# POSIX Lexing with Bitcoded Derivatives

## Chengsong Tan ✉
King's College London

## Christian Urban ✉
King's College London

──── **Abstract** ────

Sulzmann and Lu described a lexing algorithm that calculates Brzozowski derivatives using bitcodes annotated to regular expressions. Their algorithm generates POSIX values which encode the information of *how* a regular expression matches a string—that is, which part of the string is matched by which part of the regular expression. The purpose of the bitcodes is to generate POSIX values incrementally while derivatives are calculated. They also help with designing an "aggressive" simplification function that keeps the size of derivatives small. Without simplification the size derivatives can grow exponentially resulting in an extremely slow lexing algorithm. In this paper we describe a variant of Sulzmann and Lu's algorithm: Our algorithm is a recursive functional program, whereas Sulzmann and Lu's version involves a fixpoint construction. We *(i)* prove in Isabelle/HOL that our program is correct and generates unique POSIX values; we also *(ii)* establish a polynomial bound for the size of the derivatives. The size can be seen as a proxy measure for the efficiency of the lexing algorithm: because of the polynomial bound our algorithm does not suffer from the exponential blowup in earlier works.

## 1 Introduction

In the last fifteen or so years, Brzozowski's derivatives of regular expressions have sparked quite a bit of interest in the functional programming and theorem prover communities. The beauty of Brzozowski's derivatives [3] is that they are neatly expressible in any functional language, and easily definable and reasoned about in theorem provers—the definitions just consist of inductive datatypes and simple recursive functions. A mechanised correctness proof of Brzozowski's matcher in for example HOL4 has been mentioned by Owens and Slind [9]. Another one in Isabelle/HOL is part of the work by Krauss and Nipkow [6]. And another one in Coq is given by Coquand and Siles [4].

The notion of derivatives [3], written $r \backslash c$, of a regular expression give a simple solution to the problem of matching a string $s$ with a regular expression $r$: if the derivative of $r$ w.r.t. (in succession) all the characters of the string matches the empty string, then $r$ matches $s$ (and *vice versa*). The derivative has the property (which may almost be regarded as its specification) that, for every string $s$ and regular expression $r$ and character $c$, one has $cs \in L\, r$ if and only if $s \in L\, (r \backslash c)$.

If a regular expression matches a string, then in general there is more than one way of how the string is matched. There are two commonly used disambiguation strategies to generate a unique answer: one is called GREEDY matching [5] and the other is POSIX matching [1, 7, 8, 10, 11]. For example consider the string $xy$ and the regular expression $(x + y + xy)^\star$. Either the string can be matched in two 'iterations' by the single letter-regular expressions $x$ and $y$, or directly in one iteration by $xy$. The first case corresponds to GREEDY matching, which first matches with the left-most symbol and only matches the next symbol in case of a mismatch (this is greedy in the sense of preferring instant gratification to delayed repletion). The second case is POSIX matching, which prefers the longest match.
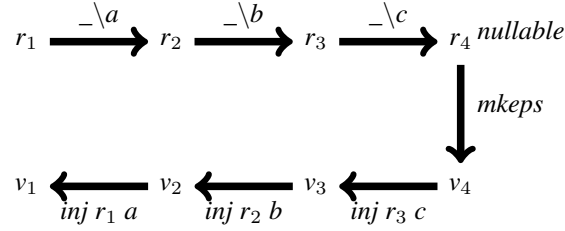
**Figure 1** The two phases of the algorithm by Sulzmann & Lu [10], matching the string $[a, b, c]$. The first phase (the arrows from left to right) is Brzozowski's matcher building successive derivatives. If the last regular expression is *nullable*, then the functions of the second phase are called (the top-down and right-to-left arrows): first *mkeps* calculates a value $v_4$ witnessing how the empty string has been recognised by $r_4$. After that the function *inj* "injects back" the characters of the string into the values.

