# Intrinsic Verification of a Regular Expression Matcher

### Joomy Korkut, Maksim Trifunovski, Daniel R. Licata<sup>\*</sup> Wesleyan University

January 15, 2016

#### Abstract

Harper's 1999 Functional Pearl on regular expression matching is a strong example of the interplay between programming and proof, and has been used for many years in introductory functional programming classes. In this paper, we revisit this algorithm from the point of view of dependently typed programming. In the process of formalizing the algorithm and its correctness using the Agda proof assistant, we found three interesting variations. First, defunctionalizing the matcher allows Agda to see termination without an explicit metric, and provides a simple first-order matcher with a clear relationship to the original, giving an alternative to a later Educational Pearl by Yi. Second, intrinsically verifying the soundness of the algorithm has useful computational content, allowing the extraction of matching strings from the parse tree. Third, while Harper uses a negative definition of *standard* regular expressions (no starred subexpression accepts the empty string), using a syntactic definition of standardness simplifies the staging of the development. These variations provide a nice illustration of the benefits of thinking in a dependently typed language, and have some pedagogical value for streamlining and extending the presentation of this material.

<sup>\*</sup>This material is based on research sponsored in part by by The United States Air Force Research Laboratory under agreement number FA9550-15-1-0053. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the United States Air Force Research Laboratory, the U.S. Government or Carnegie Mellon University.

# Contents

1	Introduction	<b>2</b>	
	1.1 Agda Definitions	5	
<b>2</b>	Syntactically standard regular expressions	6	
3	Defunctionalized intrinsic matcher	8	
	3.1 Definition $\ldots$	9	
	3.1.1 Base cases $\ldots$	9	
	3.1.2 Concatenation $\ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots$	10	
	3.1.3 Alternation $\ldots$	11	
	3.1.4 Kleene plus	11	
	3.1.5 Acceptance	13	
	3.2 Completeness	13	
4	Higher-order intrinsic matcher	15	
-	4.1 Definition	15	
	4.1.1 Base cases	16	
	4.1.2 Concatenation	16	
	4.1.3 Alternation	17	
	4.1.4 Kleene plus	17	
	4.1.5 Acceptance	18	
	4.2 Completeness	18	
5	Matching non-standard regular expressions	20	
0	5.1 Capturing groups	23	
		20	
6	Conclusion	<b>24</b>	
Re	References		

# 1 Introduction

Regular expression matching is a venerable problem, studied in many programming and theory of computation courses. Harper (Harper, 1999) presents a higher-order algorithm for regular expression matching, using continuationpassing to store the remainder of a matching problem when a concatenation is encountered, while using the ordinary control stack to represent the branching when an alternation is encountered. The code for the matcher is quite short, but also quite subtle; the emphasis of Harper's paper is on how the correctness proof for the matcher informs the reader's understanding of the code. For example, the first matcher presented has a termination bug, which is revealed when the induction in the correctness proof breaks down. The problem can be fixed by restricting the domain of the function to *standard* regular expressions, which have no Kleene-stared subexpressions that accept the empty string, and then using a preprocessor translation to solve the original problem. Harper's algorithm has been used in first- and second-year functional programming courses at Carnegie Mellon for around 20 years, as a high-water example of integrating programming and program verification. A later paper by Yi (Yi, 2006) revisits the example, motivated by the author's sense that the higher-order matcher is too difficult for students in their introductory programming course, and gives a first-order matcher based on compilation to a state machine.

Motivated by its strong interplay between programming and proof and its pedagogical usefulness, we set out to formalize Harper's algorithm using the Agda proof assistant (Norell, 2007), believing that it could be a pedagogically useful example of dependently typed programming. The process of mechanizing the algorithm led us to a few new observations that streamline and extend its presentation—which was quite surprising to the third author, who has previously taught this material several times. Our goal in this paper is to document these variations on the matching algorithm, and how the process of programming it in Agda led us to them.

The three variations are as follows. First, because Agda is a total language, we must make the termination of the matcher evident in the code itself. Harper's original algorithm can be shown to terminate using lexicographic induction on first the regular expression and second well-founded induction on the string being matched, but in Agda the latter requires passing an explicit termination measure. We discovered that *defunctionalizing* (Reynolds, 1972) the matcher avoids the explicit termination measure, because the problematic recursive call is moved from the Kleene star case to the character literal base case, where it is clear that the string is getting smaller. The defunctionalized matcher is of interest not only because programming it in Agda is simpler: it also achieves Yi's goal of a first-order matcher in a simple way, which has a clear relationship to the higher-order functions have been introduced, and it could be used as a stepping stone to the more sophisticated higher-order matcher.

Second, the matcher discussed in Harper's paper, whose Agda type is

 $\_accepts\_$  : RegExp  $\rightarrow$  String  $\rightarrow$  Bool

determines whether or not a string is accepted by a regular expression. However, for most applications of regexp matching (and for making compelling homework assignments), it is useful to allow a "bracket" or "grouping" construct that allows the user to specify sub-regular-expressions whose matching strings should be reported— for example, AG [.\*] TC [(G $\oplus$ C) \*] GA for extracting the parts of a DNA string surrounded by certain signal codes. When coding a program/proof in a dependently typed language, there is a choice between "extrinsic" verification (write the simply-typed code, and then prove it correct) and "intrinsic verification" (fuse the program and proof into one function, with a dependent type). We have formalized both a straightforward extrinsic verification, and an intrinsically *sound* verification<sup>1</sup>, which has the dependent type

#### accepts-intrinsic : (r : RegExp) $\rightarrow$ (s : List Char) $\rightarrow$ Maybe (s $\in$ L r)

When this matcher succeeds, it returns the derivation that the string is in the language of the regexp; completeness, which says that the matcher does not improperly reject strings, is still proved separately. The reason for this choice is that the *computational content* of the soundness proof is relevant to the above problem: the derivation gives a parse tree, which allows reporting the matching strings for each specified sub-expression. Indeed, even for the extrinsic matcher, one can run the separate soundness proof to produce the matching strings—but running the matcher and then its soundness proof (which has success of the matcher as a precondition) duplicates work, so we present the intrinsic version in the paper. Though we were led to this variation by coding the soundness proof of the matcher using dependent types, analogous code could be used in a simply typed language, with the less informative result type Maybe Derivation which does not say what string and regexp it is a derivation for.

