

**Definition 1. Strings**

For a set  $X$ , we define strings over  $X$  to be elements of the free monoid over  $X$ , called  $X^*$ , with  $\varepsilon \in X^*$  denoting the 0 element. The set  $X^*$  will be ranged over by the metavariable  $xs$ .

We use the notation  $\cdot : (X^* \times X^*) \rightarrow X^*$  to mean the concatenation of strings.

The operator  $\text{fold} : ((X \times A \rightarrow A) \times A \times X^*) \rightarrow A$  is defined  
 $\text{fold}(f, a, \varepsilon) = a, \text{fold}(f, a, (x \in X) \cdot xs) = f(x, \text{fold}(f, a, xs))$ .

Strings are equivalent to finite sequences and so will be indexed. For an  $i$  in  $i \in \mathbb{Z}^+ \cup 0$ ,  $xs_i$  is equal to the  $i$ th index and is otherwise undefined.

**Definition 2. Alphabet**

Let  $\Sigma$  be an arbitrary, finite set. This document will call  $\Sigma$  the Alphabet. We will let  $c$  be a metavariable ranging over  $\Sigma$ , and  $s$  be a metavariable ranging over  $\Sigma^*$ , ie, strings of  $\Sigma$ .

**Definition 3. Parsing Expression Grammars**

For a set of labels  $\text{Lab}$ , we define the set  $\text{Peg}(\text{Lab})$  via the following bnf formula, overloading the operator  $\cdot$ .

$$\begin{array}{lcl} p \in \text{Peg}(\text{Lab}) & ::= & \text{lab: } c \\ & | & \text{lab: } p_1 \cdot p_2 \\ & | & \text{lab: } p_1 / p_2 \\ & | & \text{lab: } p^* \end{array}$$

Figure 1: Parsing Expression Grammar Syntax

These rules are called character, sequence, choice, and possessive star respectively.

**Definition 4. Regular Expressions**

We define the set  $\text{Reg}(\text{Lab})$  via the following bnf formula, again overloading the operator  $\cdot$ .  $\text{Lab}$  is a set of labels, and  $l$  will be a metavariable ranging over these.

$$\begin{array}{lcl} r \in \text{Reg}(\text{Lab}) & ::= & \text{lab: } \varepsilon \\ & | & \text{lab: } c \\ & | & \text{lab: } . \\ & | & \text{lab: } \perp \\ & | & \text{lab: } r_1 \cdot r_2 \\ & | & \text{lab: } r_1 \cup r_2 \\ & | & \text{lab: } r_1 \cap r_2 \\ & | & \text{lab: } \neg r \\ & | & \text{lab: } r^* \end{array}$$

Figure 2: Regular Expression Grammar Syntax

These rules are empty, character, wildcard, empty, sequence, union, intersection, negation, and kleene star respectively.

**Definition 5. Parsing Expression Grammar Matching**

Letting  $\text{Matched} = \Sigma^*$  and  $\text{Remainder} = \Sigma^*$ ,

We inherit the partial function  $\text{pegMatch} : (\text{Peg} \times \Sigma^*) \rightarrow (\text{Matched} \times \text{Remainder}) \uplus \{\perp\}$

from the paper “Towards a Typed Semantics for Parsing Expression Grammars”.

**Definition 6.** *Regular Expression Matches*

When matched against a string, a regular expression  $r \in \text{Reg}(\text{Lab})$  will give a structured result,  $\text{Match}(\text{Lab})$ .

$$\begin{aligned}
 m \in \text{Match}(\text{Lab}) &::= \text{Emp}(\text{lab}: l, \text{string}: \varepsilon) \\
 &| \text{Char}(\text{lab}: l, \text{string}: c) \\
 &| \text{Any}(\text{lab}: l, \text{string}: c) \\
 &| \text{Seq}(\text{lab}: l, \text{string}: s, \text{right}: m_1, \text{left}: m_2) \\
 &| \text{JoinL}(\text{lab}: l, \text{string}: s, \text{subexpr}: m) \\
 &| \text{JoinR}(\text{lab}: l, \text{string}: s, \text{subexpr}: m) \\
 &| \text{Meet}(\text{lab}: l, \text{string}: s, \text{left}: m_1, \text{right}: m_1) \\
 &| \text{Neg}(\text{lab}: l, \text{string}: s) \\
 &| \text{StarBase}(\text{lab}: l) \\
 &| \text{StarRec}(\text{lab}: l, \text{string}: s, \text{left}: m, \text{right}: )
 \end{aligned}$$

Figure 3: Regular Expression Matches

Fields, e.g. `.lab`, are maps from  $\text{Match}(\text{Lab})$  with obvious semantics. The field `.submatch` is identity where otherwise undefined.

