

Progress Report - 9 Months

Meshal Binnasban

October 29, 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 marked-based regular expression matcher with the goal of POSIX value extraction. The report reviews the derivative-based approach, which appears to have a correspondence with the marked approach—a connection that we aim to investigate. It also reviews the marked algorithms of Fischer et al. and Asperti et al., which provide matchers only and do not support value extraction. Our work develops several versions of the marked algorithm in Scala aimed at extracting POSIX values. During this process, challenges in handling certain cases led us to refine the algorithm through successive versions, each improving on the last. Future directions include extending the matcher to additional operators, refining disambiguation for repetitions, and formally proving correctness.

Synopsis

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

Derivative-based methods, though elegant, suffer from severe size explosion and remain sensitive to syntactic form. Sulzmann and Lu [?] extended derivatives with bitcodes to record lexing information, but the size problem persists. The marked approach, described in the works of Fischer et al. [?] and Asperti et al. [?], 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-based matcher that consistently yields the POSIX-preferred 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
2	Derivatives	4
2.1	Derivative Extension	5
2.2	Size Explosion.	8
3	Marked Approach	9
3.1	Motivation for a Marked Approach	10
3.2	Versions and Implementation	10
3.2.1	Bit-Annotated POINT, version 1	13
3.2.2	Bit-Annotated POINT, version 2	17
3.2.3	Input-Carrying Marks	17
4	Future Work	17

1 Introduction

The notion of derivatives in regular expressions is well established but has gained renewed attention in the last decade [?, ?]. Their simplicity and compatibility with functional programming have encouraged further study. However, derivatives suffer from growth issues, since each step of taking the derivative can increase the size of subexpressions, as will be reviewed later. The marked approach propagates marks through the regex without creating new subexpressions and offers an attractive potential replacement for derivatives. At present, only matchers based on this method exist, while our work aims to extend it to provide POSIX value extraction. This report first reviews Brzozowski derivatives and the bitcoded variant by Sulzmann and Lu [?], followed by background on the marked approach algorithms described by Fischer et al. [?] and Asperti et al. [?]. We then present our work so far, including several versions of the marked algorithm developed during this period.

2 Derivatives

Brzozowski's derivatives offer an elegant way for string matching. By successively taking the derivative of a regular expression with respect to each input character, one obtains a sequence of derivatives. If the final expression can match the empty string, then the original input is accepted.

The notion of derivatives in regular expressions is well established, but have gained attention in the last decade [?, ?]. Their simplicity and compatibility with functional programming have renewed interest in their use.

To decide whether a string: $a_1a_2\dots a_n$ is in the language of a regular expression: r , we successively compute its derivatives:

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

The string matches r if and only if the final expression r_n accepts the empty string. Here $\text{der}_{a_n}(r)$ stands for the derivative of r with respect to the character a_n , as introduced by Brzozowski [?].

To illustrate how derivatives can be used to match a regex against a string, consider the regular expression $(ab + ba)$. The derivative can check the matching of a string **ba** by taking the derivative of the regex with respect

to **b**, then **a**.

$$\begin{aligned} \text{der}_b r &= \text{der}_b(ab + ba) \\ &= \text{der}_b(ab) + \text{der}_b(ba) \\ &= (\text{der}_b a) \cdot b + (\text{der}_b b) \cdot a \\ &= \emptyset \cdot b + \varepsilon \cdot a \end{aligned}$$

$$\begin{aligned} \text{der}_a(\text{der}_b r) &= \text{der}_a(\emptyset \cdot b + \varepsilon \cdot a) \\ &= \text{der}_a(\emptyset \cdot b) + \text{der}_a(\varepsilon \cdot a) \\ &= \text{der}_a(\emptyset) \cdot b + (\text{der}_a(\varepsilon) \cdot a + \text{der}_a a) \\ &= \emptyset \cdot b + (\emptyset \cdot a + \varepsilon) \end{aligned}$$

Since the final derivative expression contains ε , it matches the empty string. Because no input characters remain, this confirms that the original string **ba** is in the language of the regular expression.

Some subexpressions that arise during the computation of derivatives may be redundant. For example, $\emptyset \cdot b$ expresses matching the empty language, then **b**. Since \emptyset is the empty language, no string can satisfy this, so $\emptyset \cdot b$ simplifies to \emptyset . Such simplifications are sometimes necessary in the derivative method. Although the construction is elegant and recursively defined, it may duplicate large parts of the expression and is highly sensitive to the syntactic form of the regular expression. Simplification reduces the number of generated expressions, but it does not solve the underlying problem of size explosion.

