

**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

BACHELOR THESIS

Arkadiusz Martin Antoniewicz

Tokenization-aware Diff and Patch

Department of Software Engineering

Supervisor of the bachelor thesis: RNDr. Miroslav Kratochvíl

Study programme: Computer Science

Study branch: Programming and Software
Development

Prague 2021

I declare that I carried out this bachelor thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date

Author's signature

First and foremost, I would like to thank my supervisor Miroslav Kratochvíl for his patient guidance and professional advice. Next, I would like to express gratitude to Simona Zálešáková and my family for boundless support and faith in me. Last but not least, I am very grateful to Martina Matuláková and Jaromír Šimonek for their help with scientific English.

Title: Tokenization-aware Diff and Patch

Author: Arkadiusz Martin Antoniewicz

Department: Department of Software Engineering

Supervisor: RNDr. Miroslav Kratochvíl, Department of Software Engineering

Abstract: File comparison algorithms and utilities 'diff', 'patch' and 'diff3' are widely used in programming for the purpose of code comparison, and in many version control systems. Despite the usefulness, the differences and patches produced by the tools are strictly line-oriented, which complicates processing of differently formatted data, such as free flowing text, markup, and various other formats where line breaks are not crucial. This thesis describes and implements a customizable version of these tools, which allows the user to specify an arbitrary tokenization of the input, thus allowing easy diffing, patching and change-merging of content not supported by the traditional diff. Additionally, the thesis describes a newly appearing challenge of managing the whitespace in the patches, and demonstrates the functionality on a practical use-case that can not be performed with the current diff utilities.

Keywords: editing distance, three-way merge, text algorithms, version control

Contents

Introduction	3
1 Algorithms for comparing text	7
1.1 Diff implementation	8
1.1.1 Tokenization	8
1.1.2 Edit distance	9
1.1.3 Wagner and Fischer algorithm	9
1.1.4 Backtrack algorithm to find edit operations in matrix	9
1.1.5 Other algorithms used for diff purpose	10
1.2 Merging and applying changes	12
1.2.1 Patch	12
1.2.2 Three-way merge	14
2 Custom tokenization support	17
2.1 Lexing specification	17
2.1.1 Specifying REDFA using strings	21
2.2 Whitespace handling	22
2.2.1 Resolving whitespace conflicts	23
3 Implementation and results	27
3.1 Program structure	27
3.1.1 Shared library	27
3.1.2 TDiff	29
3.1.3 TPatch	29
3.1.4 TDiff3	30
3.2 Data formats	30
3.2.1 Tokenizer specification	30
3.2.2 Patch file format	30
3.3 Performance and use cases	31
3.3.1 Tokenization performance	31
3.3.2 Use case	32

3.4	GIT integration	33
	Conclusion	37
	Bibliography	39
A	Using tdiff	41

Introduction

Text comparison is an essential part of working with computers. Not only programmers but also many other professions use text comparison tools on a daily basis. There exist many different file and text comparison tools [2]. The text comparison tools include finding and showing differences, as well as revisiting, modifying and applying changes, comparing and showing differences in three files, merging three files together and many other cases of usage.

There is a problem with all common existing solutions. They do not allow a user to choose the unit of comparison, whether it may be a section, sentence, word, cells in a table or anything else. By considering a text with multiple changes in a long section we can present why common existing solutions fail to provide a convenient way of working with them. GNU Diff [8] (shown in Figure 1) simply shows that the lines are different but does not point to a place in the section where the difference is. Longer sections would make diff inapplicable. GNU Wdiff [5] (shown in Figure 2) demonstrates the difference in a more profound matter than the diff, however it lacks tools to patch and merge. Other solutions such as Beyond Compare [1] or the git diff [4] with `word-diff=color` option (shown in Figure 3) are capable of patching and showing the difference well. However, there is still a problem with merging. Neither of those tools would be able to three way merge if one file had changes at the beginning of a section and the other file had changes at the end of the same section.

The aim of this thesis is to design and implement tools that can work with various file formats, print readable differences and apply them. To be capable of working with many different formats, the user needs to be able to divide the text into sections of their own accord, which then they compare to each other. The process of text division is called tokenization and the results are called tokens.

The implemented utilities should be able to tokenize texts using rules defined by the user, work with the tokenized text effectively and show readable differences between them. This results in a multipurpose tool of comparing any text file format.

Notably, the custom tokenization creates a problem not present in other diff implementations. When a part of the text is left untokenized, it is considered

```

*** t1          2020-07-12 11:26:02.268930863 +0200
--- t2          2020-07-12 11:26:01.728660862 +0200
*****
*** 1,2 ****
! In mathematical theory, linguistics and computer science, the Levenshtein
distance is a string metric for measuring the difference between two sequences.
Informally, the Levenshtein distance between two words is the minimum number of
single-character edits (insertions, deletions or substitutions) required to
change one word into the other. It is named after the Russian mathematician
Vladimir Levenshtein, who considered this distance in 1965.
! Levenshtein distance may also be referred as edit distance, although that
term may also denote a larger family of distance metrics known collectively as
edit distance. It is not closely related to pairwise string alignments.
--- 1,2 ----
! In information theory, linguistics and computer science, the Levenshtein
distance is a string metric for measuring the difference between two sequences.
Informally, the Levenshtein distance between two words is the minimum number of
single-character edits (insertions, deletions or substitutions) required to
change one word into the other. It is named after the Soviet mathematician
Vladimir Levenshtein, who considered this distance in 1965.
! Levenshtein distance may also be referred to as edit distance, although that
term may also denote a larger family of distance metrics known collectively as
edit distance. It is closely related to pairwise string alignments.

```

Figure 1 GNU diff used on a text with multiple changes in the same section.

as a whitespace which is not significant for comparing. However, a whitespace around the tokens may sometimes carry information that is relevant for the result, and thus needs to be handled separately.

Layout of this Thesis

This thesis is structured as follows: the first chapter provides a detailed overview on handling text differences. There is described how to compare text, how to apply patches and how to compare and merge three files. The second chapter outlines the proposed solution for user-specified tokenization, how to implement it and what the lexers are in general. Whitespace handling is also addressed in the second chapter. The attention of the third chapter is focused on the program structure, tokenizer itself, patch format specification and the performance of the implemented tools.

In [-mathematical-] +information+ theory, linguistics and computer science, the Levenshtein distance is a string metric for measuring the difference between two sequences. Informally, the Levenshtein distance between two words is the minimum number of single-character edits (insertions, deletions or substitutions) required to change one word into the other. It is named after the [-Russian-] +Soviet+ mathematician Vladimir Levenshtein, who considered this distance in 1965.

Levenshtein distance may also be referred +to+ as edit distance, although that term may also denote a larger family of distance metrics known collectively as edit distance. It is [-not-] closely related to pairwise string alignments.

Figure 2 GNU wdiff used on a text with multiple changes in the same section.

```
@@ -1,2 +1,2 @@
In mathematicalinformation theory, linguistics and computer science, the
Levenshtein distance is a string metric for measuring the difference between
two sequences. Informally, the Levenshtein distance between two words is the
minimum number of single-character edits (insertions, deletions or
substitutions) required to change one word into the other. It is named after
the RussianSoviet mathematician Vladimir Levenshtein, who considered this
distance in 1965.
Levenshtein distance may also be referred to as edit distance, although that
term may also denote a larger family of distance metrics known collectively as
edit distance. It isnot closely related to pairwise string alignments.
```

Figure 3 GIT diff with enabled word diff color option used on a text with multiple changes in the same section.

