

POSIX Lexing with Derivatives of Regular Expressions (Proof Pearl)

Fahad Ausaf¹, Roy Dyckhoff², and Christian Urban³

¹ King's College London

fahad.ausaf@icloud.com

² University of St Andrews

roy.dyckhoff@st-andrews.ac.uk

³ King's College London

christian.urban@kcl.ac.uk

Abstract. Brzozowski introduced the notion of derivatives for regular expressions. They can be used for a very simple regular expression matching algorithm. Sulzmann and Lu cleverly extended this algorithm in order to deal with POSIX matching, which is the underlying disambiguation strategy for regular expressions needed in lexers. Sulzmann and Lu have made available on-line what they call a “rigorous proof” of the correctness of their algorithm w.r.t. their specification; regrettably, it appears to us to have unfillable gaps. In the first part of this paper we give our inductive definition of what a POSIX value is and show (i) that such a value is unique (for given regular expression and string being matched) and (ii) that Sulzmann and Lu’s algorithm always generates such a value (provided that the regular expression matches the string). We also prove the correctness of an optimised version of the POSIX matching algorithm. Our definitions and proof are much simpler than those by Sulzmann and Lu and can be easily formalised in Isabelle/HOL. In the second part we analyse the correctness argument by Sulzmann and Lu and explain why it seems hard to turn it into a proof rigorous enough to be accepted by a system such as Isabelle/HOL.

Keywords: POSIX matching, Derivatives of Regular Expressions, Isabelle/HOL

1 Introduction

Brzozowski [1] introduced the notion of the *derivative* $r \setminus c$ of a regular expression r w.r.t. a character c , and showed that it gave 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 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 $s \in L(r \setminus c)$. The beauty of Brzozowski’s derivatives 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

completely formalised correctness proof of this matcher in for example HOL4 has been given in [4].

One limitation of Brzozowski’s matcher is that it only generates a YES/NO answer for whether a string is being matched by a regular expression. Sulzmann and Lu [5] extended this matcher to allow generation not just of a YES/NO answer but of an actual matching, called a [lexical] *value*. They give a simple algorithm to calculate a value that appears to be the value associated with POSIX lexing [3,6]. The challenge then is to specify that value, in an algorithm-independent fashion, and to show that Sulzmann and Lu’s derivative-based algorithm does indeed calculate a value that is correct according to the specification.

The answer given by Sulzmann and Lu [5] is to define a relation (called an “Order Relation”) on the set of values of r , and to show that (once a string to be matched is chosen) there is a maximum element and that it is computed by their derivative-based algorithm. This proof idea is inspired by work of Frisch and Cardelli [2] on a GREEDY regular expression matching algorithm. Beginning with our observations that, without evidence that it is transitive, it cannot be called an “order relation”, and that the relation is called a “total order” despite being evidently not total⁴, we identify problems with this approach (of which some of the proofs are not published in [5]); perhaps more importantly, we give a simple inductive (and algorithm-independent) definition of what we call being a *POSIX value* for a regular expression r and a string s ; we show that the algorithm computes such a value and that such a value is unique. Proofs are both done by hand and checked in Isabelle/HOL. The experience of doing our proofs has been that this mechanical checking was absolutely essential: this subject area has hidden snares. This was also noted by Kuklewitz [3] who found that nearly all POSIX matching implementations are “buggy” [5, Page 203].

If a regular expression matches a string, then in general there are 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 [2] and the other is POSIX matching [3,5]. For example consider the string xy and the regular expression $(x + y + xy)^*$. 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.

In the context of lexing, where an input string needs to be separated into a sequence of tokens, POSIX is the more natural disambiguation strategy for what programmers consider basic syntactic building blocks in their programs. These building blocks are often specified by some regular expressions, say r_{key} and r_{id} for recognising keywords and identifiers, respectively. There are two underlying rules behind tokenising a string in a POSIX fashion:

- The Longest Match Rule (or “maximal munch rule”):
The longest initial substring matched by any regular expression is taken as next token.

⁴ We should give an argument as footnote

- **Rule Priority:**
For a particular longest initial substring, the first regular expression that can match determines the token.

