# Efficient Regular Pattern Matching avoiding Denial of Service

## Daniel Afonso de Resende

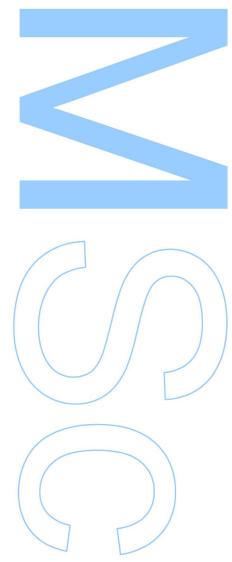
Mestrado em Segurança Informática Departamento de Ciência de Computadores 2025

#### Orientador

Nelma Resende Araújo Moreira, Professora Associada Faculdade de Ciências da Universidade do Porto

#### Orientador

Rogério Ventura Lages dos Santos Reis, Professor Auxiliar Faculdade de Ciências da Universidade do Porto

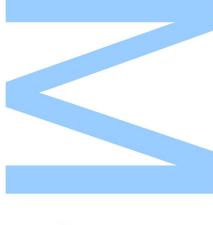


U. PORTO	
FC	FACULDADE DE CIÊNCIAS UNIVERSIDADE DO PORTO

Todas as correções determinadas pelo júri, e só essas, foram efetuadas.

O Presidente do Júri,

Porto, \_\_\_\_/\_\_\_/\_\_\_\_





## **Abstract**

Regular expressions (*regex*) are a foundational tool in modern software, widely used for string matching, input validation, and text parsing. Despite their utility, improperly constructed regular expressions can introduce serious security vulnerabilities. One of the most critical threats is the Regular Expression Denial of Service (*ReDoS*) attack, in which carefully crafted inputs cause the regex engine to perform excessive and redundant processing. This results in dramatic slowdowns or even complete unresponsiveness of the system. ReDoS poses a significant risk to web applications, APIs, and other input-facing systems, where user-controlled input is matched against vulnerable patterns.

In this work, we propose a system to address ReDoS by transforming regular expressions into a modified *position automata*, a form of nondeterministic finite automaton (NFA) that tracks the exact start and end positions of all matches within an input string. This structure enables a matching function that computes *all* match positions, including overlapping ones, without relying on backtracking. By exhaustively and efficiently exploring the automaton's transitions, our approach avoids the exponential blowup typical of vulnerable engines, while preserving a somewhat full regex expressiveness.

Furthermore, we also review and compare this approach with existing solutions present in state-of-the-art programming languages and libraries, such as *RE#*. **Palavras-chave:** regular expressions, ReDoS, position automata, nondeterministic finite automata, pattern matching.

# Resumo

O teu resumo COOL, its me TEST WOWIEESSS

Palavras-chave: palavra, chave..



# Acknowledgements

First of all, I would like to thank my family, etc, etc

Dedico à minha mãe ...

# **Contents**

# **List of Figures**



# Listings



# Acronyms



## Introduction

In this chapter, the problem is overviewed, the study's importance is explained along with goals for the proposed solution.

#### 1.1 Background

Regular expressions are a foundational tool in computer science, widely used in pattern matching, lexical analysis, input validation, and string processing. Their expressiveness and concise syntax make them a powerful language for describing regular languages.

A regular expression R is used (along with an input W) in regex matching engines. The matching engines will verify if W is fully matched by R, meaning that the entire input is a match - or they will verify if a substring of W is matched by R.

## 1.2 Regular Expression Denial of Service

One such vulnerability is known as *Regular Expression Denial of Service* (ReDoS). ReDoS exploits the pathological worst-case behavior of certain regular expressions, causing exponential time complexity during matching. In typical backtracking matchers—such as those found in JavaScript, Java, and many scripting environments—ambiguous or nested expressions (especially involving repetition, such as (a+)+) can lead the engine to explore an exponential number of paths for certain crafted inputs. This behavior allows an attacker to intentionally supply inputs that force excessive computation, effectively rendering a service unavailable or degraded.

The root of the ReDoS problem lies not in regular expressions as a theoretical model but in how they are operationalized in software. While deterministic finite automata (DFAs) evaluate regular expressions in linear time, many real-world engines opt for backtracking NFAs due to their flexibility and ease of implementation. Unfortunately, these NFAs are susceptible to exponential

blow-up in ambiguous or unguarded patterns.

