality.

Replacing triangle inequality by birthday triangle inequality is particularly challenging since birthday triangle inequality only holds on average. In contrast, triangle inequality holds for any tuple of arbitrary strings. Most of our technical discussions is dedicated to proving that we can algorithmically use birthday triangle inequality to obtain a solution for the worst-case scenarios.
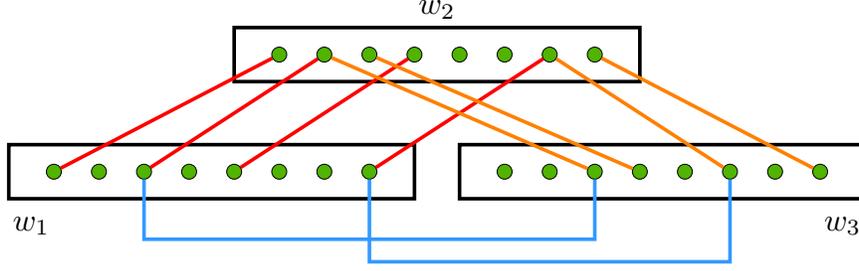
Figure 1: Birthday paradox for triangle inequality: let $w_1, w_2, w_3$ be three windows of length $d = 8$ and assume $\lambda = 1/2$. The LCS between $w_1$ and $w_2$ is $\lambda d = 4$ and the LCS between $w_2$ and $w_3$ is $\lambda d = 4$. Finally due to birthday paradox, we expect that the LCS between $w_1$ and $w_3$ is $\lambda^2 d = 2$.

The most inconvenient step of our analysis is to show that our algorithm estimates the LCS of most window pairs in the sparsification phase. While this is straightforward for edit distance, birthday triangle inequality requires a deeper analysis of the underlying graph. In particular, we need to prove that if the undiscovered edges are too many, then birthday triangle inequality can be applied to certain neighborhoods of the graph.

There are two difficulties that we face here. On one hand, in order to apply birthday triangle inequality to a subgraph, we need to have enough structure for that subgraph to show the implication can be made. On the other hand, our assumptions cannot be too strong, otherwise such neighborhoods may not cover the edges of the graph. Therefore, the first challenge that we need to overcome is characterizing subgraphs in which birthday triangle inequality is guaranteed to be applicable. Our suggestion is the *bi-cliques* structure. Although combinatorial techniques seem unlikely to prove this, we use the Blakley-Roy inequality to show that in a large enough bi-clique, we can use birthday triangle inequality to imply a bound on the LCS of certain pairs. The second challenge is to prove that if the underlying graph is dense enough, the graph contains many bi-cliques that cover almost all the edges that we plan to discover. This is again a challenging graph theoretic problem. We leverage extremal graph theory tools such as Turan's theorem for cliques and bi-cliques to obtain this bound.

Similar to edit distance, we construct a set $W = W_A \cup W_B$ of $k$ windows in Step 0 and aim to sparsify the edges of the lcs-graph in Step 1. Our construction ensures that $kd \simeq \Theta(n)$ and that knowing the LCS of the window pairs suffices to approximate the LCS of the two strings. For a threshold $0 \leq \lambda \leq 1$, define a matrix $\mathsf{O} : [k] \times [k] \to \{0, 1\}$ to be a matrix which identifies whether $||\mathsf{lcs}(w_i, w_j)|| \geq \lambda$. In other words, $\mathsf{O}[i][j] = 1 \iff ||\mathsf{lcs}(w_i, w_j)|| \geq \lambda$. For an $0 < \alpha \leq 1$, we call a matrix $\mathsf{O}_\alpha$ an $\alpha$ approximation of $\mathsf{O}$ if it meets the following two conditions:

$$\mathsf{O}_\alpha[i][j] = 0 \implies ||\mathsf{lcs}(w_i, w_j)|| < \lambda \qquad \text{and} \qquad \mathsf{O}_\alpha[i][j] = 1 \implies ||\mathsf{lcs}(w_i, w_j)|| \geq \alpha \cdot \lambda$$

Notice that $\mathsf{O}_\alpha$ gives more flexibility than $\mathsf{O}$ for the cases that $\lambda\alpha \leq \mathsf{lcs}(w_i, w_j) < \lambda$. That is, both 0 and 1 are acceptable in these cases. Indeed an $\alpha$ approximation algorithm for the above problem is enough to obtain an $\alpha$ approximation algorithm for LCS. However, this is not necessary as Step 2 allows for incorrect approximation for up to $k^{2-\Omega(1)}$ many window pairs. Therefore, the problem of approximating LCS essentially boils down to approximating $\mathsf{O}$ for a given basket of windows $W = W_a \cup W_b$ and a fixed $\lambda$ by allowing sufficiently small error in the output. A naive solution is to iterate over all pairs $w_i$ and $w_j$ and compute $\mathsf{lcs}(w_i, w_j)$ in time $O(d^2)$ and determine $\mathsf{O}$ accordingly. However, this amounts to a total running time of $O(k^2 d^2)$ which is quadratic and not desirable. In order to save time, we need to compute the LCS of fewer than $k^2$ pairs of windows.

6

To make this possible, we allow our algorithm to miss up to $O(k^{2-\Omega(1)})$ edges of the graph. Step 2 ensures that this does not hurt the approximation factor significantly.

We construct a graph from the windows wherein each vertex corresponds to a window and each edge identifies a pair with a large LCS (in terms of $\lambda$). Let us call this graph the lcs-graph and denote it by $\mathsf{G}_\lambda$. The goal is to detect the edges of the graph by allowing false-positive. As we discussed earlier, the hard instances of the problem are the cases where the lcs-graph is dense for which we need a sparsifier. Roughly speaking, in our sparsification technique, our algorithm constructs another graph $\widehat{\mathsf{G}}_\lambda$ such that $\widehat{\mathsf{G}}_\lambda$ is valid in the sense that the edges of $\widehat{\mathsf{G}}_\lambda$ correspond to pairs of windows with large enough LCS. In addition to this, our algorithm guarantees that after the removal of the edges of $\widehat{\mathsf{G}}_\lambda$ from $\mathsf{G}_\lambda$ the remainder is sparse. In other words, $|E(\mathsf{G}_\lambda) \setminus E(\widehat{\mathsf{G}}_\lambda)| = O(|V(\mathsf{G}_\lambda)|^{2-\Omega(1)})$. Of course, if the overall running time of the sparsification phase is truly subquadratic, the error of undiscovered edges can be addressed by the techniques of [CDG+18] in Step 2. Below, we bring a formal definition for sparsification.

---

sparsification

input: Windows $w_1, w_2, \ldots, w_k$, parameters $\lambda$, and $\alpha$.
solution: A matrix $\widehat{\mathsf{O}}_\alpha \in \{0,1\}^{k \times k}$ such that:

- $\widehat{\mathsf{O}}_\alpha[i][j] = 1 \implies ||\mathsf{lcs}(w_i, w_j)|| \geq \alpha \cdot \lambda$

- $\left| \left\{ (i,j) \mid ||\mathsf{lcs}(w_i, w_j)|| \geq \lambda \text{ and } \widehat{\mathsf{O}}_\alpha[i][j] = 0 \right\} \right| = k^{2-\Omega(1)}$

---

We present two sparsification techniques for LCS. The first one (Section 3.1), has an approximation factor of $(1-\epsilon) \cdot \lambda^2$. In Section 3.2 we present another sparsification technique that has a worse approximation factor $\Omega(\lambda^3)$ but leaves fewer edges behind. Although the second sparsification technique has a worse approximation factor, it has the advantage that the number of edges that remain in the sparse graph is truly subquadratic regardless of the value of $\lambda$ and therefore it extends our solution to the case that $\lambda = o(1)$ (see Section 3.2 for a detailed discussion).

Let us note one last algorithmic challenge to keep in mind before we begin to describe our sparsification techniques. For edit distance, if window pairs $(w_1, w_2)$ and $(w_2, w_3)$ are close, we are *guaranteed* that $w_1$ and $w_3$ are also close; for longest common subsequence, we will argue that $(w_1, w_3)$ are *likely* to be have a long LCS (for a "random" choice of $(w_1, w_3)$). Nonetheless, in order to add $(w_1, w_3)$ as an edge to our graph we have to *verify* that their LCS is indeed long. If we were to verify an edge naively, we would need as much time as computing the LCS between $(w_1, w_3)$ from scratch!

**Sparsification 1, $(1-\epsilon)\lambda^2$-approximation** Similar to edit distance, applying the birthday variant of triangle inequality to paths of length 2 for LCS does not improve the running time significantly. Therefore, we need to use birthday triangle inequality for paths of length 3. To this end, we define the notion of constructive tuples as follows: a tuple $\langle w_i, w_a, w_b, w_j \rangle$ is an $(\epsilon, \lambda)$-constructive tuple, if we have $||\mathsf{lcs}(w_i, w_a)|| \geq \lambda$, $||\mathsf{lcs}(w_i, w_j)|| \geq \lambda$, $||\mathsf{lcs}(w_b, w_j)|| \geq \lambda$ and by taking the intersection of the three LCS matchings, we are able to imply $||\mathsf{lcs}(w_a, w_b)|| \geq (1-\epsilon)\lambda^3$ (see Figure 2 for an example). Taking the intersection of the matchings can be done in linear time which is faster than computing the LCS.

Our sparsification technique here is simple but the analysis is very intricate. We subsample a set $S$ of windows and compute the LCS of every window in $S$ and all other windows. We set $|S| = k^\gamma \log k$, where $\gamma \in (0,1)$. At this point, for some pairs, we already know their LCS. However, if neither $w_i$ nor $w_j$ is in $S$, we do not know if $||\mathsf{lcs}(w_i, w_j)|| \geq \lambda$ or not. Therefore, for such pairs,
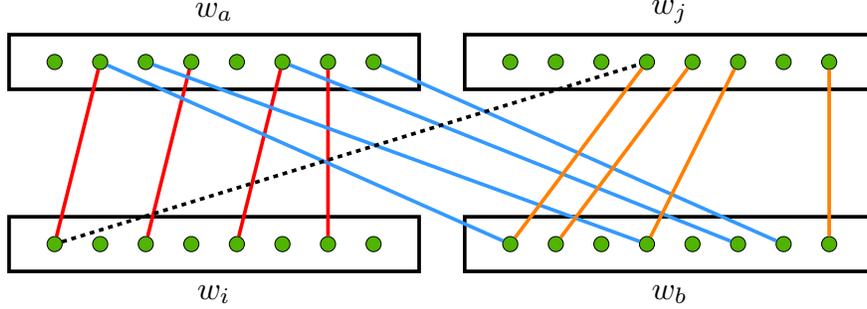
Figure 2: Let $w_i, w_a, w_b$ and $w_j$ denote four windows and each of them has length $d = 8$. This figure shows how the intersection of the edges of three windows are taken in order to construct a solution for the LCS of $w_i$ and $w_j$. If the size of the intersection is large, then such a tuple is called constructive. The solid lines represent LCS between two strings, and the dashed line represents the intersection of the three LCSs.

we try to find windows $w_a, w_b \in S$ such that $\langle w_i, w_a, w_b, w_j \rangle$ is constructive. If such a constructive tuple is found for a pair of windows, then we conclude that their normalized LCS is at least $(1-\epsilon)\lambda^3$.

All that remains is to argue that this method discovers almost all the edges of the lcs-graph $\mathsf{G}_\lambda$ and the number of undiscovered edges is $k^{2-\Omega(1)}$. This is the most difficult part of the analysis. We note that even proving the existence of **only one** constructive tuple is already non-trivial **even when $\mathsf{G}_\lambda$ is complete**. However, our goal is to show that almost all the edges are discovered via constructive tuples when $\mathsf{G}_\lambda$ is dense (and of course not necessarily complete).

Define a pair of windows $w_i$ and $w_j$ to be *well-connected*, if there are at least $k^{2-\gamma}$ different $(w_a, w_b)$ pairs such that $\langle w_i, w_a, w_b, w_j \rangle$ is $(\epsilon, \lambda)$-constructive. Since each window appears in $S$ with probability $k^{\gamma-1} \log k$, for each well-connected pair we find one constructive tuple via our algorithm with high probablity. Therefore, we need to prove that the total number of pairs $(w_i, w_j)$ such that $(w_i, w_j)$ is not well-connected but $\|\mathsf{lcs}(w_i, w_j)\| \geq \lambda$ is subquadratic. Let us put these edges in a new graph $\mathsf{NG}_\lambda$ whose vertices are all the windows.

We first leverage the Blakley-Roy inequality and a double counting technique to prove that if $\mathsf{NG}_\lambda$ has a large complete bipartite subgraph, then there is one constructive tuple which includes only the vertices of this subgraph (Lemma 3.12). Next, we apply the Turan's theorem to show that if $\mathsf{NG}_\lambda$ is dense, then it has a lot of large complete bipartite subgraphs. Finally, we use a probabilistic method to conclude that $\mathsf{NG}_\lambda$ cannot be too dense otherwise there are a lot of constructive tuples in the graph which implies that at least one edge $(w_i, w_j)$ in $\mathsf{NG}_\lambda$ is well-connected. This is not possible since all the well-connected pairs are detected in our sparsification algorithm with high probability.

The above argument proves that if we sparsify our graph using our sparsification algorithm, the remainder graph would have a subquadratic number of edges. Therefore after plugging Step 2 into the algorithm, the running time remains subquadratic. However, since Turan theorem gives us a weak bound, the running time of the algorithm using this sparsification is $O(n^{2-\Omega(\lambda)})$ and is only truly subquadratic when $1/\lambda$ is constant.

**Sparsification 2, $\Omega(\lambda^3)$-approximation**    In Section 3.1, we present another sparsification method that although gives us a slightly worse approximation factor $\Omega(\lambda^3)$ it always leaves a truly subquadratic number of edges behind and therefore the running time of the algorithm would be truly subquadratic regardless of the parameter $\lambda$. This sparsification is based on a novel data structure.

8

We assume for simplicity in the following that all the windows are of the same length even though this does not hold in general and having windows of different length adds more complication to the algorithm. We discuss the details in Section 3.1.

Let $\mathsf{opt}_{i,a}$ denote the longest common subsequence of $w_i$ and $w_a$ (with some fixed tie-breaking rule, e.g. lexicographically first). Define $\mathsf{lcs}_{w_a}(w_i, w_j)$ to be the size of the longest common subsequence between $\mathsf{opt}_{i,a}$ and $w_j$. Notice that this definition is no longer symmetric. Let $\|\mathsf{lcs}_{w_a}(w_i, w_j)\|$ denote the relative value, i.e., $\|\mathsf{lcs}_{w_a}(w_i, w_j)\| = \mathsf{lcs}_{w_a}(w_i, w_j)/\sqrt{|w_i| \cdot |w_j|}$. The first ingredient of the algorithm is a data-structure, namely $\mathsf{lcs\text{-}cmp}$. After a preprocess of time $O(|w_a| \sum_{i \in S} |w_i|)$, $\mathsf{lcs\text{-}cmp}$ is able to answer queries of the following type in time $O(|w_i| + |w_j|)$:

- "for a $0 \leq \widetilde{\lambda} \leq 1$ either certify that $\|\mathsf{lcs}_{w_a}(w_i, w_j)\| \geq \Omega(\widetilde{\lambda}^2)$ or report that

$$\|\mathsf{lcs}_{w_a}(w_i, w_j)\| < O(\widetilde{\lambda})"$$

.

In our sparsification, we repeat the following procedure $k^\gamma$ times, where $\gamma \in (0, 1)$. We sample a window $w_a$ uniformly at random and construct $\mathsf{lcs\text{-}cmp}(w_a, S)$ for $S = \{w_i | i \neq a \text{ and } |w_i| \geq |w_a|\}$. After the preprocessing step, we make a query for every pair of windows $(w_i, w_j)$ such that $w_i, w_j \in S$ and determine if $\mathsf{lcs}_{w_a}(w_i, w_j)$ is at least $\Omega(\lambda^4)$ or upper bounded by $O(\lambda^2)$ (here $\widetilde{\lambda} = \lambda^2/2$). If their $\mathsf{LCS}$ is at least $\Omega(\lambda^4)$ we report this pair as an edge in our $\mathsf{lcs}$-graph. Finally, we use the Turan theorem to prove that the number of remaining edges in our graph is small.

To be more precise, we first construct a graph $\mathsf{NG}_\lambda$ that reflects the edges that are not detected via our sparsification. If $\mathsf{NG}_\lambda$ is dense enough, then there is one vertex $v$ in $\mathsf{NG}_\lambda$ with a large enough degree. We use the neighbors of $v$ to construct another graph $\mathsf{NF}_\lambda$ with vertex set $N(v)$. An edge exists in $\mathsf{NF}_\lambda$ if $\max\{\|\mathsf{lcs}_{w_v(w_i, w_j)}\|, \|\mathsf{lcs}_{w_v(w_j, w_i)}\|\} \geq \Omega(\lambda^2)$. We prove that $\mathsf{NF}_\lambda$ has no large independent set. In other words, if we select a large enough set of vertices in $\mathsf{NF}_\lambda$, then there is at least one edges between them. Next, we apply the Turan theorem to prove that $\mathsf{NF}_\lambda$ is dense. Finally, we imply that since $\mathsf{NF}_\lambda$ is dense, there is one vertex $u$ in the neighbors of $v$ such that there are a lot of 2-paths between $v$ and $u$. This implies that the edge $(u, v)$ should have been detected in our sparsification and therefore must not exist in $\mathsf{NG}_\lambda$. This contradiction implies that $\mathsf{NG}_\lambda$ is sparse in the first place.

### 1.3.3 LIS

In this section, we present our result for longest increasing subsequence. More precisely, we show that when the solution size is lower bounded by $n\lambda$ ($\lambda \in [0, 1]$), one can approximate the solution within a factor $\Omega(\lambda^3)$ in time $\widetilde{O}(\sqrt{n}/\lambda^7)$. This married with a simple sampling algorithm for the cases that $\lambda < n^{-\Omega(1)}$, provides an $\Omega(\lambda^3)$-approximate algorithm with running time of $\widetilde{O}(n^{0.85})$ (without further dependence on $\lambda$). We further extend this result to reduce the running time to $\widetilde{O}(n^\kappa \operatorname{poly}(1/\lambda))$ for any $\kappa > 0$ by imposing a multiplicative factor of $\operatorname{poly}(1/\lambda)$[10] to the approximation.

Our algorithm heavily relies on sampling random elements of the array for which longest increasing subsequence is desired. Denote the input sequence by $A = \langle a_1, a_2, \ldots, a_n \rangle$. A naive approach to approximate the solution is to randomly subsample the elements of $A$ to obtain a smaller array $B$ and then compute the longest increasing subsequence of $B$ to estimate the solution size for $A$. Let us first show why this approach alone fails to provide a decent approximation factor. First, consider

---

[10]The exponent of $1/\lambda$ depends exponentially on $1/\kappa$.

an array $A = \langle 1, 2, \ldots, n \rangle$ which is strictly increasing. Based on $A$, we construct two inputs $A'$ and $A''$ in the following way:

- $A'$ is exactly equal to $A$ except that a $p$ fraction of the elements in $A'$ are replaced by 0.

- $A''$ is exactly equal to $A$ except that every block of length $\sqrt{n}$ is reversed in $A''$. In other words, $A'' = \langle \sqrt{n}, \sqrt{n}-1, \sqrt{n}-2, \ldots, 1, 2\sqrt{n}, 2\sqrt{n}-1, \ldots, \sqrt{n}+1, \ldots, n, n-1, n-2, \ldots, n-\sqrt{n}+1 \rangle$.

We subsample the two arrays $A'$ and $A''$ with a rate of $1/\sqrt{n}$ to obtain two smaller arrays $B'$ and $B''$ of size roughly $O(\sqrt{n})$. It is easy to prove that $\mathsf{lis}(B') = \Omega(\sqrt{n})$ and $\mathsf{lis}(B'') = \Omega(\sqrt{n})$, yet $\mathsf{lis}(A') = \Omega(n)$ but $\mathsf{lis}(A'') = O(\sqrt{n})$. By setting $p = 1/e$[11] we can also make sure that $\mathsf{lis}(B')$ and $\mathsf{lis}(B'')$ are within a small multiplicative range even though the gap between $\mathsf{lis}(A')$ and $\mathsf{lis}(A'')$ is substantial.

The above observation shows that the problem is very elusive when random sampling is involved. We bring a remedy to this issue in the following. Divide the input array into $\sqrt{n}$ subarrays of size $\sqrt{n}$. We denote the subarrays by $\mathsf{sa}_1, \mathsf{sa}_2, \ldots, \mathsf{sa}_{\sqrt{n}}$ and fix an optimal solution $\mathsf{opt}$ for the longest increasing subsequence of $A$. Define $\mathsf{sm}(\mathsf{sa}_i)$ to be the smallest number in $\mathsf{sa}_i$ that contributes to $\mathsf{opt}$ and $\mathsf{lg}(\mathsf{sa}_i)$ to be the largest number in $\mathsf{sa}_i$ that contributes to $\mathsf{opt}$. Moreover, define $\mathsf{lis}^{[\ell,r]}$ to be the longest increasing subsequence of an array subject to the elements whose values lie within the interval $[\ell, r]$. This immediately implies

$$\mathsf{lis}(A) = \sum_{i=1}^{\sqrt{n}} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i).$$

Another observation that we make here is that since we assume $\|\mathsf{lis}(A)\| \geq \lambda$ and the size of each subarray is bounded by $\sqrt{n}$, then we have

$$\frac{\mathsf{lis}(A)}{\max_i \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i)} \geq \sqrt{n}\lambda$$

which means that in order to approximate $\mathsf{lis}(A)$ it suffices to compute $\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i)$ for $\widetilde{O}(1/\lambda)$ many randomly sampled subarrays. This is quite helpful since this shows that we only need to sample $\widetilde{O}(1/\lambda)$ many subarrays and solve the problem for them. However, we do not know the values of $\mathsf{sm}(\mathsf{sa}_i)$ and $\mathsf{lg}(\mathsf{sa}_i)$ in advance. Therefore, the main challenge is to predict the values of $\mathsf{sm}(\mathsf{sa}_i)$ and $\mathsf{lg}(\mathsf{sa}_i)$ before we sample the subarrays.

Indeed, one needs to read the entire array to correctly compute $\mathsf{sm}(\mathsf{sa}_i)$ and $\mathsf{lg}(\mathsf{sa}_i)$ for each of the subarrays. However, we devise a method to approximately guess these values without losing too much in the size of the solution. Roughly speaking, we show that if we sample $k = O(1/(\lambda\epsilon))$ different elements from a subarray $\mathsf{sa}_i$ for some constant $\epsilon$ and denote them by $a_{j_1}, a_{j_2}, \ldots, a_{j_k}$, then for at least one pair $(\alpha, \beta)$, $[a_{j_\alpha}, a_{j_\beta}]$ is approximately close to $[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]$ up to a $(1-\epsilon)$ factor.

The above argument provides $O((1/(\lambda\epsilon))^2)$ candidate domain intervals for each $\mathsf{sa}_i$. This does not provide a solution since we do not know which candidate domain interval approximates $[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]$ for each $\mathsf{sa}_i$. Of course, if we were to randomly choose one candidate interval for every subarray, we would make a correct guess for at least $O(\sqrt{n}(\lambda\epsilon)^2)$ subarrays which provides an approximation guarantee of $\Omega(\lambda^2)$ for our algorithm. However, our assignments have to be monotone too. More precisely, let $[\widetilde{\mathsf{sm}}(\mathsf{sa}_i), \widetilde{\mathsf{lg}}(\mathsf{sa}_i)]$ be the guesses that our algorithm makes, then we should have
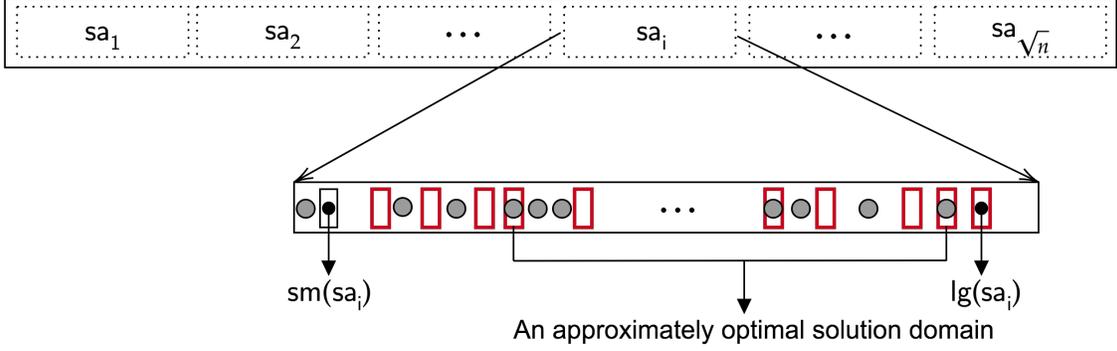
---

[11] $e \simeq 2.7182$.

