

# Locally Consistent Decomposition of Strings with Applications to Edit Distance Sketching

Sudatta Bhattacharya\* sudatta@iuuk.mff.cuni.cz Computer Science Institute of Charles University Prague, Czech Republic Michal Koucký<sup>†</sup> koucky@iuuk.mff.cuni.cz Computer Science Institute of Charles University Prague, Czech Republic

# ABSTRACT

In this paper we provide a new locally consistent decomposition of strings. Each string x is decomposed into blocks that can be described by grammars of size  $\tilde{O}(k)$  (using some amount of randomness). If we take two strings x and y of edit distance at most k then their block decomposition uses the same number of grammars and the *i*-th grammar of x is the same as the *i*-th grammar of y except for at most k indexes *i*. The edit distance of x and y equals to the sum of edit distances of pairs of blocks where x and y differ. Our decomposition can be used to design a sketch of size  $\tilde{O}(k^2)$ for edit distance, and also a rolling sketch for edit distance of size  $\tilde{O}(k^2)$ . The rolling sketch allows to update the sketched string by appending a symbol or removing a symbol from the beginning of the string.

# **CCS CONCEPTS**

- Theory of computation  $\rightarrow$  Sketching and sampling.

# **KEYWORDS**

Edit distance, sketching, string decomposition, locally consistent parsing

#### ACM Reference Format:

Sudatta Bhattacharya and Michal Koucký. 2023. Locally Consistent Decomposition of Strings with Applications to Edit Distance Sketching. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing (STOC '23), June 20–23, 2023, Orlando, FL, USA*. ACM, New York, NY, USA, 14 pages. https://doi.org/10.1145/3564246.3585239

# **1 INTRODUCTION**

Edit distance is a measure of similarity of two strings. It measures how many symbols one has to insert, delete or substitute in a string x to get a string y. The measure has many applications from text processing to bioinformatics. The edit distance ED(x, y) of two strings x and y can be computed in time  $O(n^2)$  by a classic dynamic

<sup>&</sup>lt;sup>†</sup>Partially supported by the Grant Agency of the Czech Republic under the grant agreement no. 19-27871X. This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 823748 (H2020-MSCA-RISE project CoSP).



This work is licensed under a Creative Commons Attribution 4.0 International License.

STOC '23, June 20–23, 2023, Orlando, FL, USA © 2023 Copyright held by the owner/author(s). ACM ISBN 978-1-4503-9913-5/23/06. https://doi.org/10.1145/3564246.3585239 programming algorithm [29]. Save for poly-log improvements in the running time [15, 25], the best known running time for edit distance computation is  $O(n+k^2)$  [22], where k = ED(x, y). Assuming Strong Exponential Time Hypothesis (SETH) this running time cannot be substantially improved [2]. The conditional lower bound does not exclude some approximation algorithms, though, and there was a recent progress on computing edit distance in almost-linear time to within some constant factor approximation [1, 7, 8, 20].

Another problem for edit distance that saw a major progress in recent years is sketching. In sketching we want to map a string x to a short sketch  $sk_{n,k}^{\text{ED}}(x)$  so that from sketches  $sk_{n,k}^{\text{ED}}(x)$  and  $sk_{n,k}^{\text{ED}}(y)$  of two strings x and y we can compute their edit distance, either exactly or approximately. Apriori it is not even obvious that short sketches for edit distance exist. In a surprising construction, Belazzougui and Zhang [4] gave an exact edit distance sketch of size  $O(k^8 \log^5 n)$  bits. The sketch size was then improved to  $O(k^3 \log^2(\frac{n}{\delta}) \log n)$  bits by Jin, Nelson and Wu [16], where the ED(x, y) was computed exactly from the sketches with probability at least  $1 - \delta$ , if  $\text{ED}(x, y) \leq k$ . The current best sketch is of size  $O(k^2 \log^3 n)$  bits and was given by Kociumaka, Porat and Starikovskaya [19]. [16] gives a lower bound  $\Omega(k)$  on the size of a sketch for exact edit distance.

The major problem in edit distance computation as well as in sketching is how to align the matching parts of two strings x and y. Finding an optimal alignment of two strings is the crux in the computation of edit distance and its sketching. In sketching finding a good alignment is even more challenging as we do not have both strings in our hands simultaneously to look for the matching. To the best of our knowledge, to resolve this issue all edit distance sketches use *CGK random walk* on strings [9] which allows to embed the edit distance metrics into Hamming distance metrics with distortion O(k). The walk implicitly fixes some reasonably good matching between the two strings. Going from the CGK random walk to a sketch is non-trivial undertaking and all three sketch results rely on sophisticated machinery to achieve it.

In this paper we provide a new technique to align two strings x and y in oblivious manner. In nutshell, we provide a decomposition procedure that breaks x and y into the same number of "short" blocks so that at most k pairs of blocks in the decomposition of x and y differ, and all other pairs of blocks are matching in an optimal alignment. So the edit distance of x and y is the sum of edit distances of the differing blocks. To be more specific our blocks are not short in their length but they are short in the sense that each of them can be described by a context-free grammar of size  $\tilde{O}(k)$ . Our decomposition algorithm constructs the grammars. Our decomposition is based on *locally consistent parsing* of strings a

<sup>\*</sup>Partially supported by the Grant Agency of the Czech Republic under the grant agreement no. 19-27871X.

technique similar to the one used in [3, 6, 17, 28] and hash based partitioning similar to [30]. Our main technical result is:

THEOREM 1.1 (STRING DECOMPOSITION). There is an algorithm running in time O(nk) that for each string x of length at most n produces grammars  $G_1^x, \ldots, G_s^x$  such that with probability at least  $1 - O(1/\sqrt{n}), x = eval(G_1^x) \cdots eval(G_s^x)$  and each of the grammars is of size  $\widetilde{O}(k)$ . Furthermore, for any two strings x and y of edit distance at most k with grammars  $G_1^x, \ldots, G_s^x$  and  $G_1^y, \ldots, G_{s'}^y$ , resp., that are produced by the algorithm using the same randomness, the following is true simultaneously with probability at least 4/5:

(1) 
$$s = s'$$
,

- (2)  $G_i^x = G_i^y$ , for all  $i \in \{1, ..., s\}$  except for at most k indices i, and (3)  $ED(x, y) = \sum_i ED(eval(G_i^x), eval(G_i^y)).$

Here, for a grammar G, eval(G) denotes its evaluation. Our decomposition can be used immediately to give an embedding of edit distance into Hamming distance with distortion O(k). It also readily yields a sketch for exact edit distance of size  $O(k^2)$ :

THEOREM 1.2 (SKETCH FOR EDIT DISTANCE). There is a randomized sketching algorithm  $sk_{n,k}^{ED}$  that on an input string x of length at most *n* produces a sketch sk<sup>ED</sup><sub>n,k</sub>(x) of size  $\widetilde{O}(k^2)$  in time  $\widetilde{O}(nk)$ , and a comparison algorithm running in time  $\tilde{O}(k^2)$  such that given two sketches  $\mathrm{sk}_{n,k}^{\mathrm{ED}}(x)$  and  $\mathrm{sk}_{n,k}^{\mathrm{ED}}(y)$  for two strings x and y of length at most n obtained using the same randomness of the sketching algorithm outputs with probability at least 1 - 1/n (over the randomness of the sketching and comparison algorithms) the edit distance of x and y if it is less than k and  $\infty$  otherwise.

Furthermore, we can also provide a rolling sketch, a sketch in which we can update the stored string by appending a symbol or removing its first symbol.

THEOREM 1.3 (ROLLING SKETCH FOR EDIT DISTANCE). There are algorithms Append( $sk_x$ , a), Remove( $sk_{ax}$ , a), and Compare( $sk_x$ ,  $sk_y$ ) such that for integer parameters  $k \leq m$ :

- (1) Given a sketch  $sk_x$  representing a string x and a symbol a, Append $(sk_x, a)$  outputs a sketch  $sk_{xa}$  for the string xa in time  $\widetilde{O}(k^2)$ .
- (2) Given a sketch  $sk_{ax}$  representing a string ax for a symbol a, Remove( $sk_{ax}$ , a) outputs a sketch  $sk_x$  for the string x in time  $O(k^2)$ .
- (3) Given two sketches  $sk_x$  and  $sk_y$  representing strings x and y obtained from the same random sketch for empty string using two sequences of at most m operations Append and Remove, Compare $(sk_x, sk_y)$  calculates the edit distance of x and y if it is less than k, and outputs  $\infty$  otherwise. The algorithm Compare $(sk_x, sk_y)$  runs in time  $\widetilde{O}(k^2)$ .

All the sketches are of size  $\tilde{O}(k^2)$ . The probability that any of the algorithms fails or produces incorrect output is at most 1/m over the initial randomness of the sketch for empty string and internal randomness of the algorithms.

We remark that we did not attempt to optimize the running time of either of our algorithms, or poly-log factors in the sketch sizes, and we believe that both parameters can be readily improved by

Sudatta Bhattacharya and Michal Koucký

usual amortization techniques of processing symbols in batches of size  $\tilde{O}(k)$ . We believe that building the sketch in the first theorem can be done in time O(n) using fast multi-point polynomial evaluation for  $\widetilde{O}(k)$ -wise independent hash functions, the update time in the last theorem can be improved to  $\widetilde{O}(1)$  by buffering  $\widetilde{O}(k)$ symbols that shall be inserted or removed without affecting the other parameters of the algorithm.

Another distinguishing feature of our decomposition procedure compared to the technique of CGK random walks is its parallelizability. CGK random walk seems inherently sequential whereas our decomposition procedure can be easily parallelized. We believe that our decomposition will allow for further applications beyond our simple sketches.

### 1.1 Related Work

The problem of embedding edit distance to other distance measures, like Hamming distance,  $\ell_1$ , etc. has been studied extensively. In [9], the authors have given a randomized embedding from edit distance to Hamming distance, where any string  $x \in \{0, 1\}^n$  can be mapped to a string  $f(x) \in \{0,1\}^{3n}$ , given a random string  $r \in \{0,1\}^{\log^2 n}$ , such that,  $ED(x, y)/2 \le Ham(f(x), f(y)) \le O(ED(x, y)^2)$  with probability at least 2/3. Batu, Ergun and Sahinalp [3] have introduced a dimensionality reduction technique, where any string x of length n can be mapped to a string f(x) of length at most n/r, for any parameter r, with a distortion of O(r). They used the locally consistent parsing technique for their embedding. Ostrovsky and Rabani [26] gave an embedding from edit distance to  $\ell_1$  distance with a distortion of  $O(\sqrt{\log n \log \log n})$ . Jowhari [17] also gave a randomized embedding from edit distance to  $\ell_1$  distance with a distortion of  $O(\log n \log^* n)$ . He used the embedding given by Cormode and Muthukrishnan [12] who showed that any string x of length n can be mapped to a vector f(x) of length  $m = O(2^{n \log n})$ , such that for any pair of strings x, y of length neach,  $\operatorname{ED}(x, y)/2 \le ||f(x) - f(y)||_{\ell_1} \le O(\log n \log^* n) \cdot \operatorname{ED}(x, y).$ Since the size of the vector was too large, [17] used random hashing to get his final embedding.

#### **Our Techniques** 1.2

We first provide the intuition for our technique. We would like to break a string x into small blocks obliviously so that when a string y is broken by the same procedure, the difference between x and ycaused by the edit operations is confined within the corresponding blocks of x and y, and the overall decomposition is not affected by them. For random binary strings x and y this could be done fairly easily: look on all the (overlapping) windows of log n consecutive bits in each of the strings and for each window decide at random whether to make a break at that window or not. To make it consistent between x and y use some random hash function  $H: \{0,1\}^{\log n} \to \{0,\ldots,D-1\}$  so that if the hash function evaluates to 0 on a given window then start a next block of the decomposition. If we chose *D* suitably, say  $D \ge 10k \log n$ , then we are unlikely to start a new block in any window which is affected by the the at most k edit operations on x and y. In that case we obtain the desired decomposition. Hence, decomposing random strings x and y is easy.

The issue is what to do with non-random strings. Consider for example strings *x* and *y* that are very sparse, so they contain  $\sqrt{n}$  ones sprinkled within the vast ocean of zeros. The hash function *H* will see mostly windows of 0's and occasionally a window of the form  $0^i 10^{\log(n)-i-1}$ . The decomposition will have no effect on such strings despite the fact that the string might contain  $\Omega(\sqrt{n})$  bits of entropy.

However, we can compress such sparse strings: replace stretches of zeros by some binary encoded information about their length, and try to break the strings again. Still, this will fail if in our example the stretches of zeros are replaced by stretches of some repeated pattern such as  $(01)^*$ . So we need slightly more general compression which will compress any log *n* bits into  $\log(n)/2$  bits. By repeating the sequence of steps: split and compress, we will eventually get the desired decomposition of each string.

Our actual algorithm mimics the above intuition. It is technically easier to work with a larger alphabet, so we extend the input alphabet  $\Sigma$  by adding special compression symbols into the work alphabet  $\Gamma$ . (Without loss of generalization we can assume that  $\Sigma$ is of size  $O(n^3)$  otherwise we can hash each symbol of our input strings using some perfect hash function into an alphabet of size  $O(n^3)$  without affecting the edit distance of a given pair of strings.) To split a string we will use a random  $\widetilde{O}(k)$ -wise independent hash function  $H : \Gamma^2 \rightarrow \{0, \ldots, D-1\}$ , for  $D = \Theta(k \log n)$ . If the hash function is zero on a pair of consecutive symbols in a string, we start a new block of the decomposition on the first symbol in the pair.