A third variation is that, while Harper uses a negative semantic definition of standard regular expressions ("no subexpression of the form  $r^*$ accepts the empty string"), it is often more convenient in Agda to use positive/inductive criteria. While formalizing the notion of standard, we realized that it is possible to instead use a syntactic criterion, generating standard regular expressions by literals, concatenation, alternation, and Kleene *plus* (one or more occurrences) instead of Kleene *star* (zero or more occurrences), and omitting the empty string regexp  $\varepsilon$ . While the syntactic criterion omits some semantically standard expression (such as  $(\varepsilon \cdot r)$ \*, where r does not

<sup>&</sup>lt;sup>1</sup>All formalizations are available from http://github.com/joom/regexp-agda. Use Agda version 2.4.2.2 with standard library version 0.11

accept the empty string), it still suffices to define a matcher for all regular expression. In addition to simplifying the Agda formalization, this observation has the pedagogical benefit of allowing a self-contained treatment of these syntactically standard regular expressions.

Though dependently typed programming led us to these new insights into a problem that has been very thoroughly studied from a very similar point of view, they can all be ported back to simply-typed languages. Thus, in addition to being a strong pedagogical example of dependently typed programming, these variations on regular expression matching could be used in introductory programming courses to offer a streamlined treatmente.g. using the defunctionalized matcher for only syntactically standard regular expressions, which still captures the basic interplay between programming and proof—which scales to higher-order matching and more interesting homework assignments—e.g. by computing matching strings. Therefore, even though there is existing work on verified parsing of regular expressions and context-free grammars (see (Danielsson, 2010; Ridge, 2011; Firsov & Uustalu, 2013) for work using other algorithms; Wouter Swierstra and collaborators also worked on formalizing Harper's algorithm, but had unresolved termination issues<sup>2</sup>), we believe these new variations on Harper's algorithm will be of interest to the dependent types and broader functional programming communities.

The remainder of this paper is organized as follows. In Section 2, we define syntactically standard regular expressions. In Section 3, we give an intrinsically sound defunctionalized matcher, with no explicit termination measure. In Section 4, we give an intrinsically sound higher-order matcher, which uses an explicit termination measure, and explain the correspondence with the defunctionalized matcher. In Section 5, we show that these matchers suffice to match all regular expressions by translation. In Section 5.1 we discuss how to extract matching strings.

### 1.1 Agda Definitions

We assume the reader is familiar with Agda (Norell, 2007) and the notation from the Agda standard library. Though the Agda library differentiates between strings (String) and lists of characters (List Char), we will ignore this distinction in the paper. Because the intrinsically verified matcher returns a Maybe/option type, it will be useful to use monadic notation for Maybe (in Haskell terminology, \_||\_\_ is mplus and map is fmap):

<sup>&</sup>lt;sup>2</sup>personal communication

## 2 Syntactically standard regular expressions

Rather than Harper's negative semantic criterion (no starred subexpression accepts the empty string), we use an inductive definition of standard regular expressions. Compared with the standard grammar of regular expressions, which includes the regexp matching the empty string ( $\epsilon$ ), the regexp matching the empty language ( $\emptyset$ ), character literals, concatenation ( $r_1 \cdot r_2$ ), alternation ( $r_1 \oplus^{s} r_2$ ), and Kleene star/repetition ( $r^*$ ), we omit  $\epsilon$ , and replace Kleene star with Kleene plus, which represents repetition one or more times. As an operator on regular languages, Kleene plus ( $\Sigma^+$ ) is equivalent to  $\Sigma \cdot \Sigma^*$ , where  $\Sigma^*$  is the Kleene star.

In Agda, we define a type of StdRegExp as follows:

 $\emptyset^{s}$  is the regular expression matching no strings, Lit<sup>s</sup> is the character literal, \_.\*\_\_ is concatenation, \_ $\oplus^{s}$ \_\_ is alternation, and \_\_+s is Kleene plus. We use the <sup>s</sup> superscript to differentiate standard regular expressions from the the full language, which is defined in Section 5. Informally, Kleene plus is defined as the least language closed under the following rules:

$$\frac{s \in L(\mathsf{r})}{s \in L(\mathsf{r}^{+\mathsf{s}})} \qquad \frac{s_1 s_2 = s \quad s_1 \in L(\mathsf{r}) \quad s_2 \in L(\mathsf{r}^{+\mathsf{s}})}{s \in L(\mathsf{r}^{+\mathsf{s}})}$$

In Agda, we represent sets of strings by something analogous to their membership predicates. For standard regular expressions, we define a binary relation  $s \in L^s r$ , which means s is in the language of r, by recursion on r, illustrating how it is possible to compute types based on values in a dependently typed language. This uses an auxiliary, inductively definition relation  $s \in L^+ r$ , represented by an Agda inductively defined datatype family, which corresponds to the inference rules above.

#### mutual

$$\begin{array}{l} \_ \in L^s \_: \mbox{List Char} \to \mbox{StdRegExp} \to \mbox{Set} \\ \_ \in L^s \oslash^s = \bot \\ s \in L^s \ (\mbox{List}^s \ c) = s \equiv [c] \\ s \in L^s \ (r_1 \oplus^s r_2) = (s \in L^s \ r_1) \uplus (s \in L^s \ r_2) \\ s \in L^s \ (r_1 \cdot^s \ r_2) = \\ \Sigma \ (\mbox{List Char} \times \mbox{List Char}) \ (\lambda \ \{(p,q) \to (p \ \# \ q \equiv s) \times (p \in L^s \ r_1) \times (q \in L^s \ r_2)\}) \\ s \in L^s \ (r^{+s}) = s \in L^+ \ r \\ \hline \mbox{data} \ \_ \in L^+ \ \_: \ \mbox{List Char} \to \mbox{StdRegExp} \to \mbox{Set where} \\ S+ : \ \forall \ \{s \ r\} \to s \in L^s \ r \to s \in L^+ \ r \\ \hline \ C+ : \ \forall \ \{s \ s_1 \ s_2 \ r\} \to s_1 \ \# \ s_2 \equiv s \to s_1 \in L^s \ r \to s_2 \in L^+ \ r \to s \in L^+ \ r \end{array}$$