Figure 3: Red rectangles show the elements of $\mathsf{sa}_i$ that contribute to $\mathsf{lis}(A)$ and gray circles show the elements of $\mathsf{sa}$ that are sampled via our algorithm.

$$\widetilde{\mathsf{sm}}(\mathsf{sa}_1) \leq \widetilde{\mathsf{lg}}(\mathsf{sa}_1) \leq \widetilde{\mathsf{sm}}(\mathsf{sa}_2) \leq \widetilde{\mathsf{lg}}(\mathsf{sa}_2) \leq \ldots \leq \widetilde{\mathsf{sm}}(\mathsf{sa}_{\sqrt{n}}) \leq \widetilde{\mathsf{lg}}(\mathsf{sa}_{\sqrt{n}}).$$

Random sampling does not guarantee that the sampled intervals are monotone. To address this issue, we introduce the notion of *pseudo-solutions*. A pseudo-solution is an assignment of monotone intervals to subarrays in order to approximate $\mathsf{sm}(\mathsf{sa}_i)$ and $\mathsf{lg}(\mathsf{sa}_i)$. The *quality* of a pseudo-solution with intervals $[\ell_1, r_1], [\ell_2, r_2], \ldots, [\ell_{\sqrt{n}}, r_{\sqrt{n}}]$ is equal to $\sum_i \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i)$. For a fixed pseudo-solution, this can be easily approximated via random sampling. Thus, our goal is to construct a pseudo-solution whose quality is at least an $\Omega(\lambda^3)$ approximation of the size of the optimal solution. To this end, we present a greedy method in Section 6.2 to construct the desired pseudo-solution.

Finally, in Section 6.4, we show how the above ideas can be generalized to improve the running time down to $\widetilde{O}(n^\kappa \operatorname{poly}(1/\lambda))$ for any arbitrarily small $\kappa > 0$ by imposing a factor $\operatorname{poly}(1/\lambda)$[12] to the approximation guarantee.

---

[12] The exponent of $1/\lambda$ exponentially depends on $1/\kappa$.

# 2 Organization of the Paper

Our algorithm for LCS is explained in Section 4 (Step 0), Section 3 (Step 1), and Section 5 (Step 2).

In both our results for LCS and LIS, we assume that the goal is to find approximate solutions, provided that the solution size is at least $\lambda_0 n$. After the algorithms terminate, if the output is smaller than what we expect, we realize that the solution is smaller than $\lambda_0 n$. Therefore, we begin by setting $\lambda_0 = 1$ and iteratively multiply $\lambda_0$ by a $1 - \epsilon$ factor until we obtain a solution. This only adds a multiplicative factor of $\log 1/\lambda$ to the running time and a multiplicative factor of $1 - \epsilon$ to the approximation. Since we present two different sparsification techniques, we obtain two theorems: one is Theorem 2.1 and the other is Theorem 2.2.

**Theorem 2.1.** *Given strings $A, B$ of length $|A| = |B| = n$ with $\|\mathsf{lcs}(A, B)\| = \lambda$, we can approximate the length of the LCS between the two strings within a factor $\Omega(\lambda^3)$ in time $\widetilde{O}(n^{39/20})$.*

*Proof.* Fix a sufficiently small constant $\epsilon$ (e.g. $\epsilon = 1/10000$). Since we do not know the value of $\lambda$, we start with $\lambda = 1$ and iteratively try to solve the problem within a factor $\Omega(\lambda^3)$. Each time we are not able to find a solution, we multiply the value of $\lambda$ by $1 - \epsilon$ and proceed. This imposes a $1 - \epsilon$ factor to the approximation and a logarithmic factor to the runtime that can be hidden in the $\widetilde{O}$ and $\Omega$ notations. Thus, in what follows, we assume that we fix a $\lambda$ and we know that the solution size is $\lambda n$. Also, we define $\kappa = n^{-1/140}$.

If $\lambda \leq \kappa$ we run the following algorithm: we choose $n\lambda^3$ characters of $A$ uniformly at random to obtain a string $A'$. With high probability, the LCS of $A'$ and $B$ is at least $(1 - \epsilon)\lambda^4 n$. Therefore, we set our aim to find such a solution. To this end, we spend a preprocessing time of $O(n \log n)$ on string $B$ and then we will be able to find such a solution in time $O(|A'|\lambda^4 n \log n)$ (see Theorem A.8). Thus, the overall runtime would be bounded by

$$\widetilde{O}(|A'|\lambda^4 n) \leq \widetilde{O}(n\lambda^7 n) \leq \widetilde{O}(n^{2-7/140}) = \widetilde{O}(n^{39/20}).$$

Moreover, the approximation factor is also $\Omega(\lambda^3)$ as desired.

When $\lambda \geq \kappa$ we run the three steps of our algorithm (Step 0 stated in Fact 4.6, Step 1 stated in Theorem 3.21, and Step 2 stated in Lemma 5.1) by setting $d = \sqrt{n}\lambda$, $\gamma = 2/3$, $\mathsf{w}_{\mathsf{max}} = \sqrt{n}$ and $k = \widetilde{O}(\sqrt{n}/\lambda) \leq \widetilde{O}(n^{71/140})$. After constructing the windows in Step 0 (Fact 4.6), we run the algorithm of Theorem 3.21 (Step 1) for every $\lambda' \in \{\epsilon\lambda, \epsilon(1 + \epsilon)\lambda, \epsilon(1 + \epsilon)^2\lambda, \ldots, 1\}$. If for a pair of windows $w_i, w_j$ our algorithm in Step 1 detects an edge at $\lambda'$ then we update the solution size for such a pair to $\max\{|w_i|, |w_j|\}\lambda'^4/16$. We then run the algorithm of Step 2 (Lemma 5.1) to find a solution. In what follows, we bound the approximation factor and the runtime of the algorithm.

**Approximation factor:** For now, we only consider the multiplicative and additive approximation losses that are incurred in Steps 0, 1, and 2 but we assume that for each $\lambda$, the output of Step 1 is without any errors. We then incorporate those errors to bound the overall approximation factor. These errors are listed below:

- We lose a multiplicative constant factor in Step 0 followed by an additive error of $\epsilon\lambda n$ (see Lemma 4.7).

- We lose a multiplicative factor $1 - 2\epsilon$ in Step 1 due to the fact that we ignore window pairs whose LCS sizes drops below a threshold $\epsilon\lambda$ times the maximum window size.

- We also lose a multiplicative factor $1 - \epsilon$ in Step 3 (Lemma 5.1).

12

Since $\epsilon$ is very small and all the multiplicative loss factors are constant, we can assume that all the above errors amount to an overall $1/c$ for a **constant** factor $c$. Now we are ready to discuss the loss in the approximation incurred in Step 1. If Step 1 did not incur any error, we would find a solution of size $1/c\lambda n$ for the two strings. Suppose that $(w_1, w_1'), (w_2, w_2'), \ldots$ is a such a window-compatible solution that provides a solution of total size $1/c\lambda n$ if the estimations were correct. We put these pairs in two sets based on whether the $A$ side window is larger or the $B$ side window is larger. At least one set gives us a solution of size $1/(2c)\lambda n$. Without loss of generality, we assume that these are the pairs whose $B$ side is at least as large as their $A$ side and in the rest of the proof we restrict ourselves to such pairs.

For each pair $(w_i, w_i')$, define $\lambda_i$ to be the ratio of the actual $\mathsf{LCS}$ of $w_i$ and $w_i'$ divided by $|w_i'|$. Since $\sum |w_i'| \leq n$ and $\sum \lambda_i |w_i'| \geq 1/(2c)\lambda n$ then it follows from Lemma A.7 that

$$\sum |w_i'|\lambda_i^4 \geq \left(\frac{1/(2c)n\lambda}{\sum |w_i'|}\right)^4 \sum |w_i'| = (1/(2c)\lambda)^4 n(n/\sum |w_i'|)^3 \geq (1/(2c)\lambda)^4 n,$$

and therefore $\sum |w_i'|\lambda_i^4/16 \geq \frac{\lambda^4 n}{256c^4} = \Omega(\lambda^4 n)$. Thus, via the estimations we find in Step 1, we would be able to find a solution that loses a factor of at most $\Omega(\lambda^3)$.

**Runtime:** The runtime of Step 0 is $\widetilde{O}(k) = \widetilde{O}(\sqrt{n}/\lambda) = \widetilde{O}(n^{71/140})$ which is negligible. By Theorem 3.21, the runtime of Step 1 is equal to

$$\widetilde{O}(k^{1+\gamma}\mathsf{w}_{\mathsf{max}}^2 + k^{2+\gamma}\mathsf{w}_{\mathsf{max}}) \leq \widetilde{O}(n^{1+71/140(5/3)} + n^{1/2+71/140(8/3)}) \leq \widetilde{O}(n^{1.86}).$$

Moreover, the number of remaining edges in the graph is bounded by (see Theorem 3.21)

$$\widetilde{O}(k^{2-\gamma}/\lambda) \leq \widetilde{O}(k^{2-2/3+1/70}) \simeq \widetilde{O}(k^{1.35}).$$

This implies a runtime of (see Theorem 5.2)

$$\widetilde{O}(k^{2-(2-1.35)/2}\mathsf{w}_{\mathsf{max}}^2/\lambda^6) = \widetilde{O}(k^{1.675}\mathsf{w}_{\mathsf{max}}^2/\lambda^6) \leq \widetilde{O}(n^{1+6/140+71/140(1.675)}) \leq \widetilde{O}(n^{1.9})$$

for Step 2. $\qquad\square$

**Theorem 2.2.** *Given strings $A, B$ of length $|A| = |B| = n$ with $||\mathsf{lcs}(A,B)|| = \lambda$ where $\lambda$ is constant, we can approximate the length of the $\mathsf{LCS}$ between the two strings within a factor $(1-\epsilon)\lambda^2$ in $n^{2-\Omega_{\epsilon\lambda}(1)}$ time for any constant $0 < \epsilon < 1$.*

*Proof.* Let $\epsilon' = \epsilon/200$. Since we do not know the value of $\lambda$, we start with $\lambda = 1$ and iteratively try to solve the problem within a factor $(1-\epsilon')\lambda^2$. Each time we are not able to find a solution, we multiply the value of $\lambda$ by $1-\epsilon'$ and proceed. This imposes a $1-\epsilon'$ factor to the approximation and a constant factor to the runtime. Thus, in what follows, we assume that we fix a $\lambda$ and we know that the solution size is at least $\lambda n$.

We run the three steps of our algorithm (Step 0 stated in Fact 4.3, Step 1 stated in Theorem 3.13, and Step 2 stated in Lemma 5.1) by setting $d = \sqrt{n}$, $\mathsf{w}_{\mathsf{max}} = \sqrt{n}$ and $k = \widetilde{O}(\sqrt{n})$ (Notice that here $\lambda$ is constant). After constructing the windows in Step 0 (Fact 4.3), we run the algorithm of Theorem 3.13 (Step 1) for every $\lambda' \in \{\epsilon\lambda, \epsilon(1+\epsilon')\lambda, \epsilon(1+\epsilon')^2\lambda, \ldots, 1\}$. If for a pair of windows $w_i, w_j$ our algorithm in Step 1 detects an edge at $\lambda'$ then we update the solution size for such a pair to $\sqrt{|w_i||w_j|}(1-\epsilon')\lambda'^3$. We then run the algorithm of Step 2 (Lemma 5.1) to find a solution. In what follows, we bound the approximation factor and the runtime of the algorithm.

**Approximation factor:** For now, we only consider the multiplicative and additive approximation losses that are incurred in Steps 0, 1, and 2 but we assume that for each $\lambda$, the output of Step 1 is without any errors. We then incorporate those errors to bound the overall approximation factor. These errors are listed below:

- We lose an additive error of $8\epsilon'\lambda n$ in Step 0 (see Lemma 4.4).

- We lose a multiplicative factor $1 - \epsilon'$ in Step 1 due to the fact that we ignore window pairs whose normalized LCS sizes drops below a threshold $\epsilon'\lambda$.

- We also lose a multiplicative factor $1 - \epsilon'$ in Step 3 (Lemma 5.1).

Since $\epsilon' = \epsilon/200 < 1/200$, we can assume that all the above errors amount to an overall $1 - 20\epsilon'$ factor. Now we are ready to discuss the loss in the approximation incurred in Step 1. If Step 1 did not incur any error, we would find a solution of size $(1 - 20\epsilon')\lambda n$ for the two strings. Suppose that $(w_1, w_1'), (w_2, w_2'), \ldots$ is a such a window-compatible solution that provides a solution of total size $(1 - 20\epsilon')\lambda n$ if the estimations were correct. For each pair $(w_i, w_i')$, define $\lambda_i$ to be $||\mathsf{lcs}(w_i, w_i')||$. Since $\sum |w_i| \leq n$, $\sum |w_i'| \leq n$ and $\sum \lambda_i \sqrt{|w_i| \cdot |w_i'|} \geq (1 - 20\epsilon')\lambda n$, it follows from Lemma A.7 and Cauchy–Schwarz inequality that

$$\sum \sqrt{|w_i| \cdot |w_i'|}(1 - \epsilon')\lambda_i^3 \geq (1 - \epsilon') \sum \sqrt{|w_i| \cdot |w_i'|} \left( \frac{\sum \lambda_i \sqrt{|w_i| \cdot |w_i'|}}{\sum \sqrt{|w_i| \cdot |w_i'|}} \right)^3$$
$$\geq (1 - \epsilon') \frac{((1 - 20\epsilon')\lambda n)^3}{\left( \sum \sqrt{|w_i| \cdot |w_i'|} \right)^2}$$
$$\geq (1 - \epsilon') \frac{((1 - 20\epsilon')\lambda n)^3}{n^2}$$
$$\geq (1 - \epsilon')(1 - 20\epsilon')^3 \lambda^3 n$$
$$\geq (1 - \epsilon')(1 - 60\epsilon')\lambda^3 n.$$

Thus, via the estimations we find in Step 1, we would be able to find a solution that loses a factor of at most $(1 - \epsilon')(1 - 60\epsilon')\lambda^3 \geq (1 - \epsilon)\lambda^3$.

**Runtime:** The runtime of Step 0 is $\widetilde{O}(k) = \widetilde{O}(\sqrt{n})$ which is negligible. By Theorem 3.13 and setting $\gamma = 1/10$, the runtime of Step 1 is equal to

$$\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{1.1} + \mathsf{w}_{\mathsf{max}} k^{2.2}) = \widetilde{O}(n^{1.6}).$$

Moreover, the number of remaining edges in the graph is bounded by (see Theorem 3.13)

$$\widetilde{O}(k^{2 - \epsilon\lambda^3/(800\mathsf{w}_{\mathsf{layers}}\mathsf{w}_{\mathsf{gap}})}) = \widetilde{O}(k^{2 - \Omega(\epsilon\lambda^6/\log(1/\lambda))}) = \widetilde{O}(k^{2 - \Omega(\epsilon\lambda^7)})$$

This implies a runtime of (see Theorem 5.3)

$$\widetilde{O}(k^{2 - \Omega(\epsilon\lambda^7)}\mathsf{w}_{\mathsf{max}}^2/\lambda^6) = \widetilde{O}(n^{2 - \Omega(\epsilon\lambda^7/2)}) = n^{2 - \Omega_{\lambda, \epsilon}(1)}$$

for Step 2. $\qquad \qquad \square$

As an immediate corollary of Theorem 2.1, we present an algorithm that beats the $1/|\Sigma|$ approximation factor in truly subquadratic time, when the strings are balanced.

**Corollary 2.3.** *Given a pair of strings $(A, B)$ of length $n$ over alphabet $\Sigma$ that satisfy the balance condition, we can approximate their* LCS *within an $O(|\Sigma|^{3/4})$ factor in time $O(n^{39/20})$.*

*Proof.* Since $A$ and $B$ are balanced, there is a character $\sigma \in \Sigma$ that appears at least $n/|\Sigma|$ times in both strings. Indeed, finding a solution of size $n/|\Sigma|$ by restricting our attention to only character $\sigma$ can be done in time $O(n)$. If $\mathsf{lcs}(A, B) \leq n/|\Sigma|^{1/4}$ this already gives us an $O(|\Sigma|^{3/4})$ approximate solution. Otherwise, $||\mathsf{lcs}(A, B)|| > 1/|\Sigma|^{1/4}$ and the approximation factor of our $\Omega(\lambda^3)$-approximation algorithm would be bounded by $O(|\Sigma|^{3/4})$. $\qquad\square$

Finally, we bring our results for $\mathsf{LIS}$ in Section 6. We show that

**Theorem 2.4.** *Given a length-$n$ sequence $A$ with $\mathsf{lis}(A) = n\lambda$. We can approximate the length of the $\mathsf{LIS}$ within a factor of $\Omega(\lambda^3)$ in time $\widetilde{O}(n^{17/20})$.*

*Proof.* If $\lambda < n^{-1/20}$ we sample the array with a rate of $n^{-3/20}$ and compute the $\mathsf{LIS}$ for the sampled array. The running time of the algorithm is $\widetilde{O}(n^{17/20})$. The approximation factor is $O(n^{-3/20}) \geq \Omega(\lambda^3)$. Otherwise, by Theorem 6.8, we estimate the size of $\mathsf{LIS}$ up to an $\Omega(\lambda^3)$ approximation factor in time $\widetilde{\Omega}(\lambda^{-7}\sqrt{n}) \leq \widetilde{O}(n^{17/20})$. $\qquad\square$

# 3   LCS Step 1: Sparsification via Birthday Triangle Inequality

Recall that we are given two sets of windows $W_A$ and $W_B$ for the strings and our goal is to approximate the LCS of all but a few pairs of windows from $W_A \times W_B$. For simplicity, we put all the windows in the same basket $W = W_A \cup W_B$ and denote the windows by $w_1, w_2, \ldots, w_k$ where $k$ is the total number of windows. Since the windows have different lengths, we define $\mathsf{w_{max}} = \max_{i \in [k]} |w_i|$ to be the maximum length of the windows. Similarly, we also define $\mathsf{w_{min}} = \min_{i \in [k]} |w_i|$ to be the minimum length of the windows. Let $\mathsf{w_{gap}} = \mathsf{w_{max}}/\mathsf{w_{min}}$. Let $\mathsf{w_{layers}}$ denote the number of different window sizes. Notations $\mathsf{w_{gap}}$ and $\mathsf{w_{layers}}$ will be used in the later analysis.

In order to approximate the LCS's we fix a $\lambda \in \{\epsilon\lambda_0, (1+\epsilon)\epsilon\lambda_0, (1+\epsilon)^2\epsilon\lambda_0, \ldots, 1\}$ and sparsify graph $\mathsf{G}_\lambda$. In Section 3.1, we present a sparsification algorithm (Algorithm 1) which provides $(1-\epsilon)\lambda^2$-approximation when $\lambda$ is constant. The formal guarantee of the algorithm is provided in Theorem 3.13. In Section 3.2, we present a sparsification which provides $\Omega(\lambda^3)$-approximation for any (potentially sub-constant) $\lambda$.

## 3.1 Sparsification for constant $\lambda$ (Step 1 of Theorem 2.2)

Fix an arbitrary $\mathsf{LCS}$ for every pair of windows and refer to that as $\mathsf{opt}_{i,j}$ for two windows $w_i$ and $w_j$. Note that we do not explicitly compute $\mathsf{opt}_{i,j}$ in our algorithm. Let us for simplicity, think of each $\mathsf{opt}_{i,j}$ as a matching between the characters of the two windows. Also, denote by $\big(\mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j}\big)$ a solution which is constructed for windows $w_i$ and $w_j$ by taking pairs of characters $(x, y)$ such that:

$$(x, y) \in \big(\mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j}\big) \Longleftrightarrow \exists x', y' \text{ such that}$$
$$(x, y') \in \mathsf{opt}_{i,a} \text{ and}$$
$$(y', x') \in \mathsf{opt}_{a,b} \text{ and}$$
$$(x', y) \in \mathsf{opt}_{b,j}.$$

Let $\left\| \big(\mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j}\big) \right\| = \frac{\left| \big(\mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j}\big) \right|}{\sqrt{|w_i||w_j|}}$.

Our analysis is based on a notion which roughly reflects "in how many ways a desirable solution can be made for a pair of windows $(w_i, w_j)$ by taking the intersection of the $\mathsf{LCS}$ for other pairs". Below, we provide a definition for this notion.

**Definition 3.1** (($\epsilon, \lambda$)-constructive). Let $0 < \epsilon, \lambda < 1$ be fixed values. We call a tuple $\langle w_i, w_a, w_b, w_j \rangle$ ($w_i \neq w_a \neq w_b \neq w_j$) an ($\epsilon, \lambda$)-*constructive* tuple, if

$$\left\| \big(\mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j}\big) \right\| \geq (1 - \epsilon)\lambda^3.$$

The advantage of a constructive tuple is that if $\mathsf{opt}_{i,a}, \mathsf{opt}_{a,b}$, and $\mathsf{opt}_{b,j}$ are provided, one can construct a desirable solution for $\mathsf{opt}_{i,j}$ in linear time by taking the intersection of the given matchings.

We parametrize our algorithm by a value $0 < \gamma < 1$ to be set later. One may optimize the runtime of the algorithm by setting the value of $\gamma$ in terms of the number of windows and the length of the windows. We first sample a set $S$ of $O(k^\gamma \log k)$ windows. Next, we compute $\mathsf{opt}_{i,j}$ of every window $w_i \in S$ and every other window $w_j$ (not necessarily in $S$). Finally, we find all tuples $\langle w_i, w_a, w_b, w_j \rangle$ such that $w_a, w_b \in S$ and they satisfy the following property:

$$\left\| \big(\mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j}\big) \right\| \geq (1 - \epsilon)\lambda^3.$$

Recall that we call such tuples ($\epsilon, \lambda$)-*constructive* and update $\widehat{\mathsf{O}}_{\lambda^2}[i][j], \widehat{\mathsf{O}}_{\lambda^2}[j][i] \leftarrow 1$ accordingly. This is shown in Algorithm 1.

The running time of our algorithm is equal to $O(k|S|\mathsf{w}_{\mathsf{max}}^2 + k^2|S|^2\mathsf{w}_{\mathsf{max}})$. The rest of this section is dedicated to proving that what remains in the $\mathsf{lcs}$-graph is sparse; this is formalized in Lemma 3.4.

**Definition 3.2** ($\mathsf{NG}_\lambda$). Define a graph $\mathsf{NG}_\lambda$ with $k$ vertices and the following edges:

$$E(\mathsf{NG}_\lambda) = \left\{ (i, j) \mid \||\mathsf{lcs}(w_i, w_j)\| \geq \lambda \quad \text{and} \quad \widehat{\mathsf{O}}_{\lambda^2}[i][j] = 0 \right\}.$$

In other words, $\mathsf{NG}_\lambda$ contains all of the edges that are not detected in our algorithm. We extend the notion of constructive tuples to the undetected edges in our algorithm. We call such tuples *undetected-constructive*.

**Definition 3.3** (undetected-constructive tuple). We say a tuple $\langle w_i, w_a, w_b, w_j \rangle$ is ($\epsilon, \lambda$)-undetected-constructive if it is ($\epsilon, \lambda$)-constructive and also

$$(i, a), (a, b), (b, j), (i, j) \in E(\mathsf{NG}_\lambda).$$

**Algorithm 1** Sparsification for constant $\lambda$ (Step 1 of Theorem 2.2)

---

1: **procedure** QUADRATICSPARSIFICATION$(w_1, w_2, \ldots, w_k, \lambda, \epsilon)$          ▷ Theorem 3.13
2:      $S \leftarrow 40k^\gamma \log k$ i.i.d. samples of $[k]$
3:      $\widehat{\mathsf{O}}_{\lambda^2} \leftarrow \{0\}^{k \times k}$
4:      **for** $w_i \in S$ **do**                            ▷ Takes $k|S|\mathsf{w}_{\mathsf{max}}^2$ time
5:          **for** $j \leftarrow 1$ to $k$ **do**
6:              $\mathsf{opt}_{i,j}, \mathsf{opt}_{j,i} \leftarrow \mathsf{lcs}(w_i, w_j)$
7:          **end for**
8:      **end for**
9:      **for** $w_i \in S$ **do**                            ▷ Takes $k|S|\mathsf{w}_{\mathsf{max}}$ time
10:          **for** $j \leftarrow 1$ to $k$ **do**
11:              **if** $\|\mathsf{opt}_{i,j}\| \geq \lambda$ **then**
12:                  $\widehat{\mathsf{O}}_{\lambda^2}[i][j] \leftarrow 1$
13:                  $\widehat{\mathsf{O}}_{\lambda^2}[j][i] \leftarrow 1$
14:              **end if**
15:          **end for**
16:      **end for**
17:      **for** $i \leftarrow 1$ to $k$ **do**                      ▷ Takes $k^2|S|^2\mathsf{w}_{\mathsf{max}}$ time
18:          **for** $j \leftarrow 1$ to $k$ **do**
19:              **for** $w_a \in S$ **do**
20:                  **for** $w_b \in S$ **do**
21:                      **if** $\|(\mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j})\| \geq (1-\epsilon)\lambda^3$ **then**
22:                          $\widehat{\mathsf{O}}_{\lambda^2}[i][j] \leftarrow 1$
23:                          $\widehat{\mathsf{O}}_{\lambda^2}[j][i] \leftarrow 1$
24:                      **end if**
25:                  **end for**
26:              **end for**
27:          **end for**
28:      **end for**
29:      **return** $\widehat{\mathsf{O}}_{\lambda^2}$
30: **end procedure**

---

**Lemma 3.4.** *With probability at least* $1 - 2/k^3$, $\mathsf{NG}_\lambda$ *has at most* $k^{2-\gamma\epsilon\lambda^3/(80\mathsf{w}_{\mathsf{layers}}\mathsf{w}_{\mathsf{gap}})}$ *edges.*

Instead of arguing directly about $\mathsf{NG}_\lambda$, our proof uses another graph $\mathsf{NF}_\lambda$, which is sub-sampled from $\mathsf{NG}_\lambda$:

**Definition 3.5** ($\mathsf{NF}_\lambda$)**.** Let $\mathsf{NF}_\lambda$ be the graph constructed in the following way: for each node $v \in \mathsf{NG}_\lambda$ we keep $v$ in $\mathsf{NF}_\lambda$ with probability

$$p := k^{-1+\gamma/4}/4.$$

For each $u, v \in \mathsf{NF}_\lambda$, if $(u, v) \in \mathsf{NG}_\lambda$ then we also draw an edge $u, v$ in $\mathsf{NF}_\lambda$. We identify between the vertices and edges of $\mathsf{NF}_\lambda$ and the corresponding vertices and edges of $\mathsf{NG}_\lambda$, i.e. $V(\mathsf{NF}_\lambda) \subseteq V(\mathsf{NG}_\lambda)$ and $E(\mathsf{NF}_\lambda) \subseteq E(\mathsf{NG}_\lambda)$.

