

Interactive Near Duplicate Search in Software Documentation

D. V. Luciv^{a,*}, D. V. Koznov^{a,**}, A. A. Shelikhovskii^{a,***}, K. Yu. Romanovsky^{a,****},
G. A. Chernishev^{a,*****}, A. N. Terekhov^{a,*****}, D. A. Grigoriev^{a,*****},
A. N. Smirnova^{a,*****}, D. V. Borovkov^{a,*****}, and A. I. Vasenina^{a,*****}

^aSt. Petersburg University, St Petersburg, 199034 Russia

* e-mail: d.lutsiv@spbu.ru

** e-mail: d.koznov@spbu.ru

*** e-mail: tshel231@gmail.com

**** e-mail: k.romanovsky@spbu.ru

***** e-mail: g.chernyshev@spbu.ru

***** e-mail: a.terekhov@spbu.ru

***** e-mail: d.a.grigoriev@spbu.ru

***** e-mail: anna.en.smirnova@gmail.com

***** e-mail: danila_yumsh@mail.ru

***** e-mail: saibrog@yandex.ru

Received June 28, 2018; revised September 5, 2018; accepted September 5, 2018

Abstract—Various software features such as classes, methods, requirements, and tests often have similar functionality. This can lead to emergence of duplicates in their descriptive documentation. Uncontrolled duplicates created via copy/paste hinder the process of documentation maintenance. Therefore, the task of duplicate detection in software documentation is of importance. Solving it makes planned reuse possible, as well as creating and using templates for unification and automatic generation of documentation. In this paper, we present an approach for interactive detection of near duplicates that involves the user in order to conduct meaningful search. It includes a new formal definition of a near duplicate, a pattern-based, and the proof of its completeness. Moreover, we demonstrate the results of experimenting on a collection of documents of several industrial projects.

DOI: 10.1134/S0361768819060045

1. INTRODUCTION

Software documentation quality issues have been studied since the 1970's [1], and continue to be addressed nowadays [2]. Besides, the documentation, similarly to software, becomes increasingly complex, requires more and more resources for its development and maintenance.

Copy/paste is commonly used in creating and modifying documents: a fragment of text is copied multiple times, and then edited and expanded to suit the subject. The use of this technique is justified by the fact that many software features described in documentation reuse functionality — this is true for requirements, user interface elements, tests, source code, etc. However, without the aid of specialized tools, repeatedly copied fragments create additional difficulties during maintenance because they require extensive synchronisation of changes in corresponding software features. Software reuse is a considerably more advanced research area than software documentation reuse. There is a multitude of studies on soft-

ware reuse and a part of them has been adopted by the industry (see surveys [3–5]). However, the problem of software documentation reuse largely remains a subject of academic research only [6–10]. We should also mention the problem of documentation unification — if there is a large volume of similar information, it is only reasonable to have it presented in a consistent manner. Thus, duplicate detection and duplicate examination are important for setting up documentation reuse and unification.

Duplicates in software documentation have been extensively studied during the last decade [6, 11–17]. At the same time, there are no specialized tools for duplicate detection. Generally, various text search tools are used for this purpose. However, these tools cannot be employed in detection of near duplicates, i.e. text fragments with a substantial common part and certain variations. For this reason, we have created Duplicate Finder [18, 19]. In [16], we have presented a near duplicate search algorithm which considers near duplicates as a combination of exact duplicates.

However, the output quality of this algorithm was low because it does not assess the meaningfulness of found duplicates.

In this paper we present an approach for interactive detection of near duplicates. We involve the user in order to provide meaningfulness of the search process. In short, this process is organized as follows. At first, we automatically create a map of exact duplicates of the document using Clone Miner [20].

The next step relies on the following assumption: an accumulation of exact duplicates in a certain place of the document points to a possible emergence of a near duplicate. At this step, the user moves to the most duplicate-populated document section using this map as a clue. After that, they select the text fragment that contains the most frequently used exact duplicates. Then, the user transforms this fragment into a full description of a certain software feature by extending its bounds. This way, the user ensures the meaningfulness of the fragment. Next, this description (textual string) is used as a pattern for further search. The user can also edit the search results by filtering out false positives and ensuring the meaningfulness of found fragments by expanding or narrowing their bounds in the document.

This article is organized as follows. In Section 2 we present an overview of existing studies that are similar to the current work, and in Section 3 we describe the approaches, methods, and technologies we have used in our study. Section 4 contains the description of the interactive near duplicate search technique. In Section 5 we propose a new formal definition of a near duplicate, and in Section 6 we present the pattern-based near duplicate search algorithm that forms the basis of the proposed technique. In addition, we formulate the criterion of completeness for the algorithm and prove that the proposed algorithm is complete. In Section 7 we provide complexity estimates for the algorithm, and in Section 8 we demonstrate the results of an experimental evaluation.

2. RELATED WORK

Duplicates in software documentation have been extensively studied during the last decade. Horie et al. [6] consider duplicates in API documentation of Java projects, extending JavaDoc with reuse tools. This approach is expanded with consideration of near duplicates by Nosál' and Porubän in [12]. Similarly to approaches presented in references [8–10], the authors of this paper employ parametrization for defining variative parts of duplicates. However, this study does not consider the task of near duplicate search itself, and their definition of a near duplicate is informal.

In their further research [13], they examine exact duplicates in embedded documentation of several

open-source projects, but do not consider near duplicates.

Wingkvist et al. [14] use duplicates to evaluate the quality of documentation, with no consideration of near duplicates.

