

# APPLICATION OF SEMI-LOCAL SA TO APPROXIMATE PATTERN MATCHING

NIKITA MISHIN\* AND DANIIL BEREZUN†

**Abstract.** In the paper we study an application of semi-local sequence alignment (sa) algorithms to approximate pattern matching problem. We both developed two new algorithms as well as improved the existing near duplicate search algorithm (Programming and Computer Software'19). The key idea behind the algorithms is a usage of the underlying algebraic structure of semi-local sa (Tiskin, 2007) together with a novel data structure for submatrix maximum queries in Monge matrices (TALG'20). We also show that the improved near duplicate search algorithm not only has a better complexity but also preserves all declared properties. We show that the presented algorithms running time and space complexity are  $O(\max(|t||p|, \frac{|t|\log^2|t|}{\log\log|t|}))$  and  $O(|t|)$  for the first one and  $O(\max(|t||p|, |t|\log|t|))$  and  $O(|t|\log|t|)$  for the last two, respectively, where  $t$  is a text,  $p$  — pattern, and  $v = O(1)$  is denominator of normalized mismatch score for semi-local sequence alignment.

**Key words.** semi-local lcs, monge matrix, range queries, approximate matching, near-duplicate detection

**AMS subject classifications.** 68Q25, 68R10, 68U05

**1. Introduction.** Approximate string matching is an important task in many fields such as computational biology, signal processing, text retrieval and etc. It also refers to a duplicate detection subtask.

In general form it formulates as follows: Given some pattern  $p$  and text  $t$  need to find all occurrences of pattern  $p$  in text  $t$  with some degree of similarity.

There are many algorithms that solve the above problem. Nonetheless, the number of algorithms sharply decreases when the algorithm needs to meet some specific requirements imposed by running time, space complexity or specific criterion for the algorithm itself. For example, recently there was developed an approach for interactive duplicate detection for software documentation [2]. The core of this approach is an algorithm that detects approximate clones of a given user pattern with a specified degree of similarity. The main advantage of the algorithm is that it meets a specific requirement of completeness. Nonetheless, it has an unpleasant time complexity.

The algorithm for approximate detection utilizes mainly algorithm for solving the longest commons subsequence ( $LCS$ ) problem. The longest common subsequence is a well-known fundamental problem in computer science that also has many applications of its own. The major drawback of it that it shows only the global similarity for given input strings. For many tasks, it's simply not enough. The approximate matching is an example of it.

There exist generalization for  $LCS$  called *semi-local LCS* [] which overcome this constraint. The effective theoretical solutions for this generalized problem found applications to various algorithmic problems such as bla bla add cited. For example, there has been developed algorithm for approximate matching in the grammar-compresed strings[].

Although the algorithms for *semi-local LCS* have good theoretical properties, there is unclear how they would behave in practice for a specific task and domain.

To show the applicability of semi-local lcs on practice we developed several algorithms based mainly on it and the underlying algebraic structure. As well as devel-

---

\*Saint Petersburg State University, Russia (mishinnikitam@gmail.com).

†IntelliJ Labs Co. Ltd., Saint Petersburg, Russia (daniil.berezun@jetbrains.com).

opening new algorithms we improve and significantly outperform the existing one for interactive duplicate detection for software documentation [1]. It should be noted that improvement preserves all properties of this algorithm. **Do we need to state that ant algo is slow for current strucute of algorithm**

The paper is organized as follows. Blablabla [2], our new algorithm is in [3], experimental results are in [4], and the conclusions follow in [5].

## 2. Preliminaries.

**2.1. Approximate pattern matching.** The approximate pattern matching problem (*AMatch*) defined as follows. Given text  $t$ , pattern  $p$  and some threshold  $h$  the *approximate pattern matching* problem ask for all substrings from text  $t$  that have similarity score with given pattern  $p$  at least  $h$  according to some similarity function  $g$ .

There exist different kinds of extensions and particular cases of this problem. For example, *complete approximate pattern matching* (*CompleteMatch*) that ask for substrings of text  $t$  that are exact clones of pattern  $p$ . The approach for this special case of *AMatch* is usage of well-know algorithms such as Aho-Korasic, BouerMurr, Knuth-Morris-Pratt, and so on. The latter one have optimal running time complexity of  $O(|p| + |t|)$  for *CompleteMatch* problem [6]. *Approximate pattern matching with  $k$  mismatches* is an another example of special case [7]. The search of *pattern with wild-card symbols* or search set of patterns in text  $t$  [8], multidimensional *AMatch* [9], search with lenght constraint of detected duplicates [10] are examples of such extension. There exits many more examples of constraints, extensions and special cases of *AMatch* problem [11].

The one of the common approach to solve approximate pattern matching is the usage of solution of string similarity problem. Latter represent a set of fundamental problems such as *edit distance*, *longest common subsequence*, *sequence alignment*. In this paper we primarily focuses on the usage of latter two when developing algorithms.

Recently there have been developed algorithm for solving interesting extension of *AMatch* problem with length constarint [12]. Although their algorithm have poor result in terms of running time complexity, the proposed solution possesses a completnessess proprety i.e it founds *all* non-intersected clones of pattern  $p$  with specified similarity threshold and length constraint on matching substrings. Thus, this algorithm is an subject of interest in this paper. The complete description of algorithm and its improved version may be found in section [13] respectively.