The proof has two main ingredients. The first part of the proof uses the details of our algorithm to show that, w.h.p. over the algorithm's randomness, $\mathsf{NG}_\lambda$ contains few undetected-constructive

tuples (Lemma 3.7). Whenever this is the case, w.h.p. over the sub-sampling procedure, $\mathsf{NF}_\lambda$ does not contain any undetected-constructive tuples (Claim 3.8).

The second part of the proof assumes by contradiction that $\mathsf{NG}_\lambda$ is dense. It uses this assumption to conclude that $\mathsf{NF}_\lambda$ is also dense (w.h.p. over the sub-sampling procedure, see Claim 3.9), and therefore by Turan's Theorem contains a large bi-clique (Claim 3.11). Finally, we use the large bi-clique to show that $\mathsf{NF}_\lambda$ does contain an undetected-constructive tuple (Lemma 3.12) - a contradiction!

### 3.1.1 Proof of Lemma 3.4, part 1: $\mathsf{NF}_\lambda$ does not contain an undetected-constructive tuple

We first introduce the notion of "well-connected pairs".

**Definition 3.6** (well-connected pair). We say that a pair of windows $(w_i, w_j)$ is *well-connected* if there are at least $k^{2-\gamma}$ pairs of windows $(w_a, w_b)$ such that $\langle w_i, w_a, w_b, w_j \rangle$ is $(\epsilon, \lambda)$-constructive.

We argue that well-connected pairs are detected in our algorithm with high probability.

**Lemma 3.7.** *Let $\widehat{\mathsf{O}}_{\lambda^2} \in \{0,1\}^{k \times k}$ denote the output of Algorithm 1. With probability at least $1 - 2/k^3$, for all $(i,j) \in [k] \times [k]$ such that $(w_i, w_j)$ is a well-connected pair (Definition 3.6), we have $\widehat{\mathsf{O}}_{\lambda^2}[i][j] = 1$.*

*Proof.* We consider a fixed $(i,j)$ such that pair $(w_i, w_j)$ is well-connected. Let

$$Q_{i,j} = \left\{ (a,b) \ \middle| \ \left\| \mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j} \right\| \geq (1-\epsilon)\lambda^3 \right\}.$$

Conceptually, we divide the process of sampling $S$ into two phases: we sample $20k^\gamma \log k$ windows in the first phase, and then we sample $20k^\gamma \log k$ more windows in the second phase.

For each $a \in [k]$, let $Q_{i,j,a} = \{b : (a,b) \in Q_{i,j}\}$. Since $\sum_{a \in [k]} |Q_{i,j,a}| = |Q_{i,j}| \geq k^{2-\gamma}$, there are at least $\frac{k^{1-\gamma}}{2}$ different number $a$'s in $[k]$ such that $|Q_{i,j,a}| \geq \frac{k^{1-\gamma}}{2}$. Hence, in the first phase, there is a sampled number $q$ such that $|Q_{i,j,q}| \geq k^{1-\gamma}/2$ with probability at least

$$1 - \left(1 - \frac{k^{1-\gamma}/2}{k}\right)^{20k^\gamma \log k} > \frac{k^5 - 1}{k^5}.$$

We fix such a $q$. In the second phase, there is a sampled number $r$ such that $r \in Q_{i,j,q}$ with probability at least

$$1 - \left(1 - \frac{k^{1-\gamma}/2}{k}\right)^{20k^\gamma \log k} > \frac{k^5 - 1}{k^5}.$$

Since the number of $(i,j)$ pairs is at most $k^2$, the lemma is obtained by a union bound on all the well-connected pairs $(w_i, w_j)$. □

We complete the first part of the proof by showing that with probability at least 0.99, $\mathsf{NF}_\lambda$ does not have an undetected-constructive tuple.

**Claim 3.8** ($\mathsf{NF}_\lambda$ has no undetected-constructive tuple). *With probability at least 0.99, there is no undetected-constructive tuple in $\mathsf{NF}_\lambda$.*

19

*Proof.* By Lemma 3.7, we assume that no edge $(i, j)$ in $\mathsf{NG}_\lambda$ is well-connected. In other words, for all pairs of windows $(w_i, w_j)$ that are connected in $\mathsf{NG}_\lambda$, there are at most $k^{2-\gamma}$ pairs of windows $(w_a, w_b)$ such that $\langle w_i, w_a, w_b, w_j \rangle$ is $(\epsilon, \lambda)$-constructive.

Recall that for a tuple $\langle w_i, w_a, w_b, w_j \rangle$ to be undetected-constructive, edge $(i, j)$ should belong to $\mathsf{NG}_\lambda$. For a fixed $(i, a, b, j)$, the probability that we keep it in $\mathsf{NF}_\lambda$ is $p^4$. Therefore, the expected number of undetected-constructive tuples in $\mathsf{NF}_\lambda$ is bounded by

$$\mathbb{E}[\#\text{undetected-constructive tuple}] \le k^2 k^{2-\gamma} \cdot p^4 = 1/256$$

where the last step follows from $p = k^{-1+\gamma/4}/4$.

By Markov's inequality, we have

$$\Pr[\#\text{undetected-constructive tuple} \ge 1/2] \le \frac{\mathbb{E}[\#\text{undetected-constructive tuple}]}{1/2} < 1/100.$$

Thus, with probability 0.99, there is no undetected-constructive tuple in $\mathsf{NF}_\lambda$. $\qquad\square$

### 3.1.2 Proof of Lemma 3.4, part 2: $\mathsf{NF}_\lambda$ contains an undetected-constructive tuple

Now we assume by contradiction that the $\mathsf{NG}_\lambda$ is dense. To simplify the notation, we introduce another parameter,

$$\beta := \gamma\epsilon\lambda^3/(80\mathsf{w}_{\mathsf{layers}}\mathsf{w}_{\mathsf{gap}}).$$

Note that $\beta < \gamma - \Omega(1)$.

We are able to show that if $\mathsf{NG}_\lambda$ is dense, so is $\mathsf{NF}_\lambda$.

**Claim 3.9** ($\mathsf{NF}_\lambda$ is a dense graph)**.** *If for some $0 < \beta < \gamma/4 - \Omega(1)$, $\mathsf{NG}_\lambda$ contains is at least $k^{2-\beta}$ edges then with probability at least $0.98$ we have*

$$|E(\mathsf{NF}_\lambda)| \ge \frac{|V(\mathsf{NF}_\lambda)|^{2-4\beta/\gamma}}{128}.$$

The proof of Claim 3.9 uses the following elementary graph theory fact:

**Fact 3.10.** *In every graph $G = (V, E)$, we have*

$$\sum_v \deg(v)^2 \le 2|V||E|.$$

*Proof.*

$$\sum_v \deg(v)^2 \le |V| \sum_v \deg(v) \le 2|V||E|.$$

$\qquad\square$

*Proof of Claim 3.9.* Based on the sampling rate, we know that the following holds in expectation:

$$\mathbb{E}[|V(\mathsf{NF}_\lambda)|] = p|V(\mathsf{NG}_\lambda)| = k^{\gamma/4}/4, \tag{1}$$

and

$$\mathbb{E}[|E(\mathsf{NF}_\lambda)|] \ge p^2|E(\mathsf{NG}_\lambda)| \ge k^{\gamma/2-\beta}/16.$$

Using standard Chernoff bound, we have with probability 0.99,

$$|V(\mathsf{NF}_\lambda)| \leq 2\mathbb{E}[|V(\mathsf{NF}_\lambda)|] = k^{\gamma/4}/2.$$

In the rest of this proof, we show that with probability 0.99,

$$|E(\mathsf{NF}_\lambda)| \geq k^{\gamma/2-\beta}/128,$$

and then the claim follows.

It is known that any graph can be made bipartite by removing at most half of its edges [Wes01]. Thus, for the sake of this proof, we only consider a bipartite subgraph of $\mathsf{NG}_\lambda$ with parts $\mathcal{P}$ and $\mathcal{Q}$ that contains at least $k^{2-\beta}/2$ edges. We refer to this graph by $\widehat{\mathsf{NG}_\lambda}$. We consider a two-step construction of $\widehat{\mathsf{NF}_\lambda}$ based on sampling in $\widehat{\mathsf{NG}_\lambda}$, where we first only sub-sample the vertices on the $\mathcal{P}$-side, and then also sub-sample the vertices on the $\mathcal{Q}$-side. Let $\widehat{\mathsf{NF}_\lambda}(\mathcal{P})$ denote the graph after subsampling only the $\mathcal{P}$-side, and $V(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}), \mathcal{P})$ denote the set of vertices in $\widehat{\mathsf{NF}_\lambda}$ on the $\mathcal{P}$-side.

The number of edges in $\widehat{\mathsf{NF}_\lambda}(\mathcal{P})$ is the sum of degrees of surviving $\mathcal{P}$-vertices, aka a sum of i.i.d. random variables bounded in $[0, k]$. It's expectation is given by:

$$\mathbb{E}[|E(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}))|] = p|E(\widehat{\mathsf{NG}_\lambda}(\mathcal{P}))| \geq k^{2-\beta}/2,$$

and by Hoeffding's inequality this sum concentrates around its expectation with high probability:

$$\Pr\left[|E(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}))| \leq pk^{2-\beta}/4\right] \leq \Pr\left[|E(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}))| \leq \mathbb{E}[|E(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}))|] - pk^{2-\beta}/4\right]$$
$$\leq \exp\left(-\frac{2|\mathcal{P}|^2 p^2 k^{4-2\beta}/16}{|\mathcal{P}|k^2}\right)$$
$$\leq \exp\left(-\Theta\left(k^{\gamma/2-2\beta}\right)\right). \hspace{2cm} \text{By } |P| \geq 1$$

By Chernoff bound,

$$\Pr[|V(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}), \mathcal{P})| > 2pk] \leq [|V(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}), \mathcal{P})| > 2p|\mathcal{P}|] \leq \exp\left(-p|\mathcal{P}|/3\right) \leq \exp(-\Theta(k^{\gamma/4-\beta})),$$

where the last inequality is obtained by $|\mathcal{P}| \leq k$. We henceforth fix any realization of $\widehat{\mathsf{NF}_\lambda}(\mathcal{P})$ conditioned on

$$|E(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}))| \geq pk^{2-\beta}/4 \text{ and } |V(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}), \mathcal{P})| \leq 2pk.$$

We now consider the second step of the sub-sampling, which transforms $\widehat{\mathsf{NF}_\lambda}(\mathcal{P})$ to $\widehat{\mathsf{NF}_\lambda}$. The expected number of edges in $\widehat{\mathsf{NF}_\lambda}$ satisfies:

$$\mathbb{E}[|E(\widehat{\mathsf{NF}_\lambda})|] = p|E(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}))| \geq p^2 k^{2-\beta}/4 = k^{\gamma/2-\beta}/64.$$

We now argue about concentration. Notice that the degree of each vertex in $\widehat{\mathsf{NF}_\lambda}(\mathcal{Q})$ is at most the number of vertices on $\mathcal{P}$-side in $\widehat{\mathsf{NF}_\lambda}(\mathcal{P})$ is at most $|V(\widehat{\mathsf{NF}_\lambda}(\mathcal{P}), \mathcal{P})| \leq 2pk$. Therefore, by Hoeffding's inequality,

$$\Pr\left[|E(\widehat{\mathsf{NF}_\lambda})| < k^{\gamma/2-\beta}/128\right] \leq \Pr\left[|E(\widehat{\mathsf{NF}_\lambda})| < \mathbb{E}[|E(\widehat{\mathsf{NF}_\lambda})|] - k^{\gamma/2-\beta}/128\right]$$
$$\leq \exp\left(-\Theta\left(\frac{|\mathcal{Q}|^2 k^{\gamma-2\beta}}{|\mathcal{Q}|p^2 k^2}\right)\right)$$
$$\leq \exp\left(-\Theta\left(k^{\gamma/2-2\beta}\right)\right) \hspace{2cm} \text{By } |Q| \geq 1.$$

By taking a union bound on the low probability events and the fact that $E(\widehat{\mathsf{NF}_\lambda})$ is a subset of $E(\mathsf{NF}_\lambda)$, we have $|E(\mathsf{NF}_\lambda)| \geq k^{\gamma/2-\beta}/128$. $\hspace{1cm} \square$

We now use the fact that $\mathsf{NF}_\lambda$ is dense to argue that it contains a large bi-clique.

**Claim 3.11** ($\mathsf{NF}_\lambda$ contains a large bi-clique)**.** *With probability at least* $0.98$ $\mathsf{NF}_\lambda$*, contains a complete bipartite subgraph* $K_{\gamma/(5\beta),\gamma/(5\beta)}$*.*

*Proof.* Using Turán's Theorem (Lemma A.4), we know, for any integer $s \geq 2$, that a graph $G$ with $n$ vertices and $n^{2-1/s}$ edges has at least one $K_{s,s}$ subgraph. We apply this to graph $\mathsf{NF}_\lambda$. Since $\lambda$ is constant then $\beta/\gamma$ becomes constant and thus $|V(\mathsf{NF}_\lambda)|^{\beta/\gamma}$ is super constant. It follows from Claim 3.9 that

$$|E(\mathsf{NF}_\lambda)| \geq \frac{|V(\mathsf{NF}_\lambda)|^{2-4\beta/\gamma}}{128} > |V(\mathsf{NF}_\lambda)|^{2-5\beta/\gamma}$$

holds with probability at least $0.98$. This in turn implies that with probability at least $0.98$ $\mathsf{NF}_\lambda$ contains a complete bipartite subgraph $K_{\gamma/(5\beta),\gamma/(5\beta)}$. $\qquad\square$

The last step in the proof is to use the large bi-clique to exhibit an undetected-constructive tuple.

**Lemma 3.12.** *Let $X$ and $Y$ be two sets of windows such that for every $w_i \in X$ and $w_j \in Y$ there is an edge $(i,j)$ in $E(\mathsf{NG}_\lambda)$. If $|X| \geq 16\mathsf{w}_{\mathsf{layers}}\mathsf{w}_{\mathsf{gap}}/(\epsilon\lambda^3)$ and $|Y| \geq 16\mathsf{w}_{\mathsf{layers}}\mathsf{w}_{\mathsf{gap}}/(\epsilon\lambda^3)$, then there exist $w_i, w_a, w_b, w_j \in X \cup Y$ such that $\langle w_i, w_a, w_b, w_j \rangle$ is $(\epsilon, \lambda)$-constructive, i.e.,*

$$\left\| \mathsf{opt}_{i,a} \cap \mathsf{opt}_{a,b} \cap \mathsf{opt}_{b,j} \right\| \geq (1-\epsilon)\lambda^3$$

*and either $w_i, w_b \in X$ and $w_a, w_j \in Y$ or $w_i, w_b \in Y$ and $w_a, w_j \in X$.*

*Proof.* Let $\epsilon' = \epsilon/4$. By assumption, we know that $|X|, |Y| \geq 4\mathsf{w}_{\mathsf{layers}}\mathsf{w}_{\mathsf{gap}}/(\epsilon'\lambda^3)$. Moreover the total number of different window sizes is bounded by $\mathsf{w}_{\mathsf{layers}}$. Thus, there exist two sets $\widehat{X} \subseteq X$ and $\widehat{Y} \subseteq Y$ such that $|\widehat{X}|, |\widehat{Y}| \geq 4\mathsf{w}_{\mathsf{gap}}/(\epsilon'\lambda^3)$ and the windows within $\widehat{X}$ are of the same size and the windows within $\widehat{Y}$ are of the same size (though the windows of $\widehat{X}$ and $\widehat{Y}$ may have different sizes). Let $d_x$ denote the window size for each window in $\widehat{X}$, and $d_y$ denote the window size for each window in $\widehat{Y}$.

We select $X' \subseteq \widehat{X}$ and $Y' \subseteq \widehat{Y}$ such that

1. $|X'| \geq 4/(\epsilon'\lambda^3)$.

2. $|Y'| \geq 4/(\epsilon'\lambda^3)$.

3. $|X'|d_x \leq (1+\epsilon')|Y'|d_y$.

4. $|Y'|d_y \leq (1+\epsilon')|X'|d_x$.

To do this, if $|\widehat{X}|d_x > (1+\epsilon')|\widehat{Y}|d_y$, then we set $Y' = \widehat{Y}$ and select $X'$ as an arbitrary subset of $\widehat{X}$ with size $\left\lceil \frac{|Y'|d_y}{d_x} \right\rceil$. Otherwise, we set $X' = \widehat{X}$ and select $Y'$ an arbitrary subset of $\widehat{Y}$ with size $\left\lceil \frac{|X'|d_x}{d_y} \right\rceil$.

We define a character-based (bipartite) graph $G_C = (V_C, E_C)$ as follows: Each window of $X'$ has $d_x$ nodes in the character-based graph such that each node represents a character of the window. Similarly, each window of $Y'$ has $d_y$ nodes in the character-based graph. Two nodes $x, y$ in the character-based graph are adjacent if and only if $(x,y) \in \mathsf{opt}_{i,j}$ where $w_i$ is the window containing character $x$ and $w_j$ is the window containing character $y$. Let $l_x = |X'|$ and $l_y = |Y'|$.
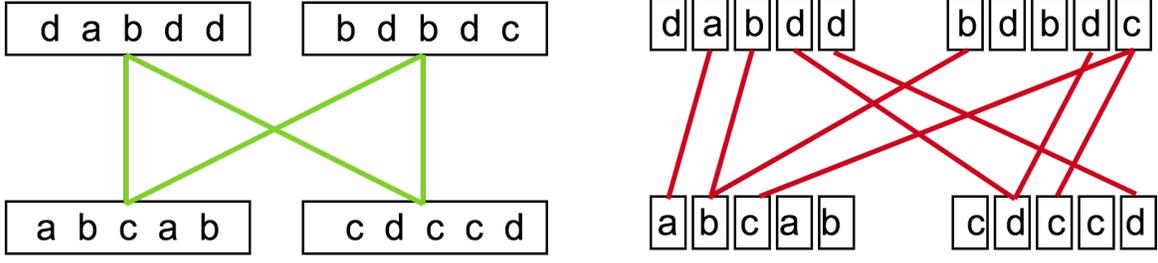
Figure 4: The graph on the left is an example of the string-based graph and the graph on the right is an example of the character-based graph.

The total number of nodes in the character-based graph is $|V_C| = l_x d_x + l_y d_y$. By the last two desiderata of definition of $X', Y'$ that, each side of the character based-graph has approximately the same number of nodes:

$$l_x d_x, l_y d_y \approx_{(1 \pm \epsilon')\text{-factor}} |V_C|/2. \tag{2}$$

Similarly, we also have that:

$$\max\{l_x, l_y\} \min\{d_x, d_y\} \approx_{(1 \pm \epsilon')\text{-factor}} |V_C|/2. \tag{3}$$

The total number of edges in the character-based graph satisfies

$$|E_C| \leq l_x l_y \min\{d_x, d_y\}, \tag{4}$$

and, by (2),

$$|E_C| \geq l_x l_y \lambda \sqrt{d_x d_y} \approx_{\sqrt{1 \pm \epsilon'}\text{-factor}} \sqrt{l_x l_y} \lambda |V_C|/2. \tag{5}$$

By Blakley-Roy inequality (Lemma A.3), the number of walks of length 3 in the character-based graph is at least

$$\#3\text{-walks} \geq |V_C| \cdot \left( 2 \frac{|E_C|}{|V_C|} \right)^3$$
$$\approx_{(1 \pm \epsilon')^{1.5}\text{-factor}} \lambda^3 (l_x l_y)^{1.5} |V_C|.$$

We are interested in the number of 3-walks that are not degenerate. Thus, we need to exclude such degenerate walks from the total count. The number of such walks is upper bounded by

$$\#\text{degenerate 3-walks} \leq (\text{max degree of } G_C) \cdot 4|E_C|$$
$$\leq \max\{l_x, l_y\} \cdot 4(l_x l_y \min\{d_x, d_y\}) \qquad (\text{Eq. (4)})$$
$$\approx_{(1 \pm \epsilon')\text{-factor}} 2 l_x l_y \cdot |V_C| \qquad (\text{Eq. (3)})$$

Therefore the ratio of the number of degenerate 3-walks and the total number of 3-walks is bounded

23

by

$$\frac{\#\text{degenerate 3-walks}}{\#\text{3-walks}} \approx_{(1 \pm O(\epsilon'))\text{-factor}} \frac{2 l_x l_y \cdot |V_C|}{\lambda^3 (l_x l_y)^{1.5} |V_C|}$$

$$= \frac{2}{\lambda^3 \sqrt{l_x l_y}}$$

$$\leq \frac{2}{4/\epsilon'} \qquad \text{(Def. of } X', Y')$$

$$= \epsilon'/2.$$

Hence, the total number of 3-paths (3-walks that are not degenerate) is at least

$$\#\text{3-paths} \gtrsim_{(1 \pm \epsilon')^{2.5}\text{-factor}} \lambda^3 (l_x l_y)^{1.5} |V_C|.$$

On the other hand, the number of 4-tuples of windows that contribute to those 3-paths of characters is $8\binom{l_x}{2}\binom{l_y}{2} \leq 2 l_x^2 l_y^2$. Thus, there must exist a 4-tuple containing at least

$$\frac{\#\text{3-paths}}{\#\text{4-tuples}} \gtrsim_{(1 \pm \epsilon')^{2.5}\text{-factor}} \frac{\lambda^3 (l_x l_y)^{1.5} |V_C|}{2 l_x^2 l_y^2}$$

$$= \frac{\lambda^3 |V_C|}{\sqrt{l_x l_y}}$$

$$\approx_{\sqrt{1 \pm \epsilon'}\text{-factor}} \frac{\lambda^3 \sqrt{l_x d_x l_y d_y}}{\sqrt{l_x l_y}} \qquad \text{(Eq. (2))}$$

$$= \lambda^3 \sqrt{d_x d_y}$$

many 3-walks. This means that such a 4-tuple is $(\epsilon, \lambda)$-constructive.

$\square$

### 3.1.3 Guarantees of the sparsification algorithm

**Theorem 3.13** (Sparsification for constant $\lambda$ (Step 1 of Theorem 2.2)). *Given $k$ windows $w_1, \cdots, w_k$. Let $\mathsf{w}_{\max} = \max_{i \in [k]} |w_i|$, $\mathsf{w}_{\min} = \min_{i \in [k]} |w_i|$ and $\mathsf{w}_{\mathsf{gap}} = \mathsf{w}_{\max}/\mathsf{w}_{\min}$. Let the number of different window sizes be $\mathsf{w}_{\mathsf{layers}}$. For any constant $\lambda \in (0,1)$ and $0 < \epsilon < 1$, there is a randomized algorithm (Algorithm 1) that runs in time*

$$O(\mathsf{w}_{\max}^2 k^{1.1} \log k + \mathsf{w}_{\max} k^{2.2} \log^2 k),$$

*outputs a table $\widehat{\mathsf{O}}_{\lambda^2} \in \{0,1\}^{k \times k}$ such that*

$$||\mathsf{lcs}(w_i, w_j)|| \geq (1 - \epsilon)\lambda^3, \text{ if } \widehat{\mathsf{O}}_{\lambda^2}[i][j] = 1$$

*and*

$$\left| \left\{ (i,j) \,\middle|\, ||\mathsf{lcs}(w_i, w_j)|| \geq \lambda, \quad \text{and} \quad \widehat{\mathsf{O}}_{\lambda^2}[i][j] = 0 \right\} \right| \leq k^{2 - \epsilon \lambda^3 / (800 \mathsf{w}_{\mathsf{layers}} \mathsf{w}_{\mathsf{gap}})}.$$

*The algorithm has success probability at least $1 - 1/\operatorname{poly}(k)$.*

*Proof.* We set $\gamma = 1/10$. The overall running time is

$$O(k|S|w_{\mathsf{max}}^2 + k^2|S|^2 w_{\mathsf{max}})$$
$$= O(k \cdot k^\gamma \log k \cdot w_{\mathsf{max}}^2 + k^2 \cdot k^{2\gamma} \log^2 k \cdot w_{\mathsf{max}})$$
$$= O(k^{1.1} \log k \cdot w_{\mathsf{max}}^2 + k^{2.2} \log^2 k \cdot w_{\mathsf{max}})$$

where the first step follows from $|S| = O(k^\gamma \log k)$, the second step follows from $\gamma = 0.1$.

