

Progress Report - 9 Months

Meshal Binnasban

November 19, 2025

Abstract

This is the nine-month progress report for my PhD at King's College London. It summarises the motivation, challenges, and progress made in developing a regular expression matcher based on marked regular expressions, with the goal of POSIX value extraction. The report first reviews regular expression matching using derivatives, then reviews marked regular expression matchers, which presently provide matching but not POSIX value extraction—our primary focus. Our work developed several versions of the marked algorithm in Scala. We have also begun formally proving the correctness of one of our versions.

Synopsis

This research investigates the marked approach-based method for regular expression matching, with the goal of extending it to provide POSIX value extraction.

Derivative-based methods, though elegant, suffer from size explosion. Sulzmann and Lu [8] extended derivatives with bitcodes to record lexing information, but the size explosion persists. The marked approach, described in the works of Fischer et al. [3] and Asperti et al. [1], offers an alternative by propagating marks directly through the expression. These works provide matchers only, without value extraction.

Our work explores how the marked approach can be extended to recover POSIX values. We have implemented several versions of the marked algorithm in Scala, making use of bitcoded annotations to track lexing values. Challenges in handling constructs such as sequences and repetitions required refining the algorithm through successive versions, each improving on the last. Large-scale testing against a derivative-based reference matcher has been used to confirm correctness and uncover edge cases.

The long-term aim of the project is to establish a marked approach-based matcher that consistently yields the POSIX value for all regular expressions, and to formally prove its correctness. Future work also includes extending the matcher to additional operators such as intersection and negation.

Contents

1	Introduction	4
1.1	Overview and Motivation	4
1.2	Background	4
2	Derivatives	5
2.1	Size Explosion	7
2.2	Derivative Extension	11
3	Marked Approach	14
3.1	Motivation for a Marked Approach	15
4	Our Approach	16
4.1	String-Carrying Marks - Matcher	17
5	Future Work	24
	References	25
A	Appendix	25
A.1	Scala Implementation of Fischer et al.'s Marked Approach . .	25
A.2	Bit-Annotated Version 1	27
A.3	Bit-Annotated Version 2	33
A.4	String-Carrying Marks — Lexer (so far)	34

1 Introduction

1.1 Overview and Motivation

The notion of derivatives in regular expressions was introduced by Brzozowski in 1964 [2], but it was rediscovered in the 2000s and gained renewed attention—for example in the work of Owens [7] and Might [4]. This renewed interest is largely due to the method’s conceptual simplicity and the ease with which it can be implemented in functional programming languages.

However, derivatives suffer from “size explosion”, since taking derivatives may increase the size of subexpressions. This means the algorithm has to traverse larger and larger regular expressions, which results in slower performance.

In contrast, the approach based on marks leaves regular expressions unchanged during matching but moves annotations through them. At present, only matching algorithms using this approach exist, while our work aims to extend them to provide POSIX value extraction—the value representing the longest leftmost match.

This report first reviews Brzozowski’s derivatives and the bit-coded variant of Sulzmann and Lu [8], followed by background on the regular expression matchers based on marks described by Fischer et al. [3] and Asperti et al. [1]. We then present our own work, including several versions of matchers based on marks developed during the first nine months.

1.2 Background

Regular expressions are a way to describe languages of strings over an alphabet. They provide a declarative way to specify such sets of strings. Each regular expression r is associated with a language $L(r)$, which is the set of strings that match r . The regular expressions we consider are the standard regular expression constructors, including bounded repetitions, namely:

$$r ::= \mathbf{0} \mid \mathbf{1} \mid c \mid r_1 + r_2 \mid r_1 \cdot r_2 \mid r^* \mid r^n$$

We will use $\mathbf{0}$ for the regular expression that cannot match any string, $\mathbf{1}$ for the regular expression that can match the empty string, and c for the regular expression that can match a character from the alphabet. The operator $+$ is for alternation of r_1 and r_2 , while the dot \cdot is for concatenation or sequencing. The Kleene star r^* is for arbitrary repetition of r , and the

power r^n is for exactly n repetitions of r . The meaning of these regular expressions is given by the language function L (Figure 1).

$L(\mathbf{0})$	$\stackrel{\text{def}}{=}$	$\mathbf{0}$
$L(\mathbf{1})$	$\stackrel{\text{def}}{=}$	$\{[]\}$
$L(c)$	$\stackrel{\text{def}}{=}$	$\{[c]\}$
$L(r_1 + r_2)$	$\stackrel{\text{def}}{=}$	$L(r_1) \cup L(r_2)$
$L(r_1 \cdot r_2)$	$\stackrel{\text{def}}{=}$	$L(r_1) \cdot L(r_2)$
$L(r^*)$	$\stackrel{\text{def}}{=}$	$(L(r))^*$
$L(r^n)$	$\stackrel{\text{def}}{=}$	$(L(r))^n$

Figure 1: the function L , which gives meaning to a regular expression, i.e. the set of strings that belong to its language.

For example, if $r = a + b$, then

$$L(r) = L(a) \cup L(b) = \{a\} \cup \{b\} = \{a, b\}.$$

Similarly, $L(a \cdot b) = \{ab\}$, and $L(a^*)$ is the set of all strings consisting of zero or more a 's.

2 Derivatives

Brzozowski's derivatives offer an elegant way for regular expression matching (Figure 2). What is beautiful about derivatives is that they can be easily implemented in a functional programming language and are easy to reason about in theorem provers. By successively taking the derivative of a regular expression with respect to each input character, one obtains a sequence of derivatives. If the final derivative can match the empty string, then the original input is accepted. For example, to decide whether a string $a_1 a_2 \dots a_n$ is matched by a regular expression r , we successively compute derivatives:

$$r_0 = r, \quad r_1 = \text{der}_{a_1}(r_0), \quad r_2 = \text{der}_{a_2}(r_1) \dots, r_n = \text{der}_{a_n}(r_{n-1}).$$

The string matches r if and only if the final expression r_n accepts the empty string, which can be easily checked by a separate function. Here $der_a(r)$ stands for the derivative of r with respect to the character a .

$$\begin{aligned}
der_a(\mathbf{0}) &\stackrel{\text{def}}{=} \mathbf{0} \\
der_a(\mathbf{1}) &\stackrel{\text{def}}{=} \mathbf{0} \\
der_a(c) &\stackrel{\text{def}}{=} \begin{cases} \mathbf{1} & \text{if } a = c \\ \mathbf{0} & \text{if } a \neq c \end{cases} \\
der_a(r_1 + r_2) &\stackrel{\text{def}}{=} der_a(r_1) + der_a(r_2) \\
der_a(r_1 \cdot r_2) &\stackrel{\text{def}}{=} der_a(r_1) \cdot r_2 + \begin{cases} der_a(r_2) & \text{if } \varepsilon \in L(r_1) \\ \mathbf{0} & \text{otherwise} \end{cases} \\
der_a(r^*) &\stackrel{\text{def}}{=} der_a(r) \cdot r^*
\end{aligned}$$

Figure 2: Brzozowski's Derivative.

To illustrate how derivatives can be used to match a regular expression against a string, consider the regular expression $(ab + ba)$. The derivative can check the matching of the string \mathbf{ba} by taking the derivative of the regular expression with respect to b , then with respect to a (see Figure 3). Since the final derivative expression can match the empty string with no input characters remaining, this confirms that the string \mathbf{ba} belongs to the language of the regular expression.

Note, however, that some subexpressions that arise during the computation of derivatives may be redundant. For example, $\mathbf{0} \cdot b$ cannot match any string, so we could simplify it to $\mathbf{0}$, which can match the same strings—namely none (see Figure 3). Other simplification rules may also be applied. These include associativity $(r + s) + t \equiv r + (s + t)$, commutativity $r + s \equiv s + r$, and idempotence $r + r \equiv r$. Rules for $\mathbf{0}$ and $\mathbf{1}$ can also be applied; e.g. $\mathbf{0} \cdot r \equiv r \cdot \mathbf{0} \equiv \mathbf{0}$, $r + \mathbf{0} \equiv r$, and $r \cdot \mathbf{1} \equiv \mathbf{1} \cdot r \equiv r$.

These simplifications preserve the language of the regular expression while reducing its overall size. Such simplifications are sometimes necessary in the derivative method to keep the size manageable, but even with a best effort, the size can still grow significantly.

$$\begin{aligned}
der_b r &= der_b(ab + ba) \\
&= der_b(ab) + der_b(ba) \\
&= (der_b a) \cdot b + (der_b b) \cdot a \\
&= \mathbf{0} \cdot b + \mathbf{1} \cdot a \\
\\
der_a(der_b r) &= der_a(\mathbf{0} \cdot b + \mathbf{1} \cdot a) \\
&= der_a(\mathbf{0} \cdot b) + der_a(\mathbf{1} \cdot a) \\
&= der_a(\mathbf{0}) \cdot b + (der_a(\mathbf{1}) \cdot a + der_a a) \\
&= \mathbf{0} \cdot b + (\mathbf{0} \cdot a + \mathbf{1})
\end{aligned}$$

Figure 3: Derivatives of $(ab + ba)$ with respect to string ba .