$$\frac{}{([], \mathbf{1}) \rightarrow Empty} P1 \qquad \frac{}{([c], c) \rightarrow Char\ c} Pc$$

$$\frac{(s, r_1) \rightarrow v}{(s, r_1 + r_2) \rightarrow Left\ v} P+L \qquad \frac{(s, r_2) \rightarrow v \qquad s \notin L\ r_1}{(s, r_1 + r_2) \rightarrow Right\ v} P+R$$

$$\frac{\begin{array}{c}(s_1, r_1) \rightarrow v_1 \qquad (s_2, r_2) \rightarrow v_2 \\ \nexists\ s_3\ s_4.\ s_3 \neq [] \wedge s_3\ @\ s_4 = s_2 \wedge s_1\ @\ s_3\ \in\ L\ r_1 \wedge s_4\ \in\ L\ r_2\end{array}}{(s_1\ @\ s_2, r_1 \cdot r_2) \rightarrow Seq\ v_1\ v_2} PS$$

$$\frac{}{([], r^\star) \rightarrow Stars\ []} P[]$$

$$\frac{\begin{array}{c}(s_1, r) \rightarrow v \qquad (s_2, r^\star) \rightarrow Stars\ vs \qquad |v| \neq [] \\ \nexists\ s_3\ s_4.\ s_3 \neq [] \wedge s_3\ @\ s_4 = s_2 \wedge s_1\ @\ s_3\ \in\ L\ r \wedge s_4\ \in\ L\ (r^\star)\end{array}}{(s_1\ @\ s_2, r^\star) \rightarrow Stars\ (v :: vs)} P\star$$

**Figure 2** Our inductive definition of POSIX values.

$$\begin{aligned}
\mathbf{0} \backslash c &\overset{def}{=} \mathbf{0} \\
\mathbf{1} \backslash c &\overset{def}{=} \mathbf{0} \\
d \backslash c &\overset{def}{=} if\ c = d\ then\ \mathbf{1}\ else\ \mathbf{0} \\
(r_1 + r_2) \backslash c &\overset{def}{=} (r_1 \backslash c) + (r_2 \backslash c) \\
(r_1 \cdot r_2) \backslash c &\overset{def}{=} if\ nullable\ r_1 \\
&\qquad then\ (r_1 \backslash c) \cdot r_2 + (r_2 \backslash c) \\
&\qquad else\ (r_1 \backslash c) \cdot r_2 \\
(r^\star) \backslash c &\overset{def}{=} (r \backslash c) \cdot r^\star
\end{aligned}$$

$$\begin{aligned}
nullable\ (\mathbf{0}) &\overset{def}{=} False \\
nullable\ (\mathbf{1}) &\overset{def}{=} True \\
nullable\ (c) &\overset{def}{=} False \\
nullable\ (r_1 + r_2) &\overset{def}{=} nullable\ r_1 \vee nullable\ r_2 \\
nullable\ (r_1 \cdot r_2) &\overset{def}{=} nullable\ r_1 \wedge nullable\ r_2 \\
nullable\ (r^\star) &\overset{def}{=} True
\end{aligned}$$

## 2   Background

Sulzmann-Lu algorithm with inj. State that POSIX rules. metion slg is correct.

$$\begin{aligned}
mkeps\ \mathbf{1} &\overset{def}{=} Empty \\
mkeps\ (r_1 \cdot r_2) &\overset{def}{=} Seq\ (mkeps\ r_1)\ (mkeps\ r_2) \\
mkeps\ (r_1 + r_2) &\overset{def}{=} if\ nullable\ r_1\ then\ Left\ (mkeps\ r_1)\ else\ Right\ (mkeps\ r_2) \\
mkeps\ (r^\star) &\overset{def}{=} Stars\ []
\end{aligned}$$