The guarantee of table $\widehat{O}_{\lambda^2}$ follows from properties of graph $\mathsf{NG}_\lambda$ (Algorithm 1 provides the first property of table in Theorem statement, Lemma 3.4 provides the second property of table in Theorem statement). $\qquad\square$

## 3.2 Sparsification for arbitrary $\lambda$ (Step 1 of Theorem 2.1)

One shortcoming of Lemma 3.4 is that the number of remaining edges is only truly subquadratic if $\lambda$ is constant. As we discuss in Section 3.1, the overall running time of the algorithm depends on the number of edges in the remaining graph; and in order for the running time to be truly subquadratic, we need to reduce the number of edges to truly subquadratic. In this section, we show how one can obtain this bound even when $\lambda$ is sub-constant. However, instead of losing a factor $\lambda^2$ in the approximation, our technique loses a factor of $\Omega(\lambda^3)$.

From here on, we assume that all of the windows are of the same length and we show at the end of the section that this is almost without loss of generality. We begin by giving a definition:

**Definition 3.14** ($\mathsf{lcs}_{w_a}(w_i, w_j)$)**.** For two windows $w_i$ and $w_j$, and a window $w_a$, define $\mathsf{lcs}_{w_a}(w_i, w_j)$ as the length of the $\mathsf{LCS}$ of $\mathsf{opt}_{i,a}$ and $w_j$, where $\mathsf{opt}_{i,a}$ denotes a fixed $\mathsf{LCS}$ of $w_i$ and $w_a$.

Notice that unlike $\mathsf{lcs}$, this new definition is not symmetric. That is, the size of $\mathsf{lcs}_{w_a}(w_i, w_j)$ may be different from the size of $\mathsf{lcs}_{w_a}(w_j, w_i)$. Throughout this section and Section 3.2.1, $\mathsf{opt}_{i,a}$ refers to a fixed (e.g., lexicographically smallest) longest common subsequence of $w_i$ and $w_a$. We moreover assume that $\mathsf{opt}_{i,a}$ and $\mathsf{opt}_{a,i}$ refer to the same matching.

Similar to $\mathsf{lcs}$, we also normalize the size of $\mathsf{lcs}_s$ by the geometric mean of the lengths of the two windows. Our assumption from here on, until the statement of Theorem 3.21 is that all the windows are of the same length which implies that

$$||\mathsf{lcs}_{w_a}(w_i, w_j)|| = \mathsf{lcs}_{w_a}(w_i, w_j)/\sqrt{|w_i||w_j|} = \mathsf{lcs}_{w_a}(w_i, w_j)/|w_i| = \mathsf{lcs}_{w_a}(w_i, w_j)/|w_j|.$$

We bring a reduction in Theorem 3.21 to make our solution work for windows of arbitrary length. In what follows, we first give an algorithm for detecting *close pairs* (a notion that we introduce later in the section) of windows, and then prove that the number of remaining pairs whose normalized $\mathsf{lcs}$ is at least $\lambda$ is truly subquadratic.

Our algorithm is based on a data structure which we call $\mathsf{lcs\text{-}cmp}(w_a, S, \widetilde{\lambda})$. The reader can think of $\widetilde{\lambda} = \lambda^2/2$. Let us fix a threshold $\widetilde{\lambda}$ and a window $w_a$. $\mathsf{lcs\text{-}cmp}(w_a, S, \widetilde{\lambda})$ receives a set $S$ of windows as input and preprocesses the windows in time $\widetilde{O}(|w_a|^2|S|)$. Next, $\mathsf{lcs\text{-}cmp}(w_a, S, \widetilde{\lambda})$ would be able to answer each query of the following form in almost linear time ($O(w_{\mathsf{max}})$):

> Given two windows $w_i, w_j \in S$, either certify that $||\mathsf{lcs}_{w_a}(w_i, w_j)|| < \widetilde{\lambda}$ or find a solution for $||\mathsf{lcs}_{w_a}(w_i, w_j)||$ of size at least $\widetilde{\lambda}^2/4$.

We first show in Section 3.2.1, how $\mathsf{lcs\text{-}cmp}$ gives us a sparsification in truly subquadratic time and then discuss the algorithm for $\mathsf{lcs\text{-}cmp}$ in Section 3.2.2.

### 3.2.1 $\Omega(\lambda^3)$ Sparsification using lcs-cmp

We use a constant parameter $\gamma$ in our algorithm, and in the end we adjust $\gamma$ to minimize the total running time. Our algorithm repeats the following procedure $O(k^\gamma \log k)$ times: sample a window $w_a$ uniformly at random and let $S$ be the set of all the other windows. Next, by setting $\widetilde{\lambda} = \lambda^2/2$, we obtain lcs-cmp$(w_a, S, \lambda^2/2)$ via running the preprocessing step. Finally, for each pair of windows $w_i, w_j \in S$, we make a query to lcs-cmp to verify one of the following two possibilities:

- $||\mathsf{lcs}_{w_a}(w_i, w_j)|| < \lambda^2/2$;

- $||\mathsf{lcs}_{w_a}(w_i, w_j)|| \geq \lambda^4/16$.

If the latter is verified we set $\widehat{\mathsf{O}}_{\lambda^3}[i][j]$ and $\widehat{\mathsf{O}}_{\lambda^3}[j][i]$ to 1 otherwise we take no action. In what follows, we prove that after the above sparsification, the number of edges in the remaining graph is truly subquadratic.

---

**Algorithm 2** Sparsification for arbitrary $\lambda$ (Step 1 of Theorem 2.1)

---

1: **procedure** CUBICSPARSIFICATION$(w_1, w_2, \ldots, w_k, \lambda)$             $\triangleright$ Theorem 3.21
2:      $\widetilde{\lambda} \leftarrow \lambda^2/2$
3:      **for** counter $= 1 \rightarrow 10k^\gamma \log k$ **do**
4:          Sample $a \sim [k]$ uniformly at random
5:          $S \leftarrow \emptyset$
6:          **for** $i = 1 \rightarrow k$ **do**
7:              **if** $i \neq a$ **then**
8:                  $S \leftarrow S \cup \{w_i\}$
9:              **end if**
10:          **end for**
11:          lcs-cmp.INITIAL$(w_a, S, \widetilde{\lambda})$                  $\triangleright$ Algorithm 3, Lemma 3.22
12:          **for** $w_i \in S$ **do**
13:              **for** $w_j \in S$ **do**
14:                  **if** lcs-cmp.QUERY$(w_i, w_j)$ outputs accept **then**      $\triangleright$ Alg. 3, Lemma 3.23
15:                      $\widehat{\mathsf{O}}_{\lambda^3}[i][j] \leftarrow 1$               $\triangleright ||\mathsf{lcs}_{w_a}(w_i, w_j)|| \geq \widetilde{\lambda}^2/8$
16:                  **end if**
17:              **end for**
18:          **end for**
19:      **end for**
20:      **return** $\widehat{\mathsf{O}}_{\lambda^3}$
21: **end procedure**

---

Our goal for the rest of this section is to prove (Lemma 3.19) an upper bound on the number of edges that remain in the following graph:

**Definition 3.15** ($\mathsf{NG}_\lambda$). Define a graph $\mathsf{NG}_\lambda$ with $k$ vertices and the following edges:

$$E(\mathsf{NG}_\lambda) = \left\{ (i, j) \;\middle|\; ||\mathsf{lcs}(w_i, w_j)|| \geq \lambda \quad \text{and} \quad \widehat{\mathsf{O}}_{\lambda^3}[i][j] = 0 \right\}.$$

We define a notation called "close" which is similar to the notion of "well-connected" vertices in Section 3.1.

**Definition 3.16** (close). Let $\gamma \in (0,1)$ and $\lambda \in (0,1)$. We say a pair $(w_i, w_j)$ of windows is *close*, if there are at least $k^{1-\gamma}$ windows $w_a$ such that $||\mathsf{lcs}_{w_a}(w_i, w_j)|| \geq \lambda^2/2$.

Our first observation is that Algorithm 2 detects all the close pairs with high probability.

**Lemma 3.17.** *Let $\widehat{\mathsf{O}}_{\lambda^3} \in \{0,1\}^{k \times k}$ be the output of Algorithm 2. For each $(i,j) \in [k] \times [k]$, if $(w_i, w_j)$ is close (Definition 3.16), then $\widehat{\mathsf{O}}_{\lambda^3}[i][j] = 1$ holds with probability at least $1 - 1/k^3$.*

*Proof.* We consider a fixed $(i,j)$ such that $(w_i, w_j)$ is close. By Definition 3.16, there are at least $k^{1-\gamma}$ windows $w_a$ such that $||\mathsf{lcs}_{w_a}(w_i, w_j)|| \geq \lambda^2/2$. If such a $w_a$ window is sampled in our algorithm, then we detect the edge between the close pair. The probability that none of these windows is sampled is at most

$$\left(1 - \frac{k^{1-\gamma}}{k}\right)^{10k^\gamma \log k} = \left(1 - \frac{1}{k^\gamma}\right)^{10k^\gamma \log k} \leq \frac{1}{k^5}.$$

Taking a union over at most $k^2$ pairs completes the proof. $\qquad\square$

Before we proceed to Lemma 3.19, we bring Lemma 3.18 as an auxiliary observation.

**Lemma 3.18** (existence of a correlated pair). *Let $\lambda \in (0,1)$. Given a window $w_a$ and a set $T$ containing at least $2/\lambda$ windows. If for each $w_i \in T$ we have $||\mathsf{lcs}(w_a, w_i)|| \geq \lambda$, then there exist two windows $w_i, w_j \in T$ such that both $||\mathsf{lcs}_{w_a}(w_i, w_j)|| \geq \lambda^2/2$ and $||\mathsf{lcs}_{w_a}(w_j, w_i)|| \geq \lambda^2/2$ hold.*

*Proof.* Let $d$ be the size of the windows. We consider a character-based bipartite graph: On one side, it has $d$ nodes, and on the other side it has $d|T|$ nodes. Two nodes $x, y$ in the character-based graph are adjacent iff $(x,y) \in \mathsf{opt}_{a,i}$ where $w_a$ is the window containing character $x$ and $w_i$ is the window containing character $y$. Since $||\mathsf{lcs}(w_a, w_i)|| \geq \lambda$ for every $w_i \in T$, the total number of edges in character-based bipartite graph is at least $\lambda|T|d$.

For $\ell$-th character in window $w_a$, we use $D_\ell$ to denote the degree of the corresponding node in the character-based bipartite graph. The number of 2-walks between pairs of nodes on the side with $d|T|$ nodes is at least

$$
\begin{aligned}
\sum_{\ell=1}^{d} D_\ell(D_\ell - 1) &= \left(\sum_{\ell=1}^{d} D_\ell^2 - \sum_{\ell=1}^{d} D_\ell\right) \\
&\geq \left(\frac{1}{d}(\sum_{\ell=1}^{d} D_\ell)^2 - \sum_{\ell=1}^{d} D_\ell\right) && \text{by Cauchy-Schwarz inequality} \\
&\geq \left(\frac{1}{d}(\lambda|T|d)^2 - (\lambda|T|d)\right) && \text{by Eq. (8) explained below} \\
&= (\lambda^2 d|T|^2 - \lambda|T|d) \\
&\geq 1/2\lambda^2 d|T|^2 && \text{by } |T| \geq 2/\lambda
\end{aligned}
$$

The number of 3-tuples $(w_i, w_a, w_j)$ is at most $|T|^2$. Thus, there must exist a pair $(w_i, w_j)$ such that

$$||\mathsf{lcs}_{w_a}(w_i, w_j)|| \geq \lambda^2/2.$$

This also means that

$$||\mathsf{lcs}_{w_a}(w_j, w_i)|| \geq \lambda^2/2$$

which is the other conclusion of our lemma.

It remains to show Eq. (8). Since the number of edges in the character-based bipartite graph is at least $\lambda|T|d$, we have

$$\sum_{\ell=1}^{d} D_\ell \geq \lambda|T|d$$

$$\geq \lambda(2/\lambda)d > d. \qquad (|T| \geq 2/\lambda \text{ by lemma premise}) \qquad (6)$$

We also require the following simple fact:

$$\forall x \geq y \geq z \geq 0, \quad x(x-z) \geq y(y-z). \qquad (7)$$

Combining the last two inequalities, we finally have:

$$\frac{1}{d}(\sum_{\ell=1}^{d} D_\ell)^2 - \sum_{\ell=1}^{d} D_\ell = \frac{\sum_{\ell=1}^{d} D_\ell}{d}(\sum_{\ell=1}^{d} D_\ell - d)$$

$$\geq \frac{1}{d}(\lambda|T|d)^2 - (\lambda|T|d) \qquad (\text{Ineq. (6) and (7)}). \qquad (8)$$

$\square$

Now, we are ready to prove that the remaining graph is sparse.

**Lemma 3.19** (upper bound on $|E(\mathsf{NG}_\lambda)|$)**.**

$$|E(\mathsf{NG}_\lambda)| \leq \frac{2k^{2-\gamma}}{\lambda}$$

*holds with probability at least* $1 - 1/k^3$.

In the proof of this lemma we will construct an auxiliary graph $\mathsf{NF}_\lambda$ in the following way:

**Definition 3.20** ($\mathsf{NF}_\lambda$)**.** Let $a$ denote the node in $V(\mathsf{NG}_\lambda)$ that has the highest degree. The vertices of $\mathsf{NF}_\lambda$ would be the set of neighbors of $a$ in $\mathsf{NG}_\lambda$. We add an edge between vertices $(i, j)$ in $\mathsf{NF}_\lambda$ if both $||\mathsf{lcs}_{w_a}(w_i, w_j)|| \geq \lambda^2/2$ and $||\mathsf{lcs}_{w_a}(w_j, w_i)|| \geq \lambda^2/2$ hold.

*Proof of Lemma 3.19.* We prove the lemma by contradiction. Suppose

$$|E(\mathsf{NG}_\lambda)| > \frac{2k^{2-\gamma}}{\lambda}.$$

Since the number of vertices in $\mathsf{NF}_\lambda$ is equal to the maximum degree of $\mathsf{NG}_\lambda$ we have

$$|V(\mathsf{NF}_\lambda)| > \frac{2k^{1-\gamma}}{\lambda}. \qquad (9)$$

Using Lemma 3.18, we have for each set $T \subseteq V(\mathsf{NF}_\lambda)$ with $|T| \geq 2/\lambda$, there exist two nodes $u$ and $v$ in $T$ such that $(u, v)$ is an edge of $\mathsf{NF}_\lambda$.

If we look at the complement of graph $\mathsf{NF}_\lambda$, we know there is no clique $K_r$ where $r = 2/\lambda$. Using Turan's theorem (Lemma A.5) we know that the complement of graph $\mathsf{NF}_\lambda$ has at most $(1 - \frac{1}{r-1})\frac{q^2}{2} < (1 - \frac{1}{r})\frac{q^2}{2}$ edges, where $q := |V(\mathsf{NF}_\lambda)|$. Then we have

$$|E(\mathsf{NF}_\lambda)| > \frac{q(q-1)}{2} - (1 - \frac{1}{r})\frac{q^2}{2} = \frac{q^2}{2r} - \frac{q}{2}.$$

This implies that there exists a vertex $b$ whose degree in $\mathsf{NF}_\lambda$ is more than

$$\frac{|V(\mathsf{NF}_\lambda)|}{r} - 1 = \frac{|V(\mathsf{NF}_\lambda)|}{2/\lambda} - 1 \geq k^{1-\gamma} - 1,$$

where the inequality follows by Eq. (9). Thus, the degree of $b$ is at least $k^{1-\gamma}$. For each vertex $c$ of $\mathsf{NF}_\lambda$ which is adjacent to $b$ we have

$$||\mathsf{LCS}_{w_c}(w_a, w_b)|| = ||\mathsf{LCS}_{w_a}(w_c, w_b)|| \geq \lambda^2/2.$$

Since the degree of $b$ in $\mathsf{NF}_\lambda$ is at least $k^{1-\gamma}$, this means that pair $(w_a, w_b)$ is a close pair and the edge $(w_a, w_b)$ should not have existed between the vertices in $\mathsf{NG}_\lambda$ in the first place. $\qquad\square$

Now, we are ready to bring our main theorem of this section. In Theorem 3.21 we assume that the windows may have different length but the total number of distinct window sizes is bounded by $\mathsf{w}_{\mathsf{layers}}$.

**Theorem 3.21** (Sparsification for arbitrary $\lambda$ (Step 1 of Theorem 2.1)). *Given $k$ windows $w_1, \cdots, w_k$. Let $\mathsf{w}_{\mathsf{layers}}$ denote the number of different sizes for windows. For any $\lambda \in (0,1)$ and $\gamma \in (0,1)$, there is a randomized algorithm (Algorithm 2) that runs in time*

$$O(\mathsf{w}_{\mathsf{layers}}^2(k^{1+\gamma}\mathsf{w}_{\mathsf{max}}^2 \log^2 k + k^{2+\gamma}\mathsf{w}_{\mathsf{max}} \log k))$$

*and outputs a table $\widehat{\mathsf{O}}_{\lambda^3} \in \{0,1\}^{k \times k}$ such that*

$$\mathsf{lcs}(w_i, w_j)/\max\{|w_i|, |w_j|\} \geq \lambda^4/16, \ \ if \ \widehat{\mathsf{O}}_{\lambda^3}[i][j] = 1$$

*and*

$$\left|\left\{(i,j) \mid \mathsf{lcs}(w_i, w_j)/\max\{|w_i|, |w_j|\} \geq \lambda, \ \ \text{and} \ \ \widehat{\mathsf{O}}_{\lambda^3}[i][j] = 0\right\}\right| = O(k^{2-\gamma}\mathsf{w}_{\mathsf{layers}}^2/\lambda).$$

*The algorithm has success probability $1 - 1/k^3$.*

*Proof.* We run the algorithm explained above $\mathsf{w}_{\mathsf{layers}} + \binom{\mathsf{w}_{\mathsf{layers}}}{2}$ times. Each of the first $\mathsf{w}_{\mathsf{layers}}$ runs considers one set of windows with equal sizes. Each of the next $\binom{\mathsf{w}_{\mathsf{layers}}}{2}$ runs on one pair of window sizes. For such runs, we pad the smaller windows by dummy characters to make sure all windows have equal lengths. Our estimation for each pair of windows is obtained in one of the runs (which focuses on those particular lengths). Thus, both the runtime and the number of false-negatives of our algorithm are multiplied by a factor of $O(\mathsf{w}_{\mathsf{layers}}^2)$.

The running time of each round of Algorithm 2 is

$$\begin{aligned}
&= O(\mathsf{lcs\text{-}cmp.}\textsc{Initial} \text{ time} + |S|^2 \cdot (\mathsf{lcs\text{-}cmp.}\textsc{Query} \text{ time})) \\
&= O(|S|\mathsf{w}_{\mathsf{max}}^2 \log k + |S|^2\mathsf{w}_{\mathsf{max}}) \\
&= O(k\mathsf{w}_{\mathsf{max}}^2 \log k + k^2\mathsf{w}_{\mathsf{max}})
\end{aligned}$$

Since the algorithm repeats the sampling procedure $O(k^\gamma \log k)$ rounds, thus the overall running time is

$$O(k^\gamma \log k) \cdot O(k\mathsf{w}_{\mathsf{max}}^2 \log k + k^2\mathsf{w}_{\mathsf{max}}) = O(k^{1+\gamma}\mathsf{w}_{\mathsf{max}}^2 \log^2 k + k^{2+\gamma}\mathsf{w}_{\mathsf{max}} \log k).$$

$\qquad\square$

### 3.2.2 Implementation of lcs-cmp

The goal of this section is to design a data-structure that will be used in our subquadratic time algorithm for LCS. The data structure receives a window $w_a$, a threshold $\widetilde{\lambda}$, and a set $S$ of windows $w_1, w_2, \ldots, w_s$. (To avoid confusing, one can assume $a = 0$ and $w_a = w_0$ but we refer to this special window by $w_a$ to be consistent with the previous sections.) All the windows $w_a, w_1, w_2, \ldots, w_s$ have equal sizes. We present an $O(s|w_a|^2 \log n)$ time preprocessing algorithm and an $O(|w_a|)$ time query algorithm for lcs-cmp such that for any $i, j \in [s]$

1. if $\|\mathsf{lcs}_{w_a}(w_i, w_j)\| \geq \widetilde{\lambda}$, then the query algorithm outputs accept;

2. if $\|\mathsf{lcs}_{w_a}(w_i, w_j)\| < \widetilde{\lambda}^2/4$, then the query algorithm outputs reject.

We bring the high level idea in the following (see also Algorithm 3 for pseudocode). We first compute $\mathsf{opt}_{a,i}$ for every $i \in [s]$. Next, for each window $w_i$, we find a set of at most $2/\widetilde{\lambda}$ common subsequences between $w_a$ and $w_i$. Let us put the corresponding indices in the $w_a$ side of all such common subsequences in a set $Y_{a,i}$ for each window $w_i$. Our algorithm maintains the property that if we remove all characters corresponding to $Y_{a,i}$ from $w_a$, the LCS of the remainder with $w_i$ is smaller than $\frac{|w_a|\widetilde{\lambda}}{2}$. For each pair $w_i, w_j \in S$, if $\mathsf{LCS}_{w_a}(w_i, w_j) \geq |w_a|\widetilde{\lambda}$ then at least half of the characters that contribute to such a solution lie in $Y_{a,j}$. Thus, the intersection of at least one common subsequence in $Y_{a,j}$ and $\mathsf{opt}_{a,i}$ should have size $\frac{|w_a|\widetilde{\lambda}}{2}/(2/\widetilde{\lambda}) = |w_a|\widetilde{\lambda}^2/4$.

In what follows, we discuss the details.

**Lemma 3.22** (INITIAL). *Given a parameter $\widetilde{\lambda}$, a window $w_a$, and a set of windows $w_1, \cdots, w_s$, the* INITIAL *procedure of data structure* lcs-cmp *(Algorithm 3) takes $O(s|w_a|^2 \log n)$ time, and outputs $\{\mathsf{opt}_{a,i}\}_{i \in [s]}$ and $\{Y_{a,i}\}_{i \in [s]}$ such that the following hold.*

1. $\mathsf{opt}_{a,i}$ *corresponds to the $w_a$ side indices of a fixed longest common subsequence between $w_a$ and $w_i$ for every $i \in [s]$.*

2. $Y_{a,i}$ *corresponds to the $w_a$ side indices of at most $2/\widetilde{\lambda}$ common subsequence between $w_a$ and $w_i$ for every $i \in [s]$.*

3. *For any common subsequence between $w_a$ and $w_i$, if none of its $w_a$ side indices are in $Y_{a,i}$, then the length of this common subsequence is less than $\widetilde{\lambda}|w_a|/2$.*

*Proof.* In order to construct set $Y_{a,i}$ for each window $w_i$ we run the following algorithm: We start by setting $Y_{a,i}$ to an empty set. Next, we set $w'_a = w_a$ and iteratively compute the LCS of $w'_a$ and $w_i$. Each time we find a solution, we compare its size to $\widetilde{\lambda}|w_a|/2$. If the solution size is smaller, we terminate the algorithm. Otherwise, we add the original positions of the common subsequence in $w_a$ to set $Y_{a,i}$ and continue by removing these characters from $w'_a$. We stop when the solution size drops below $\widetilde{\lambda}|w_a|/2$. It immediately follows that since each time the size of $Y_{a,i}$ is increased by at least $\widetilde{\lambda}|w_a|/2$ then we repeat this procedure at most $2/\widetilde{\lambda}$ times.