Chapter 1

Algorithms for comparing text

Utilities for comparing texts are used by programmers on a daily basis. Probably the biggest application of these utilities are version control systems. For the proper development of large projects, it is important to store all versions of previous projects as a form of communication among the programmers. Every additional version of a project is called a revision. Every revision, except for the first one, originates in the previous one. When the revision needs to be checked — what has changed — the differences between current revision and the one it originated from need to be shown.

The three main tools used in comparing text are diff, patch and merge. The diff serves as a data comparison tool which displays the differences between two files. The changes made in a standard format, so that both humans and machines can understand them, are displayed by the diff. An example of how colored side-by-side comparison of two files looks like can be seen in Figure 1.1. The patch utility takes a comparison output produced by the diff and applies the differences to a copy of the original file, producing a patched version. Diff3 is used when two people make changes to the same base file. It can produce a merged output that contains changes from both files and warnings when conflict appears.

```
The Way that can be told of is not the eternal Way;  
The name that can be named is not the eternal name.  
The Nameless is the origin of Heaven and Earth;  
The Named is the mother of all things.  
  
Therefore let there always be non-being,  
so we may see their subtlety,  
And let there always be being,  
so we may see their outcome.  
The two are the same,  
But after they are produced,  
they have different names.
```

```
<  
<  
> The Nameless is the origin of Heaven and Earth;  
| The named is the mother of all things.  
>  
> Therefore let there always be non-being,  
so we may see their subtlety,  
And let there always be being,  
so we may see their outcome.  
The two are the same,  
But after they are produced,  
they have different names.  
> They both may be called deep and profound.  
> Deeper and more profound,  
> The door of all subtleties!
```

Figure 1.1 The colored diff in a side-by-side format.

```

*** Lao 2019-07-14 12:21:15.548156075 +0200
--- tzu 2019-07-14 12:21:34.348156075 +0200
*****
*** 1,7 ****
- The Way that can be told of is not the eternal Way;
- The name that can be named is not the eternal name.
  The Nameless is the origin of Heaven and Earth;
! The Named is the mother of all things.
  Therefore let there always be non-being,
  so we may see their subtlety,
  And let there always be being,
--- 1,6 ----
  The Nameless is the origin of Heaven and Earth;
! The named is the mother of all things.
!
  Therefore let there always be non-being,
  so we may see their subtlety,
  And let there always be being,
*****
*** 9,11 ****
--- 8,13 ----
  The two are the same,
  But after they are produced,
  they have different names.
+ They both may be called deep and profound.
+ Deeper and more profound,
+ The door of all subtleties!

```

Figure 1.2 The colored diff in a context format.

1.1 Diff implementation

The diff produces differences between two files. One way to achieve this is by computing the edit distance (as seen in Section 1.1.2) using the Wagner-Fischer algorithm (as seen in Section 1.1.3) and using the output of the Wagner-Fischer to find a sequence of insertions, substitutions and deletions to get from one text to another (as seen in Section 1.1.4).

1.1.1 Tokenization

Before we describe algorithms, we need to explain what tokenization and tokens are. Tokenization is a process of demarcating sections of a string of input characters. Tokenizers are usually designed to use a regular grammar (although it usually can not be achieved). An output is a list of tokens. Unlike parsing, which is usually a context-free grammar, the output is an abstract syntax tree. The parsing results in obtaining more information about the input and it is, certainly, more complex. We are going to analyze the tokenization more thoroughly in the Chapter 2.

Parts which are compared in the text are tokens. In edit distance, each character is a single token. In diff, tokenization is done by splitting text with delimiters being newlines. Each line is a single token. Tokens are comparable — it is possible to determine whether they are equal or not.

1.1.2 Edit distance

To be able to tell whether two strings are similar and by how much we need a metric that results in a number for two given strings. Many different types of edit distance exist. The most common being Levenshtein distance [12]. It uses 3 edit operations: changing one character to another single character, deleting one character from a given string and inserting a single character into the given string. Another type of edit distance is called longest common subsequence where only insertions and deletions are allowed, substitutions are not. The amount of such operations needed to change one string into another is called edit distance.

1.1.3 Wagner and Fischer algorithm

The Wagner-Fischer algorithm [12] is an algorithm used for finding edit distance. There are two strings as an input (can be applied to any two lists of items that can be compared). The computing is based on the following observation. If we reserve a matrix to hold edit distances between all the prefixes of the first string and all the prefixes of the second one, then the values in the matrix can be computed by flood filling the matrix, and thus the distance between the two full strings can be determined as the last value computed. An example of such implementation can be observed in Algorithm 1.

Definition 1 (Notation). *Let A and B be arrays of tokens and a and b be the tokens. Define $A\langle i \rangle = A\langle 1 : i \rangle$, $B\langle j \rangle = B\langle 1 : j \rangle$, and $D(i, j) = \gamma(A\langle i \rangle, B\langle j \rangle)$, $0 \leq i \leq |A|$, $0 \leq j \leq |B|$.*

$\gamma(a \rightarrow b)$ is 0 if a equals b otherwise is 1

$$D(i, j) = \min\{D(i-1, j-1) + \gamma(A\langle i \rangle \rightarrow B\langle j \rangle), \\ D(i-1, j) + \gamma(A\langle i \rangle \rightarrow \Lambda), \\ D(i, j-1) + \gamma(\Lambda \rightarrow B\langle j \rangle)\}$$

for all i, j , $1 \leq i \leq |A|$, $1 \leq j \leq |B|$.

$$D(0, 0) = 0; D(i, 0) = \sum_{r=1}^i \gamma(A\langle r \rangle \rightarrow \Lambda); D(0, j) = \sum_{r=1}^j \gamma(\Lambda \rightarrow B\langle r \rangle)$$

1.1.4 Backtrack algorithm to find edit operations in matrix

We can apply edit distance matrix of substrings filled by the Wagner Fischer algorithm to find edit operations. An example of such implementation can be

```

--- a/Testdata/t1.c
+++ b/Testdata/t2.c
@@ -1,9 +1,9 @@
-int add(int a, int b)
+int add(int a, int b, int c)
+  {
-   return a + b;
+   return a + b + c;
+ }

-int add(int a, int b, int c)
+int add(int a, int b)
+  {
-   return a + b + c;
+   return a + b;
+ }

--- a/Testdata/t1.c
+++ b/Testdata/t2.c
@@ -1,9 +1,9 @@
-int add(int a, int b)
-  {
-   return a + b;
- }
-int add(int a, int b, int c)
+  {
+   return a + b + c;
+ }
+int add(int a, int b)
+  {
+   return a + b;
+ }

```

Figure 1.3 Myers algorithm (left) and histogram algorithm (right). In this case histogram produce better arranged output, but has more edit operations.

seen in Algorithm 1 as a backtrack function. The algorithm starts at the right bottom cell of the matrix (edit distance between the two full strings). It finds a path in which the last cell was taken to fill. The path is always nondecreasing and ambiguous. As an example solution for strings 'ac' and 'b' are deletion 'a' and substitution 'b' for 'c'. The second possible solution is substitution 'b' for 'a' and deletion 'c'. Both solutions are of length 2 and are correct.

1.1.5 Other algorithms used for diff purpose

As edit operations are not uniquely determined, other diff algorithms can result in different outputs. This leads to that on various text formats different algorithms could provide better user readable outputs than the other ones. Sometimes it also could be beneficial not to find the smallest edit distance. This leads to better performance on some use cases. Myers' diff is one example of such algorithm [9]. Its time complexity is $O((m+n)*d)$ where m and n are lengths and d is number of edits. As we can see in scenarios where there are small amounts of edits in large files, this provides much better execution time. This is used in GIT version control system as default diff algorithm. Other example is Histogram algorithm which is derived from patience algorithm. It creates histogram of occurrences for each element and tries to match positions recursively [7].

