# Efficient Regular Pattern Matching avoiding Denial of Service
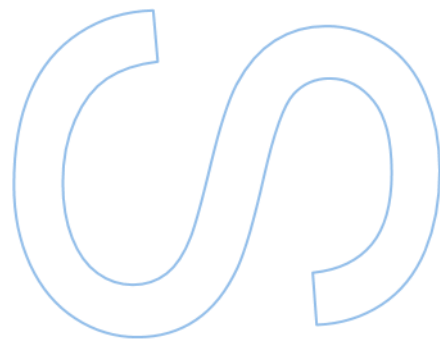
Daniel Afonso de Resende

Master in Information Security
Department of Computer Science
Faculty of Sciences of the University of Porto
2025

# Efficient Regular Pattern Matching avoiding Denial of Service
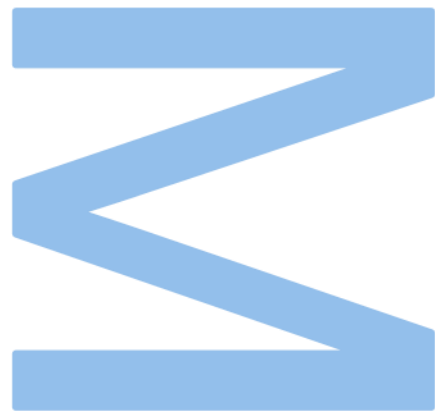
Daniel Afonso de Resende
Dissertation carried out as part of the Master in Information Security
Department of Computer Science
2025

Supervisor
Rogério Ventura Lages dos Santos Reis, Associate Professor, Faculty of Sciences of the University of Porto

Co-supervisor
Nelma Resende Araújo Moreira, Associate Professor, Faculty of Sciences of the University of Porto

Dedicado aos meus pais

# Acknowledgements

First and foremost, I want to thank professors Nelma Moreira and Rogério Reis, who have assisted me this last year. Their help, guidance, willingness and friendliness were fundamental for this work and for my development. Words will never be enough to describe what I feel for my friends and family, who have helped me carry my woes and solaces along this journey. I am especially grateful to my mom and dad for taking care of me. Lastly, I would also like to thank my work colleagues, who have been very thoughtful and patient with me along the way.

# Abstract

Regular expressions (*regex*) formalize a class of patterns definable over strings, corresponding to the family of regular languages in automata theory. They provide a declarative mechanism for specifying sets of strings, making them central both to theoretical models of computation and to practical applications such as string matching, input validation, and text parsing. Despite their utility, certain implementations of regular expression engines (particularly those based on backtracking) can introduce serious security vulnerabilities when combined with specific patterns. 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, bounded quantifiers (counting) were implemented in *FAdo*. Building upon this, we also propose a system to address ReDoS by transforming regular expressions into a modified position automaton, a Non-deterministic 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.

We also review existing state-of-the-art solutions implemented in modern programming languages and libraries, such as *RE#* [20] and *Hyperscan* [22].

**Keywords:** regular expressions, ReDoS, position automata, nondeterministic finite automata, pattern matching.

# Resumo

As expressões regulares (*regex*) formalizam uma classe de padrões definíveis sobre cadeias de caracteres, correspondendo à família das linguagens regulares na teoria dos autómatos. Fornecem um mecanismo declarativo para especificar conjuntos de cadeias, sendo, por isso, centrais tanto para modelos teóricos de computação como para aplicações práticas, tais como procura de padrões, validação de dados e análise de texto. Apesar da sua utilidade, certas implementações de motores de expressões regulares (particularmente aqueles baseados em *backtracking*) podem introduzir vulnerabilidades de segurança graves quando utilizam padrões específicos. Uma das ameaças mais críticas é o ataque Regular Expression Denial of Service (ReDoS), onde certas entradas levam o motor de regex a realizar processamento excessivo e redundante. Isto resulta em lentidões significativas ou mesmo na completa inoperacionalidade do sistema. O ReDoS representa um risco significativo para aplicações web, APIs e outros sistemas voltados para entrada de dados, onde a entrada controlada pelo utilizador é comparada com padrões vulneráveis.

Neste trabalho, foram implementados quantificadores limitados (contagem) no *FAdo*. Com isto, propomos também um sistema para mitigar o ReDoS, ao transformar expressões regulares num autómato de posições modificado — um Non-deterministic Finite Automaton (NFA) que regista exatamente as posições de início e fim de todas as ocorrências dentro de uma cadeia de entrada. Esta estrutura permite que uma função de *matching* calcule todas as posições de ocorrência, incluindo as sobrepostas, sem recorrer a *backtracking*. Ao explorar de forma exaustiva e eficiente as transições do autómato, a nossa abordagem evita a explosão exponencial típica de motores vulneráveis, mantendo, ao mesmo tempo, uma expressividade relativamente completa das expressões regulares.

Analisamos também soluções existentes em linguagens e bibliotecas de programação de ponta, como o *RE#* [20] e o *Hyperscan* [22].

**Palavras-chave:** expressões regulares, ReDoS, autómatos de posições, autómatos finitos não deterministas, correspondência de padrões.

# Contents

# List of Tables

# List of Figures

# Listings

# Acronyms

**DFA**    Deterministic Finite Automaton

**FCUP**   Faculdade de Ciências da
         Universidade do Porto

**NFA**     Non-deterministic Finite Automaton

**ReDoS** Regular Expression Denial of Service

**Regex**   Regular Expression

# Chapter 1

# Introduction

Regular expressions are a powerful and ubiquitous tool in software development, enabling concise and expressive definitions of complex text patterns. However, their utility comes with a hidden cost: when implemented naively or misused, certain regular expression constructs can expose systems to catastrophic performance degradation. This vulnerability, known as Regular Expression Denial of Service (ReDoS), arises particularly in backtracking-based matching engines and has been responsible for real-world outages across major platforms.

This chapter introduces the ReDoS problem by first providing foundational background on regular expression engines and their operational behavior. It then outlines the nature of ReDoS vulnerabilities, illustrating their practical impact through two real-world case studies. Finally, it discusses the reasons for continued reliance on legacy engines despite known risks. By establishing the technical and practical motivations behind this research, this chapter sets the stage for the proposed approach introduced in the subsequent chapters.

## 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 running

time complexity during matching.

In commonly used backtracking matchers (such as those found in JavaScript, Java, and many scripting environments) ambiguous or nested expressions (such as '(a+)+', which involves repetition) 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 implemented in many real-world software systems—particularly through backtracking-based matchers. While deterministic finite automata (DFAs) evaluate regular expressions in linear time, backtracking matchers explore multiple paths through the input, making them susceptible to exponential blow-up in the presence of ambiguous or nested patterns.

## 1.3   Case Studies

In this section, two case studies are presented. These serve as a form of introduction to getting to know and understand the ReDoS problem.

### 1.3.1   Stack Overflow

On the 20th of July 2016, a user published a malformed post on the online information exchange forum *Stack Overflow.* A couple minutes after the post, at around 14:44 UTC, the entire website became unavailable for around 34 minutes, after which there was an update that fixed the underlying issue.

On the forum, there is an automatic text formatter that runs every time someone posts something. This tool will trim any group of leading whitespace or invisible characters at both the beginning and end of a post. The regular expression that does so is the following:

$$\hat{\ }[\backslash s\backslash u200c]+|[\backslash s\backslash u200c]+\$$$

- ^ will anchor the following expression to the start of a string

- $[\backslash s\backslash u200c]+$ will match one or more of either whitespace characters (space, tab, newline, etc..) or the unicode characters U+200C (zero-width non-joiner)

- $ will anchor the match to the end of that string

The aforementioned tool was an automatic text formatter which contained a matcher that tried to match the regular expression described above against an input that contained around 20,000 consecutive whitespace characters on one line that started with $--$. The backtracking matcher in place worked as follows:

Given an input string $M$ of length $\text{len}(M)$, let $n_k$ denote the character at position $k$, where $0 \leq k < \text{len}(M)$. For each possible starting position $p$ such that $0 \leq p < \text{len}(M)$, perform the following steps:

1. Check whether the character $n_k$ (for $k \geq p$) is either a whitespace or a zero-width non-joiner (U+200C).

2. Continue checking characters until it reaches a character that is neither a whitespace nor a zero-width non-joiner.

3. If the end of the string is reached and all characters from position $p$ onward matched, the pattern succeeds.

4. If a non-matching character is encountered before the end of the string, the match fails at this position; increment $p$ and repeat from step 1.

For a 20,000 whitespace-character (both whitespace and zero-width non-joiner characters) input, the sum of computations is given as follows:

$$\sum_{k=1}^{20,000} k = \frac{20,000 \cdot (20,000) + 1}{2} = 200,010,000$$

This means that the matching algorithm ran in $O(n^2)$ complexity, and this blow up was to be expected.

The engineers quickly fixed this issue by switching to a substring replacing method.

### 1.3.2   minimatch

**minimatch** is a minimal matching utility, used internally by the **Node Package Manager**, better known as **npm**. [23] This utility works by converting glob expressions into JavaScript's *RegExp* objects, supporting the following glob features:

- Brace Expansion

- Extended glob matching

- "Globstar" (**) matching

- POSIX character classes

The utility has millions of downloads, as it is an essential component of **npm**.

On the 18th of October 2022, a new CVE was introduced: **CVE-2022-3517**. The report showed that all of minimatch's versions below 3.0.5 were vulnerable to a ReDoS attack similar to the one described in Subsection 1.3.1. The culprit for this was a function called *braceExpand*, which is responsible for expanding brace patterns in glob strings. This is commonly known as brace expansion and is often seen in Unix shells (such as bash). The function contained a regular expression that would match against given patterns and decide if a brace expansion was in order. The expression used was "/\{.*\}/", which matches any string containing a single pair of curly braces with any characters inside. But this expression poses an issue:

For example, the following text:

$$\texttt{"\{\{\{\{\{\{\{\{\{\{\{\{\{\{\{...X"}}$$

with the opening brace { repeated over 30,000 times and no closing brace }, can cause a significant CPU spike or even hang the process due to catastrophic backtracking.

To fix this issue, the developer decided to switch to a safer regular expression: \{(?:(?!\{).)*\}/

This new pattern prevents the exponential blow-up because it removes the ambiguity introduced by the symbol ., which is the greedy quantifier. In the original expression, the greedy quantifier could match arbitrarily long substrings and repeatedly backtrack in search of a closing brace symbol, resulting in exponential-time behavior when no match was found.

In contrast, the new expression uses a *negative lookahead*, which ensures that each symbol is consumed only if it is not followed by another opening brace, thereby eliminating redundant backtracking.

## 1.4   Reluctance Toward Changing Legacy Matchers

Despite well-documented vulnerabilities such as ReDoS, there remains significant reluctance in the software engineering community to replace or refactor legacy regex engines—particularly those built into performance-critical or widely adopted tools such as grep, sed, and many scripting languages.

These tools often rely on matching engines that prioritize speed and simplicity of implementation over safety. For example, grep and similar UNIX utilities implement regex matchers using finite automata, but their behavior with extended features (like bounded repetitions) can still lead to performance degradation in edge cases. While these engines are generally immune to the exponential blow-up typical of backtracking matchers, they may suffer from linear but high-cost processing when automata grow excessively large due to poorly constructed patterns.

The situation is more severe in environments that rely on backtracking matchers, such as

JavaScript, Java, and many shell-based text processors. In these ecosystems, regular expressions are both expressive and dangerously permissive, allowing patterns that trigger catastrophic backtracking without warning.

Refactoring or replacing these engines is often resisted for several reasons:

- **Backwards compatibility**: Legacy codebases and systems expect specific regex semantics, and changing the underlying engine could break existing behavior.

- **Perceived performance cost**: DFA-based matchers may require significant memory or preprocessing, which is viewed as a performance risk in lightweight tools.

- **Lack of awareness**: Many developers are unaware that regex matchers can introduce denial-of-service vulnerabilities, especially when ReDoS exploits are subtle or input-driven.

- **Cultural inertia**: Tools like `grep` are deeply embedded in developer workflows and scripts, making any modification to their behavior or performance profile controversial.

Even modern matchers that are designed to avoid ReDoS—such as Google's RE2 or Rust's regex crate—are often underutilized due to these legacy constraints.

These factors highlight the need for not only technical solutions, such as safer regex engines and static analysis tools, but also a cultural shift in how regular expressions are authored, reviewed, and validated in production systems.

The next chapter will give some more insight into the basics of regular expressions and automata.