Then in each resulting block we replace stretches of repeated symbols by a special compression symbol from  $\Gamma$  representing the block, and we use a pair-wise independent hash function  $C: \Gamma^2 \rightarrow$  $(\Gamma \setminus \Sigma)$  to compress non-overlapping pairs of symbols into one symbol. This latter step requires some care as we have to make sure that we select non-overlapping pairs in the same way in *x* and *y*. For the selection of non-overlapping pairs we use the locally consistent coloring of Cole and Vishkin [11, 23, 24] where the selection of pairs depends only on the context of  $O(\log^* n)$  symbols. The compression reduces the size of each block by a factor of 2/3. We repeat the compress and split process for  $O(\log n)$  iterations until each compressed block of *x* is of size at most 2. *Decompression* of each block then gives us the desired decomposition of *x*. (See Fig. 1 for an illustration.)

It is natural and convenient to represent each of the blocks by a context-free grammar which corresponds to the compression process. We can argue that the grammars will be of size  $O(D \log n)$ with high probability. So we can represent each string by a sequence of small grammars so that if x and y are at edit distance at most kthen at most k pairs of their grammars will differ, and the sum of the edit distances of differing pairs is the edit distance of x and y. Note, that edit distance of two strings represented by context-free grammars can be computed efficiently [14]. These are the main ideas behind our decomposition algorithm, and we provide more details in Section 3

Building a sketch from the string decomposition is straightforward: We encode each grammar in binary using fixed number of bits, and we use off-the-shelf sketch for Hamming distance to sketch the sequence of grammars. As the Hamming distance sketch does not recover identical bits but only the mismatched bits we make sure that if two grammars differ then their binary encoding differ in every bit. Over binary alphabet this might be impossible but over large alphabets one could use error-correcting codes to achieve the desired effect of recovering the differing grammars; for simplicity we use the Karp-Rabin fingerprint of the whole grammar to encode the binary 0 and 1 distinctly. See Section 3.3 for the details of our encoding and Section 3.4 for details of the sketch for edit distance.

To design a rolling sketch for edit distance where we can extend the represented string by a new symbol or repeatedly remove the first symbol of the represented string we will employ our decomposition technique together with the rolling sketch for Hamming distance of Clifford, Kociumaka, and Porat [10]. We will argue that appending a new symbol to a string affects only some fixed number of grammars in the decomposition of a string. There is a certain threshold T so that except for the last T grammars the decomposition of a string stays the same regardless of how many other symbols are appended. Hence, we will keep a buffer of at most Tactive grammars corresponding to the recently added symbols, and upon addition of a new symbol we will only update those grammars. We are guaranteed that the grammars before this threshold will stay the same forever, so we can commit them into the rolling Hamming sketch (in the form of their binary encoding.) Similarly, we will keep a buffer of up-to T active grammars that capture the symbols that were deleted from the sketch most recently. Once they become "mature" enough we can commit them by removing their binary encoding from the rolling Hamming sketch. (See Fig. 3 for an illustration.) This allows to maintain a rolling sketch for edit distance.

Evaluation of an edit distance query on two rolling sketches will use their Hamming sketch to recover differing committed grammars. Together with the active grammars of inserted and deleted symbols this provides enough information for evaluating the edit distance query. Technical details are explained in Section 4. In Section 5 we give a table of parameters used throughout the paper.

#### 2 NOTATIONS AND PRELIMINARIES

For any string  $x = x_1x_2x_2...x_n$  and integers p, q, x[p] denotes  $x_p, x[p, q]$  represents substring  $x' = x_p...x_q$  of x, and x[p,q) = x[p,q-1]. If q < p, then x[p,q] is the empty string  $\varepsilon$ . x[p,...] represents x[p,|x|], where |x| is the length of x. "."-operator is used to denote concatenation, e.g  $x \cdot y$  is the concatenation of two strings x and y. Dict $(x) = \{x[i, i+1], i \in [n-1]\}$ , is the dictionary of string x, which stores all pairs of consecutive symbols that appear in x. For strings x and y, ED(x, y) is the minimum number of modifications (*edit operations*) required to change x into y, where a single modification can be adding a character, deleting a character or substituting a character in x. All logarithms are based-2 unless stated otherwise. For integers p > q,  $\sum_{i=p}^{q} a_i = 0$  by definition regardless of  $a_i$ 's.

#### 2.1 Grammars

Let  $\Sigma \subseteq \Gamma$  be two alphabets and  $\# \notin \Gamma$ . A grammar *G* is a set of *rules* of the type  $c \to ab$  or  $c \to a^r$ , where  $c \in (\Gamma \cup \{\#\}) \setminus \Sigma$ ,  $a, b \in \Gamma$  and  $r \in \mathbb{N}$ . *c* is the *left hand side* of the rule, and *ab* or  $a^r$  is the *right hand side* of the rule. # is the starting symbol.

The size |G| of the grammar is the number of rules in *G*. We only consider grammars where each  $a \in \Gamma \cup \{\#\}$  appears on the left hand side of at most one rule of G, we call such grammars deterministic. (We assume that rules of the form  $c \rightarrow a^r$  are stored in implicit (compressed) form.) The eval(G) is the string from  $\Sigma^*$  obtained from # by iterative rewriting of the intermediate results by the rules from G. If the rewriting process never stops or stops with a string not from  $\Sigma^*$ , eval(*G*) is undefined. Observe, that we can replace each rule of the type  $c \rightarrow a^r$  by a collection of at most  $2\lceil \log r \rceil$  new rules of the other type using some auxiliary symbols. Hence, for each grammar G there is another grammar G' using only the first type of the rules such that eval(G) = eval(G') and  $|G'| \leq$  $|G| \cdot 2\lceil \log |eval(G)| \rceil$ . Using a depth-first traversal of a deterministic grammar G we can calculate its *evaluation size* |eval(G)| in time O(|G|). Given a deterministic grammar G and an integer m less or equal to its evaluation size, we can construct in time O(|G|) another grammar G' of size O(|G|) such that eval(G') = eval(G)[m, ...]. G' will use some new auxiliary symbols. Given a deterministic grammar G, using a depth-first traversal on symbols reachable from the starting symbol # we can identify in time O(|G|) the smallest sub-grammar  $G' \subseteq G$  with the same evaluation.

We will use the following observation of Ganesh, Kociumaka, Lincoln and Saha [14]:

PROPOSITION 2.1 ([14]). There is an algorithm that on input of two grammars  $G_x$  and  $G_y$  of size at most m computes the edit distance k of  $eval(G_x)$  and  $eval(G_y)$  in time  $O((m + k^2) \cdot poly(\log m + n))$ , where  $n = |eval(G_x)| + |eval(G_y)|$ .

# 2.2 Rolling Hamming Distance Sketch

For two strings *x* and *y* of the same length, we define their *mismatch information* MIS(*x*, *y*) = {(*i*, *x*[*i*], *y*[*i*]); *i*  $\in$  {1, . . . , |*x*|} *and x*[*i*]  $\neq$  *y*[*i*]}. The Hamming distance of *x* and *y* is Ham(*x*, *y*) = |MIS(*x*, *y*)|.

There exist various sketches for Hamming distance, which allow to compute Hamming distance with low error probability [13, 21]. Moreover, [10, 27] also allow to retrieve the mismatch information. For our purposes we will use the sketch given by Clifford, Kociumaka, and Porat [10].

Let  $k \leq n$  be integers and  $p \geq n^3$  be a prime. [10] give a randomized sketch for Hamming distance  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}} : \{1, \ldots, p-1\}^* \rightarrow \{0, \ldots, p-1\}^{k+4}$  computable in time  $\widetilde{O}(n)$  with the following properties.<sup>1</sup>

PROPOSITION 2.2 ([10]). There is a randomized algorithm working in time  $O(k \log^3 p)$  that given sketches  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(x)$  and  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(y)$  of two strings x and y of length  $\ell \leq n$  constructed using the same randomness decides whether  $\operatorname{Ham}(x, y) \leq k$ , and if so returns  $\operatorname{MIS}(x, y)$ , with probability of error at most 1/n over the randomness of the sketches and the internal randomness of the algorithm.

They also construct the following update procedures for their sketch. We will use them to construct a rolling sketch for edit distance.

PROPOSITION 2.3 (LEMMA 2.3 OF [10]). For  $x \in \{1, ..., p\}^*$  of length less than n and  $a \in \{1, ..., p\}$ , in time  $O(k \log p)$  we can compute:

(1) 
$$\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(xa)$$
 and  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(ax)$ , given  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(x)$  and a.  
(2)  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(x)$  given  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(xa)$  or  $\operatorname{sk}_{n,k,p}^{\operatorname{Ham}}(ax)$ , and a.

Corollary 2.5 of [10] states that appending a character to a sketch of *x* can be done even faster namely in amortized time  $O(\log p)$ .

#### 2.3 Locally Consistent Coloring

The following color reduction procedure allows for locally consistent parsing of strings. The technique was originally proposed by Cole and Vishkin [11] and further studied by Linial [23, 24].

PROPOSITION 2.4 ([11, 23, 24]). There exists a function  $F_{\text{CVL}}$ :  $\Gamma^* \rightarrow \{1, 2, 3\}^*$  with the following properties. Let  $R = \log^* |\Gamma| + 20$ . For each string  $x \in \Gamma^*$  in which no two consecutive symbols are the same:

- (1)  $|F_{\text{CVL}}(x)| = |x|$  and  $F_{\text{CVL}}(x)$  can be computed in time  $O(R \cdot |x|)$ .
- (2) For  $i \in \{1, ..., |x|\}$ , the *i*-th symbol of  $F_{CVL}(x)$  is a function of symbols of x only in positions  $\{i R, i R + 1 ..., i + R\}$ .
- (3) No two consecutive symbols of  $F_{\text{CVL}}(x)$  are the same.
- (4) Out of every three consecutive symbols of F<sub>CVL</sub>(x) at least one of them is 1.
- (5) If |x| = 1 then  $F_{CVL}(x) = 3$ , and otherwise  $F_{CVL}(x)$  starts by 1 and ends by either 2 or 3.

The first three items are standard for  $R = \log^* |\Gamma| + 10$ . The other two can be obtained by a simple modification of the output of the standard function.

# **3 DECOMPOSITION ALGORITHM**

In this section we describe our main technical tool that we have developed. It is a randomized procedure that splits a string *x* into blocks  $B_1^x, B_2^x, \ldots, B_s^x$  and for each block it produces a grammar of size at most  $S = \widetilde{O}(k)$ . Furthermore, if  $B_1^x, B_2^x, \ldots, B_s^x$  is the decomposition for a string *x* and  $B_1^y, B_2^y, \ldots, B_{s'}^x$  is the decomposition for a string *y*, obtained using the same randomness, where  $ED(x, y) \leq k$  then with good probability, s = s' and  $B_i^x = B_i^y$  for all but *k* indices *i*. The edit distance of *x* and *y* can be calculated as  $ED(x, y) = \sum_i ED(B_i^x, B_i^y)$  where *i* ranges over the differing blocks.

First we provide an overview of the algorithm, specific details are given in the next sub-section. The decomposition procedure proceeds in  $O(\log n)$  rounds. In each round, the algorithm maintains a decomposition of x into *compressed* blocks. In each round each block of size at least two is first *compressed* and then *split*. The compression is done by compressing pairs of consecutive symbols into one using a randomly chosen pair-wise independent hash function  $C_{\ell} : \Gamma^2 \to \Gamma$ , where  $\ell$  is the round number (*level*). Nonoverlapping pairs of symbols are chosen for compression using a *locally consistent coloring* so that every three symbols shrink to at most two. Prior to the compression of pairs we replace each repeated sequence  $a^r$  of a symbol  $a, r \ge 2$ , by a special character  $\Gamma_{a,r}$ .

The splitting procedure uses a  $\widetilde{O}(k)$ -wise independent hash function  $H_{\ell}: \Gamma^2 \to \{0, \dots, D-1\}$  to select places where to subdivide

<sup>&</sup>lt;sup>1</sup>Clifford, Kociumaka and Porat have the sketch size only k + 3 elements but we include as an extra item the randomness of the sketch, which is a single element from  $\{0, \ldots, p-1\}$  used to compute Karp-Rabin fingerprint.

each block into sub-blocks, where  $D = \widetilde{O}(k)$  is a suitable parameter. We start a new block at each consecutive pair of symbols *ab*, where  $H_{\ell}(ab) = 0$ .

After  $O(\log n)$  rounds, each block is compressed into at most two symbols and we output a grammar that can generate the block.

For the correctness of the algorithm we will need to establish several properties of the algorithm. Some of these properties are related to behaviour on a single string x, others analyze the behaviour of the procedure on a pair of strings x and y of edit distance at most k.

The properties we want from the algorithm when it runs on *x* are the following: In each round, each block should be compressed by factor at least 2/3 while the size of the required grammar capturing the compression should be  $\tilde{O}(k)$ . The former is achieved by the design of the compression procedure. The latter goal is provided by the property of the splitting procedure which makes sure that each block  $B = b_1b_2 \cdot b_m$  resulting from a split has small dictionary Dict(B) = { $b_ib_{i+1}$ , i = 1, ..., m - 1}. In particular, we require |Dict(B)| =  $\tilde{O}(k)$ . The grammar size will be proportional to this dictionary.

For the compression procedure we require that it preserves information so the function  $C_{\ell}$  is one-to-one on each Dict(B). Since the total size of all dictionaries is bounded by  $\widetilde{O}(n)$  this can be easily achieved by picking  $C_{\ell}$  at random provided that its range size is  $\Omega(n^3)$ .

Additionally, we need the following property to hold on a pair of strings x and y of edit distance at most k with good probability: The splitting procedure should never split x or y in a *region* which is affected by edit operations that transform x to y (for some canonical choice of those operations.) The total size of those regions will be again  $\widetilde{O}(k)$  so we can satisfy this property if each pair of symbols has probability at most  $1/\widetilde{O}(k)$  to start a new block. This constrains the choice of the range size for the splitting function  $H_{\ell}$ .

In the next section we describe the decomposition algorithm fully, and then we establish its properties.

