

POSIX Lexing with Derivatives of Regular Expressions

Fahad Ausaf Roy Dyckhoff Christian Urban

December 14, 2021

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 [1] cleverly extended this algorithm in order to deal with POSIX matching, which is the underlying disambiguation strategy for regular expressions needed in lexers. In this entry 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.

Contents

1	Values	2
2	The string behind a value	2
3	Relation between values and regular expressions	2
4	Sulzmann and Lu functions	3
5	Mkeys, injval	4
6	Our Alternative Posix definition	4
7	The Lexer by Sulzmann and Lu	5
8	Lexer including simplifications	6

```
theory Lexer  
  imports Regular–Sets.Derivatives  
begin
```

1 Values

```
datatype 'a val =  
  Void  
| Atm 'a  
| Seq 'a val 'a val  
| Right 'a val  
| Left 'a val  
| Stars ('a val) list
```

2 The string behind a value

```
fun  
  flat :: 'a val  $\Rightarrow$  'a list  
where  
  flat (Void) = []  
| flat (Atm c) = [c]  
| flat (Left v) = flat v  
| flat (Right v) = flat v  
| flat (Seq v1 v2) = (flat v1) @ (flat v2)  
| flat (Stars []) = []  
| flat (Stars (v#vs)) = (flat v) @ (flat (Stars vs))
```

```
lemma flat-Stars [simp]:  
  flat (Stars vs) = concat (map flat vs)  
(proof)
```

3 Relation between values and regular expressions

```
inductive  
  Prf :: 'a val  $\Rightarrow$  'a rexp  $\Rightarrow$  bool ( $\vdash$  - : - [100, 100] 100)  
where  
  [[ $\vdash$  v1 : r1;  $\vdash$  v2 : r2]]  $\Longrightarrow$   $\vdash$  Seq v1 v2 : Times r1 r2  
|  $\vdash$  v1 : r1  $\Longrightarrow$   $\vdash$  Left v1 : Plus r1 r2  
|  $\vdash$  v2 : r2  $\Longrightarrow$   $\vdash$  Right v2 : Plus r1 r2  
|  $\vdash$  Void : One  
|  $\vdash$  Atm c : Atom c  
|  $\vdash$  Stars [] : Star r  
| [[ $\vdash$  v : r;  $\vdash$  Stars vs : Star r]]  $\Longrightarrow$   $\vdash$  Stars (v # vs) : Star r
```

```
inductive-cases Prf-elim:  
   $\vdash$  v : Zero  
   $\vdash$  v : Times r1 r2  
   $\vdash$  v : Plus r1 r2  
   $\vdash$  v : One  
   $\vdash$  v : Atom c
```

lemma *Prf-flat-lang*:
assumes $\vdash v : r$ **shows** $\text{flat } v \in \text{lang } r$
 $\langle \text{proof} \rangle$

lemma *Prf-Stars*:
assumes $\forall v \in \text{set } vs. \vdash v : r$
shows $\vdash \text{Stars } vs : \text{Star } r$
 $\langle \text{proof} \rangle$

lemma *Star-string*:
assumes $s \in \text{star } A$
shows $\exists ss. \text{concat } ss = s \wedge (\forall s \in \text{set } ss. s \in A)$
 $\langle \text{proof} \rangle$

lemma *Star-val*:
assumes $\forall s \in \text{set } ss. \exists v. s = \text{flat } v \wedge \vdash v : r$
shows $\exists vs. \text{concat } (\text{map } \text{flat } vs) = \text{concat } ss \wedge (\forall v \in \text{set } vs. \vdash v : r)$
 $\langle \text{proof} \rangle$

lemma *L-flat-Prf1*:
assumes $\vdash v : r$ **shows** $\text{flat } v \in \text{lang } r$
 $\langle \text{proof} \rangle$

lemma *L-flat-Prf2*:
assumes $s \in \text{lang } r$ **shows** $\exists v. \vdash v : r \wedge \text{flat } v = s$
 $\langle \text{proof} \rangle$

lemma *L-flat-Prf*:
 $\text{lang } r = \{\text{flat } v \mid v. \vdash v : r\}$
 $\langle \text{proof} \rangle$