# Chapter 2

# 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, Words, and Languages

### Alphabets

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

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

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

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

### Words

A *word* over an alphabet $\Sigma$ is a finite sequence of symbols from $\Sigma$. Words are typically denoted by $w$, and the *length* of a word $w$ is denoted by $|w|$. The unique word of length zero is called the *empty word*, denoted by $\varepsilon$. The concatenation of two words $w$ and $v$ over $\Sigma$ is denoted by $wv$ (or, more explicitly, $w \cdot v$) and results in a separate word composed by the singular word $v$ right after $w$. For instance, let $w = ab$ and $v = cd$, the concatenation $wv$ is $abcd$. We can define the

concatenation of two words $w$ and $v$ (which are composed of any amount of symbols) as follows:

$$w \cdot \varepsilon = \varepsilon \cdot w = w,$$

$$(wa) \cdot v = w \cdot (av), a \in \Sigma.$$

The set of all words over $\Sigma$ is denoted by $\Sigma^\star$ and is defined as:

$$\varepsilon \in \Sigma^\star,$$

$$\forall_{a \in \Sigma}, \forall_{w \in \Sigma^\star} \, aw \in \Sigma^\star$$

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 word is, as mentioned above, denoted by $\varepsilon$ and also belongs to $\Sigma^*$.

Given a word $w = a_1 a_2 \cdots a_n$, the reversal of $w$, denoted by $w^R$, is the word formed by reversing the order of its characters such that $w^R = a_n a_{n-1} \cdots a_1$. We can define the reversal of a word recursively like so:

$$w^R = v^R a, \quad \text{where } w = av, a \in \Sigma \text{ and } v \in \Sigma^*,$$

$$\varepsilon^R = \varepsilon.$$

## Languages

A *language* over an alphabet $\Sigma$ is a set of words over $\Sigma^\star$,

$$L \subseteq \Sigma^\star$$

That is, a language is any subset of $\Sigma^*$, and can be either finite or infinite.
Since a language is a set of words, 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*: $A \setminus B = \{x \mid x \in A \text{ and } x \notin B\}$

- *Complement*: $\overline{A} = \{x \mid x \in \Sigma^\star \text{ and } x \notin A\}$

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*: $L^\star = \bigcup_{n=0}^{\infty} L^n$, 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 a word $x$.

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 this instance, the class of regular languages over an alphabet $\Sigma$ includes the finite languages and is closed under *union* ($\bigcup$), *concatenation* ($\cdot$) 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

A *finite automaton* is a model of computation used to recognize regular languages. It processes input words symbol by symbol and determines whether the word 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.

- **Non-deterministic 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,

- $Q_0 \subseteq Q$ is the set of initial states,

- $F \subseteq Q$ is the set of accepting (final) states.

A word $w \in \Sigma^*$ 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 contains an example of a NFA.
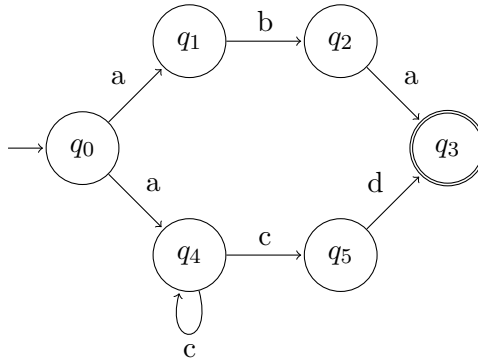
Figure 2.1: Example of an NFA that accepts $L = \{aba, ac^\star d\}$

An NFA is *deterministic* (also known as DFA) if $|\delta(q, \sigma)| \leq 1$, for any $(q, \sigma) \in Q \times \Sigma$ and $|Q_0| = 1$. Figure 2.2 contains an example of a DFA.
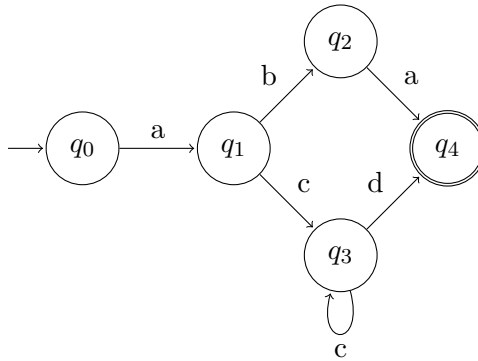


Figure 2.2: A DFA whose language is the same as the NFA from Figure 2.1

For a DFA, membership testing of a word $w$ in a language $L \subseteq \Sigma^\star$ can be decided in time $O(|w|)$.

## 2.3   Regular Expressions

Let $\Sigma$ be a finite alphabet. Let $L \subseteq \Sigma$. The set of *regular expressions* over $\Sigma$, denoted by $\mathrm{RegExp}(\Sigma)$, 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 word: $L(\varepsilon) = \{\varepsilon\}$.

- 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:

- $(r_1 \mid r_2)$ or $(r_1 + r_2)$, denoting the union: $L(r_1 \mid r_2) = L(r_1) \cup L(r_2)$.
- $(r_1 \cdot r_2)$, denoting concatenation: $L(r_1 \cdot r_2) = L(r_1) \cdot L(r_2)$.
- $(r_1)^*$, denoting Kleene star: $L(r_1^*) = (L(r_1))^*$.

For each $r \in \text{RegExp}(\Sigma)$, 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 $(a^*)$, followed by concatenation (e.g. $a \cdot b$, whose operator "$\cdot$" is omitted for convenience), and finally union $(+)$.

We denote by $\alpha = \beta$ if two regular expressions $\alpha$ and $\beta$, both over $\Sigma$, represent the same language $(L(\alpha) = L(\beta))$.

If the language of a regular expression $\alpha$ contains the empty word, then that regular expression possesses the *empty word property*.

**Definition 2.3.1** (Empty Word Property)**.** The empty word property is characterized by the function $\varepsilon : \text{RegExp} \rightarrow \{\varepsilon, \emptyset\}$ and is defined recursively as follows, given $\alpha$ and $\beta$ regular expressions defined over $\Sigma$:

$$\varepsilon(\emptyset) = \emptyset$$
$$\varepsilon(\varepsilon) = \varepsilon$$
$$\varepsilon(a) = \emptyset, \ a \in \Sigma$$
$$\varepsilon(\alpha + \beta) = \begin{cases} \emptyset & \text{if } \varepsilon(\alpha) = \varepsilon(\beta) = \emptyset \\ \varepsilon & \text{if } \varepsilon(\alpha) = \varepsilon \text{ or } \varepsilon(\beta) = \varepsilon \end{cases}$$
$$\varepsilon(\alpha\beta) = \begin{cases} \emptyset & \text{if } \varepsilon(\alpha) = \emptyset \text{ or } \varepsilon(\beta) = \emptyset \\ \varepsilon & \text{if } \varepsilon(\alpha) = \varepsilon(\beta) = \varepsilon \end{cases}$$
$$\varepsilon(a^*) = \varepsilon$$

To effectively compare regular expressions and discuss their equivalence in automata construction, it is useful to define a notion of similarity based on basic algebraic properties. The following definition, adapted from Brzozowski's work in [4], formalizes this concept:

**Definition 2.3.2** ([4])**.** Two regular expressions are similar if one can be transformed into the other using only the following rules:

$$r + r = r,$$
$$p + q = q + p,$$
$$(p + q) + r = p + (q + r)$$
$$\varepsilon r = r\varepsilon = r$$
$$\emptyset r = r\emptyset = \emptyset$$

Two regular expressions are dissimilar if they are not similar.

### 2.3.1   Extended Regular Expressions

In addition to the basic operations, some other 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 $n \geq 0$, 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 $0 \leq m < n$, the bounded repetition $r^{[m,n]}$ denotes the language containing all words formed by concatenating between $m$ and $n$ copies of words from $L(r)$:
$$L(r^{[m,n]}) = \bigcup_{k=m}^{n-1} L(r^k).$$

These 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. Derivatives can be used to decide whether a word belongs to $L(r)$, for a given regular expression $r$; to check whether two regular expressions are equivalent; or even to construct a DFA equivalent to the expression itself.

**Definition 2.3.3** (Derivative of a Regular Expression [4])**.** The derivative of a regular expression $r$ with respect to $\sigma \in \Sigma$ is itself a regular expression $d_\sigma(r)$ such that $L(d_\sigma(r)) = \{w \mid \sigma w \in L(r)\}$

and is defined as:

$$d_\sigma(\emptyset) = \emptyset,$$
$$d_\sigma(\varepsilon) = \emptyset,$$
$$d_\sigma(r') = \begin{cases} \varepsilon & \text{if } \sigma' = \sigma \\ \emptyset & \text{otherwise,} \end{cases}$$
$$d_\sigma(r + r') = d_\sigma(r) + d_\sigma(r'),$$
$$d_\sigma(rr') = \begin{cases} d_\sigma(r)r' & \text{if } \varepsilon(r) = \emptyset \\ d_\sigma(r)r' + d_\sigma(r') & \text{otherwise,} \end{cases}$$
$$d_\sigma(r^*) = d_\sigma(r)r^*,$$
$$d_\sigma(r^+) = d_\sigma(r)r^*$$

This definition can be extended to a word $w \in \Sigma^*$ like so:

$$d_\varepsilon(r) = r,$$
$$d_{wa}(r) = d_a(d_w(r)), \quad a \in \Sigma$$

With this extension:

$$L(d_x(r) = \{w \mid xw \in L(r)\}), \quad x \in \Sigma^*$$

Originally, Brzozowski defined the derivative for a concatenation of two regular expressions $r$ and $r'$:

$$d_\sigma(rr') = d_\sigma(r)r' + \varepsilon(r)d_\sigma(r'),$$

This definition will lead to an infinite number of derivatives for a regular expression. For instance, let $r = a^*$ and $x \in \Sigma^*$:

If $x = a$:

$$d_x(r) = d_a(a^*)$$
$$= \varepsilon a^*$$

If $x = aa$:

$$d_x(r) = d_{aa}(a^*)$$
$$= \varepsilon\varepsilon a^*$$

As seen, each of these derivatives is syntactically different and dissimilar. This pattern repeats, meaning that the derivative set can be infinite. The derivatives related to the concatenation

of two different regular expressions, defined in Definition 2.3.3, were proposed by Antimirov [2] and they were introduced as a way to address this issue. This modification avoids explicit multiplication by multiple $\varepsilon$ terms, ensuring that only a finite number of dissimilar derivatives are generated. Consequently, this modified definition guarantees the termination of derivative-based automaton constructions.

### 2.3.3 Partial Derivatives

The notion of *partial derivative* was introduced by Antimirov [2]. Contrary to Brzozowski's derivatives, *partial derivatives* will lead to the construction of an NFA.

Before defining partial derivatives, we extend the notion of concatenation to sets of regular expressions. Let $\cdot : 2^{RegExp(\Sigma)} \times RegExp(\Sigma) \to 2^{RegExp(\Sigma)}$ be a binary operation for an alphabet $\Sigma$ defined recursively as follows:

$$S \cdot \emptyset = \emptyset,$$
$$S \cdot \varepsilon = S,$$
$$\emptyset \cdot r' = \emptyset,$$
$$\{\varepsilon\} \cdot r' = \{r'\},$$
$$\{r\} \cdot r' = \{rr'\},$$
$$(S \cup S') \cdot r' = (S \cdot r') \cup (S' \cdot r'),$$

where $S, S' \subseteq RegExp(\Sigma), r \in RegExp(\Sigma) \setminus \{\emptyset\}$, and $r' \in RegExp(\Sigma) \setminus \{\emptyset, \varepsilon\}$.

Let $r \in RegExp(\Sigma)$ be a *regular expression* over $\Sigma$. The set of partial derivatives of $r$ with respect to a symbol $b \in \Sigma$, represented as $\partial_b(r)$, is defined as follows:

$$\partial_b(\emptyset) = \emptyset, \qquad\qquad \partial_b(r + r') = \partial_b(r) \cup \partial_b(r'), \ r' \neq r$$
$$\partial_b(\varepsilon) = \emptyset, \qquad\qquad \partial_b(rr') = \partial_b(r)r' \cup \partial_b(r'), \ \text{ if } \varepsilon(r) = \varepsilon$$
$$\partial_b(b) = \varepsilon, \qquad\qquad \partial_b(rr') = \partial_b(r)r', \ \text{ if } \varepsilon(r) = \emptyset$$
$$\partial_b(c) = \emptyset, \ b \neq c \text{ and } b \in \Sigma \qquad \partial_b(r^*) = \partial_b(r)r^*$$

The set of partial derivatives of $r$ with respect to a word $w \in \Sigma^*$, denoted by $\partial_w(r)$, is defined recursively as:

$$\partial_\varepsilon(r) = \{r\},$$
$$\partial_{wa}(r) = \bigcup_{r' \in \partial_w(r)} \partial_a(r'), \quad \text{for } a \in \Sigma, w \in \Sigma^*$$

For instance, given the regular expression $r = (ad + d^*)db$, where $\Sigma = \{a, b, d\}$, the set of partial derivatives of $r$ for the symbol $d$ is given by

$$\begin{aligned} \partial_d(r) &= \partial_d((ad + d^*)db) \\ &= \partial_d(ad + d^*)db \cup \partial_d(db) \\ &= (\partial_d(ad) \cup \partial_d(d^*))db \cup \{b\} \\ &= (\partial_d(a)d) \cup \partial_d(d)d^*)db \cup \{b\} \\ &= (\emptyset \cup \{d^*\})db \cup \{b\} \\ &= \{d^*db\} \cup \{b\} \end{aligned}$$

resulting in $\partial_d(r) = \{d^*db, b\}$.