### 3.1 Algorithm Description

Let *n* be an upper bound on the length of the input string and  $k \le n$  be given. Set  $L = \lceil \log_{3/2} n \rceil + 3$  to be an upper bound on the decomposition depth. Let  $\Sigma$  be an input alphabet of size at most  $n^3$ ,  $\Sigma_c = \{c_1, c_2, \ldots, c_{Ln^3}\}$  and  $\Sigma_r = \{r_{a,r}, a \in \Sigma \cup \Sigma_c, r \in \{2, 3, \ldots, n\}\}$  be auxiliary pair-wise disjoint alphabets. Let  $\Gamma = \Sigma \cup \Sigma_c \cup \Sigma_r$  be the working alphabet, and # be a symbol not in  $\Gamma$ . Notice  $|\Gamma| = O(n^5 + |\Sigma|)$ . We call symbols from  $\Sigma_c^0 = \Sigma$  *level-0 compression symbols*, and for  $\ell \ge 1$ , symbols from  $\Sigma_c^0 = \{c_i, (\ell - 1)n^3 < i \le \ell n^3\}$  are *level-l compression symbols*. Additionally, symbols from  $\Sigma_r^\ell = \{r_{a,r} \in \Sigma_r, a \text{ is a level-}(\ell - 1) \text{ compression symbol}\}$  are also *level-l compression symbols*.

Let  $R = \log^* |\Gamma| + 20$ , D = 110R(L + 1)k and  $S = 30DL \log n + 6$ be parameters. The algorithm is a recursive algorithm of depth at most *L*. It starts by selecting at random several hash functions: For  $\ell = 1, ..., L$ , it selects at random a compression hash function  $C_{\ell} : \Gamma^2 \to \Sigma_c^{\ell}$  from a pair-wise independent hash family, and for  $\ell = 0, ..., L$ , it selects at random a splitting function  $H_{\ell} : \Gamma^2 \to \{0, ..., D-1\}$  from a (5*D* log *n*)-wise independent hash family. Main building blocks of the algorithm are functions Compress and Split. The first one compresses strings by a factor of 2/3, and the other splits strings at random points. Their pseudo-code is provided as Algorithm 1 and 2. We describe them next.

Compress. The function Compress( $B, \ell$ ) takes as input a string *B* over alphabet  $\Gamma$  of length at least two, and an integer  $\ell \geq 1$ , which denotes the level number. Divide B into minimum number of blocks  $B_1, \ldots, B_m, B = B_1 B_2 B_3 \ldots B_m$ , so that in each  $B_i$  either all the characters are the same, i.e.  $B_i = a^r$  for some  $a \in \Gamma$  and  $r \ge 2$ , or no two adjacent characters are the same. The first step is to compress the  $B_i$ 's which contain repeated characters by simply replacing the whole  $B_i$  with the symbol  $r_{a,|B_i|}$ , where *a* is the repeated character. Then for the remaining blocks, the following compression is applied: Let  $B_i$  be an uncompressed block. Each character of  $B_i$  is colored by applying  $F_{\text{CVL}}(B_i)$ . Divide  $B_i$  into blocks  $B_i = B'_1 B'_2 \dots B'_s$ , such that for each  $B'_i$  only the first character is colored 1. Now, according to Proposition 2.4, length of each  $B'_i$  is either 2 or 3. If  $B'_i = ab$ , replace it with  $C_{\ell}(ab)$  else if  $B'_{i} = abc$ , replace it with  $C_{\ell}(ab) \cdot c$ , where  $a, b, c \in \Gamma$ . The actual pseudo-code given below performs the compression of blocks of repeats in two stages, where in the first stage we replace the repeated sequence  $a^r$  by  $r_{a,r} \cdot \#$ , and then in the next stage we remove the extra symbol #. This simplifies analysis in Lemma 3.10. Assuming that  $C_{\ell}$  can be evaluated in time O(1), the running time of Compress $(B, \ell)$  is dominated by the time needed to compute  $F_{\text{CVL}}$ -coloring of blocks which is  $O(R \cdot |B|)$  in total.

Algorithm 1: Compress( <i>B</i> , <i>l</i> )					
<b>Input:</b> String <i>B</i> over alphabet $\Gamma$ of length at least two, and					
level number $\ell$ .					
<b>Output:</b> String $B''$ over alphabet $\Gamma$ .					
1 Divide $B = B_1 B_2 B_3 \dots B_m$ into minimum number of blocks					
so that each maximal subword $a^r$ of $B$ , for $a \in \Gamma$ and $r \ge 2$ ,					
is one of the blocks.					
2 <b>for</b> each $i \in \{1,, m\}$ <b>do</b>					
3 <b>if</b> $B_i = a^r$ , where $r \ge 2$ then Set $B'_i = r_{a,r} \cdot \#$ and color					
$r_{a,r}$ by 1 and # by 2. <sup>2</sup> ;					
4 <b>else</b> Set $B'_i = B_i$ and color each symbol of $B'_i$ according					
to $F_{\text{CVL}}(B_i)$ ;					
5 end					
6 Set $B' = B'_1 B'_2 \cdots B'_m$ , $B'' = \varepsilon$ , and $i = 1$ .					
7 while $i <  B' $ do					
8 <b>if</b> $B'[i+1] = #$ <b>then</b> $B'' = B'' \cdot B'[i];$					
9 <b>else</b> $B'' = B'' \cdot C_{\ell}(B'[i, i+1]);$					
10 $i = i + 2.$					
11 <b>if</b> $i \le  B' $ and $B'[i]$ is not colored 1 <b>then</b>					
$B^{\prime\prime} = B^{\prime\prime} \cdot B^{\prime}[i], i = i + 1;$					
12 end					
13 Return $B''$ .					

<sup>&</sup>lt;sup>2</sup>If  $a = r_{b,s}$  for some  $b \in \Gamma$  and  $s \in \mathbb{N}$ , then set  $B'_i = r_{b,rs} \cdot \#$ . However, such a situation should never happen during the execution of the algorithm as level- $\ell$  compression symbol can be introduced only at level  $\ell$ .

Split. The function takes as input a string *B* over alphabet  $\Gamma$  of length at least two, and an integer  $\ell \geq 1$ . The function splits the string *B* into smaller blocks. The algorithm works as follows: For each  $i \in \{2, ..., |B| - 1\}$ , if  $H_{\ell}(B[i, i + 1]) = 0$ , start a new block at position *i*. The running time of Split(*B*,  $\ell$ ) is dominated by the time to evaluate  $H_{\ell}$  at |B| - 2 points.

ł	Algorithm 2: $Split(B, \ell)$				
	<b>Input:</b> String <i>B</i> over alphabet $\Gamma$ of length at least two, and				
	level number $\ell$ .				
	<b>Output:</b> A sequence of strings $(B_0, B_1, \ldots, B_s)$ over				
	alphabet Γ.				
1	Let $i_1 < \dots < i_s$ be all $i \in \{2, \dots,  B  - 1\}$ where				
	$H_{\ell}(B[i, i+1]) = 0$ . Set $s = 0$ if no such <i>i</i> exists.				
2	Let $i_0 = 1$ and $i_{s+1} =  B  + 1$ .				
3	For $j = 0,, s$ , set $B_j = B[i_j, i_{j+1})$ .				
4	Return $(B_0, B_1,, B_s)$ .				

The main recursive step of the algorithm is encompassed in function Process. The function gets a block  $B \in \Gamma^*$  as its input. The block might have already been compressed previously, so the function also gets dictionaries that allow decompression of the block. If the block is already of length at most two, then the function outputs the block. Otherwise it compresses the block *B* using Compress, then it subdivides the compressed block using Split, and invokes itself recursively on each sub-block. For the output, each block is represented by a grammar. The grammar is reconstructed from the compressed block and its dictionaries by a simple bread-first search algorithm provided in the function Grammar.

<b>Algorithm 3:</b> $Process(B, (D_1, D_2,, D_{\ell-1}), \ell)$				
<b>Input:</b> String $B \in \Gamma^*$ , a sequence of dictionaries $D_i \subseteq \Gamma^2$ for decompressing <i>B</i> , and level number $\ell$ .				
<b>Output:</b> A sequence of blocks of <i>B</i> each encoded by a				
grammar.				
1 <b>if</b> $ B  \le 2$ <b>then</b> Output				
Grammar( $B, (D_1, D_2,, D_{\ell-1}), \ell - 1$ ) and return ;				
$A = \text{Compress}(B, \ell).$				
$(B_0, B_1,, B_s) = \text{Split}(A, \ell).$				
4 For $i = 0,, s$ , Process $(B_i, (D_1,, D_{\ell-1}, \text{Dict}(B)), \ell + 1)$ .				

To decompose an input string x into blocks, we first apply function Split(x, 0) to x and then invoke Process(B, (), 1) on each of the obtained blocks B. Breaking the string x into sub-blocks guarantees that each block passed to Process has small dictionary whereas the dictionary of x could have been arbitrarily large.

# 3.2 Correctness of the Decomposition Algorithm

Our goal is to establish the following theorem which is a stronger version of Theorem 1.1:

Algorithm 4: Grammar $(B, (D_1, D_2, \ldots, D_\ell), \ell)$						
	<b>Input:</b> String $B \in \Gamma^*$ , a sequence of dictionaries $D_i \subseteq \Gamma^2$ for					
	decompressing B.					
	<b>Output:</b> The smallest grammar <i>G</i> for <i>B</i> based on the					
	dictionaries $D_i$ and hash functions $C_1, \ldots, C_{\ell}$ .					
1	Let $C = \{c \in \Sigma_c :$					
	<i>c</i> appears in <i>B</i> or $r_{c,r}$ appears in <i>B</i> for some <i>r</i> }. // Symbols					
	needed to decompress B					
2	$2  G = \{\# \to B\}.$					
3	<b>3</b> for $j = \ell,, 1$ do					
4	4 <b>for</b> each $ab \in D_i$ <b>do</b>					
5	<b>if</b> $C_j(ab) \in C$ <b>then</b> $G = G \cup \{C_j(ab) \rightarrow ab\},$					
6	$C = C \cup \{c \in \Sigma_c; c \in \{a, b\} \text{ or } r_{c,r} \in \{a, b\} \}$					
	$\{a, b\}$ for some $r\}$ ;					
7	end					
8	8 end					
9	For each $r_{a,r}$ appearing in any of the rules in <i>G</i> , add					
	$r_{a,r} \rightarrow a^r$ to G.					

10 Return G.

THEOREM 3.1. Let x and y be a pair of strings of length at most n with  $ED(x, y) \le k$ . Let  $G_1^x, \ldots, G_s^x$  and  $G_1^y, \ldots, G_{s'}^y$  be the sequence of grammars output by the decomposition algorithm on input x and y respectively, using the same choice of random functions  $C_1, \ldots, C_L$ and  $H_0, \ldots, H_L$ . The following is true for n large enough:

- (1) With probability at least 1 2/n,  $x = eval(G_1^x) \cdots eval(G_s^x)$ and  $y = eval(G_1^y) \cdots eval(G_{s'}^y)$ .
- (2) With probability at least  $1 2/\sqrt{n}$ , for all  $i \in \{1, \ldots, s\}$  and  $j \in \{1, \ldots, s'\}, |G_i^x|, |G_j^y| \le S$ .
- (3) With probability at least 9/10, s = s', G<sub>i</sub><sup>x</sup> = G<sub>i</sub><sup>y</sup>, for all i ∈ {1,...,s} except for at most k indices i, and ED(x, y) = ∑<sub>i</sub> ED(eval(G<sub>i</sub><sup>x</sup>), eval(G<sub>i</sub><sup>y</sup>)).

By union bound, all three parts happen simultaneously with probability at least  $9/10 - 2/n - 1/\sqrt{n}$  which is  $\ge 4/5$  for *n* large enough.

To prove the theorem we make some simple observations about the algorithm, first.

LEMMA 3.2. For any string B of length at least two, and  $\ell \ge 1$ ,  $|\text{Compress}(B, \ell)| \le \frac{2}{3}|B| + 1$  and  $|\text{Compress}(B, \ell)| < |B|$ .

PROOF. Let  $B = B_1B_2B_3...B_m$  be as in the procedure. Every block  $B_i$  that equals to  $a^r$ , for some a and  $r \ge 2$ , is reduced to one symbol by the compression. The other blocks are colored using  $F_{\text{CVL}}(\cdot)$  and compressed. Unless a block  $B_i$  is of size one, the coloring induces division of the block  $B_i$  into subwords of size two or three, where the former is compressed into one symbol and the latter into two symbols. Hence, each such a block is compressed to at most 2/3 of its size. So the only blocks  $B_i$  that do not shrink are of size one, and are sandwiched between blocks of repeated symbols (that shrink by a factor of at least two). The worst-case situation is when m is odd, blocks  $B_i$  are of size one for odd i, and of size two for even i. In that case the original string B shrinks to size  $\lfloor \frac{2}{3} |B| \rfloor + 1$ . This proves the first inequality. The second inequality is also clear from the analysis above: The only time the string does not shrink is if it is of size one.  $\hfill \Box$ 

COROLLARY 3.3. On a string B of length at most n, the depth of the recursive calls of Process is at most L.

Indeed, from the previous lemma it follows that each block after  $\ell$  compressions and splits is of size at most  $(2/3)^{\ell}|B|+3$ . Hence, after  $L = \lceil \log_{3/2} n \rceil + 3$  recursive calls Process must stop the recursion.

LEMMA 3.4. Let  $B \in \Gamma^*$  be of length at most n, and  $\ell \in \{0, \ldots, L\}$ . Let  $(B_0, B_1, \ldots, B_s) = \text{Split}(B, \ell)$  where  $H_\ell : \Gamma^2 \to \{0, \ldots, D-1\}$  is chosen at random from  $(5D \log n)$ -wise independent hash family. Then with probability at least  $1 - 1/n^3$ , for all  $j \in \{0, \ldots, s\}$ ,  $|\text{Dict}(B_j)| \leq 5D \log n$ .