1.3

## **Preliminaries**

Theory builds upon theory, therefore it is essential to establish a solid foundation by understanding the basic concepts and terminology that compose the core topics of formal languages and automata theory. In this chapter we begin by formally defining what a language is and then move on to describe the class of languages known as regular languages. Along the way, we will also introduce various concepts such as finite automata (DFA, NFA) and regular expressions.

#### 2.1 Alphabets, Strings and Languages

#### **Alphabets**

An *alphabet* is a finite, non-empty set of symbols, typically denoted by the Greek letter  $\Sigma$ . That is,

$$\Sigma = \{a_1, a_2, \dots, a_n\}$$

where each  $a_i$  is a symbol in the alphabet.

For example, one can represent the binary alphabet as  $\Sigma = \{0,1\}$ , or the English alphabet as  $\Sigma = \{a,b,c,\ldots,z\}$ .

#### Strings

A *string* over an alphabet  $\Sigma$  is a finite sequence of symbols from  $\Sigma$ . Strings are typically denoted by w, and the *length* of a string w is denoted by |w|.

The set of all strings over the alphabet  $\Sigma$  is denoted by  $\Sigma^*$  and defined as:

 $\Sigma^* = \{ w \mid w \text{ is a finite sequence of symbols from } \Sigma \}$ 

The unique string of length zero is called the *empty string*, denoted by  $\varepsilon$ . It is important to note that  $\varepsilon \in \Sigma^*$ .

For example, if  $\Sigma = \{0, 1\}$ , then we have that:

$$\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, 101, 110, 111, \ldots\}$$

Where the empty string is, as mentioned above, denoted by  $\varepsilon$  and also belongs to  $\Sigma^*$ .

#### Languages

A *language* over an alphabet  $\Sigma$  is a set of strings over  $\Sigma$ .

$$L\subseteq \Sigma^*$$

That is, a language is any subset of  $\Sigma^*$ , possibly infinite, finite, or even empty. Since a language is a set of strings, the following standard set operations can be applied (assuming A and B are languages over the same alphabet  $\Sigma$ ):

- Intersection:  $A \cap B = \{x \mid x \in A \text{ and } x \in B\}$
- Union:  $A \cup B = \{x \mid x \in A \text{ or } x \in B\}$
- Difference:  $A B = \{x \mid x \in A \text{ and } x \notin B\}$

Furthermore, we can also operate specifically over languages with the following operations (assuming  $L_1$  and  $L_2$  are languages over the same alphabet  $\Sigma$ ):

- Concatenation:  $L_1 \cdot L_2 = \{xy \mid x \in L_1 \text{ and } y \in L_2\}$
- Kleene Star:  $L^* = \bigcup_{n=0}^{\infty} L^n$ , where  $L^0 = \{\varepsilon\}$  and  $L^n = L \cdot L^{n-1}$  for n > 0.
- Reversal:  $L^R = \{x^R \mid x \in L\}$ , where  $x^R$  denotes the reversal of string x.
- Complement:  $\overline{L} = \Sigma^* L$ , i.e., the set of all strings over  $\Sigma$  that are not in L.

These operations form the basis for reasoning about the expressiveness and closure properties of language classes such as regular, context-free, and context-sensitive languages. In particular, regular languages are closed under all the operations listed above, including union, intersection, concatenation, and Kleene star. This robustness makes them especially amenable to algorithmic manipulation, as seen in finite automata and regular expression engines.

2.2. Finite Automata 5

#### 2.2 Finite Automata

A *finite automaton* is a theoretical machine used to recognize regular languages. It processes input strings symbol by symbol and determines whether the string belongs to the language defined by the automaton. There are two main types of finite automata:

- **Deterministic Finite Automaton (DFA)**: An automaton where, for each state and input symbol, there is exactly one possible next state.
- Nondeterministic Finite Automaton (NFA): An automaton that allows multiple possible transitions for a given state and input symbol, including transitions without consuming any input (called  $\varepsilon$ -transitions).

Formally defined, an NFA is a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$  where:

- Q is a finite set of states,
- $\Sigma$  is the input alphabet,
- $\delta: Q \times (\Sigma \cup \{\varepsilon\}) \to 2^Q$  is the transition function,
- $q_0 \in Q$  is the initial state,
- $F \subseteq Q$  is the set of accepting (final) states.