Algorithm 1 Wagner Fischer algorithm to fill a matrix with edit distances of substrings. The backtrack algorithm to determine edit operations.

```

1: procedure WAGNER AND FISCHER( $D$ )
2:    $D[0, 0] \leftarrow 0$ 
3:   for  $i \leftarrow 1, |A|$  do
4:      $D[i, 0] \leftarrow D[i - 1, 0] + \gamma(A\langle i \rangle \rightarrow \Lambda)$ 
5:   end for
6:   for  $j \leftarrow 1, |B|$  do
7:      $D[0, j] \leftarrow D[0, j - 1] + \gamma(\Lambda \rightarrow B\langle j \rangle)$ 
8:   end for
9:   for  $i \leftarrow 1, |A|$  do
10:    for  $j \leftarrow 1, |B|$  do
11:       $m_1 \leftarrow D[i - 1, j - 1] + \gamma(A\langle i \rangle \rightarrow B\langle j \rangle)$ 
12:       $m_2 \leftarrow D[i - 1, j] + \gamma(A\langle i \rangle \rightarrow \Lambda)$ 
13:       $m_3 \leftarrow D[i, j - 1] + \gamma(\Lambda \rightarrow B\langle j \rangle)$ 
14:       $D[i, j] \leftarrow \min(m_1, m_2, m_3)$ 
15:    end for
16:  end for
17: end procedure
18: procedure BACKTRACK( $D$ )
19:    $i \leftarrow |A|$ 
20:    $j \leftarrow |B|$ 
21:   while  $i \neq 0$  and  $j \neq 0$  do
22:     if  $D[i, j] = D[i - 1, j] + \gamma(A\langle i \rangle \rightarrow \Lambda)$  then
23:        $i \leftarrow i - 1$ 
24:       print("Addition : ",  $A\langle i \rangle$ )
25:     else if  $D[i, j] = D[i, j - 1] + \gamma(\Lambda \rightarrow B\langle j \rangle)$  then
26:        $j \leftarrow j - 1$ 
27:       print("Deletion : ",  $B\langle j \rangle$ )
28:     else
29:        $i \leftarrow i - 1$ 
30:        $j \leftarrow j - 1$ 
31:       if  $A\langle i \rangle \neq B\langle j \rangle$  then
32:         print("Substitution : ",  $A\langle i \rangle, B\langle j \rangle$ )
33:       end if
34:     end if
35:   end while
36: end procedure

```

1.2 Merging and applying changes

1.2.1 Patch

The output of the diff is not only for people but also for other programs as well. A patch is a program that takes an output of a diff and applies it. The utility that updates text files according to instructions is called a patch. The instructions are produced by the diff. This may seem to have no use. Comparing files and then applying changes to the first file, results in forming of the second file. The use of the patch is to be able to review the output of the diff, adjust or remove some of the changes and apply them afterwards.

The GNU Patch manual [8] describes the patch algorithm as follows:

"The patch reads instructions and applies them to the file (as seen in Figure 1.4). As for context diffs, patch can detect when the line numbers mentioned in the patch are incorrect, and it attempts to find the correct place to apply each hunk of the patch. A hunk is a sequence of lines common to both files, interspersed with groups of differing lines. As a first guess, it takes the line number mentioned in the hunk, plus or minus any offset used in applying the previous hunk. If that is not the correct place, the patch makes a forward and a backward scan for a set of lines to match the context given in the hunk.

At first, the patch looks for a place where all lines of the context match. If it cannot find such place, and it reads a context or a unified diff and the maximum fuzz factor is set to 1 or more, then the patch makes another scan, ignoring the first and the last line of the context. If that fails, and the maximum fuzz factor is set to 2 or more, it makes a scan again, ignoring the first two and the last two lines of context. It behaves similarly if the maximum fuzz factor is larger.

If the patch cannot find a place to install a hunk of the patch, it writes the hunk out to a reject file. The line numbers on the hunks in the reject file may be different from those in the patch file: they show the approximate location where the patch thinks the failed hunks belong in the new file rather than in the old one.

The patch usually produces correct results, even when it makes many guesses. However, the results are guaranteed only when the patch is applied to an exact copy of the file that the patch was generated from."

Algorithm 2 Patch pseudocode.

```
1: hunks ← parseContext()
2: file ← readFile()
3: for all h ← hunks do
4:   p ← position(h)
5:   if h.match(file.atPosition(p)) then
6:     applyContext(h, file)
7:   else if h.match(neighborhood(file, p)) then
8:     applyContext(h, file)
9:   else
10:    saveRejectedHunk(h)
11:  end if
12: end for
```

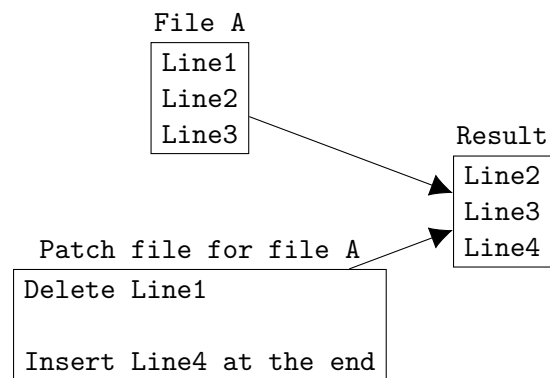


Figure 1.4 Simplified patching.

1.2.2 Three-way merge

A more interesting phenomenon than patching is merging. One possible way of merging is a three-way merge algorithm. Two files that are merged and their common ancestor (base) are considered. The result is a single file containing both sets of changes. Let us call chunk, fragments of text from all 3 files. Stable chunk is a chunk where all 3 fragments of text in files are the same. Unstable chunk is when fragments differ. Possibilities of unstable chunks are, changes in one of the merged files, falsely conflicting chunk – changes in both merged files, but the change is the same and conflicting chunk – all 3 fragments differ.

To sum it up, the result consists of:

- The parts matching in all 3 files (stable chunk)
- The parts not matching the base in one of the merging files (unstable chunk)
- The parts not matching the base but matching each other (falsely conflicting unstable chunk)
- Placeholders for parts where all of them are different marked as conflict to be resolved (conflicting chunk)

An example with all types of chunks can be seen in Figure 1.5

Firstly it finds differences between the base and each merging file. It starts applying both changes on the base file from the beginning. If two differences which are not the same (unstable chunk) exist, one in A file and one in B, that should be applied on the same line (or adjacent to each other), algorithm can not decide which change to apply or how to merge them together. It marks the place where differences should be applied as conflict and it is left to the user to be resolved. All other differences, that do not interfere with each other or are of the same change on the same line can be applied on the base file.

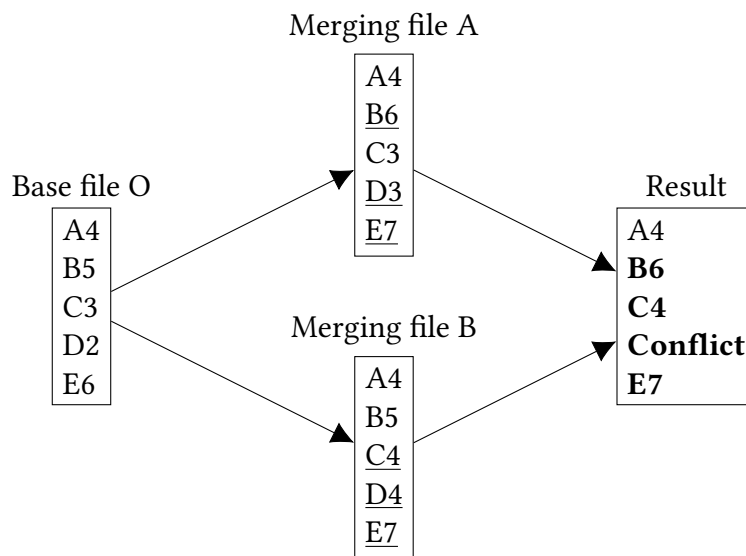


Figure 1.5 The three-way merge with the base file O and two files A and B. The result is composed of: A4, a stable chunk. B6 changed in A. C4 changed in B. Next, there is a truly conflicting chunk. E7 is a falsely conflicting chunk.