**2.2. Semi-local lcs.** First of all we give definition of *lcs* and *sa*.

**DEFINITION 2.1.** *Given two strings  $a$  and  $b$  the longest common subsequence (LCS) problem ask for the maximal length of the longest common subsequence of  $a$  and  $b$  ( $lcs(a, b)$ ).*

In other words, *LCS* problem asks about maximal *lcs* score of two given string  $a$  and  $b$  ( $lcs(a, b)$ ).

**DEFINITION 2.2.** *Given two strings  $a$  and  $b$  and scoring scheme  $w = (w_+, w_0, w_-)$  the sequence alignment (SA) problem ask for the maximal alignment score between  $a$  and  $b$  ( $sa(a, b)$ ).*

Scoring scheme determines how calulate alignment score of two aligned sequences. If pair of character in aligned sequences are matches (equals) then this pair contributes to final alignemnt score  $w_+$ , if their mismatch it contributes  $w_0$ . If symbol  $\alpha$  of one of the sequences is not aligned with any other symbol from other sequence it means that

92  $\alpha$  is aligned with *gap*. Thus, this pair contributes  $w_0$ . The scoring scheme calculates  
 93 as follows:

$$94 \quad (2.1) \quad sa(a, b, w) = w_+k^+ + w_0k^0 + w_-(|a| + |b| - 2k^+ - 2k^0) = \\ k^+(w_+ - 2w_-) + k^0(w_0 - 2w_-) + w_-(|a| + |b|)$$

95 The  $k^+$  states for the number of matching symbols,  $k^-$  — mismatched symbols.

96 Note that *LCS* is a special case of *SA* when scoring scheme is  $(1, 0, 0)$ .

97 Both described problems are solved by classical dynamic programming algorithm  
 98 and have running time complexity  $O(|a||b|)$ . *LCS* and *SA* allow you to find how much  
 99 whole given strings are similar i/e how similar two string in a global sense.

100 In many cases, this is not enogh. There also exist fully local version of these  
 101 problems and semi-local one. The last one is in sight of this paper due to natural  
 102 applicability to approximate pattern matching.

103 **2.3. Semi-local lcs.** Given two strings  $a$  and  $b$  the semi-local lcs is asks about  
 104 lcs scores for following:

- 105 1. *string-substring*: whole  $a$  against every substring of  $b$
- 106 2. *substring-string*: whole  $b$  against every substring of  $a$
- 107 3. *prefix-suffix*: every prefix of  $a$  against every suffix of  $b$
- 108 4. *suffix-prefix*: every prefix of  $b$  against every suffix of  $a$

109 The following *semi-local lcs matrix* associated with the defined *semi-local lcs*.

110 DEFINITION 2.3. The semi-local lcs matrix  $H_{a,b}$  for strings  $a, b$  defined as follows:

$$111 \quad (2.2) \quad H_{a,b}[i, j] = if(j \leq i)j - ielse lcs(a, b^{pad}[i, j])$$

112 where  $i \in [-|a| : |b|]$ ,  $j \in [0 : |a| + |b|]$  and  $b^{pad} = ?^{|a|}b?^{|a|}$ ,  $?$  — wildcard symbol that  
 113 matches any other symbol.

114 The semi-local lcs matrix  $H_{a,b}$  comprises from four quadrant associated with described  
 115 subproblems:

$$116 \quad (2.3) \quad H_{a,b} = \begin{bmatrix} H_{a,b}^{suf-pre} & H_{a,b}^{sub-str} \\ H_{a,b}^{str-sub} & H_{a,b}^{pre-suf} \end{bmatrix}$$

117 DEFINITION 2.4. Matrix  $H$  called (anti) Monge matrix if

$$118 \quad H[i, j] + H[i', j'] (\geq) \leq H[i, j'] + H[i', j], \forall i \leq i', j \leq j'$$

119 DEFINITION 2.5. Let  $H[0 : m, 0 : n]$  be a matrix.  $H^\square[0 : m - 1, 0 : n - 1]$   
 120 constructed as a result of taken cross difference between secondary and first diagonal  
 121 for all adjacent 2 by 2 squares called cross-difference matrix of  $H$

122 DEFINITION 2.6. Matrix  $H$  called unit anti Monge matrix if  $H$  is (anti) Monge  
 123 matrix and its cross-difference matrix  $(- )H^\square$  is permutation matrix.