### 2.3.3.1   Linear Form

It is possible to calculate the partial derivatives for every symbol $a \in \Sigma$ of a regular expression $r$ over $\Sigma$ in a single pass. To do this, Antimirov [2] defined the *linear form*.

A linear form $lf$ of a regular expression $r$ is a finite set of symbol–continuation pairs. These pairs are called *monomials* and they can encode all the possible one-symbol prefixes of that regular expression. Given $r \in RegExp(\Sigma)$ and $a \in \Sigma$, a pair $(a, r)$ is called a monomial, whose head is $a$ and tail is $r$.

Formally, given an alphabet $\Sigma$ and the set of all regular expressions over $\Sigma$ denoted $RegExp$, a monomial is an element of $\Sigma \times RegExp$. We denote the set of all monomials for a given alphabet $\Sigma$ by $M_\Sigma$.

For $r \in RegExp(\Sigma) \setminus \{\emptyset\}$, $t \in RegExp(\Sigma) \setminus \{\varepsilon, \emptyset\}$, $a \in \Sigma$ and $S, S' \subseteq 2^{M_\Sigma}$, the operator of concatenation for monomials is given by $\cdot : 2^{M_\Sigma} \times RegExp(\Sigma) \to 2^{RegExp(\Sigma)}$ and is defined as:

$$\begin{aligned} S \cdot \emptyset &= \emptyset \\ S \cdot \varepsilon &= \varepsilon \\ (S \cup S') \cdot t &= (S \cdot t) \cup (S' \cdot t) \\ \{(a, \emptyset)\} \cdot t &= \{(a, \emptyset)\} \\ \{(a, \varepsilon)\} \cdot t &= \{(a, t)\} \\ \{(a, r)\} \cdot t &= \{(a, rt)\} \end{aligned}$$

A *linear form* is then a finite set of monomials:

$$lf(r) \subseteq 2^{M_\Sigma}.$$

For instance, a linear form $l = \{(a_1, \alpha_1), \ldots, (a_n, \alpha_n)\}$ corresponds to the regular expression

$$\kappa_l = a_1 \cdot \alpha_1 + \cdots + a_n \cdot \alpha_n.$$

The function $lf : RegExp \to 2^{M_\Sigma}$ is defined recursively as follows, where $\alpha, \beta \in RegExp$:

$$lf(\emptyset) = \emptyset$$
$$lf(\varepsilon) = \emptyset$$
$$lf(a) = \{(a, \varepsilon)\}, \quad a \in \Sigma$$
$$lf(\alpha + \beta) = lf(\alpha) \cup lf(\beta)$$
$$lf(\alpha \cdot \beta) = (lf(\alpha) \cdot \beta) \cup \begin{cases} lf(\beta) & \text{if } \varepsilon \in L(\alpha) \\ \emptyset & \text{otherwise} \end{cases}$$
$$lf(\alpha^*) = lf(\alpha) \cdot \alpha^*$$

Antimirov proved in [2, Proposition 2.5] through induction that one can decompose a regular expression into all of its one-symbol branches plus the possible empty-word case. For any regular expression $\alpha$ over $\Sigma$, the following holds:

$$\alpha = \sum_{(a, \alpha') \in lf(\alpha)} a\alpha' + \varepsilon(\alpha)$$

Equivalently, we also have

$$\partial_a(\alpha) = \{\alpha' \mid (a, \alpha') \in lf(\alpha)\}$$

## 2.4   From Regular Expressions to Automata

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

### 2.4.1   Thompson's Algorithm

In 1968, Thompson provided a method capable of transforming a regular expression into an equivalent nondeterministic finite automaton with $\varepsilon$-transitions (an $\varepsilon$-NFA) [19]. The main idea was to associate each regular-expression operator with a small automaton fragment, and then combine these fragments recursively until the entire expression is represented.

Each basic symbol $a \in \Sigma$ is represented by two states with a transition labeled $a$ between them. The empty word $\varepsilon$ is represented by two states connected by a single $\varepsilon$-transition.

More complex expressions are built by combining smaller fragments:

- Concatenation $(\alpha\beta)$ is obtained by connecting the accept state of $\alpha$ to the start state of $\beta$ with an $\varepsilon$-transition.

- Union $(\alpha \mid \beta)$ is built by adding a new start state with $\varepsilon$-transitions to the start states of $\alpha$ and $\beta$, and a new accept state with $\varepsilon$-transitions from the accept states of $\alpha$ and $\beta$.

- Kleene star $(\alpha^*)$ is represented by adding a new start and accept state. The new start connects by $\varepsilon$ both to the new accept (for the empty repetition) and to the start of $\alpha$ (for one or more repetitions). The accept state of $\alpha$ connects by $\varepsilon$ back to its own start and to the new accept.

Thompson described this as a compiler-like process: the expression is first parsed into a convenient form (such as reverse Polish notation), and then a stack-based procedure builds the automaton fragment by fragment.

In the end, a single fragment represents the entire regular expression. The resulting automaton can be simulated efficiently by maintaining two lists of active states: the current list and the next list. For each input symbol, all $\varepsilon$-reachable states are added to the current list, transitions on the current symbol are followed, and the next list becomes the current list for the following step. A word is accepted if the accept state becomes active at the end of the input. An example of an automaton constructed using Thompson's method is shown in Figure 2.3.
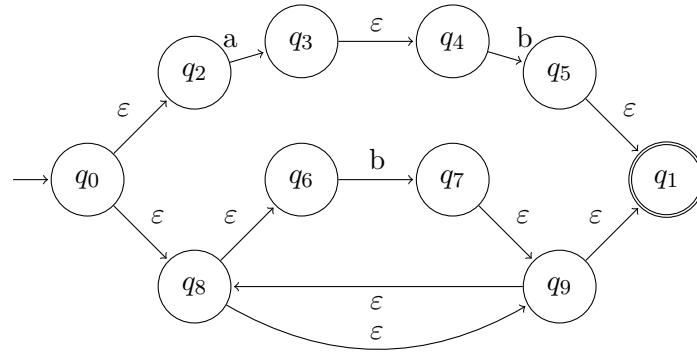


Figure 2.3: Diagram of a $\varepsilon$-NFA built using Thompson's construction for the regular expression $(ab + b^*)$.

### 2.4.2  Brzozowski's Derivatives

In 1964, Brzozowski [4] proposed a method of constructing deterministic finite automata using derivatives, where consecutive derivations of a regular expression would result in the states and their transitions with respect to the symbol derived.

**Definition 2.4.1** (Brzozowski's Automaton [4])**.** Let $\alpha \in RegExp(\Sigma)$. Brzozowski's automaton for the regular expression $\alpha$ is defined as the deterministic finite automaton $\mathcal{D} = (Q, \Sigma, q_0, \delta, F)$. $Q$ is the set of all dissimilar derivatives of $\alpha$, $q_0$ is the initial state $(q_0 \in Q)$, $\delta$ is the transition function defined by $\delta : Q \times \Sigma \to Q$ such that $\delta(q, a) = d_a(q)$ for all $a \in \Sigma$ and $q \in Q$, and the set of final states is $F = \{q \in Q \,|\, \varepsilon(q) = \varepsilon\}$.

For instance, let $\alpha = (ad + d^*)db$ over $\Sigma = \{a, b, d\}$. The transition function $\delta$ is defined for a

Brzozowski's automaton $\mathcal{D}$ of $\alpha$ as:

$$\delta((ad + d^*)db, a) = d_a((ad + d^*)db)$$
$$= d_a(addb + d^*db)$$
$$= d_a(addb) + d_a(d^*db)$$
$$= d_a(a)ddb + \emptyset + \emptyset$$
$$= \varepsilon ddb$$

$$\delta((ad + d^*)db, b) = d_b((ad + d^*)db)$$
$$= d_b(addb + d^*db)$$
$$= d_b(addb) + d_b(d^*db)$$
$$= \emptyset + \emptyset = \emptyset$$

$$\delta((ad + d^*)db, d) = d_d((ad + d^*)db)$$
$$= d_d(addb + d^*db)$$
$$= d_d(addb) + d_d(d^*db)$$
$$= \emptyset + d_d(d^*)db + d_d(db)$$
$$= d_d(d)d^*db + d_d(d)b + d_d(b)$$
$$= \varepsilon d^*db + \varepsilon b + \emptyset$$
$$= d^*db + b$$

$$\delta(ddb, a) = d_a(ddb)$$
$$= d_a(d)db$$
$$= \emptyset$$

$$\delta(ddb, b) = d_b(ddb)$$
$$= d_b(d)db$$
$$= \emptyset$$

$$\delta(ddb, d) = d_d(ddb)$$
$$= d_d(d)db$$
$$= \varepsilon db$$
$$= db$$

$$\delta(d^*db + b, a) = d_a(d^*db + b)$$
$$= d_a(d^*db) + d_a(b)$$
$$= \emptyset + \emptyset = \emptyset$$

$$\delta(d^* db + b, b) = d_b(d^* db + b)$$
$$= d_b(d^* db) + d_b(b)$$
$$= \emptyset + \varepsilon = \varepsilon$$

$$\delta(d^* db + b, d) = d_d(d^* db + b)$$
$$= d_d(d^* db) + d_d(b)$$
$$= d_d(d^*)db + d_d(db) + \emptyset$$
$$= d_d(d)d^* db + d_d(d)b + d_d(b)$$
$$= \varepsilon d^* db + \varepsilon b + \emptyset = d^* db + b$$

$$\delta(db, a) = d_a(db)$$
$$= d_a(d)b$$
$$= \emptyset$$

$$\delta(db, b) = d_b(db)$$
$$= d_b(d)b$$
$$= \emptyset$$

$$\delta(db, d) = d_d(db)$$
$$= d_d(d)b$$
$$= \varepsilon b = b$$

$$\delta(b, a) = d_a(b)$$
$$= \emptyset$$

$$\delta(b, b) = d_b(b)$$
$$= \varepsilon$$

$$\delta(b, d) = d_d(b)$$
$$= \emptyset$$

Finally, the set of states for $\mathcal{D}$ is $Q = \{(ad + d^*)db, ddb, d^* db + b, db, b, \varepsilon, \emptyset\}$ and the set of final states is $F = \{\varepsilon\}$. The diagram for this automaton is shown in Figure 2.4.
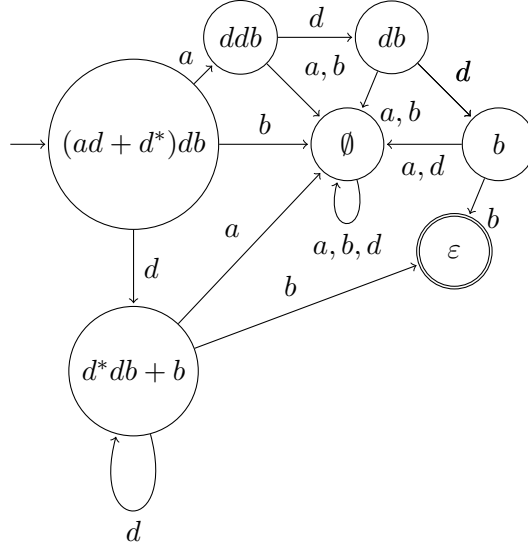
Figure 2.4: Diagram of a Brzozowki DFA for the regular expression $(ad + d^*)db$.

### 2.4.3   Antimirov's Partial Derivatives

Proposed by Valery Antimirov in 1996 [2], the partial derivatives construction generalizes Brzozowski'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.

**Definition 2.4.2** (Partial Derivatives Automaton)**.** Let $\alpha \in RegExp$ over $\Sigma$. The partial derivatives automaton for $\alpha$ is defined as the non-deterministic finite automaton $\mathcal{N} = (Q_{\mathcal{PD}}, \Sigma, q_0, \delta, F)$. $q_0$ is the initial state and belongs to the set of states $Q_{\mathcal{PD}}$. It is important to note that $q_0 = \alpha$. The transition function is defined as $\delta(q, a) = \partial_a(q)$ for all symbols $a \in \Sigma$ and states $q \in Q_{\mathcal{PD}}$. The set of final states is defined as $F = \{q \in Q_{\mathcal{PD}} \mid \varepsilon(q) = \varepsilon\}$.