Chapter 2

Custom tokenization support

There are many different versions of diffs. They differ in the application and the way of displaying changes. Some of them are designed for finding differences in specific file formats. HTML diff tries to compare not only the source codes but also the appearance of the final webpage. XML diffs compare the hierarchical structure of XML documents. There is a word comparing option in the Gits implementation of the diff (as seen in Figure 2.1) but it lacks the patch. None of the tools mentioned above allow the user to specify the tokenization. The user specifiable tokenization has an advantage in its wide variety of applications. It can result in a universal diff that could then handle any programming language and text format.

There are many possible existing solutions for the user-specifiable implementation of the tokenization process: regular expressions with capture groups, lexical analyzer generators such as the Lex and the Flex. In this thesis we are going to design and implement our own solution.

Terms such as automaton, regular grammar, (non)deterministic finite state machine etc. are used in this section. Their definitions can be found in a book by Hopcroft, Motwani, and Ullman [6].

2.1 Lexing specification

Most tokenizers are designed to use a regular grammar. Tokenizers are sometimes referred to as lexers. Although they share very similar properties, the difference between them is that a lexer usually attaches an extra context to the tokens. We are going to consider regular expression [6] and try to simplify the defining of a more complex tokenization using a deterministic finite state machine with regular expression as edges. Let us call it Regex Edge Deterministic Finite Automaton (REDFA). The definition of REDFA is similar to the definition

```

diff --git a/lao b/tzu
index 635ef2c..5af88a8 100644
--- a/lao
+++ b/tzu
@@ -1,7 +1,6 @@
[-The Way that can be told of is not the eternal Way;-]
[-The name that can be named is not the eternal name.-]
The Nameless is the origin of Heaven and Earth;
The [-Named-]{+named+} is the mother of all things.

Therefore let there always be non-being,
    so we may see their subtlety,
And let there always be being,
@@ -9,3 +8,6 @@ And let there always be being,
The two are the same,
But after they are produced,
    they have different names.
{+They both may be called deep and profound.+}
{+Deeper and more profound,+}
{+The door of all subtleties!+}

```

Figure 2.1 The git diff with enabled word option.

of the deterministic finite automaton but with an elaborate transition function. An example of REDFA can be seen in Figure 2.2.

Definition 2. *Definition of Regex Edge Deterministic Finite Automaton (REDFEA)* M is a 6-tuple, $(Q, \Sigma, R, \gamma, q_0, F)$, consisting of:

- a finite set of states Q
- a finite set of input symbols called the alphabet Σ
- a finite set of regular expressions, search patterns over alphabet R
- a finite set of tuples $(q_1, r, q_2) \in \gamma$, where $q_1, q_2 \in Q$ and $r \in R$
- an initial state $q_0 \in Q$
- a set of accept states $F \subseteq Q$

Let w, a_1, a_2, \dots, a_n be strings over the alphabet Σ . $w = a_1 a_2 \dots a_n$. The automaton M accepts the string w if a sequence of states, s_0, s_1, \dots, s_n , exists in Q with the following conditions:

1. $s_0 = q_0$
2. $r \in R : r$ accepts (is matching) $a_i \wedge (a_i, r, a_{i+1}) \in \gamma$
3. $s_n \in F$

Theorem 1. *Any language accepted by REDFA is a regular language.*

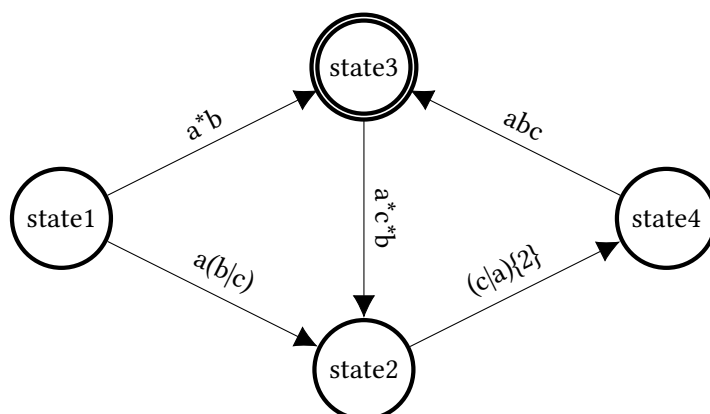


Figure 2.2 The state machine where edges are regular expressions.

Proof. At first, let us prove that a finite state machine with the edges as regular expression still fulfills the criteria for being a finite state machine.

We start with converting the regular expression to a NFA (nondeterministic finite automaton). This is called the Thompson algorithm [10]. The algorithm works recursively by splitting an expression into its constituent subexpressions, from which the NFA will be constructed using a set of rules. The constants and operations, which define a basis for the construction of the regular expression, are going to be used. The elementary constants are an empty expression ϵ and an expression with one symbol of alphabet. Operations are a union expression $|$, a concatenation expression and The Kleene star expression $*$. Firstly we convert elementary constants (as seen in Figure 2.3) and expand the constants with regular expression operations (as seen in Figure 2.4). Now we have an NFA instead of a regular expression. We can replace all regular expression edges in our REDFA with the NFA. Then we create an epsilon edge (empty expression ϵ) from the starting node of REDFA edge to the starting node of NFA created from the regular expression and also create epsilon edges starting in all the ending nodes of the created NFA to the ending node of the REDFA edge.

As for the next step, ordered edges are transformed to be a part of a nondeterministic finite-state machine, starting from the highest priority to the lowest priority edge for each node, unioning the complement with all edges with lower priority. The complement of state machine is done by reversing accepting and non accepting states. The complement of a regular expression $a(b|c)$ is shown in Figure 2.6. \square

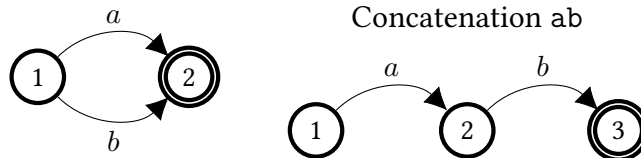
A regular finite state machine with ordered edges and edges as regular expressions (REDFA) is proven to accept regular language.

NFA representing an empty string NFA representing a



Figure 2.3 Elementary constants – empty and one character long.

The union operator $a | b$



The Kleene closure a^*

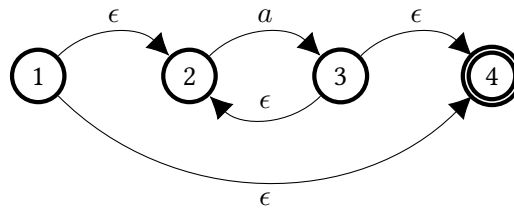


Figure 2.4 Converting a basic regular expression operator into NFA.

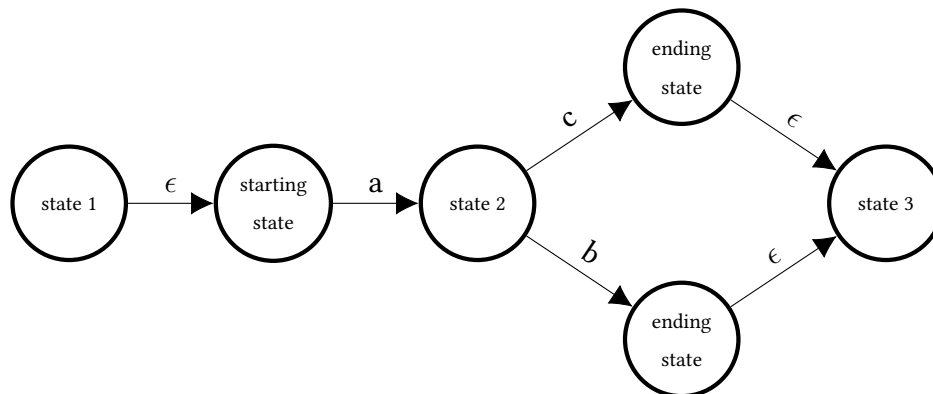


Figure 2.5 Inserting a regular expression $a(b | c)$ into state machine.