Consider for example r_{key} recognising keywords such as *if*, *then* and so on; and r_{id} recognising identifiers (a single character followed by characters or numbers). Then we can form the regular expression $(r_{key} + r_{id})^*$ and use POSIX matching to tokenise strings, say *iffoo* and *if*. In the first case we obtain by the longest match rule a single identifier token, not a keyword followed by identifier. In the second case we obtain by rule priority a keyword token, not an identifier token—even if r_{id} matches also.

Not Done Yet

Contributions:

Derivatives as calculated by Brzozowski’s method are usually more complex regular expressions than the initial one; various optimisations are possible, such as the simplifications of $\mathbf{0} + r$, $r + \mathbf{0}$, $\mathbf{1} \cdot r$ and $r \cdot \mathbf{1}$ to r . One of the advantages of having a simple specification and correctness proof is that the latter can be refined to allow for such optimisations and simple correctness proof.

An extended version of [5] is available at the website of its first author; this includes some “proofs”, claimed in [5] to be “rigorous”. Since these are evidently not in final form, we make no comment thereon, preferring to give general reasons for our belief that the approach of [5] is problematic rather than to discuss details of unpublished work.

2 Preliminaries

Strings in Isabelle/HOL are lists of characters with the empty string being represented by the empty list, written $[]$, and list-cons being written as $_::_$. Often we use the usual bracket notation for strings; for example a string consisting of a single character is written $[c]$. By using the type *char* for characters we have a supply of finitely many characters roughly corresponding to the ASCII character set. Regular expressions are defined as usual as the following inductive datatype:

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

where $\mathbf{0}$ stands for the regular expression that does not match any string and $\mathbf{1}$ for the regular expression that matches only the empty string. The language of a regular expression is again defined routinely by the recursive function L with the clauses:

$$\begin{aligned} L(\mathbf{0}) &\stackrel{\text{def}}{=} \emptyset \\ L(\mathbf{1}) &\stackrel{\text{def}}{=} \{[]\} \\ L(c) &\stackrel{\text{def}}{=} \{[c]\} \\ L(r_1 \cdot r_2) &\stackrel{\text{def}}{=} L(r_1) @ L(r_2) \\ L(r_1 + r_2) &\stackrel{\text{def}}{=} L(r_1) \cup L(r_2) \\ L(r^*) &\stackrel{\text{def}}{=} (L(r))^* \end{aligned}$$

In the fourth clause we use $_ @ _$ for the concatenation of two languages. We use the star-notation for regular expressions and sets of strings (in the last clause). The star on sets is defined inductively as usual by two clauses for the empty string being in the star of a language and if s_1 is in a language and s_2 and in the star of this language then also $s_1 @ s_2$ is in the star of this language.

Semantic derivatives of sets of strings are defined as

$$Der\ c\ A \stackrel{\text{def}}{=} \{s \mid [c] @ s \in A\}$$

where the second definition lifts the notion of semantic derivatives from characters to strings.

The nullable function

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

The derivative function for characters and strings

$$\begin{aligned} \mathbf{0} \setminus c &\stackrel{\text{def}}{=} \mathbf{0} \\ \mathbf{1} \setminus c &\stackrel{\text{def}}{=} \mathbf{0} \\ c \setminus c &\stackrel{\text{def}}{=} \text{if } c = c' \text{ then } \mathbf{1} \text{ else } \mathbf{0} \\ (r_1 + r_2) \setminus c &\stackrel{\text{def}}{=} (r_1 \setminus c) + (r_2 \setminus c) \\ (r_1 \cdot r_2) \setminus c &\stackrel{\text{def}}{=} \text{if nullable } r_1 \text{ then } (r_1 \setminus c) \cdot r_2 + (r_2 \setminus c) \text{ else } (r_1 \setminus c) \cdot \\ &\quad r_2 \\ (r^*) \setminus c &\stackrel{\text{def}}{=} (r \setminus c) \cdot r^* \end{aligned}$$

It is a relatively easy exercise to prove that

$$\begin{aligned} \text{nullable } r &= (\mathbf{1} \in L(r)) \\ L(r \setminus c) &= Der\ c\ (L(r)) \end{aligned}$$

3 POSIX Regular Expression Matching

4 The Argument by Sulzmann and Lu

5 Conclusion

Nipkow lexer from 2000

Values

$$v := \text{Void} \mid \text{Char } c \mid \text{Left } v \mid \text{Right } v \mid \text{Seq } v_1 v_2 \mid \text{Stars } vs$$