124 The example of unit anti Monge matrix is following:

$$125 \quad (2.4) \quad \begin{bmatrix} 0 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 1 & 1 \end{bmatrix}^\square = \begin{bmatrix} (2+0) - (1+0) & (3+1) - (2+2) \\ (1+0) - (1+0) & (2+1) - (1+1) \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

126 DEFINITION 2.7. Let  $H[0 : m - 1, 0 : n - 1]$  be a matrix.  $H^{\nearrow}[0 : m, 0 : n]$   
 127 constructed as sum of element that lies below and left given cell  $i, j$  in matrix  $H$  called  
 128 dominance-sum matrix of  $H$

The example dominance sum matrix:

$$(2.5) \quad \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^{\nearrow} = \begin{bmatrix} 0+0+0 & 1 & 1+1 \\ 0+0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 2 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

In [?] is proved that  $H_{a,b}$  is unit anti Monge. Also it is proved that this matrix may be decomposed to permutation matrix i.e into *cross-difference* matrix. It allows to store  $H_{a,b}$  implicitly and query any element of  $H_{a,b}$  via dominance sum query (orthogonal range queries). Thus, there may be several ways to storing matrix  $H_{a,b}$  or one of it quadrant implicitly. A simple storing of two list of permutation gives  $O(|a| + |b|)$  space and time complexity with  $O(|a| + |b|)$  orthogonal range queries (need to check how many points dominated by given point), whereas more sophisticated approach requires  $O(|a| + |b|)$  space with  $O((|a| + |b|)\sqrt{\log(|a| + |b|)})$  preprocessing time and allows to query any point of  $H$  in  $O(\frac{\log(|a| + |b|)}{\log \log(|a| + |b|)})$  time.

The one useful proposition of  $H$  is following.

PROPOSITION 2.8. *Given a permutation matrix  $P$  and the value  $P^{\nearrow}[i; j]$ , the values  $P^{\nearrow}[i + -1; j]$ ,  $P^{\nearrow}[i; j + -1]$ , where they exist, can be queried in time  $O(1)$ .*

We particularly interesting in lower left quadrant that refers to string substring problem:

$$(2.6) \quad H_{a,b}^{str-sub}[i, j] = lcs(a, b[i, j]), i, j \in [0, |b|]$$

There exists several algorithms **second on, recursive not described as i see** that solve *semi-local lcs*. Both have the optimal running time  $O(|a||b|)$  for given dynamic problem **Impossibility faster then pt.**

**2.4. Semi-local sa.** The semi-local sequence alignment (sa) is a generalization of semi-local lcs in same sense as sequence alignment is generalization of lcs.

Given two strings  $a$  and  $b$  and scoring scheme  $w = (w_+, w_0, w_-)$  the semi-local sa asks about sa scores for following:

1. *string-substring*: whole  $a$  against every substring of  $b$
2. *substring-string*: whole  $b$  against every substring of  $a$
3. *prefix-suffix*: every prefix of  $a$  against every suffix of  $b$
4. *suffix-prefix*: every prefix of  $b$  against every suffix of  $a$

The associated matrix for *semi-local sa* is defined analogously as for *semi-local lcs*.

The approach for solving *semi-local sa* is as follows. The problem reduced to *semi-local lcs*. First, note that scoring scheme in 2.1 may be simplified by so called normalization:

$$(2.7) \quad w = (w_+, w_0, w_-) \rightarrow (w_+ + 2x, w_0 + 2x, w_- + x) = \left( \frac{w_+ + 2x}{w_+ + 2x}, \frac{w_0 + 2x}{w_+ + 2x}, \frac{w_- + x}{w_+ + 2x} \right)_{x=-w_-} = \left( 1, \frac{\mu}{v}, 0 \right)$$

The resulted scoring scheme  $w_{normalized} = (1, \frac{\mu}{v}, 0)$  called normalized scoring scheme.

Then to query initial score  $sa$  for scoring scheme  $w$  knowing  $sa_{normalized}$  for  $w_{normalized}$  you need to apply reverse regularization:

$$(2.8) \quad sa(a, b, w) = sa_{normalized}(w_+ - 2w_-) + w_-(|a| + |b|)$$

The blown-up technique is applied after reducing scoring scheme which increases both input strings in  $v$  times. Nonetheless, only one of the described algorithm time complexity increases in  $v^2$  times, the second one only  $v$ . **Bad sentecnce**. The space complexity also increses by factor  $v$ .

For detailed description we refer readers to TISKIN BOOK[].

**2.5. Range maximum/minimum queries.** Range maximum/minimum queries (rmq) (submatrix query) refers to search maximum/minimum element in submatrix  $[i_1 : i_2] \times [j_1 : j_2]$  of given matrix  $M$  of size  $n \times n$ . The associated data strucutre that can report maximum/minimum element in any submatrix query called *range maximum/minimum data structure*.

For the generic case of Matrix  $M$  it is not possible to achieve running time faster then  $O(n^2)$  due to fact that storing matrix  $M$  requires  $O(n^2)$ .

Nonetheless, the situation is changed if we consider special cases such as Monge matrices. There have been several researches over several decades about rmq on monge matrices [].

The recent research achives following result[].

**THEOREM 2.9.** [] *Given an  $n \times n$  Monge matrix  $M$ , a data structure of size  $O(n)$  can be constructed in  $O(n \log n)$  time to answer submatrix maximum queries in  $O(\log \log n)$  time when random access to Monge matrix is  $O(1)$ .*

**THEOREM 2.10.** [] *Given an  $n \times n$  staircase<sup>1</sup> Monge matrix  $M$ , a data structure of size  $O(n)$  can be constructed in  $O(n \log n)$  time to answer submatrix maximum queries in  $O(\log \log n)$  time when random access to Monge matrix is  $O(1)$ .*

**THEOREM 2.11.** [] *Given an  $n \times n$  partial Monge matrix<sup>2</sup>  $M$ , a data structure of size  $O(n)$  can be constructed in  $O(n \log n)$  time to answer submatrix maximum queries in  $O(\log \log n)$  time when random access to Monge matrix is  $O(1)$ .*

The above results applies both to range minimum queries and to monge matrices with non-constant access  $O(\beta)$  to queries. The latter one, costs in increased construction time and query time by factor  $\beta$ .