The paper [11] by Juergens et al. is an examination of duplicates in requirement specifications: the authors have analyzed 28 industrial documents, manually filtering and classifying found duplicates. The meaning of these duplicates was discussed (with emphasis on duplicates that corresponded to code clones). They also have studied the influence of redundancy on the speed of reading and understanding texts. This work does not consider other types of software documentation, as well as near duplicates, although the authors do mention their existence.

Ouzmazis et al. [7] analyze API documentation of several well-known open source projects, classify detected duplicates, and consider the problem of documentation reuse. They do not consider near duplicates, however, they note that those duplicates occur quite often and are important in practice.

Rago et al. [21] present near duplicate search in textual use case descriptions with the use of natural language processing methods. However, they consider a highly specific type of requirement specifications, which is rarely used in practice. Moreover, it is unclear how to apply this method to other types of documentation.

Concluding our overview, we should note that most existing approaches except [21] use token-based tools of code clone analysis. This fact seriously complicates near duplicate detection. However, some authors acknowledge the existence and importance of near duplicates in documentation redundancy analysis and documentation reuse [7, 11, 12].

The approach presented in this work is largely based on pattern matching. This problem is well-studied and has been solved in multitude of ways for different contexts. Let us provide a short overview of this problem in the context of text search.

The algorithm proposed by Ukkonen [22] makes efficient matching approximate occurrences of pattern in text possible, but it requires a costly preprocessing of the pattern.

Broder [23] describes a method for matching approximate occurrences of pattern in text using information retrieval methods; we should note that this approach also requires expensive preprocessing of input data (both document and pattern).

Algorithms presented in [22, 24–26] are efficient but quite sophisticated, which complicates their use and modification, as well as proving their formal properties.

Ukkonen's algorithm [22] is suited for working with an immutable pattern, and Broder's approach [23] operates on an immutable document. Both of

these situations are irrelevant to our task. Studies by Landau and Vishkin [25], Myers [26] describe algorithms for detecting text fragments for which Levenshtein distance [27] does not exceed a pre-defined threshold. We should emphasize the high computational complexity of Levenshtein distance calculation, which, as shown by our experiments, makes this approach unsuitable for duplicate detection. A more detailed review of approximate pattern matching in text can be found in [28]. Using ideas proposed in this research area, we have created our own pattern matching algorithm while adhering to the following requirements:

- (i) the algorithm needs to perform near duplicate detection in accordance with our definition of near duplicates;
- (ii) we wanted to formally prove a number of properties of this algorithm;
- (iii) the run time has to be adequate because the algorithm is run in an interactive mode;
- (iv) the algorithm needs to yield as few false positives as possible.

3. BACKGROUND

3.1. Edit Distance

We use *edit distance* [27] to determine the degree of similarity of two text fragments (strings of text). This distance is essentially the number of string editing operations required to convert one string into another: the less operations, the more similar the strings. Different definitions of edit distance differ by their admissible operations. In our work, we use *longest common subsequence distance* [29, 30] that uses only insertion and deletion of a symbol due to its suitability for near duplicate model described below, and, consequently, the convenience of further proofs. The authors of [31] prove that this type of edit distance has metric properties. Further on, we will denote the longest common subsequence distance between two strings s_1 and s_2 as $d(s_1, s_2)$.

Computing edit distance is a resource-consuming task. The complexity of the algorithm we have selected is estimated [32] as $\mathcal{O}(|s_1| * |s_2|)$ in the average case. Furthermore, the authors of reference [33] show that it is impossible to design an algorithm that can provide better complexity estimates for the worst case. In this work, we use the difflib library [34], which is included in the standard Python package (performance-critical parts of which are implemented in C).

3.2. Detecting Exact Duplicates with Clone Miner

We have employed exact duplicate detection to create a duplicate map, using which the user could select a pattern for matching. We have selected Clone Miner [20] for this task, which is a software code clone detec-

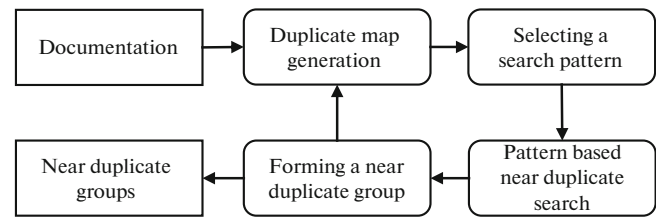


Fig. 1. Process overview.

tion tool. Clone Miner is a token-based tool, and it does not employ an abstract syntax tree. A token is a stand-alone word (sequence of symbols) in a document, separated from adjacent words by such delimiters as “.”, “;”, “:”, “(”, “)”, etc. For example, the fragment “FM registers” consists of two tokens. Clone Miner considers text as an ordered collection of tokens and detects duplicate fragments (clones) using algorithms based on suffix trees [35]. We have chosen Clone Miner because it is easily integrated with other tools through command line interface and supports the Russian language.

4. THE PROCESS

The general purpose of our process is to ensure meaningfulness of duplicate detection via user interaction. A diagram which describes the workflow of the process is presented in Fig. 1. Let us describe the process in detail.

Generating a duplicate map. Using Clone Miner [20], all exact duplicate groups in the document are detected. Every token (word) t is assigned a color from an RGB interval from white to red: $color(t) = (h(t)/T_m) * R + (1 - (h(t))/T_m) * W$, where $R = [1, 0, 0]$, $W = [1, 1, 1]$, $h(t)$ is the exact duplicate group that has the maximum cardinality and contains t (further on called token temperature), and T_m is the maximum cardinality of an exact duplicate group in the document (maximum token temperature)¹. The closer a token’s color to red, the “warmer” it is. This metaphorical representation is called a heat map [36], an example of which can be seen in Fig. 2

The generated heat map provides an overview of the potential near duplicate occurrences. The areas most likely to contain near duplicates are represented by red areas of different hue. Tokens that occur in the reddest areas are repeated the same or roughly the same number of times. Therefore, the probability that these tokens would form a meaningful near duplicate is quite high. This way, one may hope to obtain mean-

¹ We only consider groups that consist of fragments longer than four tokens, because, according to our experiments [15], this particular constraint filters out many false positives.