The language of a regular expression

$$\begin{aligned} L(\mathbf{0}) &\stackrel{\text{def}}{=} \emptyset \\ L(\mathbf{1}) &\stackrel{\text{def}}{=} \{\emptyset\} \\ L(c) &\stackrel{\text{def}}{=} \{\{c\}\} \\ L(r_1 \cdot r_2) &\stackrel{\text{def}}{=} L(r_1) @ L(r_2) \\ L(r_1 + r_2) &\stackrel{\text{def}}{=} L(r_1) \cup L(r_2) \\ L(r^*) &\stackrel{\text{def}}{=} (L(r))^* \end{aligned}$$

The *flat* function for values

$$\begin{aligned} |\text{Void}| &\stackrel{\text{def}}{=} \emptyset \\ |\text{Char } c| &\stackrel{\text{def}}{=} [c] \\ |\text{Left } v| &\stackrel{\text{def}}{=} |v| \\ |\text{Right } v| &\stackrel{\text{def}}{=} |v| \\ |\text{Seq } v_1 v_2| &\stackrel{\text{def}}{=} |v_1| @ |v_2| \\ |\text{Stars } []| &\stackrel{\text{def}}{=} \emptyset \\ |\text{Stars } (v :: vs)| &\stackrel{\text{def}}{=} |v| @ |\text{Stars } vs| \end{aligned}$$

The *mkeps* function

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

The *inj* function

$$\begin{aligned} \text{inj } (d) c \text{ Void} &\stackrel{\text{def}}{=} \text{Char } d \\ \text{inj } (r_1 + r_2) c \text{ (Left } v_1) &\stackrel{\text{def}}{=} \text{Left } (\text{inj } r_1 c v_1) \\ \text{inj } (r_1 + r_2) c \text{ (Right } v_2) &\stackrel{\text{def}}{=} \text{Right } (\text{inj } r_2 c v_2) \\ \text{inj } (r_1 \cdot r_2) c \text{ (Seq } v_1 v_2) &\stackrel{\text{def}}{=} \text{Seq } (\text{inj } r_1 c v_1) v_2 \\ \text{inj } (r_1 \cdot r_2) c \text{ (Left (Seq } v_1 v_2)) &\stackrel{\text{def}}{=} \text{Seq } (\text{inj } r_1 c v_1) v_2 \\ \text{inj } (r_1 \cdot r_2) c \text{ (Right } v_2) &\stackrel{\text{def}}{=} \text{Seq } (\text{mkeps } r_1) (\text{inj } r_2 c v_2) \\ \text{inj } (r^*) c \text{ (Seq } v \text{ (Stars } vs)) &\stackrel{\text{def}}{=} \text{Stars } ((\text{inj } r c v) :: vs) \end{aligned}$$

The inhabitation relation:

$$\begin{array}{c}
\frac{\vdash v_1 : r_1 \quad \vdash v_2 : r_2}{\vdash \mathit{Seq} v_1 v_2 : (r_1 \cdot r_2)} \\
\frac{\vdash v_1 : r_1}{\vdash (\mathit{Left} v_1) : (r_1 + r_2)} \quad \frac{\vdash v_2 : r_1}{\vdash (\mathit{Right} v_2) : (r_2 + r_1)} \\
\frac{}{\vdash \mathit{Void} : (\mathbf{1})} \quad \frac{}{\vdash (\mathit{Char} c) : (c)} \\
\frac{}{\vdash \mathit{Stars} [] : (r^*)} \quad \frac{\vdash v : r \quad \vdash \mathit{Stars} vs : (r^*)}{\vdash \mathit{Stars} (v :: vs) : (r^*)}
\end{array}$$

We have also introduced a slightly restricted version of this relation where the last rule is restricted so that $|v| \neq []$. This relation for *non-problematic* is written $\models v : r$.