**2.6. Near-duplicate detection algorithm.** First, we denote several parameters, that is used in algorithm [].  $k$ — constant in interval  $[\frac{1}{\sqrt{3}}, 1]$  that set similarity measure. A window  $w$  of size  $L_w = |p|/k$  is to process text  $t$  with sliding window of one symbol step.  $k_{di} = |p| * (\frac{1}{k} + 1)(1 - k^2)$  is threshold value for edit distance.  $I$  — interval of size  $[|p|k, \frac{|p|}{k}]$  that set boundaries for lenght of matching substrings.  $d_{di}$  — function that measure similarity between two strings.

The algorithm comprises of three phases.

At the first phrase text  $t$  is processed with sliding window of size  $L_w$  with one symbol step. Further, substrings that corsepond to window  $w$  compared using edit distance<sup>3</sup> and if  $d_{di}(p, t_w) \leq k_{di}$  i.e close enough, then they saved to set  $W_1$  to be further proceeded.

On the second phase each of the detected substrings in  $W_1$  are shrunk i.e they lenght could be decreased. More preciesely, within each of the element of  $W_1$  the largest one substring with legnth fall in  $I$  that most similar to pattern  $p$  according to  $d_{di}$  is selected. The set  $W_2$  is a result of this phase.

<sup>1</sup>Defintion

<sup>2</sup>Definition

<sup>3</sup>Authors of [] used lcs edit distance — where operations substituiou,removoal, addition of one symbol costs 2,1,1 respectively

At the final third phase set  $W_2$  iterated over to remove elements that fully contains in other elements of  $W_2$  or duplicates.

*Running time analysis. 1st phase.* The first phase requires at most  $O(|t||p|^2)$  due to fact that computing cost of edit distance is  $|p|^2$  for strings of size  $O(|p|)$  and we need to process  $O(|t| - \frac{|p|}{k}) = O(|t|)$  windows. *2nd phase.* The cardinality of set  $W_1$  at worst case be  $O(|t|)$ . Thus, first loop will be iterated over  $O(|t|)$  times. The second loop refers to iteration over  $I$  set with cardinality  $O(\frac{|p|}{k} - |p|k) = O(|p|)$ . The third loop requires at most  $O(|p|)$  since we again goes with sliding window. The compare operation requires at most  $O(|p|^2)$  since both string of size  $O(|p|)$ . Thus, the total running time complexity of second phase is  $O(|p|^4|t|)$  at worst case.

*3rd phase.* Note that cardinality of set  $W_2 == W_1$  and thus at worst case is  $O(|t|)$ . Then, this phase requires  $O(|t| \log |t|)$  time.

Thus, the total time complexity of algorithms estimates as  $O(\max(|t||p|^4, |t| \log |t|))$

THEOREM 2.12. [?] *Completeness theorem*

Add luciv pseudocode

### 3. Related work. ?????

could mention about approximation. Need discuss

**4. Algorithm for near duplicate detection.** We now describe an improved version of Luciv et.al. algorithm [2] by utilizing a *semi-local sa* solution. Then we present proof that improved version preserves completeness property. It is achieved by imitating all phases of the algorithm.

**4.1. Algorithm description.** The algorithm comprises three phases as in [2]. At phase one (Lines 1-3) *semi-local sa* problem is solved for the pattern  $p$  against the whole text  $t$ . This solution provides access to the string-substring matrix  $H_{p,t}^{str-sub}$  which allows performing fast queries of *sa* score for pattern  $p$  against every substring of text  $t$ . We apply implicitly transposition and inverse operation on  $H_{p,t}^{str-sub}$ :

$$(4.1) \quad M[j, i] := -H_{p,t}^{str-sub}[i, j]$$

Note that, transposition operation preserves (*anti*) *Monge* property whereas inverse operation make *anti Monge* matrix *Monge* and vice versa. So, matrix  $M$  is *Monge* matrix.

The second phase comprises several steps (Lines 4-6). First, we want to get for each prefix of the text  $t$  the longest suffix that has the highest similarity with the given pattern  $p$  with the following constraint. The lengths of obtained suffixes should be in  $|p| * k \dots \frac{|p|}{k}$  interval where  $k \in [\frac{1}{\sqrt{3}}, 1]$ . It could be done in several ways. For example, direct pass through diagonal with width  $w := \frac{|p|}{k} - |p| * k = |p|(\frac{1}{k} - k)$  in  $H_{p,t}^{str-sub}$  (see fig) or in  $M$  (see fig). The other approach is the following. Note that in  $M$  is *Monge matrix* and indices are swapped. It allows us to descry this diagonal as approximately  $|t|$  square windows of size  $w \times w$  i.e a sliding window of step 1 that goes diagonally. Because of length constraint we only interesting in elements that lie in the main diagonal and below it (remember, transposition) in these submatrices  $w \times w$ . Each of these  $W := w \times w$  matrix is *Monge matrix* by definition (as a submatrix of *Monge matrix*). This implies that  $W$  also totally monotone. If we set to  $+\inf$  for elements that lie above the main diagonal that result matrix will remain totally monotone. Thus, we can apply *SMAWK* algorithm to this matrix to find a leftmost element that has a minimum in a given row with a corresponding column position.