More simplification rules may be applied after each derivative to provide a finite bound on the number of intermediate expressions [?]. These include associativity $(r + s) + t \equiv r + (s + t)$, commutativity $r + s \equiv s + r$, and idempotence $r + r \equiv r$. Rules for the empty language and the empty string can also be applied, such as the one mentioned earlier; e.g. $\emptyset \cdot r \equiv r \cdot \emptyset \equiv \emptyset$, $r + \emptyset \equiv r$, and $r \cdot \varepsilon \equiv \varepsilon \cdot r \equiv r$. These simplifications preserve the language accepted by the regular expression while reducing the number of intermediate expressions. They help to mitigate—though not eliminate—the size explosion in derivatives noted by Sulzmann and Lu [?]; see Section 1.2.

2.1 Derivative Extension

Sulzmann and Lu [?] extend Brzozowski’s derivatives to produce lexing values in addition to deciding whether a match exists. These values record how the match occurred: which part of the regular expression corresponded to which

part of the string, which alternative branch was taken, and how sequences and Kleene stars were matched.

They provide two variants: a bitcode-based construction and an injection-based construction [?]. In the bitcode variant, bit sequences encoding lexing choices are embedded during derivative construction and, after acceptance, decoded to the value [?]. In the injection variant, an *inj* function "injects back" the consumed characters into the value; it reverts the derivative steps to obtain the POSIX value [?]. We focus on the bitcode variant, which more directly inspired our marked approach.

Bitcodes are sequences over $\{0, 1\}$ that encode the choices made in a match. These bits encode branch choices in alternatives and repetitions made during matching.

The following example illustrates how bitcoding works with regular expression: $(a+ab)(b+\varepsilon)$ and string: ab . We proceed by constructing bitcoded derivatives step by step and tracking how the bitcode grows.

Initially, the algorithm internalizes the regular expression. Internalization only adds bit annotations to the alternative constructors found in the expression, where the left branch is annotated with bitcode 0 and the right branch with bitcode 1 [?]. This annotation process is implemented by the function *fuse*, whose definition is given at the end of this section. After internalization, the algorithm proceeds by taking the derivative with respect to the first character of the input string, which in this case is a .

1. Step 1: Internalizing the regex

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

2. Step 2: input a

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

3. Step 3: input b

$$\begin{aligned}
der_b(der_a(r')) &= der_b(({}_0\varepsilon + {}_1\varepsilon \cdot b) \cdot ({}_0b + {}_1\varepsilon)) \\
&= der_b({}_0\varepsilon + {}_1\varepsilon \cdot b) \cdot ({}_0b + {}_1\varepsilon) + der_b(fuse(mkeps(r_1), ({}_0b + {}_1\varepsilon))) \\
&= (der_b({}_0\varepsilon) + der_b({}_1\varepsilon \cdot b)) \cdot ({}_0b + {}_1\varepsilon) + {}_0(der_b({}_0b) + der_b({}_1\varepsilon)) \\
&= (\emptyset + der_b({}_1\varepsilon \cdot b)) \cdot ({}_0b + {}_1\varepsilon) + {}_0({}_0\varepsilon + \emptyset) \\
&= (der_b({}_1\varepsilon) \cdot b + der_b({}_1b)) \cdot ({}_0b + {}_1\varepsilon) + {}_0({}_0\varepsilon + \emptyset) \\
&\rightarrow (\emptyset \cdot b + {}_1\varepsilon) \cdot ({}_0b + {}_1\varepsilon) + {}_0({}_0\varepsilon + \emptyset)
\end{aligned}$$

The result of Step 2 shows that the ε symbols indicate a successful match, while the bitcode records how that match was obtained. One match is obtained by taking the left branch to match the string a , reflected in the bitcode [0]. The other match arises by taking the right branch, which matches the regular expression: (ab) . The function $mkeps$, as defined by Sulzmann and Lu [?], extracts the bitcode from nullable derivatives; its definition is given at the end of this section.

In Step 3, the resulting derivative expands into an alternative. This occurs because the concatenation has become nullable, which means the first component, r_1 , may be skipped while matching. To account for this, the derivative expands into an alternative: the left branch assumes r_1 is not skipped, while the right branch assumes it is skipped and instead takes the derivative of r_1 directly. This illustrates why derivatives tend to grow in size. Sulzmann and Lu [?] use $fuse(mkeps(r_1))$ to include the bit annotations needed when r_1 is skipped; these bits, extracted by $mkeps$, indicate how r_1 matched the empty string.