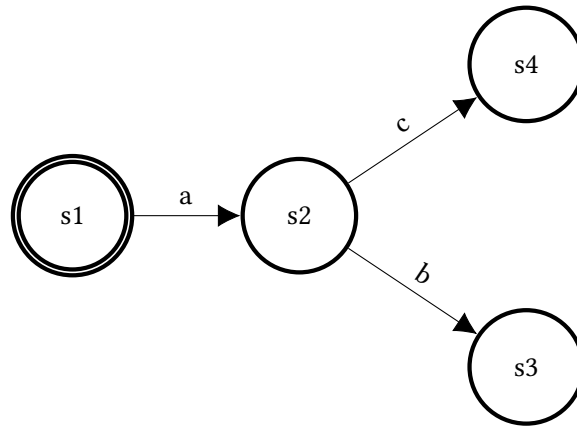


Figure 2.6 The complement of a regular expression $a(b|c)$ converted into an NFA.

2.1.1 Specifying REDFA using strings

Let us show how REDFA can be used to make user-specifiable tokenization easier. The usage is similar to using a simple regular expression. Capture groups are used to define tokens and anything that is not in any capture group is considered a whitespace.

Before showing the differences of a simple regular expression and REDFA we need to decide how to define REDFA. We can easily define REDFA by specifying a set of 3-tuples γ , 3-tuples consisting of (starting node, regular expression, ending node). This is sufficient to define REDFA:

- Q – All nodes in rules
- Alphabet Σ – same as the alphabet regular expression use
- R – All regular expressions in rules
- γ – It is the same as rules
- q_0 – Starting node of the first rule
- F – All nodes ($F = Q$)

Simple regular expression is shorter but it becomes unreadable in more complicated rules. On the other hand, REDFA definitions are easily readable and extendable. It is possible to use REDFA as a simple regular expression – an automaton with only one node and an edge going into itself.

Examples can be seen in Figure 2.7.

- Text definition of REDFA for words
`word, (\s*), whitespace`
`whitespace, \S*, word`
- Text definition of REDFA for CSV with header line
`header, .*\\r\\n, token`
`token, ([^\\r\\n,]*) , separator`
`separator, (?:\\r\\n|,) , token`

Figure 2.7 Examples of text defined REDFA.

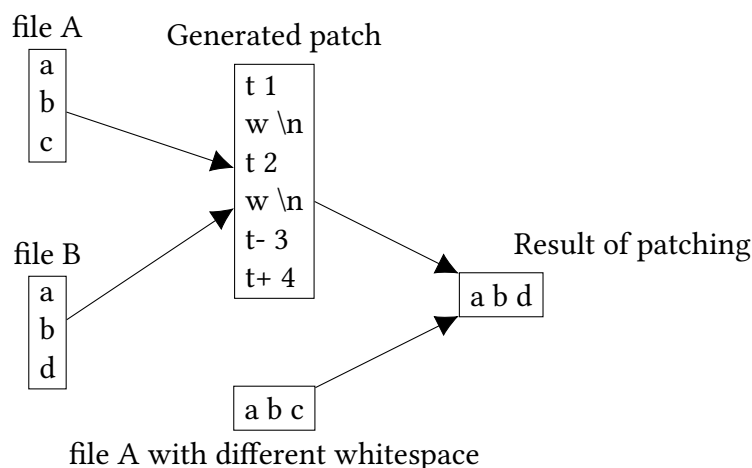


Figure 2.8 Example of patch, not failing when whitespaces changed in original file.

2.2 Whitespace handling

Whitespaces in this context are everything that is not compared in the text. The whitespaces in line-oriented diff are new lines, in word diff the whitespaces are the actual whitespaces (spaces, tabs, newlines etc.). In the user-defined tokenization the whitespaces are parts of the text that are between the tokens. In the line-oriented diff, whitespaces do not need to be handled because all the whitespaces are always the same. In the user-specified tokenization whitespaces can be anything, thus they need to be handled. A simple example why whitespaces needs to be handled can be seen in Figure 2.8.

Whitespace changes between tokens which are not changed nor shown within the context are not found by the diff and the patch. The whitespace from

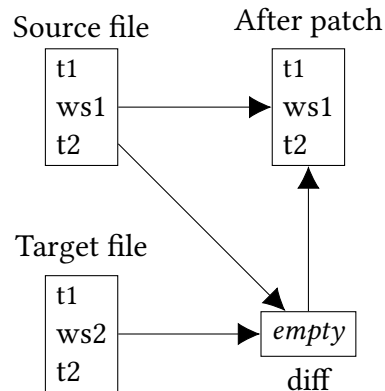


Figure 2.9 The change of a whitespace between non changing tokens. The diff is not able to find such change and the patch cannot patch it.

the source file is going to be used as shown in Figure 2.9 and the change is not detected. Only whitespace changes that are around the token (in the context) changes are found. When inserting a token, the whitespaces around the token being inserted, are inserted as well (as seen in Figure 2.10). When deleting a token, the whitespaces around the token, that is being deleted are deleted too and the whitespace from the target file is inserted as shown in Figure 2.11. During tokenizing, this needs to be considered. It is advised not to leave crucial parts of the text as whitespace because the program is not able to determine which whitespaces are to be used or deleted.

When the whitespace change is shown in a context but it is not directly located next to a token change, it is considered the same way as whitespace changes on a different place where they are not a part of the patch file. Therefore, the change is not applied. An example is shown in Figure 2.13.

2.2.1 Resolving whitespace conflicts

A conflict occurs when a whitespace in a patch file does not match the whitespace of the file that patch is applied to. When tokens in context match but whitespaces do not, the hunk can not be rejected as whitespaces are not significant. The whitespaces in patch files are used in the result. This example can be seen in Figure 2.12 – for deletion case when the diff is running, the source file has ws3 between t1 and t2 and ws4 between t2 and t3. So that the information about the change from ws1 to ws2 is not lost, the whitespaces are saved into a separate whitespace file, for insertion case when the diff is running, the source file has ws4 between t1 and t3, ws4 was changed to ws1 before the patch was running, so ws1 is saved into a whitespace file not to lose the information about having it.

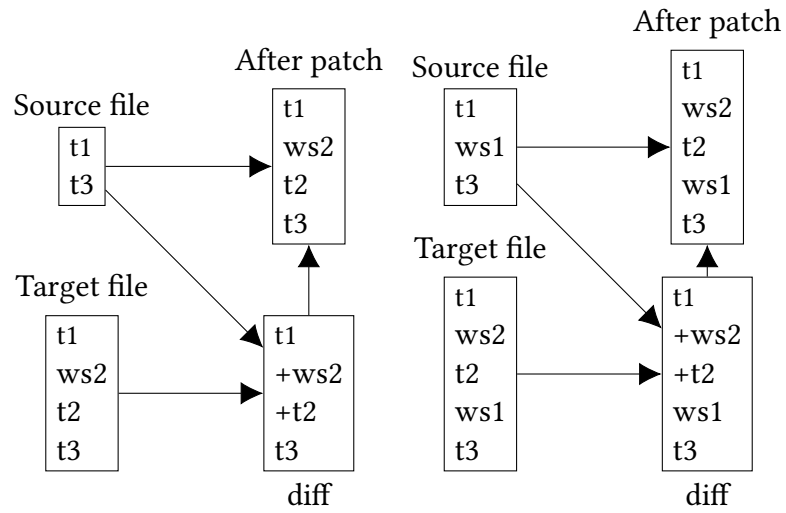


Figure 2.10 Insertion of token `t2`.