2.1 Size Explosion

Although the construction of derivatives is elegant, it may duplicate large parts of the expression and increase the size of the regular expression in intermediate steps (see Figure 2, the sequence case can create additional subexpressions, as can the star case). Sulzmann and Lu [8] also note this problem, describing it as the well-known issue that the size and number of derivatives may explode. The size of derivatives can become very large even for simple expressions. This happens because each derivative is usually more complicated than the original expression, and even when standard simplifications are applied, they may not remove all duplications. Consider, for example, the regular expression $(a + aa)^*$ and a derivative with respect to a single a .

$$\begin{aligned}
der_a(a + aa)^* &= (der_a(a + aa)) \cdot (a + aa)^* \\
&= (der_a a + der_a(aa)) \cdot (a + aa)^* \\
&= (\mathbf{1} + der_a(a \cdot a)) \cdot (a + aa)^* \\
&= (\mathbf{1} + (\mathbf{1} \cdot a)) \cdot (a + aa)^* \\
&= (\mathbf{1} + a) \cdot (a + aa)^*
\end{aligned}$$

Each derivative step may introduce new structure, unfolding all possible ways in which r^* can match the input.

It is already known that the simplification reduces the number of gen-

erated expressions and helps to provide a finite bound on the number of intermediate expressions [8, 9], but it does not completely solve the underlying problem. Even with aggressive simplifications—as in Sulzmann and Lu’s bit-coded approach, and in the variant described by Tan and Urban—the resulting derivatives can still be large even if finitely bounded. For example, expressions of the form $((a)^* + (aa)^* + (aaa)^* + \dots)^*$ cannot be simplified completely, since the duplicates occur at different levels of the derivatives and across sequences and alternatives. Thus, even under aggressive simplification, the size of the derivatives remains large and continues to grow across derivative steps. If we consider the regular expression

$$((a)^* + (aa)^* + (aaa)^* + (aaaa)^* + (aaaaa)^*)^*$$

with a small input string of the form $\underbrace{a \dots a}_n$, each derivative step introduces subexpressions that are copies of the original, as shown in Figure 4.

Figure 5 illustrates the “size explosion” of derivatives by showing the runtime difference between derivatives with aggressive simplifications and Fischer’s marked algorithm. We conducted a naive test using a regular expression of the kind mentioned above to demonstrate the impact of the “size explosion” of derivatives, even under aggressive simplifications. Although the test is naive, it clearly shows a considerable difference in behaviour. The test was designed for inputs of up to 100,000 characters; however, the derivatives implementation ran out of memory after 10,000, with the last successful case taking almost 22 seconds, while the marked algorithm finished executing all inputs, with each case under 0.0035 seconds. Even Antimirov’s partial derivatives, which do not track POSIX values, may still produce cubic growth in the worst case: there can be linearly many distinct partial derivatives, and each derivative can already be quadratic in the size of the original expression, resulting in an overall cubic upper bound.

$$\begin{aligned}
der_a r &= der_a (((a)^* + (aa)^* + (aaa)^* + (aaaa)^* + (aaaaa)^*)^*) \\
&= der_a ((a)^* + (aa)^* + (aaa)^* + (aaaa)^* + (aaaaa)^*) \cdot r^* \\
&= ((der_a (a)^*) + (der_a (aa)^*) + (der_a (aaa)^*) + \dots) \cdot r^* \\
&= ((der_a a \cdot (a)^*) + (der_a aa \cdot (aa)^*) + (der_a aaa \cdot (aaa)^*) \\
&\quad + \dots) \cdot r^* \\
&= ((\mathbf{1} \cdot (a)^*) + ((\mathbf{1} \cdot a) \cdot (aa)^*) + ((\mathbf{1} \cdot aa) \cdot (aaa)^*) \\
&\quad + \dots) \cdot r^* \\
&= ((a)^* + (a \cdot (aa)^*) + (aa \cdot (aaa)^*) + \dots) \cdot r^* \\
\\
der_a (der_a r) &= der_a ((a)^* + (a \cdot (aa)^*) + (aa \cdot (aaa)^*) + \dots) \cdot r^* \\
&\quad + der_a r^* \\
&= (\mathbf{1} \cdot (a)^* + (\mathbf{1} \cdot (aa)^*) + ((\mathbf{1} \cdot a) \cdot (aaa)^*) + \dots) \cdot r^* \\
&\quad + ((a)^* + (a \cdot (aa)^*) + (aa \cdot (aaa)^*) + \dots) \cdot r^* \\
&= ((a)^* + (aa)^* + (a \cdot (aaa)^*) + \dots) \cdot r^* \\
&\quad + ((a)^* + (a \cdot (aa)^*) + (aa \cdot (aaa)^*) + \dots) \cdot r^* \\
\\
der_a (der_a (der_a r)) &= der_a (((a)^* + (aa)^* + (a \cdot (aaa)^*) + \dots) \cdot r^* \\
&\quad + der_a (((a)^* + (a \cdot (aa)^*) + (aa \cdot (aaa)^*) + \dots) \cdot r^*) \\
&= (der_a (((a)^* + ((aa)^*) + (a \cdot (aaa)^*) + \dots) \cdot r^*) \\
&\quad + der_a r^*) + (der_a (((a)^* + (a \cdot (aa)^*) + (aa \cdot (aaa)^*) \\
&\quad + \dots) \cdot r^*) + der_a r^*)
\end{aligned}$$

Figure 4: Derivatives of $((a)^* + (aa)^* + (aaa)^* + (aaaa)^* + (aaaaa)^*)^*$ with respect to the string **aaa**. Each step introduces additional copies of the original regular expression at different levels of the derivative, where simplifications cannot be applied. For example, after the second step, $(a)^*$ appears in both alternatives, but, since they occur under different concatenations with r , the two occurrences cannot be merged.

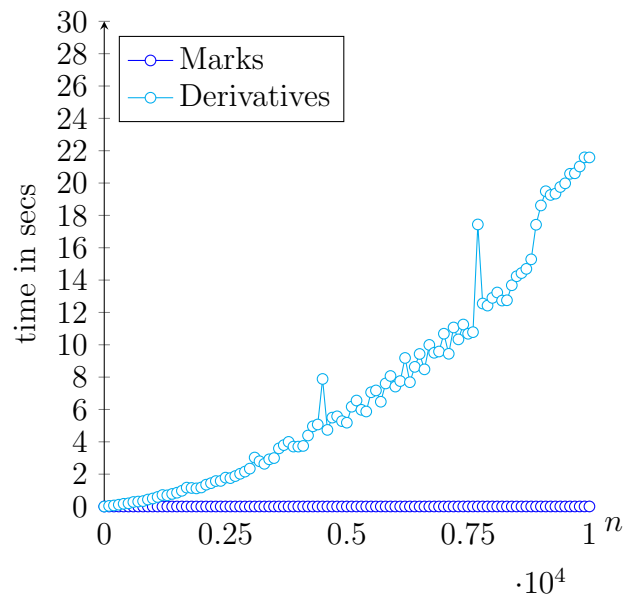


Figure 5: Runtime comparison of derivative and marked regular expression matchers for $((a)^* + (aa)^* + (aaa)^* + (aaaa)^* + \dots)^*$ with input strings of length n .

2.2 Derivative Extension

Sulzmann and Lu [8] extend Brzozowski’s derivatives to produce POSIX values in addition to deciding whether a match exists. These values record how the match occurred, for example by showing which part of the regular expression matched the input string, whether the left or right branch was taken, and how sequences and Kleene stars were matched.

Sulzmann and Lu provide two variants: a bit-coded construction and an injection-based construction [8]. In the bit-coded variant, bitsequences encode lexing choices during derivative construction and, after acceptance, are decoded to the value. In the injection-based variant, an *inj* function “injects back” the consumed characters into the value; the injection function essentially reverts the derivative steps to obtain a POSIX value. We focus here on the bit-coded variant, which more directly inspired our marked approach.

Bitsequences are sequences over $\{0, 1\}$ that encode the choices made in a match. These bits record branch choices in alternatives and repetitions during matching. Figure 6 illustrates how bitsequences work with the regular expression $(a + ab)(b + \mathbf{1})$ on the string **ab**, proceeding step by step by constructing derivatives and tracking how the bitsequence grows. Initially, the algorithm internalizes the regular expression by applying the function *intern* (Figure 7), adding bit annotations to the alternative constructors in the expression: 0 for left branches and 1 for right branches [8]. This annotation process is implemented by the function *fuse* (Figure 8). The bitsequence is updated while taking the derivatives to reflect the choices made during matching. The function *mkeps* (Figure 9), as defined by Sulzmann and Lu [8], extracts the bitsequence from nullable derivatives. A regular expression is nullable if it can match the empty string.

Derivatives can cause the regular expression to expand. In particular, a sequence case expands into an alternative when the first part is nullable, as shown in Figure 6, Step 3, where the concatenation becomes nullable. This accounts for the situation in which the first part, r_1 , may be skipped during matching: the left branch assumes r_1 is not skipped, while the right branch assumes it is skipped and instead takes the derivative of r_2 directly. This behaviour also illustrates why derivatives tend to increase in size. Sulzmann and Lu [8] use $\text{fuse}(\text{mkeps}(r_1))$ to include the bit annotations (bitsequences) needed when r_1 is skipped; these bits, extracted by *mkeps*, indicate how r_1 matched the empty string.