Recalling from Chapter 2.3.3.1, the linear form can be used as an efficient way to compute partial derivatives. Let $\alpha \in RegExp$ over $\Sigma$ and $a \in \Sigma$:

$$\delta(\alpha, a) = \partial_a(\alpha) = \{\beta \mid (a, \beta) \in lf(\alpha)\}.$$

That is, to compute the set of partial derivatives of $\alpha$ with respect to a symbol $a$, one inspects the linear form of $\alpha$ and collects all continuations $\beta$ that follow an $a$ in the prefix position. This characterization is not only concise but also algorithmically useful: it gives a direct recipe for building the transition relation of the partial derivatives automaton.

**Theorem 2.4.1** ([2])**.** For every regular expression $\alpha$, the partial derivatives automaton $\mathcal{N}_\alpha$ accepts exactly the language $L(\alpha)$. Moreover, the number of states of $\mathcal{N}_\alpha$ is bounded by $|\alpha| + 1$.

### 2.4.4   Position Automata

The *position automaton*, also known as the *Glushkov automaton*, is a $\varepsilon$-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. [3]

Given a regular expression $r$, the Glushkov automaton $M_r$ is defined based on four key position-based functions:

- **first**$(r)$: the set of positions that can appear first in some word of the language $L(r)$.

- **last**$(r)$: the set of positions that can appear last in some word of $L(r)$.

- **follow**$(r, x)$: for each position $x$, the set of positions that can immediately follow $x$ in some word of $L(r)$.

- **symbol**$(x)$: for each position $x$, symbol$(x)$ represents the symbol on that position.

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 $(a_1 + b_2)^*a_3(b_4 + a_5)^*$, where the positions are $\{1, 2, 3, 4, 5\}$. Each marked symbol corresponds to a unique position and becomes a distinct state in the automaton.

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

- $Q$ is the set of positions in $r$ (i.e., the marked symbols), plus an initial state $q_0$.

- For each symbol $a \in \Sigma$:

  - $\delta(q_0, a) = \{x \in \text{first}(r) \mid \text{symbol}(x) = a\}$
  - $\delta(x, a) = \{y \in \text{follow}(r, x) \mid \text{symbol}(y) = a\}$

- The set of final states $F$ is last$(r)$; if $\varepsilon \in L(r)$, then $q_0$ is also final.

This automaton captures the structural flow of $r$ 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 $r$.

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 [16] 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`.
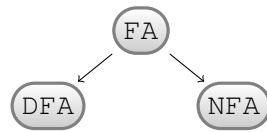
### 2.5.1  Finite Automata



Figure 2.5: Class inheritance organization in FAdo's finite automata code file, `fa.py`.

The class `FA` defines the abstract interface for finite automata, introducing the fundamental components: the set of states, the alphabet, the transition function, the initial state(s), and the set of final states. This class also establishes common methods for manipulating these elements and for interoperability between automata.

Two subclasses will then inherit from `FA`:

- `NFA` — implements nondeterministic finite automata, possibly with $\varepsilon$-transitions. It provides methods for epsilon-closure, elimination of $\varepsilon$-moves, subset construction into a DFA, and regular operations such as union, concatenation, and Kleene star.

- `DFA` — implements deterministic finite automata. It enforces determinism by construction and provides efficient word evaluation and minimization routines.

### 2.5.2  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. ($\emptyset$)

- **CEpsilon**: The empty word. ($\varepsilon$)

- **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 $\cdot$ operation

- **CStar**: The Kleene closure over a set of symbols (e.g. `CStar(CDisj(CAtom(a),CAtom(b)))` represents 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* [18] is a parsing toolkit for Python. It can parse any context-free language.

## 2.6 Single Matching

When referring to matching in the context of regular expressions, one typically means simply determining whether a given word is a member of the language defined by a specific regular expression. We can define this problem formally as:

**Definition 2.6.1.** Given $e \in RegExp(\Sigma)$ and a word $w \in \Sigma^\star$, we say that $w$ matches $e$ if $w \in L(e)$.

In practice, it is also useful to check if $w \in L(\Sigma^\star e \Sigma^\star)$, for a given word $w \in \Sigma^\star$ and $e \in RegExp(\Sigma)$. This process is known as *sub-matching.*

To look for matches using *FAdo*, we first need to choose a regular expression. For the sake of this example, let $r = ab + ba$. Informally, this regular expression will match any *ab* or *ba* character sequences in an input text. To do this, we first need to import the grammar using Lark and parse the regular expression to get a parsed syntax tree. With this tree, we can use *FAdo*'s `BuildRegexp` class and invoke its `transform` method, which will iteratively transform each of the trees nodes into abstract *FAdo* class objects:

```
r = "ab+ba"
regGrammar = lark.Lark.open("lang/regexp_test.lark", start="rege",
    parser="lalr")
tree = regGrammar.parse(r)
reg : RegExp = BuildRegexp(context={"sigma": None}).transform(tree)
reg.setSigma(reg.setOfSymbols())
```

We now have a regular expression object. We can use it to build a DFA using Brzozowski's derivatives construction:

```
dfa_z = reg.dfaBrzozowski()
```

Finally, having a DFA object, we can then use the `evalWordP` function to verify if $w \in L(r)$ for a given $w \in \Sigma^\star$, which can be, for this example, $w = ab$:

```
w = "ab"
if dfa_z.evalWordP(w):
```

```
        print("w belongs to the language of r")
    else:
        print("w does not belong to the language of r")
```

With the result being:

```
    w belongs to the language of r
```

# Chapter 3

# State of the Art

## 3.1 Introduction

Regular expressions (regex) remain one of the most powerful and widely adopted tools for string matching across programming languages, search tools, and data processing pipelines. While their theoretical foundation lies in formal language theory, the practical implementation of regex engines often diverges from the idealized models. This chapter summarizes the state of the art in regex engine designs and regex language features, with a particular focus on performance trade-offs, security concerns, and evolving capabilities.

## 3.2 Regex Syntax Variants

Beyond the theoretical framework of regular expressions lies a diverse set of practical syntax variants. Different tools, programming languages, and libraries adopt different interpretations and extensions of regex syntax, leading to a rich but inconsistent ecosystem. Two of the most influential and widely adopted syntax families are POSIX and PCRE.

### 3.2.1 POSIX Regular Expressions

The POSIX (Portable Operating System Interface) standard defines a family of regular expression syntaxes intended to promote consistency across UNIX-like systems. POSIX defines two primary levels of syntax: *Basic Regular Expressions* (BRE) and *Extended Regular Expressions* (ERE). These are used in many command-line utilities such as `grep` [7], `sed` [8], and `awk` [9].

In BRE, several metacharacters (such as parentheses, curly braces, and the plus sign) must be escaped to be interpreted as operators. For example:

- `a\{3\}` matches exactly three consecutive a's (`aaa`, for instance).

- `\(a|b\)` matches either a or b using grouping and alternation.

In contrast, ERE relaxes these constraints and allows a more expressive and readable syntax without requiring escapes for most metacharacters. For example:

- `a{3}` matches aaa (three consecutive a's).

- `a+` matches one or more a's.

- `a|b` matches either a or b.

- `(abc)?` matches either abc or $\varepsilon$.

A defining feature of POSIX regular expressions is their *leftmost-longest match* semantics. When multiple matches are possible, the engine must return the one that starts at the earliest position in the input and is the longest among those starting at that position. This requirement ensures deterministic behavior but can increase the complexity of the matching algorithm. [11]

POSIX regular expressions do not support many advanced features common in other regex dialects. For example, they lack support for backreferences, lookahead/lookbehind assertions, non-greedy quantifiers, and conditional expressions. This makes POSIX regexes more limited in expressiveness but often easier to reason about and implement efficiently.

### 3.2.2   PCRE: Perl-Compatible Regular Expressions

PCRE (Perl-Compatible Regular Expressions) define an expressive and flexible regex syntax, modeled after the pattern language used in the Perl programming language. It has become the standard for regular expression syntax in many modern programming environments, including PHP, Python, JavaScript, and others. PCRE introduces a wide range of features that go far beyond the capabilities of POSIX regexes and theoretical regular languages. These include:

- **Backreferences**, allowing a pattern to refer to previously matched groups.
  Example: `(a)b\1` matches aba (the `\1` refers to the first captured a).

- **Lookahead and lookbehind assertions**, enabling matches based on surrounding context without consuming it.
  Examples:

  - Lookahead: `foo(?=bar)` matches foo only if followed by bar.

  - Lookbehind: `(?<=foo)bar` matches bar only if preceded by foo.

- **Non-greedy quantifiers**, which match as little as possible.
  Example: `<.*?>` matches the shortest HTML-like tag, such as `<b>` in `<b>bold</b>`, instead of greedily matching the entire string.

- **Named capture groups**, for improved readability and maintainability.
  Example: `(?<year>\d{4})-(?<month>\d{2})` captures `2025-09` with named groups `year` and `month`.

- **Conditional expressions**, recursion, atomic groups, and other advanced constructs.
  Examples:

  - Conditional: `(a)?b(?(1)c|d)` matches `abc` if `a` is present, `bd` otherwise.
  - Recursion: `(a|b(?1))*` matches balanced nested patterns using self-reference.
  - Atomic group: `(?>a+)` prevents backtracking inside the group.

The added power of PCRE comes at the cost of increased complexity in both the syntax, generative expressivity and the engine. To support these features, most PCRE-compatible engines rely on backtracking-based evaluation, which explores all potential matching paths through a pattern. While this approach enables support for advanced constructs, it also introduces the risk of exponential-time behavior in certain patterns—especially those involving nested quantifiers or ambiguous alternations.

## 3.3 Engine Architectures

Regular expressions are most commonly used to look for a pattern in a string. This is generally called *matching*. The basis for this operation are the regular expression engines. At a high level, regex engines can be grouped by how they traverse the implicit nondeterministic automaton of a pattern.

### 3.3.1 Backtracking NFA

The classic backtracking NFA matcher is driven by the pattern: the regular expression acts like a small procedural program that dictates how the engine explores matches and handles failure. The engine begins at the start of the text and attempts to match the pattern from that position; if it fails, it "bumps along," advancing one character and trying again. Once the first tokens of the pattern match the text, the engine proceeds through the regex. At any choice point—such as an alternation, an optional, or a quantifier—it selects one alternative to try and records the others, together with the current input position. If the chosen path later fails, the engine backtracks to the most recent choice point and resumes with a saved alternative. If all saved alternatives are exhausted, the attempt from that starting position fails and the engine bumps along to the next character.

If the engine reaches the end of the pattern with all constraints satisfied, it declares success and discards any remaining, unexplored alternatives. A key consequence is that the order of alternatives matters: typical backtracking engines implement leftmost-first semantics rather than leftmost-longest, so the first workable alternative can win even if a longer match exists.

Because the engine literally follows the structure of the regex, you control its search by how you write the pattern: place safer or more selective alternatives first, avoid ambiguous constructs under repetition, and structure patterns to minimize backtracking and to fail fast when appropriate.

Despite the name, a backtracking NFA engine is not an NFA in the formal sense used in automata theory. Theoretical NFAs are memoryless machines that recognize only regular languages, whereas real-world backtracking engines—such as those used in PCRE, Java, and Python—go beyond regularity (Câmpeanu et al. proved this using a pumping lemma in [5]). They simulate nondeterministic behavior by exploring multiple paths through the pattern, but they do so using a stack and internal memory to track backtracking points, group captures, and even previous input matches. This allows them to support powerful features like backreferences and conditional expressions, which cannot be recognized by finite automata. The term "NFA" here refers to the engine's matching strategy, inspired by the branching behavior of NFAs, rather than to its computational limitations.

### 3.3.2   POSIX NFA

Similarly to the classic NFA engine, the POSIX NFA engine will match in the same way but memorize and continue when a successful match is found. This is done to see if a longer, leftmost match can be found later on.

### 3.3.3   DFA

A DFA-based matcher is driven by the input text rather than the structure of the regular expression. The pattern is first compiled into a deterministic finite automaton (DFA), which enables the engine to scan the input in a single pass. As each character is read, the engine transitions deterministically from one state to another, maintaining only a single active state at any point during execution.

Each DFA state represents a set of possible continuations in the pattern, allowing the matcher to recognize valid strings without the need for backtracking. This results in highly predictable performance, with matching taking linear time in the length of the input.

In search mode, DFA-based engines typically return the leftmost match. Many implementations are also designed to return the leftmost-longest match by continuing the scan while remembering the last encountered accepting state. Because the DFA encodes all matching possibilities in advance, the matching process is entirely deterministic and independent of the syntactic order of alternatives in the pattern.