255 For our case leftmost means that for each prefix algorithm will detect longest suffix  
 256 (remember that  $M$  is transposed  $H_{p,t}^{str-sub}$ ).

257 The second step, it is one-way pass through these suffixes with a sliding window  
 258 of size  $\frac{|p|}{t}$  to find for each window most similar suffix with the longest length. Then  
 259 the resulting set is filtered out that remaining suffixes have a score greater or equal  
 260 to given threshold  $-k_{di}$ .

261 The third phase is the same as in [2] (Lines 8-12).

---

**Algorithm 4.1** PATTERN BASED NEAR DUPLICATE SEARCH ALGORITHM  
 VIA SEMI-LOCAL SA

---

Input: pattern  $p$ , text  $t$ , similarity measure  $k \in [\frac{1}{\sqrt{3}}, 1]$

Output: Set of non-intersected clones of pattern  $p$  in text  $t$

---

$$(4.2) \quad k_{di} = |p| * (\frac{1}{k} + 1)(1 - k^2)$$

$$(4.3) \quad L_w = \frac{|p|}{k}$$

$$(4.4) \quad w = |p|(\frac{1}{k} - k)$$

Pseudocode:

```

1:  $W = semilocalsa(p, t)$  {1st phase}
2:  $H_{p,t}^{str-sub} = semilocalsa(p, t).stringSubstringMatrix$ 
3:  $M[j, i] = -H_{p,t}^{str-sub}[i, j]$ 
4:  $sufixes = processDiagonal(M, L)$  {2d phase}
5:  $W_2 = SuffixMaxForEachWindow(sufixes, L_w)$ 
6:  $filter(W_2, k_{di})$ 
7:  $W_3 = UNIQUE(W_2)$  {3rd phase unchanged}
8: for  $w \in W_3$  do
9:   if  $\exists w' \in W_3 : w \subset w'$  then
10:     remove  $w$  from  $W_3$ 
11:   end if
12: end for
13: return  $W_3$ 
```

---

262 **THEOREM 4.1.** *Algorithm 4.1 could be solved in  $\max(O(|t| * |p|), O(|t| * \log |t|))$*   
 263 *time with  $O(|t| \log |t|)$  additional space where  $p$  is pattern,  $t$  is text when  $|p| \leq |t|$ ,*  
 264  *$v = O(1)$  where  $v$  is denominator of normalized mismatch score for semi-local sa*  
 265  *$w_{normalized} = (1, \frac{\mu}{v}, 0)$ .*

266 For each phases of algorithm we provide it's time and space bounds.

267 *First phase.* . We will store solution  $H$  of *semi-local sa* by decomposing it to  
 268 permutation matrix  $P$  of size  $O(v * |t| \times v * |t|)$  (Lines 1-3) (ref add). The permutation  
 269 matrix can be stored via two permutations of size  $v * |t|$  for column and rows. It is  
 270 simply two lists of size  $v * |t|$ . Then, to random access query in specific position  $i, j$  of  
 271 matrix  $H$  we need to check how many points dominated by  $i, j$ . It is just pass through  
 272 permutations that requires  $O(v * |t|)$ . Thus, the total time and space complexity of



1st phase is  $O(v * |p| * |t|)$  (time complexity for solving *semi-local sa*) and  $O(v * |t|)$ . Given  $v = O(1)$  we have  $O(|p| * |t|)$  and  $O(|t|)$  respectively. Also random access query for our case is  $O(|t|)$ .

*Second phase.* We omit  $k$  factor in analysis because when  $k \in [\frac{1}{\sqrt{3}}, 1]$   $O(k) = 1$

We will use the first approach described in the algorithm description for this phase. First, although the random access query to a matrix element requires  $O(|t|)$ . We only need one such query to step on the diagonal. Precisely, to the cell that represents substring  $t_{0, |p| * k}$ , starting at zero position and ending in  $|p| * k$  position. Further we use **Theorem about adjacent cell query** that allows us to perform  $O(1)$  access to adjacent elements for given  $i, j$  cell in matrix  $M$ . Thus, we can visit each cell in the desired diagonal of size at most  $O(|t|) * O(|p|)$  in  $O(|t| * |p|)$  time in the following way. Process row  $i'$  with starting  $j'$  (recall it cell by  $M[i', j']$ ) position (go right i.e increment  $j'$ ) until  $i' - j' \geq |p| * k$ . Then shift by one  $i'$  down and  $j'$  to right by one if needed (see picture **This about the top left corner**).

When we pass through a slice of the specific column, we also will find the longest suffix with the highest similarity simply by checking elements twice. First for detect maximum score, second for detect the longest suffix among those who have this score. Thus, for storing for each prefix its longest suffix we need additionally  $O(|t|)$  space. Also for each substring of length  $\frac{p}{k}$  we store similarity score by querying them during diagonal passage because they lie also on this diagonal. Let's denote it by  $C$ . At the end of *processDiagonal* we will have  $O(t)$  suffixes that require  $O(t)$  space for storing them. Then, *processDiagonal* requires  $O(|t| + |t| * |p|) = O(|t| * |p|)$  time for processing diagonal with  $O(|t| + |p|) = O(|t|)$  additional space.