4 Sulzmann and Lu functions

fun

$mkeps :: 'a \text{ rexp} \Rightarrow 'a \text{ val}$

where

$mkeps(\text{One}) = \text{Void}$

$| mkeps(\text{Times } r1 \ r2) = \text{Seq } (mkeps \ r1) \ (mkeps \ r2)$

$| mkeps(\text{Plus } r1 \ r2) = (\text{if } \text{nullable}(r1) \ \text{then } \text{Left } (mkeps \ r1) \ \text{else } \text{Right } (mkeps \ r2))$

$| mkeps(\text{Star } r) = \text{Stars } []$

fun $\text{injval} :: 'a \text{ rexp} \Rightarrow 'a \Rightarrow 'a \text{ val} \Rightarrow 'a \text{ val}$

where

$\text{injval } (\text{Atom } d) \ c \ \text{Void} = \text{Atm } d$

$| \text{injval } (\text{Plus } r1 \ r2) \ c \ (\text{Left } v1) = \text{Left}(\text{injval } r1 \ c \ v1)$

$| \text{injval } (\text{Plus } r1 \ r2) \ c \ (\text{Right } v2) = \text{Right}(\text{injval } r2 \ c \ v2)$

$| \text{injval } (\text{Times } r1 \ r2) \ c \ (\text{Seq } v1 \ v2) = \text{Seq } (\text{injval } r1 \ c \ v1) \ v2$

$| \text{injval } (\text{Times } r1 \ r2) \ c \ (\text{Left } (\text{Seq } v1 \ v2)) = \text{Seq } (\text{injval } r1 \ c \ v1) \ v2$

$| \text{injval } (\text{Times } r1 \ r2) \ c \ (\text{Right } v2) = \text{Seq } (mkeps \ r1) \ (\text{injval } r2 \ c \ v2)$

| $injval (Star\ r)\ c\ (Seq\ v\ (Stars\ vs)) = Stars\ ((injval\ r\ c\ v)\ \#)\ vs$

5 Mkeys, injval

lemma *mkeys-nullable*:

assumes *nullable r*
shows $\vdash mkeys\ r : r$

<proof>

lemma *mkeys-flat*:

assumes *nullable r*
shows $flat\ (mkeys\ r) = []$

<proof>

lemma *Prf-injval*:

assumes $\vdash v : deriv\ c\ r$
shows $\vdash (injval\ r\ c\ v) : r$

<proof>

lemma *Prf-injval-flat*:

assumes $\vdash v : deriv\ c\ r$
shows $flat\ (injval\ r\ c\ v) = c\ \# (flat\ v)$

<proof>

6 Our Alternative Posix definition

inductive

Posix :: 'a list \Rightarrow 'a rexp \Rightarrow 'a val \Rightarrow bool (- \in - \rightarrow - [100, 100, 100] 100)

where

Posix-One: $[] \in One \rightarrow Void$

| *Posix-Atom*: $[c] \in (Atom\ c) \rightarrow (Atm\ c)$

| *Posix-Plus1*: $s \in r1 \rightarrow v \Longrightarrow s \in (Plus\ r1\ r2) \rightarrow (Left\ v)$

| *Posix-Plus2*: $\llbracket s \in r2 \rightarrow v; s \notin lang\ r1 \rrbracket \Longrightarrow s \in (Plus\ r1\ r2) \rightarrow (Right\ v)$

| *Posix-Times*: $\llbracket s1 \in r1 \rightarrow v1; s2 \in r2 \rightarrow v2;$

$\neg(\exists s3\ s4. s3 \neq [] \wedge s3 @ s4 = s2 \wedge (s1 @ s3) \in lang\ r1 \wedge s4 \in lang\ r2) \rrbracket \Longrightarrow$
 $(s1 @ s2) \in (Times\ r1\ r2) \rightarrow (Seq\ v1\ v2)$

| *Posix-Star1*: $\llbracket s1 \in r \rightarrow v; s2 \in Star\ r \rightarrow Stars\ vs; flat\ v \neq [];$

$\neg(\exists s3\ s4. s3 \neq [] \wedge s3 @ s4 = s2 \wedge (s1 @ s3) \in lang\ r \wedge s4 \in lang\ (Star\ r)) \rrbracket$
 $\Longrightarrow (s1 @ s2) \in Star\ r \rightarrow Stars\ (v\ \#)\ vs$

| *Posix-Star2*: $[] \in Star\ r \rightarrow Stars\ []$

inductive-cases *Posix-elim*s:

$s \in Zero \rightarrow v$

$s \in One \rightarrow v$

$s \in Atom\ c \rightarrow v$

$s \in Plus\ r1\ r2 \rightarrow v$

$s \in Times\ r1\ r2 \rightarrow v$

$s \in \text{Star } r \rightarrow v$

lemma *Posix1*:

assumes $s \in r \rightarrow v$

shows $s \in \text{lang } r \text{ flat } v = s$

<proof>

lemma *Posix1a*:

assumes $s \in r \rightarrow v$

shows $\vdash v : r$

<proof>

lemma *Posix-mkeys*:

assumes *nullable* r

shows $\square \in r \rightarrow \text{mkeys } r$

<proof>

lemma *Posix-determ*:

assumes $s \in r \rightarrow v1$ $s \in r \rightarrow v2$

shows $v1 = v2$

<proof>

lemma *Posix-injval*:

assumes $s \in (\text{deriv } c \ r) \rightarrow v$

shows $(c \# s) \in r \rightarrow (\text{injval } r \ c \ v)$

<proof>

7 The Lexer by Sulzmann and Lu

fun

lexer :: 'a rexp \Rightarrow 'a list \Rightarrow ('a val) option

where

lexer $r \ \square = (\text{if nullable } r \text{ then } \text{Some}(\text{mkeys } r) \text{ else } \text{None})$

| *lexer* $r \ (c \# s) = (\text{case } (\text{lexer } (\text{deriv } c \ r) \ s) \text{ of}$

$\text{None} \Rightarrow \text{None}$

| $\text{Some}(v) \Rightarrow \text{Some}(\text{injval } r \ c \ v)$)

lemma *lexer-correct-None*:

shows $s \notin \text{lang } r \iff \text{lexer } r \ s = \text{None}$

<proof>

lemma *lexer-correct-Some*:

shows $s \in \text{lang } r \iff (\exists v. \text{lexer } r \ s = \text{Some}(v) \wedge s \in r \rightarrow v)$

<proof>

lemma *lexer-correctness*:
shows $(\text{lexer } r \ s = \text{Some } v) \longleftrightarrow s \in r \rightarrow v$
and $(\text{lexer } r \ s = \text{None}) \longleftrightarrow \neg(\exists v. s \in r \rightarrow v)$
 $\langle \text{proof} \rangle$

end

theory *Simplifying*
imports *Lexer*
begin

8 Lexer including simplifications

fun *F-RIGHT* **where**
 $F\text{-RIGHT } f \ v = \text{Right } (f \ v)$

fun *F-LEFT* **where**
 $F\text{-LEFT } f \ v = \text{Left } (f \ v)$

fun *F-Plus* **where**
 $F\text{-Plus } f_1 \ f_2 \ (\text{Right } v) = \text{Right } (f_2 \ v)$
 $| F\text{-Plus } f_1 \ f_2 \ (\text{Left } v) = \text{Left } (f_1 \ v)$
 $| F\text{-Plus } f_1 \ f_2 \ v = v$

fun *F-Times1* **where**
 $F\text{-Times1 } f_1 \ f_2 \ v = \text{Seq } (f_1 \ \text{Void}) \ (f_2 \ v)$

fun *F-Times2* **where**
 $F\text{-Times2 } f_1 \ f_2 \ v = \text{Seq } (f_1 \ v) \ (f_2 \ \text{Void})$

fun *F-Times* **where**
 $F\text{-Times } f_1 \ f_2 \ (\text{Seq } v_1 \ v_2) = \text{Seq } (f_1 \ v_1) \ (f_2 \ v_2)$
 $| F\text{-Times } f_1 \ f_2 \ v = v$

fun *simp-Plus* **where**
 $\text{simp-Plus } (\text{Zero}, f_1) \ (r_2, f_2) = (r_2, F\text{-RIGHT } f_2)$
 $| \text{simp-Plus } (r_1, f_1) \ (\text{Zero}, f_2) = (r_1, F\text{-LEFT } f_1)$
 $| \text{simp-Plus } (r_1, f_1) \ (r_2, f_2) = (\text{Plus } r_1 \ r_2, F\text{-Plus } f_1 \ f_2)$