Here,  $\perp$  is the empty Agda type,  $\uplus$  is disjoint union (Either), and  $\times$  is the pair type.  $\Sigma \land (\lambda \times \rightarrow B)$  is an existential/dependent pair, where the type of the second component depends on the value of the first—because this is not built in in Agda, it takes a type  $\land$  and a function from  $\land$  to Set as arguments. The notation  $\lambda \lbrace p \rightarrow e \rbrace$  allows a pattern-matching anonymous function. The function # appends lists, and  $\equiv$  is Agda's propositional equality. Thus, in full, the clause for alternation means "there exist strings p and q such that appending p and q gives s, where p is in the language of  $r_1$  and q is in the language of  $r_2$ ".

However, is is important to note that these membership "predicates" land in Set, the Agda type of types, and thus may have computational content. For example, a witness that  $s \in L^s (r_1 \oplus r_2)$  includes a bit  $(inj_1 \text{ or } inj_2)$  that tells which possibility was taken, and a witness  $s \in L^+ r$  is a non-empty list of strings matching r, which concatenate to s. Thus, there can be different witnesses that a string matches a regular expression, such as

Derivation A := 
$$\frac{a^{"} \in L(\mathsf{Lit}^{\mathsf{s}}, \mathsf{a}^{\mathsf{s}})}{a^{"} \in L(\mathsf{Lit}^{\mathsf{s}}, \mathsf{a}^{\mathsf{s}})}$$
Derivation B := 
$$\frac{a^{"} + a^{"} = aaa^{"} + aa^{"} \in L(\mathsf{Lit}^{\mathsf{s}}, \mathsf{a}^{\mathsf{s}}) \quad \text{Derivation A}}{aa^{"} \in L(\mathsf{Lit}^{\mathsf{s}}, \mathsf{a}^{\mathsf{s}})}$$
Derivation C := 
$$\frac{a^{"} + aa^{"} = aaa^{"} \quad \text{Derivation A}}{aaa^{"} \in L(\mathsf{Lit}^{\mathsf{s}}, \mathsf{a}^{\mathsf{s}})}$$

 $\text{Derivation D} := \frac{``aa" + ``a" = ``aaa"}{``aaa"} \xrightarrow{\text{Derivation B}} \xrightarrow{\text{Derivation A}} \xrightarrow{\text{Derivation A}} \xrightarrow{(\texttt{Lits'}, \texttt{a'}, \texttt{+s})} \xrightarrow{(\texttt{Lits'}, \texttt{a'}, \texttt{+s})} \xrightarrow{(\texttt{Lits'}, \texttt{a'}, \texttt{+s})}$ 

We will exploit this in Section 5.1 to extract matching strings from such derivations.

# 3 Defunctionalized intrinsic matcher

In this section, we define a first-order matcher for standard regular expressions. This is a defunctionalization of Harper's algorithm, though we will describe it from first principles. The idea is to generalize from matching a string s against a regular expression r by adding an additional stack of regular expressions k that need to be matched against the suffix of s if some prefix of s matches r. We represent the stack by a list, and say that a string is in the language of a stack if it splits into strings in the language of each stack element:

$$\begin{array}{l} \_\in L^{k}\_: \mbox{ List Char} \rightarrow \mbox{ List StdRegExp} \rightarrow \mbox{ Set} \\ s \in L^{k} [] = s \equiv [] \\ s \in L^{k} (r :: rs) = \\ \Sigma (\mbox{ List Char} \times \mbox{ List Char}) \\ (\lambda \left\{ (p,s') \rightarrow (p \ \mbox{ + } s' \equiv s) \times (p \in L^{s} r) \times (s' \in L^{k} rs) \right\} ) \end{array}$$

If the stack is empty, the string also has to be empty. If the stack has a head, then a prefix of the string should match the head of the list and the rest of the string should match with the rest of the list.

#### 3.1 Definition

The soundness part of informal description of the matcher given above translates into the following dependent type/specification:

$$\begin{array}{l} \mathsf{match} \ : \ (\mathsf{r} \ : \ \mathsf{StdRegExp}) \ (\mathsf{s} \ : \ \mathsf{List} \ \mathsf{Char}) \ (\mathsf{k} \ : \ \mathsf{List} \ \mathsf{StdRegExp}) \\ & \to \mathsf{Maybe} \ (\Sigma \ (\mathsf{List} \ \mathsf{Char} \times \ \mathsf{List} \ \mathsf{Char}) \\ & \qquad (\lambda \ \{(\mathsf{p},\mathsf{s}') \to (\mathsf{p} \ + \ \mathsf{s}' \equiv \mathsf{s}) \times (\mathsf{p} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r}) \times \mathsf{s}' \in \mathsf{L}^{\mathsf{k}} \ \mathsf{k}\})) \end{array}$$

This says that the matcher takes a regular expression r, a string s and a stack k, and, if matching succeeds, returns a splitting  $(p + s' \equiv s)$ , a derivation that p is in the language of r and a derivation that s' is in the language of the stack k. This specifies the soundness of successful matching, and has computational content that will let us extract matching strings, but leaves open the possibility that the matcher fails incorrectly (for example, the constantly nothing function has this type)—this will be addressed in the completeness proof below.

The complete code is in Figure 1. Agda can verify that this function terminates by induction on first the string s, and then the regular expression r. We now discuss the code case by case.