To bound the runtime, we state the following fact: when the size of the LCS between a string $Q$ and a string $R$ is bounded by $\ell$, one can preprocess one of the strings ($R$ in this case) in time $O(|R| \log n)$ and then compute the LCS in time $O(|Q| \ell \log n)$ (see Theorem A.8). Notice that in order to construct $Y_{a,i}$ for a window $w_i$, the total size of the LCS's that we find is bounded by $|w_a|$. Therefore, for a fixed window $w_i$, if we preprocess $w_i$ once and use it for all computations, the total runtime is bounded by

$$O(|w_a||Y_{a,i}| \log n + |w_i| \log n) = O(|w_a|^2 \log n + |w_i| \log n).$$

Thus, the total preprocessing time of our algorithm over all windows is $O(s|w_a|^2 \log n)$. $\qquad \square$

---
**Algorithm 3** lcs-cmp data structure
---
1: **data structure** lcs-cmp
2:
3: **members**
4:     $\mathsf{opt}_{a,1}, \cdots, \mathsf{opt}_{a,s}$
5:     $Y_{a,1}, \cdots, Y_{a,s}$
6:     $\widetilde{\lambda} \in (0, 1)$
7: **end members**
8:
9: **procedure** INITIAL($w_a, \{w_i\}_{i \in [s]}, \widetilde{\lambda}$)                          ▷ Lemma 3.22
10:     **for** $i \in [s]$ **do**
11:         compute $\mathsf{opt}_{a,i}$, a fixed (e.g. lexicographically smallest) LCS between $w_a$ and $w_i$
12:     **end for**
13:     **for** $i \in [s]$ **do**
14:         $Y_{a,i} \leftarrow \emptyset$
15:         $w'_a \leftarrow w_a$
16:         **while** $\mathsf{LCS}(w'_a, w_i) \geq |w_a|\widetilde{\lambda}/2$ **do**
17:             $Y_{a,i} \leftarrow Y_{a,i} \cup$ characters in the LCS of $w'_a$ and $w_i$
18:             $w'_a \leftarrow w'_a \setminus$ characters in the LCS of $w'_a$ and $w_i$
19:         **end while**
20:     **end for**
21:     **return** $\{\mathsf{opt}_{a,i}\}_{i \in [s]}, \{Y_{a,i}\}_{i \in [s]}$
22: **end procedure**
23:
24: **procedure** QUERY($w_i, w_j$)                          ▷ Lemma 3.23
25:     **if** $|\mathsf{opt}_{a,i} \cap Y_{a,j}| > \widetilde{\lambda}|w_a|/2$ **then**
26:         **return** accept
27:     **else**
28:         **return** reject
29:     **end if**
30: **end procedure**
31:
32: **end data structure**
---

**Lemma 3.23** (QUERY). *For any $(i, j) \in [s] \times [s]$, the* QUERY *of data structure* lcs-cmp *(Algorithm 3) runs in time $O(|w_a|)$ with the following properties:*

1. *If $\mathsf{lcs}_{w_a}(w_i, w_j) \geq \widetilde{\lambda}|w_a|$, then it outputs accept.*

2. *If $\mathsf{lcs}_{w_a}(w_i, w_j) < \widetilde{\lambda}^2|w_a|/4$, then it outputs reject.*

*Proof.* Our algorithm takes the intersection of the characters that contribute to $\mathsf{opt}_{a,i}$ with $Y_{a,j}$. As explained earlier, if $\mathsf{lcs}_{w_a}(w_i, w_j) \geq \widetilde{\lambda}|w_a|$ then at least half of these characters belong to $Y_{a,j}$ and thus the intersection of one of the common subsequences contributing to $Y_{a,j}$ with $\mathsf{opt}_{a,i}$ is at least $\widetilde{\lambda}^2|w_a|/4$. Indeed, this can only happen if $\mathsf{lcs}_{w_a}(w_i, w_j) \geq \widetilde{\lambda}^2|w_a|/4$ in the first place. Otherwise we certainly output reject. □

# 4   LCS Step 0: Window-compatible Solutions

The goal of this section is to show that window-compatible solutions provide very accurate estimates for LCS. Roughly speaking, in this section, we construct a set of windows for $A$ and a set of windows for $B$ and we show that it suffices to only compute the LCS of pairs of windows. If that information is available, then we can estimate the LCS of the original two strings very accurately.

We use the same definition of window-compatible transformation as [BEG$^+$18], but we will use slightly different sets of windows. Every window is a (contiguous) substring of each string. Notice that windows may not have the same length. We use $W_A$ to denote the set of windows constructed for $A$ and $W_B$ for the set of the windows that we construct for $B$.

Let us first state our definition of window-compatible solutions for LCS.

**Definition 4.1** (Window-compatible common subsequence).    • Let $S = \langle w_1, \cdots, w_\alpha \rangle$ and $S' = \langle w_1', \cdots w_\alpha' \rangle$ be two sequences of non-overlapping windows from $W_A$ and $W_B$, respectively. We call a common subsequence of $(A, B)$ *window-compatible with respect to $S$ and $S'$*, if it is a union of $k$ common subsequences of $(w_1, w_1'), \ldots, (w_\alpha, w_\alpha')$

   • We call a common subsequence *window-compatible*, if it is window-compatible with respect to some pair of sequences of non-overlapping windows from $W_A$ and $W_B$, ordered by their respective starting indices.

$W_A$ will simply be a partitioning of $A$ into disjoint windows of length $d$ (for parameter $d$ to be finalized later). Setting $d = n^x$ leads to a truly subquadratic time solution for any $0 < x < 1$, however, one can play with the value of $d$ to optimize the running time. Below we consider two constructions of windows $W_B$ for string $B$.

## 4.1   Window construction for constant $\lambda$ (Step 0 of Theorem 2.2)

For clarity, we use $B^{(j,l)}$ to denote a length-$l$ substring of $B$ which starts at index $j$ and ends at index $j + l - 1$. Our construction has multiple layers, where each layer contains windows with equal lengths. For parameter $\epsilon_0 \in (0, 1)$ to be defined later, let

$$f := \lceil \log_{1+\epsilon_0}(1/\epsilon_0) \rceil + 1 = \Theta \left( \frac{1}{\epsilon_0} \log(\frac{1}{\epsilon_0}) \right).$$

For each $i \in \{-f, \cdots, f\}$, we define $d_i = d(1 + \epsilon_0)^i$, $g_i = \epsilon_0 d_i$, and $t_i = n/d_i$. For each $i$, let $W_B^i$ denote the set of all the windows in this layer. In the $i$-th layer, all the windows have the same size which is $d_i$. $g_i$ is the shift size, i.e. for each window (except the leftmost and rightmost windows), if we shift by $\pm g_i$ we get another window in $W_B^i$. $t_i$ is the number of windows in the $i$-th layer.

**Definition 4.2** (Window construction for constant $\lambda$)**.**

$$W_B^i := \{ B^{(\text{left}, \text{len})} \mid \text{left} = x \cdot g_i + y \cdot d_i, \text{len} = d_i,$$
$$\forall x \in \{0, 1, \cdots, d_i/g_i - 1\}, \forall y \in \{0, 1, \cdots, t_i - 1\} \}$$

where we use $B^{(\text{left}, \text{len})}$ denotes a (contiguous) substring of string $B$ that starts at left and has length len. Finally $W_B$ is $\cup_{i \in \{-f, \cdots, f\}} W_B^i$.

**Fact 4.3** (Parameters for constant $\lambda$ (Theorem 2.2))**.**

   *Recall that $\epsilon$ is a small constant and $\lambda$ is the relative length of the true LCS between $A$ and $B$. In order to make sure the loss in the approximation is negligible, we use*

$$\epsilon_0 := \epsilon \lambda.$$

**Algorithm 4** Window construction algorithm for constant $\lambda$

---

1: **procedure** QUADRATICWINDOWS$(A, B, n, d, \epsilon_0)$
2:                                                                                              ▷ Generate $W_A$:
3:      $W_A \leftarrow \emptyset$
4:      $t \leftarrow n/d$
5:      **for** $y = 0 \to t - 1$ **do**
6:         left $\leftarrow y \cdot d$
7:         $W_A \leftarrow W_A \cup \left\{ A^{(\text{left}, d)} \right\}$
8:      **end for**
9:                                                                                              ▷ Generate $W_B$:
10:     $f \leftarrow \lceil \log_{1+\epsilon_0}(1/\epsilon_0) \rceil + 1$                      ▷ $d \cdot (1 + \epsilon_0)^f = \Theta(d/\epsilon_0)$
11:     $W_B \leftarrow \emptyset$
12:     **for** $i = -f \to f$ **do**
13:        $W_B^i \leftarrow \emptyset$
14:        $d_i \leftarrow d(1 + \epsilon_0)^i$, $t_i \leftarrow n/d_i$, $g_i \leftarrow \epsilon_0 d_i$     ▷ $d_i$ is the length, $g_i$ is the shift
15:        **for** $x = 0 \to d_i/g_i - 1$ **do**
16:           **for** $y = 0 \to t_i - 1$ **do**
17:              left $\leftarrow x \cdot g_i + y \cdot d_i$
18:              len $\leftarrow d_i$
19:              $W_B^i \leftarrow W_B^i \cup \left\{ B^{(\text{left}, \text{len})} \right\}$
20:           **end for**
21:        **end for**
22:        $W_B \leftarrow W_B \cup W_B^i$
23:     **end for**
24:     **return** $W_A, W_B, f$
25: **end procedure**

---

- *The total number of windows is equal to $k = |W_A| + |W_B| = O((1/\lambda)^3 n/d)$.*

- *The maximum size of the windows is equal to $\mathsf{w}_{\mathsf{max}} = \Theta(d/\lambda)$.*

- *The minimum size of the windows is equal to $\mathsf{w}_{\mathsf{min}} = \Theta(d\lambda)$.*

- *The ratio of the maximum window size over the minimum window size is bounded by $\mathsf{w}_{\mathsf{gap}} = \mathsf{w}_{\mathsf{max}}/\mathsf{w}_{\mathsf{min}} = \Theta(1/\lambda^2)$.*

- *The number of different layers which is equal to the number of different window sizes is equal to $\mathsf{w}_{\mathsf{layers}} = O((1/\lambda)\log(1/\lambda))$*

*Proof.* We derive the number of windows (the rest of the parameters follow immediately from the construction). In each layer of $W_B$, the number of windows is

$$|W_B^i| = O(t_i d_i / g_i) = O(n/(d_i \epsilon_0)).$$

Thus the total number of windows is given by

$$
\begin{aligned}
|W_B| &= \sum_{i=-f}^{f} |W_B^i| \\
&= O\left( n/\epsilon_0 \sum_{i=-f}^{f} 1/d_i \right) && \text{(Def. 4.2)} \\
&= O\Big( n/(d\epsilon_0) \underbrace{\sum_{i=-f}^{f} 1/(1+\epsilon_0)^i}_{=O(1/\epsilon_0^2)} \Big) \\
&= O(n/(d\epsilon_0^3)) && \text{(Geometric sequence).}
\end{aligned}
$$

$\square$

### 4.1.1 Near-optimality of window-compatible common subsequence

We present a
structural lemma for window-compatible common subsequences. This essentially, shows that the problem of computing $\mathsf{LCS}$ between the two strings reduces to the problem of computing the $\mathsf{LCS}$'s between the windows.

**Lemma 4.4** (Window-compatible structural lemma (constant $\lambda$)).
*For any two strings $A, B$ of length $n$: If $\|\mathsf{lcs}(A, B)\| \geq \lambda$, there exists a window-compatible common subsequence of $A, B$ of length at least $\lambda n - 8\epsilon_0 n$ for the window sets constructed in Definition 4.2.*

*Proof.* Fix a LCS of $A, B$.
Let $A_j$ be the $j$-th window of $A$, and let $\widetilde{B}_j$ be the minimal substring of $B$ containing all characters common to $A_j$ and the LCS. We consider three cases:

**Case 1:** $|\widetilde{B}_j| \in [\epsilon_0 d, d/\epsilon_0]$**:** In this case we can keep this pair.

34

**Case 2:** $|\widetilde{B}_j| \geq d/\epsilon_0$**:** Then we throw out this pair. We throw out at least $d/\epsilon_0$ characters, but we only decrease the entire lcs by at most $d$.

**Case 3:** $|\widetilde{B}_j| \leq d\epsilon_0$ **:** Then we also throw out this pair. We throw out at least $d$ characters, but we only decrease the entire lcs by at most $\epsilon_0 d$.

Let $Z$ denote the set of $A$-windows that we don't throw out. Whenever we throw out a pair, the amount of lcs got decreased is always at most $\epsilon_0$ fraction of the characters we throw out. There are $2n$ characters in total. Then we decrease the lcs by at most $2\epsilon_0 n$, i.e.,

$$\sum_{j \in Z} \mathsf{lcs}(\widetilde{A}_j, \widetilde{B}_j) \geq \mathsf{lcs}(A, B) - 2\epsilon_0 n, \tag{10}$$

We now derive, for each $j \in Z$, a *window* $B_j \subseteq \widetilde{B}_j$ of approximately the same length. Let $i \in \{-f+1, \ldots, f\}$ be the layer such that $|\widetilde{B}_j| \in (d(1+\epsilon_0)^i, d(1+\epsilon_0)^{i+1})$. Similarly, $\widetilde{B}_j$'s starting index is in $\big((y-1) \cdot \epsilon_0 d(1+\epsilon_0)^{j-1}, y \cdot \epsilon_0 d(1+\epsilon_0)^{j-1}\big]$ for some $y$. Let $B_j \in W_B$ be the window of length $d(1+\epsilon_0)^{j-1}$ with starting index in $y \cdot \epsilon_0 d(1+\epsilon_0)^{j-1}$. Notice that $B_i$ ends at

$$y \cdot \epsilon_0 d(1+\epsilon_0)^{j-1} + d(1+\epsilon_0)^{j-1} = y \cdot \epsilon_0 d(1+\epsilon_0)^{j-1} + d(1+\epsilon_0)^j - \epsilon_0 d(1+\epsilon_0)^{j-1}$$
$$= (y-1) \cdot \epsilon_0 d(1+\epsilon_0)^{j-1} + d(1+\epsilon_0)^j,$$

which is a lower bound on $\widetilde{B}_j$'s end index, hence indeed $B_j \subseteq \widetilde{B}_j$. Finally, notice that

$$\begin{aligned}
|\widetilde{B}_j \setminus B_j| &\leq d(1+\epsilon_0)^{i+1} - d(1+\epsilon_0)^{i-1} \\
&\leq d(1+\epsilon_0)^{i-1} 3\epsilon_0 \\
&< 3\epsilon_0 |\widetilde{B}_j|.
\end{aligned} \tag{11}$$

Therefore, in total the number of characters we lose by restricting to windows $B_j$ is bounded by $3\epsilon_0 n$:

$$\begin{aligned}
\sum_{i \in Z} |\mathsf{lcs}(A_j, B_j)| &\geq \sum_{j \in Z} (|\mathsf{lcs}(A_j, \widetilde{B}_j)| - 3\epsilon_0 |\widetilde{B}_j|) && \text{(Eq. (11))} \\
&\geq \sum_{j \in Z} |\mathsf{lcs}(A_j, \widetilde{B}_j)| - 3\epsilon_0 n && (\widetilde{B}_j\text{'s are disjoint}) \\
&\geq |\mathsf{lcs}(A, B)| - 8\epsilon_0 n && \text{(Eq. (10))}.
\end{aligned}$$

$\qquad\square$

## 4.2 Window construction for arbitrary $\lambda$ (Step 0 of Theorem 2.1)

Our construction for this case is similar to Section 4.1, with the exception that when we're OK with losing constant factors in approximation, we can afford less shifts and less layers, which will improve the running time.

Let

$$f := \log(1/\epsilon_0).$$

For each $i \in \{0, \cdots, f\}$, we define $d_i = d \cdot 2^i$, and $t_i = n/d_i$. For each $i$, let $W_B^i$ denote the set of all the windows in this layer. We consider window set $W_B^i$, in $i$-th layer, all the window have sizes that are multiples of $d$, and in the range $(d_i/2, d_i]$. The shift for each layer will be $d_i$.

**Definition 4.5** (Window construction for arbitrary $\lambda$)**.**

$$W_A, W_B^0 := \big\{ B^{(\text{left,len})} \mid \text{left} = y \cdot d,$$
$$\forall y \in \{0, 1, \cdots, t_0 - 1\} \big\}$$

$$W_B^i := \big\{ B^{(\text{left,len})} \mid \text{left} = y \cdot d_i, \text{len} = d_i/2 + x \cdot d,$$
$$\forall x \in \{1, \cdots, d_i/(2d)\} \forall y \in \{0, 1, \cdots, t_i - 1\} \big\}$$

where we use $B^{(\text{left,len})}$ to denote a (contiguous) substring of string $B$ that starts at left and has length len. Finally $W_B$ is $\cup_{i \in \{0,1,\cdots,f\}} W_B^i$.

In other words, the window construction in Definition 4.5 considers all the substrings that start at any multiple of $d \cdot 2^i$ and whose length is a multiple of $d$, and at most $d \cdot 2^i - 1$.

---

**Algorithm 5** Window construction algorithm for arbitrary $\lambda$

---

1: **procedure** CUBICWINDOWS$(A, B, n, d, \epsilon_0)$
2:                                                      $\triangleright$ Generate $W_A$ and $W_B^0$:
3:     $W_A, W_B^0 \leftarrow \emptyset$
4:     $t \leftarrow n/d$
5:     **for** $y = 0 \rightarrow t - 1$ **do**
6:         left $\leftarrow y \cdot d$
7:         $W_A \leftarrow W_A \cup \big\{ A^{(\text{left},d)} \big\}$
8:         $W_B^0 \leftarrow W_B^0 \cup \big\{ B^{(\text{left},d)} \big\}$
9:     **end for**
10:                                                       $\triangleright$ Generate $W_B$:
11:     $f \leftarrow \log(1/\epsilon_0)$                                           $\triangleright d \cdot 2^f = d/\epsilon_0$
12:     $W_B \leftarrow W_B^0$
13:     **for** $i = 1 \rightarrow f$ **do**
14:         $W_B^i \leftarrow \emptyset$
15:         $d_i \leftarrow d \cdot 2^i$, $t_i \leftarrow n/d_i$
16:         **for** $x = 1 \rightarrow d_i/(2d)$ **do**
17:             **for** $y = 0 \rightarrow t_i - 1$ **do**
18:                 left $\leftarrow y \cdot d_i$
19:                 len $\leftarrow d_i/2 + x \cdot d$
20:                 $W_B^i \leftarrow W_B^i \cup \big\{ B^{(\text{left,len})} \big\}$
21:             **end for**
22:         **end for**
23:         $W_B \leftarrow W_B \cup W_B^i$
24:     **end for**
25:     **return** $W_A, W_B, f$
26: **end procedure**

---

**Fact 4.6** (Parameters for arbitrary $\lambda$ (Theorem 2.1))**.**

*Let $\epsilon > 0$ be a small constant and $\lambda$ is the relative length of the true LCS between $A$ and $B$. In order to make sure the loss in the approximation is small, we use*

$$\epsilon_0 := \epsilon \lambda.$$

- *The total number of windows is equal to $k = |W_A| + |W_B| = \widetilde{O}(n/d)$.*

- *The maximum size of the windows is equal to $\mathsf{w_{max}} = \Theta(d/\lambda)$.*

- *The minimum size of the windows is equal to $\mathsf{w_{min}} = \Theta(d)$.*

- *The ratio of the maximum window size over the minimum window size is bounded by $\mathsf{w_{gap}} = \mathsf{w_{max}}/\mathsf{w_{min}} = \Theta(1/\lambda)$.*

- *The number of different layers which is equal to the number of different window sizes is equal to $\mathsf{w_{layers}} = O(\log(1/\lambda))$.*

*Proof.* We derive the number of windows (the rest of the parameters follow immediately from the construction). In each layer of $W_B$, the number of windows is

$$|W_B^i| = O(t_i \cdot (d_i/d)) = O(n/d).$$

Therefore the total number of windows is

$$|W_B| = O(fn/d) = \widetilde{O}(n/d).$$

$\square$

### 4.2.1 Up-to-constant-factor-optimality of window-compatible common subsequence

We present another structural lemma for window-compatible common subsequences. It is similar to Lemma 4.7, but loses a constant factor in length of common subsequence.

**Lemma 4.7** (Window-compatible structural lemma (arbitrary $\lambda$)).
    *For any two strings $X, Y$ of length $n$, it is possible to map them into $A, B$ in one of the following ways:*

- $A = X, B = Y$

- $A = Y, B = X$

- $A = \mathrm{reverse}(X), B = \mathrm{reverse}(Y)$, *or*

- $A = \mathrm{reverse}(Y), B = \mathrm{reverse}(X)$,

*such that the following holds: If $\|\,\mathsf{lcs}(A,B)\| \geq \lambda$, there exists a window-compatible common subsequence of $A, B$ of length at least $\Omega(\lambda n) - 2\epsilon_0 n$ for the window sets constructed in Definition 4.5.*

*Proof.* Fix any LCS of $X, Y$, and suppose we first try $A = X, B = Y$. For each $j$, let $A_j$ be the $j$-th window of $A$, and let $\widetilde{B}_j$ be the minimal substring of $B$ containing all characters common to $A_j$ and the LCS.

We will analyze a few different cases of $\widetilde{B}_j$; the trickiest case is when $\widetilde{B}_j$ and $\widetilde{B}_{j+1}$ are both contained in the same $W_B^0$-window of length $d$. Denote this window by $\widehat{B}_{j,j+1}$. We claim that in this case the opposite cannot also be true: The minimal substring of $A$ containing all characters common to $\widehat{B}_{j,j+1}$ and the LCS is not contained in any $W_A$-window — this is because it intersects both $A_j$ and $A_{j+1}$. We will throw out all such $j$'s; wlog their contribution is at most half of the LCS (otherwise we can consider the reverse assignment $A = Y, B = X$).

---

**Algorithm 6** Dynamic programming algorithm for block-based lcs problem

---

1: **procedure** DP-LCS($W_A, W_B, M$)                                    ▷ Lemma 4.8
2:     Note that $M$ is table s.t. $M(k_1, k_2) \leq \mathsf{lcs}(A_{k_1}, B_{k_2})$.
3:     Note that $W_A$ and $W_B$ are sets of windows.
4:
5:     Let $S_1$ denote the sorted list of right-indices of $W_A$-windows.
6:     Let $S_2$ denote the sorted list of right-indices of $W_B$-windows.
7:     **for** $i_1 \in S_1$ **do**
8:         **for** $i_2 \in S_2$ **do**
9:             $C[i_1][i_2] \leftarrow 0$
10:            $x_1 \leftarrow 0$
11:            **for** $A_{k_1}, B_{k_2}$ have right index $i_1, i_2$ **do**          ▷ Case 1
12:                $\text{left}_1 \leftarrow$ left index of $A_{k_1}$
13:                $\text{left}_2 \leftarrow$ left index of $B_{k_2}$
14:                $\text{tmp} \leftarrow C[\text{left}_1][\text{left}_2] + M(A_{k_1}, B_{k_2})$
15:                **if** $\text{tmp} > x_1$ **then**
16:                    $x_1 \leftarrow \text{tmp}$
17:                **end if**
18:            **end for**
19:            Let $\text{prev}(i_1)$ denote the first index $\in S_1$ that is earlier than $i_1$
20:            Let $\text{prev}(i_2)$ denote the first index $\in S_2$ that is earlier than $i_2$
21:            $x_2 \leftarrow C[\text{prev}(i_1)][i_2]$                          ▷ Case 2
22:            $x_3 \leftarrow C[i_1][\text{prev}(i_2)]$                          ▷ Case 3
23:            $C[i_1][i_2] \leftarrow \max(x_1, x_2, x_3)$
24:        **end for**
25:    **end for**
26: **end procedure**

---

We also throw out all the $\widetilde{B}_j$'s such that $|\widetilde{B}_j| > d/(2\epsilon_0)$. Throwing out each such $\widetilde{B}_j$ can decrease the LCS by at most $d$, but removes $d/(2\epsilon_0)$ characters; hence the total loss to LCS is at most $2\epsilon_0 \cdot n$.

Finally, we further throw out either all the $\widetilde{B}_j$'s that are contained in some $W_B^0$-window of length $d$, or all the ones that aren't; again this step loses at most half of the LCS.

**Case 1: we keep all the $\widetilde{B}_j$'s that are contained in some $W_B^0$-window of length $d$.** By the previous paragraph, each of them is the only $\widetilde{B}_j$ contained in that window, so we can just match $A_j$ to that window.

**Case 2: we keep $\widetilde{B}_j$'s that are not contained in any $W_B^0$-window of length $d$.** For each $\widetilde{B}_j$, consider the index that is a multiple of $d \cdot 2^i$ for the largest possible $i$ (notice that this is unique). We truncate $\widetilde{B}_j$ to the left of this index; wlog the total contribution of truncated characters to the LCS is at most half (otherwise we reverse both strings). Each $\widetilde{B}_j$ is now contained in a disjoint $W_B^i$-window, so we can match $A_j$ to that window.