One practical consideration in DFA-based matching is the potential growth of the automaton during the compilation phase. In the worst case, the number of DFA states can grow exponentially with respect to the size of the regular expression, particularly for patterns involving complex

nesting or alternation. While this does not affect runtime performance during matching — which remains linear with respect to the size of the input — the memory and time costs of building the DFA can be significant. To address this, many modern implementations adopt strategies such as lazy DFA construction or hybrid evaluation models that balance determinism with scalability.

### 3.3.4  Hybrid

Hybrid engines try to take the best of both worlds in both NFAs and DFAs. They perform a depth-first search through the space of matches. At each point of nondeterminism—due to alternation, optional constructs, or unbounded quantifiers — they choose one option to continue and save the others as choice points on an internal stack. If a later step fails, the engine pops a choice point and resumes from there. This approach yields intuitive behavior and supports rich features, including capturing groups, backreferences, and look-around. However, in the presence of ambiguous subpatterns under repetition, the number of explored paths can grow exponentially, making these engines particularly susceptible to ReDoS.

Spencer's classic backtracking implementation is perhaps the best known and most used, embodying closely the architecture described above. Conceptually, its state comprises the current regex node, the current input index, and a stack of saved alternatives. Consider the pattern `(a|aa)+b$` against the input `aaaa...a` without a trailing `b`. The matcher repeatedly chooses the left alternative `a`, consuming a single character while saving a choice point for the `aa` branch at each iteration. At the end of input it fails to satisfy the end anchor, so it begins to backtrack, popping choice points and trying to repartition the run of `a`s using `aa` segments. The number of such partitions grows exponentially with the length of the input, and the engine must examine many of them before determining there is no match. This example illustrates how overlapping alternatives inside a quantifier, combined with a terminal failure, can lead directly to catastrophic backtracking. The tools GNU egrep and awk use a hybrid approach for this.

## 3.4  Engines and Libraries

Modern regular expression engines vary widely in their design goals, trade-offs, and supported features. While traditional engines prioritize expressiveness (often at the cost of performance guarantees) there has been a growing interest in implementations that prioritize safety, predictability, and efficiency, especially in systems that process untrusted or large-scale input. This section presents a selection of such libraries and engines.

### 3.4.1  RE2

RE2 is a regular expression engine developed by Russ Cox at Google, designed to provide predictable performance and strong safety guarantees. Unlike backtracking engines, RE2 avoids

features that lead to exponential-time behavior, such as backreferences and arbitrary lookbehinds. By restricting itself to regular languages, RE2 ensures that all matches can be performed in linear time. [6] However, being restricted to regular languages is not enough. It is also good practice to use DFAs or NFAs accordingly, instead of resorting to backtracking.

The engine compiles patterns into automata using a combination of DFA and NFA techniques. For simple match queries, RE2 prefers deterministic finite automata (DFA) for their speed and predictability. When submatch extraction (identifying not just whether a string matches a regular expression, but also which specific substrings of the input match particular subexpressions) is needed, it may fall back to an NFA simulation, provided the regex satisfies certain properties (e.g., being one-pass). [6]

By trading off advanced features like backreferences for performance and safety, RE2 provides a robust alternative to PCRE-style engines, especially in environments where worst-case behavior must be avoided.

RE2 can also operate in either POSIX mode (accepting POSIX syntax regular expressions) or Perl mode (accepting PCRE syntax regular expressions). [10]

### 3.4.2   PCRE2

PCRE2 is a modern and widely adopted regular expression library that implements the rich, feature-complete syntax mentioned above (PCRE). [15]

At its core, PCRE2 relies on a backtracking-based matcher, which simulates nondeterminism by exploring alternative execution paths through the pattern. Although this resembles the behavior of nondeterministic finite automata (NFAs), it does not correspond to an NFA in the theoretical sense. Instead, the engine maintains an explicit control stack and memory for backtracking, enabling it to handle features like backreferences and complex group captures. This matching strategy enables high expressiveness but comes with the risk of exponential time complexity in certain pathological cases, particularly when ambiguous repetition and nested alternations are involved. [15]

PCRE2 also offers an alternative, limited matching mode based on deterministic finite automata (DFA), which guarantees linear-time performance but does not support the full range of Perl-compatible features. Despite these limitations, PCRE2 remains widely used in scripting languages and developer tools due to its flexibility and familiar syntax. [15]

### 3.4.3   Hyperscan

Hyperscan is Intel's regular expression matching engine, designed specifically for high-throughput and low-latency applications. It serves as a core component in several security and networking tools, including intrusion detection systems like Suricata and firewalls.

Traditional regex engines (e.g., backtracking-based ones like PCRE) often struggle with performance bottlenecks due to their sequential nature and vulnerability to ReDoS attacks. Hyperscan addresses these limitations by combining multiple automata models—particularly NFAs and DFAs—with a hybrid execution strategy that leverages Single Instruction, Multiple Data (SIMD) parallelism and tiled execution on modern CPUs .

According to Wang et al. ([21]), Hyperscan divides regexes into multiple subgraphs, such as anchored DFAs for simple patterns and NFAs for complex constructs. This hybrid approach enables it to process large volumes of data streams efficiently without the exponential-time risks associated with backtracking engines.

However, as noted in [21], Hyperscan will also enforce syntactic restrictions when compiling. Regexes that are deemed vulnerable or considered ambiguous, therefore limiting expressiveness and versatility, especially when the user is looking for nested repetition (e.g. $(a+)+$, looking to match one or more of one or more $a$'s) or greedy alternation (e.g. $(a|aa)+$, matching a sequence of $a$ or $aa$, repeated, resulting in three matches for a string such as $aa$).

### 3.4.4   Rust's `regex` library

Rust is a systems programming language focused on safety, concurrency, and performance, offering memory safety without garbage collection through its unique ownership model. The `regex` crate implements a hybrid DFA/NFA model and guarantees linear-time performance. All regex searches in Rust's `regex` library have worst case $O(m*n)$ time complexity, where $m$ is the size of the regular expression and $n$ is the size of the input string. [1]

To uphold this guarantee, the pattern syntax is intentionally restricted: features that are not known to admit efficient matching—most notably *backreferences* and general *look-around assertions*—are omitted. The engine's design combines automata techniques: a lazy DFA is used to accelerate finding overall matches, while a PikeVM/Thompson-NFA simulation recovers submatch boundaries when needed; these strategies avoid recursive backtracking and keep stack usage bounded. The crate provides first-class Unicode support (including Unicode character classes and word boundaries) and exposes its parser and automata components as reusable libraries (`regex_syntax`, `regex-automata`). In practice, this places `regex` in the same family as RE2-style engines: it favors regular-language constructs and predictable performance over Perl-style expressiveness.

## 3.5   Tools for ReDoS Detection and Exploitation

### 3.5.1   Revealer

REVEALER is a hybrid static–dynamic analysis tool designed to detect and exploit ReDoS vulnerabilities, particularly in regex engines that support extended, non-regular features such

as backreferences and lookahead assertions. Unlike traditional ReDoS analyzers that operate under the assumption that regular expressions conform to classical finite automata semantics, REVEALER targets the realities of modern backtracking-based engines like Java's regex library, where non-regular constructs are common and matching behavior is often implementation-specific [14]. The tool begins with a static analysis phase, modeling regexes as extended nondeterministic finite automata. It detects vulnerable structural patterns (such as nested loops or ambiguous branching) by analyzing the internal state graph used by the Java regex engine. These structures are known to induce exponential or polynomial backtracking behavior under certain input conditions. To confirm whether a suspected pattern can cause denial-of-service, REVEALER proceeds with a dynamic validation phase that simulates matching and attempts to synthesize concrete attack strings. This dual-phase approach minimizes false positives and ensures that reported vulnerabilities are exploitable in practice. REVEALER was evaluated on over 29,000 Java regexes and successfully identified all previously known vulnerabilities while discovering 213 new ones [14].

### 3.5.2   ReScue

ReScue is a technique for detecting ReDoS vulnerabilities in modern regex engines. It can generate special input strings that, when submitted as an input to a pre-defined regex engine, will lead to ReDoS. To do this, it starts by exploring the regex engine's behaviors under various inputs, which is called the *seeding phase*. These inputs are selected and evaluated on how many of the $\varepsilon$-NFA states they reach, regardless of the time it takes to reach them. After that, it will use the selected strings from the seeding phase and apply "mutations" to them, this phase is called the *incubating phase*. Lastly, the *pumping phase* will enhance the strings by finding which parts of the string take the longest to process, and copy and paste them one after the other. Parallel to this, it will also trim unnecessary parts of the string in order to retain the attack core (the minimal input responsible for the performance hit). [17]

ReScue was evaluated on 29088 real-world regexes and outperformed prior tools including RXXR2, Rexploiter, and SlowFuzz, finding 49% more exploitable inputs and uncovering ten previously unknown vulnerabilities in popular open-source projects [17].

### 3.5.3   ReDoSHunter

ReDoSHunter is a vulnerability detection framework that combines static analysis with dynamic validation to identify ReDoS vulnerabilities in modern regex engines. Unlike prior approaches that suffer from a trade-off between precision and recall, ReDoSHunter is designed to achieve both high accuracy and coverage. The framework begins by transforming real-world regexes into a simplified form (standardized regex) while preserving ReDoS-relevant semantics. It then uses pattern-specific static diagnosis algorithms to identify potential vulnerabilities and synthesize attack strings. Each candidate is subsequently validated dynamically against real

regex engines to ensure the vulnerability is exploitable in practice. ReDoSHunter also detects multiple vulnerabilities in a single regex, a feature lacking in most prior tools. Empirically, ReDoSHunter outperforms seven state-of-the-art detection tools in both precision and recall. It achieved 100% precision and 100% recall across multiple large-scale datasets containing over 37,000 regexes and discovered 28 previously unknown vulnerabilities in widely used software projects, 26 of which were assigned CVEs [13]. This level of precision and recall mean that the tool triggers very few false alarms, and misses very few real vulnerabilities.

## 3.6  Reluctance

Theory for automata has been around for very long, and its theory has been checked and developed by many. However, software is often built from demand. Most of the time, someone builds something because it could be useful. Regular expressions entered practical programming through Ken Thompson's implementation in the QED text editor on CTSS (and later, in Dennis Ritchie's version for GE-TSS). When Thompson and Ritchie developed Unix, they carried regular expressions with them to the new ecosystem, adding them to tools such as `awk` and `sed`. By the end of the 1970s, regular expressions were a fundamental part of Unix. Nowadays, regular expressions illustrate the risks of straying from theoretical principles. Contemporary regex engines (like backtracking, for example) are often far less efficient than the real automata-based implementations used in the early Unix tools. Despite known limitations and security risks associated with traditional backtracking-based regular expression engines, there remains significant resistance to replacing them in many modern and legacy systems. This reluctance doesn't stem only from the fact that people don't want to change stuff if it "works", especially when it is so embedded into a trustworthy system:

- **Backwards compatibility**: Changing old code could break countless projects, many of them supporting running projects today

- **Pattern rewritting**: While tools like RE2 offer guaranteed linear-time matching, they lack features such as backreferences or lookaround assertions, which make them unsuitable for existing expressions that depend on those features

- **The risk**: Many organizations fear regressions or downtime, and often ponder if changing this old code is worth it

For instance, Irani et. al ([12]) performed a study within the public administration domain but it can be applied here. They found that there are multiple layers of complexity that prevent modernization, ranging from outdated architectures to deep interdependencies across institutions. The systems they studied were often tied to statutory processes, policy constraints and siloed governance, making them difficult to refactor or retire. Irani et. al found that overcoming a legacy system's inertia to change often requires more than just technical solutions. It demands

stakeholder buy-in, incremental migration strategies, semantic data interoperability and, most importantly, alignment between policy ambitions and organizational capacity.

# Chapter 4

# Counting

In this chapter, we will explore the concept of counting in the context of formal languages and automata theory as well as explain the implementation for fixed and bounded repetition in FAdo.

## 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. Current tools are already able to do this, including some non-backtracking matchers. However, even for tools based on NFA and DFA, these operators are problematic due to the complexity of their manipulation. The aim of the work here described is to build derivative-based tools that may be more efficient and/or secure.

## 4.2 Derivatives of fixed and bounded repetition operators

Given a word $w$ and a regular expression $e$, $w \in L(e)$ if $\varepsilon(d_w(e)) = \varepsilon$. In order to match with fixed and bounded repetition using derivatives, one must first define them for those operations.