The function  $\cdot_* : \text{Match}(\text{Lab}) \rightarrow \text{Match}(\text{Lab}) \rightarrow \text{Match}(\text{Lab})$  concatenates a star recursion to a match and is defined

$$\cdot_*(m_1, m_2) = \text{StarRec}(\text{lab}: m_2.\text{lab}, \text{string}: m_1.\text{string} \cdot m_2.\text{string}, \text{left}: m_1, \text{right}: m_2)$$

**Definition 7.** *Character Of*

For a set  $X$  with elements  $x \in X$ , we say that  $x$  is a character of the string  $xs \in X^*$  iff

$$\exists(\text{prf}, \text{suf} \in X^*), xs = \text{prf} \cdot x \cdot \text{suf}$$

We use the notation  $x \text{ char-of } xs$  to mean  $x$  is a character of the string  $xs$ .

**Definition 8.** *Regular Expression Matching*

We define  $\text{regMatch} : (\text{Reg}(\text{Lab}) \times \Sigma^*) \rightarrow \text{Pow}(\text{Match}(\text{Lab}))$  recursively:

$$\text{regMatch}(l: \varepsilon, \varepsilon) = \{\text{Emp}(\text{lab}: l, \text{string}: \varepsilon)\}$$

$$\text{regMatch}(l: \varepsilon, c \cdot s) = \emptyset$$

$$\text{regMatch}(l: c, c) = \{\text{Char}(\text{lab}: l, \text{string}: c)\}$$

$$\text{regMatch}(l: c, s) \text{ where } s \neq c = \emptyset$$

$$\text{regMatch}(l: \cdot, c) = \{\text{Any}(\text{lab}: l, \text{string}: c)\}$$

$$\text{regMatch}(l: \cdot, c \cdot c' \cdot s) = \emptyset$$

$$\text{regMatch}(l: \perp, s) = \emptyset$$

$$\text{regMatch}(l: r_1 \cdot r_2, s) = \{\text{Seq}(\text{lab}: l, \text{string}: s, \text{right}: m_1, \text{left}: m_2)$$

$$| \exists(s', s'' \in \Sigma^*), s' \cdot s'' = s \wedge m_1 \in \text{regMatch}(r_1, s') \wedge m_2 \in \text{regMatch}(r_2, s'')\}$$

$$\text{regMatch}(l: r_1 \cup r_2, s) =$$

$$\{\text{JoinL}(\text{lab}: l, \text{string}: s, \text{subexpr}: m) \mid m \in \text{regMatch}(r_1, s)\} \cup$$

$$\{\text{JoinR}(\text{lab}: l, \text{string}: s, \text{subexpr}: m) \mid m \in \text{regMatch}(r_2, s)\}$$

$$\text{regMatch}(l: r_1 \cap r_2, s) =$$

$$\{\text{Meet}(\text{lab}: l, \text{string}: s, \text{left}: m_1, \text{right}: m_1) \mid m_1 \in \text{regMatch}(r_1, s) \wedge m_2 \in \text{regMatch}(r_2, s)\}$$

$$\text{regMatch}(l: \neg r, s) = \left\{ \begin{array}{l|l} \text{match}(\text{regMatch}(r, s)) \equiv \emptyset & \{\text{Neg}(\text{lab}: l, \text{string}: s)\} \\ \text{otherwise} & \emptyset \end{array} \right.$$

$$\begin{aligned} \text{regMatch}(l: r^*, s) = \\ \{ \text{fold}(\cdot, *, \text{StarBase}(\text{lab}: l), ms \in \text{Match}(\text{Lab})^*) \\ \mid \exists ss \in \Sigma^{**}, \bullet ss = s \wedge \forall i, ms_i \in \text{regMatch}(r, ss_i) \} \end{aligned}$$

While most rules are self explanatory, the rule for concatenation and star may not be.

Concatenation splits the string into two halves, the first of which matches to the left regex, and the second to the right.