□

## 4.3   Dynamic Programming for window-compatible LCS

**Lemma 4.8** (Dynamic Programming for Longest Common Sequence). *Given two strings $A, B$ of length $n$, respective sets of windows $W_A, W_B$, and estimates-table $M$ on the pairwise LCS, Algorithm 6 computes the longest common subsequence of $A, B$ that is both: (i) window-compatible, and*

*(ii) on each pair of windows uses the common subsequence length from $M$.*

*The algorithm runs in time $|W_A||W_B|$.*

*Proof.* The proof is based on a classic DP. Define $D[i][j]$ which stores the size of the solution ending at windows $w_i \in A$ and $w_j \in B$. We then can update the size of the solution for each pair in time $O(1)$ which gives us a solution in time $O(|W_A||W_B|)$. □

# 5  LCS Step 2: Nearby Searching

The arguments of this section follow from the work of [CDG$^+$18]. However, for the sake of completeness, we restate the theorems here. This section is dedicated to presenting a method NonMetricEstimation($W_A, W_B$) to be used for estimating the longest common subsequence. The output of this method is a matrix $\widehat{M} : [|W_A|] \times [|W_B|] \to \mathbb{Z}^+$ where $\widehat{M}[i][j]$ estimates the LCS of windows $w_i \in W_A$ and $w_j \in W_B$. We allow for errors in a few elements of $\widehat{M}$ but we prove at the end of this section that with high probability, the error does not affect the solution significantly.

Our algorithm is pretty simple: We call a pair $(w_i, w_j)$ of windows *underestimated* if their LCS is large enough (at least $\epsilon\lambda$ when normalizing the LCS size) but it is not computed within the desired approximation factor in Step 1.

We are guaranteed by the bounds of Step 1 that the total number of underestimated pairs of windows is sublinear in the number of pairs. Thus, we show their count by $k^{2-\eta}$. Define $\mathcal{W} = k^{\eta/2}\mathsf{w_{max}}$ and let two windows of the $A$ side or two windows of the $B$ side be *nearby* if and only if the distance of the starting indices of the windows is bounded by $\mathcal{W}$. Fix $\epsilon_{\mathsf{nbs}}$ to be a small error rate. In our algorithm, we sample the windows of the $A$ side with a rate $p = (10 \log n)k^{-\eta/2}/\epsilon_{\mathsf{nbs}}$ and discover all the underestimated pairs whose $A$ side is subsampled via a naive brute-force. We then recompute the LCS of each pair of windows $(w_i, w_j)$ such that $w_i \in W_A, w_j \in W_B$ and we detect an underestimated pair $(w_{i'}, w_{j'})$ such that $w_i$ and $w_{i'}$ are nearby and $w_j$ and $w_{j'}$ are also nearby. Finally, starting from the distance matrix provided in Step 1, we update the solution for each pair of windows that we compute their LCS from scratch and output it as matrix $\widehat{M}$. We prove in Lemma 5.1 that the error of our estimations is bounded with high probability.

---
**Algorithm 7** Nearby Searching
---
1: $p \leftarrow (10 \log n)k^{-\eta/2}/\epsilon_{\mathsf{nbs}}$
2: $\mathcal{W} \leftarrow k^{\eta/2}\mathsf{w_{max}}$
3: $\widehat{M} \leftarrow M$
4: **for** $w \in W_A$ **do**
5:     Skip the following commands with probability $1 - p$
6:     **for** $w' \in W_B$ **do**
7:         Compute the LCS of $w$ and $w'$
8:         **if** $(w, w')$ is underestimated in $M$ **then**
9:             **for** any window $w_1$ in $W_A$ which is nearby to $w$ **do**
10:                 **for** any window $w_2$ in $W_B$ which is nearby to $w'$ **do**
11:                     Update $\widehat{M}$ by computing the LCS of $w_1$ and $w_2$ from scratch.
12:                 **end for**
13:             **end for**
14:         **end if**
15:     **end for**
16: **end for**
17: **return** $\widehat{M}$
---

Let $M$ be the output of Step 1 of our algorithm and $\widehat{M}$ be the output of Step 2. We refer to the size of the optimal LCS made by matrices $M$ and $\widehat{M}$ by LCS($M$) and LCS($M'$) respectively. We moreover, define a matrix $M'$ which is the same as $M$ except that the LCS values of all the underestimated pairs are corrected in $M'$. Similarly, we denote by LCS($M'$) the size of the longest common subsequence made by running DP on matrix $M'$. We prove below that the error between LCS($M'$) and LCS($\widehat{M}$) is bounded.

**Lemma 5.1** (follows from [CDG$^+$18]). *After running Algorithm 7, $\mathsf{LCS}(\widehat{M}) \geq \mathsf{LCS}(M') - 2n\epsilon_{\mathsf{nbs}}$ holds with probability at least $1 - 1/n^4$.*

*Proof.* It follows from the definition that the error can only come from the underestimated pairs of windows; for all other pairs, the corresponding entries of matrix $\widehat{M}$ are at least as large as the respective entries of $M'$. We prove below that the error from underestimated pairs is bounded.

To this end, fix an optimal window-compatible solution with respect to $M'$ and let

$$(x_1, y_1), (x_2, y_2), \ldots$$

be the first indices of the underestimated pairs of windows that contribute to this optimal solution. Let the pairs be sorted by their index. That is $x_1 < x_2 < x_3 < \ldots$ and $y_1 < y_2 < y_3 < \ldots$ hold. The proof is based on the following intuitive argument: with high probability, all but very few pairs of underestimated windows that contribute to the optimal window-compatible solution of $M'$ are detected in Algorithm 7 and thus their values are corrected.

To see this, consider one pair $(x_i, y_i)$. This means that there is a window $w \in W_A$ and a window $w' \in W_B$ such that the starting index of $w$ is equal to $x_i$ and the starting index of $w'$ is equal to $y_i$. Moreover, the value of the pair $(w, w')$ is underestimated in $M$. Let $X$ be the set of all $(x_j, y_j)$'s in $(x_1, y_1), (x_2, y_2), \ldots$ such that $|x_j - x_i| \leq \mathcal{W}$ and $|y_j - y_i| \leq \mathcal{W}$. Notice that each elements of $X$ in addition to being an underestimated pair of windows also contributes to the optimal window-compatible solution of $M'$. If $|X| \geq \epsilon_{\mathsf{nbs}}k^{\eta/2}$, then with high probability the $A$-side window of one of these underestimated pairs will be sampled in Algorithm 7 and therefore after detecting the underestimated pair, we compute the $\mathsf{LCS}$ of $(x_i, y_i)$ from scratch. If $|X| < \epsilon_{\mathsf{nbs}}k^{\eta/2}$ then at least one of the following two constraints should hold:

- The number of $x_j$'s in $(x_1, y_1), (x_2, y_2), \ldots$ such that $x_i - \mathcal{W} \leq x_j \leq x_i$ is smaller than $\epsilon_{\mathsf{nbs}}k^{\eta/2}$.

- The number of $y_j$'s in $(x_1, y_1), (x_2, y_2), \ldots$ such that $y_i - \mathcal{W} \leq y_j \leq y_i$ is smaller than $\epsilon_{\mathsf{nbs}}k^{\eta/2}$.

Keep in mind that in the constrains above, each $(x_j, y_j)$ is an underestimated pair which contributes to the optimal solution of $M'$. If the former holds, we call $x_i$ *lonely*. Similarly, if the latter holds, we call $y_i$ *lonely*. A simple argument implies that the number of lonely $x_i$'s (and similarly the number of lonely $y_i$'s) is bounded by $(\epsilon_{\mathsf{nbs}}k^{\eta/2})n/\mathcal{W}$: Divide the string into blocks of size $\mathcal{W}$. In each block, at most $\epsilon_{\mathsf{nbs}}k^{\eta/2}$ $x_j$'s are lonely. Thus, the total number of lonely $x_j$'s is bounded by $(\epsilon_{\mathsf{nbs}}k^{\eta/2})n/\mathcal{W}$. A similar argument also holds for the number of lonely $y_j$'s.

Therefore, the size of $|X|$ may be smaller than $\epsilon_{\mathsf{nbs}}k^{\eta/2}$ for at most $2(\epsilon_{\mathsf{nbs}}k^{\eta/2})n/\mathcal{W}$ underestimated pairs in $(x_1, y_1), (x_2, y_2), \ldots$. This implies that with high probability, the error of $\widehat{M}$ in comparison to $M'$ is bounded by

$$2\mathsf{w}_{\mathsf{max}}(\epsilon_{\mathsf{nbs}}k^{\eta/2})n/\mathcal{W} = 2\mathsf{w}_{\mathsf{max}}(\epsilon_{\mathsf{nbs}}k^{\eta/2})n/(k^{\eta/2}\mathsf{w}_{\mathsf{max}})$$
$$= 2n\epsilon_{\mathsf{nbs}}$$

$\square$

**Theorem 5.2** (nearby searching for $\Omega(\lambda^3)$ approximation). *Given that the number of underestimated pairs in Step 1 is bounded by $k^{2-\eta}$, for an arbitrary small constant $\epsilon > 0$, Step 2 takes time $\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^6)$ and has approximation factor $1 - \epsilon$ with probability at least $1 - 1/n^3$.*

*Proof.* To make sure the multiplicative factor in the approximation remains $1 - \epsilon$, we set $\epsilon_{\mathsf{nbs}} = \Theta(\lambda^4)$. Thus, the runtime of the algorithm would be as follows: In the sampling process, we

subsample each window with probability $p = (10 \log n) k^{-\eta/2}/\epsilon_{\mathsf{nbs}} = \widetilde{O}(k^{-\eta/2}/\lambda^4)$. Therefore, the number of sampled windows would be $\widetilde{O}(k^{1-\eta/2}/\lambda^4)$ with high probability and therefore the runtime for detecting the underestimated pairs of windows would be

$$\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^4).$$

Since the total number of underestimated pairs is bounded by $k^{2-\eta}$ and we sample the $A$-side windows with probability $\widetilde{O}(k^{-\eta/2}/\lambda^4)$, the expected number of underestimated pairs we detect in the first step is $\widetilde{O}(k^{2-3\eta/2}/\lambda^4)$. Moreover, for each underestimated pair of windows that we detect, we recompute the LCS of all the nearby pairs of windows. The number of nearby windows on the $A$-side is $O(k^{\eta/2}/\lambda)$ and the number of nearby windows on the $B$ side is $\widetilde{O}(k^{\eta/2}/\lambda)$. This makes a total runtime of $\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^6)$.

Notice that the runtime of the algorithm is in expectation since it depends on the number of underestimated edges we detect in the sampling phase. To make the runtime fixed, we observe that with probability at least $1/2$, the algorithm terminates after $\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^6)$ operations. Thus, we run the algorithm in parallel $10 \log n$ times and report the output of any instance that terminates before $\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^6)$ operations. This comes with a constant factor overhead in the error rate of the algorithm which remains smaller than $1/n^3$. $\qquad\square$

**Theorem 5.3** (nearby searching for $(1-\epsilon)\lambda^2$ approximation). *Given that the number of underestimated pairs in Step 1 is bounded by $k^{2-\eta}$, for an arbitrary small constant $\epsilon > 0$, Step 2 takes time $\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^6)$ and has approximation factor $1 - \epsilon$ with probability at least $1 - 1/n^3$.*

*Proof.* To make sure the multiplicative factor in the approximation remains $1 - \epsilon$, we set $\epsilon_{\mathsf{nbs}} = \Theta(\lambda^3)$. Thus, the runtime of the algorithm would be as follows: In the sampling process, we subsample each window with probability $p = (10 \log n) k^{-\eta/2}/\epsilon_{\mathsf{nbs}} = \widetilde{O}(k^{-\eta/2}/\lambda^3)$. Therefore, the number of sampled windows would be $\widetilde{O}(k^{1-\eta/2}/\lambda^3)$ with high probability and therefore the runtime for detecting the underestimated pairs of windows would be

$$\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^3).$$

Since the total number of underestimated pairs is bounded by $k^{2-\eta}$ and we sample the $A$-side windows with probability $\widetilde{O}(k^{-\eta/2}/\lambda^3)$, the expected number of underestimated pairs we detect is $\widetilde{O}(k^{2-3\eta/2}/\lambda^3)$. Moreover, for each underestimated pair of windows that we detect, we recompute the LCS of all the nearby pairs of windows. The number of nearby windows on the $A$-side is $O(k^{\eta/2}/\lambda)$ and the number of nearby windows on the $B$ side is $\widetilde{O}(k^{\eta/2}/\lambda^2)$. This makes a total runtime of

$$\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^6).$$

Similar to Theorem 5.2, one can make sure the runtime of the algorithm is

$$\widetilde{O}(\mathsf{w}_{\mathsf{max}}^2 k^{2-\eta/2}/\lambda^6)$$

and the error remains bounded by $1/n^3$. $\qquad\square$

# 6 Longest Increasing Subsequence

We outlined the algorithm in Section 1.3.3. Here, we bring the details for each step of the algorithm. In Section 6.1 we discuss the solution domains and show how we construct them. Next, in Section 6.2 we discuss the details of constructing pseudo-solutions and finally in Section 6.3 we show how we can obtain an approximate solution from the pseudo-solutions. Also, in Section 6.4 we bring an improvement to the running time at the expense of having a larger approximation factor for the algorithm.

## 6.1 Solution Domains

We assume from now on that $\mathsf{lis}(A) \geq n\lambda$ holds. As mentioned earlier, we divide the input array into $\sqrt{n}$ subarrays of size $\sqrt{n}$ and denote them by $\mathsf{sa}_1, \mathsf{sa}_2, \ldots, \mathsf{sa}_{\sqrt{n}}$. For a fixed optimal solution $\mathsf{opt}$, our goal is to approximate the smallest and the largest number of each subarray that contribute to $\mathsf{opt}$. Let us refer to these numbers as the *domain* of each subarray. Let $\epsilon := 1/1000$ be the accuracy parameter. For a subarray $\mathsf{sa}_i$, we sample $k$ (will be decided later) different elements and refer to them by $a_{j_1}, a_{j_2}, \ldots, a_{j_k}$.

---

**Algorithm 8** Constructing the candidate domains

1: **procedure** CONSTRUCTCANDIDATEDOMAINS($\mathsf{sa}_i$)                    ▷ Lemma 6.1
2:                                                           ▷ Given random access to a subarray $\mathsf{sa}_i$
3:     $k \leftarrow 20 \log(1/\delta)/(\lambda\epsilon^2)$
4:     Sample $k$ elements from $\mathsf{sa}_i$, and denote the sampled elements by $a_{j_1}, a_{j_2}, \ldots, a_{j_k}$
5:     **for** $\alpha$ in $[k]$ **do**
6:         **for** $\beta$ in $[k]$ **do**
7:             If $a_{j_\alpha} \leq a_{j_\beta}$, then construct a candidate domain $[a_{j_\alpha}, a_{j_\beta}]$
8:         **end for**
9:     **end for**
10:     **return** all the constructed candidate domains
11: **end procedure**

---

We first prove that,

**Lemma 6.1** (constructing candidate domains). *Let $\lambda \in (0, 1)$, $\epsilon \in (0, 1/2)$ and $\delta \in (0, 1/10)$. Let $\mathsf{sa}_i$ be a length-$\sqrt{n}$ subarray whose contribution to the optimal solution is at least $\epsilon\sqrt{n}\lambda$, i.e., $\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) \geq \epsilon\sqrt{n}\lambda$. If we uniformly sample $k = 20 \log(1/\delta)/(\lambda\epsilon^2)$ elements $a_{j_1}, a_{j_2}, \ldots, a_{j_k}$ from $\mathsf{sa}_i$, then with probability at least $1 - \delta$, there exists a pair $(\alpha, \beta) \in [k] \times [k]$ such that the following two conditions hold*

1. *$\mathsf{sm}(\mathsf{sa}_i) \leq a_{j_\alpha} \leq a_{j_\beta} \leq \mathsf{lg}(\mathsf{sa}_i)$,*

2. *$\mathsf{lis}^{[a_{j_\alpha}, a_{j_\beta}]}(\mathsf{sa}_i) \geq (1-\epsilon)\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i)$.*

*Proof.* At least $\epsilon\sqrt{n}\lambda$ elements of $\mathsf{sa}_i$ appear in $\mathsf{opt}$. Let us put all these elements in an array $\mathsf{b}$ in the same order that they appear in $\mathsf{sa}_i$. Then it is obvious that $\mathsf{b}$ has at least $\epsilon\sqrt{n}\lambda$ elements. To prove the lemma, we bound the probability that none of the first $\epsilon/2$ fraction of the elements of $\mathsf{b}$

43

are sampled in our algorithm.

$$\Pr[\text{none of the elements in the first } \epsilon/2 \text{ fraction of } \mathsf{b} \text{ is sampled}]$$

$$\leq \left(1 - \frac{\epsilon}{2} \cdot \epsilon\sqrt{n}\lambda \cdot \frac{1}{\sqrt{n}}\right)^k$$

$$= \left(1 - \frac{\epsilon^2\lambda}{2}\right)^{\frac{2}{\epsilon^2\lambda} \cdot 10\log(1/\delta)}$$

$$\leq e^{-10\log(1/\delta)}$$

$$\leq \delta/2,$$

where the first step follows from the fact that $\mathsf{b}$ contains at least $\epsilon\lambda\sqrt{n}$ elements, the second step follows from $k = 20\log(1/\delta)/(\lambda\epsilon^2)$, and the third step follows from the fact that $(1 - 1/x)^x \leq 1/e$ for $\forall x \geq 4$. Hence, with probability at least $1 - \delta/2$, at least one of the elements in the first $\epsilon/2$ fraction of $\mathsf{b}$ are sampled.

With the same analysis, one can prove that with probability at least $1 - \delta/2$ at least one of the elements in the last $\epsilon/2$ fraction of $\mathsf{b}$ are also sampled.

Taking a union bound of two events, with probability at least $1 - \delta$, at least one of the elements in the first and at least one of the elements in the last $\epsilon/2$ fraction of $\mathsf{b}$ are sampled.

Therefore, the $\mathsf{lis}$ of $\mathsf{sa}_i$ subject to this interval is at least a $1 - \epsilon$ fraction of $\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i)$.

$\square$

Notice that the average contribution of each subarray to $\mathsf{opt}$ is $\sqrt{n}\lambda$ and Lemma 6.1 applies to a subarray if its contribution to $\mathsf{opt}$ is at least an $\epsilon$ fraction of this value. Therefore Lemma 6.1 implies that a considerable fraction of the solution is covered by the candidate domains.

**Corollary 6.2** (existence of a desirable solution)**.** *Let $\lambda \in (0,1)$ such that $\mathsf{lis}(A) \geq n\lambda$ and $\epsilon \in (0,1/4)$. If we run Algorithm 8 with parameter $\delta = \epsilon$ on every subarray independently, then with probability at least $1 - \exp(-\Omega(\epsilon^2\sqrt{n}\lambda))$, there exists a set $T \subseteq [\sqrt{n}]$ and elements $\alpha_i$ and $\beta_i$ sampled from $\mathsf{sa}_i$ for each $i \in T$ such that the following conditions hold:*

1. *For any $i \in T$, $\alpha_i \leq \beta_i$.*

2. *For any $i, j \in T$ satisfying $i < j$, $\beta_i < \alpha_j$.*

3. $\sum_{i \in T} \mathsf{lis}^{[\alpha_i,\beta_i]}(\mathsf{sa}_i) \geq (1 - 4\epsilon)\mathsf{lis}(A)$.

*Proof.* Lemma 6.1 holds for all subarrays whose contribution to $\mathsf{opt}$ is at least $\epsilon\sqrt{n}\lambda$. Let $S \subseteq [\sqrt{n}]$ denote the set of coordinates such that for each $i \in S$

$$\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) \geq \epsilon\sqrt{n}\lambda.$$

Since $\sum_{i=1}^{\sqrt{n}} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) = \mathsf{lis}(A)$ and $\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) \leq \sqrt{n}$, we have

$$\sum_{i \in S} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) \geq \sum_{i=1}^{\sqrt{n}} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) - \sqrt{n} \cdot \epsilon\sqrt{n}\lambda \geq \mathsf{lis}(A) - \epsilon n\lambda. \qquad (12)$$

Let $T \subseteq S$ denote the set of coordinates such that for each $i \in T$,

$$\mathsf{lis}^{[\alpha_i,\beta_i]}(\mathsf{sa}_i) \geq (1 - \epsilon)\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i), \quad \text{and} \quad \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i),\mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) \geq \epsilon\sqrt{n}\lambda.$$

44

Now we show that with probability at least $1 - \exp(-\Omega(\epsilon^2\sqrt{n}\lambda))$,

$$\sum_{i \in T} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) \geq (1 - 2\epsilon) \sum_{i \in S} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i). \tag{13}$$

For each $i \in S$, let $X_i$ denote a random variable such that

$$X_i = \begin{cases} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i), & \text{with probability of } i \in T; \\ 0, & \text{with probability of } i \notin T, \end{cases}$$

and $X = \sum_{i \in S} X_i$. By Lemma 6.1 (with $\delta = \epsilon$), We have

$$\mathbf{E}[X] \geq (1 - \epsilon) \sum_{i \in S} \mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i).$$

By Hoeffding bound (Theorem A.2),

$$\Pr[X - \mathbf{E}[X] \geq \epsilon\, \mathbf{E}[X]] \leq 2\exp\left( -\frac{2\epsilon^2(\mathbf{E}[X])^2}{\sum_{i \in S}\left(\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i)\right)^2} \right)$$

$$\leq 2\exp\left( -\frac{2\epsilon^2(1-\epsilon)^2\left(\sum_{i \in S}\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i)\right)^2}{n^{3/2}} \right)$$

$$\leq 2\exp(-2\epsilon^2(1-\epsilon)^4\sqrt{n}\lambda)$$

$$\leq \exp(-\Omega(\epsilon^2\sqrt{n}\lambda)).$$

Hence, Equation (13) holds with probability at least $1 - \exp(-\Omega(\epsilon^2\sqrt{n}\lambda))$.

Conditioned on Equation (13), we have

$$\sum_{i \in T} \mathsf{lis}^{[\alpha_i, a_{\beta_i}]}(\mathsf{sa}_i) \geq \sum_{i \in T}(1-\epsilon)\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i)$$

$$\geq (1-\epsilon)(1-2\epsilon)\sum_{i \in S}\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]}(\mathsf{sa}_i) \tag{14}$$

$$\geq (1-\epsilon)(1-2\epsilon)(\mathsf{lis}(A) - \epsilon n\lambda)$$

$$\geq (1-4\epsilon)\mathsf{lis}(A)$$

where the first inequality follows from the definition of $T$, the second inequality follows from Equation (13), the third inequality follows from Equation (12) and the last inequality follows from $\mathsf{lis}(A) \geq n\lambda$.

Finally, by Equation (13) and (14), we have

$$\Pr\left[\sum_{i \in T} \mathsf{lis}^{[\alpha_i, a_{\beta_i}]}(\mathsf{sa}_i) \geq (1-4\epsilon)\mathsf{lis}(A)\right] \geq 1 - \exp(-\Omega(\epsilon^2\sqrt{n}\lambda)).$$

$\square$

---

**Algorithm 9** Constructing the pseudo solutions

---

1: **procedure** CONSTRUCTPSEUDOSOLUTIONS($\mathsf{cdi}_1, \ldots, \mathsf{cdi}_{\sqrt{n}}$)  ▷ Lemma 6.3,6.4,6.5
2:   ▷ $\{\mathsf{cdi}_i\}_{i \in [\sqrt{n}]}$ is $\sqrt{n}$ sets of candidate domain intervals
3:   pseudo-solutions $\leftarrow \emptyset$
4:   **while true do**
5:     assg $\leftarrow$ largest assigment of candidate domain intervals to subarrays which is monotone
6:     **if** assg contains less than $\epsilon\sqrt{n}\lambda$ non-empty candidate domain intervals **then**
7:       **break**
8:     **else**
9:       Add assg to pseudo-solutions
10:       **for** $i \leftarrow 1$ to $\sqrt{n}$ **do**
11:         **if** assg contains a candidate domain interval for subarray $\mathsf{sa}_i$ **then**
12:           remove the corresponding candidate domain interval from $\mathsf{cdi}_i$
13:         **end if**
14:       **end for**
15:     **end if**
16:   **end while**
17:   **return** pseudo-solutions $\mathsf{ps}_1, \mathsf{ps}_2, \ldots, \mathsf{ps}_t$
18: **end procedure**

---

## 6.2 Constructing Approximately Optimal Pseudo-solutions

We call a sequence of $\sqrt{n}$ intervals $[\ell_1, r_1], [\ell_2, r_2], \ldots, [\ell_{\sqrt{n}}, r_{\sqrt{n}}]$ a pseudo-solution if all of the intervals are monotone. That is $\ell_1 \leq r_1 < \ell_2 \leq r_2 < \ell_3 \leq r_3 \leq \ldots \leq \ell_{\sqrt{n}} \leq r_{\sqrt{n}}$. These intervals denote solution-domains for the subarrays. We also may decide not to assign any solution domain to a subarray in which case we show the corresponding interval by $\emptyset$. We define monotonicity of a pseudo-solution such that it is not violated by $\emptyset$. The quality of a pseudo-solution is defined as $\sum_i \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i)$. We denote the quality of a pseudo-solution $\mathsf{ps}$ by $\mathsf{q}(\mathsf{ps})$.

Another way to interpret Corollary 6.2 is that one can construct a pseudo-solution using the sampled elements whose quality is at least a $1 - \epsilon$ fraction of $\mathsf{lis}(A)$. In this section, we present an algorithm to construct a small set of pseudo-solutions with the promise that at least one of them has a quality of at least $\mathsf{lis}(A)/t$, where $t$ is the number of pseudo-solutions. Finally, in Section 6.3, we present a method to approximate the size of the optimal solution using pseudo-solutions.