Further (Line 5), we need to find longest suffix within  $O(|p|)$  window with step one in list of size  $|t|$  with additional condition that within each window of size  $O(\frac{|p|}{k} - |p| * k) = O(|p|)$  the suffix with length  $\frac{p}{k}$  have similarity score at least  $-k_{di}$ . It is simply a one-way pass-through list of suffixes where the processing of each window requires at most  $O(|p| + 1) = O(|p|)$ . More precisely, first, we check that for current window of size  $O(|p|)$  associated suffix has similarity not less then given threshold  $k_{di}$ . It is simply lookup for a specific element in  $C$  with  $O(1)$ . If that true, then we need  $O(p)$  lookups within *suffixes* to query the most similar and longest one. The total number of such windows at most  $O(|t|)$ . Thus, *SuffixMaxForEachWindow* requires  $O(|t|) * O(|p|) = O(|t| * |p|)$  time with  $O(|t|)$  space for storing suffix for each window.

The filtering process (Line 6) is a one-way pass through a list of suffixes  $W_2$ . It requires at most  $O(t)$  time.

As we see, the total running time and space complexity of the second phase is  $O(|t| * |p|)$  and  $O(|t|)$  respectively.

*Third phase.* The third phase remains unchanged, thus have the same time and space-bound. Note that it possible to perform this phase in-place during a second phase which make the algorithm even faster i.e decrease space and time complexity to  $O(|t|)$  and  $O(|t| * |p|)$ . The third phase is  $O(|t| \log |t|)$  at most both for space and running time complexity.

Thus, the total running time is  $\max(O(tp), O(t \log t))$  and space complexity  $t \log t$ . **It be good if we also improve third phase))**

**THEOREM 4.2.** *Algorithm 4.1 preserves completeness property of algorithm [2].*

First we show, equivalence between similarity functions, then we show that set  $W_2$  from algorithm ?? equals to set  $W_2$  from algorithm 4.1. Let be  $A_1$  a set  $W_2$  from algorithm ?. Let be  $A_2$  a set  $W_2$  from algorithm 4.1. We will show that  $A_2 = A_1$ .



322 *First part. Take from dimpla*

323 *Second part.* At first algorithm ?? pass through text  $t$  with sliding window to  
 324 detect those fragments which has similarity above given threshold  $k_{di}$  with size  $\frac{p}{k}$ .  
 325 Then within these fragments algorithm detects longest suffixes most similar to pattern  
 326  $p$  with size within  $pk \dots \frac{p}{k}$  interval. That how  $A_1$  constructed.

327 The second algorithm 4.1 proceed in similar way but it first detects longest suffixes  
 328 with size in  $pk \dots \frac{p}{k}$  interval for each prefix of text  $t$ . Then it proceeds in a such way  
 329 that for each window the longest suffix it detected that have similarity above given  
 330 threshold  $h$  for current window of size  $\frac{p}{k}$ . That how  $A_1$  constructed.

331 Thus,  $A_1 = A_2$  by resulting equivalence of construction.

332 Note that set  $A_1$  contains only those fragments of size  $\frac{p}{k}$  from text  $t$  that close  
 333 enough to pattern  $p$  i.e

334 The fragment from  $W_1$  then shrunked. It means that after second phase set  $W_2$   
 335 will have size of  $W_1$ .

336 **5. CutMax a new approximate mathing algorithm.** We now describe sev-  
 337 eral algorithms that heavily based on semi-local lcs and it's underlying algebraic  
 338 structure.

339 The first algorithm 5.2 refers to following constraint. There should be found all  
 340 non-intersected clones  $\tau_i$  of pattern  $p$  from text  $t$  that has the highest similarity score  
 341 on the uncovered part of the text  $t$  i.e algorithm should perform greedy choice at each  
 342 step. This is a more intuitive approach i.e like looking for the most similar one every  
 343 time. *Formally:*

$$344 \quad (5.1) \quad \tau_i = \arg \max_{l, r \in (t \cap (\cup_{j=1}^{i-1} \tau_j), l < r, t_l, r \cap (\cup_{j=1}^{i-1} \tau_j) = \emptyset)} sa(t_{l,r}, p)$$

345 The algorithm proceeds as follows. First, upon string-substring Monge matrix  
 346  $M$  of semi-local solution is built data structure for performing range queries on it  
 347 denoted by *rmq2D* (Lines 1).

348 Second, algorithm make recursive call to subroutine *greedy*. The *greedy* routine  
 349 perfoms greedy choice of  $\tau_i$  with maximal alignment within the current uncovered  
 350 part of the text  $t_{i,j}$ . More precisely, it refers to searching maximum value with  
 351 corresponding position (row and column) in matrix  $M$  within  $t_{i,j}$  (starting at  $i$ th  
 352 position and ending at  $j$ th position of text  $t$ . It is solved via range queries. When  
 353 detected *interval* has alignment score less then threshold it means that no clones  
 354 of pattern  $p$  are presented in this part of text  $t_{i,j}$ , and further processing should be  
 355 skipped. Otherwise, the founded clone is added to final result and the current part  
 356 of the text splits on two smaller parts and processed in the same way. Finally, the  
 357 algorithm outputs a set of the non-intersected intervals of clones of pattern  $p$  in text  
 358  $t$ .