Fig. 2. Duplicate map.

ingful near duplicates that appear a significant number of times in the document.

Selecting a search pattern. The user moves to the reddest (the “warmest”) area on the heat map (Fig. 2), zooms in on it, and selects a fragment (pattern) for further search (see Fig. 3). During this process, the user does not only consider the color of the selected fragment, but aims to select a fragment that describes a software feature in full. To achieve this, the user can either include a white-colored text fragment in the pattern, or not include a red-colored one. Consider an example. Following the information from Fig. 3, we select this fragment:

“To alter the owner, you must also be a direct or indirect member of the new owning role, and that role must have CREATE privilege on the table’s schema. (These restrictions enforce that altering the owner doesn’t do anything you couldn’t do by dropping and recreating table. However, a superuser can alter ownership of any table anyway.)”

This fragment describes an integral software feature concerning the administration of the PostgreSQL DBMS: to alter the owner of a database table, you must also have specific rights or be an administrator.

Near duplicate search. The user selects a similarity measure for the highlighted fragment, which is a number from $1/\sqrt{3}$ to 1, and launches the pattern matching algorithm².

Forming a near duplicate group. Having received the algorithm’s output, the user modifies it. During the process the user deletes elements (near duplicate

occurrences) that only resemble the pattern syntax-wise but not meaning-wise. Furthermore, for each occurrence the user can modify the bounds of the fragments to ensure the meaningfulness of each.

5. DEFINING A NEAR DUPLICATE GROUP

In this section we generalize the definition of a near duplicate group that we have proposed earlier in [16, 37]. In contrast with the previous definition, here we use a parameter instead of a constant to define the similarity measure, and we allow to place extension points at the ends of duplicates. We will consider a document D as a finite sequence of symbols, denoting its length as $\text{length}(D)$.

Definition 1. A **text fragment** is an occurrence of a certain symbol string in document D .

Therefore, for every text fragment g of document D there is an integer interval $[b, e]$, where b is the position of the first symbol of the fragment, and e is the position of the last. By $g \in D$, we denote a text fragment g of document D . Next, let $[g]$ be the function that maps a text fragment g to its interval, and let $\text{str}(g)$ be the function that maps a text fragment g to its textual content. By $b(g)$ and $e(g)$ we denote the positions of the beginning and the end of g . Next, $|g|$ is a function that takes a text fragment g and returns its length as $|g| = 1 + e(g) - b(g)$. Finally, we introduce a two-place predicate $\text{Before}(g_1, g_2)$, which is true if and only if $e(g_1) < b(g_2)$.

Definition 2. **Near duplicate group.** Consider a collection of text fragments g_1, \dots, g_M of document D . We will call this collection a near duplicate group with the similarity measure $k \in (1/\sqrt{3}, 1]$ (or simply a near duplicate group) if the following conditions are satisfied.

- (1) $\forall i \in \{1, \dots, M - 1\}$ holds $\text{Before}(g_i, g_{i+1})$;
- (2) There exists a ordered collection of strings (I_1, \dots, I_N) such as there is an occurrence of this collection in every text fragment, i.e. $\forall j \in \{1, \dots, M\} \forall i \in \{1, \dots, N\} I_i \subset \text{str}(g_j)$ and $\forall i \in \{1, \dots, N - 1\}$ holds $\text{Before}(I_i^j, I_{i+1}^j)$, where I_i^j is an occurrence of I_i in g_j , and the following condition is satisfied:

$$\forall j \in \{1, \dots, M\} \frac{\sum_{i=1}^N |I_i|}{|g_j|} \geq k.$$

We define the *archetype* of a given group as a collection of strings (I_1, \dots, I_N) . It is easy to show that the definition proposed above generalizes the definition

² The $1/\sqrt{3} \approx 0.577$ value was selected for the convenience of the following proofs; following our experiments, we have concluded that if the similarity measure is less than $1/2$, then, for smaller patterns (up to 15–20 tokens), the algorithm produces many non-meaningful matches; the lower bound we have selected is insignificantly larger than $1/2$.

```

large tables, since only one pass over the table need be made.
You must own the table to use ALTER TABLE. To change the schema or
tablespace of a table, you must also have CREATE privilege on the new
schema or tablespace. To add the table as a new child of a parent
table, you must own the parent table as well. To alter the owner, you
must also be a direct or indirect member of the new owning role, and
that role must have CREATE privilege on the table's schema. (These
restrictions enforce that altering the owner doesn't do anything you
couldn't do by dropping and recreating the table. However, a superuser
can alter ownership of any table anyway.) To add a column or alter a
column type or use the OF clause, you must also have USAGE privilege on
the data type.
PARAMETERS

```

Fig. 3. Selecting a search pattern (PostgreSQL documentation).

given in [16, 37]. If G is a near duplicate group, then by $|G|$ denote the number of elements of this group.

Definition 3. Consider a text fragment p of document D ($p \in D$) and $g \in D$. We say that g is a near duplicate of p with similarity k , if g and p form a near duplicate group with similarity k defined according to 2.