We construct the pseudo-solutions via Algorithm 9. The input of Algorithm 9 is the set of candidate domain intervals obtained by Algorithm 8 on every subarray. We first find an assignment of candidate solution domains to the subarrays which is monotone and has the largest number of candidate domain intervals. (This step can be implemented by dynamic programming or solved with an algorithm similar to activity selection algorithm.) We make a pseudo-solution out of this assignment and update the set of candidate intervals by removing the ones which are used in our pseudo-solution. We then repeat the same procedure to construct the second pseudo-solution and update the candidate solution domains accordingly. We continue on, until the number of solution domains used in our pseudo-solution drops below $\epsilon\lambda\sqrt{n}$ in which case we stop.

We first prove in Lemma 6.3 that the number of pseudo-solutions constructed in Algorithm 9 is bounded by $O(k^2/(\lambda\epsilon))$. Next, we show in Lemma 6.4 that at least one of the pseudo-solutions constructed by Algorithm 9 has a quality of at least $\Omega(\mathsf{lis}(A)/t)$ where $t$ is the number of pseudo-solutions. Finally we prove in Lemma 6.5 that the running time of Algorithm 9 is $O(tk^2\sqrt{n}\log n)$.

**Lemma 6.3** (number of pseudo-solutions). *For each $i \in [\sqrt{n}]$, let $\mathsf{cdi}_i$ be a set of at most $k^2$ candidate domain intervals. Let $t$ denote the number of pseudo-solutions constructed in Algorithm 9. Then, we have $t \leq k^2/(\lambda\epsilon)$.*

*Proof.* Note that for each subarray we have at most $k^2$ candidate domain intervals. Since there are $\sqrt{n}$ subarrays, then in total we have $\sqrt{n}k^2$ candidate domain intervals. Each time we construct a pseudo-solution, the total number of the candidate domain intervals is decreased by at least $\epsilon\sqrt{n}\lambda$. Thus, the total number of pseudo-solutions $t$ can be upper bounded,

$$t \leq \frac{\sqrt{n}k^2}{\epsilon\sqrt{n}\lambda} \leq \frac{k^2}{\lambda\epsilon}.$$

□

**Lemma 6.4** (quality of pseudo-solutions). *Let $\mathsf{ps}_1, \mathsf{ps}_2, \ldots, \mathsf{ps}_t$ be the pseudo-solutions constructed by Algorithm 9. If $\mathsf{lis}(A) \geq n\lambda$ holds, then with probability at least $1 - \exp(-\Omega(\sqrt{n}\lambda))$, there exists an $i \in [t]$ such that*

$$\mathsf{q}(\mathsf{ps}_i) \geq \frac{\mathsf{lis}(A)}{2t}.$$

*Proof.* Let us focus again on the actual solution domains

$$[\mathsf{sm}(\mathsf{sa}_1), \mathsf{lg}(\mathsf{sa}_1)], [\mathsf{sm}(\mathsf{sa}_2), \mathsf{lg}(\mathsf{sa}_2)], \ldots, [\mathsf{sm}(\mathsf{sa}_{\sqrt{n}}), \mathsf{lg}(\mathsf{sa}_{\sqrt{n}})].$$

We define set $S \subseteq [\sqrt{n}]$ such that

$$\mathsf{lis}^{[\mathsf{sm}(\mathsf{sa}_i), \mathsf{lg}(\mathsf{sa}_i)]} \geq \epsilon\sqrt{n}\lambda, \forall i \in S.$$

Using Corollary 6.2 with $\epsilon \leq 1/10$, we know that there is a monotone pseudo-solution $[\alpha_i, \beta_i]_{i \in T}$ ($T \subseteq S$) such that $[\alpha_i, \beta_i]$ are candidate domain intervals and

$$\sum_{i \in T} \mathsf{lis}^{[\alpha_i, \beta_i]}(\mathsf{sa}_i) \geq (1 - 4\epsilon)\mathsf{lis}(A).$$

Denote this pseudo-solution as $\mathsf{sol}$. At the time we terminate Algorithm 9, there are at most $\epsilon\sqrt{n}\lambda$ candidate domain intervals of $\mathsf{sol}$ that do not belongs to any pseudo-solution of the pseudo-solution set. Also, since each candidate domain interval contributes at most $\sqrt{n}$ to the quality of the pseudo-solution containing the interval, we have

$$\sum_{i=1}^{t} \mathsf{q}(\mathsf{ps}_i) \geq (1 - 4\epsilon)\mathsf{lis}(A) - \epsilon\sqrt{n}\lambda \cdot \sqrt{n} \geq (1 - 4\epsilon)\mathsf{lis}(A) - \epsilon\mathsf{lis}(A) = (1 - 5\epsilon)\mathsf{lis}(A).$$

Thus, there exists an $i \in [t]$ such that

$$\mathsf{q}(\mathsf{ps}_i) \geq (1 - 5\epsilon)\mathsf{lis}(A)/t \geq \mathsf{lis}(A)/(2t).$$

□

**Lemma 6.5** (running time). *For each $i \in [\sqrt{n}]$, let $\mathsf{cdi}_i$ be a set of at most $k^2$ candidate domain intervals. Let $t$ denote the number of pseudo-solutions. The running time of Algorithm 9 is bounded by $O(tk^2\sqrt{n}\log n)$.*

47

---

**Algorithm 10** Evaluate the pseudo solutions

---

1: **procedure** EVALUATEPSEUDOSOLUTIONS($\mathsf{ps}_1, \ldots, \mathsf{ps}_t$)                    ▷ Lemma 6.6
2:                                                      ▷ $\{\mathsf{ps}_i\}_{i \in [t]}$ is a set of pseudo solutions
3:      $p \leftarrow \frac{1000 t \log^4 n}{\epsilon^4 \lambda \sqrt{n}}$
4:      Randomly sample each $i \in [\sqrt{n}]$ with probability p, and put all the samples in a set $W$
5:      **for** each $\mathsf{ps}_j$ **do**
6:           $\widetilde{\mathsf{q}}(\mathsf{ps}_j) \leftarrow 0$
7:           **for** each interval $[\ell_i, r_i]$ in $\mathsf{ps}_j$ **do**
8:                **if** $i \in W$ **then**
9:                     $\widetilde{\mathsf{q}}(\mathsf{ps}_j) \leftarrow \widetilde{\mathsf{q}}(\mathsf{ps}_j) + \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i)/p$
10:               **end if**
11:          **end for**
12:     **end for**
13:     **return** largest $\widetilde{\mathsf{q}}(\mathsf{ps}_j)$ for all $j \in [t]$
14: **end procedure**

---

*Proof.* Lemma 6.3 states that Algorithm 9 terminates after constructing $t$ pseudo-solutions. Now we show that constructing each pseudo-solution takes time $O(k^2 \sqrt{n} \log n)$. Our solution is based on a dynamic programming technique. Let $D : [\sqrt{n}] \times [k^2] \to \mathbb{N}$ be an array such that $D[i][j]$ stores the size of the largest monotone pseudo-solution for the first $i$ subarrays which ends with the $j$'th candidate domain interval of $\mathsf{sa}_i$. Using classic segment-tree data structure (this data structure can be found in many textbooks, e.g. [CLRS09]), one can compute the value of $D[i][j]$ in time $O(\log n)$ from the previously computed elements of the array.

Thus, the total running is bounded by $O(tk^2 \sqrt{n} \log n)$. $\qquad\qquad\qquad\qquad\qquad\square$

## 6.3    Evaluating the Pseudo-solutions

We finally use a concentration bound to show that the quality of a pseudo-solution can be approximated well by sampling a small number of subarrays. Since a pseudo-solution specifies the range of the numbers used in every subarray, the quality of a pseudo-solution, or in other words, the size of the corresponding increasing subsequence of a pseudo-solution can be formulated as

$$\mathsf{q}(\mathsf{ps}) := \sum_{i=1}^{\sqrt{n}} \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i)$$

where $[\ell_i, r_i]$ denotes the corresponding solution domain of $\mathsf{ps}$ for $\mathsf{sa}_i$.

In Lemma 6.6, we prove that by sampling $O(\log^{O(1)} n/\lambda^4)$ many subarrays and computing $\mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i)$ for them, one can approximate the quality of a pseudo-solution pretty accurately.

**Lemma 6.6** (the quality of pseudo-solution)**.** *Let* $\lambda \in (0, 1)$ *and* $\epsilon$ *be a constant in* $(0, 1/100)$. *Let* $\mathsf{ps}_1, \mathsf{ps}_2, \cdots, \mathsf{ps}_t$ *be a set of t pseudo-solutions. With probability at least* $1 - \exp(-\Omega(\log^2 n))$, *Algorithm 10 runs in time* $O(t^2 \sqrt{n} \log^{O(1)} n/\lambda)$ *such that,*

1. *If there exists an* $i \in [t]$, $\mathsf{q}(\mathsf{ps}_i) \geq \frac{\lambda n}{2t}$, *then the algorithm outputs an estimation at least* $\frac{\lambda n}{4t}$.

2. *If* $\mathsf{q}(\mathsf{ps}_i) < \frac{\lambda n}{8t}$ *for all* $i \in [t]$, *then the algorithm outputs an estimation smaller than* $\frac{\lambda n}{4t}$.

*Proof.* We iterate over all $t$ pseudosolutions, and for each one we estimate their quality separately. In the end we output the largest estimated value over all pseudosolutions. More precisely, we define $p = \frac{1000t \log^4 n}{\epsilon^4 \lambda \sqrt{n}}$ and for each pseudosolution, we sample each of its subarrays with probability $p$. We then compute the LIS of the sampled subarrays subject to the range which corresponds to the pseudosolution. We then estimate the quality of the pseudosolution by scaling the sum of LIS's by a factor $1/p$. In what follows, we prove that the estimation error is negligible.

By Chernoff bound, for each pseudosolution, with probability $1 - \exp(-\Omega(\log^3 n))$ at most $2p\sqrt{n} = \frac{2000t \log^4 n}{\epsilon^4 \lambda}$ subarrays are sampled. For each pseudo-solution, Algorithm 10 computes the value of longest increasing subsequence subject to the corresponding range once for every sampled subarray. Hence, the total running time of the algorithm is $O(t^2 \sqrt{n} \log^{O(1)} n/\lambda)$.

Consider an arbitrary pseudo-solution $\mathsf{ps} \in \{\mathsf{ps}_1, \ldots, \mathsf{ps}_t\}$ and let $\widetilde{\mathsf{q}}(\mathsf{ps})$ denote its estimated LIS. We show that

$$\Pr\left[\left(1 - \frac{\epsilon}{4}\right)^2 \left(\mathsf{q}(\mathsf{ps}) - \frac{\epsilon \lambda^4 n}{100t}\right) \leq \widetilde{\mathsf{q}}(\mathsf{ps}) \leq \left(1 + \frac{\epsilon}{4}\right)^2 \mathsf{q}(\mathsf{ps}) + \frac{\epsilon \lambda^4 n}{100t}\right] \leq 1 - \exp(-\Omega(\log^3 n)). \quad (15)$$

Then the lemma holds by a union bound on all the pseudo-solutions.

Let $T$ be the set of subarray indices such that $i \in T$ iff there is a non-empty interval $[\ell_i, r_i]$ of $\mathsf{ps}$ corresponding to subarray $\mathsf{sa}_i$. Let $p = \frac{1000t \log^4 n}{\epsilon^4 \lambda \sqrt{n}}$ be the probability of sampling a subarray. For each $i \in T$, let $X_i$ denote a random variable such that

$$X_i = \begin{cases} 1, & \text{with prob. } p; \\ 0, & \text{with prob. } 1 - p. \end{cases}$$

and

$$X = \sum_{i \in T} \frac{1}{p} \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i) X_i.$$

We have

$$\mathbf{E}[X] = \mathbf{E}\left[\sum_{i \in T} \frac{1}{p} \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i) X_i\right] = \sum_{i \in T} \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i) = \mathsf{q}(\mathsf{ps}).$$

Let

$$T_j = \left\{i \in T : (1 + \epsilon/4)^{j-1} \leq \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i) < (1 + \epsilon/4)^j\right\}$$

for integer $1 \leq j \in \lceil \log_{1+\epsilon} \sqrt{n} \rceil$.

Let $\Delta = \frac{\epsilon^2 \lambda \sqrt{n}}{1000t \log n}$. Consider $T_j$'s such that $|T_j| \geq \Delta$. The contribution of $\mathsf{q}(\mathsf{ps})$ mostly comes from $T_j$ in the following sense

$$\sum_{j:|T_j| \geq \Delta} \sum_{i \in T_j} \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i) \leq \mathsf{q}(\mathsf{ps}) \quad (16)$$

and

$$\sum_{j:|T_j| \geq \Delta} \sum_{i \in T_j} \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i) = \mathsf{q}(\mathsf{ps}) - \sum_{j:|T_j| < \Delta} \sum_{i \in T_j} \mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i)$$
$$\geq \mathsf{q}(\mathsf{ps}) - \Delta\sqrt{n} \cdot \lceil \log_{1+\epsilon} \sqrt{n} \rceil \quad (17)$$
$$\geq \mathsf{q}(\mathsf{ps}) - \frac{\epsilon \lambda n}{100t}.$$

49

Now we bound the random variable $\sum_{j:|T_j|\geq\Delta}\sum_{i\in T_j}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)X_i/p$. By Chernoff bound, for each $j$ such that $|T_j|\geq\Delta$, we have

$$\Pr\left[\left(1-\frac{\epsilon}{4}\right)p|T_j|\leq\sum_{i\in T_j}X_i\leq\left(1+\frac{\epsilon}{4}\right)p|T_j|\right]\geq1-\exp(-\Omega(\epsilon^2 p|T_j|))=1-\exp(-\Omega(\log^3 n)). \quad (18)$$

Notice that for set $T_j$,

$$\frac{(1+\epsilon/4)^{j-1}}{p}\sum_{i\in T_j}X_i\leq\sum_{i\in T_j}\frac{1}{p}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)X_i\leq\frac{(1+\epsilon/4)^j}{p}\sum_{i\in T_j}X_i. \quad (19)$$

We have

$$\Pr\left[\left(1-\frac{\epsilon}{4}\right)^2\sum_{i\in T_j}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)\leq\sum_{i\in T_j}\frac{1}{p}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)X_i\leq\left(1+\frac{\epsilon}{4}\right)^2\sum_{i\in T_j}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)\right]$$

$$\geq\Pr\left[\left(1-\frac{\epsilon}{4}\right)^2\sum_{i\in T_j}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)\leq\frac{(1+\epsilon/4)^{j-1}}{p}\sum_{i\in T_j}X_i\right.$$

$$\left.\bigwedge\frac{(1+\epsilon/4)^j}{p}\sum_{i\in T_j}X_i\leq\left(1+\frac{\epsilon}{4}\right)^2\sum_{i\in T_j}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)\right]$$

$$\geq\Pr\left[\left(1-\frac{\epsilon}{4}\right)^2\left(1+\frac{\epsilon}{4}\right)^j|T_j|\leq\frac{(1+\epsilon/4)^{j-1}}{p}\sum_{i\in T_j}X_i\right. \qquad\qquad (20)$$

$$\left.\bigwedge\ \frac{(1+\epsilon/4)^j}{p}\sum_{i\in T_j}X_i\leq\left(1+\frac{\epsilon}{4}\right)^{j+1}|T_j|\right]$$

$$\geq\Pr\left[\left(1-\frac{\epsilon}{4}\right)p|T_j|\leq\sum_{i\in T_j}X_i\ \bigwedge\ \sum_{i\in T_j}X_i\leq\left(1+\frac{\epsilon}{4}\right)p|T_j|\right]$$

$$=1-\exp(-\Omega(\log^3 n)),$$

where the first inequality follows from Equation (19), the second inequality follows from the definition of $T_j$ and the last inequality follows from Equation (18). By Equation (16), (17), (21) and union bound, we have

$$\Pr\left[\left(1-\frac{\epsilon}{4}\right)^2\left(\mathsf{q}(\mathsf{ps})-\frac{\epsilon\lambda n}{100t}\right)\leq\sum_{j:|T_j|\geq\Delta}\sum_{i\in T_j}\frac{1}{p}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)X_i\leq\left(1+\frac{\epsilon}{4}\right)^2\mathsf{q}(\mathsf{ps})\right]$$

$$\geq\Pr\left[\left(1-\frac{\epsilon}{4}\right)^2\sum_{j:|T_j|\geq\Delta}\sum_{i\in T_j}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)\leq\sum_{j:|T_j|\geq\Delta}\sum_{i\in T_j}\frac{1}{p}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)X_i\right.$$

$$\qquad\qquad\qquad (21)$$

$$\left.\leq\sum_{j:|T_j|\geq\Delta}\left(1+\frac{\epsilon}{4}\right)^2\sum_{i\in T_j}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)\right]$$

$$\geq1-O\left(\frac{\log n}{\epsilon}\right)\exp(-\Omega(\log^3 n))$$

$$=1-\exp(-\Omega(\log^3 n)).$$

Now we bound random variable $\sum_{j:|T_j|<\Delta} \sum_{i\in T_j} \mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)X_i/p$. If $|T_j|<\Delta$, then by Chernoff bound we have

$$\Pr\left[\sum_{i\in T_j} X_i \leq \frac{2\log^3 n}{\epsilon^2}\right] \geq 1 - \exp(-\Omega(\log^3 n)).$$

By union bound we have

$$\Pr\left[\sum_{j:|T_j|<\Delta}\sum_{i\in T_j}\frac{1}{p}\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i)X_i \leq \frac{2\log^3 n}{\epsilon^2}\cdot\frac{\sqrt{n}}{p}\cdot\lceil\log_{1+\epsilon}\sqrt{n}\rceil < \frac{\epsilon\lambda n}{100t}\right] \geq \tag{22}$$

$$1 - \exp(-\Omega(\log^3 n))). \tag{23}$$

By Equation (21) and Equation (22), we obtain Equation (15).

$\square$

Putting all the previous lemmas together gives the following result:

**Corollary 6.7** (algorithm for LIS decision problem). *Given a length-n sequence A, let $\lambda \in [1/n, 1]$. There is a randomized algorithm that runs in time $O(\lambda^{-7}\sqrt{n}\log^{O(1)} n)$ such that with probability $1 - 1/\operatorname{poly}(n)$*

- *The algorithm accepts if $\mathsf{lis}(A) \geq n\lambda$.*

- *The algorithm rejects if $\mathsf{lis}(A) = O(n\lambda^4)$.*

*Proof.* The correctness follows from Lemma 6.6, Lemma 6.4, and Lemma 6.5.
    **Running time:** The running time is

$$\text{time} = \underbrace{O(tk^2\sqrt{n}\log^{O(1)} n)}_{\text{Lemma 6.5}} + \underbrace{O(t^2\lambda^{-1}\sqrt{n}\log^{O(1)} n)}_{\text{Lemma 6.6}}$$

$$= O(t^2\lambda^{-1}\sqrt{n}\log^{O(1)} n) \qquad\qquad\qquad \text{since } t \leq O(k^2/\lambda)$$

$$= O(k^4\lambda^{-3}\sqrt{n}\log^{O(1)} n) \qquad\qquad\qquad \text{since } k \leq O(1/\lambda)$$

$$= O(\lambda^{-7}\sqrt{n}\log^{O(1)} n)$$

Thus, we complete the proof. $\square$

Finally, by starting with $\lambda = 1$ and iteratively multiplying $\lambda$ by a $1/(1+\epsilon)$ factor until a solution is found, we can approximate $\mathsf{lis}(A)$ within an approximation factor of $O(\lambda^3)$.

**Theorem 6.8** (polynomial approximation for LIS). *Given a length-n sequence A such that $\mathsf{lis}(A) = n\lambda$ where $\lambda \in [1/n, 1]$ is unknown to the algorithm. There is an algorithm that runs in time $\widetilde{O}(\lambda^{-7}\sqrt{n})$ and outputs a number est such that*

$$\Omega(\mathsf{lis}(A)\lambda^3) \leq \mathsf{est} \leq O(\mathsf{lis}(A)).$$

*with probability at least $1 - 1/\operatorname{poly}(n)$.*

We remark that one can turn Theorem 6.8 into an algorithm with running time $\widetilde{O}(n^{17/20})$ by considering two cases separately. If $\lambda < n^{-1/20}$ we sample the array with a rate of $n^{-3/20}$ and compute the LIS for the sampled array. Otherwise, the running time of the algorithm is already bounded by $\widetilde{O}(n^{17/20})$.

## 6.4 An $O(n^\kappa)$ Time Algorithm via Bootstrapping

In this section, we present a general framework to reduce the running time for LIS approximation at the expense of worse approximation factor.

Let us move a step backward and analyze the previous algorithm for obtaining an $O(\lambda^3)$ approximate solution. We first divide the input array into $\sqrt{n}$ subarrays of size $\sqrt{n}$ and after constructing the pseudo-solutions, we sample $O(\lambda^4)$ subarrays to estimate the size of the solution for pseudo-solutions. The reason we set the size of the subarrays to $\sqrt{n}$ is that there is a trade-off between the first and the last steps of the algorithm. More precisely, if we have more than $\sqrt{n}$ subarrays then the number of samples we draw in the beginning would exceed $O_\lambda(\sqrt{n})$. On the other hand, having fewer than $\sqrt{n}$ subarrays results in larger subarrays which would be costly in the last step.

If one favors the running time over the approximation factor, the following improvement can be applied to the algorithm: In the last step of the algorithm, instead of sampling the entire subarrays and computing lis for every pseudo-solution, we recursively call the same procedure to approximate the size of the solution for each subarray. This way, having large subarrays would no longer be an issue and therefore we can have fewer subarrays to improve the number of samples we draw in the first step of the algorithm.

More formally, in order to obtain a running time of $O(\text{poly}(\lambda)n^\kappa)$ for any constant $0 < \kappa < 1$, we set the size of each subarray to $n^{1-\kappa}$ and therefore after constructing the pseudo-solutions, the problem boils down to approximating the solution for $\text{poly}(\lambda)$ many subarrays of length $n^{1-\kappa}$. By running the same algorithm, we would have $n^\kappa$ subarrays of length $n^{1-2\kappa}$ in the second iteration. After $1/\kappa - 1$ iterations, the subarrays are small enough and we can access all their elements in time $O(\text{poly}(\lambda)n^\kappa)$. Of course, this imposes a factor of $\text{poly}(\lambda)$ to the approximation.
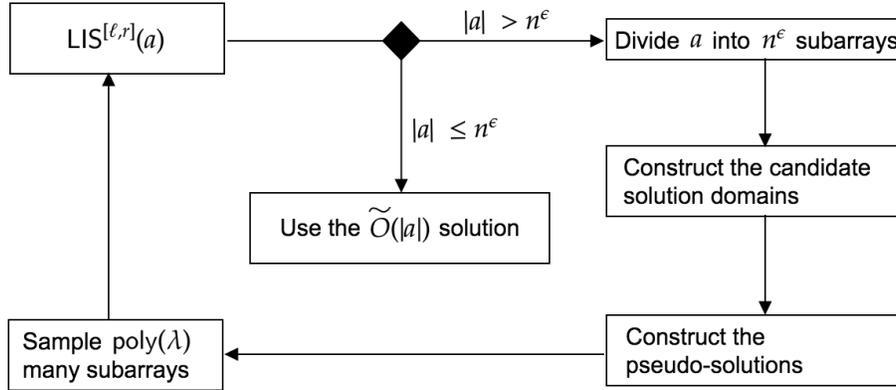


Figure 5: The flowchart of the $O_\lambda(n^\epsilon)$ time algorithm is shown.

By generalizing the ideas from previous subsections, we show that if there is an algorithm for LIS with approximation factor $f(\lambda)$, then we can get a $\left(f\left(\frac{\lambda^4}{2^{32}}\right) \cdot \frac{\lambda^4}{2^{33}}\right)$-approximate LIS algorithm with better running time using the $f(\lambda)$-approximate algorithm as a subroutine.