**Theorem 4.2.1.** Given $e \in RegExp(\Sigma)$, the following holds:

$$\varepsilon(e) + e = e$$

*Proof.* Let $e \in RegExp(\Sigma)$. The proof considers the cases where $\varepsilon \in L(e)$ and $\varepsilon \notin L(e)$. For $\varepsilon \in L(e)$, we have that $\varepsilon(e) = \varepsilon$. And then,

$$\varepsilon(e) + e = e$$
$$\Rightarrow \varepsilon + e = e$$
$$\Rightarrow e = e$$

holds. On the other hand, for $\varepsilon \notin L(e)$, we have that $\varepsilon(e) = \emptyset$. Therefore,

$$\varepsilon(e) + e = e$$
$$\Rightarrow \emptyset + e = e$$
$$\Rightarrow e = e$$

also holds. $\qquad\qquad\square$

### 4.2.1 Fixed Repetition

**Theorem 4.2.2.** Given $e \in RegExp(\Sigma)$ and $n \in \mathbb{N}^+$, we have:

$$d_a(e^n) = d_a(e)\, e^{n-1}.$$

*Proof.* We proceed by induction on $n \in \mathbb{N}^+$. For the base case $n = 1$,

$$d_a(e^1) = d_a(e) = d_a(e)e^0 = d_a(e)\varepsilon = d_a(e).$$

Since the identity holds for $n = 1$, lets assume that the identity holds for some $n - 1 \geq 1$, i.e.,

$$d_a(e^{n-1}) = d_a(e)e^{n-2}.$$

We want to show it holds for $n$. Using the product rule:

$$\begin{aligned}
d_a(e^n) &= d_a(ee^{n-1}) \\
&= d_a(e)e^{n-1} + \varepsilon(e)d_a(e^{n-1}) \\
&= d_a(e)e^{n-1} + \varepsilon(e)d_a(e)e^{n-2} \\
&= d_a(e)\left(e^{n-1} + \varepsilon(e)e^{n-2}\right) \\
&= d_a(e)\left(ee^{n-2} + \varepsilon(e)e^{n-2}\right) \\
&= d_a(e)(e + \varepsilon(e))e^{n-2} \\
&= d_a(e)e^{n-1}.
\end{aligned}$$

This completes the induction. $\qquad\qquad\square$

As an example, let $e = ab$ such that $e^2 = abab$. Using the theorem above:

$$\begin{aligned}
d_a(e^2) &= d_a(e)e \\
&= d_a(ab)ab \\
&= bab
\end{aligned}$$

We can also define partial derivatives and the linear form for this new operator. We have

$$\partial_a(e^n) = \partial_a(e)e^{n-1}$$
$$lf(e^n) = lf(e)e^{n-1}$$

for $e \in RegExp(\Sigma)$, $a \in \Sigma$, and $n \in \mathbb{N}^+$.

### 4.2.2 Bounded Repetition

**Theorem 4.2.3.** Given $e \in RegExp(\Sigma)$, $n \in \mathbb{N}^+$, $m > n$, and $a \in \Sigma$, we have:

$$d_a(e^{n,m}) = \begin{cases} d_a(e)e^{n-1,m-1}, & \text{if } n > 1, \\ d_a(e) + d_a(e)e^{1,m-1}, & \text{if } n = 1. \end{cases}$$

*Proof.* Let $e \in RegExp(\Sigma)$ and let $n, m \in \mathbb{N}^+$ with $m > n$. By definition of bounded iteration:

$$e^{n,m} = \sum_{i=n}^{m} e^i,$$

and by the linearity of derivatives over sums:

$$d_a(e^{n,m}) = \sum_{i=n}^{m} d_a(e^i).$$

Using the known identity $d_a(e^i) = d_a(e)e^{i-1}$, we have:

$$d_a(e^{n,m}) = \sum_{i=n}^{m} d_a(e)e^{i-1}$$
$$= d_a(e) \sum_{i=n}^{m} e^{i-1}.$$

Letting $j = i - 1$, this becomes:

$$d_a(e^{n,m}) = d_a(e) \sum_{j=n-1}^{m-1} e^j = d_a(e)e^{n-1,m-1}.$$

Considering both cases, for $n > 1$, the expression is already in its final form:

$$d_a(e^{n,m}) = d_a(e)e^{n-1,m-1}.$$

And for $n = 1$, we have:

$$d_a(e^{1,m}) = d_a(e)e^{0,m-1}$$
$$= d_a(e) + d_a(e)e^{1,m-1}.$$

This completes the proof.                                                                 □

Furthermore, we can include the infinite upper bound: $e^{n,\infty}$.

$$e^{n,\infty} = e^n e^\star$$

We can simplify it even more, depending on the value of $n$. Assuming $n = 0$:

$$e^{0,\infty} = e^0 e^\star = \varepsilon e^\star = e^\star$$

Lastly, for $n = 1$:

$$e^{1,\infty} = e^1 e^\star = ee^\star = e^+$$

These relations allow to substitute $r^\star$ and $r^+$ with counting expressions. However, in the *FAdo* package, we consider also the classic iterators.

**Theorem 4.2.4.** Given $e \in RegExp(\Sigma)$, $n \in \mathbb{N}_0$ and $a \in \Sigma$, we have:

$$d_a(e^{n,\infty}) = \begin{cases} d_a(e)e^{n-1,\infty}, & \text{if } n > 0, \\ d_a(e)e^{0,\infty}, & \text{otherwise.} \end{cases}$$

*Proof.* Let $e \in RegExp(\Sigma)$, $n \in \mathbb{N}_0$, and $a \in \Sigma$. We have,

$$\begin{aligned} d_a(e^{n,\infty}) &= d_a(e^n e^\star) \\ &= d_a(e^n)e^\star + \varepsilon(e)d_a(e)e^\star \\ &= (d_a(e^n) + \varepsilon(e)d_a(e))e^\star \\ &= (d_a(e)e^{n-1} + d_a(e)\varepsilon(e))e^\star \\ &= (d_a(e)(e^{n-1} + \varepsilon(e)))e^\star \\ &= d_a(e)e^{n-1}e^\star \\ &= d_a(e)e^{n-1,\infty}. \end{aligned}$$

Otherwise, if $n = 0$, we can just replace accordingly:

$$d_a(e^n e^\star) = d_a(e^0 e^\star) = d_a(e^\star) = d_a(e)e^\star = d_a(e)e^{0,\infty}.$$

This concludes the proof.                                                                                      $\square$

For instance, given $e = ab$ such that $e^2 = abab$ and $e^3 = ababab$, we have:

$$\begin{aligned} d_a(e^{2,4}) &= d_a(e)e^{[1,3[} \\ &= b(ab)^{[1,3[} \\ &= bab + babab \end{aligned}$$

Extending to the infinity case, $e^{2,\infty}$:

$$\begin{aligned} d_a(e^{2,\infty}) &= d_a(e)e^{1,\infty} \\ &= b(ab)e^{1,\infty} \\ &= bab + b(ab)^+ \end{aligned}$$

Lastly, we can define the partial derivatives and linear form for these operators too. We have for bounded repetition (excluding infinity cases),

$$\partial_a(e^{n,m}) = \begin{cases} \partial_a(e)e^{n-1,m-1}, & \text{if } n > 1, \\ \partial_a(e) + \partial_a(e)e^{1,m-1}, & \text{if } n = 1. \end{cases}$$

$$lf_a(e^{n,m}) = \begin{cases} lf_a(e)e^{n-1,m-1}, & \text{if } n > 1, \\ lf_a(e) + lf_a(e)e^{1,m-1}, & \text{if } n = 1. \end{cases}$$

and for its infinity cases,

$$\partial_a(e^{j,\infty}) = \begin{cases} \partial_a(e)e^{j-1,\infty}, & \text{if } j > 0, \\ \partial_a(e)e^{0,\infty}, & \text{otherwise.} \end{cases}$$

$$lf_a(e^{j,\infty}) = \begin{cases} lf_a(e)e^{j-1,\infty}, & \text{if } j > 0, \\ lf_a(e)e^{0,\infty}, & \text{otherwise.} \end{cases}$$

for bounded repetition, with $e \in RegExp(\Sigma)$, $a \in \Sigma$, $n \in \mathbb{N}^+$, $m > n$, $j \in \mathbb{N}_0$.

## 4.3 Implementation in FAdo

In order to implement these new syntactic constructs of extended regular expressions in FAdo, we added exact-power and counted-repetition nodes to the regular expression abstract syntax tree and define their derivatives, partial derivatives and linear form cases.

### 4.3.1 CPower

*CPower* is the class responsible for the construction of regular expressions using the fixed repetition operator.

```
class CPower(Unary):
    def __init__(self, arg, n, sigma=None):
    self.arg = arg
    self.n = n
    self.Sigma = sigma
    self._ewp = False if self.n > 0 else True


    ...
```

As shown above, the class inherits from the *Unary* class, which only defines the *Unary.Sigma* and *Unary.arg* properties, letting the class hold a alphabet set (representing $\Sigma$) and the symbol used to construct this operator. In order to integrate *CPower* fully with *FAdo*, some more methods had to be implemented.

```
    def linearForm(self):
        arg_lf = self.arg.linearForm()
        lf = dict()
        for head in arg_lf:
            lf[head] = set()
            for tail in arg_lf[head]:
                lf[head].add(self.derivative(head))
        return lf

    def partialDerivatives(self, sigma):
```

```
        return self.arg.partialDerivatives(sigma)

    def derivative(self, sigma):
        if str(sigma) in str(self.Sigma):
            if self.n == 0:
                return CEmptySet(sigma)
            elif self.n == 1:
                return self.arg.derivative(sigma)
            else:
                return CConcat(self.arg.derivative(sigma), CPower(self.arg,
                    self.n-1, self.Sigma))
        else:
            return CEmptySet(sigma)
```

### 4.3.2   CCount

*CCount* is the class responsible for the construction of regular expressions using the bounded repetition operator.

```
class CCount(Unary):
    def __init__(self, arg, min, max = None, sigma = None):
        self.arg = arg
        self.min = int(min)
        self.max = "inf" if max == -1 or max == "inf" else int(max)-1
        self.Sigma = sigma


        if self.min==0:
            self._ewp = True
```

As shown above, and much like *CPower*, the class also inherits from *Unary*. The constructor adjusts the bounds accordingly and sets the empty word property based on whether the minimum repetition is zero.

To fully integrate the class with *FAdo*, several methods are implemented, including string representation, partial derivatives, linear form computation, and marking functions:

```
    def derivative(self, sigma):
        if str(sigma) in str(self.Sigma):
            if self.max == "inf" or self.max == None:
                if self.min == 0:
                    return CConcat(self.arg.derivative(sigma), CStar(self.arg,
                        self.Sigma))
                else:
                    return CConcat(self.arg.derivative(sigma), CCount(self.arg,
                        self.min-1, self.max, sigma), self.Sigma)
            else:
                if self.max > 1:
                    if self.min == 0:
```

```python
                return CConcat(self.arg.derivative(sigma), CCount(self.arg,
                    self.min, int(self.max), self.Sigma))
            else:
                return CConcat(self.arg.derivative(sigma), CCount(self.arg,
                    self.min-1, int(self.max), self.Sigma))
        elif self.max == 1:
            return self.arg.derivative(sigma)
        else:
            CEpsilon(sigma)
    else:
        return CEmptySet(sigma)

def partialDerivatives(self, sigma):
    arg_pdset = self.arg.partialDerivatives(sigma)
    pds = set()
    for pd in arg_pdset:
        if pd.emptysetP():
            pds.add(CEmptySet(self.Sigma))
        elif pd.epsilonP():
            pds.add(self)
        else:
            pds.add(CConcat(pd, self, self.Sigma))
    return pds
```

The `linearForm` method is exactly the same as `CPower`'s.