After taking the derivatives with respect to the entire input string, the

1. Step 1: Internalising the regular expression

$$r' = \text{intern}(r) = ({}_0a + {}_1ab) \cdot ({}_0b + {}_1\mathbf{1})$$

2. Step 2: input a

$$\begin{aligned} \text{der}_a(r') &= \text{der}_a(({}_0a + {}_1ab) \cdot ({}_0b + {}_1\mathbf{1})) \\ &= \text{der}_a({}_0a + {}_1ab) \cdot ({}_0b + {}_1\mathbf{1}) \\ &= (\text{der}_a({}_0a) + \text{der}_a({}_1ab)) \cdot ({}_0b + {}_1\mathbf{1}) \\ &= ({}_0\mathbf{1} + \text{der}_a({}_1a) \cdot b) \cdot ({}_0b + {}_1\mathbf{1}) \\ &= ({}_0\mathbf{1} + {}_1\mathbf{1} \cdot b) \cdot ({}_0b + {}_1\mathbf{1}) \end{aligned}$$

3. Step 3: input b

$$\begin{aligned} \text{der}_b(\text{der}_a(r')) &= \text{der}_b(({}_0\mathbf{1} + {}_1\mathbf{1} \cdot b) \cdot ({}_0b + {}_1\mathbf{1})) \\ &= \text{der}_b({}_0\mathbf{1} + {}_1\mathbf{1} \cdot b) \cdot ({}_0b + {}_1\mathbf{1}) \\ &\quad + \text{der}_b(\text{fuse}(\text{mkeps}(r_1), ({}_0b + {}_1\mathbf{1}))) \\ &= (\text{der}_b({}_0\mathbf{1}) + \text{der}_b({}_1\mathbf{1} \cdot b)) \cdot ({}_0b + {}_1\mathbf{1}) \\ &\quad + {}_0(\text{der}_b({}_0b) + \text{der}_b({}_1\mathbf{1})) \\ &= (\mathbf{0} + \text{der}_b({}_1\mathbf{1} \cdot b)) \cdot ({}_0b + {}_1\mathbf{1}) + {}_0({}_0\mathbf{1} + \mathbf{0}) \\ &= (\text{der}_b({}_1\mathbf{1}) \cdot b + \text{der}_b({}_1b)) \cdot ({}_0b + {}_1\mathbf{1}) + {}_0({}_0\mathbf{1} + \mathbf{0}) \\ &= (\mathbf{0} \cdot b + {}_1\mathbf{1}) \cdot ({}_0b + {}_1\mathbf{1}) + {}_0({}_0\mathbf{1} + \mathbf{0}) \end{aligned}$$

Figure 6: Step-by-step construction of derivatives with bit-sequences for $(a + ab) \cdot (b + \mathbf{1})$ with input **ab**. Step 1 shows the result of internalization using *fuse*, where alternatives are annotated with bits. Step 2 takes the derivative with respect to *a*; the two branches of the result are nullable at this stage: the left branch of r_1 and the right branch of r_2 can both match the empty string after consuming the first character, corresponding to the bitsequences [0] and [1] respectively. Step 3 takes the derivative with respect to *b*, expanding the expression due to the nullable concatenation. Finally, *mkeps* extracts the bitsequence [1, 1].

algorithm checks whether the result is nullable. If so, it calls *mkeps* to extract the bitsequence indicating how the match was obtained. In Figure 6, there

are two possible matches, but the algorithm prefers the left one, $[1, 1]$, which encodes the choice of the right branch in the alternative of r_1 , matching the string **ab**, followed by the right branch in the alternative of r_2 , which matches the empty string. Decoding this against the original expression yields the POSIX value

$$Seq(Right(Seq(a, b)), Right(Empty))$$

$intern(\mathbf{0})$	$\stackrel{\text{def}}{=}$	$\mathbf{0}$
$intern(\mathbf{1})$	$\stackrel{\text{def}}{=}$	$\mathbf{1}$
$intern(c)$	$\stackrel{\text{def}}{=}$	c
$intern(r_1 + r_2)$	$\stackrel{\text{def}}{=}$	$fuse([0], intern(r_1)) + fuse([1], intern(r_2))$
$intern(r_1 \cdot r_2)$	$\stackrel{\text{def}}{=}$	$intern(r_1) \cdot intern(r_2)$
$intern(r^*)$	$\stackrel{\text{def}}{=}$	$(intern(r))^*$

Figure 7: The *intern* function, which internalises regular expressions by adding bit annotations to alternatives, as defined by Sulzmann and Lu [8].

$fuse\ bs\ (c_{bs'})$	$\stackrel{\text{def}}{=}$	$c\ (bs\ @\ bs')$
$fuse\ bs\ ((r_1 + r_2)_{bs'})$	$\stackrel{\text{def}}{=}$	$(r_1 + r_2)_{(bs\ @\ bs')}$
$fuse\ bs\ ((r_1 \cdot r_2)_{bs'})$	$\stackrel{\text{def}}{=}$	$(r_1 \cdot r_2)_{(bs\ @\ bs')}$
$fuse\ bs\ (r^*_{bs'})$	$\stackrel{\text{def}}{=}$	$(r^*)_{(bs\ @\ bs')}$
$fuse\ bs\ (r^n_{bs'})$	$\stackrel{\text{def}}{=}$	$(r^n)_{(bs\ @\ bs')}$

Figure 8: The *fuse* function, which adds bit annotations to expressions, as defined by Sulzmann and Lu [8].

$$\begin{aligned}
mkeps(\mathbf{1}_{bs}) &\stackrel{\text{def}}{=} bs \\
mkeps((r_1 + r_2)_{bs}) &\stackrel{\text{def}}{=} \begin{cases} bs @ mkeps(r_1) & \text{if } nullable(r_1) \\ bs @ mkeps(r_2) & \text{otherwise} \end{cases} \\
mkeps((r_1 \cdot r_2)_{bs}) &\stackrel{\text{def}}{=} bs @ mkeps(r_1) @ mkeps(r_2) \\
mkeps((r^*)_{bs}) &\stackrel{\text{def}}{=} bs @ [1]
\end{aligned}$$

Figure 9: The *mkeps* function, which extracts bit annotations from nullable expressions, as defined by Sulzmann and Lu [8].

3 Marked Approach

The marked approach is a method for regular expression matching that tracks matching progress by inserting and moving marks into the expression. As noted by Nipkow and Traytel [5], the idea can be traced to the earlier work, in particular to Yamada and Glushkov, who identify positions in regular expressions by marking their atoms. Nipkow and Traytel cite Fischer et al., [3] and Asperti et al. [1] as reviving this work and developing it in a modern setting.

This approach aims to enable more efficient matching, particularly in complex expressions such as $(a^* \cdot b^*)^*$, in large alternations like $(a + b + c + \dots + z)^*$, or in nested choices such as $((a)^* + (aa)^* + (aaa)^* + (aaaa)^* + (aaaaa)^*)^*$. The regular expression itself does not increase in size; rather, progress is tracked through the placement of marks. With each input character, these marks traverse the regular expression according to a set of rules. At the end of the input, the expression is evaluated to determine whether the marks are in positions that make the expression “final”—that is, whether the expression accepts the string.

There is a slight difference in how marks are interpreted in the works of Fischer et al. [3] and Asperti et al. [1]. In Fischer et al.’s approach, marks are inserted after a character has been matched, thereby recording the matched character or subexpression. In contrast, Asperti et al. interpret the positions of marks as indicating the potential to match a character: as input is consumed, the marks move through the regular expression to indicate

the next character that can be matched. In both approaches, acceptance is determined from the final marked structure of the expression. For Fischer et al., this corresponds to having matched characters, whereas for Asperti et al. it requires that the marks end in positions where the expression can accept the empty string, ensuring that the entire input has been consumed. If the marks do not reach such positions, some characters remain unmatched and the matching problem fails.

Nipkow and Traytel, in their work on regular expression equivalence, build a unified framework for decision procedures on regular expressions. They describe McNaughton–Yamada–Glushkov and Fischer’s work as mark-after-atom, meaning that when a character is read, the mark is placed immediately after it. In their framework, this is expressed by two operations on marks, namely *follow* and *read*. The function *follow* is, as they describe, similar to an epsilon-closure, in that it moves all marks in a regular expression to the next atom to be read, and then *read* marks all atoms matching the current input symbol. The main transition function is therefore written as

$$\textit{shift } m \ r \ x = \textit{read } x \ (\textit{follow } m \ r).$$

By contrast, Nipkow and Traytel show that the work of Asperti et al., although presented as McNaughton–Yamada, is in fact a dual construction: mark-before-atom. While Asperti et al. described their approach as following the McNaughton–Yamada framework, their work actually corresponds to the opposite variant, where marks are placed before reading the next atom. In this version, the marks indicate atoms that are about to be read rather than atoms that have just been consumed. The main transition function is therefore the reverse of the previous one:

$$\textit{move } c \ r \ m = \textit{follow } m \ (\textit{read } c \ r).$$

3.1 Motivation for a Marked Approach

The marked approach offers a promising alternative to derivatives for regular expression matching by addressing some of the limitations found in the derivative method. The derivative method provides a more direct way of matching compared to constructing an NFA and then a DFA, which can suffer from exponential state explosion, up to 2^n . Engines that construct DFAs, such as those used in Rust, attempt to mitigate this problem by building the

automaton dynamically: they start from the initial state and construct additional DFA states only as required by the input string. This avoids generating the full automaton up front, but even so, it can still lead to large memory usage for expressions such as r^n , which can effectively require n copies of the automaton in the worst case. Other engines, such as those in Python, use NFAs instead of DFAs. This design allows for additional features, such as backreferences and lookarounds, but may result in catastrophic backtracking when the engine needs to explore a large portion of the automaton—or even the entire automaton—before finding a match or rejecting an input. The derivative approach, albeit elegant in design and supported by formal proofs—unlike many existing engines that lack such guarantees—can still grow quickly in size, as each derivative may add new subexpressions at each step. Even with aggressive simplifications, these subexpressions cannot always be simplified when they occur at different levels of the expression tree.