#### 3.1.1 Base cases

First, the empty language does not accept any string:

match  $\varnothing^{s} s k = fail$ 

For character literals

$$\begin{array}{ll} \mathsf{match}\;(\mathsf{Lit}^{s}\;c)\;[]\;k\;=\;\mathsf{fail}\\ \mathsf{match}\;(\mathsf{Lit}^{s}\;c)\;(x::xs)\;k\;=\\ (\mathsf{isEqual}\;x\;c) \gg \\ & (\lambda\;p\to\mathsf{map}\;(\lambda\;\mathsf{pf}\to((([c],xs),\mathsf{cong}\;(\lambda\;x\to x::xs)\;(\mathsf{sym}\;p),\mathsf{refl},\mathsf{pf})))\\ & (\mathsf{match-helper}\;k\;xs)) \end{array}$$

in the first case, if we are trying to match an empty list with a regular expression that requires a character, the matcher fails. In the next case, the isEqual x c call has type Maybe ( $x \equiv c$ )—i.e. it optionally shows that x is equal to c. Thus, by the monad bind, when x is not c, the matcher fails, and when x is c, we try to match the stack k against the suffix xs using match-helper.

The function **match-helper** is mutually recursive with our matcher and is defined as follows:

 $\begin{array}{lll} {\rm match-helper}\,:\,(k\,:\,List\,StdRegExp)\rightarrow(s\,:\,List\,Char)\rightarrow Maybe\,(s\in L^k\,k)\\ {\rm match-helper}\,[\,]\,[\,]\,=\,return\,refl\\ {\rm match-helper}\,[\,]\,(x\,::\,s)\,=\,fail\\ {\rm match-helper}\,(r\,::\,k')\,s\,=\,match\,r\,s\,k' \end{array}$ 

It succeeds when matching the empty string against the empty stack, fails when matching a non-empty string against the empty stack, and otherwise refers back to **match**. Relative to Harper's algorithm, this is the application function for the defunctionalized continuation.

Returning to the character literal case, if the first character in our list matches the character literal c, then we call match-helper on the continuation and the rest of the list, which will produce a derivation of  $s \in L^k$  k if the rest of the list indeed matches the rest of the StdRegExps in k. The remainder of the code packages this as a pair showing that therefore the list x :: xs splits as  $[x] \in L^s$  (Lit c) and  $xs \in L^k$  k.

Agda's termination checker is able to verify that, for the call from match to match-helper and back to match, the string x :: xs becomes xs, and is therefore smaller, which justifies termination in this case.

#### 3.1.2 Concatenation

$$\begin{array}{l} \mathsf{match} \ (\mathsf{r}_1 \cdot {}^{\mathsf{s}} \ \mathsf{r}_2) \ \mathsf{s} \ \mathsf{k} \ = \\ \mathsf{map} \ (\mathsf{reassociate-left} \ \{\mathsf{R} \ = \ \_ {}^{\mathsf{s}}\_\} \ (\lambda \ \mathsf{inL} \ \mathsf{inL'} \rightarrow \_, \mathsf{refl}, \mathsf{inL}, \mathsf{inL'})) \\ (\mathsf{match} \ \mathsf{r}_1 \ \mathsf{s} \ (\mathsf{r}_2 \ :: \ \mathsf{k})) \end{array}$$

In the case for concatenation  $r_1 \cdot {}^s r_2$ , we match against  $r_1$  first, and add  $r_2$  to the stack k. Calling match  $r_1 s (r_2 :: k)$  will give us a split  $xs + ys \equiv s$  and derivations of  $xs \in L^s r_1$  and  $ys \in L^k (r_2 :: k)$ . If we unpack the second derivation (using the definition of  $\in L^k$  and the fact that our continuation list contains at least one element,  $r_2$ ), we will have another split  $as + bs \equiv ys$  and derivations  $as \in L^s r_2$  and  $bs \in L^k k$ . The helper function reassociate-left states that if we have such a situation, we can reassociate it to show that xs + as matching the entire regular expression  $r_1 \cdot {}^s r_2$ . Because we will want to do similar reasoning in the Kleene plus case below, we use a higher-order function that says that if R is a binary operation on regular expressions that respects splitting, then given the first kind of splitting, we can produce the second:

```
\begin{array}{l} \mathsf{reassociate-left} \ : \ \forall \ \{\mathsf{r}_1 \ \mathsf{r}_2 \ \mathsf{s} \ \mathsf{k}\} \ \{\mathsf{R} \ : \ \mathsf{StdRegExp} \rightarrow \mathsf{StdRegExp} \rightarrow \mathsf{StdRegExp} \\ \rightarrow (\mathsf{f} \ : \ \forall \ \{\mathsf{xs} \ \mathsf{as}\} \rightarrow \mathsf{xs} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r}_1 \rightarrow \mathsf{as} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r}_2 \rightarrow ((\mathsf{xs} \ \ \mathsf{+} \ \mathsf{as}) \in \mathsf{L}^{\mathsf{s}} \ \mathsf{R} \ \mathsf{r}_1 \ \mathsf{r}_2)) \end{array}
```

$$\begin{split} & \rightarrow \Sigma_{-} \left( \lambda \left\{ (xs, ys) \rightarrow (xs \ \# \ ys \equiv s) \times xs \in L^{s} r_{1} \right. \\ & \times \Sigma_{-} \left( \lambda \left\{ (as, bs) \rightarrow (as \ \# \ bs \equiv ys) \times as \in L^{s} r_{2} \times bs \in L^{k} k \right\} \right) \right\} ) \\ & \rightarrow \left( \Sigma_{-} \left( \lambda \left\{ (p, s') \rightarrow (p \ \# \ s' \equiv s) \times (p \in L^{s} \ R \ r_{1} \ r_{2}) \times s' \in L^{k} \ k \right\} \right) \right) \end{split}$$

#### 3.1.3 Alternation

```
match (r_1 \oplus^s r_2) s k =
(map (change-\in L inj_1) (match r_1 s k)) ||
(map (change-\in L inj_2) (match r_2 s k))
```

In the alternation case, we match the string with  $r_1$ , and if that fails match with  $r_2$  (recall that || handles failure of its first disjunct by trying the second). If the call match  $r_1 \ s \ k$  succeeds, it will produce a triple splitting s as  $p \ \# \ s'$ , with a derivation of  $p \in L^s \ r_1$  and  $s' \in L^k \ k$ . However the return type for the alternation case should contain a derivation of type  $p \in L^s \ (r_1 \oplus^s r_2)$ , so we use the helper function change- $\in L$ , which applies a function to this position of the result triple, to make the appropriate modification:

 $\begin{array}{l} \mathsf{change-} \in \mathsf{L} \ : \ \{ \mathsf{a} \ \mathsf{b} \ \mathsf{d} \ : \ \mathsf{List} \ \mathsf{Char} \to \mathsf{Set} \} \ \{ \mathsf{c} \ : \ \mathsf{List} \ \mathsf{Char} \to \mathsf{List} \ \mathsf{Char} \to \mathsf{Set} \} \\ \to (\forall \ \{ \mathsf{s} \} \to \mathsf{a} \ \mathsf{s} \to \mathsf{b} \ \mathsf{s}) \\ \to (\Sigma_{-} (\lambda \ \{ (\mathsf{p},\mathsf{s}') \to (\mathsf{c} \ \mathsf{p} \ \mathsf{s}') \times (\mathsf{a} \ \mathsf{p}) \times (\mathsf{d} \ \mathsf{s}') \})) \\ \to (\Sigma_{-} (\lambda \ \{ (\mathsf{p},\mathsf{s}') \to (\mathsf{c} \ \mathsf{p} \ \mathsf{s}') \times (\mathsf{b} \ \mathsf{p}) \times (\mathsf{d} \ \mathsf{s}') \})) \\ \mathsf{change-} \in \mathsf{L} \ \mathsf{f} \ (\mathsf{x},\mathsf{eq},\mathsf{inL},\mathsf{rest}) = \ (\mathsf{x},\mathsf{eq},\mathsf{finL},\mathsf{rest}) \end{array}$ 

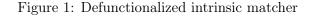
We use change- $\in L$  to apply inj<sub>1</sub> or inj<sub>2</sub> to the derivation, depending on which part of the alternation successfully matched the string.

#### 3.1.4 Kleene plus

```
 \begin{array}{l} \mathsf{match}\;(\mathsf{r}^{+\mathsf{s}})\;\mathsf{s}\;\mathsf{k}\;=\;\\ (\mathsf{map}\;(\mathsf{change-}\!\in\!\mathsf{L}\;\mathsf{S+})\;(\mathsf{match}\;\mathsf{r}\;\mathsf{s}\;\mathsf{k}))\;||\\ (\mathsf{map}\;(\mathsf{reassociate-left}\;\{\mathsf{R}\;=\;\lambda\;\mathsf{r}\;\_\to\;\mathsf{r}^{+\mathsf{s}}\}\;(\lambda\;\mathsf{inL}\;\mathsf{inL'}\to\mathsf{C+}\;\mathsf{refl}\;\mathsf{inL}\;\mathsf{inL'}))\\ (\mathsf{match}\;\mathsf{r}\;\mathsf{s}\;((\mathsf{r}^{+\mathsf{s}})\;::\;\mathsf{k}))) \end{array}
```

In the Kleene plus case, we first try to match s with just r, and if that succeeds we apply change- $\in L S+$  to the derivation since we matched from the single r case. If this fails, then similar to the  $\cdot^{s}$  case, we try to match a prefix of the string to r and then the suffix that follows with the continuation which now includes  $r^{+s}$ . Just like in the  $\cdot^{s}$  case, we use reassociate-left in order to get that our splitting of s matches the entire starting r.

```
mutual
   match : (r : StdRegExp) (s : List Char) (k : List StdRegExp)
      \rightarrow Maybe (\Sigma (List Char \times List Char)
                         (\lambda \{ (p, s') \rightarrow (p + s' \equiv s) \times (p \in L^{s} r) \times s' \in L^{k} k \}))
   match \varnothing^s s k = fail
   match (Lit<sup>s</sup> c) [] k = fail
   match (Lit<sup>s</sup> c) (x :: xs) k =
      (isEqual x c) \gg
          (\lambda p \rightarrow map \ (\lambda pf \rightarrow ((([c], xs), cong \ (\lambda x \rightarrow x :: xs) \ (sym p), refl, pf)))
             (match-helper k xs)
   match (r_1 \cdot s r_2) s k =
      map (reassociate-left { R = \_.^{s} } (\lambda inL inL' \rightarrow _, refl, inL, inL'))
          (match r_1 s (r_2 :: k))
   match (r_1 \oplus^s r_2) s k =
      (map (change \in L inj_1) (match r_1 s k)) \parallel
      (map (change-\in L inj_2) (match r_2 s k))
   match (r^{+s}) s k =
      (map (change \in L S+) (match r s k)) \parallel
      (\text{map (reassociate-left } \{ \mathsf{R} \ = \ \lambda \ r \ \_ \rightarrow r \ ^{+s} \} \ (\lambda \ \text{inL inL'} \rightarrow \mathsf{C} + \text{refl inL inL'}))
          (match r s ((r + s) :: k)))
   match-helper : (k : List StdRegExp) \rightarrow (s : List Char) \rightarrow Maybe (s \in L<sup>k</sup> k)
   match-helper [] [] = return refl
   match-helper [](x :: s) = fail
   match-helper (r :: k') s = match r s k'
```



#### 3.1.5 Acceptance

From the outside, we can call the generalized matcher with the empty stack to check membership:

```
accepts<sup>s</sup>-intrinsic : (r : StdRegExp) \rightarrow (s : List Char) \rightarrow Maybe (s \inL<sup>s</sup> r) accepts<sup>s</sup>-intrinsic r s = map \inL-empty-stack (match r s [])
```

When the **match** function succeeds on an empty stack, the suffix is in the language of an empty stack and is therefore an empty string, so we use the following lemma to change the result of the function call **match** r s [] into a derivation over the entire string s.

```
 \begin{array}{l} \in \mathsf{L\text{-empty-stack}} \ : \ \{r \ : \ \mathsf{StdRegExp}\} \ \{s \ : \ \mathsf{List} \ \mathsf{Char}\} \\ \rightarrow \Sigma \ \_ \ (\lambda \ \{(\mathsf{p}, \mathsf{s}') \rightarrow (\mathsf{p} \ + \ \mathsf{s}' \equiv \mathsf{s}) \times (\mathsf{p} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r}) \times (\mathsf{s}' \equiv [])\}) \\ \rightarrow \mathsf{s} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r} \end{array}
```