$$
\begin{array}{lll}
(1) & inj\ d\ c\ (Empty) & \stackrel{\text{def}}{=} & Char\ d \\
(2) & inj\ (r_1 + r_2)\ c\ (Left\ v_1) & \stackrel{\text{def}}{=} & Left\ (inj\ r_1\ c\ v_1) \\
(3) & inj\ (r_1 + r_2)\ c\ (Right\ v_2) & \stackrel{\text{def}}{=} & Right\ (inj\ r_2\ c\ v_2) \\
(4) & inj\ (r_1 \cdot r_2)\ c\ (Seq\ v_1\ v_2) & \stackrel{\text{def}}{=} & Seq\ (inj\ r_1\ c\ v_1)\ v_2 \\
(5) & inj\ (r_1 \cdot r_2)\ c\ (Left\ (Seq\ v_1\ v_2)) & \stackrel{\text{def}}{=} & Seq\ (inj\ r_1\ c\ v_1)\ v_2 \\
(6) & inj\ (r_1 \cdot r_2)\ c\ (Right\ v_2) & \stackrel{\text{def}}{=} & Seq\ (mkeps\ r_1)\ (inj\ r_2\ c\ v_2) \\
(7) & inj\ (r^\star)\ c\ (Seq\ v\ (Stars\ vs)) & \stackrel{\text{def}}{=} & Stars\ (inj\ r\ c\ v :: vs)
\end{array}
$$

## 3  Bitcoded Regular Expressions and Derivatives

Sulzmann and Lu describe another algorithm that generates POSIX values but dispences with the second phase where characters are injected "back" into values. For this they annotate bitcodes to regular expressions, which we define in Isabelle/HOL as the datatype

$$
\begin{array}{lll}
breg & ::= & ZERO \\
& | & ONE\ bs \\
& | & CHAR\ bs\ c \\
& | & ALTs\ bs\ rs \\
& | & SEQ\ bs\ r_1\ r_2 \\
& | & STAR\ bs\ r
\end{array}
$$

where *bs* stands for a bitsequences; $r$, $r_1$ and $r_2$ for annotated regular expressions; and *rs* for a list of annotated regular expressions. In contrast to Sulzmann and Lu we generalise the alternative regular expressions to lists, instead of just having binary regular expressions. The idea with annotated regular expressions is to incrementally generate the value information by recording bitsequences. Sulzmann and Lu then define a coding function for how values can be coded into bitsequences.

$$
\begin{array}{lll}
code\ (Empty) & \stackrel{\text{def}}{=} & [] \\
code\ (Char\ c) & \stackrel{\text{def}}{=} & [] \\
code\ (Left\ v) & \stackrel{\text{def}}{=} & Z :: code\ v \\
code\ (Right\ v) & \stackrel{\text{def}}{=} & S :: code\ v \\
code\ (Seq\ v_1\ v_2) & \stackrel{\text{def}}{=} & code\ v_1\ @\ code\ v_2 \\
code\ (Stars\ []) & \stackrel{\text{def}}{=} & [S] \\
code\ (Stars\ (v :: vs)) & \stackrel{\text{def}}{=} & Z :: code\ v\ @\ code\ (Stars\ vs)
\end{array}
$$

There is also a corresponding decoding function that takes a bitsequence and generates back a value. However, since the bitsequences are a "lossy" coding (*Seq*s are not coded) the decoding function depends also on a regular expression in order to decode values.