The marked approach relies on propagating marks within the regular expression rather than constructing new subexpressions. Our main motivation is that this method could support efficient and high-performance matching with POSIX value extraction, since it handles matching in a way that avoids some of the limitations of the derivative-based approach. Inspired by the works of Fischer et al. and Asperti et al. [3, 1], we aim to extend the marked approach to extract POSIX values, as well as to handle extended constructors such as bounded repetitions and intersections. Our work is primarily based on the algorithm described by Fischer et al. [3]. As shown by Nipkow and Traytel [5], the pre-mark algorithm of Asperti et al. is a dual construction to the post-mark algorithm of Fischer et al.; moreover, they prove that the mark-before-atom automaton is a quotient of the mark-after-atom automaton. Consequently, we adopt Fischer et al.’s approach as the foundation for extending the marked approach to matching with POSIX value extraction.

4 Our Approach

We began our work by implementing the marked approach described by Fischer et al. [3] in Scala; details are provided in Appendix A.1. This initial implementation of the algorithm included only acceptance checking, without any value construction. Over the course of this work, we developed several versions of the algorithm, each addressing specific challenges.

The first version extended the marked approach with bit annotations,

producing values but not always the POSIX value. This occurs when marks are overwritten after a character is matched more than once; further details are provided in Appendix A.2.

The second version was developed to address two main issues: the overwriting of marks and the absence of a mechanism to order them so as to always produce the POSIX value. Initially, we modified the previous version to accumulate all possible paths to a match, retaining every bitsequence that could lead to acceptance. Through reasoning and testing, we refined so that it seems that with every input it produces the possible values that the input can match the regular expression. One resulting difficulty was the proliferation of matches in the *Star* case, where every possible way of matching was generated. Another difficulty, which also existed in the first version, was the absence of an ordering for marks during the shifting process. We implemented an ordering after shifting to evaluate the results of the algorithm, following the work of Okui and Suzuki [6], although this ordering is not embedded within the shifting process itself. Instead, all possible values are generated first, and then the ordering function is applied to select the POSIX value; further details are provided in Appendix A.3.

This line of thought eventually led us to the third version of the algorithm, which uses string-annotated marks. In this version, each mark carries the suffix of the input string that is still to be matched. This allows us to compare marks based on their remaining suffixes. We also began annotating these marks with bitsequences to record the choices made during matching. Although this version is not yet complete, the initial results are promising. The remainder of this section describes the third (latest) version; the earlier versions are presented in the Appendix.

4.1 String-Carrying Marks - Matcher

We initially followed the approach of Fischer et al. [3] and Asperti et al. [1], in which shifting is performed character by character and marks are propagated at each step. However, as in our previous versions, this method made it difficult to maintain an ordering on marks and then to extract the POSIX value.

To address this, we instead let each mark carry the suffix of the input string that is still to be matched, starting with an initial mark carrying the full input string. This gives a more straightforward way of ordering the marks, based on the remaining suffixes. It also appears to offer a closer

correspondence with derivatives, in that each mark could be associated with subexpressions created by the derivative, though this point is not yet fully established. The trade-off is that more information must be stored. However, with string-carrying marks, we currently have a better handle on extracting POSIX values.

We define *shifts* in place of *shift*, as the function now operates on the full input string rather than character by character, as in the previous versions. The full definition of *shifts* is given in Figure 10. As discussed in Appendix A.1, the empty-string expression **1** is not marked. This choice ensures termination of the *shifts* function, particularly in the *Star* case: if **1** were allowed to produce marks, it would lead to infinite unfolding and non-termination.

The shifting process starts with a list of marks containing one initial mark carrying the full input string. As the mark is shifted through the expression, new marks may be produced depending on its structure, such as in sequences, stars, or alternatives. For example, in the case of an alternative, the mark is passed to each branch, and the resulting lists of marks are combined. Each time a character matches, it is stripped from the corresponding string carried by the mark. If a character does not match the prefix of the string carried by a mark, that mark is deleted. A match occurs at the end of shifting when an empty mark has been produced, indicating that all characters have been matched. When the input string itself is empty, *shifts* is not applied; instead, the *nullable* function (with the same definition as in Section A.1, shown in Figure 16) is used to check whether the regular expression accepts the empty string.

The number of marks when using lists is, in most cases, linear in the input length. The exception is the *Star* case, where the total number of marks can grow exponentially or more. This proliferation arises when the inner *r* contains alternatives and sequences that may add or combine lists of marks, which are then fed back into the *Star* through repeated unfoldings. Even though the length of individual marks becomes shorter with each successful shift, their number still grows in this way. In such cases, if the regular expression has m branch choices, then on an input of length n , the number of marks can reach $\mathcal{O}(m^n)$. For example, for the regular expression $(a + a)^*$ and an input of the form $\underbrace{a \cdots a}_n$, we have $m = 2$. For $n = 3$, the first three

unfoldings are as follows, where @ stands for list concatenation.

$$\begin{aligned}
\text{step 1 : } & [\bullet aa] @ [\bullet aa] \quad (2) \\
\text{step 2 : } & [\bullet a, \bullet a] @ [\bullet a, \bullet a] \quad (4) \\
\text{step 3 : } & \underbrace{[\bullet [], \dots] @ [\bullet [], \dots]}_{8 \text{ marks}}
\end{aligned}$$

Using sets, on the other hand, removes duplicates, which are the reason lists can grow exponentially (or more). The number of distinct marks therefore remains linear in the input string length. In the same example, we obtain only $\{\bullet aa \rightarrow \bullet a \rightarrow \bullet []\}$. When marks are represented only by their remaining suffixes, there can be at most n distinct marks for an input string of length n , since each mark corresponds to a unique suffix of the input and there are at most n such suffixes. Even in the presence of nested stars, the outer unfolding reuses the same set of suffix marks generated by the inner star, so although the number of recursive calls can increase, the number of distinct marks at each stage remains bounded by n .

We started with a lexer implementation of this version, and from initial testing, it produces the POSIX value, yet this remains to be proved. In that version, we introduce reshuffling, where marks are ordered based on their remaining strings. At present, the reordering is carried out in the sequence case after shifting through the first part, and in the *Star* and r^n cases after shifting into the inner regular expression.

The behaviour of each case of the *shifts* function is discussed in detail below.

Character case (c):

- If the head of the string in a mark matches c , it is stripped and the remainder retained; otherwise, the mark is dropped.

Alternative case ($r_1 + r_2$):

- The list of marks ms is shifted into both subexpressions r_1 and r_2 . The two resulting lists are then concatenated, corresponding to matching with either r_1 or r_2 .

$$\begin{aligned}
shifts(ms, 0) &\stackrel{\text{def}}{=} [] \\
shifts(ms, 1) &\stackrel{\text{def}}{=} [] \\
shifts(ms, d) &\stackrel{\text{def}}{=} [\bullet s \mid \bullet d :: s \in ms] \\
shifts(ms, r_1 + r_2) &\stackrel{\text{def}}{=} shifts(ms, r_1) @ shifts(ms, r_2) \\
shifts(ms, r_1 \cdot r_2) &\stackrel{\text{def}}{=} \text{let } ms' = shifts(ms, r_1) \text{ in} \\
&\quad \begin{cases} shifts(ms' @ ms, r_2) @ ms' & \text{if } nullable(r_1) \wedge nullable(r_2), \\ shifts(ms' @ ms, r_2) & \text{if } nullable(r_1), \\ shifts(ms', r_2) @ ms' & \text{if } nullable(r_2), \\ shifts(ms', r_2) & \text{otherwise} \end{cases} \\
shifts(ms, r^*) &\stackrel{\text{def}}{=} \text{let } ms' = shifts(ms, r) \text{ in} \\
&\quad \text{if } ms' = [] \text{ then } [] \text{ else } shifts(ms', r^*) @ ms'
\end{aligned}$$

Figure 10: The *shifts* function. Marks are annotated with strings. A mark is represented as $\bullet s$, where s is a string, and ms is a list of such marks.

Sequence case ($r_1 \cdot r_2$):

- This is the most elaborate case. It begins by shifting the marks into the first part of the sequence, r_1 , producing a new list of marks ms' . The result then depends on several conditions, namely:
 1. $nullable(r_1) \wedge nullable(r_2)$: both r_1 and r_2 are nullable, so both can be skipped, and the shifting must account for this by combining the results appropriately. To account for r_1 being skipped, the original marks ms are appended to ms' and then shifted into r_2 . In this way, skipping r_1 has the same effect as shifting ms directly into r_2 . Finally, ms' itself is appended to the result to account for skipping r_2 , which corresponds to matching with r_1 only.
 2. $nullable(r_1)$: if r_1 is nullable, it may be skipped. In this case, the original marks ms are appended to ms' , and the result is then shifted into r_2 , which corresponds to matching directly with r_2 .

3. $nullable(r_2)$: if r_2 is nullable, it may be skipped. Here, ms' is shifted into r_2 , and the result is then appended with ms' to account for the case where r_2 is skipped, which corresponds to matching with r_1 only.
4. $\neg nullable(r_1) \wedge \neg nullable(r_2)$: if neither r_1 nor r_2 is nullable, the marks ms' are shifted directly into r_2 .