Star considers all possible splits of the string, and requires that the constituent strings of the split are each matched by the subexpression.

It follows from the definitions alone that for arbitrary  $(r, s)$ , and for an arbitrary match  $m$  in  $\text{regMatch}(r, s)$ ,  $m.\text{string} = s$ .

### Definition 9. Well Founded Pairs

We define  $\text{labelToExpr} : \text{Peg}(\text{Lab}_1) \cup \text{Reg}(\text{Lab}_2) \rightarrow (\text{Lab}_1 \cup \text{Lab}_2 \rightarrow \text{Pow}(\text{Peg}(\text{Lab}_1) \cup \text{Reg}(\text{Lab}_2)))$  recursively on the syntax of  $\text{Peg}(\text{Lab}_1)$  and  $\text{Reg}(\text{Lab}_2)$ , collecting the parsing expression grammar corresponding to the label.

A parsing expression grammar  $p$  or regular expression  $r$  is in the set of well formed expressions,  $\text{WF}(\text{Lab})$ , iff for every label  $l \in \text{Lab}$ ,  $\text{labelToExpr}(p)(l)$  (resp  $\text{labelToExpr}(r)(l)$ ) has no more than one element.

A pair  $(p, r) \in \text{Peg}(\text{Lab}_1) \times \text{Reg}(\text{Lab}_2)$  is in the set of well formed pairs  $\text{WF2}(\text{Lab}_1 \cup \text{Lab}_2)$  if for every label, the pointwise union  $\text{labelToExpr}(p)$  and  $\text{labelToExpr}(r)$  applied to the label has no more than one element.

WF and WF2 informally mean expressions and pairs where each label is different.

### Definition 10. Parsing Expression Grammar Translation

On input  $(p, r)$ , the function  $\text{pegreg}$  produces a regular expression with the characters consumed by  $p$  removed. It is defined in mutual recursion with the  $\text{negate}$  function, which on input  $p$ , generates a regular expression corresponding to the set of characters that cause  $p$  to match  $\perp$ .

In order to define the output set, we need to have a labelling scheme for combining these; so the labelling scheme will first be detailed, followed by the definition of the function.

$$\begin{array}{lcl}
l_{pr} \in \text{PegregLab}(\text{Lab}) & ::= & (l_1 \in \text{Lab}) \\
& | & \text{PRLabChar}(l \in \text{Lab}) \\
& | & \text{PRLabSeq}(l \in \text{Lab}) \\
& | & \text{PRLabUnionWildcard}(l \in \text{Lab}) \\
& | & \dots
\end{array}$$

Figure 4: Labelling Scheme

Below, the labels in the output will be given implicitly, with examples shown in more detail to make it clear how the labels are formed.

$$\text{pegreg} : (\text{Peg}(\text{Lab}_1) \times \text{Reg}(\text{Lab}_2)) \rightarrow \text{Reg}(\text{PegregLab}(\text{Lab}_1 \cup \text{Lab}_2))$$

$$\text{pegreg}(c, r) = c \cdot r \tag{1}$$

$$\text{pegreg}(p_1 \cdot p_2, r) = \text{pegreg}(p_1, \text{pegreg}(p_2, r)) \tag{2}$$

$$\text{pegreg}(p_1 / p_2, r) = \text{pegreg}(p_1, r) \cup (\text{pegreg}(p_2, r) \cap \text{negate}(p_1) \cdot .^*) \tag{3}$$

$$\text{pegreg}(p^*, r) = \text{pegreg}(p, \varepsilon)^* \cdot (r \cap \text{negate}(p) \cdot .^*) \tag{4}$$

$$\text{negate} : \text{Peg}(\text{Lab}_1) \rightarrow \text{Reg}(\text{PegregLab}(\text{Lab}_1 \cup \text{Lab}_2))$$

$$\text{negate}(c) = . \cap \neg c \tag{5}$$

$$\text{negate}(p_1 \cdot p_2) = \text{negate}(p_1) \cup \text{pegreg}(p_1, \text{negate}(p_2)) \tag{6}$$

$$\text{negate}(p_1 / p_2) = \text{negate}(p_1) \cap \text{negate}(p_2) \tag{7}$$

$$\text{negate}(p^*) = \perp \tag{8}$$

To make it clear how the labels are formed, let  $\text{.lab}$  be, by abuse of notation, a projection from regular expressions to their labels, and see the case for char in more detail:

$$\text{pegreg} : (\text{Peg}(\text{Lab}_1) \times \text{Reg}(\text{Lab}_2)) \rightarrow \text{Reg}(\text{PegregLab}(\text{Lab}))$$

$$\text{pegreg}(c, r) = \text{PRLabSeq}(r.\text{lab}) : (\text{PRLabChar}(r.\text{lab}) : c) \cdot r$$