PROOF. If for some  $j \in \{0, ..., s\}$ ,  $|\operatorname{Dict}(B_j)| > 5D \log n$ , then there exists  $1 < r < t \le |B|$  such that  $|\operatorname{Dict}(B[r, t])| = 5D \log n$ and for all  $i \in \{r, ..., t - 1\}$ ,  $H_{\ell}(B[i, i + 1]) \ne 0$ . (Pick r to be the position in B of the second symbol of  $B_j$  and r some later position in  $B_j$ .) For a fixed r and t with  $|\operatorname{Dict}(B[r, t])| = 5D \log n$ ,  $\Pr_{H_{\ell}}[\forall i \in \{r, ..., t - 1\}, H_{\ell}(B[i, i + 1]) \ne 0] \le \left(1 - \frac{1}{D}\right)^{5D \log n}$  by the  $(5D \log n)$ -wise independence of  $H_{\ell}$ . Hence,  $\Pr_{H_{\ell}}[\exists 1 < r < t \le |B|, |\operatorname{Dict}(B[r, t])| = 5D \log n$  and  $\forall i \in \{r, ..., t - 1\}, H_{\ell}(B[i, i + 1]) \ne 0] \le |B|^2 \left(1 - \frac{1}{D}\right)^{5D \log n} \le n^2 e^{-5 \log n} \le 1/n^3$ .  $\Box$ 

LEMMA 3.5. For  $B \in \Gamma^*$ ,  $\ell \leq L, D_1, D_2, \ldots, D_\ell \subseteq \Gamma^2$ , Grammar $(B, (D_1, \ldots, D_\ell), \ell)$  outputs a grammar G of size at most  $3|B| + 6\sum_i |D_i|$ , and runs in time  $\widetilde{O}(|B| + \sum_i |D_i|)$ .

**PROOF.** The main loop of the algorithm iterates over all the pairs from  $D_j$ . In each iteration we can add a rule of the type  $c \rightarrow ab$  to G. Hence, the number of such rules in G is at most  $|B| + 2\sum_i |D_i|$ . Last, we add to G rules for symbols from  $\Sigma_r$  that appear on right hand sides of rules in G. This increases the size of G by at most factor of 3. If C is stored using some efficient data structure such as binary search trees or hash tables, each iteration takes  $\widetilde{O}(1)$  time. (We assume that evaluation of  $C_j(\cdot)$  takes O(1).) Hence, the total running time is bounded by claimed bound.

During processing of a string x, there are at most Ln calls to the function Split. (The actual number of calls is O(n) as the strings shrink exponentially but our simple upper bound suffices.) The probability that any one of them would produce a block with dictionary larger than  $5D \log n$  is at most  $Ln/n^3$ . We can conclude the next corollary which implies the second item of Theorem 3.1.

COROLLARY 3.6. For n large enough, on a string x of length at most n, processing the string x produces a sequence of grammars each of size at most  $S = 30DL \log n + 6$  with probability at least 1 - 1/n.

For the grammars produced by the algorithm to be deterministic, we need that each  $C_{\ell}$  is one-to-one on Dict(*B*) for each block *B* on which Compress(*B*,  $\ell$ ) is invoked. That will happen with high probability by a standard argument:

LEMMA 3.7. Let  $B \in \Gamma^*$  be of length at most n and  $\ell \in \{1, \ldots, L\}$ . Let  $C_{\ell} : \Gamma^2 \to \{c_i, (\ell-1)n^3 < i \leq \ell n^3\}$  be chosen at random from a pair-wise independent family of hash functions. Then with probability at least  $1 - |B|/n^2$ ,  $C_{\ell}$  is one-to-one on Dict(B). PROOF. For two distinct elements from Dict(B), the probability of a collision for randomly chosen  $C_{\ell}$  is at most  $1/n^3$ . By the union bound, the probability that  $C_{\ell}$  is not one-to-one on  $\text{Dict}(B_j)$  is at most  $|\text{Dict}(B)|^2/n^3 \le |B|/n^2$  as  $|\text{Dict}(B)| \le |B| \le n$ .

During processing of a string x, there are at most Ln calls to the function Compress. For a fixed level  $\ell \in \{1, \ldots, L\}$ , the total size of blocks B for which Compress $(B, \ell)$  is invoked is at most n. By the previous lemma and the union bound, the probability that during any of those calls Compress $(B, \ell)$  uses a function  $C_{\ell}$  that is not one-to-one on Dict(B) is at most 1/n. If all the hash functions  $C_1, C_2, \ldots, C_L$  that are used to compress blocks of x are one-to-one on their respective blocks then the grammars that Grammar produces will be deterministic, and they will evaluate to their respective blocks of x. (We can actually conclude a stronger statement that each  $C_{\ell}$  will be one-to-one on the union of all blocks at level  $\ell$  with high probability.) We can conclude the next corollary which implies the first item of Theorem 3.1.

COROLLARY 3.8. For *n* large enough, on a string *x* of length at most *n*, with probability at least 1 - L/n, processing the string *x* produces a sequence of grammars  $G_1, \ldots, G_s$  such that  $x = eval(G_1) \cdots eval(G_s)$ .

At this point we can estimate the running time of the decomposition algorithm. We can let the algorithm fail, and produce some trivial decomposition of x, whenever Split produces a block with dictionary larger than 5D log n. If it does not fail, then all grammars are of size at most S which is O(k). There are at most n of them so time spent in Grammar(...) is bounded by O(nk). The total time spent in Compress(...) is proportional to the sum of sizes of all nontrivial blocks over all levels of recursion which is O(nL) = O(n). (A more accurate estimate on the total size of blocks is O(n) since the blocks are shrinking geometrically in each iteration.) This means that the time to execute all calls to Compress is O(nLR) = O(n). The time spent in Split(...) is dominated by the time needed to evaluate  $H_{\ell}$ . The number of evaluation points at a given level  $\ell$  is proportional to the total size of all blocks at that level. Since  $H_{\ell}$  can be evaluated at a single point in time  $O(D \log n) = O(k)$ , we get a trivial upper bound  $O(nLD \log n) = \widetilde{O}(nk)$  on time spent in Split. Hence, in total the decomposition procedure runs in time  $\tilde{O}(nk)$ . (We believe that the total running time can be improved to O(n)on average. One could argue that in expectation the number of grammars the procedure produces is  $\widetilde{O}(n/k)$  as the average block size a string *x* is decomposed into should be at least  $\Omega(D/\log n)$ . So we believe that the total running time of calls to Grammar is O(n). Using multi-point evaluation of  $(5D \log n)$ -wise independent hash functions we could reduce the time for evaluation of  $H_{\ell}$  on a given level to  $\widetilde{O}(n)$ .)

PROPOSITION 3.9. Given  $k \le n$ , the running time of the decomposition algorithm on a string x of length at most n is  $\widetilde{O}(nk)$  with probability at least 1 - 1/n.

It remains to address the properties of the algorithm run on a pair of strings x and y of edit distance at most k to establish Theorem 3.1. For the pair of strings x and y we fix a *canonical decomposition of* x and y to be a sequence of words  $w_0, w_1, \ldots, w_k, u_i, \ldots, u_k, v_1, \ldots, v_k \in \Gamma^*$  such that  $x = w_0 u_1 w_1 u_2 w_2 \cdots u_k w_k, y = w_0 v_1 w_1 \cdots v_k w_k$ and  $|u_i|, |v_i| \le 1$  for all *i*. By the definition of edit distance such a decomposition exists: each pair  $(u_i, v_i)$  represents one edit operation, and we fix one such decomposition to be *canonical*. Observe, if we now partition x into blocks  $B_1^x, \ldots, B_s^x$  so that each  $B_i^x$  starts within one of the  $w_j$ 's, and we partition y into blocks  $B_1^y, \ldots, B_s^y$ so that each block  $B_i^y$  starts at the corresponding location in  $w_j$  as  $B_s^x$ , then  $\text{ED}(x, y) = \sum_i \text{ED}(B_s^x, B_s^y)$ .

We need to understand what happens with the decomposition of *x* and *y* when we apply the Compress function. Let x = uwvand  $x' = \text{Compress}(x, \ell) = u'w'v'$ , for some  $u, w, v, u'w'v' \in \Gamma^*$ . We say that a symbol c in w' comes from the compression of w if either it is directly copied from *w* by Compress, or it is the image  $c = C_{\ell}(ab)$  of a pair of symbols *ab* where *a* belongs to *w*, or  $c = r_{a,r}$ replaced a block  $a^r$  where the first symbol of  $a^r$  belongs to w. w' is the compression of w if it consists precisely of the symbols that come from the compression of *w*. Furthermore, we say a symbol *c* in w' comes weakly from the compression of w if either it is directly copied from *w* by Compress, or it is the image  $c = C_{\ell}(ab)$  of a pair of symbols *ab* where *a* or *b* belong to *w*, or  $c = r_{a,r}$  replaced a block  $a^r$  where some symbol of  $a^r$  belongs to w. w' is the weak compression of w if it consists precisely of the symbols that come weakly from the compression of w. Notice, a weak compression of w might contain and extra symbol at the beginning compared to the compression of *w*.

The following lemma captures what compression does to the canonical decomposition of x and y. (See Fig. 2 for illustration.)

LEMMA 3.10. Let x, y be strings over  $\Gamma$ , and let  $x' = \text{Compress}(x, \ell)$ and  $y' = \text{Compress}(y, \ell)$ . Let  $x = w_0 u_1 w_1 u_2 w_2 \cdots u_q w_q$  and  $y = w_0 v_1 w_1 v_2 w_2 \cdots v_q w_q$  for some strings  $w_i$ ,  $u_i$  and  $v_i$  where for  $i \in \{1, \ldots, q\}, |u_i|, |v_i| \le 4R + 24$ .

Then there are  $w'_0, w'_1, \ldots, w'_q, u'_1, \ldots, u'_q, v'_1, \ldots, v'_q \in \Gamma^*$  such that for  $i \in \{1, \ldots, q\}, |u'_i|, |v'_i| \leq 4R + 24, x' = w'_0 u'_1 w'_1 u'_2 w'_2 \cdots u'_q w'_q$ and  $y' = w'_0 v'_1 w'_1 v'_2 w'_2 \cdots v'_q w'_q$ . Moreover, each  $w'_i$  is the compression of the same subword of  $w_i$  in both x and y.

For each  $x = w_0 u_1 w_1 u_2 w_2 \cdots u_q w_q$ ,  $y = w_0 v_1 w_1 v_2 w_2 \cdots v_q w_q$ and  $\ell$  we fix one choice of  $w'_0, \ldots, w'_q, u'_0, \ldots, u'_q, v'_0, \ldots, v'_q$  satisfying the lemma. We will refer to it as the *canonical decomposition* of x' and y' induced by the decomposition of x and y as given by the lemma.

**PROOF.** The first stage of Compress replaces maximal blocks of repeated symbols by shortcuts. To simplify our analysis first we will reassign blocks of repeated symbols among neighboring blocks of  $w_i$ ,  $u_i$  and  $v_i$ , resp., so each maximal block of symbols in x and y is fully contained in one of the words  $w_i$ ,  $u_i$  or  $v_i$ .

For i = 1, ..., q - 1 we define words  $w_i^{(1)}$  and parameters  $a_i, b_i \in \Gamma$  and  $k_i, k'_i \in \mathbb{N}$  as follows: If  $w_i$  contains at least two distinct symbols let  $w_i = a_i^{k_i} w_i^{(1)} b_i^{k'_i}$  so that  $k_i$  and  $k'_i$  are maximum possible, otherwise  $w_i = a_i^{k_i}$  for some  $a_i$  and  $k_i$  ( $k_i$  might be zero), and we set  $w_i^{(1)} = \varepsilon$ ,  $b_i = a_i$  and  $k'_i = 0$ . Let  $w_0 = w_0^{(1)} b_0^{k'_0}$  for maximum possible  $k_0$  and some symbol  $b_0$ . Let  $w_q = a_q^{k_q} w_q^{(1)}$  for maximum possible  $k_q$  and some symbol  $a_q$ . For i = 1, ..., q, we let  $u_i^{(1)} = b_{i-1}^{k'_{i-1}} u_i a_i^{k_i}$ . Similarly,  $v_i^{(1)} = b_{i-1}^{k'_{i-1}} v_i a_i^{k_i}$ . Hence,  $x = w_0^{(1)} u_1^{(1)} w_1^{(1)} \cdots u_q^{(1)} w_q^{(1)}$  and  $y = w_0^{(1)} v_1^{(1)} w_1^{(1)} \cdots v_q^{(1)} w_q^{(1)}$ .

Next, if there is a maximal block of symbols  $a^r$  contained in  $u_s^{(1)} w_s^{(1)} \cdots u_t^{(1)}$  starting in  $u_s^{(1)}$  and ending in  $u_t^{(1)}$ ,  $s \neq t$ , we add all the symbols of the  $a^r$  to the end of  $u_s^{(1)}$  and remove them from the other  $u_i^{(1)}$ ,  $i = s + 1, \ldots, t$ . (Notice,  $w_i^{(1)} = \varepsilon$  for s < i < t because of the definition of  $w_i^{(1)}$ , and  $u_i^{(1)}$  will become empty for s < i < t.) We do this for all maximal blocks of repeated symbols that span multiple  $u_i^{(1)}$ . We perform similar moves on  $v_i^{(1)}$ , s. After all of those moves we denote the resulting subwords by  $w_i^{(2)}$ ,  $u_i^{(2)}$ , and  $v_i^{(2)}$ . (Notice,  $w_i^{(2)} = w_i^{(1)}$  for all *i*.) We have:  $x = w_0^{(2)} u_1^{(2)} w_1^{(2)} \cdots u_q^{(2)} w_q^{(2)}$  and  $y = w_0^{(2)} v_1^{(2)} w_1^{(2)} \cdots v_q^{(2)} w_q^{(2)}$ . At this stage, each maximal block of repeated symbols in x or y is contained in one of the subwords  $w_i^{(2)}$ ,  $u_i^{(2)}$ , and  $v_i^{(2)}$ .