6. PATTERN BASED NEAR DUPLICATE SEARCH ALGORITHM

6.1. Algorithm Description

The algorithm consists of three phases. At **phase 1 (scanning)**, document D is scanned by a sliding window w of size $L_w = \frac{|p|}{k}$ with a one symbol step³. The text fragment that corresponds to the current window position is compared to pattern p using edit distance, and if they are close, i.e. $d(p, w) \leq k_{di}$, then this fragment is saved in the set W_1 . The threshold value k_{di} is defined as follows:

$$k_{di} = |p| \left(\frac{1}{k} + 1 \right) (1 - k^2). \quad (1)$$

This choice will be explained below.

At **phase 2 ("shrinking")**, we search for the largest text fragment that is closest to pattern p in every element of W_1 . Essentially, during this phase lengths of elements of W_1 decrease, i.e. text fragments are "shrunk." This is reasonable since the window (and consequently, all elements of W_1) is of maximum possible size of a near duplicate of p (see Lemma. 1). During "shrinking" for every $w_2 \in W_1$, all of its internal fragments are iterated over, starting with fragments of $|p| * k$ length up to fragments of $\frac{|p|}{k}$ length. The one that is closest to the pattern in terms of edit distance is selected. If there are

³ Here and further we do not round the lengths of the intervals to integers to save up space. Nevertheless, all proofs can be performed with rounded values as well.

several such fragments, the longest one should be taken. This phase results in the set W_2 .

Algorithm 1: Pattern based near duplicate search algorithm

```

Input data:  $D$  – document,
 $p$  – pattern,  $k$  – similarity measure
Result:  $R$ 
// Phase 1 (scanning)
1  $W_1 \leftarrow \emptyset$ 
2 for  $\forall w_1 : w_1 \in D \wedge |w_1| = L_w$  do
3   if  $d_{di}(w_1, p) \leq k_{di}$  then
4      $\hookrightarrow$  add  $w_1$  to  $W_1$ 
// Phase 2 ("shrinking")
5  $W_2 \leftarrow \emptyset$ 
6 for  $w \in W_1$  do
7    $w'_2 \leftarrow w$ 
8   for  $l \in I$  do
9     for  $\forall w_2 : w_2 \subseteq w \wedge |w_2| = l$  do
10      if  $\text{Compare}(w_2, w'_2, p)$  then
11         $w'_2 \leftarrow w_2$ 
12       $\hookrightarrow$  add  $w'_2$  to  $W_2$ 
// Phase 3 (filtering)
13  $W_3 \leftarrow \text{Unique}(W_2)$ 
14 for  $w_3 \in W_3$  do
15   if  $\exists w'_3 \in W_3 : w_3 \subset w'_3$  then
16      $\hookrightarrow$  remove  $w_3$  from  $W_3$ 
17  $R \leftarrow W_3$ 

```

At the **phase 3 (filtering)**, duplicate elements in W_2 are eliminated. They emerge at the previous phase because W_1 can contain text fragments that differ by a window shift of several symbols. Furthermore, ele-

ments that are fully contained in other elements of W_2 are filtered out. This phase results in the set W_3 which is the output of the algorithm, i.e. the set R .

Let us describe the auxiliary functions used in Algorithm 1. The Compare function is used during phase 2 to identify the text fragment which is closer to the pattern p in terms of edit distance. If the distance from both fragments to the pattern is the same, the longest fragment is selected:

$$\text{Compare}(w_1, w_2, p) = \begin{cases} \text{true} & d(w_1, p) < d(w_2, p) \\ \text{false} & d(w_1, p) > d(w_2, p) \\ |w_1| > |w_2| & d(w_1, p) = d(w_2, p) \end{cases}$$

The Unique function receives a collection of text fragments, iterates over it and discards duplicate fragments.

6.2. Algorithm Completeness

The criterion of completeness for our pattern based near duplicate search algorithm is defined as follows. The algorithm is complete if for arbitrary D , $p \in D$, output of the algorithm R , and for any near duplicate group G of fragment p with similarity k (see def. 2), the following condition holds true:

$$\forall g \in G \quad \exists w \in R: |g \cap w| \geq O_{\min}(k), \quad (2)$$

where $O_{\min}(k) = \frac{|p|}{2} \left(3k - \frac{1}{k} \right)$. This criterion can be explained as follows: for any fragment of document D that is a near duplicate of pattern p , the set R will contain a text fragment that significantly intersects with this near duplicate, allowing the user to easily recognise this duplicate in the output. The ratio of the intersecting portion to the whole pattern is bounded from below by the $O_{\min}(k)$ function. $O_{\min}\left(\frac{1}{\sqrt{3}}\right) = 0$, and for larger values of k $O_{\min}(k) > 0$. This is true since the function increases with increasing k – its derivative is $\frac{|p|}{2} \left(3 + \frac{1}{k^2} \right)$ and it is obviously positive for all $\frac{1}{\sqrt{3}} < k \leq 1$. In practice, the best results are achieved for $k \geq 0.77$: for these values $O_{\min}(k) > \frac{|p|}{2}$, i.e. all elements of R intersect all near duplicates at least by half of the pattern's length. Note that the lower estimate $O_{\min}(k)$ is pessimistic: the experimental results demonstrate a larger overlap of the output and the near duplicates contained in the document. Let us continue on to the completeness of the proposed algorithm, proving several auxiliary propositions first.