Our Posix relation $s \in r \rightarrow v$

$$\begin{array}{c}
\frac{}{\overline{[] \in (\mathbf{1}) \rightarrow \mathit{Void}}} \quad \frac{}{\overline{[c] \in (c) \rightarrow (\mathit{Char} c)}} \\
\frac{s \in r_1 \rightarrow v}{s \in (r_1 + r_2) \rightarrow (\mathit{Left} v)} \quad \frac{s \in r_2 \rightarrow v \quad s \notin L(r_1)}{s \in (r_1 + r_2) \rightarrow (\mathit{Right} v)} \\
\frac{\begin{array}{c} s_1 \in r_1 \rightarrow v_1 \quad s_2 \in r_2 \rightarrow v_2 \\ \nexists s_3 s_4. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r_1) \wedge s_4 \in L(r_2) \end{array}}{(s_1 @ s_2) \in (r_1 \cdot r_2) \rightarrow \mathit{Seq} v_1 v_2} \\
\frac{\begin{array}{c} s_1 \in r \rightarrow v \quad s_2 \in (r^*) \rightarrow \mathit{Stars} vs \\ |v| \neq [] \quad \nexists s_3 s_4. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r) \wedge s_4 \in L(r^*) \end{array}}{(s_1 @ s_2) \in (r^*) \rightarrow \mathit{Stars} (v :: vs)} \\
\frac{}{\overline{[] \in (r^*) \rightarrow \mathit{Stars} []}}
\end{array}$$

Our version of Sulzmann's ordering relation

$$\begin{array}{c}
\frac{v_1 \succ_{r_1} v_1' \quad v_1 \neq v_1'}{\text{Seq } v_1 v_2 \succ (r_1 \cdot r_2) \text{Seq } v_1' v_2'} \quad \frac{v_2 \succ_{r_2} v_2'}{\text{Seq } v_1 v_2 \succ (r_1 \cdot r_2) \text{Seq } v_1 v_2'} \\
\frac{\text{len}(|v_1|) \leq \text{len}(|v_2|)}{(\text{Left } v_2) \succ (r_1 + r_2) (\text{Right } v_1)} \quad \frac{\text{len}(|v_2|) < \text{len}(|v_1|)}{(\text{Right } v_1) \succ (r_1 + r_2) (\text{Left } v_2)} \\
\frac{v_2 \succ_{r_2} v_2'}{(\text{Right } v_2) \succ (r_1 + r_2) (\text{Right } v_2')} \quad \frac{v_1 \succ_{r_1} v_1'}{(\text{Left } v_1) \succ (r_1 + r_2) (\text{Left } v_1')} \\
\\
\frac{}{\text{Void} \succ (\mathbf{1}) \text{Void}} \quad \frac{}{(\text{Char } c) \succ (c) (\text{Char } c)} \\
\frac{|\text{Stars } (v :: vs)| = []}{\text{Stars } [] \succ (r^*) \text{Stars } (v :: vs)} \quad \frac{|\text{Stars } (v :: vs)| \neq []}{\text{Stars } (v :: vs) \succ (r^*) \text{Stars } []} \\
\frac{v_1 \succ_r v_2 \quad v_1 \neq v_2}{\text{Stars } (v_1 :: vs_1) \succ (r^*) \text{Stars } (v_2 :: vs_2)} \\
\frac{\text{Stars } vs_1 \succ (r^*) \text{Stars } vs_2}{\text{Stars } (v :: vs_1) \succ (r^*) \text{Stars } (v :: vs_2)} \quad \frac{}{\text{Stars } [] \succ (r^*) \text{Stars } []}
\end{array}$$

A prefix of a string s

$$s_1 \sqsubseteq s_2 \stackrel{\text{def}}{=} \exists s_3. s_1 @ s_3 = s_2$$

Values and non-problematic values

$$\text{Values } r s \stackrel{\text{def}}{=} \{v \mid \vdash v : r \wedge (|v|) \sqsubseteq s\}$$

The point is that for a given s and r there are only finitely many non-problematic values.

Some lemmas we have proved:

$$\begin{array}{l}
L(r) = \{|v| \mid \vdash v : r\} \\
L(r) = \{|v| \mid \models v : r\} \\
\text{If nullable } r \text{ then } \vdash \text{mkeps } r : r. \\
\text{If nullable } r \text{ then } |\text{mkeps } r| = []. \\
\text{If } \vdash v : (r \setminus c) \text{ then } \vdash (\text{inj } r c v) : r. \\
\text{If } \vdash v : (r \setminus c) \text{ then } |\text{inj } r c v| = c :: (|v|). \\
\text{If nullable } r \text{ then } [] \in r \rightarrow \text{mkeps } r. \\
\text{If } s \in r \rightarrow v \text{ then } |v| = s. \\
\text{If } s \in r \rightarrow v \text{ then } \models v : r. \\
\text{If } s \in r \rightarrow v_1 \text{ and } s \in r \rightarrow v_2 \text{ then } v_1 = v_2.
\end{array}$$