**Algorithm 5.1** Greedy subroutine

Input:  $rmq2D$ — range maximum query data structure for performing range queries on monge matrix  $M$ ,  $h$  — threshold value,  $i, j$  — start and end positions of current text  $t_{i,j}$

Output: Set of non-intersected intervals from  $t_{i,j}$

Pseudocode:

$greedy(rmq2D, h, i, j, t_{i,j}) :$

```

1:  $interval = rmq2D.query(i, j, i, j)$ 
2:  $result = \emptyset$ 
3: if  $interval.score < h$  then
4:   return  $result$ 
5: end if
6: if  $interval.i - i \geq 1$  then
7:    $cl = greedy(rmq2D, h, t_{i,interval.i})$ 
8:    $result.add(cl)$ 
9: end if
10: if  $j - interval.j \geq 1$  then
11:    $cl = greedy(rmq2D, h, t_{j,interval.j})$ 
12:    $result.add(cl)$ 
13: end if
14: return  $result$ 
```

**Algorithm 5.2** GREEDY-PATTERN BASED NEAR DUPLICATE SEARCH ALGORITHM

Input: monge matrix  $M$  that correspond to string-substring matrix for pattern  $p$  and text  $t$ , threshold value  $h$

Output: Set of non-intersected clones of pattern  $p$  in text  $t$

Pseudocode:

$GreedyMatching(M, h, t)$

```

1:  $rmq2D = buildRMQStructure(M)$ 
2:  $result = greedy(rmq2D, 0, |t|, t)$ 
3: return  $result$ 
```

The second algorithm 5.3 uses a less sophisticated approach and a more lightweight one but found fewer duplicates of pattern  $p$  (see example ??). The algorithm also follows a greedy approach but instead of looking at the uncovered part of text  $t$  at each step it looks at the text  $t$  and chooses the first available substring with the highest score that doesn't intersect with already taken substrings. More formally, it approximates algorithm 5.2.

*Algorithm description.* First, the *semi-local sa* problem is solved (Line 1). Then we solve *complete approximate matching problem* (Line 3) i.e for each prefix of text  $t$  we find the shortest suffix that has the highest similarity score with pattern  $p$  (Line 3):

$$(5.2) \quad a[j] = \max_{i \in 0..j} sa(p, t[i, j]), j \in 0..|t|$$

Further, we remove suffixes whose similarity is below the given threshold  $h$  (Line 4). Then remaining suffixes are sorted in descending order (Line 5) and the interval

tree is built upon them (Lines 7-11). The building process comprises from checking that current substring *candidate* not intersected with already added substrings to *tree* and adding it to *tree*. Finally, algorithm output set of non-intersected substrings (clones) of pattern *p* in text *t*.

---

**Algorithm 5.3** Greedy approximate

---

Input: pattern *p*, text *t*, threshold value *h*

Output: Set of non-intersected clones of pattern *p* in text *t*

Pseudocode:

```

1: sa = semilocalsa(p, t)
2: matrix = sa.getStringSubstringMatrix()
3: colmax = smawk(matrix)
4: colmax = colmax.filter(it.score >= h)
5: colmax = colmax.sortedByDescending(it.score)
6: tree = buildIntervalTree()
7: for candidate ∈ colmax do
8:   if candidate ∩ tree = ∅ then
9:     tree.add(candidate)
10:  end if
11: end for
12: result = tree.toList()
13: return result

```

---

THEOREM 5.1. Algorithm 5.3 could be solved in  $\max(O(|p|*|t|*|v|), O(|t|*\log^2|t|v))$  time with  $O(|t|*v*\log|t|*v)$  space when  $|p| < |t|$  where *p* is pattern, *t* is text and *v* is denominator of normalized mismatch score for semi-local sequence alignment  $w_{\text{normalized}} = (1, \frac{\mu}{v}, 0)$  assuming we are storing solution matrix implicitly.

First phase. As shown in section 2 the time complexity of solving semi-local is  $O(|p|*|t|*|v|)$ . The space complexity of storing monge matrix of semi-local solution is  $O(|t|*v*\log|t|*v)$  at most due to fact that *v* – subbistochasticmatrix has at most *v* non-zeros in each row and upon these  $v*|t|$  points we build two dimensional range tree data structure with  $|t|*v*\log|t|*v$  nodes that have report range sum queries in  $O(\log^2|t|v)$  time.

Second phase. SMAWK algorithm requires  $O(|t|*q)$  time where *q* stands for time complexity of random access of monge matrix. Thus, the total time complexity of line 3 is  $O(|t|*\log^2|t|v)$ . Filtering and sorting have at most  $O(|t|)$  and  $O(|t|*\log|t|)$  time complexity. In Line 6 simple initialization of interval tree is performed that requires  $O(1)$ .

Third phase *colmax* array has as worst case  $O(|t|)$  elements when filtering does not eliminate any substrings. Thus, adding to interval tree (both operation at most require  $O(\log|t|)$  time) as well as intersection in (Lines 8-9) will be performed at most  $O(|t|)$ . Thus, the total complexity of last phase is  $O(|t|*\log t)$ .