A string  $w \in \Sigma^*$  is accepted by the NFA if there exists a sequence of transitions (possibly including  $\varepsilon$ -moves) that consumes w and ends in a state q such that  $q \in F$ .

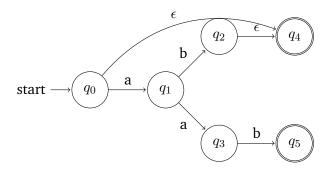


Figure 2.1: An example of a NFA that accepts  $L = \{\epsilon, ab, aab\}$ 

A NFA is deterministic (also known as a DFA) if  $|\delta(q,\sigma)| \leq 1$ , for any  $(q,\sigma) \in Q \times \Sigma$ .

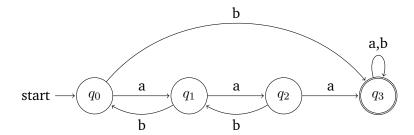


Figure 2.2: An example of a DFA that accepts any strings that start with the symbol 'b' or contain the substring aaa

#### 2.3 Regular Expressions

Let  $\Sigma$  be a finite alphabet. Let  $L \subseteq \Sigma$ . The set of *regular expressions* over  $\Sigma$ , denoted by RegExp( $\Sigma$ ), is defined inductively as follows:

- $\emptyset$  is a regular expression denoting the empty language:  $L(\emptyset) = \emptyset$ .
- $\varepsilon$  is a regular expression denoting the language containing only the empty string:  $L(\varepsilon)=\{\varepsilon\}.$
- For each symbol  $a \in \Sigma$ , a is a regular expression denoting the singleton language:  $L(a) = \{a\}$ .
- If  $r_1$  and  $r_2$  are regular expressions, then so are:
  - $(r_1 | r_2)$ , denoting the union:  $L(r_1 | r_2) = L(r_1) \cup L(r_2)$ .
  - $(r_1 \cdot r_2)$ , denoting concatenation:  $L(r_1 \cdot r_2) = L(r_1) \cdot L(r_2)$ .
  - $(r_1)^*$ , denoting Kleene star:  $L(r_1^*) = (L(r_1))^*$ .

We write  $\operatorname{RegExp}(\Sigma)$  to denote the set of all such syntactic expressions, and for each  $r \in \operatorname{RegExp}(\Sigma)$ , the function L(r) yields the language defined by r.

Parentheses are used to disambiguate expressions and enforce precedence; by convention, Kleene star binds most tightly, followed by concatenation, and finally union.

#### 2.3.1 Extended Regular Expressions

In addition to the basic operations, some extended operators are often used for convenience. These include:

• Kleene plus: Given a regular expression r, the expression  $r^+$  denotes one or more repetitions of r:

$$L(r^+) = L(r) \cdot L(r)^*.$$

• Fixed repetition (power): For a regular expression r and integer  $n \ge 0$ , the expression  $r^n$  denotes n consecutive concatenations of r:

$$L(r^0) = \{\varepsilon\}, \quad L(r^n) = L(r) \cdot L(r^{n-1}) \text{ for } n > 0.$$

• Bounded repetition: For a regular expression r and integers m, n with  $0 \le m < n$ , the bounded repetition  $r^{[m,n[}$  denotes the language containing all strings formed by concatenating between m and n copies of strings from L(r):

$$L(r^{[m,n[}) = \bigcup_{k=m}^{n} L(r^k).$$

These extended forms do not increase the expressive power of regular expressions but are useful for readability and practical applications. They can always be rewritten using the fundamental operators: union and concatenation.

#### 2.3.2 Derivatives

The derivative of a regular expression was first introduced in 1962 by Janusz Brzozowski. It is a powerful concept used to define the behavior of regular expressions in a more operational manner. They can be used as a means of verifying equivalence of regular expressions, for example. The derivative of a regular expression r with respect to a symbol a is another regular expression  $D_a(r)$  that describes the set of strings that can be obtained by taking the derivative of r with respect to a.

The derivative can be defined based on the structure of the regular expression:

- If  $r = \varepsilon$ , then  $D_a(r) = \emptyset$  for all  $a \in \Sigma$ .
- If r = a, then  $D_a(r) = \varepsilon$ .
- If  $r = r_1 \cdot r_2$ , then  $D_a(r) = D_a(r_1) \cdot L(r_2) \cup \varepsilon \cdot D_a(r_2)$ .
- If  $r = r_1 + r_2$ , then  $D_a(r) = D_a(r_1) + D_a(r_2)$ .
- If  $r = r^*$ , then  $D_a(r) = D_a(r) \cdot r^*$ .