Star case (r^*):

- The list of marks ms is first shifted into r , producing a new list ms' . If ms' is empty, this corresponds to the end of the iterations of the *Star*, and the result is the empty list. Otherwise, ms' records one completed iteration and is appended to the result of shifting ms' into the *Star*, corresponding to adding a new iteration.

The matcher is defined as follows.

$$matcher(s, r) \stackrel{\text{def}}{=} \begin{cases} nullable(r) & \text{if } s = [], \\ \bullet [] \in shifts([\bullet s], r) & \text{otherwise.} \end{cases}$$

To illustrate the behaviour of this version, we consider three examples. The first uses the regular expression $(a + (a \cdot c))$ and the string **ac**. We show, step by step, how the marks are reduced as characters are stripped and non-matching paths are dropped. The second example considers the *Star* case, highlighting how the marks are unfolded. A third example, $(a^* \cdot a^*)$ with the string **aa**, shows that, in this version, all possible marks must be generated, since one of them may be the only one required to complete the match when the *Star* expression is part of a larger expression.

In Example 1 (Figure 11), we consider the regular expression $(a + (a \cdot c))$ and the string **ac**. Matching begins with the initial mark carrying the full string, which is then propagated into the two subexpressions: the left subexpression a and the right subexpression $(a \cdot c)$. In the left subexpression a , the initial character matches, so it is stripped from the mark, leaving the string **c**. In the right subexpression $(a \cdot c)$, the sequence matches both characters, so the string in the mark is reduced to the empty string. At this point, the list of marks $[\bullet_{[c]}, \bullet_{[]}]$ is obtained, and with no further calls to *shifts* and one mark reduced to the empty string, the expression is accepted.

In Example 2 (Figure 12), we consider the regular expression (a^*) and the string **aaa**. For the *Star* case, the process starts by shifting the received

marks into the inner expression; here, the initial mark is shifted into a . When the character a matches, the leading a is stripped from the string carried by the mark, and the inner *shifts* call produces a new mark with this shorter string. If the result is empty—meaning that no character was consumed—no further iteration takes place. This mechanism has two effects. First, the strings carried by the marks either remain unchanged or become shorter, in the sense that their length decreases. Second, it enumerates all the ways the *Star* expression can match the input, from no consumption to consuming the entire string. The algorithm therefore returns a list of marks representing these alternatives, and, as Example 3 will show, one of these marks may be the only one required to complete the match when the *Star* expression is part of a larger expression. This observation also suggests a possible optimisation: rather than generating all marks, it may be sufficient to generate only those needed in a given case.

In Example 3 (Figure 13), we consider the regular expression $((a^* \cdot a) \cdot a)$ and the string `aaa`. The full string is initially carried by the mark, as before. The mark is then passed to the first part of the outer sequence, which is itself a sequence. In step 2, the mark enters the first part of the inner sequence, r_1 , which is a^* , and produces the list of marks $[\bullet aa, \bullet a, \bullet []]$, as shown in Example 2. However, since r_1 is a *Star* and thus nullable, the original mark $[\bullet aaa]$ is appended to the result list of the first part to account for the possibility of skipping r_1 . This list of marks is then passed, in step 3, to the second part of the inner sequence, resulting in the consumption of one character from each of the marks (and dropping the empty-list mark from r_1 , since it cannot match the character in r_2). In step 4, the resulting list is passed to the second part of the outer sequence, consuming another character from the remaining marks (and again dropping the empty-list mark), resulting in the final list of marks $[\bullet [], \bullet a]$.

This example shows that the *Star* part of the expression produces multiple marks, and one of them is necessary to complete the match for the entire expression—in particular, the mark $\bullet aa$, which represents consuming only one character by the *Star* constructor.

- String **ac**, regular expression: $(a + (a \cdot c))$

1. $|_{[\bullet ac]} (a + (a \cdot c))$
2. $(|_{[\bullet ac]} a + |_{[\bullet ac]} (a \cdot c))$
3. $(a |_{[\bullet c]} + (a \cdot c) |_{[\bullet []]})$

Figure 11: Example 1: String-carrying shifts for $(a + (a \cdot c))$ with the string **ac**.

- String **aaa**, regular expression: (a^*)

1. $|_{[\bullet aaa]} (a^*)$
2. $(|_{[\bullet aa]} a^*)$
3. $(|_{[\bullet aa, \bullet a]} a^*)$
4. $(a^* |_{[\bullet aa, \bullet a, \bullet []]})$

Figure 12: Example 2: String-carrying shifts for (a^*) with the string **aaa**.

- String **aaa**, regular expression: $((a^* \cdot a) \cdot a)$

1. $|_{[\bullet aaa]} ((a^* \cdot a) \cdot a)$
2. $((a^* |_{[\bullet aa, \bullet a, \bullet [], \bullet aaa]} \cdot a) \cdot a)$
3. $((a^* \cdot a |_{[\bullet a, \bullet [], \bullet aa]}) \cdot a)$
4. $((a^* \cdot a) \cdot a) |_{[\bullet [], \bullet a]}$

Figure 13: Example 3: String-carrying shifts for $((a^* \cdot a) \cdot a)$ with the string **aaa**.

5 Future Work

The development of the marked matcher is still ongoing, and several directions for future work have been identified:

- **Continue the latest version to produce POSIX values.** Work has started on POSIX value extraction for the latest version of the matcher. Initial results indicate that it produces POSIX values. Further work is needed to fix the handling of exact repetition r^n .
- **Extend support for additional operators.** Future work will focus on extending the matcher to handle intersection and negation operators. This will require extending both the shifting function and the bitcoding scheme so that these operators are handled correctly during matching.
- **Formal verification in Isabelle/HOL.** Work has started on formally proving the correctness of the latest version of the algorithm in Isabelle/HOL. The next step is to prove that it always produces the POSIX value.
- **Performance and optimisation.** We are exploring ways to make the algorithm more efficient. We plan to investigate strategies to control or reduce mark proliferation in complex expressions and to handle bitcode annotations more efficiently.
- **Parsing with marks.** We plan to explore extending the marked approach to parsing by generalising from regular expressions to context-free grammars, in analogy to Might et al.'s [4] generalisation of derivatives to obtain a functional parsing algorithm.

References

- [1] A. Asperti, C. S. Coen, and E. Tassi. Regular Expressions, au point. *arXiv*, <http://arxiv.org/abs/1010.2604>, 2010.
- [2] J. A. Brzozowski. Derivatives of regular expressions. *J. ACM*, 11(4):481–494, Oct. 1964.

- [3] S. Fischer, F. Huch, and T. Wilke. A Play on Regular Expressions: Functional Pearl. In *Proc. of the 15th ACM SIGPLAN International Conference on Functional Programming (ICFP)*, pages 357–368, 2010.
- [4] M. Might, D. Darais, and D. Spiewak. Parsing with derivatives: A functional pearl. In *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming (ICFP)*, pages 189–195. ACM, 2011.
- [5] T. Nipkow and D. Traytel. Unified decision procedures for regular expression equivalence. In G. Klein and R. Gamboa, editors, *Interactive Theorem Proving*, pages 450–466, Cham, 2014. Springer International Publishing.
- [6] S. Okui and T. Suzuki. Disambiguation in regular expression matching via position automata with augmented transitions. Technical Report Technical Report 2013-002, The University of Aizu, Language Processing Systems Laboratory, June 2013.
- [7] S. Owens, J. Reppy, and A. Turon. Regular-expression derivatives re-examined. *Journal of Functional Programming*, 19(2):173–190, 2009.
- [8] M. Sulzmann and K. F. Lu. Posix regular expression parsing with derivatives. *Science of Computer Programming*, 89:179–193, 2014.
- [9] C. Tan and C. Urban. *POSIX Lexing with Bitcoded Derivatives*, pages 26:1–26:18. Apr. 2023.

A Appendix

A.1 Scala Implementation of Fischer et al.’s Marked Approach

In Fischer et al.’s approach, the marks are shifted through the regular expression with each input character. The process starts with an initial mark inserted at the beginning, which is then moved step by step as the input is consumed. This behaviour is implemented by the function *shift*, which performs the core logic of the algorithm. The initial specification of this function is given below, as we have developed several versions throughout our work.