#### 3.2 Completeness

Though the above matcher is intrinsically sound, it is not intrinsically complete for example, the function that always fails has the above type. One way to resolve this would be to write a matcher that is intrinsically both sound and complete — ignoring the stack, given **s** and **r**, we would like to know  $(\mathbf{s} \in \mathsf{L}^{\mathsf{s}} \mathsf{r}) \uplus \neg (\mathbf{s} \in \mathsf{L}^{\mathsf{s}} \mathsf{r})$ , a type that expresses decidability of matching. However, while there is an efficiency reason to intrinsically compute the the derivation of  $\mathsf{s} \in \mathsf{L}^{\mathsf{s}} \mathsf{r}$ —we will use it to extract matching strings in Section 5.1—there is no computational content to  $\neg (\mathsf{s} \in \mathsf{L}^{\mathsf{s}} \mathsf{r})$ . Because of this, and because it keeps the matcher code itself simpler, we choose to make completeness extrinsic. In full, completeness says that if we have **r**, **s**, **k** and a split of  $\mathsf{p} \ + \ \mathsf{s}' \equiv \mathsf{s}$  such that there are derivations of type  $\mathsf{p} \in \mathsf{L}^{\mathsf{s}} \mathsf{r}$  and  $\mathsf{s}' \in \mathsf{L}^{\mathsf{k}} \mathsf{k}$ , then we know that our match function does not fail:

```
 \begin{array}{l} \mathsf{match-completeness} \ : \ (\mathsf{r} \ : \ \mathsf{StdRegExp}) \ (\mathsf{s} \ : \ \mathsf{List} \ \mathsf{Char}) \ (\mathsf{k} \ : \ \mathsf{List} \ \mathsf{StdRegExp}) \\ & \rightarrow \Sigma_{-} \left( \lambda \ \{ (\mathsf{p},\mathsf{s}') \rightarrow (\mathsf{p} \ + \ \mathsf{s}' \equiv \mathsf{s}) \times (\mathsf{p} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r}) \times (\mathsf{s}' \in \mathsf{L}^{\mathsf{k}} \ \mathsf{k}) \} \right) \\ & \rightarrow \mathsf{isJust} \ (\mathsf{match} \ \mathsf{r} \ \mathsf{s} \ \mathsf{k}) \end{array}
```

That is, when there is a way for the matcher to succeed, it does. We cannot make a stronger claim and say that it returns the same derivation that is given as input, because as we showed before, there can be different derivations for the same s and r, and the given one may not be the one the matcher finds.

The full proof is in the companion code. The base cases  $\emptyset^s$  and Lit<sup>s</sup> are easy. Since the Kleene plus case captures the essence of both concatenation and alternation cases, we will explain only this case. The proof proceeds by cases, depending on how the given derivation of  $p \in L(r^+)$  was constructed. For the first,

match-completeness (r <sup>+s</sup>) s k ((xs, ys), eq, S+ inL, rest) with match r s k | match-completeness r s k ((xs, ys), eq, inL, rest) ... | nothing | () ... | just \_ | \_ = tt

if the given derivation of  $xs \in L(r^+)$  was by the constructor S+, then sitting under the constructor is a derivation of  $inL : xs \in L^s r$ , and the result follows from the inductive hypothesis on D.

For the second, where the derivation was constructed by C+,

```
\begin{array}{l} \text{match-completeness } (r^{+s}) \circ ((s_1 \ + \ s_2) \ + \ ys) \ k \\ ((.\_, ys), \text{refl}, C+ \{ .\_ \} \ \{ s_1 \} \ \{ s_2 \} \ \text{refl inL1 inL2, rest}) \\ \text{with match } r \ ((s_1 \ + \ s_2) \ + \ ys) \ k \\ \dots \ \mid \ \text{just} \ \_ \ = \ tt \\ \dots \ \mid \ \text{nothing} \\ \text{with match } r \ ((s_1 \ + \ s_2) \ + \ ys) \ ((r^{+s}) \ :: \ k) \\ \ \mid \ \text{match-completeness } r \ ((s_1 \ + \ s_2) \ + \ ys) \ ((r^{+s}) \ :: \ k) \\ (\_, \text{append-assoc } s_1 \ s_2 \ ys, \text{inL1}, (\_, ys), \text{refl, inL2, rest}) \\ \dots \ \mid \ \text{nothing} \ \mid () \\ \dots \ \mid \ \text{just} \ \_ \ \mid \ = \ tt \end{array}
```

we already had a split  $xs + ys \equiv s$ , so we can replace s with xs + ys by pattern matching on the equality proof. The constructor C+ gives us another split  $s_1 + s_2 \equiv xs$ , so we can also replace xs with  $s_1 + s_2$ . This means now we have to show the goal for  $(s_1 + s_2) + ys$  instead of s. The definition of the Kleene plus case uses  $||_{-}$ , which has to try and fail the first case to return the second case. To verify completeness, we first check whether match  $r((s_1 + s_2) + ys) k$  succeeds or fails. If the call succeeds, then we satisfy the first case of  $||_{-}$ , so the matcher succeeds. If the call fails, then we checker whether the second disjunct fails or succeeds. If it fails, then we obtain a contradiction by the inductive hypothesis, which shows that the recursive call should have succeeded because  $s_1$  matches r and  $s_2 + ys$  matches the continuation (using the associative property of appending lists to show that it is the same string). If it succeeds, then the matcher succeeds, so we have the result.