After taking the derivative with respect to the entire input string, the algorithm checks whether the result is nullable. If so, it calls $mkeps$ to extract the bitcode indicating how the match was obtained. In this example, there are two possible matches, but the algorithm prefers the left one. Consequently, $mkeps$ returns the bitcode [1, 1], which encodes the choice of the right branch in the first alternative of r_1 (matching the string: ab), followed by the right branch in the second part of the concatenation (matching the empty string).

The final bitcode after calling $mkeps$ is: [1, 1]. Decoding this against the original expression yields the POSIX value:

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

As mentioned earlier, the formal definitions of the auxiliary functions $intern$, $fuse$, and $mkeps$ are given below, as defined by Sulzmann and Lu [?].

- $fuse : bs, r \rightarrow r'$

$$\begin{aligned}
fuse\ bs\ (\emptyset) &\stackrel{\text{def}}{=} \emptyset \\
fuse\ bs\ (\varepsilon_{bs'}) &\stackrel{\text{def}}{=} \varepsilon_{(bs@bs')} \\
fuse\ bs\ (c_{bs'}) &\stackrel{\text{def}}{=} c_{(bs@bs')} \\
fuse\ bs\ ((r_1 + r_2)_{bs'}) &\stackrel{\text{def}}{=} (r_1 + r_2)_{(bs@bs')} \\
fuse\ bs\ ((r_1 \cdot r_2)_{bs'}) &\stackrel{\text{def}}{=} (r_1 \cdot r_2)_{(bs@bs')} \\
fuse\ bs\ (r_{bs'}^*) &\stackrel{\text{def}}{=} (r^*)_{(bs@bs')} \\
fuse\ bs\ (r_{bs'}^n) &\stackrel{\text{def}}{=} (r^n)_{(bs@bs')}
\end{aligned}$$

- $intern : r \rightarrow r'$

$$\begin{aligned}
intern(\emptyset) &\stackrel{\text{def}}{=} \emptyset \\
intern(\varepsilon) &\stackrel{\text{def}}{=} \varepsilon \\
intern(c) &\stackrel{\text{def}}{=} c \\
intern(r_1 + r_2) &\stackrel{\text{def}}{=} fuse([0], intern(r_1)) + fuse([1], intern(r_2)) \\
intern(r_1 \cdot r_2) &\stackrel{\text{def}}{=} intern(r_1) \cdot intern(r_2) \\
intern(r^*) &\stackrel{\text{def}}{=} (intern(r))^*
\end{aligned}$$

- $mkeps : r \rightarrow bs$

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

2.2 Size Explosion.

Even with aggressive simplifications—as shown by Sulzmann and Lu and another variant of the algorithm in *POSIX Lexing with Bitcoded Derivatives* [?—the simplifications keep the number of derivatives finitely bounded, the number can still grow significantly, making practical usage difficult. The size of derivatives can arise to very large numbers even for some simple expressions. Consider the case $(a + aa)^*$. The number of expressions can grow

significantly because each derivative may introduce new structure. Each step unfolds all possible ways the $\{*\}$ expression can match the input. Urban and Tan showed that even under simplifications such as removing redundant subterms and collapsing identical alternatives, the size remains only finitely bounded, but still grows to drastically large numbers even with said simplifications. This is because derivatives reintroduce the same sublanguage in different syntactic forms (e.g. $a + aa$ versus $a \cdot (1 + a)$), which these rules and duplication removal do not identify as equal [?].

Even Antimirov’s partial derivatives [?], which do not track POSIX values, may produce cubic growth in worst cases. [?]

*** check Chensong example of size explosion even with simplifications

3 Marked Approach

The marked approach is a method for regular expression matching that tracks progress by inserting marks into the expression. As noted by Nipkow and Traytel [?], the idea can be traced to earlier work; they point to Fischer et al. [?] and Asperti et al. [?] as reviving and developing it in a modern setting.

This approach allows for more efficient matching, particularly in complex expressions. The regular expression itself does not grow in size; instead, marks are inserted into it. With each input character, these marks move (or shift) 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, meaning whether the expression accepts the string.

There is a slight difference in how marks are interpreted in the works of Fischer et al. [?] and Asperti et al. [?]. 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 regex to indicate the next character that can be matched. In both cases, acceptance is determined by evaluating the final state of the regex. 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 regex 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 expression is rejected.