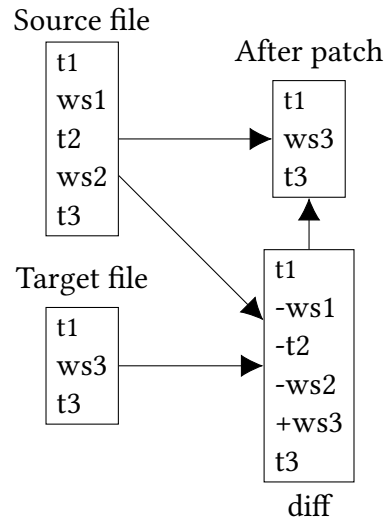


Figure 2.11 Deletion of token `t2`.

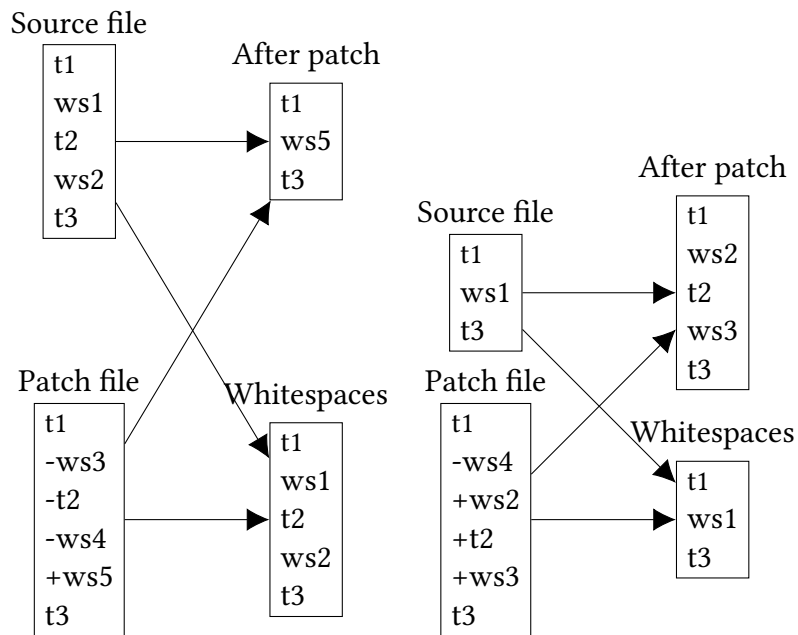


Figure 2.12 Deletion and insertion with different whitespaces in the patch file.

If the whitespaces do not change, the information about their deletion is saved in the patch file, and thus there is no need for them to be stored again.

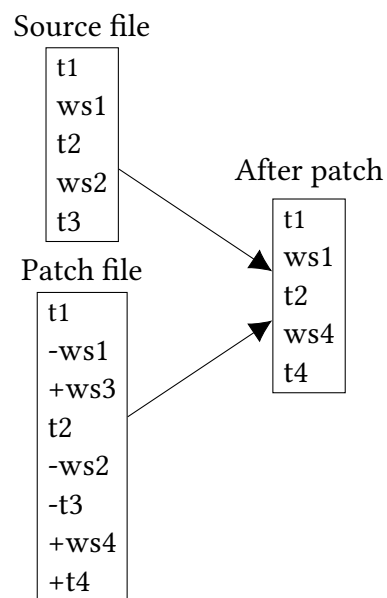


Figure 2.13 Differing whitespace in the context but not directly next to the token change. ws1 is not changed into ws3 despite the fact that the patch knows about the change.

Chapter 3

Implementation and results

Three utilities, TDiff, TPatch and TDiff3, have been implemented. They were programmed using the C++ programming language and the C++20 standard. Proof of concept of patching files without collisions in whitespaces is implemented in TPatch. TDiff3 supports token-by-token patching which is sufficient in most cases. In this chapter, the structure of the program and its implementation details are going to be introduced. After that we are going to show the results of a benchmark to see if it can handle larger files in reasonable time. At the end, we are going to show many different applications where having a user specifiable tokenization in the diff is superior to almost all the other diff utilities.

3.1 Program structure

The program is divided into 3 standalone executables and one library which is used by all the projects. In every project there is a Main file and an InputOutput file. The Main is used for parsing arguments and calling appropriate methods from the InputOutput file. In the InputOutput file there is a logic of the program and it is calling the shared library methods. The program structure can be seen in Figure 3.1.

3.1.1 Shared library

The shared library contains 3 header files and their implementation.

- Tokenizer

The Tokenizer is used for parsing the text into tokens. There are classes used for the definition of REDFA and its edges. It also contains the definition of the parsed text which consists of the tokens and the text that

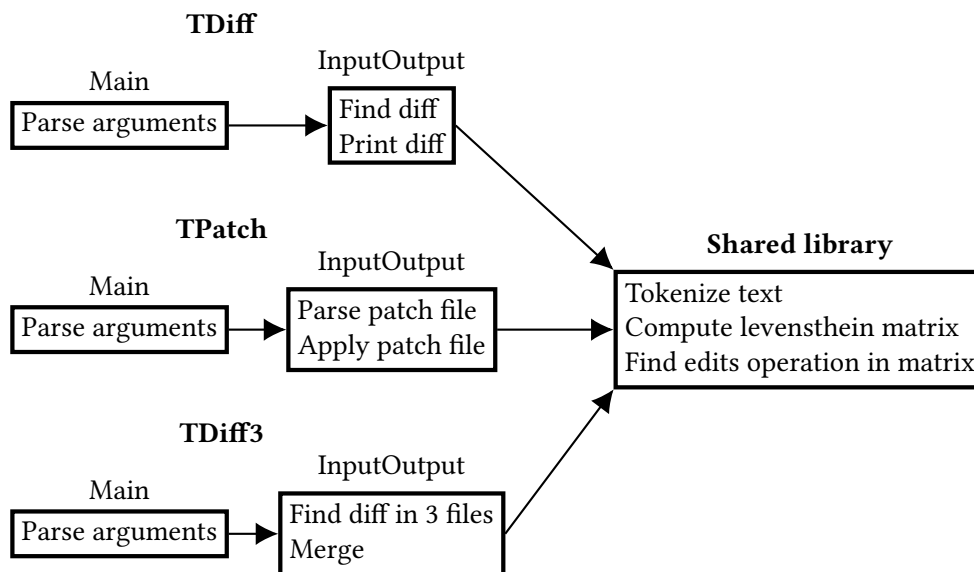


Figure 3.1 The program structure.

was tokenized. Methods for building the REDFA from the rules and for tokenizing the text with the REDFA are implemented as well.

- Differentiate

The Differentiate is used for finding the differences in two tokenized texts. It contains the definition of edit distance matrix, file differences and the difference. The file difference contains information about differentiating two files, paths to those files, their contents, the tokenized content and differences between them. One difference contains its type and indexes of tokens that are compared. There are methods for the Wagner and Fischer algorithm to create a matrix from the tokens and to backtrack the matrix to find the differences.

- Utils

The Utils contains other methods and classes – a method that takes two paths to files and an REDFA definition, creating an REDFA, reading files, comparing files and returning them filled into the file difference class, a method for replacing whitespaces (tab, newline) for printable characters and vice versa, a method for splitting one argument which contains all the rules of REDFA into separate rules and a method to group the diffs that are near each other within the context.