fun *simp-Times* **where**
 $\text{simp-Times } (\text{One}, f_1) \ (r_2, f_2) = (r_2, F\text{-Times1 } f_1 \ f_2)$
 $| \text{simp-Times } (r_1, f_1) \ (\text{One}, f_2) = (r_1, F\text{-Times2 } f_1 \ f_2)$
 $| \text{simp-Times } (r_1, f_1) \ (r_2, f_2) = (\text{Times } r_1 \ r_2, F\text{-Times } f_1 \ f_2)$

lemma *simp-Times-simps*[*simp*]:

$\text{simp-Times } p1 \ p2 = (\text{if } (\text{fst } p1 = \text{One}) \text{ then } (\text{fst } p2, \text{F-Times1 } (\text{snd } p1) (\text{snd } p2))$
 $\text{else } (\text{if } (\text{fst } p2 = \text{One}) \text{ then } (\text{fst } p1, \text{F-Times2 } (\text{snd } p1) (\text{snd } p2))$
 $\text{else } (\text{Times } (\text{fst } p1) (\text{fst } p2), \text{F-Times } (\text{snd } p1) (\text{snd } p2))))$
 <proof>

lemma *simp-Plus-simps*[simp]:
 $\text{simp-Plus } p1 \ p2 = (\text{if } (\text{fst } p1 = \text{Zero}) \text{ then } (\text{fst } p2, \text{F-RIGHT } (\text{snd } p2))$
 $\text{else } (\text{if } (\text{fst } p2 = \text{Zero}) \text{ then } (\text{fst } p1, \text{F-LEFT } (\text{snd } p1))$
 $\text{else } (\text{Plus } (\text{fst } p1) (\text{fst } p2), \text{F-Plus } (\text{snd } p1) (\text{snd } p2))))$
 <proof>

fun
 $\text{simp} :: 'a \text{ rexp} \Rightarrow 'a \text{ rexp} * ('a \text{ val} \Rightarrow 'a \text{ val})$
where
 $\text{simp } (\text{Plus } r1 \ r2) = \text{simp-Plus } (\text{simp } r1) (\text{simp } r2)$
 $| \text{simp } (\text{Times } r1 \ r2) = \text{simp-Times } (\text{simp } r1) (\text{simp } r2)$
 $| \text{simp } r = (r, \text{id})$

fun
 $\text{slexer} :: 'a \text{ rexp} \Rightarrow 'a \text{ list} \Rightarrow ('a \text{ val}) \text{ option}$
where
 $\text{slexer } r \ [] = (\text{if } \text{nullable } r \text{ then } \text{Some}(\text{mkeys } r) \text{ else } \text{None})$
 $| \text{slexer } r \ (c\#\#s) = (\text{let } (rs, fr) = \text{simp } (\text{deriv } c \ r) \text{ in}$
 $\text{case } (\text{slexer } rs \ s) \text{ of}$
 $\text{None} \Rightarrow \text{None}$
 $| \text{Some}(v) \Rightarrow \text{Some}(\text{injval } r \ c \ (fr \ v))))$

lemma *slexer-better-simp*:
 $\text{slexer } r \ (c\#\#s) = (\text{case } (\text{slexer } (\text{fst } (\text{simp } (\text{deriv } c \ r)))) \ s) \text{ of}$
 $\text{None} \Rightarrow \text{None}$
 $| \text{Some}(v) \Rightarrow \text{Some}(\text{injval } r \ c \ ((\text{snd } (\text{simp } (\text{deriv } c \ r))) \ v)))$
 <proof>

lemma *L-fst-simp*:
shows $\text{lang } r = \text{lang } (\text{fst } (\text{simp } r))$
 <proof>

lemma *Posix-simp*:
assumes $s \in (\text{fst } (\text{simp } r)) \rightarrow v$
shows $s \in r \rightarrow ((\text{snd } (\text{simp } r)) \ v)$
 <proof>

lemma *slexer-correctness*:
shows $\text{slexer } r \ s = \text{lexer } r \ s$
 <proof>

end

References

- [1] 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.