The following presents the Scala implementation of the shifting behaviour as originally defined by Fischer et al. [3]. The *shift* function takes as input a regular expression to match against, a flag m , and a character c , and returns a *marked regular expression*—that is, a regular expression annotated with marks. We write a marked regular expression as $\bullet r$, where the preceding dot indicates that the expression r has been annotated with marks to record the progress of matching.

$$\text{shift}(m, c, r) = \begin{cases} \bullet r & \text{if } c \text{ matches in } r, \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

The flag m indicates the mode of operation: when set to true, a new mark is introduced; otherwise, the function shifts the existing marks. This was realised in our first implementation by adding a Boolean attribute to the character constructor to represent a marked character. In later versions, we instead introduced a wrapper constructor around the character constructor to explicitly represent a marked character. Shifting marks for the base cases $\mathbf{0}$

$$\begin{aligned} \text{shift}(m, c, \mathbf{0}) & \stackrel{\text{def}}{=} \mathbf{0} \\ \text{shift}(m, c, \mathbf{1}) & \stackrel{\text{def}}{=} \mathbf{1} \\ \text{shift}(m, c, d) & \stackrel{\text{def}}{=} \begin{cases} \bullet d & \text{if } c = d \wedge m \\ d & \text{otherwise} \end{cases} \\ \text{shift}(m, c, r_1 + r_2) & \stackrel{\text{def}}{=} \text{shift}(m, c, r_1) + \text{shift}(m, c, r_2) \\ \text{shift}(m, c, r_1 \cdot r_2) & \stackrel{\text{def}}{=} \begin{cases} \text{shift}(m, c, r_1) \cdot \text{shift}(\text{true}, c, r_2) & \text{if } m \wedge \text{nullable}(r_1) \\ \text{shift}(m, c, r_1) \cdot \text{shift}(\text{true}, c, r_2) & \text{if } \text{fin}(r_1) \\ \text{shift}(m, c, r_1) \cdot \text{shift}(\text{false}, c, r_2) & \text{otherwise} \end{cases} \\ \text{shift}(m, c, r^*) & \stackrel{\text{def}}{=} \begin{cases} \text{shift}(\text{true}, c, r^*) & \text{if } \text{fin}(r) \\ \text{shift}(m, c, r^*) & \end{cases} \end{aligned}$$

Figure 14: The *shift* function in the Scala implementation of Fischer’s Marked Approach.

and **1** is straightforward: **0** cannot be marked, and **1**—for now—will not carry a mark, since it matches only the empty string. In the initial algorithm, **1** was not marked, and this choice is carried over into later versions. The reason is to avoid complications and ensure termination of the *shift* function—some of which become apparent in the later versions, as discussed in subsequent sections.¹ The behaviour of the remaining cases is described next.

- **Character case** (d): If the input character c matches d and the flag m is true, a mark is added and stored in the character constructor; otherwise, the character remains unmarked.
- **Alternative case** ($r_1 + r_2$): Marks are shifted into both subexpressions, since either branch may match the input character.
- **Sequence case** ($r_1 \cdot r_2$):
 - If r_1 is neither nullable nor in a final position, marks are shifted only into r_1 , indicating that matching proceeds with the first component.
 - If r_1 is nullable and may be skipped, marks are shifted into both r_1 and r_2 , so that either component can begin matching.
 - If $\text{fin}(r_1)$ holds, meaning r_1 has finished matching, marks are shifted into r_2 to continue matching with that component.
- **Star case** (r^*): Marks are shifted into the subexpression if m is true or if $\text{fin}(r)$ holds.

The formal definitions of the auxiliary functions *fin* and *nullable* appear in Figures 15 and 16, as defined by Fischer et al. [3].

A.2 Bit-Annotated Version 1

In this version, we extended the marked approach with bit sequencing. Inspired by Sulzmann and Lu [8], we introduced bitcodes in the form of lists attached to each mark, which are incrementally built as the marks are shifted.

¹Our choice follows Fischer et al. [3], where **1** is left unmarked (‘shift EPS = EPS’). Asperti et al. [1], on the other hand, use pointed regular expressions (pREs), where acceptance of the empty string is represented by the trailing point being set to true.

$fin(\mathbf{0})$	$\stackrel{\text{def}}{=}$	false
$fin(\mathbf{1})$	$\stackrel{\text{def}}{=}$	false
$fin(c)$	$\stackrel{\text{def}}{=}$	false
$fin(\bullet c)$	$\stackrel{\text{def}}{=}$	true
$fin(r_1 + r_2)$	$\stackrel{\text{def}}{=}$	$fin(r_1) \vee fin(r_2)$
$fin(r_1 \cdot r_2)$	$\stackrel{\text{def}}{=}$	$(fin(r_1) \wedge nullable(r_2)) \vee fin(r_2)$
$fin(r^*)$	$\stackrel{\text{def}}{=}$	$fin(r)$

Figure 15: The *fin* function, which checks whether a marked expression is in a final state, as defined by Fischer et al. [3].

$nullable(\mathbf{0})$	$\stackrel{\text{def}}{=}$	false
$nullable(\mathbf{1})$	$\stackrel{\text{def}}{=}$	true
$nullable(c)$	$\stackrel{\text{def}}{=}$	false
$nullable(r_1 + r_2)$	$\stackrel{\text{def}}{=}$	$nullable(r_1) \vee nullable(r_2)$
$nullable(r_1 \cdot r_2)$	$\stackrel{\text{def}}{=}$	$nullable(r_1) \wedge nullable(r_2)$
$nullable(r^*)$	$\stackrel{\text{def}}{=}$	true

Figure 16: The *nullable* function, which checks whether an expression can match the empty string, as defined by Fischer et al. [3].

This version produces a value, though not necessarily the POSIX value, because when a character is matched more than once at the same point, the associated bitsequence may be overwritten. This can cause value erasure and, in some cases, the loss of the POSIX value, as illustrated in Example 2 below. We use the bit annotations 0 and 1, similar to the bitcoded derivatives described by Sulzmann and Lu [8].

In contrast to the original work by Fischer et al., our *shift* function takes an additional argument, *Bits*. As *shift* is applied, bits are appended to the bitsequence. For example, when shifting through an alternative, 0 is appended to the bitsequence and passed to the left subexpression, and 1 to the right subexpression. If the input character matches a leaf character node, it is wrapped by the newly defined *POINT* constructor, which represents the mark. The associated bitsequence is stored inside this constructor together with the character.

We define the auxiliary function *mkfin* to extract the bitsequence representing the path, in bits, describing how the expression matched. We adapt *mkeps* from Sulzmann and Lu [8] and from Tan and Urban [9], who define it to construct a value tree and a bitsequence, respectively, for how a nullable expression matches the empty string. Our version also returns a bitsequence: 0 for the left branch and 1 for the right branch in choices; and for the *Star* case, 0 indicates the start of an iteration and 1 its end, with the single bit 1 used for the empty-star case. The auxiliary functions *mkfin* and *mkeps* are defined in Figures 18 and 19, while the definition of *fin* is the same as in Section A.1. We define *shift* for this version as follows:

$$\text{shift}(m, c, bs, r) = \begin{cases} \bullet r_{bs'} & \text{if } c \text{ matches in } r, \\ 0 & \text{otherwise.} \end{cases}$$

where *bs'* is the updated bitsequence. The complete definition appears in Figure 17. The behaviour of each case of the *shift* function is discussed in detail below.

Character case (*d*):

- If the input character *c* matches *d* and the flag *m* is true, a *POINT* wraps the constructor *c*, with the updated bitsequence *bs* stored inside the *POINT*. Otherwise, the character remains unmarked.

Alternative case (*r*₁ + *r*₂):

- Marks are shifted as before, and the direction of the match is annotated: 0 is added to the bitsequence passed to the left subexpression, and 1 is added to the bitsequence passed to the right subexpression.

$$\begin{aligned}
\text{shift}(m, bs, c, \mathbf{0}) & \stackrel{\text{def}}{=} \mathbf{0} \\
\text{shift}(m, bs, c, \mathbf{1}) & \stackrel{\text{def}}{=} \mathbf{1} \\
\text{shift}(m, bs, c, d) & \stackrel{\text{def}}{=} \begin{cases} \bullet d_{bs} & \text{if } m \wedge d = c \\ d & \text{otherwise} \end{cases} \\
\text{shift}(m, bs, c, r_1 + r_2) & \stackrel{\text{def}}{=} \text{shift}(m, bs \oplus 0, c, r_1) + \text{shift}(m, bs \oplus 1, c, r_2) \\
\text{shift}(m, bs, c, r_1 \cdot r_2) & \stackrel{\text{def}}{=} \begin{cases} \text{shift}(m, bs, c, r_1) \cdot \text{shift}(\text{true}, bs @ mkeps(r_1), c, r_2) & \text{if } m \wedge \text{nullable}(r_1) \\ \text{shift}(m, bs, c, r_1) \cdot \text{shift}(\text{true}, bs @ mkfin(r_1), c, r_2) & \text{if } fin(r_1) \\ \text{shift}(m, bs, c, r_1) \cdot \text{shift}(\text{false}, [], c, r_2) & \text{otherwise} \end{cases} \\
\text{shift}(m, bs, c, r^*) & \stackrel{\text{def}}{=} \begin{cases} (\text{shift}(m, bs \oplus 0, c, r))^* & \text{if } m \\ (\text{shift}(\text{true}, bs @ (mkfin(r) \oplus 1), c, r))^* & \text{if } m \wedge fin(r) \\ (\text{shift}(\text{true}, mkfin(r) \oplus 0, c, r))^* & \text{if } fin(r) \\ (\text{shift}(\text{false}, [], c, r))^* & \text{otherwise} \end{cases}
\end{aligned}$$

Figure 17: Bit-Annotated Version 1 *shift*. Here, \oplus stands for appending a bit to a bitsequence, while $@$ stands for concatenation of two bitsequences.

Sequence case ($r_1 \cdot r_2$):

- If r_1 is nullable, a mark is shifted to both r_1 and r_2 : bs is passed to r_1 (representing the path to this expression), and $bs @ mkeps(r_1)$ is passed to r_2 , where *mkeps* returns the bits for an empty-string match. This corresponds to the case where the first part of the sequence is skipped.
- If r_1 is in a final position (that is, it has finished matching), a mark is shifted to r_2 with the bitsequence describing how r_1 was matched, extracted using the *mkfin* function, i.e. $bs @ mkfin(r_1)$.
- Otherwise, marks are shifted only into r_1 , with bs representing the current path to r_1 .