### 4.3.3  Grammar

For ease of use and protoyping, we utilized *FAdo*'s ability to parse a regular expression using a common Python *string*. To do this, *FAdo* utilizes a library called Lark. Lark is a modern, fast and flexible parsing toolkit for Python that supports both LALR(1) and Earley parsers, enabling the parsing of all context-free grammars with features like EBNF, AST generation, and grammar composition [18]. The new grammar can be seen in Appendix A and the most notable changes are the following:

```
| rep _CARET _LEFBR digit _RIGBR -> pow_min
| rep _CARET _LEFBR digit _COMMA digit _LEFBR -> pow_minmax
| rep _CARET _LEFBR digit _COMMA _INFTY _LEFBR -> pow_inf
```

These rules (along with some more utilitary definitions such as the `_LEFBR` and `_RIGBR`, which are the symbols '[' and ']') will route the following cases (code in Appendix B):

- **pow_min**: For the cases $e^n$ and $e^{n,\infty}$, where $e \in RegExp(\Sigma)$ and $n \in \mathbb{N}^+$

- **pow_minmax**: For the case $e^{n,m}$, where $e \in RegExp(\Sigma)$ and $n \in \mathbb{N}^+, m > n$

- **pow_inf**: Calls *pow_min* with a special argument that selects the infinity case

### 4.3.4  Matching

To demonstrate how matching is performed using the implemented operators, we can consider a concrete example. For instance, let $r = (ab)^{1,3}$. The following code snippet shows how this expression can be parsed, transformed into an internal representation, and converted into a DFA using Brzozowski's construction:

```
r = "(ab)^[1,3["
tree = regGrammar.parse(r)
reg : RegExp = BuildRegexp(context={"sigma": None}).transform(tree)
reg.setSigma(reg.setOfSymbols())
dfa_z = reg.dfaBrzozowski()
dfa_z.display()
```
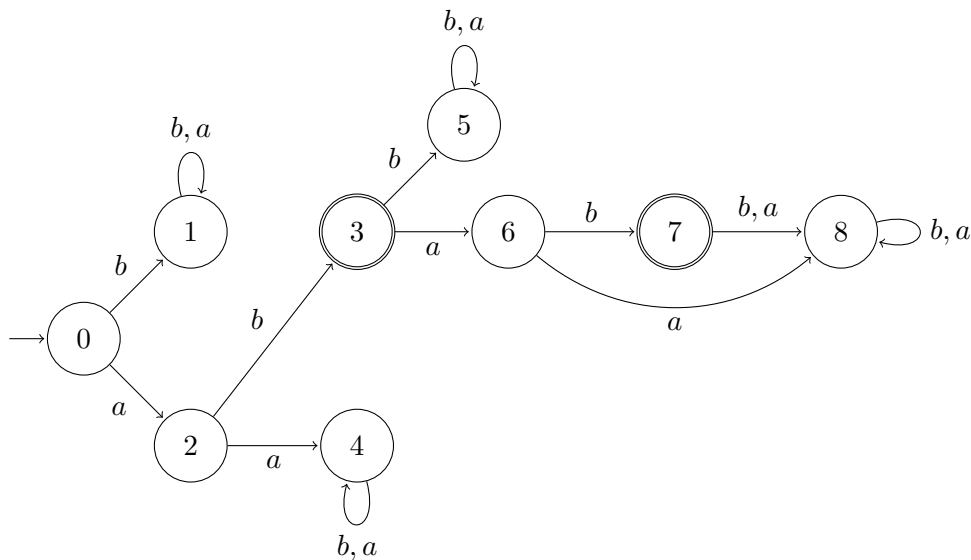


Figure 4.1: Diagram of the DFA built using Brzozowski's method for the regular expression $r = (ab)^{1,3}$.

Then, try to match $w = ab$ using the same code from the snippet 2.6, which yields the following:

```
w belongs to the language of r
```

Trying again with a text that does not match, $w = ababab$, yields:

```
w does not belong to the language of r
```

# Chapter 5

# Matching

*Matching* is the process of checking whether a piece of text fits a specific pattern described using a regular expression (regex). In this chapter, we will discuss the different approaches to matching regular expressions and their differences, the implications of using them, and the performance considerations that arise from these choices. Furthermore, due to the importance of multiple matching and the lack of tools that accomplish it, we will also contribute with a novel approach based on a modified position automaton and matcher, with the focus of mitigating the performance issues associated with traditional regex engines while preserving some of the extended expressiveness of regex patterns.

## 5.1    Multiple Matching

The definition of single matching was introduced in Definition 2.6.1. Now, the concept of multiple matching is simply to get all the matches of a given regular expression within an input text. We can define it formally like so,

**Definition 5.1.1.** Let $r \in RegExp(\Sigma)$, $(i, j) \in \mathbb{N}^{+2}$ where $j > i$, and $t \in \Sigma^\star$.

$$\text{matches}(r, t) = \{(i, j) \mid \exists w \in L(r), \ t \in \Sigma^i w \Sigma^j\}$$

Informally, we want to capture the positions in $t$ where a matching substring is between a prefix of length $i$ and a suffix of length $j$. In this chapter, we will use the term *matching* to refer to *multiple matching* or *all matching*.

## 5.2    Overlapped versus Non-Overlapped Matching

In the context of regular expression matching, two distinct paradigms exist: *overlapped matching* and *non-overlapped matching*. Understanding their differences is crucial when designing matching engines, especially when completeness or performance is a concern.

**Overlapped Matching**

Overlapped matching refers to finding all possible matches of a pattern in an input string. This is typically achieved by attempting a match starting at every index of the input. It is more exhaustive and useful in domains where no potential match should be missed, such as bioinformatics (DNA pattern searching).

For example, the indexed string $w = a_0a_1a_2a_3$, when matched against using the pattern $r = aa$, will yield the following matched substrings:

- $w_{[0,2[} = \textcolor{red}{aa}aa$

- $w_{[1,3[} = a\textcolor{red}{aa}a$

- $w_{[2,4[} = aa\textcolor{red}{aa}$

**Non-Overlapped Matching**

Non-overlapped matching (also referred to as *disjoint*, *standard*, or in some engines, *greedy* matching) finds matches sequentially from left to right, and once a match is found, it advances the input pointer beyond the match. Unlike overlapped matchers, where overlapping is a feature by default, greedy matchers will often depend on the lookahead assertions to do so.

Using the same pattern $r = aa$ on $w = a_0a_1a_2a_3$, a non-overlapping matcher may return:

- $w_{[0,2[} = \textcolor{red}{aa}aa$

- $w_{[2,4[} = aa\textcolor{red}{aa}$

As an example, this is the exact behaviour that PCRE2 demonstrates when trying to match for *aa* on an input text *aaaa*.

To summarize, overlapped matching provides a more complete view of potential matches, but at a higher computational cost. It is particularly well-suited to automata-based approaches like the modified position automaton described in this work.

## 5.3    Modified Position Automaton

In order to perform multiple matching, we had to use a different automaton. A *position automaton* is a construction that transforms a regular expression into a nondeterministic finite automaton (NFA) (2.4.4). The key component to this construction is the fact that each state corresponds to a position in the regular expression. For instance, let $r = a^{2,4}$. This expression can be built using Brzozowski's derivatives construction. The resulting DFA can be seen below:
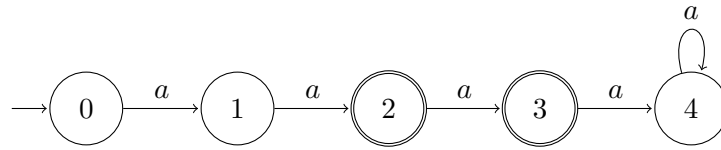
Figure 5.1: DFA resulting from Brzozowski's derivatives construction for the expression $r = a^{2,4}$.

For single matching, a typical matcher should be able to match in $O(n)$ using this DFA, where $n$ is the length of the input string. The matcher would need to be restarted after each starting position, leading to a running time complexity of $O(n^2)$. So, in order to perform overlapped matching in linear time, the automaton must be constructed using a modified position automaton construction. Using Algorithm 1 to build an NFA for $r$ results in the NFA represented below:
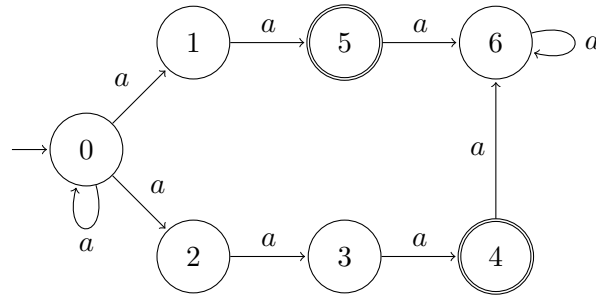


Figure 5.2: NFA resulting from the modified position automaton construction for the expression $r = a^{2,4}$.

This NFA now contains two different branches, one that captures all the matches for $aa$ and another one that captures all the matches for $aaa$. Also, it can continuously match due to the loops on states 0 and 6. To demonstrate the matcher logic, it is better to implement a more complex regular expression. Consider the following regular expression $s$ and input word $w$:

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

We can separate $R$ into two matching groups:

- The first group $(aa + aaa)$ will match either two or three $a$ symbols (e.g. $aaa$ will yield three overlapped matches: $aaa$, $aaa$ and $aaa$).

- The second group $(aaa + aa)$ will also match either two or three $a$ symbols, much like the first group.

The normal position automaton construction for $s$ is as follows, considering the marked version $(a_1a_2 + a_3a_4a_5)(a_6a_7a_8 + a_9a_{10})$:

Figure 5.3: Unmodified position automaton construction of $s$

Meanwhile, the modified position automaton for this regular expression can be constructed using the Algorithm 1, resulting in the following:



Figure 5.4: Modified position automaton construction of $s$

When we match using Algorithm 2 to an input string, it will traverse the automaton and record all positions where matches occur. This approach ensures that we can find all possible matches, including those that overlap, without falling into the exponential blowup trap of backtracking.

The result is a table that maps each accepting state to a set of index pairs, each indicating the start and end of a successful match. For instance, applying this process to $R = (aa+aaa)(aaa+aa)$ and $w = aaaaabaaaaa$ will yield the following table:

Table 5.1: Match positions table using the regular expression $R = (aa + aaa)(aaa + aa)$ and input string $w = aaaaabaaaaa$

|   | $q_0$ | $q_1$ | $q_2$ | $q_3$ | $q_4$ | $q_5$ | $q_6$ | $q_7$ | $q_8$ | $q_9$ | $q_{10}$ | $q_{11}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\varepsilon$ | (0,0) | | | | | | | | | | | |
| a | (1,1) | (0,1) | (0,1) | | | | | | | | | |
| a | (2,2) | (1,2) | (1,2) | (0,2) | | | | | | | (0,2) | |
| a | (3,3) | (2,3) | (2,3) | (1,3) | (0,3) | (0,3) | (0,3) | | | | (1,3) | |
| a | (4,4) | (3,4) | (3,4) | (2,4) | (1,4) | (0,4) (1,4) | (0,4) (1,4) | (0,4) | (0,4) | | (2,4) | |
| a | (5,5) | (4,5) | (4,5) | (3,5) | (2,5) | (1,5) (2,5) | (1,5) (2,5) | (1,5) (0,5) | (1,5) (0,5) | (0,5) | (3,5) | (0,5) |
| b | (6,6) | | | | | | | | | | | |
| a | (7,7) | (6,7) | (6,7) | | | | | | | | | |
| a | (8,8) | (7,8) | (7,8) | (6,8) | | | | | | | (6,8) | |
| a | (9,9) | (8,9) | (8,9) | (7,9) | (6,9) | (6,9) | (6,9) | | | | (7,9) | |
| a | (10,10) | (9,10) | (9,10) | (8,10) | (7,10) | (6,10) (7,10) | (6,10) (7,10) | (6,10) | (6,10) | | (8,10) | |
| a | (11,11) | (10,11) | (10,11) | (9,11) | (8,11) | (7,11) (8,11) | (7,11) (8,11) | (6,11) (7,11) | (6,11) (7,11) | (6,11) | (9,11) | (6,11) |

First, we always have to account for $\varepsilon$, since there is the possibility of having the empty word and we also want to match against it. After that, every symbol $s \in w$ is processed sequentially.

At each step, the algorithm updates a row that maps the automaton's (represented in Figure 5.4) states to sets of position intervals $(i, j)$, such that the substring $w_{ij}$ corresponds to a valid match, whether it is overlapped or not.

Transitions are computed for each input symbol using the automaton's $\delta$ function. When a final state is reached, the interval is stored as a successful match. Furthermore, during this process, only the last symbol's computed transitions and position intervals are preserved because they always carry over to the current symbol's computation.

The automaton on Figure 5.4 shows that:

$$F = \{q_7, q_9\}$$

For those states, the resulting match table is represented by Table 5.2 and Table 5.3.

Table 5.2:  Highlighted substrings of valid matches for state $q_7$

| $w_{0,4}$ | **aaaaabaaaaa** |
|-----------|-----------------|
| $w_{0,5}$ | **aaaaabaaaaa** |
| $w_{1,5}$ | **aaaaabaaaaa** |
| $w_{6,10}$ | **aaaaabaaaaa** |
| $w_{6,11}$ | **aaaaabaaaaa** |
| $w_{7,11}$ | **aaaaabaaaaa** |

Table 5.3:  Highlighted substrings of valid matches for state $q_9$

| $w_{0,5}$ | **aaaaabaaaaa** |
|-----------|-----------------|
| $w_{6,11}$ | **aaaaabaaaaa** |

As seen in the tables, we are able to get all the matches possible (even overlapped ones) for the given input word $w$, along with their positions.