3.1 Motivation for a Marked Approach

The marked approach offers an alternative to derivatives for regular expression matching. It depends on propagating markers within the regex rather than constructing new subexpressions. Our main motivation is that this method could support fast and high-performance matching with POSIX value extraction, since it handles matching in a way that avoids some of the limitations of the derivatives approach.

In the derivative method, for example in the SEQ case, the size of the expression typically increases due to the creation of new subexpressions, which contributes to the well-known size explosion problem. By contrast, in the marked approach, matching achieves a similar result by propagating marks through the regex, but without generating larger expressions, and we hope that these marks can also be used to extract POSIX values.

Inspired by the works of Fischer et al. and Asperti et al. [?, ?], we aim to extend the marked approach in order to extract POSIX values, as well as to handle complex constructors used in modern regexes, such as bounded repetitions and intersections. Our work is primarily based on the algorithm described by Fischer et al. [?]. As shown by Nipkow and Traytel [?], the pre-mark algorithm of Asperti et al. is in fact a special case of the post-mark algorithm of Fischer et al., which makes Fischer et al.’s approach the most suitable foundation for extending the marked algorithm to matching with POSIX value extraction.

3.2 Versions and Implementation

We began our work by implementing the marked approach described by Fischer et al. [?] in Scala. In this 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.

Scala Implementation of Fischer's Marked Approach

The following describes the shifting behaviour as defined by Fischer et al. [?]. The *shift* function, which forms the core of the algorithm, takes as input a regular expression to match against, a flag m , and a character c , and returns a *marked regular expression*. 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.

$$\text{shift} : (m, c, r) \mapsto \bullet r$$

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 character constructor to explicitly represent a marked character.

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

Shifting marks for the base cases `ZERO` and `ONE` is straightforward, as they cannot have any marks on them.

- **Character case** (d): A mark is set if the character matches the input character and the mode is `true`. The mode controls how marks propagate through the regular expression.
- **Alternative case** ($r_1 + r_2$): Marks are shifted into both branches, since either branch could match the same input character.
- **Sequence case** ($r_1 \cdot r_2$):
 - If neither r_1 is nullable nor in a final position, shift only into r_1 , indicating that we are matching the first part of the sequence.
 - If r_1 is nullable and can be skipped, also shift into r_2 so that both parts can begin matching.
 - If $\text{fin}(r_1)$ holds (meaning r_1 has finished matching), shift into r_2 to begin matching its part.
- **Star case** (r^*): Shift into the subexpression if the mode is `true` or if the subexpression is in a final position.

Helper Functions Definitions:

- $\text{fin}(\text{regex}) \rightarrow \text{Boolean}$

$$\begin{aligned}
 \text{fin}(\emptyset) &\stackrel{\text{def}}{=} \text{false} \\
 \text{fin}(\varepsilon) &\stackrel{\text{def}}{=} \text{false} \\
 \text{fin}(c) &\stackrel{\text{def}}{=} \text{false} \\
 \text{fin}(\bullet c) &\stackrel{\text{def}}{=} \text{true} \\
 \text{fin}(r_1 + r_2) &\stackrel{\text{def}}{=} \text{fin}(r_1) \vee \text{fin}(r_2) \\
 \text{fin}(r_1 \cdot r_2) &\stackrel{\text{def}}{=} (\text{fin}(r_1) \wedge \text{nullable}(r_2)) \vee \text{fin}(r_2) \\
 \text{fin}(r^*) &\stackrel{\text{def}}{=} \text{fin}(r)
 \end{aligned}$$

- $\text{nullable}(\text{regex}) \rightarrow \text{Boolean}$

$$\begin{aligned}
 \text{nullable}(\emptyset) &\stackrel{\text{def}}{=} \text{false} \\
 \text{nullable}(\varepsilon) &\stackrel{\text{def}}{=} \text{true} \\
 \text{nullable}(c) &\stackrel{\text{def}}{=} \text{false} \\
 \text{nullable}(r_1 + r_2) &\stackrel{\text{def}}{=} \text{nullable}(r_1) \vee \text{nullable}(r_2) \\
 \text{nullable}(r_1 \cdot r_2) &\stackrel{\text{def}}{=} (\text{nullable}(r_1) \wedge \text{nullable}(r_2)) \\
 \text{nullable}(r^*) &\stackrel{\text{def}}{=} \text{true}
 \end{aligned}$$

This first implementation of the algorithm, provides only acceptance checking without any value construction, the following subsections describe the different versions we have developed so far.