The first stage of Compress replaces each maximal block  $a^r$ ,  $r \ge 2$ , by a sequence  $r_{a,r}$ #, and we apply this procedure on each subword  $w_i^{(2)}$ ,  $u_i^{(2)}$ , and  $v_i^{(2)}$  to obtain corresponding subwords  $w_i^{(3)}$ ,  $u_i^{(3)}$ , and  $v_i^{(3)}$ . Observe, for i = 1, ..., q,  $|u_i^{(3)}|$ ,  $|v_i^{(3)}| \le 4R+28$ . This is because every  $u_i$  is transformed into  $u_i^{(3)}$  by appending or prepending possibly empty block of repeated symbols, i.e.,  $u_i^{(3)} = a^r u_i b^{r'}$  for some a, b, r, r', or removing its content entirely. Each block of repeats is reduced to two symbols so each  $u_i^{(3)}$  is longer than the original by at most 4 symbols. Similarly for  $v_i^{(3)}$ .

Next, coloring function  $F_{\text{CVL}}$  is used on parts of x and y that are not obtained from repeated symbols; the two symbols replacing each repeated block are colored by 1 and 2, resp. We refer to this as {1, 2, 3}-coloring. At most R first and last symbols of each  $w_i^{(3)}$ might be colored differently in x and y as the color of each symbol depends on the context of at most R symbols on either side of the symbol, and that context might differ in x and y. Hence, only symbols near the border of  $w_i^{(3)}$  that are in vicinity of  $u_i^{(3)}$ 's and  $v_i^{(3)}$ 's, resp., might get different colors. All the other symbols of  $w_i^{(3)}$  are colored the same in both x and y. The coloring is then used to make decisions on which pairs of symbols are compressed into one.

We will let  $u'_i$  be the symbols that come from the compression of symbols in  $u_i^{(3)}$ , the first up-to R + 2 symbols of  $w_i^{(3)}$ , and the last up-to R + 3 symbols of  $w_{i-1}^{(3)}$ . Next we specify precisely which symbols of  $w_i^{(3)}$  and  $w_{i-1}^{(3)}$  are considered to be compressed into symbols belonging to  $u'_i$ . For i = 0, ..., q, if  $|w_i^{(3)}| \ge R + 3$ , let  $s_i^x$  be the position of the first symbol in  $w_i^{(3)}$  among positions R + 1, R + 2, R + 3 which is colored 1 in x by the  $\{1, 2, 3\}$ -coloring. If  $|w_i^{(3)}| < R + 3$ , let  $s_i^x = 1$ . Next, if  $|w_i^{(3)}| \ge 2R + 3$  set  $t_i^x$  to be the first position from left colored 1 among the symbols of  $w_i^{(3)}$  at positions R + 1, R + 2, R + 3 counting from right. If  $|w_i^{(3)}| < 2R + 3$ , set  $t_i^x$  to be equal to  $s_i^x$ . For i = 0, if  $|w_0^{(3)}| \ge R + 3$  then redefine  $s_0^x = 1$ . For i = q, redefine  $t_q^x = |w_q^{(3)}| + 1$  and if  $|w_q^{(3)}| < R + 3$  then redefine  $s_q^x$  to  $t_q^x$ . Similarly, define  $s_i^y$  and  $t_i^y$  based on the  $\{1, 2, 3\}$ -coloring of y.

Notice,  $s_i^x \neq t_i^x$  iff  $s_i^y \neq t_i^y$ . Furthermore, if  $s_i^x \neq t_i^x$  then either  $i \in \{q, 0\}$  or  $|w_i^{(3)}| \ge 2R + 3$  so  $s_i^x = s_i^y$  and  $t_i^x = t_i^y$  as the symbols

*R*-away from either end of  $w_i^{(3)}$  are colored the same in *x* and *y*. We let  $u'_i$  to be the compression of  $w_{i-1}^{(3)}[t_{i-1}^x, |w_{i-1}^{(3)}|] \cdot u_i^{(3)} \cdot w_i^{(3)}[1, s_i^x)$  and similarly,  $v'_i$  to be the compression of  $w_{i-1}^{(3)}[t_{i-1}^y, |w_{i-1}^{(3)}|] \cdot v_i^{(3)} \cdot w_i^{(3)}[1, s_i^y)$ . We let  $w'_i$  be the compression of  $w_i^{(3)}[s_i^y, t_i^y)$ .

Hence,  $u'_i$  comes from the compression of at most  $|u_i^{(3)}|+2R+5 \le 6R+33$  symbols. Since each symbol after a symbol colored 1 is *removed* by the compression, and each consecutive triple of symbols contains at least one symbol colored by 1, the at most 6R+27 symbols are compressed into at most  $(6R+33) \cdot 2/3 + 2 = 4R + 24$  symbols. So  $u'_i$  is of length at most 4R + 24. Similarly for  $v'_i$ .

The following generalization of the previous lemma will be useful to design a rolling sketch. It considers situation where x and y are prefixed by some strings u and v, resp., that we want to ignore from the analysis. The proof of the lemma is a straightforward modification of the above proof.

LEMMA 3.11. Let  $x, y, u, v \in \Gamma^*$ , and let  $u'x' = \text{Compress}(ux, \ell)$ and  $v'y' = \text{Compress}(vy, \ell)$ , where x' is the weak compression of x, and y' is the weak compression of y. Let  $x = u_0w_0u_1w_1u_2w_2\cdots u_qw_q$ and  $y = v_0w_0v_1w_1v_2w_2\cdots v_qw_q$  for some strings  $w_i, u_i$  and  $v_i$  where for  $i \in \{0, \ldots, q\}, |u_i|, |v_i| \leq 4R+24$ . Then there are  $w'_0, w'_1, \ldots, w'_q, u'_0, u'_1, \ldots, u'_q, v'_0, v'_1, \ldots, v'_q \in \Gamma^*$  such that for  $i \in \{0, \ldots, q\}, |u'_i|, |v'_i| \leq 4R+24, x' = u'_0w_0u'_1w_1'u'_2w'_2\cdots u'_qw'_q$  and  $y' = v'_0w_0v'_1w_1'\cdots v'_qw'_q$ . Moreover, each  $w'_i$  is the compression of the same subword of  $w_i$  in both x and y.

Let  $x \in \Sigma^*$ . Let  $H_0, H_1, \ldots, H_L, C_1, C_2, \ldots, C_L$  be chosen. We define inductively the *trace* of the algorithm on x at level  $\ell \ge 0$  to consist of sequences  $B^x(\ell, 1), \ldots, B^x(\ell, s^x_\ell) \in \Gamma^*$ , of auxiliary sequences  $A^x(\ell, 1), \ldots, A^x(\ell, s^x_\ell) \in \Gamma^*$  and  $t^x_{\ell,1}, \ldots, t^x_{\ell,s^x_\ell+1} \in \mathbb{N}$ . Their meaning is:  $B^x(\ell, i)$  is compressed into  $A^x(\ell, i)$  and that is split into blocks  $B^x(\ell + 1, j)$  for  $t^x_{\ell+1, i} \le j < t^x_{\ell+1, i+1}$ . (See Fig. 1 for illustration.)<sup>3</sup>

Set

 $B^{x}(0, 1), \dots, B^{x}(0, s_{0}^{x}) = \text{Split}(x, 0).$ 

For  $\ell = 1, \dots, L$  we define  $B^x(\ell, 1), \dots, B^x(\ell, s^x_{\ell})$  inductively. Set  $t^x_{\ell,1} = 1$ . For  $i = 1, \dots, s^x_{\ell-1}$ , if  $|B^x(\ell-1, i)| \ge 2$ , then

 $A^{x}(\ell - 1, i) = \operatorname{Compress}(B^{x}(\ell - 1, i), \ell),$ 

and for  $(B_0, B_1, ..., B_s) = \text{Split}(A^x(\ell - 1, i), \ell)$  set

$$B^{x}(\ell, t_{\ell,i}^{x}) = B_{0}, \ B^{x}(\ell, t_{\ell,i}^{x} + 1) = B_{1}, \ \dots, \ B^{x}(\ell, t_{\ell,i}^{x} + s) = B_{s}$$

and  $t_{\ell,i+1}^x = t_{\ell,i}^x + s + 1$ . If  $|B^x(\ell - 1, i)| < 3$ , then set  $B^x(\ell, t_{\ell,i}^x)$  and  $A^x(\ell - 1, i)$  to  $B^x(\ell - 1, i)$ , and  $t_{\ell,i+1}^x = t_{\ell,i}^x + 1$ . For  $j = s_{\ell-1}^x$ , set  $s_{\ell}^x = t_{\ell,i+1}^x$ .

Furthermore, for x and  $y \in \Sigma^*$ ,  $\ell$ ,  $i \ge 0$ , define a canonical decomposition of blocks  $A^x(\ell, i)$ ,  $B^x(\ell, i)$ ,  $A^y(\ell, i)$ ,  $B^y(\ell, i)$  inductively as follows. Let  $A^x(-1, 1) = x$  and  $A^y(-1, 1) = y$ . Let  $t^x_{-1,1} = 1$ ,  $t^x_{-1,2} = 2$ ,  $s^x_{-1} = 1$ ,  $t^y_{-1,1} = 1$ ,  $t^y_{-1,2} = 2$ , and  $s^y_{-1} = 1$ . Let  $A^x(-1, 1) = w_0 u_1 w_1 u_2 w_2 \cdots u_k w_k$  &

 $A^{y}(-1,1) = w_0 v_1 w_1 v_2 w_2 \cdots v_k w_k$ 



Figure 1: The hierachical decomposition of *x*.

be the canonical decomposition of the pair x and y.

For  $\ell \ge 0$  and  $j \in \{1, \ldots, s_{\ell}^x\}$ , let *i* be such that  $t_{\ell-1,i}^x \le j < t_{\ell-1,i+1}^x$  and  $m = j - t_{\ell-1,i}^x$ . Then  $B^x(\ell, j)$  is the *m*-th block of Split $(A^x(\ell-1, i), \ell)$ . If the decomposition of  $A^x(\ell-1, i)$  is defined and is equal to  $w_0 u_1 w_1 u_2 w_2 \cdots u_q w_q$ , for some  $u_i, w_i \in \Gamma^*$ , then the decomposition of  $B^x(\ell, j)$  is the restriction of the decomposition of  $A^x(\ell-1, i)$  to symbols of the *m*-th block of Split $(A^x(\ell-1, i), \ell)$ . Otherwise the decomposition of  $B^x(\ell, j)$  is undefined. Similarly for  $B^y(\ell, j)$ . (See Fig. 2.)

For  $\ell \geq 0$  and  $j \in \{1, \ldots, s_{\ell}^{x}\}$ , if  $B^{x}(\ell, j)$  and  $B^{y}(\ell, j)$  have defined decompositions  $B^{x}(\ell, j) = w_{0}u_{1}w_{1}u_{2}\cdots u_{q}w_{q}$  and  $B^{y}(\ell, j) = w_{0}v_{1}w_{1}v_{2}w_{2}\cdots v_{q}w_{q}$  for some  $u_{i}, v_{i}, w_{i} \in \Gamma^{*}$ , then we let  $A^{x}(\ell, j) = w'_{0}u'_{1}w'_{1}u'_{2}\cdots w'_{q}$  and  $A^{y}(\ell, j) = w'_{0}v'_{1}w'_{1}v'_{2}\cdots w'_{q}$  be their canonical decomposition induced by  $B^{x}(\ell, j)$  and  $B^{x}(\ell, j)$  as given by Lemma 3.10.



Figure 2: Decomposition of  $B^{X}(\ell, i)$  after compression and split.

To conclude item 3 of Theorem 3.1 we want to argue that *x* and *y* are recursively split into sub-blocks that respect their canonical decomposition. So we want all splits of blocks to occur in matching parts of *x* and *y*. For  $A^x(\ell - 1, i)$  with canonical decomposition  $w_0u_1w_1u_2w_2\cdots u_qw_q$  we say that  $\text{Split}(A^x(\ell - 1, i), \ell)$  makes *undesirable split* if it starts a new block at a position *j* that either belongs to one of the  $u_1, u_2, \ldots, u_q$  or is the first or last symbol of one of the  $w_0, w_1, \ldots, w_q$ . Recall,  $\text{Split}(A^x(\ell - 1, i), \ell)$  starts a new block at each position *j* such that  $H_\ell(A^x(\ell - 1, i), [j, j + 1]) = 0$ . Since  $H_\ell$  is chosen at random a given position starts a new block with probability 1/D.

For  $A^{y}(\ell-1, i)$  with canonical decomposition  $w'_{0}v_{1}w'_{1}v_{2}\cdots v_{q'}w'_{q'}$ we say that Split $(A^{y}(\ell-1, i), \ell)$  makes *undesirable split* if it starts a new block at position *j* that either belongs to one of the  $v_1, v_2, \ldots, v_{q'}$ 

<sup>&</sup>lt;sup>3</sup>To avoid double and triple indexes we use our notation  $B^{x}(\ell, i)$  and  $A^{x}(\ell, i)$  instead of the usual  $B^{x}_{\ell,i}$  and  $A^{x}_{\ell,i}$ .

or is the first or last symbol of one of the  $w'_0, w'_1, \ldots, w'_{q'}$ . If  $A^x(\ell - 1, i)$  and  $A^y(\ell - 1, i)$  have *matching* canonical decomposition (that is q = q' and each  $w_j = w'_j$ ) and both  $\text{Split}(A^x(\ell - 1, i), \ell)$  and  $\text{Split}(A^y(\ell - 1, i), \ell)$  make no undesirable split then  $A^x(\ell - 1, i)$  and  $A^y(\ell - 1, i)$  are split in the same number of blocks with matching canonical decomposition as they are split at the same positions in the corresponding  $w_j$ 's.