For another example, consider the choice case. If viewed in the order from the innermost subexpression to the outermost, the recursion constructs a wildcard, then a kleene star with this subexpression, then  $\text{negate}(p_1)$ , then the concatenation of the two, etc.

So  $\text{negate}(p_1)$  is first computed, then  $\text{PRLabUnionWildcard}(\text{negate}(p_1).\text{lab})$  is used to give a label to the expression that matches any character. This pattern of computing new, unique labels from the subexpressions continues until the entire expression is labelled.

This labelling scheme will give  $\text{pegreg}$  the property that it maps WF2 pairs to WF regular expressions.

**Definition 11.** The function

$$\text{extract}_{\text{peg}} : (\text{Peg}(\text{Lab}_1) \times \text{Match}(\text{PegregLab}(\text{Lab}_1 \cup \text{Lab}_2))) \rightarrow \text{Match}(\text{PegregLab}(\text{Lab}_1 \cup \text{Lab}_2))$$

applied to a pair  $(p, m)$  extracts the part corresponding to the peg match.

For any parsing expression grammar  $p$ , any regular expression  $r$ , any string  $s$ , and any match  $m$  in  $\text{regMatch}(\text{pegreg}(p, r), s)$ ,  $\text{extract}_{\text{peg}}$  is total.

Similarly,  $\text{extract}_{\text{reg}} : (\text{Reg}(\text{Lab}_1) \times \text{Match}(\text{PegregLab}(\text{Lab}_1 \cup \text{Lab}_2))) \rightarrow \text{Match}(\text{PegregLab}(\text{Lab}_1 \cup \text{Lab}_2))$  is defined.

**Proposition 1.** *Translation Correspondance*

Let  $(p, r)$  be an arbitrary pair of parsing and regular expressions, in  $\text{WF2}(\text{Lab}_1 \cup \text{Lab}_2)$ . Then the following will hold:

1. Let  $s$  be an arbitrary string, and  $m$  be an arbitrary match in  $\text{regMatch}(\text{pegreg}(p, r), s)$ . Then  $m.\text{string} = \text{extract}_{\text{peg}}(p, m).\text{string} \cdot \text{extract}_{\text{reg}}(r, m).\text{string}$ .
2. Let  $s$  be an arbitrary string. Suppose  $\text{pegMatch}(p, s) = \perp$ . Then the set  $\text{regMatch}(\text{negate}(p), s)$  contains at least one element, and for all matches  $m$  in  $\text{regMatch}(\text{negate}(p), s)$ , the following holds:  $m.\text{string} = \text{extract}_{\text{peg}}(p, m).\text{string} = s$ .
3. Let  $s$  be an arbitrary string. Suppose  $\text{pegMatch}(p, s) = (s', k')$ , and  $\text{regMatch}(r, k)$  is nonempty. Then  $\text{regMatch}(\text{pegreg}(p, r), s)$  is nonempty.
4. Let  $s$  be an arbitrary string. Suppose  $\text{regMatch}(\text{negate}(p), s)$  is not empty. Then it contains one element, and  $\text{pegMatch}(p, s) = \perp$ .
5. Let  $s$  be an arbitrary string. For all matches  $m \in \text{regMatch}(\text{pegreg}(p, r), s)$ ,  $\text{pegMatch}(p, s) = (\text{extract}_{\text{peg}}(m), \text{extract}_{\text{reg}}(r, m))$  and  $\text{extract}_{\text{reg}}(r, m) \in \text{regMatch}(r, \text{extract}_{\text{reg}}(r, m).\text{string})$ .

In the above proposition, the first property requires that matches from *pegreg* split the string into a *peg* portion and a *reg* portion. The rest of the properties give forward and backwards correctness for *negate* and *pegreg*, in that order.

*Proof.* I think we can do a proof by induction. I'll have to take the leap and try them out. □