This is the main theorem that lets us prove that the algorithm is correct according to $s \in r \rightarrow v$:

$$\text{If } s \in (r \setminus c) \rightarrow v \text{ then } (c :: s) \in r \rightarrow (\text{inj } r c v).$$

Proof The proof is by induction on the definition of *der*. Other inductions would go through as well. The interesting case is for $r_1 \cdot r_2$. First we analyse the case where *nullable* r_1 . We have by induction hypothesis

$$\begin{aligned} (IH1) \quad & \forall s v. \text{ if } s \in (r_1 \setminus c) \rightarrow v \text{ then } (c :: s) \in r_1 \rightarrow (\text{inj } r_1 \ c \ v) \\ (IH2) \quad & \forall s v. \text{ if } s \in (r_2 \setminus c) \rightarrow v \text{ then } (c :: s) \in r_2 \rightarrow (\text{inj } r_2 \ c \ v) \end{aligned}$$

and have

$$s \in ((r_1 \setminus c) \cdot r_2 + (r_2 \setminus c)) \rightarrow v$$

There are two cases what v can be: (1) *Left* v' and (2) *Right* v' .

- (1) We know $s \in ((r_1 \setminus c) \cdot r_2) \rightarrow v'$ holds, from which we can infer that there are s_1, s_2, v_1, v_2 with

$$s_1 \in (r_1 \setminus c) \rightarrow v_1 \quad \text{and} \quad s_2 \in r_2 \rightarrow v_2$$

and also

$$\# s_3 \ s_4. \ s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r_1 \setminus c) \wedge s_4 \in L(r_2)$$

and have to prove

$$(c :: s_1 @ s_2) \in (r_1 \cdot r_2) \rightarrow \text{Seq } (\text{inj } r_1 \ c \ v_1) \ v_2$$

The two requirements $(c :: s_1) \in r_1 \rightarrow (\text{inj } r_1 \ c \ v_1)$ and $s_2 \in r_2 \rightarrow v_2$ can be proved by the induction hypotheses (IH1) and the fact above.

This leaves to prove

$$\# s_3 \ s_4. \ s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge c :: s_1 @ s_3 \in L(r_1) \wedge s_4 \in L(r_2)$$

which holds because $c :: s_1 @ s_3 \in L(r_1)$ implies $s_1 @ s_3 \in L(r_1 \setminus c)$

- (2) This case is similar.

The final case is that \neg *nullable* r_1 holds. This case again similar to the cases above.

References

1. J. A. Brzozowski. Derivatives of Regular Expressions. *Journal of the ACM*, 11(4):481–494, 1964.
2. 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.
3. C. Kuklewicz. Regex Posix. <https://wiki.haskell.org/Regex.Posix>.
4. S. Owens and K. Slind. Adapting Functional Programs to Higher Order Logic. *Higher-Order and Symbolic Computation*, 21(4):377–409, 2008.
5. 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.
6. S. Vansummeren. Type Inference for Unique Pattern Matching. *ACM Transactions on Programming Languages and Systems*, 28(3):389–428, 2006.

6 Roy's Rules

$$\begin{array}{c}
 \text{Void } \triangleleft \epsilon \quad \text{Char } c \triangleleft \text{Lit } c \\
 \\
 \frac{v_1 \triangleleft r_1}{\text{Left } v_1 \triangleleft r_1 + r_2} \quad \frac{v_2 \triangleleft r_2 \quad |v_2| \notin L(r_1)}{\text{Right } v_2 \triangleleft r_1 + r_2} \\
 \\
 \frac{v_1 \triangleleft r_1 \quad v_2 \triangleleft r_2 \quad s \in L(r_1 \setminus |v_1|) \wedge |v_2| \setminus s \in L(r_2) \Rightarrow s = \epsilon}{(v_1, v_2) \triangleleft r_1 \cdot r_2} \\
 \\
 \frac{v \triangleleft r \quad vs \triangleleft r^* \quad |v| \neq \epsilon}{(v :: vs) \triangleleft r^*} \quad \epsilon \triangleleft r^*
 \end{array}$$