Star case (r^*):

- If a new mark is introduced, bs is passed with 0 appended, representing the beginning of a new iteration of the star.
- If a new mark is introduced and r is in a final position, $bs @ mkfin(r)$ is passed with 1 appended, combining the bits describing the path to r^* with the bits showing how r reached a final position.
- If r is in a final position, $mkfin(r)$ is passed with 0 appended, representing the start of a new iteration.
- If no new mark is introduced, the existing marks are shifted with an empty bitsequence.

$$\begin{aligned}
mkfin(\bullet_{bs} r) &\stackrel{\text{def}}{=} bs \\
mkfin(\bullet_{bs} c) &\stackrel{\text{def}}{=} bs \\
mkfin(r_1 + r_2) &\stackrel{\text{def}}{=} \begin{cases} mkfin(r_1) & \text{if } fin(r_1) \\ mkfin(r_2) & \text{otherwise} \end{cases} \\
mkfin(r_1 \cdot r_2) &\stackrel{\text{def}}{=} \begin{cases} mkfin(r_1) @ mkeps(r_2) & \text{if } fin(r_1) \wedge nullable(r_2) \\ mkfin(r_2) & \text{otherwise} \end{cases} \\
mkfin(r^*) &\stackrel{\text{def}}{=} mkfin(r) \oplus 1
\end{aligned}$$

Figure 18: The $mkfin$ function, which extracts the bitsequence describing how a regular expression matched.

Next, we present two examples of matching a string and extracting a value. The first example shows how the bitsequence is constructed during matching, while the second demonstrates a case where the algorithm fails to produce the POSIX value. When shifting to a point (an already marked character) and the character matches again, the associated bitsequence is overwritten, which can lead to value erasure and, in particular, to the loss of the POSIX value, as illustrated in the second example.

In Example 1 (Figure 20), the process starts by shifting the first character, b , into the regular expression. Since the expression is an alternative, *shift* is

$$\begin{aligned}
mkeps(\mathbf{1}) &\stackrel{\text{def}}{=} [] \\
mkeps(r_1 + r_2) &\stackrel{\text{def}}{=} \begin{cases} 0 \oplus mkeps(r_1) & \text{if } nullable(r_1) \\ 1 \oplus mkeps(r_2) & \text{otherwise} \end{cases} \\
mkeps(r_1 \cdot r_2) &\stackrel{\text{def}}{=} mkeps(r_1) @ mkeps(r_2) \\
mkeps(r^*) &\stackrel{\text{def}}{=} [1]
\end{aligned}$$

Figure 19: The *mkeps* function, which returns the bitsequence describing how a nullable expression matches the empty string.

- String **ba**, regular expression: $(a \cdot b + b \cdot a)$

$$shift\ b \rightarrow (a \cdot b) + (\bullet_{[1]} b \cdot a)$$

$$shift\ a \rightarrow (a \cdot b) + (b \cdot \bullet_{[1]} a)$$

Figure 20: Example 1: Bit-Annotated Version 1 *shift* for $(a \cdot b + b \cdot a)$ with the string **ba**.

applied to both branches. However, only the right branch matches, because its first subexpression begins with *b*. A mark is therefore set with the bitcode list $[1]$. Shifting the second character, *a*, continues into the concatenation on the right branch. Here, the first part is already final, as it matched the preceding *b*, so the shift moves on to the second part, which matches *a*. With no further calls to *shift*, *mkfin* is applied, because the regular expression has reached a final position, indicated by a mark at the end of the expression. The resulting bitsequence is $[1]$, recording that the match followed the right branch of the alternative. In Example 2 (Figure 21), after the first shift on *a*, a mark is placed on the left branch with bits $[0, 0]$, indicating the start of a *Star* iteration followed by a left choice. In the right subexpression, the mark on r_1 of the $a \cdot a$ sequence carries bits $[0, 1]$, representing the start of the *Star* iteration followed by a right choice. When the second character is shifted, the right subexpression r_2 with bits $[0, 1]$ is overwritten during the third shift, when those bits should instead be preserved. These bits correspond to

- String **aaa**, regular expression: $(a + a \cdot a)^*$

$$\text{shift } a \rightarrow (\bullet_{[0,0]} a + \bullet_{[0,1]} a \cdot a)^*$$

$$\text{shift } a \rightarrow (\bullet_{[0,0,0,0]} a + \bullet_{[0,0,0,1]} a \cdot \bullet_{[0,1]} a)^*$$

$$\text{shift } a \rightarrow (\bullet_{[0,0,0,0,0,0]} a + \bullet_{[0,0,0,0,0,1]} a \cdot \bullet_{[0,0,0,1]} a)^*$$

Figure 21: Example 2: Bit-Annotated Version 1 *shift* for $(a + a \cdot a)^*$ with the string **aaa**.

the POSIX match, which begins by matching the right-hand side first and then performing another iteration to match the left-hand side. The correct bitsequence in that case would be $[0, 1, 0, 0]$, with the final 1 indicating the end of the *Star* iteration. This behaviour arises because *POINT* stores only a single bitsequence at a time, with no mechanism for preserving multiple marks.

A.3 Bit-Annotated Version 2

In this version, we updated the previous Bit-Annotated Version 1 to accumulate all possible bitsequences leading to a match, rather than storing only a single bitsequence at each marked position. This change was motivated by two issues observed in Version 1: the overwriting of marks when a character is matched more than once at the same position—which can cause the POSIX value to be lost—and the absence of a mechanism to order marks during shifting so as to always produce the POSIX value.

To address the first issue, each marked position now carries a list of bitsequences. Whenever a character is matched, instead of overwriting the existing bitsequence, the new bitsequence is appended to the list already stored at that position.

To address the second issue, ordering is applied after the shifting process, using a function inspired by the work of Okui and Suzuki [6]. This ordering is not embedded within shifting itself: all possible values are first generated and then ordered externally.

The definition of *shift* for this version is shown in Figure 22. The key

difference in this version, compared to Version 1, is that all possible paths are considered. For example, in the sequence case, the bitsequences that arise when r_1 is nullable and can be skipped are passed along together with those that arise when r_1 is in a final position. This means that the right-hand side of a sequence (r_2) receives both lists of bitsequences—those corresponding to r_1 being skipped and those where r_1 is final—ensuring that all possible paths are accounted for. Similarly, the function *mkfin* is updated to return all possible bitsequences.

The auxiliary function *mkfin* now returns lists of bitsequences, while *mkeps* remains unchanged from the previous version. The updated definition of *mkfin* is given in Figure 23.

By accumulating lists of bitsequences at each position, no potential path is lost due to overwriting. Although a formal proof is not yet complete, extensive testing and reasoning give us confidence that this version generates all possible values for matching a string against a regular expression, including the POSIX value. This improvement, however, comes at the cost of an increase in the number of values returned—particularly in the *Star* case, where every possible way of matching is generated.

A.4 String-Carrying Marks — Lexer (so far)

In this version, we extend the string-carrying-marks design by attaching a bit sequence to each mark, so that each mark consists of its remaining suffix and its corresponding bitsequence. Similar to the Bit-Annotated Versions 1 and 2, the bit sequences are built incrementally as the marks are shifted through the regular expression. The difference here is that the order of marks is determined by their remaining suffixes: shorter suffixes are prioritised over longer ones, while marks with equal suffix length preserve their original order. This reordering, which we call *reshuffle*, is applied after shifting through the first part of a sequence, after shifting into the inner regular expression of a *Star*, and during the handling of r^n . Once shifting is complete, we extract the bitsequence of the first mark whose carried suffix is empty (i.e., fully consumed).

We define *shifts* in this version in Figure 24. The behaviour of each case of the *shifts* function is discussed in detail below.

Character case (c):

$$\begin{aligned}
\text{shift}(m, bs, c, \mathbf{0}) & \stackrel{\text{def}}{=} \mathbf{0} \\
\text{shift}(m, bs, c, \mathbf{1}) & \stackrel{\text{def}}{=} \mathbf{1} \\
\text{shift}(m, bs, c, d) & \stackrel{\text{def}}{=} \begin{cases} \bullet d_{bs} & \text{if } m \wedge d = c \\ d & \text{otherwise} \end{cases} \\
\text{shift}(m, bs, c, r_1 + r_2) & \stackrel{\text{def}}{=} \text{shift}(m, bs \oplus 0, c, r_1) + \text{shift}(m, bs \oplus 1, c, r_2) \\
\text{shift}(m, bs, c, r_1 \cdot r_2) & \stackrel{\text{def}}{=} \text{shift}(m, bs, c, r_1) \cdot \\
& \begin{cases} \text{shift}(m, (bs @ mkeps(r_1) @ mkfin(r_1)), c, r_2) & \text{if } m \wedge \text{nullable}(r_1) \wedge \text{fin}(r_1) \\ \text{shift}(m, bs @ mkeps(r_1), c, r_2) & \text{if } m \wedge \text{nullable}(r_1) \\ \text{shift}(\text{true}, mkfin(r_1), c, r_2) & \text{if } \text{fin}(r_1) \\ \text{shift}(\text{false}, [], c, r_2) & \text{otherwise} \end{cases} \\
\text{shift}(m, bs, c, r^*) & \stackrel{\text{def}}{=} \\
& \begin{cases} (\text{shift}(m, (bs @ mkfin(r)) \oplus 0, c, r))^* & \text{if } m \wedge \text{fin}(r) \\ (\text{shift}(\text{true}, mkfin(r) \oplus 0, c, r))^* & \text{if } \text{fin}(r) \\ (\text{shift}(m, bs \oplus 0, c, r))^* & \text{if } m \\ (\text{shift}(\text{false}, [], c, r))^* & \text{otherwise} \end{cases}
\end{aligned}$$

Figure 22: Bit-Annotated Version 2 *shift*. Here, \oplus stands for appending a single element to each member of a list, \leq stands for appending a list to each member of a list, and $@$ stands for list concatenation.