Given  $r \in RE$  - associativity, commutativity and idempotence also apply:

- $r\emptyset = \emptyset r = \emptyset$
- $r\varepsilon = \varepsilon r = r$
- $\emptyset + r = r + \emptyset$

## 2.4 From Regular Expressions to Automata

While regular expressions provide a declarative way to specify patterns in strings, finite automata offer an operational model for recognizing such patterns.

#### 2.4.1 Thompson's Algorithm

Thompson's construction is a classic algorithm for converting a regular expression into a nondeterministic finite automaton (NFA) with  $\varepsilon$ -transitions. It was introduced in 1968 by Ken Thompson and is the foundation for many regex engines, including the lex lexical analyzer generator.

The construction proceeds recursively based on the structure of the regular expression. Each base case (such as a single symbol or the empty word  $\varepsilon$ ) and each operator ('+', concatenation, '\*') corresponds to a small NFA fragment. These fragments are then joined using  $\varepsilon$ -transitions.

Although Thompson's NFA may contain many  $\varepsilon$ -transitions, it is guaranteed to be of size linear in the length of the regular expression. The resulting NFA can be converted into a deterministic finite automaton (DFA) using the standard powerset construction, typically after removing  $\varepsilon$ -transitions.

#### 2.4.2 Brzozowski's Derivatives

Through the definition of *derivatives* mentioned in ??, one can compute the derivatives for all symbols, memoize the results and build a DFA that represents the same language as the intended regular expression.

This method can, however, lead to an exponential number of distinct derivatives in the worst case. Therefore, simplification rules and expression equivalences are critical to making the approach practical.

#### 2.4.3 Antimirov's Partial Dertivatives

Proposed by Valery Antimirov in 1996, the partial derivatives construction generalizes Brzo-zowski's derivatives to build an NFA rather than a DFA. Instead of producing a single derivative for each symbol, Antimirov's method produces a *set* of partial derivatives, reflecting the inherent nondeterminism of the regular expression.

This construction avoids  $\varepsilon$ -transitions and yields a compact  $\varepsilon$ -free NFA. Each partial derivative corresponds to a transition in the automaton, and the process naturally handles alternation and repetition.

Antimirov's approach is especially efficient for regex evaluation and analysis tasks, and forms the basis for several modern regex matchers and formal verification tools.

#### 2.4.4 Position Automata

The *position automaton*, also known as the *Glushkov automaton*, is a type of  $\varepsilon$ -free nondeterministic finite automaton (NFA) constructed directly from a regular expression. Unlike the standard Thompson construction, which introduces  $\varepsilon$ -transitions that must later be eliminated, the Glushkov construction yields an automaton in which each state corresponds uniquely to a symbol occurrence—or *position*—in the expression. [mesh-of-automata]

Given a regular expression E, the Glushkov automaton  $M_E$  is defined based on three key position-based functions:

- first(E): the set of positions that can appear first in some word of the language  $\mathcal{L}(E)$ .
- last(E): the set of positions that can appear last in some word of  $\mathcal{L}(E)$ .
- follow(E, x): for each position x, the set of positions that can immediately follow x in some word of  $\mathcal{L}(E)$ .

To distinguish different occurrences of the same symbol, the construction introduces marked symbols. For example, the expression  $(a+b)^*a(b+a)^*$  is rewritten as  $(a_1+b_2)^*a_3(b_4+a_5)^*$ . Each marked symbol corresponds to a unique position and becomes a distinct state in the automaton.

The Glushkov automaton  $M_E = (Q, \Sigma, \delta, q_0, F)$  is constructed as follows:

- Q is the set of positions in E (i.e., the marked symbols), plus an initial state  $q_0$ .
- For each symbol  $a \in \Sigma$ :

```
- \delta(q_0, a) = \{x \in \text{first}(E) \mid \text{symbol}(x) = a\}

- \delta(x, a) = \{y \in \text{follow}(E, x) \mid \text{symbol}(y) = a\}
```

• The set of final states F is last(E); if  $\varepsilon \in \mathcal{L}(E)$ , then  $q_0$  is also final.