Lemma 1. *Let G be a near duplicate group of fragment p with similarity k . Then $\forall g_1, g_2 \in G \quad k \leq \frac{|g_1|}{|g_2|} \leq \frac{1}{k}$ holds true.*

Proof. Suppose A is the archetype of group G . Then $k|g_1| \leq |A|$ and $k|g_2| \leq |A|$. Because $A \subset \text{str}(g_1)$ and $A \subset \text{str}(g_2)$, we have: $k|g_1| \leq |A| \leq |g_1|$ and $k|g_2| \leq |A| \leq |g_2|$. Therefore, $k|g_1| \leq |g_2|$ and $k|g_2| \leq |g_1|$. Dividing these inequalities by $k|g_1| \leq |g_1|$ and $k|g_2| \leq |g_2|$ respectively, we get the required result.

Lemma 2. *Let G be a near duplicate group of fragment p with similarity k . Then $\forall g \in G$ the following holds true: $d(g, p) \leq (1 - k^2)|p|$.*

Proof. Because p and g belong to the same near duplicate group, they have the same archetype and can be presented in the following way:

$$p = v_0^p I_1 v_1^p I_2 \dots v_{N-1}^p I_N v_N^p, \\ g = v_0^g I_1 v_1^g I_2 \dots v_{N-1}^g I_N v_N^g,$$

where I_1, I_2, \dots, I_N is the archetype of group G , $v_0^p, v_1^p, \dots, v_N^p$ is the variative part of p , and $v_0^g, v_1^g, \dots, v_N^g$ is the variative part of g .

Let us introduce the following notations: $v^p = v_0^p v_1^p \dots v_N^p$, $v^g = v_0^g v_1^g \dots v_N^g$, $A = I_1 I_2 \dots I_N$. Then according to (2), we have $|A|/|p| \geq k \Rightarrow |p| - |v^p| \geq |p|k \Rightarrow |p| - |p|k \geq |v^p| \Rightarrow |p|(1 - k) \geq |v^p|$, and, likewise, $|g|(1 - k) \geq |v^g|$. Moreover, g can be obtained from p by substituting v_i^g for v_i^p , i.e. $d(g, p) \leq |v^g| + |v^p| \leq (1 - k)(|p| + |g|)$. According to lemma 1 we have $|g| \leq k|p|$. Then $d(g, p) \leq (1 - k)(1 + k)|p| = (1 - k^2)|p|$.

Lemma 3. *For any $p \in D$, $k \in (1/\sqrt{3}, 1]$, near duplicate group G of fragment p with similarity k (Definition. 3), the criterion of completeness 2 is satisfied in respect to the results of phase 1.*

Proof. As mentioned above, the triangle inequality is satisfied for longest common subsequence distance: $d(fr, p) \leq d(fr, g) + d(g, p)$. According to lemma 2, $d(g, p) \leq |p|(1 - k^2)$. We also know that $g \subseteq fr$. Therefore, because we can obtain g from fr by removing all symbols that belong to $fr \setminus g$, $d(fr, g) \leq |fr| - |g|$ holds true. But because $|fr| = \frac{|p|}{k}$ and according to lemma 1, $|g| \geq k|p|$, the following also holds true: $|fr| - |g| \leq \frac{1}{k}|p| - k|p|$. Therefore, $d(fr, p) \leq |p| \left(\frac{1}{k} - k + 1 - k^2 \right) = |p| \left(1 + \frac{1}{k} \right) (1 - k^2)$. It is obvious that during the scanning on phase 1 there will be a state in which the window contains fr .

Then, according to (1), $fr \in W_1$. Therefore, the following holds true:

$$\forall g \in G : \left(|fr| = \frac{|p|}{k}, g \subseteq fr \right) \Rightarrow fr \in W_1.$$

Since for any near duplicate there is an element of W_1 that does not only intersect with this duplicate, but contains it completely, criterion 2 is satisfied for W_1 if we consider it as the set R .

Lemma 4. *For any $p \in D$, $k \in (1/\sqrt{3}, 1]$ and near duplicate group G of fragment p with similarity k (Definition. 3) the criterion of completeness 2 is satisfied in respect to the output of phase 2.*

Proof. We omit a formal proof due to its large size. The main idea here is considering the worst case where during “shrinking,” the length of $w_1 \in W_1$ elements is decreased to $k|p|$. Considering corner cases (an element positioned in the center or right at the ends of w_1) allows us to confirm that the lemma holds true.

Lemma 5. *For any $p \in D$, $k \in (1/\sqrt{3}, 1]$ and near duplicate group G of fragment p with similarity k (Definition. 3) the criterion of completeness 2 is satisfied in respect to the output of phase 3.*

Proof. Phase 3 consists of element deletion from W_2 only. The intervals of the deleted elements are contained in intervals of other elements. It is obvious that if W_2 satisfied the criterion, then W_3 will satisfy it as well.

Theorem 1. *The criterion of completeness is satisfied for any $p \in D$, $k \in (1/\sqrt{3}, 1]$, corresponding algorithm output R and any near duplicate group G of fragment p with similarity k .*

Proof. The output of phases 1–3 was proven to satisfy the criterion 2 in lemmas 3–5.

6.3. Optimizing the Algorithm

The proposed algorithm turned out to be inadequate performance-wise: its run time exceeded one hour when searching for patterns larger than 100 symbols in documents of about 2 MB in size. Furthermore, the algorithm produced many false positives – its output contained the same text fragments that were insignificantly shifted relatively to each other. As the result, a range of optimizations has been suggested.

Optimization 1 is applied during phase 1 (scanning). It allows to reduce the number of calculations of d , significantly improving the run time of the algorithm. It is based on the known Boyer-Moore algorithm, which is intended for matching a pattern in a string [38]: during the scan, a check is performed to see by how many symbols the window can be shifted without skipping the required result. Therefore, at each step of the scan we check whether $d(w, p) > k_{di} + 1$ (w is the win-

dow position) holds. If it is true, then we slide the window by $(d(w, p) - k_{di})/2$ symbols to the right. Otherwise, we slide it by one symbol.