- Similar to the previous version described in Section 4.1.

Alternative case ($r_1 + r_2$):

- The list of marks *ms* is shifted into both subexpressions r_1 and r_2 , with 0 appended to each mark's bitsequence in the left-hand side r_1 , and 1 appended in r_2 .

Sequence case ($r_1 \cdot r_2$):

$$\begin{aligned}
mkfin(\bullet_{bs} c) &\stackrel{\text{def}}{=} bs \\
mkfin(r_1 + r_2) &\stackrel{\text{def}}{=} \begin{cases} mkfin(r_1) @ mkfin(r_2) & \text{if } fin(r_1) \wedge fin(r_2), \\ mkfin(r_1) & \text{if } fin(r_1), \\ mkfin(r_2) & \text{if } fin(r_2). \end{cases} \\
mkfin(r_1 \cdot r_2) &\stackrel{\text{def}}{=} \begin{cases} mkfin(r_2) @ (mkfin(r_1) \triangleleft mkeps(r_2)) & \text{if } fin(r_1) \wedge nullable(r_2) \wedge fin(r_2), \\ mkfin(r_1) \triangleleft mkeps(r_2) & \text{if } fin(r_1) \wedge nullable(r_2), \\ mkfin(r_2) & \text{otherwise.} \end{cases} \\
mkfin(r^*) &\stackrel{\text{def}}{=} mkfin(r) \oplus 1
\end{aligned}$$

Figure 23: The *mkfin* function for Bit-Annotated Version 2, which returns a list of all bitsequences in a point, representing all possible match paths.

- As in previous versions, this is the most elaborate case, especially with the ordering, which determines how marks are shifted and combined to preserve the order of remaining suffixes. It begins by shifting the marks into the first part of the sequence, r_1 , producing a new list of marks ms' as before. We then apply *reshuffle*, which reorders the list of marks based on the remaining suffixes, giving lower priority to those with longer remaining strings. The result then depends on several conditions, namely:

1. $nullable(r_1) \wedge nullable(r_2)$: when both r_1 and r_2 can be skipped, the returned list of marks must preserve the remaining-suffix order. There are three ordered parts in the result:
 - (a) r_1 consumes and r_2 is skipped: ms' is appended with $mkeps(r_2)$ to account for the case where r_1 consumes part of the string while r_2 is skipped.
 - (b) r_1 and r_2 both consume: each mark in ms' is shifted into r_2 , accounting for both parts consuming characters. (The reason why each mark is shifted individually is explained in the final case below.)

- (c) r_1 is skipped and r_2 consumes: the original marks ms are appended with $mkeps(r_1)$ and shifted into r_2 , accounting for the case where r_1 is entirely skipped.
- 2. $nullable(r_1)$: if r_1 can be skipped, each mark in ms' is first shifted into r_2 , and then the original list of marks ms , appended with $mkeps(r_1)$, is also shifted into r_2 to account for skipping r_1 . This preserves the remaining-suffix order by prioritising marks where r_1 consumes input over those where it is skipped, similar to case (3) above.
- 3. $nullable(r_2)$: if r_2 can be skipped, ms' is first appended with $mkeps(r_2)$ (accounting for skipping r_2 entirely), and then each mark in ms' is shifted into r_2 . This preserves the remaining-suffix order by prioritising marks where r_1 consumes while r_2 is skipped over those where both consume.
- 4. $\neg nullable(r_1) \wedge \neg nullable(r_2)$: if neither r_1 nor r_2 can be skipped, we shift each mark in ms' individually into r_2 . If we moved all marks simultaneously to the second part, we would lose the remaining-suffix order in cases where a list of marks from the first part does not match all subexpressions in the second part, which would disrupt the correct ordering.

Consider the example $((a + b) + ab) \cdot (bc + (c + b))$ with the string abc :

- (a) $\mid_{[\bullet abc]} ((a + b) + ab) \cdot (bc + (c + b))$
- (b) $((a + b) + ab) \mid_{[\bullet c, \bullet bc]} \cdot (bc + (c + b))$

Here, $\bullet c$ has consumed the leftmost match via ab , while $\bullet bc$ consumed only a . If both marks were moved simultaneously, $\bullet c$ would be dropped by the left subexpression since it does not match bc , and $\bullet bc$ —which matches successfully—would produce $\bullet []$ and take priority in the result of the alternative r_2 , since the left and right mark lists are concatenated in order. To avoid this issue, we shift each mark individually into r_2 : first, $\bullet c$ is shifted, producing an empty list for the left subexpression and $\bullet []$ for the right; then $\bullet bc$ is shifted, producing $\bullet c$ for the left and $\bullet []$ for the right. Combining both results keeps the remaining-suffix order and correctly prioritises the POSIX value.

Star case (r^*):

- Similar to the previous version, but with annotating 0 to indicate the beginning of a star iteration, and applying *reshuffle* to reorder the resulting list and preserve the remaining-suffix order. The end of an iteration is indicated by appending 1 when returning the marks representing that iteration.

n-Repetitions case (r^n): The shifting behaviour depends on several cases depending on n .

- If $n = 0$, then no marks are shifted and an empty list is returned, following the choice made in the previous versions for the empty string expression **1** (see Section A.1), since $n = 0$ represents the same language as **1**.
- If $n = 1$, the marks are shifted only into the inner r , representing one iteration. The 0 is appended to indicate the beginning of the iteration, after which *reshuffle* is applied to preserve the remaining-suffix order.
- If $n > 1$, the marks are first shifted into the inner r with 0 appended to indicate the beginning of the iteration, and *reshuffle* is applied to preserve the remaining-suffix order. If the result is empty, an empty list is returned. Otherwise, we have two cases:
 - If r is nullable, we append 1 to indicate the end of the iteration without completing all n repetitions due to skipping r , and shift the result of shifting into r into r^{n-1} to represent the remaining iterations. This accounts for the possibility of skipping r , similar to the behaviour in the *Star* case.
 - If r is not nullable, we simply shift the result into r^{n-1} without appending 1.

$$\begin{aligned}
\text{shifts}(ms, 0) &\stackrel{\text{def}}{=} [] \\
\text{shifts}(ms, 1) &\stackrel{\text{def}}{=} [] \\
\text{shifts}(ms, d) &\stackrel{\text{def}}{=} [\bullet s \mid \bullet d :: s \in ms] \\
\text{shifts}(ms, r_1 + r_2) &\stackrel{\text{def}}{=} \text{shifts}(ms \oplus 0, r_1) @ \text{shifts}(ms \oplus 1, r_2) \\
\text{shifts}(ms, r_1 \cdot r_2) &\stackrel{\text{def}}{=} \text{let } ms' = \text{shifts}(ms, r_1).reshuffle \text{ in} \\
&\quad \left\{ \begin{array}{l} (ms' \triangleleft mkeps(r_2)) @ [\text{shifts}([\bullet s], r_2) \mid \bullet s \in ms']^@ @ \text{shifts}(ms \triangleleft mkeps(r_1), r_2) \\ \text{if } nullable(r_1) \wedge nullable(r_2), \\ [\text{shifts}([\bullet s], r_2) \mid \bullet s \in ms']^@ @ \text{shifts}(ms \triangleleft mkeps(r_1), r_2) \\ \text{if } nullable(r_1) \wedge \neg nullable(r_2), \\ (ms' \triangleleft mkeps(r_2)) @ [\text{shifts}([\bullet s], r_2) \mid \bullet s \in ms']^@ \\ \text{if } \neg nullable(r_1) \wedge nullable(r_2), \\ [\text{shifts}([\bullet s], r_2) \mid \bullet s \in ms']^@ \\ \text{otherwise.} \end{array} \right. \\
\text{shifts}(ms, r^*) &\stackrel{\text{def}}{=} \text{let } ms' = \text{shifts}(ms \oplus 0, r).reshuffle \text{ in} \\
&\quad \text{if } ms' = [] \text{ then } [] \text{ else } \text{shifts}(ms', r^*) @ (ms' \oplus 1) \\
\text{shifts}(ms, r^n) &\stackrel{\text{def}}{=} \text{if } n = 0 \text{ then } [] \\
&\quad \text{else if } n = 1 \text{ then } \text{shifts}(ms \oplus 0, r).reshuffle \\
&\quad \text{else let } ms' = \text{shifts}(ms \oplus 0, r).reshuffle \text{ in} \\
&\quad \text{if } ms' = [] \text{ then } [] \\
&\quad \text{else } \left\{ \begin{array}{ll} (ms' \oplus 1) @ \text{shifts}(ms', r^{n-1}) & \text{if } nullable(r), \\ \text{shifts}(ms', r^{n-1}) & \text{otherwise.} \end{array} \right.
\end{aligned}$$

Figure 24: The *shifts* function for the String-Carrying Marks — Lexer version. Here, \oplus stands for appending a single element to each member of a list, \triangleleft stands for appending a list to each member of a list, and $@$ stands for list concatenation.