This automaton captures the structural flow of E by tracing symbol sequences as state transitions. It is  $\varepsilon$ -free and has one state per symbol occurrence, which results in at most a quadratic number of transitions with respect to the size of E.

An important property of the Glushkov automaton is its relationship to unambiguity. A regular expression is *weakly unambiguous* if and only if its Glushkov automaton is unambiguous, i.e., there exists at most one accepting path for each accepted word. This makes the Glushkov construction a practical and efficient tool in applications requiring unambiguous parsing, such as syntax analysis in document type definitions (e.g., SGML and XML DTDs).

#### 2.5 FAdo

The FAdo [fado\_paper] project is an open-source implementation of several sets of tools for formal languages manipulation. In order to allow quick prototyping and testing of algorithms, these tools were developed in Python. Regular languages can be represented by regular expressions which are defined in reex.py - or by finite automata which are defined in fa.py.

#### 2.5.1 Regular Expressions

In reex.py, FAdo defines the RegExp base class for all regular expressions, declaring therein a class variable sigma, formally known as the set of symbols ( $\Sigma$ ). To represent any and all base constructions of a regular expressions, FAdo defines base classes for each of them:

- **CEmptySet**: The empty symbol set. (∅)
- **CEpsilon**: The empty string. ( $\varsigma$ )
- **CAtom**: A simple symbol. (e.g. 'a')
- CDisj: The + operation between symbols. (e.g. CDisj(CAtom(a), CAtom(b)) represents the regular expression R = a + b where  $a, b \in \Sigma_R$ )
- **CConcat**: The · operation
- CStar: The Kleene closure over a set of symbols (e.g. CStar (CDisj (CAtom(a), CAtom(b))) representes the regular expression  $R=(a+b)^*$  where  $a,b\in\Sigma_R$ )

In order to parse expressions into FAdo's classes and types, *lark* was used. *lark* is a parsing toolkit for Python. It can parse all context-free languages.

#### 2.5.2 Finite Automata

#### 2.5.3 Extended Regular Expressions

# State of the Art

The current state-of-the-art for mitigating the **ReDoS!** What!

# Counting

In this chapter, we will explore the concept of counting in the context of formal languages and automata theory as well as explain an attempt that was made towards match counting using a partial derivative automaton construction and why it didn't work.

#### 4.1 Counting in Formal Languages

In formal language theory, counting refers to the ability of a language or an automaton to enforce numeric constraints over the number of symbols or patterns within strings. Specifically, it deals with the ability to recognize whether certain elements occur a specified number of times—or in a specific numerical relationship to others.

4.2

# Matching

Regex *matching* is the process of checking whether a piece of text fits a specific pattern described using a regular expression (regex). A regex is a compact, rule-based way to describe sets of strings—like email addresses, phone numbers, or specific word formats.

In this chapter, we will discuss the different approaches to matching regular expressions, the implications of using them, and the performance considerations that arise from these choices. Furthermore, we will also present a novel approach based on position automata, which aims to mitigate the performance issues associated with traditional regex engines while preserving some of the extended expressiveness of regex patterns.

#### 5.1 Modified Position Automata

A *position automaton* is a type of nondeterministic finite automaton (NFA) that facilitates overlapped matching.

#### **Algorithm 1** NFAPOSCOUNT(R): Construct Special Position Automaton

```
Require: Regular expression R
Ensure: A NFA A
 1: A \leftarrow \text{new empty NFA}
 2: i \leftarrow A.addInitialState()
 3: A.addTransitionStar(i, i)
                                                                                          \triangleright Accept any symbol from \Sigma
 4: f_R \leftarrow R.\text{marked}()
 5: stack \leftarrow empty stack
 6: addedStates ← empty map
 7: for all p \in First(f_R) do
        q \leftarrow A.addState(p)
 9:
        addedStates[p] \leftarrow q
        stack.push((p,q))
10:
        A.addTransition(i, p, q)
12: end for
13: while stack is not empty do
        (s, s_{idx}) \leftarrow stack.pop()
14:
        for all t \in Follow(f_R, s) do
15:
             if t \in addedStates then
16:
                 q \leftarrow \mathsf{addedStates}[t]
17:
18:
             else
19:
                 q \leftarrow A.addState(t)
                 addedStates[t] \leftarrow q
20:
                 stack.push((t, q))
21:
22:
             end if
23:
             A.addTransition(s_{idx}, t, q)
        end for
24:
25: end while
26: e \leftarrow A.addState()
27: A.addTransitionStar(e, e)
28: for all p \in f_R.Last() do
29:
        if p \in addedStates then
             A.addFinal(addedStates[p])
30:
             A.addTransitionStar(addedStates[p], e)
31:
32:
        end if
33: end for
```

