

# RegEx

Fahad Ausaf<sup>1</sup>, Roy Dyckhoff<sup>2</sup>, and Christian Urban<sup>1</sup>

<sup>1</sup> King's College London, United Kingdom

<sup>2</sup> St Andrews

**Abstract.** BLA BLA Sulzmann and Lu [1]

**Keywords:**

## 1 Introduction

Regular expressions

$$r := \text{NULL} \mid \text{EMPTY} \mid \text{CHAR } c \mid \text{ALT } r_1 \ r_2 \mid \text{SEQ } r_1 \ r_2 \mid \text{STAR } r$$

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

$$\begin{aligned} L(\text{NULL}) &\stackrel{\text{def}}{=} \emptyset \\ L(\text{EMPTY}) &\stackrel{\text{def}}{=} \{\boxed{\}\!\!}\} \\ L(\text{CHAR } c) &\stackrel{\text{def}}{=} \{[c]\} \\ L(\text{SEQ } r_1 \ r_2) &\stackrel{\text{def}}{=} (L(r_1)) @ (L(r_2)) \\ L(\text{ALT } r_1 \ r_2) &\stackrel{\text{def}}{=} (L(r_1)) \cup (L(r_2)) \\ L(\text{STAR } r) &\stackrel{\text{def}}{=} (L(r))^\star \end{aligned}$$

The nullable function

$$\begin{aligned} \text{nullable } \text{NULL} &\stackrel{\text{def}}{=} \text{False} \\ \text{nullable } \text{EMPTY} &\stackrel{\text{def}}{=} \text{True} \\ \text{nullable } (\text{CHAR } c) &\stackrel{\text{def}}{=} \text{False} \\ \text{nullable } (\text{ALT } r_1 \ r_2) &\stackrel{\text{def}}{=} \text{nullable } r_1 \vee \text{nullable } r_2 \\ \text{nullable } (\text{SEQ } r_1 \ r_2) &\stackrel{\text{def}}{=} \text{nullable } r_1 \wedge \text{nullable } r_2 \\ \text{nullable } (\text{STAR } r) &\stackrel{\text{def}}{=} \text{True} \end{aligned}$$

The derivative function for characters and strings

$$\begin{aligned}
der c \text{ } NULL &\stackrel{\text{def}}{=} \text{NULL} \\
der c \text{ } EMPTY &\stackrel{\text{def}}{=} \text{NULL} \\
der c \text{ } (CHAR c') &\stackrel{\text{def}}{=} \text{if } c = c' \text{ then } EMPTY \text{ else } \text{NULL} \\
der c \text{ } (ALT r_1 r_2) &\stackrel{\text{def}}{=} ALT (der c r_1) (der c r_2) \\
der c \text{ } (SEQ r_1 r_2) &\stackrel{\text{def}}{=} \text{if nullable } r_1 \text{ then } ALT (SEQ (der c r_1) r_2) (der c r_2) \\
&\quad \text{else } SEQ (der c r_1) r_2 \\
der c \text{ } (STAR r) &\stackrel{\text{def}}{=} SEQ (der c r) (STAR r) \\
ders [] r &\stackrel{\text{def}}{=} r \\
ders (c::s) r &\stackrel{\text{def}}{=} ders s (der c r)
\end{aligned}$$

The *flat* function for values

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

The *mkeps* function

$$\begin{aligned}
mkeps \text{ } EMPTY &\stackrel{\text{def}}{=} \text{Void} \\
mkeps \text{ } (SEQ r_1 r_2) &\stackrel{\text{def}}{=} Seq (mkeps r_1) (mkeps r_2) \\
mkeps \text{ } (ALT r_1 r_2) &\stackrel{\text{def}}{=} \text{if nullable } r_1 \text{ then Left (mkeps r_1) else Right (mkeps r_2)} \\
mkeps \text{ } (STAR r) &\stackrel{\text{def}}{=} Stars []
\end{aligned}$$

The *inj* function

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

The inhabituation relation:

$$\begin{array}{c}
\frac{\vdash v_1 : r_1 \quad \vdash v_2 : r_2}{\vdash Seq\ v_1\ v_2 : SEQ\ r_1\ r_2} \\
\frac{\vdash v_1 : r_1}{\vdash (Left\ v_1) : ALT\ r_1\ r_2} \quad \frac{\vdash v_2 : r_1}{\vdash (Right\ v_2) : ALT\ r_2\ r_1} \\
\\
\frac{\vdash Void : EMPTY}{\vdash Stars\ [] : STAR\ r} \quad \frac{\vdash (Char\ c) : CHAR\ c}{\vdash Stars\ (v::vs) : STAR\ r} \\
\frac{\vdash v : r \quad \vdash Stars\ vs : STAR\ r}{\vdash Stars\ (v::vs) : STAR\ r}
\end{array}$$