$$
\begin{aligned}
decode'\ bs\ (\mathbf{1}) &\overset{\text{def}}{=} (Empty,\ bs) \\
decode'\ bs\ (d) &\overset{\text{def}}{=} (Char\ d,\ bs) \\
decode'\ []\ (r1.0 + r2.0) &\overset{\text{def}}{=} (Empty,\ []) \\
decode'\ (Z::bs)\ (r1.0 + r2.0) &\overset{\text{def}}{=} \textit{let}\ (v, y) = decode'\ bs\ r1.0\ \textit{in}\ (Left\ v, y) \\
decode'\ (S::bs)\ (r1.0 + r2.0) &\overset{\text{def}}{=} \textit{let}\ (v, y) = decode'\ bs\ r2.0\ \textit{in}\ (Right\ v, y) \\
decode'\ bs\ (r1.0 \cdot r2.0) &\overset{\text{def}}{=} \textit{let}\ (v1, bs') = decode'\ bs\ r1.0;\ (v2, y) = decode'\ bs'\ r2.0\ \textit{in}\ (Seq\ v1\ v2, y) \\
decode'\ []\ (r^\star) &\overset{\text{def}}{=} (Empty,\ []) \\
decode'\ (S::bs)\ (r^\star) &\overset{\text{def}}{=} (Stars\ [],\ bs) \\
decode'\ (Z::bs)\ (r^\star) &\overset{\text{def}}{=} \textit{let}\ (v, bs') = decode'\ bs\ r;\ (vs, y) = decode'\ bs'\ (r^\star)\ \textit{in}\ (Stars\_add\ v\ vs, y) \quad \text{fix}
\end{aligned}
$$

The idea of the bitcodes is to annotate them to regular expressions and generate values incrementally. The bitcodes can be read off from the *breg* and then decoded into a value.

$$
\begin{aligned}
breg \quad ::= \quad &ZERO \\
| \quad &ONE\ bs \\
| \quad &CHAR\ bs\ c \\
| \quad &ALTs\ bs\ rs \\
| \quad &SEQ\ bs\ r_1\ r_2 \\
| \quad &STAR\ bs\ r
\end{aligned}
$$

$$
\begin{aligned}
retrieve\ (ONE\ bs)\ (Empty) &\overset{\text{def}}{=} bs \\
retrieve\ (CHAR\ bs\ c)\ (Char\ d) &\overset{\text{def}}{=} bs \\
retrieve\ (ALTs\ bs\ [r])\ v &\overset{\text{def}}{=} bs\ @\ retrieve\ r\ v \\
retrieve\ (ALTs\ bs\ (r::rs))\ (Left\ v) &\overset{\text{def}}{=} bs\ @\ retrieve\ r\ v \\
retrieve\ (ALTs\ bs\ (r::rs))\ (Right\ v) &\overset{\text{def}}{=} bs\ @\ retrieve\ (ALTs\ []\ rs)\ v \\
retrieve\ (SEQ\ bs\ r_1\ r_2)\ (Seq\ v_1\ v_2) &\overset{\text{def}}{=} bs\ @\ retrieve\ r_1\ v_1\ @\ retrieve\ r_2\ v_2 \\
retrieve\ (STAR\ bs\ r)\ (Stars\ []) &\overset{\text{def}}{=} bs\ @\ [S] \\
retrieve\ (STAR\ bs\ r)\ (Stars\ (v::vs)) &\overset{\text{def}}{=} bs\ @\ [Z]\ @\ retrieve\ r\ v\ @\ retrieve\ (STAR\ []\ r)\ (Stars\ vs)
\end{aligned}
$$

▶ **Theorem 1.** *blexer r s = lexer r s*

bitcoded regexes / decoding / bmkeps gets rid of the second phase (only single phase) correctness

## 4    *Simplification*

Sulzmann & Lu apply simplification via a fixpoint operation; also does not use erase to filter out duplicates.

not direct correspondence with PDERs, because of example problem with retrieve
correctness

## 5    *Bound - NO*

## 6    *Bounded Regex / Not*

## 7    *Conclusion*

[2]

$$\overline{(SEQ\ bs\ ZERO\ r_2) \rightsquigarrow (ZERO)} \qquad \overline{(SEQ\ bs\ r_1\ ZERO) \rightsquigarrow (ZERO)} \qquad \overline{(SEQ\ bs_1\ (ONE\ bs_2)\ r) \rightsquigarrow fuse\ (bs_1\ @\ bs_2)\ r}$$