### 5.2 Automata-Based Matching

Given a regular expression R, one can construct an NFA A such that L(A) = L(R). Matching then reduces to verifying whether the automaton A accepts the input string s. In DFA-based engines, each character of the input leads to a deterministic transition from one state to another, resulting in a guaranteed linear-time match. In contrast, NFA-based engines may involve branching paths due to nondeterminism and can require simulating multiple transitions concurrently.

For example, consider the following regex pattern:

## 5.3 Matching with the Modified Position Automaton

One can find all matches over an input string by constructing the modified position automaton from a regular expression and then simulating the automaton's transitions over the input string. The algorithm presented in this section is designed to track the start and end positions of all matches, including overlapping ones, without relying on backtracking.

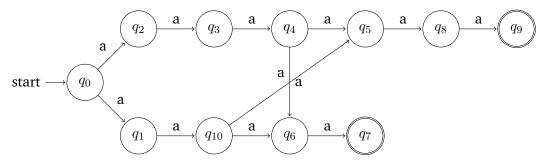
#### **Algorithm 2** TABLEMATCHER(A, s): Modified Position Automaton Multi-matcher

```
Require: A = (\Sigma, Q, \delta, I, F): NFA
Require: s: input string
Ensure: M: mapping from final states to lists of match positions
 1: symbols \leftarrow list with \varepsilon prepended to s
 2: currentRow \leftarrow empty map from states to list of position pairs
 3: finalMatches \leftarrow empty map from states to list of matches
 4: position \leftarrow 0
 5: for all sym in symbols do
        if sym = \varepsilon then
 6:
 7:
            for all q_0 \in I do
                currentRow[q_0] \leftarrow [(0,0)]
 8:
 9:
            end for
10:
        else
            nextRow \leftarrow empty map
11:
            if sym \in \Sigma then
12:
                for all q \in \text{keys}(currentRow) do
13:
                    if |\delta(q, sym)| > 0 then
14:
                        for all q' \in \delta(q, sym) do
15:
16:
                            for all (start, ) \in currentRow[q] do
                                if q' = q and q' \in I then
17:
18:
                                    append (position, position) to nextRow[q']
                                else
19:
                                    append (start, position) to nextRow[q']
20:
                                    if q' \in F then
21:
22:
                                        append (start, position) to finalMatches[q']
                                    end if
23:
                                end if
24:
25:
                            end for
                        end for
26:
27:
                    end if
28:
                end for
                                                                             \triangleright Symbol not in \Sigma; treat as fresh start
29:
            else
30:
                for all q_0 \in I do
                    nextRow[q_0] \leftarrow [(position, position)]
31:
32:
                end for
            end if
33:
34:
            currentRow \leftarrow nextRow
35:
            position \leftarrow position + 1
36:
        end if
37: end for
38: return \ final Matches
```

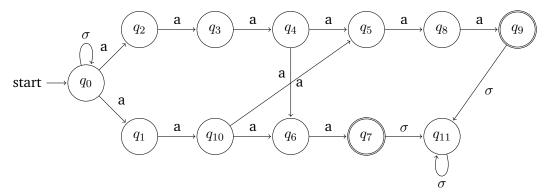
As an example, consider the following regular expression R and input string w:

$$R = (aa + aaa)(aaa + aa)$$
$$w = aaaaabaaaaa$$

The Glushkov automaton construction for R is as follows:



Meanwhile, the modified position automaton for this regular expression can be constructed using the algorithm presented in ??, resulting in the following:



One can then apply the matching algorithm (??) to an input string, resulting in the following match table:

# **Results and Discussion**

This is a test

- 6.1 Algorithm Analysis
- 6.2 Accuracy

The methods of evaluating

- 6.3 Comparison with Other Methods
- 6.4 Examples

# **Conclusion**

## 7.1 Findings Summary

This research and development project served the objective of

## 7.2 Contributions

7.2. Contributions 25

refs