3.2.1 Bit-Annotated POINT, version 1

In this version, we extended the marked approach to include bitcodes that annotate the marks that is being shifted through the regex. inspired by Sulzmann and Lu [?], we introduced a bitcode in form of a list that is attached to each mark, and it gets built as the marks are shifted.

This version produces a value but not necessarily the POSIX-preferred value. We use the bit annotations Z (or 0) and S (or 1), similar to the bitcoded derivative. In this version, the function $shift$ takes an additional argument, $Bits$, which is a list of Bit elements. When we shift through this function, bits are added to the list. If a mark is to be added to a leaf character node, this list is stored in the `POINT` constructor that wraps the marked character constructor. We define two additional functions: `mkfin`, which extracts the bit sequence of a final constructor (that is, the path, in bits, describing how the expression matched), and `mkeps`, which extracts the bit sequence of a nullable expression matching the empty string (showing the path that led to the empty-string match). The definitions are given below,

starting with the *shift* function.

$$\begin{aligned}
shift(m, bs, c, \emptyset) &\stackrel{\text{def}}{=} \emptyset \\
shift(m, bs, c, \varepsilon) &\stackrel{\text{def}}{=} \varepsilon \\
shift(m, bs, c, d) &\stackrel{\text{def}}{=} \begin{cases} \bullet_{bs} d & \text{if } m \wedge d = c \\ d & \text{otherwise} \end{cases} \\
shift(m, bs, c, r_1 + r_2) &\stackrel{\text{def}}{=} shift(m, bs \oplus Z, c, r_1) + shift(m, bs \oplus S, c, r_2) \\
shift(m, bs, c, r_1 \cdot r_2) &\stackrel{\text{def}}{=} \begin{cases} shift(m, bs, c, r_1) \cdot shift(true, bs \text{ ++ } \text{mkeps}(r_1), c, r_2) & \text{if } m \wedge \\ & \text{if } m \wedge \\ shift(m, bs, c, r_1) \cdot shift(true, \text{mkfin}(r_1), c, r_2) & \text{if } fin(r_1) \\ shift(m, bs, c, r_1) \cdot shift(false, [], c, r_2) & \text{otherwise} \end{cases} \\
shift(m, bs, c, (r)^*) &\stackrel{\text{def}}{=} \begin{cases} (shift(m, bs \oplus Z, c, r))^* & \text{if } m \\ (shift(true, bs \text{ ++ } (\text{mkfin}(r) \oplus S), c, r))^* & \text{if } m \wedge fin(r) \\ (shift(true, \text{mkfin}(r) \oplus Z, c, r))^* & \text{if } fin(r) \\ (shift(false, [], c, r))^* & \text{otherwise} \end{cases}
\end{aligned}$$

The \oplus denotes adding a bit at the end of a bit list, and ++ denotes the concatenation of two bit lists. The symbols Z (or 0) and S (or 1) are used to represent left and right choices, respectively.

Alternative case ($r_1 + r_2$):

- We shift as before and annotate the direction of the match with Z or S , adding this bit to the list of bits.

Sequence case ($r_1 \cdot r_2$):

- If r_1 is nullable, we shift a mark to both r_1 and r_2 , passing bs to r_1 (representing the current path leading to this expression) and $bs + \text{mkeps}$ to r_2 , where mkeps returns the bits for an empty-string match. This corresponds to the path for this sequence when the first part is skipped.
- If r_1 is in a final position (meaning it has finished matching), we pass a mark to r_2 with the bit list describing how r_1 was matched, extracted using the mkfin function.

- Otherwise, we pass a new mark to r_1 only, with bs representing how we reached r_1 in this sequence.

Star case (r^*):

- If only we are introducing a new mark, we pass bs , representing the bits that describe how the r^* expression was reached and append Z to indicate a new iteration has begun.
- If we are introducing a new mark and r is in a final position, we pass $bs + +\text{mkfin}(r)$, combining the bits representing the path to r^* with the bits from `mkfin` that describe how r reached its final position with S representing an end of one iteration. So to say, to append the bitcode for how r reaches a final position to the bitcode of the new mark to be introduced to the STAR.
- If r is in a final position, then we pass `mkfin` which describe how r reached its final position and then append Z representing the begining of a new iteration.
- If no new mark is introduced, the existing ones are moved by passing an empty list of bits.