$$\frac{r_1 \rightsquigarrow r_2}{(SEQ\ bs\ r_1\ r_3) \rightsquigarrow (SEQ\ bs\ r_2\ r_3)} \qquad \frac{r_3 \rightsquigarrow r_4}{(SEQ\ bs\ r_1\ r_3) \rightsquigarrow (SEQ\ bs\ r_1\ r_4)}$$

$$\overline{(ALTs\ bs\ []) \rightsquigarrow (ZERO)} \qquad \overline{(ALTs\ bs\ [r]) \rightsquigarrow fuse\ bs\ r}$$

$$\frac{rs_1 \overset{s}{\rightsquigarrow} rs_2}{(ALTs\ bs\ rs_1) \rightsquigarrow (ALTs\ bs\ rs_2)}$$

$$\overline{[] \overset{s}{\rightsquigarrow} []} \qquad \frac{rs_1 \overset{s}{\rightsquigarrow} rs_2}{r :: rs_1 \overset{s}{\rightsquigarrow} r :: rs_2} \qquad \frac{r_1 \rightsquigarrow r_2}{r_1 :: rs \overset{s}{\rightsquigarrow} r_2 :: rs}$$

$$\overline{ZERO :: rs \overset{s}{\rightsquigarrow} rs} \qquad \overline{ALTs\ bs\ rs_1 :: rs_2 \overset{s}{\rightsquigarrow} (map\ (fuse\ bs)\ rs_1\ @\ rs_2)}$$

$$\frac{r_1^{\downarrow} = r_2^{\downarrow}}{(rs_1\ @\ [r_1]\ @\ rs_2\ @\ [r_2]\ @\ rs_3) \overset{s}{\rightsquigarrow} (rs_1\ @\ [r_1]\ @\ rs_2\ @\ rs_3)}$$

■ **Figure 3** ???

───── **References** ─────

1. The Open Group Base Specification Issue 6 IEEE Std 1003.1 2004 Edition, 2004. http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd_chap09.html.
2. F. Ausaf, R. Dyckhoff, and C. Urban. POSIX Lexing with Derivatives of Regular Expressions (Proof Pearl). In *Proc. of the 7th International Conference on Interactive Theorem Proving (ITP)*, volume 9807 of *LNCS*, pages 69–86, 2016.
3. J. A. Brzozowski. Derivatives of Regular Expressions. *Journal of the ACM*, 11(4):481–494, 1964.
4. T. Coquand and V. Siles. A Decision Procedure for Regular Expression Equivalence in Type Theory. In *Proc. of the 1st International Conference on Certified Programs and Proofs (CPP)*, volume 7086 of *LNCS*, pages 119–134, 2011.
5. A. Frisch and L. Cardelli. Greedy Regular Expression Matching. In *Proc. of the 31st International Conference on Automata, Languages and Programming (ICALP)*, volume 3142 of *LNCS*, pages 618–629, 2004.
6. A. Krauss and T. Nipkow. Proof Pearl: Regular Expression Equivalence and Relation Algebra. *Journal of Automated Reasoning*, 49:95–106, 2012.
7. C. Kuklewicz. Regex Posix. https://wiki.haskell.org/Regex_Posix.
8. S. Okui and T. Suzuki. Disambiguation in Regular Expression Matching via Position Automata with Augmented Transitions. In *Proc. of the 15th International Conference on Implementation and Application of Automata (CIAA)*, volume 6482 of *LNCS*, pages 231–240, 2010.
9. S. Owens and K. Slind. Adapting Functional Programs to Higher Order Logic. *Higher-Order and Symbolic Computation*, 21(4):377–409, 2008.
10. M. Sulzmann and K. Lu. POSIX Regular Expression Parsing with Derivatives. In *Proc. of the 12th International Conference on Functional and Logic Programming (FLOPS)*, volume 8475 of *LNCS*, pages 203–220, 2014.
11. S. Vansummeren. Type Inference for Unique Pattern Matching. *ACM Transactions on Programming Languages and Systems*, 28(3):389–428, 2006.