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  $s \in r \rightarrow v$

$$\begin{array}{c}
\frac{\vdash [] : EMPTY \rightarrow Void \quad \vdash [c] : CHAR\ c \rightarrow (Char\ c)}{\vdash s \in ALT\ r_1\ r_2 \rightarrow (Left\ v)} \quad \frac{\vdash s \in r_2 \rightarrow v \quad s \notin (L\ r_1)}{\vdash s \in ALT\ r_1\ r_2 \rightarrow (Right\ v)} \\
\\
\frac{s_1 \in r_1 \rightarrow v_1 \quad s_2 \in r_2 \rightarrow v_2 \quad \#s_3\ s_4. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in (L\ r_1) \wedge s_4 \in (L\ r_2)}{(s_1 @ s_2) \in SEQ\ r_1\ r_2 \rightarrow Seq\ v_1\ v_2} \\
\\
\frac{|v| \neq [] \quad \#s_3\ s_4. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in (L\ r) \wedge s_4 \in (L\ (STAR\ r))}{(s_1 @ s_2) \in STAR\ r \rightarrow Stars\ (v::vs)}
\end{array}$$

$$\vdash [] : STAR\ r \rightarrow Stars\ []$$

Our version of Sulzmann's ordering relation

$$\begin{array}{c}
\frac{v_1 \succeq_{r_1} v_1' \quad v_1 \neq v_1'}{\text{Seq } v_1 v_2 \succeq_{SEQ r_1 r_2} \text{Seq } v_1' v_2'} \quad \frac{v_2 \succeq_{r_2} v_2'}{\text{Seq } v_1 v_2 \succeq_{SEQ r_1 r_2} \text{Seq } v_1 v_2'}
\\
\frac{\text{len } (|v_1|) \leq \text{len } (|v_2|)}{\text{Left } v_2 \succeq_{ALT r_1 r_2} \text{Right } v_1} \quad \frac{\text{len } (|v_2|) < \text{len } (|v_1|)}{\text{Right } v_1 \succeq_{ALT r_1 r_2} \text{Left } v_2}
\\
\frac{v_2 \succeq_{r_2} v_2'}{\text{Right } v_2 \succeq_{ALT r_1 r_2} \text{Right } v_2'} \quad \frac{v_1 \succeq_{r_1} v_1'}{\text{Left } v_1 \succeq_{ALT r_1 r_2} \text{Left } v_1'}
\\
\\
\frac{}{\text{Void} \succeq_{EMPTY} \text{Void}} \quad \frac{}{\text{Char } c \succeq_{CHAR c} \text{Char } c}
\\
\frac{|Stars (v::vs)| = []}{Stars [] \succeq_{STAR r} Stars (v::vs)} \quad \frac{|Stars (v::vs)| \neq []}{Stars (v::vs) \succeq_{STAR r} Stars []}
\\
\frac{v_1 \succeq_r v_2 \quad v_1 \neq v_2}{Stars (v_1::vs_1) \succeq_{STAR r} Stars (v_2::vs_2)}
\\
\frac{Stars vs_1 \succeq_{STAR r} Stars vs_2}{Stars (v::vs_1) \succeq_{STAR r} Stars (v::vs_2)} \quad \frac{}{Stars [] \succeq_{STAR r} Stars []}
\end{array}$$

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

$$\begin{aligned}
Values \ r \ s &\stackrel{\text{def}}{=} \{v \mid \vdash v : r \wedge (|v|) \sqsubseteq s\} \\
NValues \ r \ s &\stackrel{\text{def}}{=} \{v \mid \models v : r \wedge (|v|) \sqsubseteq s\}
\end{aligned}$$

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

Some lemmas we have proved:

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

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

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

Things we have proved about our version of the Sulzmann ordering

$$\text{If } s \vdash v : r \text{ then } v \succeq_r v.$$

Things we like to prove, but cannot:

$$\text{If } s \in r \rightarrow v_1, \vdash v_2 : r, \text{ then } v_1 \succeq_r v_2$$

## 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.

## 2 Roy's Rules

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