As a corollary, we get completeness of accepts<sup>s</sup>-intrinsic:

```
\begin{array}{ll} \mathsf{accepts}^{\mathsf{s}}\text{-intrinsic-completeness} \ : \ (\mathsf{r} \ : \ \mathsf{StdRegExp}) \ (\mathsf{s} \ : \ \mathsf{List} \ \mathsf{Char}) \\ & \to \mathsf{s} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r} \\ & \to \mathsf{isJust} \ (\mathsf{accepts}^{\mathsf{s}}\text{-intrinsic} \ \mathsf{r} \ \mathsf{s}) \end{array}
```

# 4 Higher-order intrinsic matcher

In this section, we show that the above intrinsic verification of the first-order matcher scales to a higher-order matcher, written using continuation-passing, which is more similar to Harper's original code. We will use this to explain why the above matcher is a defunctionalization, and why the termination reasoning is more difficult in the higher-order case. To make termination evident to Agda, we will need to use an explicit termination metric that corresponds to well-founded induction on strings/lists. This is represented in Agda by an iterated inductive definition RecursionPermission xs. Visually, you can think of RecursionPermission ys as a tree, where a node for ys has subtrees for each strict suffix of ys. Each of these subtrees is judged smaller by the termination checker, and therefore we will be allowed to recur on any suffix of ys. Such a tree type is defined by the following datatype, which has a higher-order constructor argument:

```
data RecursionPermission {A : Set} : List A \rightarrow Set where
CanRec : {ys : List A}
\rightarrow ((xs : List A) \rightarrow Suffix xs ys \rightarrow RecursionPermission xs)
\rightarrow RecursionPermission ys
```

#### 4.1 Definition

The higher-order intrinsic matcher has the following specification:

```
\begin{array}{ll} \mathsf{match} \ : \ (\mathsf{C} \ : \ \mathsf{Set}) \ (\mathsf{r} \ : \ \mathsf{StdRegExp}) \ (\mathsf{s} \ : \ \mathsf{List} \ \mathsf{Char}) \\ & \to (\mathsf{k} \ : \ \forall \ \{\mathsf{p} \ \mathsf{s}'\} \to \mathsf{p} \ \# \ \mathsf{s}' \equiv \mathsf{s} \to \mathsf{p} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r} \to \mathsf{Maybe} \ \mathsf{C}) \\ & \to \mathsf{RecursionPermission} \ \mathsf{s} \\ & \to \mathsf{Maybe} \ \mathsf{C} \end{array}
```

The type variable C stands for the output derivation computed by the matcher on success. Just as Harper's algorithm returns a **bool** and uses both the continuation and the language's control stack (i.e. it is not fully in CPS),

here both the continuation and the matcher return an option, but the success data can been chosen arbitrarily. In Harper's algorithm, the continuation takes a string that corresponds to s' in the above type. The additional arguments provided here, which are important for justifying termination, say that in a call to match  $r \ s \ k$ , the domain of the continuation is suffixes s' of s by a prefix that is in the language of r. The final argument is of type RecursionPermission s, and allows recursive calls on strict suffixes of s. The termination measure for this function is lexicographic in first the regular expression r and then the recursion permission tree. The complete code is in Figure 2.

#### 4.1.1 Base cases

The case for the empty regexp fails:

match C  $\varnothing^{s}$  s k perm = fail

The cases for character literals

 $\begin{array}{l} \mathsf{match}\ \mathsf{C}\ (\mathsf{Lit}^{\mathsf{s}}\ \mathsf{c})\ []\ \mathsf{k}\ \mathsf{perm}\ =\ \mathsf{fail}\\ \mathsf{match}\ \mathsf{C}\ (\mathsf{Lit}^{\mathsf{s}}\ \mathsf{c})\ (\mathsf{x}::\mathsf{xs})\ \mathsf{k}\ \mathsf{perm}\ =\\ (\mathsf{isEqual}\ \mathsf{x}\ \mathsf{c})\ \gg\ (\lambda\ \mathsf{p}\ \to\ \mathsf{k}\ \{[\mathsf{x}]\}\ \{\mathsf{xs}\}\ \mathsf{refl}\ (\mathsf{cong}\ (\lambda\ \mathsf{q}\ \to\ [\mathsf{q}])\ \mathsf{p})) \end{array}$ 

are as above, except that where we called match-helper to activate the stack k, here we call k itself. The packaging of the result of match-helper (the map in the above code) now happens as *input* to k, because to call k we must show that x :: xs splits as something in the language of Lit<sup>s</sup> c and some suffix.

#### 4.1.2 Concatenation

```
 \begin{array}{l} \mathsf{match}\; \mathsf{C}\;(\mathsf{r}_1 \cdot^{\mathsf{s}}\;\mathsf{r}_2)\;\mathsf{s}\;\mathsf{k}\;(\mathsf{CanRec}\;\mathsf{f}) \;= \\ \mathsf{match}\; \mathsf{C}\;\mathsf{r}_1\;\mathsf{s} \\ (\lambda\;\{\mathsf{p}\}\;\{\mathsf{s}'\}\;\mathsf{eq}\;\mathsf{inL}\;\rightarrow \\ \mathsf{match}\; \mathsf{C}\;\mathsf{r}_2\;\mathsf{s}'\;(\lambda\;\{\mathsf{p}'\}\;\{\mathsf{s}''\}\;\mathsf{eq}'\;\mathsf{inL}'\;\rightarrow \\ \mathsf{k}\;\{\mathsf{p}\;\;+\;\;\mathsf{p}'\}\;\{\mathsf{s}''\}\;(\mathsf{replace-right}\;\mathsf{p}\;\mathsf{s}'\;\mathsf{p}\;\mathsf{s}\;\mathsf{eq}\;\mathsf{eq})\;((\mathsf{p},\mathsf{p}'),\mathsf{refl},\mathsf{inL},\mathsf{inL}')) \\ (\mathsf{f}\;\mathsf{s}'\;(\mathsf{suffix-after-}\in\mathsf{L}^{\mathsf{s}}\;\mathsf{eq}\;\mathsf{inL})))\;(\mathsf{CanRec}\;\mathsf{f}) \end{array}
```

In the concatenation case of the defunctionalized version, we added  $r_2$  to the stack of continuations to be matched later. In the higher-order version, extending the stack with  $r_2$  corresponds to constructing a new continuation function which matches  $r_2$  against the suffix that results from matching  $r_1$  against a prefix—which is exactly what "applying" the stack in match-helper did in the defunctionalized version. The massaging that happens after the recursive call above (the reassociate-left) here happens in the new continuation, which repackages the given derivations inL :  $p \in L^s r_1$  and inL' :  $p' \in L^s r_2$ and the given splittings eq :  $p + s' \equiv s$  and eq' :  $p' + s'' \equiv s'$  as a splitting  $(p + p') + s'' \equiv s$  and a derivation that  $p + p' \in \in L^s (r_1 \cdot s r_2)$ . The two recursive calls pass the termination checker because the regular expressions  $r_1$  and  $r_2$  get smaller in each case. To make the inner recursive call, it is necessary to supply a recursion permission for s', i.e. to allow recursive calls on s', and to do this it suffices to show that s' is a suffix of s. The suffix-after- $\in L^s$  lemma

suffix-after- $\in L^s$  :  $\forall$  { p s' s r}  $\rightarrow$  (p + s'  $\equiv$  s)  $\rightarrow$  (p  $\in L^s$  r)  $\rightarrow$  Suffix s' s

does this: s' is a non-strict suffix of s by the equality, and because the prefix p is in the language of a *standard* regular expression r and therefore is not empty, it is a strict suffix. (In this case, it would also be sufficient to observe that s' is a non-strict suffix of s, because we do not need the string and recursion permission to get smaller to justify termination, but the argument we just gave will also be used in the Kleene plus case.)