3.1.2 TDiff

TDiff is a program which compares two files or directories according to the given rules and prints the differences to a standard output. The Main parses arguments and calls the methods from InputOutput file. The InputOutput contains a method for comparing folders and files and for printing differences.

The format for executing the diff is `tdiff from to [options]`. From and to must be paths to either files or directories.

Below, there is a summary of all the options that the Token-aware diff accepts. Standard long and short arguments can be passed to the diff. The Getopt [3] is used to parse the input arguments.

<code>--automaton, -a AT</code>	Use AT to build an automaton to tokenize the input. The rules are in a format starting state, regular expression, ending state. The rules are separated by a semicolon. The semicolon in the rules needs to be escaped with a backslash.
<code>--context, -C NUM</code>	Output NUM (default 3) lines of context.
<code>--debug, -d</code>	Tokenize file and print the tokens.
<code>--help, -h</code>	Display help.
<code>--file-automaton, -f PATH</code>	Read the rules from PATH. The rules are delimited by newlines.

3.1.3 TPatch

TPatch is a program which applies patches on a file. The input is a patch file which contains information about which file is going to be patched and what changes are going to be applied. It applies the changes, rejected changes and whitespaces that were deleted are then saved to separate files. The Main parses arguments and calls methods from the InputOutput file. The InputOutput contains a method to parse the patch file and apply the changes.

The format for running the patch is `tpatch patch-file [options]`. The patch file must be in the format described further below.

<code>--ignorews, -i</code>	Ignore whitespaces during patching. Only the tokens are inserted and deleted.
<code>--output, -o PATH</code>	Write output to PATH instead of the path specified in the patch.
<code>--help, -h</code>	Display help.

3.1.4 TDiff3

TDiff3 is a program that reads 3 files and merges them together. The Main parses arguments and calls the method from the InputOutput file. The InputOutput contains a method of comparing three files, merging them and to printing the result.

The format for running the diff3 is `tdiff3 mine base yours [options]`. Mine, base and yours are the paths to three files.

<code>--automaton, -a AT</code>	Use AT to build an automaton to tokenize the input. The rules are in a format starting state, regular expression, ending state. The rules are separated by a semicolon. The emicolon in the rules needs to be escaped with a backslash.
<code>--file-automaton, -f PATH</code>	Read the rules from PATH. The rules are delimited by newlines.
<code>--help, -h</code>	Display help.

3.2 Data formats

3.2.1 Tokenizer specification

To define the REDFA we use something that has already been mentioned in Section 2.1.1. The REDFA is going to be defined as an ordered set of rules (edges). Each rule has a starting state, regular expression and an ending state. To tokenize the input, capture groups are used. It is possible for a capture group not to end in one edge. The token starts in one edge and ends in another one. This is achieved by allowing the regular expression to have an incomplete group structure. As an example we consider rules `1, a(bc, 2; 2, d)aa, 1` and input `abcdaa`. The tokenization is going to be successful with the application of these rules and the input will result in one token `bcd`.

3.2.2 Patch file format

The output of the diff can be seen in Figure 3.2. On the first two lines there is a path to the source and the target files. On the third line there is a definition of the REDFA delimited by newlines. After that there are hunks. All hunks start with a line consisting of asterisks. It is followed by two lines with indexes of tokens of a source and a target used in the hunk. The hunk is similar to the hunk in GNU patch.

```

*** t1.txt 2019-07-17 23:21:53 +0200
--- t2.txt 2019-07-17 23:21:22 +0200
@@ 1,(.*\r\n),1
*****
*** 1,6 ****
--- 0,5 ----
t- The Way that can be told of is not the eternal way;\r\n
*****
*** 10,12 ****
--- 9,14 ----
t But after they are produced,\r\n
t they have different names.\r\n
t+ They both may be called deep and profound.\r\n
t+ Deeper and more profound,\r\n
t+ The door of all subtleties!\r\n

```

Figure 3.2 Example of a token-aware diff tokenized into lines.

When a line should be deleted there is a minus at the beginning of the line. The same is applied with a plus and addition of a line. Space means that everything is on a place where it is supposed to be and it remains as it is. In the token-aware patch there are tokens and whitespaces instead of the lines but the notions are the same in both patches. To differentiate between the whitespaces and the tokens we use letters 'w' and 't'. After 'w' or 't' we put '+', '-' or ' ' as a second character to determine the type of the operation in the patch. All these can be seen in Figure 3.2.

3.3 Performance and use cases

3.3.1 Tokenization performance

The benchmarks are done using Ubuntu 18.04.4 on the virtual machine with the host running Windows 10 1903. The tool for measurements is linux `time` [11] utility. Presented results are always mean time of ten runs with the slowest run being discarded. We are going to redirect the output to `/dev/null` as writing output can add overhead. The text will be generated lorem ipsum.

To measure the tokenization time we are going to use the diff debug option. With this option only tokenization is going to run. The results can be seen in Table 3.1. The results suggest what we expected:

- For the same REDFA, execution time depends linearly on the length of the text.
- For almost the same REDFA (only capture group changed), the execution time depends on the number of tokens.
- Regular expressions definitions can heavily affect the performance.

File size	REDFA definition	Mean time
200000 words~1.4MB	1,(\S*\s*),1	0.125s
2000000 words~13.5MB		1.274s
200000 words~1.4MB	1,(\S*\s*),2; 2,\S*\s*,1	0.088s
2000000 words~13.5MB		0.864s
200000 words~1.4MB	1,(\S+),1; 1,\s*,1	0.190s
2000000 words~13.5MB		1.901s

Table 3.1 The tokenization benchmark.

```

In [-mathematical-]
+information+ theory,
linguistics and computer
science, the Levenshtein
distance is a string
metric for measuring the
difference between two
sequences. Informally,
the Levenshtein distance
between two words is the
minimum number of
single-character edits
(insertions, deletions or
substitutions) required
to change one word into
the other. It is named
after the [-Russian-]
+Soviet+ mathematician
Vladimir Levenshtein, who
considered this distance
in 1965.
*** 1,4 ****
--- 1,4 ----
t In
t- information
t+ mathematical
t theory,
t linguistics
*****
*** 51,55 ****
--- 51,55 ----
t after
t the
t- Soviet
t+ Russian
t mathematician
t Vladimir

```

Figure 3.3 GNU wdiff and tdiff.

To sum up, the tokenization runs reasonably quickly even on larger files, however, not optimal definition of REDFA (both regular expression itself as well as DFA) can slow down a great deal of the execution time.

3.3.2 Use case

First, let us compare the long section we have already mentioned in the intro. As we can see in Figure 3.3, both GNU diff and tdiff produce precise and readable output, however only tdiff is capable of running tpatch in this format.

Next, let us show a simple case where patching using the GNU utilities fails, but tdiff and tpatch are able to handle it. Firstly, let us consider two simple c files with small changes (as can be seen in Figure 3.4). Then we run diff between this two files to produce the patching file. The outputs of diffs can be seen in Figure 3.5. As we can observe, the tdiff output is more verbose. Let us see what

```

#include <stdio.h>
int add(int a,int b)
{
    return a + b;
}

int main()
{
    int a = 5;
    int b = 4;

    printf("Hello, World!");
    printf("%d",add(a, b));
    return 0;
}

#include <stdio.h>
int add(int a,int b,int c)
{
    return a + b + c;
}

int main()
{
    int a = 5;
    int b = 4;
    int c = 6;

    printf("Hello, World!");
    printf("%d",add(a,b,c));
    return 0;
}

```

Figure 3.4 Two C codes with small changes.

happens when we change the formatting of the source file (the file we are applying the patch to). The changed file and the tpatch result can be seen in Figure 3.6. The GNU patch fails to apply anything to the changed file, on the other hand, the tpatch handles it correctly.