For given  $\ell \in \{0, \ldots, L\}$ , if no undesirable split happens during Split $(A^x(\ell'-1, i), \ell')$  and Split $(A^y(\ell'-1, i), \ell')$ , for any  $\ell' < \ell$  and *i*, then for each  $\ell' < \ell$ , the number of blocks  $B^x(\ell', i)$  and  $B^y(\ell', i)$  will be the same, i.e.,  $s_{\ell'}^x = s_{\ell'}^y$ , and blocks  $B^x(\ell', i)$  and  $B^y(\ell', i)$  will have matching canonical decomposition. The total number of  $u_j$ 's in canonical decomposition of all  $B^x(\ell', i)$ ,  $i = 1, \ldots, t_{\ell'}^x$ , will be at most k, and similarly for  $v_j$ 's. Thus, there will be at most (4R+24+2)k+2 positions where an undesirable split can happen in Split $(A^x(\ell - 1, i), \ell)$  for any *i*. Similarly, there are at most (4R+26)k+2 positions where an undesirable split can happen in Split $(A^y(\ell - 1, i), \ell)$ . By union bound, the probability that an undesirable split happens in some Split $(A^y(\ell - 1, i), \ell)$  or Split $(A^y(\ell - 1, i), \ell)$ , for some  $\ell$  and *i*, is at most  $2(4R+28)k(L+1)/D \le 11Rk(L+1)/D \le 1/10$ .

Thus, if no undesirable split happens there are at most k indices i for which the canonical decomposition of  $B^x(\ell, i)$  contains some  $u_j$ . All other blocks  $B^x(\ell, i)$  have a canonical decomposition consisting of a single block  $w_0$ , for various  $w_0$  depending on  $\ell$  and i. Similarly, the canonical decomposition of  $B^y(\ell, i)$  contains  $v_j$  if and only if  $B^x(\ell, i)$  contains  $u_j$ . Blocks  $B^y(\ell, i)$  that do not contain  $v_j$  are identical to  $B^x(\ell, i)$  so they have the same grammar.

Hence, if no undesirable split happens, item 3 of Theorem 3.1 will be satisfied.

The following theorem generalizes item 3 of Theorem 3.1 and it will be useful to construct the rolling sketch in Section 4.

THEOREM 3.12. Let  $u, v, x, y \in \Sigma^*$  be strings such that  $|ux|, |vy| \leq n$  and  $ED(x, y) \leq k$ . Let  $G_1^x, \ldots, G_s^x$  and  $G_1^y, \ldots, G_{s'}^y$  be the sequence of grammars output by the decomposition algorithm on input ux and vy respectively, using the same choice of random functions  $C_1, \ldots, C_L$  and  $H_0, \ldots, H_L$ . With probability at least 1 - 1/5 the following is true: There exist integers r, r', t, t' such that s - t = s' - t',

$$\begin{aligned} x &= \operatorname{eval}(G_t^x)[r, \dots] \cdot \operatorname{eval}(G_{t+1}^x) \cdots \operatorname{eval}(G_s^x) & \& \\ y &= \operatorname{eval}(G_{t'}^y)[r', \dots] \cdot \operatorname{eval}(G_{t'+1}^y) \cdots \operatorname{eval}(G_{s'}^y), and \\ \operatorname{ED}(x, y) &= \operatorname{ED}(\operatorname{eval}(G_t^x)[r, \dots], \operatorname{eval}(G_{t'}^y)[r', \dots]) \end{aligned}$$

+ 
$$\sum_{i>0} \text{ED}(\text{eval}(G_{t+i}^x), \text{eval}(G_{t'+i}^y)).$$

Its proof is a minor modification of the proof above and its sketch is provided in the full version [5].

# 3.3 Encoding a Grammar

We will set a parameter  $N \ge n^3$  to be a suitable integer: Let  $F_{\text{KR}}$  :  $\{0, 1\}^* \to \{1, ..., N\}$  be a hash function picked at random, such as Karp-Rabin fingerprint [18], so for any two strings  $u, v \in \{0, 1\}^*$ , if  $u \ne v$  then  $\Pr_{F_{\text{KR}}}[F_{\text{KR}}(u) = F_{\text{KR}}(v)] \le (|u| + |v|)/N$ .

Set  $M = 3S \cdot \lceil 1 + \log |\Gamma| \rceil$ . We will encode a grammar *G* over  $\Gamma$  of length at most *S* given by our decomposition algorithm by a string Enc(*G*) over alphabet  $\{1, \ldots, 2N\}$  of length *M*. The encoding is obtained as follows: First, order the rules of the grammar *G* 

lexicographically. Then encode the rules in binary one by one using  $3 \cdot \lceil 1 + \log |\Gamma| \rceil$  bits for each rule. (The extra bit allows to mark unused symbols.) This gives a binary string of length at most M, which we pad by zeros to the length precisely M. We call the resulting binary string Bin(G). Compute  $h_G = F_{KR}(Bin(G))$ . We replace each 0 in Bin(G) by  $h_G$ , and each 1 in Bin(G) by  $N + h_G$  to obtain the string Enc(G). Clearly, Enc(G) is a string over alphabet  $\{1, \ldots, 2N\}$  of length exactly M. The encoding can be computed in time O(M). For completeness, we encode any grammar G of length more than S or that uses rules with more than two symbols on the right as  $Enc(G) = 1^M$ .

By the property of  $F_{\text{KR}}$  the following holds.

LEMMA 3.13. Let G, G' be two grammars of size at most S output by our decomposition algorithm. Let  $F_{KR}$  be chosen at random.

- (1)  $\operatorname{Enc}(G) \in \{1, \dots, 2N\}^M$ .
- (2) If G = G' then Enc(G) = Enc(G').
- (3) If  $G \neq G'$  then Enc(G) = Enc(G') with probability at most 2M/N.
- (4) If Enc(G) ≠ Enc(G') then Ham(Enc(G), Enc(G')) = M, that is they differ in every symbol.

### 3.4 Edit Distance Sketch

Let *n* and  $k \le n$  be two parameters, and  $p \ge 2N + 1$  be a prime such that  $p \ge (nM)^3$ . For a string  $x \in \Sigma^*$  of length at most *n*, we compute its sketch by running first the decomposition algorithm of Theorem 3.1 to get grammars  $G_1, G_2, \ldots, G_s$ . Encode each grammar  $G_i$  by encoding  $\text{Enc}(G_i)$  from Section 3.3 using the same  $F_{\text{KR}}$  picked at random. Concatenate the encoding to get a string  $w = \text{Enc}(G_1) \cdot$  $\text{Enc}(G_2) \cdots \text{Enc}(G_s)$ . Calculate the Hamming sketch  $\text{sk}_{n',m',p}^{\text{Ham}}(w)$ on *w* for strings of length n' = nM and Hamming distance at most k' = kM from Section 2.2. Set the sketch  $\text{sk}_{n,k}^{\text{ED}}(x) = \text{sk}_{n',k',p}^{\text{Ham}}(w)$ . The calculation of  $\text{sk}_{n,k}^{\text{ED}}(x)$  can be done in time  $\widetilde{O}(nk)$  as the number of grammars is at most *n* and each grammar requires  $\widetilde{O}(k)$  time to be encoded into binary. The Hamming sketch can be constructed in time  $\widetilde{O}(nk)$ . (We believe that on average we expect only  $\widetilde{O}(n/k)$ grammars to be produced for a given string *x* so the actual running time should be  $\widetilde{O}(n)$  on average.)

THEOREM 3.14. Let  $x, y \in \Sigma^*$  be strings of length at most n such that  $ED(x, y) \leq k$ . Let  $sk_{n,k}^{ED}(x)$  and  $sk_{n,k}^{ED}(y)$  be obtained using the same randomness for the decomposition algorithm and the same choice of  $F_{\text{KR}}$ . With probability at least 2/3, we can calculate ED(x, y) from  $sk_{n,k}^{ED}(x)$  and  $sk_{n,k}^{ED}(y)$ .

Assume that the output of the decomposition algorithm on xand y satisfies all the conclusions of Theorem 3.12. In particular, for x we get  $eval(G_1^x) \cdot eval(G_2^x) \cdots eval(G_s^x)$  and for y we get  $eval(G_1^y) \cdots eval(G_s^y)$ , for some  $s \le n$ , each of the grammars is of size at most S,  $ED(x, y) = \sum_i ED(eval(G_i^x), eval(G_i^y))$ , and the number of pairs  $G_i^x$  and  $G_i^y$  where  $G_i^x \ne G_i^y$  is at most k. Assume that  $F_{\text{KR}}$  is chosen so that  $Enc(G_i^x) \ne Enc(G_i^y)$  for each of the pairs where  $G_i^x$  and  $G_i^y$  differ.

In order to determine ED(x, y), we recover the (Hamming) mismatch information between  $\text{Enc}(G_1^x) \cdot \text{Enc}(G_2^x) \cdots \text{Enc}(G_s^x)$  and  $\text{Enc}(G_1^y) \cdot \text{Enc}(G_2^y) \cdots \text{Enc}(G_s^y)$  from  $\text{sk}_{nk}^{\text{ED}}(x)$  and  $\text{sk}_{nk}^{\text{ED}}(y)$ . That gives grammars  $G_i^x$  and  $G_i^y$ , for all *i* where  $G_i^x \neq G_i^y$ . (Whenever the two grammars differ, their encoding differ in every symbol by Lemma 3.13 so we can recover them from the Hamming mismatch information.) Calculating the edit distance of each of the pair of differing grammars using the algorithm from Proposition 2.1 we recover ED(*x*, *y*) as the sum of their edit distances.

The sum is correct unless some of the assumptions fail: The probability that the grammar decomposition fails (does not have properties from Theorem 3.1) for the pair *x* and *y* is at most 1/5 for *n* large enough. The probability that the choice of  $F_{\text{KR}}$  fails (two distinct grammars have the same encoding) is at most 2kM/N < 1/n by the choice of *N*. The probability that the Hamming distance sketch fails to recover the mismatch information between all the grammars is at most 1/n. So in total, the probability that the output of the algorithm is incorrect is at most 1/3.

The running time of the comparison algorithm is  $O(k^2)$ : The Hamming mismatch information can be recovered in time  $O(kM) = O(k^2)$  (Proposition 2.2), then we build the  $\leq k$  mismatched grammars in time  $O(k^2)$ , and run the edit distance computation on the pairs of grammars in time  $\sum_{i < k} O(k + k_i^2) \leq O(k^2)$ , where  $k_i$  is the edit distance of the *i*-th pair of mismatched grammars. (We interrupt the edit distance computation if it takes more time than  $O(k^2)$  which would indicate ED(x, y) > k.)

To decide whether ED(x, y) > k we note that on input x and y, the Hamming sketch either outputs the correct mismatched places if their number is  $\leq k'$  or it outputs  $\infty$  if there are more mismatches than that or the sequences sketched by the Hamming sketch are of different length. (We assume that the Hamming sketch knows the number of symbols it is sketching.) In the  $\infty$ -case we know that there are more than k different pairs of grammars or the decomposition of x and y failed, and we can report ED(x, y) > k. In the other case we try to calculate the edit distance of the differing pairs of grammars. If we spend more than  $\widetilde{O}(k^2)$  time on it or we get a number larger than k then we report ED(x, y) > k. This correctly decides whether ED(x, y) > k with probability at least 2/3.

To prove Theorem 1.2 we build a more robust sketch by taking  $c \log n$  independent copies of the sketch  $\mathrm{sk}_{n,k}^{\mathrm{ED}}$ . To calculate the edit distance of two sketched strings we run the edit distance calculation on each of the corresponding pairs of copies, and output the majority answer. A usual application of Chernoff bound shows that the probability of correct answer is at least 1 - 1/n for suitable constant c > 0.

#### **4 ROLLING SKETCH FOR EDIT DISTANCE**

In this section we will construct the rolling sketch of Theorem 1.3. We will use two auxiliary claims. The first one addresses how much a compression of a string w might change depending on what is appended to it. Their proofs are omitted due to space limitations but can be found in the full version [5].

LEMMA 4.1. Let  $\ell \in \{0, ..., L\}$  and  $v, u, w \in \Gamma^*$ . Let w'u' =Compress $(wu, \ell)$  and let w''v' =Compress $(wv, \ell)$ , where w' is the compression of w when compressing wu and w'' is the compression of w when compressing wv. Let t = |w'| - 3(R + 1) or t = |w'u'| - |u| - 3(R + 1). Then w'[1, t] = w''[1, t].

The next lemma addresses how much the overall decomposition of a string x might change if we append a suffix z to it.

LEMMA 4.2. Let  $x, z \in \Sigma^*$ ,  $|xz| \le n$ . Let  $H_0, \ldots, H_L, C_1, \ldots, C_L$  be given. Let  $G_1^x, G_2^x, \ldots, G_s^x$  be the output of the decomposition algorithm on input x, and  $G_1^{xz}, G_2^{xz}, \ldots, G_{s'}^{xz}$  be the output of the decomposition algorithm on input xz using the given hash functions. Let T = L(3R + 6).

(1) 
$$G_i^x = G_i^{xz}$$
 for all  $i = 1..., s - T$ .  
(2)  $|x| \le \sum_{i=1}^{\min(s+T,s')} |eval(G_i^{xz})|$ .

The second part says that if x is decomposed into s grammars by itself, then it can be recovered from the first s + T grammars for xz. Hence, appending extra symbols to x cannot increase the number of grammars that cover x by more than T.

Let  $m \ge k$  and  $n \ge 10m^3$  be integers. A rolling sketch for a string obtained by up-to *m* insertions (to the right end) and *m* deletions (from the left end) from an empty word consists of three data structures: *insertion buffer*, *deletion buffer* and a Hamming distance sketch sk<sup>Ham</sup><sub>n',k',p</sub>, where k' = (4T + 1)(k + 2)M, n' = nM and  $p \ge n'^3$  is a chosen prime.