## 5.4    Implementation in FAdo

In order to support these new operations in *FAdo*, a new class was created. This new `matchNFA` class inherits from the `NFA` class. In it, several methods were implemented.

```python
class matchNFA(NFA):
    def __init__(self):
        NFA.__init__(self)


    def table_matcher(self, input_text: str):
        ...


    def enum_matches(self, match_table: dict, input_text: str):
        match_count = 0
        for entry in pos_tab:
            print(f"State {entry} matches ({len(set(pos_tab[entry]))}):")
            match_set = set(pos_tab[entry])
            match_count += len(match_set)
            for match in sorted(match_set):
                first_index = string[0:match[0]]
                second_index = string[match[1]:]
                print("Match:", colored(first_index, 'red') +
                    colored(string[match[0]:match[1]], 'green') +
                    colored(second_index, 'red'), "=>", match)
        print("Total matches:", match_count)


    def nfaPosCount(self: RegExp):
        ...


    setattr(RegExp, "nfaPosCount", nfaPosCount)
```

- `table_matcher` is the algorithm (can be seen in Table 2) that creates the matching table,

- `enum_matches` is the algorithm that displays the formatted matching table data,

- `nfaPosCount` is the modified position automaton builder (can be seen in Algorithm 1).

The `nfaPosCount` function is defined outside of `matchNFA`'s scope because, even though it can be fully integrated into `FAdo`, it still isn't. The `setattr` function is used to add this method to the `RegExp` class after it has been defined, so that it can later be used in other components of `FAdo`.

## 5.5   Limitations

There are some caveats regarding this modified position automaton construction and the table-based matcher.

### 5.5.1   Match Table Size

The match table can grow large. Practically, each input symbol will result in a new "step" in the table. Each step on the table will store data for multiple states (some can be empty, but they fill up continually). In the worst case, meaning that there are positions for each of the states, this will scale with the number of states. Ultimately, storing sets of intervals per state per position is done with space complexity of $O(m \cdot n)$, where $m$ is the total number of states and $n$ is the length of the input text.

### 5.5.2   Construction Size

Although the position automaton construction is linear in the number of symbol occurrences, bounded repetitions such as $a^{m,n}$ get unrolled into $a^m + a^{m+1} + \cdots + a^n$. Therefore, large bounds or nested counts will inflate the automaton and, therefore, slow down its construction. Having the regular expression $a^{1,50}$, for instance. Figure 5.5 represents a shortened version of this construction.

Figure 5.5: Abbreviated representation of the large modified position automaton.

Any automaton constructed for a regular expression that are constituted by a singe bounded repetition operator such as $a^{n,m}$ will have $\sum_{n=1}^{m-1} n + 2$ states. Therefore, the automaton for the expression $a^{1,50}$ contains 1227 states.

It is important to note that bounded repetitions are the worst case scenario for this construction. This blow up does not usually happen with the other operators.

# Chapter 6

# Conclusion

In this work, we contribute with a new approach to address the problem of Regular Expression Denial of Service (ReDoS), focusing on a solution outside of traditional backtracking engines. We explored the underlying causes of ReDoS vulnerabilities, particularly how certain finite automaton evaluators can lead to catastrophic performance, due to the deviation of the path of theory. To tackle this, we extended *FAdo* to support fixed and bounded repetition operators. To recognize this new grammar, the default Lark grammar present in *FAdo* was modified. The new type of automaton is capable of performing multiple overlapping matching. In security analysis, for instance, this means that one can find all possible malicious substrings in payloads (e.g., nested injection markers), because attack signatures can overlap and/or nest. Another example is the efficient search for DNA protein patterns, where biological sequences often contain repeating and/or nested patterns.

Most of the examples (regular expression pattern datasets) we found for this work did not include overlapping. The only way to do so (that we know of, currently implemented in state-of-the-art languages) is using `lookahead` assertions, and even this strategy does not work well, since they are expensive and hard to read. As a result, limited validation data exists for evaluating overlapping matchers, which makes empirical comparison with other state-of-the-art techniques challenging.

Anyhow, this work contributes towards safer and more predictable regular expression evaluation by combining formal methods and a new matching strategy.

Future work may turn towards making these algorithms more efficient by switching to a different language or integrating them into a larger ecosystem. It would also be valuable to develop tools capable of automatically generating text datasets with overlapping patterns for benchmarking and evaluation.

# Appendix A

# Regular Expression Grammar With Bounded and Fixed Repetition

```
start: rege

?rege: disjn
    ?disjn:  conjn
        | rege _UNIONT conjn -> disj

    ?conjn: shufflen
        |conjn _CONJT shufflen -> conj

    ?shufflen: concatn
        | shufflen _SHUFFLET concatn -> shuffle

    ?concatn: rep
        | concatn [_CONCATT] rep -> concat

    ?rep: lre
        | rep _START -> star
        | rep _PLUST -> plus
        | rep _OPTIONT  -> option
        | rep _CARET _LEFBR digit _RIGBR -> pow_min
        | rep _CARET _LEFBR digit _COMMA digit _LEFBR -> pow_minmax
        | rep _CARET _LEFBR digit _COMMA _INFTY _LEFBR -> pow_inf

    ?lre: base
        | _NOTT  lre  -> notn
        | _USHUFFLET  lre -> u_shuffle

    ?base: "(" rege ")" | symbol | epsilon | emptyset | sigmas | sigmap

_START: "*"
_SHUFFLET: ":"
_CONJT: "&"
_NOTT: "~"
```

```
_CONCATT: "."
_UNIONT: "+" | "|"
_OPTIONT: "-" | "?"
_TUPLET: "/"
_USHUFFLET: "!"
_CARET: "^"
_LEFBR: "["
_RIGBR: "]"
_INFTY: "..."
_COMMA: ","
_PLUST: "%"

digit: /[0-9]+/
symbol: /[a-zA-Z0-9]/

epsilon: "@epsilon"
emptyset: "@empty_set"
sigmap: "@sigmaP"
sigmas: "@sigmaS"

%ignore /[ \t\f\"]+/
```

Listing A.1: Lark grammar for regular expression parsing (with bounded/fixed repetition support)

# Appendix B

# Repetition Regular Expression Builder Cases

```python
def pow_min(self, s, inf=False):
    (arg, n_r) = s
    n = int(n_r.children[0].value)
    if inf:
        if n == 0:
            r = CStar(arg, sigma=self.sigma)
        elif n == 1:
            r = CPlus(arg, sigma=self.sigma)
        else:
            r = CCount(arg, n, -1, sigma=self.sigma)
    else:
        r = CPower(arg, n, sigma=self.sigma)

    return r
```

Listing B.1: pow_min: Fixed/Lower-bound-infinite Repetition Case

```python
def pow_minmax(self, s):
    (arg, n_mi, n_ma) = s
    n_min = n_mi.children[0].value
    n_max = n_ma.children[0].value

    r = CCount(arg, n_min, n_max, sigma=self.sigma)

    if n_max != "inf":
        if n_max == n_min:
            r = CAtom(arg, self.sigma)

    return r
```

Listing B.2: pow_minmax: Bounded Repetition Case

```python
def pow_inf(self, s):
    return self.pow_min(s, True)
```

Listing B.3: pow_inf: Lower-bound-infinite Repetition Case

# Appendix C

# Modified Position Automaton

---

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

---

**Require:** Regular expression $R$

**Ensure:** NFA $A$

 1: $A \leftarrow$ new empty NFA

 2: $i \leftarrow A.\text{addInitialState}()$

 3: $A.\text{addTransitionStar}(i, i)$                                              ▷ Accept any symbol from $\Sigma$

 4: $f_R \leftarrow R.\text{marked}()$

 5: stack $\leftarrow$ empty stack

 6: addedStates $\leftarrow$ empty map

 7: **for all** $p \in First(f_R)$ **do**

 8:     $q \leftarrow A.\text{addState}(p)$

 9:     addedStates$[p] \leftarrow q$

10:     stack.push$((p, q))$

11:     $A.\text{addTransition}(i, p, q)$

12: **end for**

13: **while** stack is not empty **do**

14:     $(s, s_{\text{idx}}) \leftarrow$ stack.pop()

15:     **for all** $t \in Follow(f_R, s)$ **do**

16:         **if** $t \in$ addedStates **then**

17:             $q \leftarrow$ addedStates$[t]$

18:         **else**

19:             $q \leftarrow A.\text{addState}(t)$

20:             addedStates$[t] \leftarrow q$

21:             stack.push$((t, q))$

22:         **end if**

23:         $A.\text{addTransition}(s_{\text{idx}}, t, q)$

24:     **end for**

25: **end while**

26: $e \leftarrow A.\text{addState}()$

27: $A.\text{addTransitionStar}(e, e)$

28: **for all** $p \in f_R.\text{Last}()$ **do**

29:     **if** $p \in$ addedStates **then**

30:         $A.\text{addFinal}(\text{addedStates}[p])$

31:         $A.\text{addTransitionStar}(\text{addedStates}[p], e)$

32:     **end if**

33: **end for**

---

# Appendix D

# Multiple Matching Evaluator

---

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

---

# Bibliography

[1] *regex — Rust crate documentation*, 2025. https://docs.rs/regex/latest/regex/.

[2] Valentin Antimirov. Partial derivatives of regular expressions and finite automaton constructions. *Theoretical Computer Science*, 155(2):291–319, 1996. ISSN: 0304-3975. doi:https://doi.org/10.1016/0304-3975(95)00182-4.

[3] S Broda, M Holzer, E Maia, N Moreira, and R Reis. A mesh of automata. *Information and Computation*, 265:94–111, 2019. doi:10.1016/j.ic.2019.01.003.

[4] Janusz A. Brzozowski. Derivatives of regular expressions. *J. ACM*, 11(4):481–494, October 1964. ISSN: 0004-5411. doi:10.1145/321239.321249.

[5] Cezar Câmpeanu, Kai Salomaa, and Sheng Yu. Regex and extended regex. In Jean-Marc Champarnaud and Denis Maurel, editors, *Implementation and Application of Automata*, pages 77–84, Berlin, Heidelberg, 2003. Springer Berlin Heidelberg. ISBN: 978-3-540-44977-5.

[6] Russ Cox. Regular expression matching in the wild, 2010.

[7] Free Software Foundation. grep, .

[8] Free Software Foundation. sed, .

[9] *GNU Awk User's Guide*. Free Software Foundation, 2024.

[10] Google. Re2, a regular expression library.

[11] *The Open Group Base Specifications Issue 8*. The IEEE and The Open Group, 2024.

[12] Zahir Irani, Raul M. Abril, Vishanth Weerakkody, Amizan Omar, and Uthayasankar Sivarajah. The impact of legacy systems on digital transformation in european public administration: Lesson learned from a multi case analysis. *Government Information Quarterly*, 40(1):101784, 2023. ISSN: 0740-624X. doi:https://doi.org/10.1016/j.giq.2022.101784.

[13] Yeting Li, Zixuan Chen, Jialun Cao, Zhiwu Xu, Qiancheng Peng, Haiming Chen, Liyuan Chen, and Shing-Chi Cheung. ReDoSHunter: A combined static and dynamic approach for regular expression DoS detection. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 3847–3864. USENIX Association, August 2021. ISBN: 978-1-939133-24-3.

[14] Yinxi Liu, Mingxue Zhang, and Wei Meng. Revealer: Detecting and exploiting regular expression denial-of-service vulnerabilities. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 1468–1484, 2021. doi:10.1109/SP40001.2021.00062.

[15] University of Cambridge. pcre2.

[16] Rogério Reis and Nelma Moreira. Fado: tools for finite automata and regular expressions manipulation. 11 2002.

[17] Yuju Shen, Yanyan Jiang, Chang Xu, Ping Yu, Xiaoxing Ma, and Jian Lu. Rescue: crafting regular expression dos attacks. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, ASE '18, page 225–235, New York, NY, USA, 2018. Association for Computing Machinery. ISBN: 9781450359375. doi:10.1145/3238147.3238159.

[18] Erez Shinan. Lark — a parsing toolkit for python, 2017. Version 1.x, Available at: https://github.com/lark-parser/lark.

[19] Ken Thompson. Programming techniques: Regular expression search algorithm. *Commun. ACM*, 11(6):419–422, June 1968. ISSN: 0001-0782. doi:10.1145/363347.363387.

[20] Ian Erik Varatalu, Margus Veanes, and Juhan Ernits. Re#: High performance derivative-based regex matching with intersection, complement, and restricted lookarounds. *Proceedings of the ACM on Programming Languages*, 9(POPL):1–32, January 2025. ISSN: 2475-1421. doi:10.1145/3704837.

[21] Xiang Wang, Yang Hong, Harry Chang, KyoungSoo Park, Geoff Langdale, Jiayu Hu, and Heqing Zhu. Hyperscan: A fast multi-pattern regex matcher for modern CPUs. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, pages 631–648, Boston, MA, February 2019. USENIX Association. ISBN: 978-1-931971-49-2.

[22] Xiang Wang, Yang Hong, Harry Chang, KyoungSoo Park, Geoff Langdale, Jiayu Hu, and Heqing Zhu. Hyperscan: A fast multi-pattern regex matcher for modern CPUs. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, pages 631–648, Boston, MA, February 2019. USENIX Association. ISBN: 978-1-931971-49-2.

[23] Isaac Z. Schlueter. minimatch, 2011.