Optimization 2 is applied during phase 2 (“shrinking”). It allows to reduce the number of computations of d as well. The approach is similar to the one used in the previous optimization. During “shrinking” of a text fragment w_1 , the window scans the fragment in a symbol-by-symbol manner. At each step $d(p, w'_2)$ is computed and, if necessary, its minimum value d_{min} is updated. If for the current window position w'_2 , $d(p, w'_2) > d_{min} + 1$ holds true, slide the window to the right by $(d(p, w'_2) - d_{min})/2$ symbols. Otherwise, slide it by one symbol. The d_{min} value is updated at the beginning of each iteration corresponding to the next value of the sliding window width.

Optimization 3 is applied during phase 3 (filtering). It allows to minimize the cardinality of W_3 . It is as follows: the set is divided into maximum subsets that are transitively closed under intersection. Further, for every such subset a w_3 fragment with the minimum value of $d(w_3, p)$ is selected, or if there are several such fragments, the one with maximum length. All remaining elements of the set are deleted.

Optimization 4 is applied during phase 3, extending all text fragments of W_3 up to complete words. The bounds of a text fragment can ignore the bounds of words, i.e. incomplete words can be included into text fragments. In order to address this, text fragments are expanded to include these words fully. This helps to decrease the number of false positives in the algorithm’s output.

Optimization 5 is applied during phases 1 and 2. Its purpose is reusing d for the same strings and parallelizing the “shrinking” of the elements of W_1 .

Let us show how these optimizations affect the algorithm’s completeness.

Theorem 2. *Optimizations 1, 2, 4, 5 preserve the completeness property.*

Proof. Consider two strings that are results of concatenation: $s_1 = ab$ and $s_2 = bc$, where $|a| = |c|$ and $d(s_1, s_2) = d$. We can easily show that $|a| + |c| \geq d$. Using this fact, it is easy to prove the completeness of optimization 1. The completeness of optimization 2 is proven in the same way. The completeness of optimization 4 can not be doubted because it only extends the elements of the output. Finally, optimization 5 is complete because it only considers the implementation of the algorithm.

Note 1. Situations where optimization 3 does not satisfy the criterion of completeness are possible, but our experiments show that their number is insignificant in practice.

7. ALGORITHM COMPLEXITY

Document length $|D|$, pattern length $|p|$, the k value, and the cardinality of the near duplicate group of the pattern $|G_p|$ are all significant parameters that influence the run time of the algorithm. Let us estimate the algorithm's complexity depending on these parameters.

The average complexity of calculating d (i.e. edit distance) is $\mathcal{O}(|p|^2)$ [32]. Consequently, the average complexity of phase 1 is proportional to $|p|^2$ and $|D|$. During phase 2 all of the internal fragments of each $w_i \in W_1$ are iterated over, and it is easy to show that their number is proportional to $|p|$. Furthermore, the cardinality of the set W_1 is proportional to $|G_p|$ and k_{di} . Finally, the complexity of phase 2 is $\mathcal{O}(|G_p|)$ and $\mathcal{O}(|p|^4)$. The complexity of phase 3 operations is $\mathcal{O}(|W_2| * \log|W_2|)$, but because $|W_2| = |W_1|$, the complexity of phase 3 is $\mathcal{O}(|G_p| * \log|G_p|)$ and $\mathcal{O}(|p| * \log|p|)$. Optimizations 1 and 2 on average lead to “skips” during iteration, the size of which is proportional to k_{di} (and hence $|p|$), making the complexity of phases 1 and 2 $\mathcal{O}(|p|)$ and $\mathcal{O}(|p|^3)$ respectively. Therefore, with $k = \text{const}$ the algorithm's run time can be estimated as $\mathcal{O}(|D|)$, $\mathcal{O}(|p|^3)$, and $\mathcal{O}(|G_p| * \log|G_p|)$.

Theorem 3. *The complexity of the algorithm with fixed D and p is estimated as $\mathcal{O}(1/k^4)$ on average.*

We omit the proof due to its large volume.

8. EVALUATION

Theoretical complexity estimates are not sufficient for determining the real run time of the proposed algorithm. These estimates were produced using certain significant parameters of the algorithm independently to simplify the proofs, while real complexity can depend on their combinations. Another argument for the necessity of experimental evaluation is the fact that theoretical estimates do not provide the real value intervals of these parameters. Finally, other properties of the algorithm need to be evaluated as well.

We have conducted our experiments to answer the following questions:

- (i) what is the run time of the pattern matching algorithm on real data;
- (ii) how large are algorithm's outputs having real data as input.

The first question is important because the algorithm is used in interactive mode, and therefore its run time should not exceed several minutes. Considering output volume, we have proven that our algorithm's output contains all existing near duplicates of a certain

pattern. It is, however, unclear, how exactly large are the real outputs of the algorithm – outputs that contain over 100 elements become more or less unfeasible for human analysis. In turn, output volume is affected by the number of false positive matches and the number of near duplicates in the document

We have experimented on 19 industrial documents both in Russian and English (described in reference [16]). The experiments were conducted on a computer with the following specifications: Intel Core i7 2600, 3.4 GHz, 16 GB RAM. The documents were converted into “flat text” (UTF-8 encoding) with Pandoc [39]. After the conversion, the size of the documents ranged from 0.04 MB to 2.5 MB (0.75 MB on average). We are inclined to think that these numbers are realistic for $|D|$. However, we should note that it is necessary to create a more representative selection of different documentation types in order to obtain more precise estimates.