Another example we can offer is a three-way merge. Considering two files already mentioned earlier in Figure 3.4, we add the third file to those 2 with changed Hello world to Hi. The GNU Diff3 fails to produce merged output of these 3 files. However, the tdiff3 is capable of doing a correct merge as can be seen in Figure 3.7.

3.4 GIT integration

To use the utilities with GIT, it is easy to configure git difftool and git mergetool. One of many ways to make it work can be seen in Figure 3.8

```

*****
*** 9,10 ****
--- 9,13 ----
w
t b
t+ ,
w+
t+ int
w+
t+ c
t )
*****
*** 15,16 ****
--- 18,21 ----
w
t b
w+
t+ +
w+
t+ c
t ;
*****
*** 32,33 ****
--- 37,43 ----
t ;
w+ \n
t+ int
w+
t+ c
w+
t+ =
w+
t+ 6
t+ ;
w \n\n
t printf
*****
*** 44,45 ****
--- 54,57 ----
w
t b
t+ ,
w+
t+ c
t )

```

```

@@ -2,5 +2,5 @@
-int add(int a, int b)
+int add(int a, int b, int c)
{
-   return a + b;
+   return a + b + c;
}
@@ -11,5 +11,6 @@
   int b = 4;
+   int c = 6;

   printf("Hello, World!");
-   printf("%d",add(5, 4));
+   printf("%d",add(5, 4, 6));
   return 0;

```

Figure 3.5 The token aware diff and the GNU diff output.

```

#include <stdio.h>

int add(int a, int b){
    return a+b;
}

int main(){
    int a=5;
    int b=4;
    printf("Hello , World!");
    printf("%d", add (a, b) );
    return 0;
}

#include <stdio.h>

int add(int a, int b, int c){
    return a+b + c;
}

int main(){
    int a=5;
    int b=4;
    int c = 6;
    printf("Hello , World!");
    printf("%d", add (a, b, c) );
    return 0;
}

```

Figure 3.6 tdiff is able to apply patches even to reformatted code. Left: Code from Figure 3.4 with changed coding style. Right: Token-aware patching is able to apply the patches from Figure 3.5 even in the reformatted code.

```

#include <stdio.h>

int add(int a,int b)
{
    return a + b;
}

int main()
{
    int a = 5;
    int b = 4;
    printf("HI!");
    printf("%d\n",add(a, b));
    return 0;
}

int add(int a,int b,int c)
{
    return a + b + c;
}

int main()
{
    int a = 5;
    int b = 4;
    int c = 6;
    printf("HI!");
    printf("%d\n",add(a,b,c));
    return 0;
}

```

Figure 3.7 Different changes in the original file from Figure 3.4 can be merged with the other patches using tdiff3. Left: The new modification. Right: Merged patches applied to the file.

```
git config --global diff.tool tdiff
git config --global difftool.tdiff.cmd "/path/to/tdiff
    \ $LOCAL \ $REMOTE -f /path/to/file/with/definitions"
git config --global difftool.tdiff.trustExitCode false

git config --global merge.tool tdiff3
git config --global mergetool.tdiff3.cmd "/path/to/tdiff3
    \ $LOCAL \ $BASE \ $REMOTE -f /path/to/file/with/definitions"
git config --global mergetool.tdiff3.trustExitCode false
```

Figure 3.8 GIT configuration with utilities.

Conclusion

In this thesis, we have designed the user-specifiable tokenization and implemented tools for differentiating text files using the user-specifiable tokenization.

In the Chapter 1 we have discussed the text difference handling. We have described the Wagner and Fischer algorithm for text comparing, algorithms for patching and merging three files.

In the Chapter 2 we have designed a solution for generic tokenization. We have created our own form of lexer. We also proposed a way of handling whitespace conflicts when working with generic tokenization.

In the Chapter 3 we have described the implementation of a program, the specification of defining the tokenizer and the tdiff output format to be able to consider whitespaces and be readable at the same time. We have also measured the performance of tdiff, verifying its sufficient fastness, and the performance scales as predicted by the asymptotic complexities of the used algorithms.

Bibliography

- [1] *Beyond Compare*. Web page. 2020. URL: <https://www.scootersoftware.com/>.
- [2] *Comparison of file comparison tools*. Web page. 2020. URL: https://en.wikipedia.org/wiki/Comparison_of_file_comparison_tools.
- [3] *Getopt manual page*. Web page. 2017. URL: https://www.gnu.org/software/libc/manual/html_node/Getopt.html.
- [4] *Git Diff*. Web page. 2020. URL: <https://git-scm.com/docs/git-diff/>.
- [5] *GNU Wdiff*. Web page. 2020. URL: <https://www.gnu.org/software/wdiff/>.
- [6] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation (3rd Edition)*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 2006. ISBN: 0321455363.
- [7] *jgit Histogram implementation*. Web page. 2020. URL: <https://github.com/eclipse/jgit/blob/ebfd62433a58d23af221adfdffed56d9274f4268/org.eclipse.jgit/src/org/eclipse/jgit/diff/HistogramDiff.java>.
- [8] D. MacKenzie, P. Eggert, and R. Stallman. *Comparing and Merging Files with Gnu Diff and Patch*. Network Theory, 2003. ISBN: 9780954161750. URL: <https://books.google.cz/books?id=oIINAAAACAAJ>.
- [9] Eugene W. Myers. “An O(ND) Difference Algorithm and Its Variations”. In: *Algorithmica* 1 (1986), pp. 251–266.
- [10] Ken Thompson. “Programming Techniques: Regular Expression Search Algorithm”. In: *Commun. ACM* 11.6 (June 1968), 419–422. ISSN: 0001-0782. DOI: 10.1145/363347.363387. URL: <https://doi.org/10.1145/363347.363387>.
- [11] *time - Linux manual page*. Web page. 2020. URL: <https://man7.org/linux/man-pages/man1/time.1.html>.

- [12] Robert A Wagner and Michael J Fischer. “The string-to-string correction problem”. In: *Journal of the ACM (JACM)* 21.1 (1974), pp. 168–173.

Appendix A

Using tdiff

To compile and run the software, you need:

1. GCC compiler with version at least 8.1 (filesystem)
2. POSIX compatible mmap in header `<sys/mman.h>`
3. POSIX compatible getopt in header `<getopt.h>`
4. External library RE2 (contained in debian package `libre2-dev`). On Debian-based Linux systems (such as Ubuntu), you may install this dependency with:

```
git clone https://code.google.com/re2
cd re2
make
make test
make install
make testinstall
```

To unpack and compile the software, proceed as follows:

```
unzip tdiff.zip
cd tdiff/tdiff
make
```

The following example shows the usage of `tdiff` and `tpatch` with tokens being printable characters delimited by whitespace characters and automaton specified in a command line:

```
tdiff file1 file2 -a '1,(\S*)\s*,1' > diffoutput
tpatch diffoutput
```

An example below presents the usage of `tdiff3` with `automaton` specified in the file `automatondef` is:

```
tdiff3 mine old your -f automatondef
```

With `automatondef` being the file with following content:

```
1, (\S*), 2
2, (\s*), 1
```

A different example displays `automaton` definition for simple C files (note that comments are not supported and for `diff3` it is necessary to put the first rule into the capture group):

```
ws, [ \r\n\t]*, wend
wend, (#include [^ \r\n\t]*), ws
wend, ([^ \r\n\t,;(){}+*=&%\-\!|^<>~\[\]\/\\"]+), ws
wend, ([,;(){}+*=&%\-\!|^<>~\[\]\/\\]), ws
wend, ("(?: [^"\\\n] |\\. |\\n)*"), ws
```