The insertion buffer maintains a buffer of *committed grammars*  $G_{s-4T+1}, G_{s-4T+2}, \ldots, G_s$  and a buffer of *active grammars*  $G_1^i, \ldots, G_t^i, t \leq T$ . The deletion buffer is similar, it maintains a buffer of *committed grammars*  $G_{r-4T+1}, G_{r-4T+2}, \ldots, G_r$  and a buffer of *active grammars*  $G_1^i, \ldots, G_{t'}^d, t' \leq T$ . The Hamming sketch is a sketch of grammars  $G_{r-2T+1}, G_{r-2T+2}, \ldots, G_{s-2T}$ , each encoded as a string of length M over the alphabet  $\{1, \ldots, 2N\}$ .

In addition to that, the sketch keeps track of the current value of r and s, and remembers a collection of pair-wise independent hash functions  $C_1, \ldots, C_L$ , a collection of  $(5D \log n)$ -wise independent hash functions  $H_0, \ldots, H_L$ , and randomness for Karp-Rabin fingerprint to compute binary encoding of grammars. The hash functions and the randomness of Karp-Rabin fingerprint are chosen at random when creating the sketch for empty string. This extra information requires  $\tilde{O}(k)$  bits to specify.

Initially, the committed grammars in the insertion and deletion buffers are all treated as empty sets, there are no active grammars in the insertion or deletion buffers so t = t' = 0 and s = r = 0.

For  $u, x \in \Sigma^*$ , if in total a string ux was inserted into the sketch then  $G_1, \ldots, G_s, G_1^i, \ldots, G_t^i$  represents ux, that is ux is the concatenation of the evaluation of the grammars. If in total the string u was deleted from the sketch, then  $G_1, \ldots, G_r, G_1^d, \ldots, G_{t^d}^d$  represents u. (See Fig. 3 for an illustration.)

Appending a symbol. When we append additional symbol a to the sketch we modify input buffers as follows: We update the active grammars  $G_1^i, \ldots, G_t^i$  by appending a as explained further below. Say the update produces grammars  $G_1'^i, \ldots, G_{t'}'^i$ . If  $t' \leq T$  then the produced grammars will become the active grammars, and no more changes are done to the sketch. Otherwise we commit the first t' - T grammars  $G_1'^i, \ldots, G_{t'-T}'^i$  one-by-one into the committed buffer as grammars  $G_{s+1}^i, \cdots, G_{s+t'-T}^i$  and we keep the remaining grammars as the active grammars.

Committing a grammar  $G_{s+1}$  into the committed buffer will trigger addition of  $G_{s-2T+1}$  into the Hamming sketch at the end of the represented sequence of grammars (if s - 2T + 1 > 0), and removing the grammar  $G_{s-4T+1}$  from the committed buffer. For insertion into

#### STOC '23, June 20-23, 2023, Orlando, FL, USA

#### Sudatta Bhattacharya and Michal Koucký





the Hamming sketch, the grammar  $G_{s-2T+1}$  is encoded into binary as in Section 3.3 and then the binary string is encoded using the Karp-Rabin fingerprint  $F_{\text{KR}}$  of *all* the grammars  $G_{s-4T+1}, \ldots, G_{s+1}$ , instead of only the grammar  $G_{s-2T+1}$ . (Thus, a change in any of the neighboring grammars will trigger a recovery of also the grammar  $G_{s-2T+1}$  when calculating a mismatch information from the Hamming sketch.) We repeat this process for each grammar being committed.

By the second part of Lemma 4.2  $t' \le t + T \le 2T$  so we will commit at most  $T = \widetilde{O}(1)$  grammars. It takes time  $O(MT) = \widetilde{O}(k)$ to prepare the binary encoding of each of the committed grammars, and  $\widetilde{O}(k^2)$  to insert it into the Hamming sketch. The update of the active grammars takes  $\widetilde{O}(k)$  time as described below. So in total this step takes  $\widetilde{O}(k^2)$  time.

*Removing a symbol.* Deletion buffer works in manner similar to insertion buffer, we add the removed symbol *a* to the active grammars, but when committing the grammar  $G_{r+1}$ , we use  $F_{\text{KR}}$ -fingerprint of all the grammars  $G_{r-4T+1}, \ldots, G_{r+1}$  to encode grammar  $G_{r-2T+1}$  which is then *removed* from the beginning of the sequence of grammars represented by the Hamming sketch (if r - 2T + 1 > 0), i.e., we update the Hamming sketch to reflect this removal. Similarly to appending a symbol, this step takes time  $\widetilde{O}(k^2)$ .

Active grammar update. The update of active grammars  $G_1^i, \ldots, G_t^i$  when appending *a* is done as follows.  $G_1, \ldots, G_s, G_1^i, \ldots, G_t^i$  represents *ux* so we need to calculate the grammars for *uxa*. We claim that only the active grammars might change: At some point,  $G_s$  became committed so at that time there was *T* active grammars following it. If at that point the grammars together represented a string *z*, by appending more symbols to *z* we cannot change grammars  $G_1, G_1, \ldots, G_s$  according to the first part of Lemma 4.2. So appending *a* to *ux* will affect only the active grammars.

From the analysis in the proof of Lemma 4.2 it follows that for  $\ell \in \{0, ..., 1\}$  if  $B^{ux}(\ell, 1), ..., B^{ux}(\ell, s_{\ell}^{xy})$  is the trace of the decomposition algorithm on ux at level  $\ell$ , and  $B^{uxa}(\ell, 1), ..., B^{uxa}(\ell, s_{\ell}^{xya})$  is the trace on uxa, then their difference spans at most  $\ell(3R + 6)$  last symbols of  $B^{ux}(\ell, 1) \cdots B^{ux}(\ell, s_{\ell}^{xy})$ .

So instead of decompressing the active grammars completely, adding *a* and recompressing them back, we only decompress the necessary part of each trace  $B^{ux}(\ell, 1) \cdots B^{ux}(\ell, s_{\ell}^{xy})$ . Let  $\# \to v_i$  be the starting rule of the active grammar  $G_i$ . Starting from the string  $v_1 \cdot v_2 \cdots v_t$ , for each  $\ell = L, \ldots, 1$ , we iteratively rewrite all level- $\ell$ symbols in the string using the appropriate grammars while only maintaining at most *T* last symbols of the resulting string. (Care has to be taken to maintain information about any sequence  $a^r$  stretching from those *T* last symbols to the left.)

We add *a* to the resulting string and re-apply compress and split procedures for levels 0, 1, ...,  $\ell - 1$  to recompress only the part of the trace affected by modifications. As we perform the compression of symbols we maintain a set *G* of all grammar rules needed for decompression. (We initialize *G* with the union of all rules from the active grammars  $G_1^i, \ldots, G_t^i$  minus the starting rules, and we iteratively add new rules coming from the recompression.) For the recompression we need to know the context of up-to R + 1 symbols preceding the modified part of the trace. On the other hand, the modification can affect the recompression of up-to R + 1 symbols to the left from the left-most modified symbol in the trace. Those R + 1 symbols all happen to be within the decompressed suffix of the trace of size at most *T*.

Eventually, we get a new level-*L* trace  $B^{uxa}(L, s_L^{xya} - t' + 1), \ldots, B^{uxa}(L, s_L^{xya})$ , for some *t'*. Each new grammar  $G_j^{\prime i}$  is obtained by taking the grammar  $G \cup \{\# \to B^{uxa}(L, s_L^{xya} - t' + j)\}$  and removing from it all useless rules. This can be done in time O(|G|). (See Section 2.1).

Overall the update of active grammars on insertion of a single symbol will require  $O(LT) = \widetilde{O}(1)$  evaluations of split hash functions  $H_0, \ldots, H_L$ ,  $O(LT) = \widetilde{O}(1)$  evaluations of compress hash functions  $C_1, \ldots, C_L$ , and  $O(T(LT + \sum_{j=1}^{t} |G_j^i|))$  time to produce the new grammars. As the total size of the grammars is  $\widetilde{O}(k)$  and the time to evaluate  $H_\ell$  at a single point is also  $\widetilde{O}(k)$ , the overall time for the update of active grammars is  $\widetilde{O}(k)$ . We provide a more detailed description of the update procedure in the full version [5].

*Edit distance evaluation.* Consider strings *x* and *y* of length at most *m* and edit distance at most *k*. Consider the rolling sketch  $\operatorname{sk}_{m,k}^{\operatorname{Rolling}}(x)$  for *x* obtained by inserting symbols *ux* and removing symbols *u*, for some  $u \in \Sigma^*$  where  $|ux| \leq m$ . Consider also the rolling sketch for *y* obtained by inserting symbols *vy* and removing symbols *v*, for some  $v \in \Sigma^*$  where  $|vy| \leq m$ . Both sketches should use the same randomness that is to start from the same sketch for empty string.

The rolling sketch for *x* consists of the insertion buffer with committed grammars  $G_{s^x-4T+1}^x, G_{s^x-4T+2}^x, \ldots, G_{s^x}^x$  and with active grammars  $G_1^{ix}, \ldots, G_{t^x}^{ix}$ , and the deletion buffer with committed grammars  $G_{r^x-4T+1}^x, G_{r^x-4T+2}^x, \ldots, G_{r^x}^x$  and active grammars  $G_1^{dx}, \ldots, G_{t^{tx}}^{dx}, t^{tx} \leq T$ . Its Hamming sketch sketches the sequence of grammars  $G_{r^x-2T+1}^x, G_{r^x-2T+2}^x, \ldots, G_{s^x-2T}^x$ . Also for *y*, we have the committed insertion grammars  $G_{s^{y}-4T+1}^y, G_{s^y-4T+2}^y, \ldots, G_{s^y}^y$ , etc.

We extend the notation so for  $j \in \{1, ..., t^x\}$ , we let  $G_{s^x+j}^x$  denote the active grammar  $G_j^{ix}$ , and similarly for y. Let  $d^x = s^x + t^x - r^x$ and  $d^y = s^y + t^y - r^y$ . We assume that the hash functions used to decompose ux and vy into grammars satisfy the probabilistic conclusion of Theorem 3.12. That means that grammars  $G_r^x, ...$ and  $G_r^y, ...$  can be aligned from the right so  $G_j^x$  corresponds to  $G_{j-d^x+d^y}^y$ , for  $j \ge r^x$  (they might not be identical because of the edit operations). Without loss of generality we assume that  $d^x \ge d^y$ .

Before proceeding with the algorithm we first observe that  $d^x - d^y < 2T$ . Let  $p^x \ge r^x + 1$  be the index of the grammar  $G_{p^x}^x$  which produces the first symbol of x when we evaluate all the grammars. Similarly,  $p^y \ge r^y + 1$  is the index of  $G_{p^y}^y$  which produces the first symbol of y. By Lemma 4.2 applied on  $x \leftarrow u$  and  $z \leftarrow x$  we get that  $p^x \le r^x + t'^x + T \le r^x + 2T$ , and similarly  $p^y \le r^y + 2T$ . By our assumption on success of Theorem 3.12,  $s^x + t^x - p^x = s^y - t^y - p^y$ . Hence,  $s^x + t^x - s^y - t^y = p^x - p^y \le r^x + 2T - r^y - 1 \le r^x - r^y + (2T - 1)$ . Thus  $d^x - d^y = s^x + t^x - r^x - s^y - t^y + r^y \le r^x - r^y + (2T - 1) - r^x + r^y \le 2T - 1$ .

If  $d^x < 10T$ , then we can recover all the grammars  $G_{r^x-2T+1}^x$ ,  $G_{r^x-2T+2}^x, \ldots, G_{s^x-2T}^x$  from their Hamming sketch by constructing an auxiliary *dummy* Hamming sketch sk' for a sequence of 1's of length  $(s^x - r^x)M$  and comparing the two sketches. (*M* is the length of the encoding of each grammar.) Their mismatch information reveals all the grammars  $G_{r^x-2T+1}^x, \ldots, G_{s^x-2T}^x$  Since  $d^y \le d^x$ , we can similarly recover all the grammars  $G_{r^y-2T+1}^y, \ldots, G_{s^y-2T}^y$  from their Hamming sketch.

We therefore know all the grammars  $G_{r^{x+1}}^x, G_{r^{x+2}}^x, \ldots, G_{s^{x+tx}}^x$ and  $G_{r^{y+1}}^x, G_{r^{y+2}}^y, \ldots, G_{s^{y+ty}}^y$ . We know grammars  $G_1^{dx}, \ldots, G_{t^{xx}}^{dx}$ and  $G_1^{dy}, \ldots, G_{t^{yy}}^{dy}$  too, that need to be *subtracted* from our grammars. As noted in Section 2.1, for each of the grammars we can calculate its evaluation size. From that information we can easily identify  $p^x$ and  $p^y$ , and shorten the grammars  $G_{p^x}^x$  and  $G_{p^y}^y$  to produce only symbols of x and y, respectively. We can combine all the grammars of x into one grammar  $G^x$ , and all the grammars of y into  $G^y$ , and run the algorithm of Ganesh, Kociumaka, Lincoln and Saha [14] to calculate the edit distance of x and y. Since  $T = \widetilde{O}(1)$ , that will take time  $\widetilde{O}(|G^x| + |G^y| + k^2) = \widetilde{O}(k^2)$ .

If  $d^x \ge 10T$  then we proceed as follows. Clearly,  $d^y \ge 8T$ , so  $s^y - r^y \ge 7T$  and  $s^x - r^x \ge 9T$ . Thus  $G^x_{r^x - 2T + 1}, G^x_{r^x - 2T + 2}, \dots, G^x_{s^x - 2T}$ and  $G^y_{r^y - 2T + 1}, G^y_{r^y - 2T + 2}, \dots, G^y_{s^y - 2T}$  consist of at least 7T grammars each, and those grammars are sketched by their Hamming sketches. Although we assume that there is a correspondence between the grammar  $G^x_j$ , for  $j \ge r^x$ , and  $G^y_{j - d^x + d^y}$  the sequences  $G^x_{r^x - 2T + 1}, \dots, G^x_{s^x - 2T}$  and  $G^y_{r^y - 2T + 1}, \dots, G^y_{s^y - 2T}$  are misaligned in their Hamming sketches by  $d^x - d^y$  grammars. To rectify this misalignment, we prepend  $(d^x - d^y)M$  copies of symbol 1 into the sketch for  $G^y_{r^y - 2T + 1}, \dots, G^y_{s^y - 2T}$ . Furthermore, if  $t^x < t^y$  then we append  $(t^y - t^x)M$  ones into the sketch for  $G^y_{r^y - 2T + 1}, \dots, G^y_{s^y - 2T}$ . The number of sketched grammars. Otherwise if  $t^x > t^y$  then we append  $(t^x - t^y)M$  ones into the sketch for  $G^y_{r^x - 2T + 1}, \dots, G^y_{s^x - 2T}$ . Now we can calculate the mismatch information from the Ham-

Now we can calculate the mismatch information from the Hamming sketches to find out the pairs of grammars  $G_j^x$  and  $G_{j-d^x+d^y}^y$ ,  $j \ge r^x + 1$ , that are different.