The experiments were conducted as follows. We have run the algorithm for patterns of length ranging from 50 to 1000 symbols with a 50-symbol step. A 1000-symbol fragment is about 25% of a page of a docx document, i.e. it is a large fragment, and following our experiments, duplicates are significantly smaller in general. We have iterated the similarity measure value k from 0.6 to 1 with 0.1 step for each selected pattern in each document. We have selected the pattern in the following way. Having a fixed pattern length, we followed our technique and selected the “warmest” area in the document of this length, calculating it automatically as a fragment where the following expression reaches its maximum value: $\sum_{t \in fr} h(t)$. In this expression t is a token of fragment fr and $h(t)$ is its temperature. The sum is calculated over all tokens of the fragment, including possibly incomplete leftmost and rightmost tokens.

Analyzing the data obtained from the experiments to answer the question whether the algorithm's run time is suitable for interactivity, we have established the following: in 38% of cases the algorithm ran for less than 5 s, in 78% cases – less than 30 s, in 90% of cases – less than 2 min. These run times are fairly adequate for interactive mode.

We have obtained the following data on the output volumes of our algorithm: 84% of outputs contained less than 100 elements, 5% outputs – from 100 to 200 elements, 5.6% – from 200 to 600 elements, 5.4% – from 600 to 1000 elements. Thus, the majority of near duplicate groups in software documentation are relatively small (containing up to 100 elements), which follows from Theorem 1 and our experimental results.

CONCLUSIONS

In this study we have presented an interactive near duplicate search process for software documentation.

This process solves the problem of meaningful extraction of near duplicates by involving the user, who can use an automatically generated heat map of exact duplicates to detect the most probable occurrences of near duplicates. We have created a pattern-based near duplicate search algorithm and provided optimizations for it. We have proven the completeness of the algorithm, meaning that all near duplicates contained in the document are present in the algorithm's output. More precisely, duplicates located in the document significantly intersect with particular elements of the output, and this is why the user can identify them with ease. Our process allows user to manually edit their bounds and to include them in the output in full. We present complexity estimates for our algorithm as well as experimental results. These results suggest that duplicate groups in software documentation generally do not exceed 100 elements, and the algorithm itself performs adequately for practical use.

In the future, we plan to study different types of software documentation in detail using our algorithm and experiment model (focusing on API documentation first). We also intend to thoroughly examine the behavior of our algorithm with varying input parameters (pattern length and similarity measure). Finally, we plan to switch to automatic methods of detecting meaningful duplicates via machine learning. A detailed analysis of different types of near duplicates in different software documentation types is required as well. Other fruitful areas for future work are integration of documentation reuse (in the context of requirement development) with automatic test development [40, 41], and visualization of duplicate structure using diagrams [42].

FUNDING

This work is partially supported by RFBR grant 16-01-00304.

REFERENCES

- Brooks, F.P., *The Mythical Man-Month: Essays on Software Engineering*, Addison-Wesley, 1975.
- Parnas, D.L., Precise documentation: The key to Better Software, *The Future of Software Engineering*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 125–148.
- Bassett, P.G., *Framing Software Reuse: Lessons from the Real World*, Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1997.
- Jarzabek, S. and Pettersson, U., Research journey towards industrial application of reuse technique, *ICSE*, 2006, pp. 608–611.
- Irshad, M., Petersen, K., and Poulding, S.M., A systematic literature review of software requirements reuse approaches, *Inform. Software Technol.*, 2018, vol. 93, pp. 223–245.
- Horie, M. and Chiba, S., Tool support for crosscutting concerns of API documentation, *Proceedings of the 9th International Conference on Aspect-Oriented Software Development*, New York, NY, USA: ACM, 2010, pp. 97–108.
- Oumaziz, M.A., Charpentier, A., Falleri, J.-R., and Blanc, X., Documentation reuse: Hot or not? An empirical study, *Mastering Scale and Complexity in Software Reuse: 16th International Conference on Software Reuse*, Springer, 2017, pp. 12–27.
- Koznov, D.V. and Romanovsky, K.Yu., DocLine: A method for software product lines documentation development, *Programming and Computer Software*, 2008, vol. 34, no. 4, pp. 216–224.
- Romanovsky, K., Koznov, D., and Minchin, L., Refactoring the documentation of Software Product Lines, *Lect. Notes Computer Science*, Berlin, Heidelberg: Springer-Verlag, 2011, vol. 4980, pp. 158–170.
- Jarzabek, S. and Dan, D., *Documentation Reuse: Managing Similar Documents*, FedCSIS, 2017, pp. 1325–1334.
- Juergens, E., Deissenboeck, F., Feilkas, M., Hummel, B., Schaetz, B., Wagner, S., Domann, C., and Streit, J., Can clone detection support quality assessments of requirements specifications?, *Proceedings of ACM/IEEE 32nd International Conference on Software Engineering*, 2010, vol. 2, pp. 79–88.
- Nosál, M. and Porubán, J., Reusable software documentation with phrase annotations, *Central Eur. J. Comp. Sci.*, 2014, vol. 4, no. 4, pp. 242–258.
- Nosál, M. and Porubán, J., Preliminary report on empirical study of repeated fragments in internal documentation, *Proceedings of Federated Conference on Computer Science and Information Systems*, 2016, pp. 1573–1576.
- Wingkvist, A., Lowe, W., Ericsson, M., and Lincke, R., Analysis and visualization of information quality of technical documentation, *Proceedings of the 4th European Conference on Information Management and Evaluation*, 2010, pp. 388–396.
- Koznov, D., Luciv, D., Basit, H.A., Lieh, O.E., and Smirnov, M., Clone Detection in Reuse of Software Technical Documentation, *International Andrei Ershov Memorial Conference on Perspectives of System Informatics (2015)*, Springer Nature, 2016, vol. 9609 of Lecture Notes in Computer Science, pp. 170–185.
- Luciv, D.V., Koznov, D.V., Chernishev, G.A., Terekhov, A.N., Romanovsky, K.Yu., and Grigoriev, D.A., Detecting near duplicates in software documentation, *Programming and Comput. Software*, 2018, vol. 44, no. 5.
- Koznov, D.V., Luciv, D.V., and Chernishev, G.A., Duplicate management in software documentation maintenance, *Proceedings of V International Conference Actual Problems of System and Software Engineering*, vol. 1989: CEUR Workshop Proceedings, 2017, pp. 195–201.
- Duplicate Finder.
<http://www.math.spbu.ru/user/kromanovsky/docline/index.html>.
- Luciv, D.V., Koznov, D.V., Chernishev, G.A., Basit, H.A., Romanovsky, K.Yu., and Terekhov, A.N., Poster: Duplicate finder toolkit, *Proceedings of the International Conference on Software Engineering (ICSE 2018)*, 2018, pp. 171–172.