**Lemma 6.9.** *Assume we partition the sequence into $\zeta$ subarrays, where $\zeta$ is polynomially related to the length of the input sequence. For parameter $\lambda \in (0,1)$, let* ORACLE *be a $f(\lambda)$-approximate* LIS *algorithm (with respect to a domain interval) with running time $g(n,\lambda)$ and success probability $1 - \exp(-\Omega(\log^2 n))$ where $n$ is the length of the input sequence. Then Algorithm 11 using* ORACLE

---

**Algorithm 11** Recursive Estimate LIS with Oracle

---

1: **procedure** RECURSIVEESTIMATIONWITHORACLE(ORACLE, $A, \lambda, \ell, r$)          ▷ Lemma 6.9
2:                                              ▷ input: sequence $A$, parameter $\lambda$, domain interval $[\ell, r]$
3:                                              ▷ assume $\mathsf{sa}_1, \mathsf{sa}_2, \ldots, \mathsf{sa}_{n^\kappa}$ are subarrays of $A$
4:              ▷ subroutine Oracle approximate LIS for subarrays with approximation factor $f(\lambda)$
5:     **for** $i \in [\zeta]$ **do**
6:         $\mathsf{cdi}_i \leftarrow$ CONSTRUCTCANDIDATEDOMAINS($\mathsf{sa}_i$)
7:         discard all the intervals which are not in $[\ell, r]$ from $\mathsf{cdi}_i$
8:     **end for**
9:     $\{\mathsf{ps}_1, \ldots, \mathsf{ps}_t\} \leftarrow$ CONSTRUCTPSEUDOSOLUTIONS($\mathsf{cdi}_1, \ldots, \mathsf{cdi}_\zeta$)
10:    $\lambda_0 \leftarrow \left(\frac{\lambda}{2^8}\right)^4$
11:    $p \leftarrow \frac{20 \log^4 n}{\lambda_0 \zeta}$
12:    randomly sample each $i \in [\zeta]$ with probability $p$, and put all the samples in a set $Q$
13:    **for** $j \in [t]$ **do**
14:        $\mathsf{c} \leftarrow 0$
15:        **for** $i \in W$  **do**
16:            **if** $\exists [\ell_i, r_i] \in \mathsf{ps}_j$ and ORACLE($\mathsf{sa}_i, \lambda_0, \ell_i, r_i$) accepts **then**
17:                $\mathsf{c} \leftarrow \mathsf{c} + 1$
18:            **end if**
19:        **end for**
20:        **if** $\mathsf{c} \geq 3\lambda_0 p\zeta/4$  **then**
21:            **return** accept
22:        **end if**
23:    **end for**
24:    **return** reject
25: **end procedure**

---

as a subroutine is a $\left( f\left(\frac{\lambda^4}{2^{32}}\right) \cdot \frac{\lambda^4}{2^{33}} \right)$-approximate LIS algorithm with

$$O\left( \lambda^{-4} g\left( \frac{n}{\zeta}, \frac{\lambda^4}{2^{32}} \right) \log^{O(1)} n + \lambda^{-7} \zeta \log^{O(1)} n \right)$$

running time and success probability $1 - \exp(-\Omega(\log^2 n))$, where $\zeta$ is the number of subarrays.

*Proof.* We first prove the correctness of the algorithm. Let $A$ be a sequence of length $n$, and $\mathsf{sa}_1, \ldots, \mathsf{sa}_\zeta$ be the subarrays.

Consider the case of $\mathsf{lis}^{[\ell,r]}(A) \geq \lambda n$. By Corollary 6.2 and Lemma 6.4 with $\epsilon = \delta = 1/10$, with probability $1 - \exp(-\Omega(\zeta\lambda))$, there exists a pseudo-solution $\mathsf{ps}_j$ within interval $[\ell, r]$ satisfying

$$\mathsf{q}(\mathsf{ps}_j) \geq \frac{\mathsf{lis}^{[\ell,r]}(A)}{2t} \geq \frac{\mathsf{lis}^{[\ell,r]}(A)\lambda\epsilon}{2k^2} \geq \frac{\mathsf{lis}^{[\ell,r]}(A)\lambda\epsilon}{2 \cdot 20^2 \log^2(1/\delta)/(\lambda^2\epsilon^4)} \geq \frac{\epsilon^5 \lambda^4}{800 \log^2(1/\delta)} n \geq \frac{\lambda^4}{2^{31}} n.$$

Let $\alpha$ denote the number of subarrays $\mathsf{sa}_i$ such that $\mathsf{lis}^{[\ell_i, r_i]}(\mathsf{sa}_i) \geq \lambda_0 n/\zeta$ where $[\ell_i, r_i]$ is the interval for subarray $\mathsf{sa}_i$ in $\mathsf{ps}_j$. We have

$$\alpha \geq \frac{\mathsf{q}(\mathsf{ps}_j) - \lambda^4 n/2^{32}}{n/\zeta} \geq \frac{\lambda^4 \zeta}{2^{32}}.$$

By Chernoff bound, Step 14 to Step 22 of Algorithm 11 accepts on $\mathsf{ps}_j$ with probability at least $1 - \exp(-\Omega(\log^2 n))$.

Consider the case of

$$\mathsf{lis}^{[\ell,r]}(A) \leq f\left(\frac{\lambda^4}{2^{32}}\right)\frac{\lambda^4}{2^{33}}n.$$

Then for any pseudo-solution $\mathsf{ps}_j$, we have

$$\mathsf{q}(\mathsf{ps}_j) \leq f\left(\frac{\lambda^4}{2^{32}}\right)\frac{\lambda^4}{2^{33}}n.$$

Let $\beta$ denote the number of subarrays $\mathsf{sa}_i$ such that $\mathsf{lis}^{[\ell_i,r_i]}(\mathsf{sa}_i) \geq f(\lambda^4/2^{32})n/\zeta$ where $[\ell_i, r_i]$ is the interval for subarray $\mathsf{sa}_i$ in $\mathsf{ps}_j$. We have

$$\beta \leq \frac{\mathsf{q}(\mathsf{ps}_j)}{f(\lambda^4/2^{32})n/\zeta} \leq \frac{\lambda^4}{2^{33}}\zeta.$$

By Chernoff bound, Step 14 to Step 22 of Algorithm 11 do not accept on $\mathsf{ps}_j$ with probability at least $1 - \exp(-\Omega(\log^2 n))$. By union bound, Algorithm 11 rejects with probability at least $1 - \exp(-\Omega(\log^2 n))$.

Hence, Algorithm 11 is a $\left(f\left(\frac{\lambda^4}{2^{32}}\right)\cdot\frac{\lambda^4}{2^{33}}\right)$-approximate algorithm for $\mathsf{LIS}$ of length $n$ with success probability at least $1 - \exp(-\Omega(\log^2 n))$.

Finally, we discuss the running time of Algorithm 11. By the definition of procedures CON-STRUCTCANDIDATEDOMAINS and CONSTRUCTPSEUDOSOLUTIONS, the running time of Step 5 to Step 9 of Algorithm 11 is $O(\lambda^{-9}\zeta\log^{O(1)}n)$. By Lemma 6.3, $t = O(\lambda^{-3})$, and by Chernoff bound with probability $1 - \exp(-\Omega(\log^2 n))$ the size of $Q$ is at most $O(\lambda^{-4}\log^4 n)$. Hence, the running time of Step 12 to Step 24 of Algorithm 11 is $O(\lambda^{-7}g(n/\zeta, \lambda^4/2^{32})\log^4 n)$. $\square$

By definition, we have the following basic fact about approximate ratio.

**Fact 6.10.** *Let $f$ and $f'$ be two functions mapping $(0,1)$ to $(0,1)$ such that $f(\lambda) \geq f'(\lambda)$ for any $\lambda \in (0,1)$. If there is a $f(\lambda)$-approximate $\mathsf{LIS}$ algorithm, then the algorithm is also $f'(\lambda)$-approximate.*

Now we present algorithm to approximate $\mathsf{LIS}$ using $\widetilde{O}(n^\kappa \mathrm{poly}(\lambda^{-1}))$ space by applying the pseudo-solution construction-evaluation framework recursively. In particular, we use the same algorithm on subarrays as an oracle and apply Lemma 6.9 recursively to approximate the entire sequence with slightly worse approximation ratio (compared with approximation ratio of the oracle).

**Lemma 6.11.** *Let $\kappa$ be a constant of $(0,1)$ and $\lambda \in (0,1)$. Algorithm 12 approximates $\mathsf{LIS}$ with approximation ratio*

$$\frac{\lambda^{2\cdot4^{(\lceil 1/\kappa\rceil-1)}}}{256^{3\cdot4^{(\lceil 1/\kappa\rceil-1)}}}$$

*and running time $O(n^\kappa \cdot \lambda^{-4^{O(1/\kappa)}}\log^{O(1)}n)$ and success probability $1 - \exp(-\Omega(\log^3 n))$.*

*Proof.* We first prove the correctness of the algorithm by induction. Without loss of generality, we assume $1/\kappa$ is an integer.

For $i \in \{2, 3, \ldots, 1/\kappa\}$, denote

$$h_i(\lambda) = \frac{\lambda^{2\cdot4^{(i-1)}-4}}{256^{2\cdot4^{(i-1)}+3\cdot4^{(i-2)}-7}}.$$

54

**Algorithm 12** Recursive Estimate LIS

---

1: **procedure** RECURSIVELIS($A, \lambda, \ell, r$)                                        ▷ Lemma 6.11
2:                                             ▷ input: sequence $A$, parameter $\lambda$, domain interval $[\ell, r]$
3:                                             ▷ assume $\mathsf{sa}_1, \mathsf{sa}_2, \ldots, \mathsf{sa}_{n^\kappa}$ are subarrays of $A$
4:     **if** the length of $A$ is greater than $n^{2\kappa}$ **then**
5:         **return** RECURSIVELISWITHORACLE(RECURSIVELIS, $A, \lambda, \ell, r$) with $\zeta = n^\kappa$
6:     **else**
7:         **for** $i \in [n^\kappa]$ **do**
8:             $\mathsf{cdi}_i \leftarrow$ CONSTRUCTCANDIDATEDOMAINS($\mathsf{sa}_i$)
9:             discard all the intervals which are not in $[\ell, r]$ from $\mathsf{cdi}_i$
10:         **end for**
11:         $\{\mathsf{ps}_1, \ldots, \mathsf{ps}_t\} \leftarrow$ CONSTRUCTPSEUDOSOLUTIONS($\mathsf{cdi}_1, \ldots, \mathsf{cdi}_{n^\kappa}$)
12:         **if** EVALUATEPSEUDOSOLUTIONS($\mathsf{ps}_1, \ldots, \mathsf{ps}_t$) $\geq \lambda |A|$ **then**
13:             **return** accept
14:         **else**
15:             **return** reject
16:         **end if**
17:     **end if**
18: **end procedure**

---

We show that Algorithm 12 is $h_i(\lambda)$-approximate if the length of the input sequence is $n^{i \cdot \kappa}$.

If the length of input sequence is $n^{2\kappa}$, then $h_2(\lambda) = \frac{\lambda^4}{2^{32}}$. By Corollary 6.2, Lemma 6.4, and Lemma 6.6, Algorithm 12 is $h_2(\lambda)$-approximate.

In the induction step, for an integer $2 \leq i < 1/\kappa$, we assume Algorithm 12 is $h_i(\lambda)$-approximate for input instance of length $n^{i \cdot \kappa}$. By Lemma 6.9, Algorithm 12 is $\left( h_i\left(\frac{\lambda^4}{2^{32}}\right) \frac{\lambda^4}{2^{33}} \right)$-approximate for input instance of length $n^{(i+1) \cdot \kappa}$. Since

$$
\begin{aligned}
h_i\left(\frac{\lambda^4}{2^{32}}\right) \frac{\lambda^4}{2^{33}} &= \frac{\lambda^{4 \cdot (2 \cdot 4^{(i-1)} - 4)} \cdot \lambda^4}{256^{4 \cdot (2 \cdot 4^{(i-1)} - 4)} \cdot 256^{2 \cdot 4^{(i-1)} + 3 \cdot 4^{(i-2)} - 7} \cdot 2^{33}} \\
&= \frac{\lambda^{2 \cdot 4^i - 12}}{256^{2 \cdot 4^i + 2 \cdot 4^{(i-1)} + 3 \cdot 4^{(i-2)} - 18.875}} \\
&> \frac{\lambda^{2 \cdot 4^i - 4}}{256^{2 \cdot 4^i + 3 \cdot 4^{(i-1)} - 7}} \\
&= h_{i+1}(\lambda),
\end{aligned}
$$

by Fact 6.10, Algorithm 12 is $h_{i+1}(\lambda)$-approximate for input instance of length $n^{(i+1) \cdot \kappa}$. Since

$$
h_{1/\kappa}(\lambda) > \frac{\lambda^{2 \cdot 4^{((1/\kappa) - 1)}}}{256^{3 \cdot 4^{((1/\kappa) - 1)}}},
$$

by Fact 6.10, Algorithm 12 is $\left( \frac{\lambda^{2 \cdot 4^{((1/\kappa) - 1)}}}{256^{3 \cdot 4^{((1/\kappa) - 1)}}} \right)$-approximate for input instance of length $n$.

By Corollary 6.2, Lemma 6.4, Lemma 6.6 and Lemma 6.9, we have the desired running time. The success probability is obtained by same corollaries/lemmas and union bound. $\qquad\square$

Finally, by starting with $\lambda = 1$ and iteratively multiplying $\lambda$ by a $1/(1 + \epsilon)$ factor until a solution is found, we can approximate $\mathsf{lis}(A)$ within an approximation factor of $\lambda^{O(4^{1/\kappa})}$.

**Theorem 6.12.** *Let $\kappa$ be a constant of $(0,1)$ and $\lambda \in (0,1)$. There exists a $\widetilde{O}(n^\kappa \cdot \lambda^{-4^{O(1/\kappa)}})$ time algorithm for* lis *with approximation factor $\lambda^{4^{O(1/\kappa)}}$ and success probability $1 - \exp(-\Omega(\log^3 n))$ .*

# 7 Acknowledgement

# References

[AB17]      Amir Abboud and Arturs Backurs. Towards hardness of approximation for polynomial time problems. In *Proceedings of the Eighth Innovations in Theoretical Computer Science Conference*, 2017.

[ABW15]    Amir Abboud, Arturs Backurs, and Virginia Vassilevska Williams. Tight hardness results for LCS and other sequence similarity measures. In *Proceedings of the Fifty-Sixth IEEE Annual Symposium on Foundations of Computer Science*, 2015.

[ACCL07]   Nir Ailon, Bernard Chazelle, Seshadhri Comandur, and Ding Liu. Estimating the distance to a monotone function. *Random Structures & Algorithms*, 31(3):371–383, 2007.

[AHWW16] Amir Abboud, Thomas Dueholm Hansen, Virginia Vassilevska Williams, and Ryan Williams. Simulating branching programs with edit distance and friends. In *Proceedings of the Forth-Eighth Annual ACM SIGACT Symposium on Theory of Computing*, 2016.

[AKO10]    Alexandr Andoni, Robert Krauthgamer, and Krzysztof Onak. Polylogarithmic approximation for edit distance and the asymmetric query complexity. In *Proceedings of the Fifty-First IEEE Annual Symposium on Foundations of Computer Science*, 2010.

[AO09]     Alexandr Andoni and Krzysztof Onak. Approximating edit distance in near-linear time. In *Proceedings of the Forty-First Annual ACM Symposium on Theory of Computing*, 2009.

[AR18]     Amir Abboud and Aviad Rubinstein. Fast and deterministic constant factor approximation algorithms for LCS imply new circuit lower bounds. In *Proceedings of the Ninth Innovations in Theoretical Computer Science Conference*, 2018.

[BBK13]    Pavle VM Blagojević, Boris Bukh, and Roman Karasev. Turán numbers for $k_{s,t}$-free graphs: Topological obstructions and algebraic constructions. *Israel Journal of Mathematics*, 197(1):199–214, 2013.

[BEG+18]   Mahdi Boroujeni, Soheil Ehsani, Mohammad Ghodsi, MohammadTaghi HajiAghayi, and Saeed Seddighin. Approximating edit distance in truly subquadratic time: Quantum and mapreduce. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, 2018.

[BES06]    Tuğkan Batu, Funda Ergun, and Cenk Sahinalp. Oblivious string embeddings and edit distance approximations. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithm*, 2006.

[BI15]      Arturs Backurs and Piotr Indyk. Edit distance cannot be computed in strongly sub-quadratic time (unless SETH is false). In *Proceedings of the Forty-Seventh Annual ACM Symposium on Theory of Computing*, 2015.

[BR65]      George R Blakley and Prabir Roy. A hölder type inequality for symmetric matrices with nonnegative entries. *Proceedings of the American Mathematical Society*, 16(6):1244–1245, 1965.

[BYJKK04]  Ziv Bar-Yossef, TS Jayram, Robert Krauthgamer, and Ravi Kumar. Approximating edit distance efficiently. In *Proceedings of the Forty-Fifth Annual IEEE Symposium on Foundations of Computer Science*, 2004.

[CDG+18]   Diptarka Charkraborty, Debarati Das, Elazar Goldenberg, Michal Koucky, and Michael Saks. Approximating edit distance within constant factor in truly sub-quadratic time. In *Proceedings of the Fifty-Ninth Annual IEEE Symposium on Foundations of Computer Science*, 2018.

[CGL+19]   Lijie Chen, Shafi Goldwasser, Kaifeng Lyu, Guy N Rothblum, and Aviad Rubinstein. Fine-grained complexity meets IP=PSPACE. *Proceedings of the Thirtieth Annual ACM-SIAM symposium on Discrete Algorithm*, 2019.

[Che52]     Herman Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pages 493–507, 1952.

[CLRS09]   Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to Algorithms*. MIT press, 2009.

[DGL+99]   Yevgeniy Dodis, Oded Goldreich, Eric Lehman, Sofya Raskhodnikova, Dana Ron, and Alex Samorodnitsky. Improved testing algorithms for monotonicity. In *Randomization, Approximation, and Combinatorial Optimization. Algorithms and Techniques*, pages 97–108. Springer, 1999.

[EKK+00]   Funda Ergün, Sampath Kannan, S Ravi Kumar, Ronitt Rubinfeld, and Mahesh Viswanathan. Spot-checkers. *Journal of Computer and System Sciences*, 60(3):717–751, 2000.

[Fis04]     Eldar Fischer. The art of uninformed decisions: A primer to property testing. *Current Trends in Theoretical Computer Science: The Challenge of the New Century*, 1:229–264, 2004.

[Fre75]     Michael L Fredman. On computing the length of longest increasing subsequences. *Discrete Mathematics*, 11(1):29–35, 1975.

[Hoe63]     Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.

[HSSS19]   MohammadTaghi Hajiaghayi, Masoud Seddighin, Saeed Seddighin, and Xiaorui Sun. Approximating LCS in linear time: Beating the barrier. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, 2019.

[KST54]     Tamás Kovári, Vera Sós, and Pál Turán. On a problem of k. zarankiewicz. In *Colloquium Mathematicum*, volume 1, pages 50–57, 1954.

[KSV13]   Peter Keevash, Benny Sudakov, and Jacques Verstraëte. On a conjecture of erdős and simonovits: Even cycles. *Combinatorica*, 33(6):699–732, 2013.

[LMS98]   Gad M Landau, Eugene W Myers, and Jeanette P Schmidt. Incremental string comparison. *SIAM Journal on Computing*, 27(2):557–582, 1998.

[MP80]   William J Masek and Michael S Paterson. A faster algorithm computing string edit distances. *Journal of Computer and System Sciences*, 20(1):18–31, 1980.

[RS20]   Aviad Rubinstein and Zhao Song. Reducing approximate longest common subsequence to approximate edit distance. In *SODA*, 2020.

[Sch61]   Craige Schensted. Longest increasing and decreasing subsequences. *Canadian Journal of Mathematics*, 13:179–191, 1961.

[SS10]   Michael Saks and C Seshadhri. Estimating the longest increasing sequence in polylogarithmic time. In *Proceedings of the Fifty-First Annual IEEE Symposium on Foundations of Computer Science*, 2010.

[Tur41]   Paul Turán. On an external problem in graph theory. *Mat. Fiz. Lapok*, 48:436–452, 1941.

[Wes01]   Douglas Brent West. *Introduction to graph theory*, volume 2. Prentice hall Upper Saddle River, 2001.

# A    Probability, combinatorial, and Graph Tools

In this section, we restate probability and graph tools that we use throughout this paper. All these theorems are proven in previous work.

**Theorem A.1** (Chernoff Bounds [Che52]). *Let $X = \sum_{i=1}^{n} X_i$, where $X_i = 1$ with probability $p_i$ and $X_i = 0$ with probability $1 - p_i$, and all $X_i$ are independent. Let $\mu = \mathbf{E}[X] = \sum_{i=1}^{n} p_i$. Then*
*1. $\Pr[X \geq (1+\delta)\mu] \leq \exp(-\delta^2\mu/3)$, $\forall \delta > 0$ ;*
*2. $\Pr[X \leq (1-\delta)\mu] \leq \exp(-\delta^2\mu/2)$, $\forall 0 < \delta < 1$.*

**Theorem A.2** (Hoeffding bound [Hoe63]). *Let $X_1, \cdots, X_n$ denote $n$ independent bounded variables in $[a_i, b_i]$. Let $X = \sum_{i=1}^{n} X_i$, then we have*

$$\Pr[|X - \mathbf{E}[X]| \geq t] \leq 2\exp\left(-\frac{2t^2}{\sum_{i=1}^{n}(b_i - a_i)^2}\right)$$

**Theorem A.3** (Blakley-Roy inequality, [BR65], see also Proposition 3.1 in [KSV13]). *Let $G$ denote a graph that has $n$ vertices and average degree $d$. The number of walks of length $k$ in graph $G$ is at least $nd^k$.*

**Theorem A.4** (Turán theorem for bipartite graphs, [KST54], see also [BBK13]). *For a graph $G$ the Turán number $\mathrm{ex}(G, n)$ is the maximum number of edges that a graph on $n$ vertices can have without containing a copy of $G$. For any $s \leq t$, $\mathrm{ex}(K_{s,t}, n) \leq \frac{1}{2}(t-1)^{1/s}n^{2-1/s} + o(n^{2-1/s})$*

In particular, when $s = t$, Theorem A.5 implies that for large enough $n$ we have $\mathrm{ex}(K_{s,s}, n) \leq n^{2-1/s}$.

**Theorem A.5** (Turán theorem for cliques [Tur41]). *Let $G$ be any graph with $n$ vertices, such that $G$ is $K_{r+1}$-free. Then the number of edges in $G$ is at most*

$$\left(1 - \frac{1}{r}\right) \cdot \frac{n^2}{2}.$$

**Corollary A.6** (of Theorem A.5). *Let $G$ be any graph with $n$ vertices, such that $G$ has no independent set of size $r + 1$. Then the number of edges in $G$ is at least*

$$\binom{n}{2} - \left(1 - \frac{1}{r}\right) \cdot \frac{n^2}{2} = \frac{nr + n^2}{2r}.$$

**Lemma A.7** (application of Jensen's inequality). *Let $n_1, n_2, \ldots, n_k$ be a sequence of integer numbers of size $k$ and $\lambda_1, \lambda_2, \ldots, \lambda_k$ be $k$ real numbers in the interval $[0, 1]$. Define $\lambda = (\sum n_i \lambda_i)/(\sum n_i)$. Then for any $y \geq 0$ we have*

$$\sum n_i \lambda_i^{1+y} \geq \sum n_i \lambda^{1+y}.$$

*Proof.* Let $\psi(x) = x^{1+y}$. Since $\psi$ is real convex function, by Jensen's inequality, we have

$$\psi\left(\frac{\sum_i n_i \lambda_i}{\sum_i n_i}\right) \leq \frac{\sum n_i \psi(\lambda_i)}{\sum n_i}$$

Applying definition $\psi$, we have

$$\left(\frac{\sum_i n_i \lambda_i}{\sum_i n_i}\right)^{1+y} \leq \frac{\sum n_i \lambda_i^{1+y}}{\sum n_i}.$$

Using definition of $\lambda$, we have

$$\lambda^{1+y} \leq \frac{\sum n_i \lambda_i^{1+y}}{\sum n_i}.$$

$\square$

**Theorem A.8** (A well-known algorithm also used in [HSSS19])**.** *Given two strings $A$ and $B$, one can with preprocessing time $O(|B| \log |B|)$ verify if the $\mathsf{LCS}$ of $A$ and $B$ is at least $q$ or not in time $O(|A| q \log |B|)$. In case the answer is positive, finding such a common subsequence can be done in time $O(|A| q \log |B|)$.*

*Proof.* In the preprocessing step, for each character $x$ in $B$, we construct a binary tree that keeps track of all the places that $x$ appears in $B$. This enables us to answer the following queries in $O(\log |B|)$ time: *Given an index $i$ of $B$ and a character $x$, what is the smallest index $i' \geq i$ such that $B_{i'} = x$?*

To find the $\mathsf{LCS}$ of $A$ and $B$, we slightly modify the conventional dynamic program for computing $\mathsf{LCS}$ and construct a two-dimensional array $T^*$ that stores the following information

$$T^*[i][j] = \begin{cases} \text{the smallest } k \text{ s.t.} & \text{if } |\mathsf{lcs}(A[1,i],B)| \geq j \\ |\mathsf{lcs}(A[1,i],B[1,k])| = j \\ \infty & \text{otherwise} \end{cases}$$

Using the above definition, we can construct table $T^*$ via the following recursive formula:

$$T^*[i][j] := \min \left\{ \begin{matrix} T^*[i-1][j], \\ f(T^*[i-1][j-1]+1, A_i) \end{matrix} \right\} \tag{24}$$

where $f(T^*[i-1][j-1]+1, A_i)$ is the index of the first occurrence of $A_i$ in $B$ after position $T^*[i-1][j-1]$ (or $\infty$ if $A_i^*$ does not appear in $B$ after position $T^*[i-1][j-1]$). Since such queries can be answered in $O(\log |B|)$ time then the overall runtime of the algorithm is $O(|A| q \log |B|)$. $\square$