If for some  $j \in \{r^x + 1, ..., r^x + 2T\}$ ,  $G_j^x$  and  $G_{j-d^x+d^y}^y$  differ then because we use the Karp-Rabin fingerprint of the two grammars to encode also the neighboring grammars up-to distance 2T, we recover from the sketch all the grammars  $G_j^x$  and  $G_{j-d^x+d^y}^y$ , for  $j = r^{x} + 1, \dots, r^{x} + 2T$ . By counting the evaluation size of each of those grammars and comparing it with the evaluation size of active grammars in deletion buffers of x and y, resp., we identify  $p^x$  and  $p^y$ , and how much the grammars  $G^x_{p^x}$  and  $G^y_{p^y}$  should be shortened to produce only symbols of x and y. After shortening  $G_{p^x}^x$ and  $G_{n^y}^y$  we calculate the edit distance of their evaluation. We sum it up with the edit distance of evaluation of each pair of grammars  $G_j^x$  and  $G_{j-d^x+d^y}^y$ , for  $j > p^x$ , that was identified as mismatch by the Hamming distance sketch or that belongs among the active grammars in insertion buffers of either x or y. There will be at most T mismatched pairs involving the active grammars, and (4T + 1)kpairs identified by the Hamming sketch.

In the remaining case when  $G_j^x$  and  $G_{j-d^x+d^y}^y$  are identical for all  $j \in \{r^x + 1, \dots, r^x + 2T\}$ , we might not be able to recover all those grammars from the Hamming sketches, and we might not be able to identify  $p^x$  and  $p^y$ . However, since  $G_{p^x}^x = G_{p^y}^y$ , we know that the part of x produced by  $G_{p^x}^x$  is either a prefix or suffix of the part of y produced by  $G_{p^y}^y$ . The difference in the size of the two parts is the edit distance of the two parts. The difference is given by the difference between the total evaluation size of active grammars in the deletion buffer of x, and the total evaluation size of active grammars in the deletion buffer of y together with grammars  $G_{r_y-j}^y$ , for  $j = 0, ..., d^x - d^y - 1$ . The latter grammars are in the committed deletion buffer of *y* and they agree with  $G_{r^x+1}^x, \ldots, G_{r^x+d^x-d^y}^x$ . Hence, the edit distance of the parts of *x* and *y* coming from  $G_{p^x}^x$ and  $G_{p^y}^y$  can be determined. All other mismatching pairs of grammars are identified by the Hamming sketch or are among active grammars of the insertion buffers. So we proceed as in the previous case to calculate their contribution to the edit distance of *x* and *y*. The edit distance of x and y is the sum of those edit distances.

We see that in both the cases we need the Hamming sketch to be able to recover at least *T* mismatched grammars at the very end caused by the dummy padding, 4*T* grammars at the beginning corresponding to  $G_{r^x-2T+1}^x, G_{r^x-2T+2}^x, \dots, G_{r^x+2T}^x, 2T$  neighbors of  $G_{r^x+2T}^x$  to the right, and at most (4T + 1)k mismatched grammars caused by the edit operations between *x* and *y*. This is less than M(4T+1)(k+2) which is the number of mismatches our Hamming sketch can recover.

The time needed to compare the sketched strings can be bounded as follows: In total the procedure generates at most O(Tk) pairs of grammars of total size  $\tilde{O}(k^2)$  on which it runs edit distance computation from Proposition 2.1. If those edit distance computations take total time more than  $\tilde{O}(k^2)$  we can terminate them as we know the overall edit distance is larger than k. Recovering differing grammars from the Hamming distance sketch takes time  $\tilde{O}(k') = \tilde{O}(k^2)$ . Their follow-up processing such as counting their evaluation size and shortening them is proportional to their total size which is  $\tilde{O}(k^2)$ . Hence, the time for comparing strings is  $\tilde{O}(k^2)$ . STOC '23, June 20-23, 2023, Orlando, FL, USA

*Failure probability.* The analysis of the failure probability is omitted due to space limitations and can be found in the full version [5].

# **5 TABLE OF PARAMETERS**

Definition	Asymptotics	Meaning	Reference
$R = \log^*  \Gamma  + 20$	log* n	compression	Sec. 2.3
		locality	
$L = \lceil \log_{3/2} n \rceil + 3$	log n	recursion	Sec. 3,
,		depth	Cor. 3.3
D = 110c - R(L+1)k	$k \log n \log^* n$	1/splitting	Sec. 3,
		probability	Lem. 3.4
$S = 30DL\log n + 6$	$k \log^3 n \log^* n$	maximum	Sec. 3,
		grammar	Thm. 3.1
		size	
$M = 3S \cdot \lceil 1 + \log  \Gamma  \rceil$	$k \log^4 n \log^* n$	grammar	Sec. 3.3
		encoding	
		size	
T = L(3R + 6)	$\log n \log^* n$	locality	Sec. 4,
		of suffix	Lem. 4.2
		changes	
$N \ge n^3$	<i>n</i> <sup>3</sup>	F <sub>KR</sub> range	Sec. 3.3
		size	

# ACKNOWLEDGEMENTS

The authors benefited greatly from discussions with Nicole Wein who took part in the initial stages of this project. The second author also benefited from many discussions on edit distance with Mike Saks. We are grateful to Tomasz Kociumaka for providing us with a reference for Proposition 2.1. We thank anonymous reviewers for their comments.

#### REFERENCES

- Alexandr Andoni and Negev Shekel Nosatzki. 2020. Edit Distance in Near-Linear Time: it's a Constant Factor. In 61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, Durham, NC, USA, November 16-19, 2020, Sandy Irani (Ed.). IEEE, 990–1001. https://doi.org/10.1109/FOCS46700.2020.00096
- [2] Arturs Backurs and Piotr Indyk. 2015. Edit Distance Cannot Be Computed in Strongly Subquadratic Time (Unless SETH is False). In Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing (Portland, Oregon, USA) (STOC '15). ACM, New York, NY, USA, 51–58.
- [3] Tuğkan Batu, Funda Ergun, and Cenk Sahinalp. 2006. Oblivious String Embeddings and Edit Distance Approximations. In Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithm (Miami, Florida) (SODA '06). Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 792–801.
- [4] Djamal Belazzougui and Qin Zhang. 2016. Edit Distance: Sketching, Streaming, and Document Exchange. In 2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS). 51–60. https://doi.org/10.1109/FOCS.2016.15
- [5] Sudatta Bhattacharya and Michal Koucký. 2023. Locally consistent decomposition of strings with applications to edit distance sketching. *CoRR* abs/2302.04475 (2023). arXiv:2302.04475
- [6] Or Birenzwige, Shay Golan, and Ely Porat. 2020. Locally Consistent Parsing for Text Indexing in Small Space. In Proceedings of the 2020 ACM-SIAM Symposium on Discrete Algorithms, SODA 2020, Salt Lake City, UT, USA, January 5-8, 2020, Shuchi Chawla (Ed.). SIAM, 607–626. https://doi.org/10.1137/1.9781611975994.37
- [7] Joshua Brakensiek and Aviad Rubinstein. 2020. Constant-factor approximation of near-linear edit distance in near-linear time. In Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2020, Konstantin Makarychev, Yury Makarychev, Madhur Tulsiani, Gautam Kamath, and Julia Chuzhoy (Eds.). ACM, 685–698. https://doi.org/10.1145/3357713.3384282
- [8] Diptarka Chakraborty, Debarati Das, Elazar Goldenberg, Michal Koucký, and Michael E. Saks. 2018. Approximating Edit Distance within Constant Factor in Truly Sub-Quadratic Time. In 59th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2018. 979–990. https://doi.org/10.1109/FOCS.2018.00096

- [9] Diptarka Chakraborty, Elazar Goldenberg, and Michal Koucký. 2016. Streaming algorithms for embedding and computing edit distance in the low distance regime. In Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016. 712–725.
- [10] Raphaël Clifford, Tomasz Kociumaka, and Ely Porat. 2019. The streaming kmismatch problem. In Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2019. SIAM, 1106–1125. https://doi.org/10.1137/1. 9781611975482.68
- [11] Richard Cole and Uzi Vishkin. 1986. Deterministic coin tossing and accelerating cascades: micro and macro techniques for designing parallel algorithms. In Proceedings of the eighteenth annual ACM symposium on Theory of computing (STOC). 206–219. https://doi.org/10.1145/12130.12151
- [12] Graham Cormode and S. Muthukrishnan. 2002. The string edit distance matching problem with moves. In Proceedings of the Thirteenth Annual ACM-SIAM Symposium on Discrete Algorithms, January 6-8, 2002, San Francisco, CA, USA. 667–676.
- [13] Joan Feigenbaum, Yuval Ishai, Tal Malkin, Kobbi Nissim, Martin J Strauss, and Rebecca N Wright. 2006. Secure multiparty computation of approximations. ACM transactions on Algorithms (TALG) 2, 3 (2006), 435–472.
- [14] Arun Ganesh, Tomasz Kociumaka, Andrea Lincoln, and Barna Saha. 2022. How Compression and Approximation Affect Efficiency in String Distance Measures. In Proceedings of the 2022 ACM-SIAM Symposium on Discrete Algorithms, SODA. 2867–2919. https://doi.org/10.1137/1.9781611977073.112
- [15] Szymon Grabowski. 2016. New tabulation and sparse dynamic programming based techniques for sequence similarity problems. *Discrete Applied Mathematics* 212 (2016), 96–103.
- [16] Ce Jin, Jelani Nelson, and Kewen Wu. 2021. An Improved Sketching Algorithm for Edit Distance. In 38th International Symposium on Theoretical Aspects of Computer Science, STACS 2021, (LIPIcs, Vol. 187). 45:1–45:16. https://doi.org/10.4230/LIPIcs. STACS.2021.45
- [17] Hossein Jowhari. 2012. Efficient Communication Protocols for Deciding Edit Distance. In Algorithms - ESA 2012 - 20th Annual European Symposium, Ljubljana, Slovenia, September 10-12, 2012. Proceedings. 648–658.
- [18] Richard M. Karp and Michael O. Rabin. 1987. Efficient randomized patternmatching algorithms. *IBM Journal of Research and Development* 31, 2 (1987), 249-260. https://doi.org/10.1147/rd.312.0249
- [19] Tomasz Kociumaka, Ely Porat, and Tatiana Starikovskaya. 2021. Small-space and streaming pattern matching with k edits. In 2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS). 885–896. https://doi.org/10.1109/ FOCS52979.2021.00090
- [20] Michal Koucký and Michael E. Saks. 2020. Constant factor approximations to edit distance on far input pairs in nearly linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2020*, Konstantin Makarychev, Yury Makarychev, Madhur Tulsiani, Gautam Kamath, and Julia Chuzhoy (Eds.). ACM, 699–712. https://doi.org/10.1145/3357713.3384307
- [21] Eyal Kushilevitz, Rafail Ostrovsky, and Yuval Rabani. 1998. Efficient search for approximate nearest neighbor in high dimensional spaces. In Proceedings of the thirtieth annual ACM symposium on Theory of computing. 614–623.
- [22] Gad M. Landau, Eugene W. Myers, and Jeanette P. Schmidt. 1998. Incremental String Comparison. SIAM J. Comput. 27, 2 (April 1998), 557–582.
- [23] Nathan Linial. 1987. Distributive Graph Algorithms-Global Solutions from Local Data. In 28th Annual Symposium on Foundations of Computer Science, FOCS. IEEE Computer Society, 331–335. https://doi.org/10.1109/SFCS.1987.20
- [24] Nathan Linial. 1992. Locality in Distributed Graph Algorithms. SIAM J. Comput. 21, 1 (1992), 193–201. https://doi.org/10.1137/0221015
- [25] William J. Masek and Michael S. Paterson. 1980. A faster algorithm computing string edit distances. J. Comput. System Sci. 20, 1 (1980), 18 – 31.
- [26] Rafail Ostrovsky and Yuval Rabani. 2007. Low distortion embeddings for edit distance. J. ACM 54, 5 (2007), 23. https://doi.org/10.1145/1284320.1284322
- [27] Ely Porat and Ohad Lipsky. 2007. Improved Sketching of Hamming Distance with Error Correcting. In Combinatorial Pattern Matching, 18th Annual Symposium, CPM, Vol. 4580. Springer, 173–182. https://doi.org/10.1007/978-3-540-73437-6\_19
- [28] Süleyman Cenk Sahinalp and Uzi Vishkin. 1994. Symmetry breaking for suffix tree construction. In Proceedings of the Twenty-Sixth Annual ACM Symposium on Theory of Computing, 23-25 May 1994, Montréal, Québec, Canada. ACM, 300–309. https://doi.org/10.1145/195058.195164
- [29] Robert A. Wagner and Michael J. Fischer. 1974. The String-to-String Correction Problem. J. ACM 21, 1 (Jan. 1974), 168–173.
- [30] Haoyu Zhang and Qin Zhang. 2019. MinJoin: Efficient Edit Similarity Joins via Local Hash Minima. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Anchorage, AK, USA) (KDD '19). Association for Computing Machinery, New York, NY, USA, 1093–1103. https://doi.org/10.1145/3292500.3330853

Received 2022-11-07; accepted 2023-02-06