20. Basit, H.A., Puglisi, S.J., Smyth, W.F., Turpin, A., and Jarzabek, S., Efficient token based clone detection with flexible tokenization, *Proceedings of the 6th Joint Meeting on European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering: Companion Papers*, New York, NY, USA: 2007, pp. 513–516.
21. Rago, A., Marcos, C., and Diaz-Pace, J.A., Identifying duplicate functionality in textual use cases by aligning semantic actions, *Software and Systems Modeling*, 2016, vol. 15, no. 2, pp. 579–603.
22. Ukkonen, E., Finding approximate patterns in strings, *J. Algorithms*, 1985, vol. 6, no. 1, pp. 132–137.
23. Broder, A.Z., On the resemblance and containment of documents, *Compression and Complexity of Sequences 1997, Proceedings*, IEEE, 1997, pp. 21–29.
24. Wu, S. and Manber, U., Fast Text Searching: Allowing Errors, *Commun. ACM*, 1992, vol. 35, no. 10, pp. 83–91.
25. Landau, G.M. and Vishkin, U., Fast string matching with k differences, *J. Comput. System Sci.* 1988, vol. 37, no. 1, pp. 63–78.
26. Myers, G., A Fast bit-vector algorithm for approximate string matching based on dynamic programming, *J. ACM*, 1999, vol. 46, no. 3, pp. 395–415.
27. Levenshtein, V., Binary codes capable of correcting spurious insertions and deletions of ones, *Problems Inform. Transmission*, 1965, vol. 1, pp. 8–17.
28. Smyth, W., *Computing Patterns in Strings*, Addison-Wesley, 2003.
29. Bergroth, L., Hakonen, H., and Raita, T., A survey of longest common subsequence algorithms, *String Processing and Information Retrieval, 2000 (SPIRE 2000), Proceedings, Seventh International Symposium on, 2000*, pp. 39–48.
30. Leskovec, J., Rajaraman, A., and Ullman, J.D., *Mining of Massive Datasets*, Cambridge: Cambridge Univ. Press, 2014.
31. Gusfield, D., *Algorithms on Strings, Trees, and Sequences*, Cambridge: Cambridge Univ. Press, 1997.
32. Ratcliff, J.W. and Metzener, D.E., Pattern matching: The Gestalt approach, *Dr. Dobbs's J.*, 1988, vol. 13, no. 7, pp. 46–72.
33. Abboud, A., Backurs, A., and Williams, V.V., Tight hardness results for LCS and other sequence similarity measures, *Foundations of Computer Science (FOCS), 2015 IEEE 56th Annual Symposium on, 2015*, pp. 59–78.
34. Python DiffLib module.
<https://docs.python.org/3/library/difflib.html>.
35. Abouelhoda, M.I., Kurtz, S., and Ohlebusch, E., Replacing suffix trees with enhanced suffix arrays, *J. Discrete Algorithms*, 2004, vol. 2, no. 1, pp. 53–86.
36. Špakov, O. and Miniotos, D., Visualization of eye gaze data using heat maps, *Elektronika ir elektrotechnika*, 2007, pp. 55–58.
37. Luciv, D.V., Detecting Near Duplicates in Software Documentation, 2017. arXiv: 1711.04705.
38. Boyer, R.S. and Moore, J.S., A fast string searching algorithm, *Commun. ACM*, 1977, vol. 20, no. 10, pp. 762–772.
39. Pandoc: A Universal Document Converter.
<https://pandoc.org/>.
40. Drobintsev, P. D. A formal approach to test scenarios generation based on guides / P. D. Drobintsev, V. P. Kotlyarov, A. A. Letichevsky // *Automatic Control and Computer Sciences*, 2014, Dec., vol. 48, no. 7, pp. 415–423.
41. Pakulin, N.V. and Tugaenko, A.N., Model-based testing of Internet Mail Protocols, *Proc. Inst. System Programming*, 2011, vol. 20, pp. 125–141.
42. Gorovoy, V.A., Bolotnikova, E.S., and Gavrilova, T.A., To a method of evaluating ontologies, *J. Comput. Systems Sci. Int.*, 2011, vol. 50, no. 3, pp. 448–461.