Next, we present two examples of matching a string and extracting a value. The first example demonstrates how the bit sequence is constructed during matching, while the second illustrates why the algorithm does not always produce the POSIX-preferred value. In version 1, when shifting to a point (an already marked character) and the character matches again, the associated bit list is overwritten. This behaviour can cause value erasure and, in certain cases, the loss of the POSIX-preferred value, as will be shown in the second example.

1. **Example 1:** ' ba' \rightarrow $(a \cdot b + b \cdot a)$

$$\text{shift } b \rightarrow (a \cdot b) + (\{S\} \bullet b \cdot a)$$

$$\text{shift } a \rightarrow (a \cdot b) + (b \cdot \{S\} \bullet a)$$

With no further calls to `shift`, `mkfin` is called because the regular expression has reached a final position, indicated by a mark at the end of the right-hand subexpression in the alternative. `mkfin` then retrieves the bit list $\{S, S\}$.

2. **Example 2:** ' aaa ' $\rightarrow (a + a \cdot a)^*$

$$\text{shift } a \rightarrow (\{Z,Z\} \bullet a + \{Z,S\} \bullet a \cdot a)^*$$

$$\text{shift } a \rightarrow (\{Z,Z,Z,Z\} \bullet a + \{Z,Z,Z,S\} \bullet a \cdot \{Z,S\} \bullet a)^*$$

$$\text{shift } a \rightarrow (\{Z,Z,Z,Z,Z,Z\} \bullet a + \{Z,Z,Z,Z,Z,S\} \bullet a \cdot \{Z,Z,Z,S\} \bullet a)$$

In this example, after the first shift on a , a mark is placed on the left branch with bits $\{Z, Z\}$, indicating the start of a **STAR** iteration followed by a left choice. In the right subexpression, the mark on r_1 of the ab sequence carries bits $\{Z, S\}$, representing the start of the **STAR** iteration followed by a right choice.

By the third shift, the bits in the right subexpression $\{Z, S\}$ are overwritten, when they should instead be preserved. These bits correspond to the POSIX-preferred match, which starts by matching the right-hand side first, then performing another iteration to match the left-hand side. The correct bit sequence in that case would be $\{Z, S, Z, Z\}$, with the final S marking the end of the **STAR** iteration. This behaviour arises because the **POINT** wrapper stores only a single bit list at a time, and in this version there is no clearly defined ordering of marks.

The annotation can become difficult to follow, which motivated us to introduce additional symbols forming sequences (transitioning from a bit-based representation to sequences) in later versions. For instance, we use N to denote the beginning of a **STAR** iteration (replacing Z) and E to denote the end of an iteration (replacing S).

*** below are more of a self note ***

3.2.2 Bit-Annotated POINT, version 2

we are fairly sure/strongly think that this version produces all possible values including the posix value.

3.2.3 Input-Carrying Marks

in this version, we modified the marks to have them carry the input string. initially, the full string is added to the initial mark which will be shifted through the regex, each time a character match, the character will be removed from the string of the mark. a match happen when there is a mark with an empty string/matching the empty string. marks are organized in order of posix value, we are fairly sure/think of that. basically, the only reordering happens at SEQ case, after shifting through the first part, this is to reorder the marks based on remaining strings meaning that the marks with shorter remaining strings will be at the front of the list.

*** from previous report ***

4 Future Work

This project focuses on implementing and validating a correct and efficient marked regular expression matcher under POSIX disambiguation. Several directions remain open and are planned for the next stages of the PhD:

- **POSIX Disambiguation for STAR.** While the current matcher correctly computes POSIX values for many expressions, disambiguation for nested or ambiguous STAR patterns is not yet complete. Ensuring that the correct POSIX-preferred value is selected in all cases involving repetition remains a primary target. The current implementation explores candidate paths, but the disambiguation logic for selecting among them requires refinement and formal confirmation.
- **Support for Additional Operators.** Beyond the basic constructs (ALT, SEQ, STAR, NTIMES), future work includes extending the matcher to handle additional regex operators such as intersection, negation, and lookahead. These additions require careful definition of how

marks behave and how disambiguation should be handled, but could significantly increase the expressiveness of the engine.

- **Formal Proof of POSIX Value Correctness.** A formal verification is planned to prove that the marked matcher always produces the correct POSIX-disambiguated value. This would involve defining the decoding function rigorously and proving its output corresponds to the POSIX-preferred parse. This direction is part of the original PhD proposal, where value extraction and correctness proofs were identified as key goals.