As we see, the third phase is dominated by the second phase in terms of running time and second phase is dominated by the space complexity of third phase. Thereby, the total time and space complexity is  $\max(O(|p|*|t|*|v|), O(|t|*\log^2|t|v))$  and  $O(|t|*v*\log|t|*v)$  respectively.

COROLLARY 5.2. Algorithm 5.3 could be solved in  $\max(O(|p|*|t|), O(|t|*\log|t|))$  when  $v = O(1)$ .

When  $v = O(1)$  we will use simple range tree for orthogonal range queries with

$O(\log|t|)$  query time.

COROLLARY 5.3. Algorithm 5.3 could be solved in  $O(|p| * |t|)$ .

When amount of clones is relatively small and threshold value is set high then after filtering out  $t$  intervals (Line 4) sorting is performed on  $s$  small set of elements. Thus, this part is dominated by calculating semi-local sa solution.

THEOREM 5.4. Algorithm 5.2 could be solved in  $\max(O(|p| * |t| * v), O(|t| * \log |t|))$  time with  $O(|t| \log |t|)$  space when  $|p| < |t|$  where  $p$  is pattern,  $t$  is text and  $v$  is denominator of normalized mismatch score for semi-local sequence alignment  $w_{normalized} = (1, \frac{p}{v}, 0)$ .

On the first phase of alg

The first phase of algorithm requires  $O(|p| * |t| * v)$  with  $O(|t| * v)$  additional space for storing monge matrix implicitly. We denote this matrix, specifically it's lower-left quadrant that refers to string-substring solution as  $M$  with size  $|t| \times |t|$ .

Theorema 3.4 First, note that

Building structure for rmq queries for staircase matrix requires Theorem 5.8. Given an  $n \times n$  partial Monge matrix  $M$ , a data structure of size  $O(n)$  can be constructed in  $O(n \log n)$  time to answer submatrix maximum queries in  $O(\log \log n)$  time.

*Proof it*

$$D = \text{diag}(d_1, \dots, d_n)$$

COROLLARY 5.5. Algorithm 5.2 could be solved in  $\max(O(|p| * |t|), O(|t| * \log |t|))$  when  $v = O(1)$ .

## 6. Evaluation.

*Research questions.* To present evaluation of algorithms we need to investigate the performance of algorithms that computes solution for *semi-local* problem first. It is justified by fact that all described algorithms heavily based on it. Thus, the following research questions have been settled by evaluation in this paper:

1. RQ1. Does both theoretical algorithms for solving *semi-local* problem applicable in practice? (perform well on practice)
2. RQ2. How differ in terms of running time computation of *semi-local lcs* and *prefix lcs*?

We had implemented algorithms and required data structures to answer RQ1 and RQ2<sup>4</sup>. Evaluation have been done in laptop machine with operation system *Ubuntu18.04* that have processor *Intel-Core i5* with *16GB* RAM.

*RQ1.* On fig. ?? the comparison between two algorithm for computing *semi-local lcs* is presented. The plot marked as *recursive* refers to the algorithm based on steady ant multiplication of associated sticky braids. The second one *reducing* refers to algorithm that based on reducing the associated unreduced sticky braid to reducing one.

Although both algorithms have the same theoretical running time, the figure completely shows that there are significant differences in practice. The complex recursive structure of the algorithm by fast multiplication of sticky braids makes it inapplicable in practice for long input. Nonetheless, such complex structure with combination of steady ant multiplication indeed allows to get rid of one  $v$  when computing *semi-local sa* (see fix ??). The recursive structure of multiplication itself is also a subject of required optimizations due to fact that it used in several theoretical algorithms.

<sup>4</sup>add link to github

For example, in solution for *Window substring* problem or *Bounded Length Smith-Waterman alignment* (implicitly).

*RQ2.* On fig ?? the comparison between computing *prefix lcs* and *semi-local lcs* is presented. More precisely, the comparison among computing prefix lcs via dynamic programming with explicit (denoted by *prefix-lcs*) and implicit (denoted by *light-prefix*) construction of 2D matrix and *semi-local lcs* via reducing approach is presented. The fig ?? show that computation of *semi-local lcs* not only applicable to large input but also comparable with computing of simple *prefix lcs*. The difference between speed computation is relatively subtle.

**7. Conclusion.** Say may be successfully be applied on practice (showed by algorithm luciv updated)

Open problem.— >

Say that need to implement with monge2020 (what we not finished)

Improve algo based on recursive steady ant. Because it's critical for algos based on it.

df[1]

**Acknowledgments.** We would like to acknowledge the assistance of volunteers in putting together this example manuscript and supplement.

## REFERENCES

- [1] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, The Johns Hopkins University Press, Baltimore, 4th ed., 2013.
- [2] D. V. LUCIV, D. V. KOZNOV, A. SHELIKHOVSKII, K. Y. ROMANOVSKY, G. A. CHERNISHEV, A. N. TEREKHOV, D. A. GRIGORIEV, A. N. SMIRNOVA, D. BOROVKOV, AND A. VASENINA, *Interactive near duplicate search in software documentation*, Programming and Computer Software, 45 (2019), pp. 346–355.