#### 4.1.3 Alternation

 $\begin{array}{l} \mathsf{match}\ \mathsf{C}\ (\mathsf{r}_1\oplus^{\mathsf{s}}\mathsf{r}_2)\ \mathsf{s}\ \mathsf{k}\ \mathsf{perm}\ =\\ \mathsf{match}\ \mathsf{C}\ \mathsf{r}_1\ \mathsf{s}\ (\lambda\ \mathsf{eq}\ \mathsf{inL}\ \to\ \mathsf{k}\ \mathsf{eq}\ (\mathsf{inj}_1\ \mathsf{inL}))\ \mathsf{perm}\ ||\\ \mathsf{match}\ \mathsf{C}\ \mathsf{r}_2\ \mathsf{s}\ (\lambda\ \mathsf{eq}\ \mathsf{inL}\ \to\ \mathsf{k}\ \mathsf{eq}\ (\mathsf{inj}_2\ \mathsf{inL}))\ \mathsf{perm}\ \end{array}$ 

The alternation case is similar to the defunctionalized version, except instead of massaging the derivations after the fact with map change- $\in L$ , we modify them before passing them to the continuation.

#### 4.1.4 Kleene plus

The structure of the Kleene plus case is similar to the defunctionalized version, except the continuations are modified analogously to the concatenation and alternation cases. The first recursive call terminates because r gets smaller. For the second recursive call, (r + s) stays the same, so it is essential that s' is a *strict* suffix of s, and that the recursion permission tree gets smaller. As in the alternation case, s' is a non-strict suffix of s by the equality  $eq : p + s' \equiv s$ , and because the prefix p is in the language of a *standard* regular expression r and therefore is not empty, s' is a strict suffix of s by the suffix-after- $\in L^{s}$  lemma. Therefore, applying f to s' and this fact selects a smaller recursion permission subtree, justifying termination.

Termination is trickier for the higher-order matcher than for the defunctionalized matcher because, here, we make the recursive call on  $r^{+s}$  in the continuation constructed in the  $r^{+s}$  case, so we must argue that whenever this continuation is applied, it will be applied to a smaller string. In the defunctionalized matcher, this recursive call is made in match-helper (the apply function for the defunctionalized continuation), which is called from the character literal case, at which point it is syntactically clear that the recursive call is being made on a smaller string.

#### 4.1.5 Acceptance

Overall, we can define

```
accepts<sup>s</sup>-intrinsic : (r : StdRegExp) \rightarrow (s : List Char) \rightarrow Maybe (s \inL<sup>s</sup> r) accepts<sup>s</sup>-intrinsic r s = match _ r s empty-continuation (well-founded s)
```

by choosing an appropriate initial continuation, and by constructing a recursion permission for s (which exists because string suffix is a well-founded relation). The initial continuation

empty-continuation :  $\forall \{p s' s r\} \rightarrow (p \# s' \equiv s) \rightarrow (p \in L^{s} r) \rightarrow Maybe (s \in L^{s} r)$ 

corresponds to the logic for the empty stack [] in match-helper in the defunctionalized version. This function takes a splitting of a string s as p + s', as well as a proof that its first part p is in the language of r. It returns either nothing if s' is not empty, or just a witness that  $(s \in L^s r)$  if s' is empty, and therefore  $s \equiv p'$ .

#### 4.2 Completeness

Completeness is similar to above, and says that the matcher succeeds whenever it should:

```
match : (C : Set) (r : StdRegExp) (s : List Char)
              (\mathsf{k} \ : \ \forall \ \{ p \ s' \} \rightarrow p \ \# \ s' \equiv s \rightarrow p \in L^s \ r \rightarrow Maybe \ \mathsf{C})
              \rightarrow RecursionPermission s
              \rightarrow Maybe C
\mathsf{match}\;\mathsf{C}\;\varnothing^{\mathsf{s}}\;\mathsf{s}\;\mathsf{k}\;\mathsf{perm}\;=\;\mathsf{fail}
match C (Lit<sup>s</sup> c) [] k perm = fail
match C (Lit<sup>s</sup> c) (x :: xs) k perm =
   (isEqual x c) \gg (\lambda p \rightarrow k \{ [x] \} \{ xs \} refl (cong (\lambda q \rightarrow [q]) p))
match C (r_1 \cdot s r_2) s k (CanRec f) =
   match C r_1 s
       (\lambda \{p\} \{s'\} eq inL \rightarrow
           match C r<sub>2</sub> s' (\lambda { p' } { s'' } eq' inL' \rightarrow
              k \{p + p'\} \{s''\} (replace-right p s' p' s'' s eq' eq) ((p, p'), refl, inL, inL'))
              (f_{-}(suffix-after-\in L^{s} eq inL))) (CanRec f)
match C (r_1 \oplus^s r_2) s k perm =
   match C r<sub>1</sub> s (\lambda eq inL \rightarrow k eq (inj<sub>1</sub> inL)) perm ||
   match C r<sub>2</sub> s (\lambda eq inL \rightarrow k eq (inj<sub>2</sub> inL)) perm
match C (r ^{+s}) s k (CanRec f) =
   match C r s (\lambda eq inL \rightarrow k eq (S+ inL)) (CanRec f) ||
   match C r s (\lambda { p } { s' } eq inL \rightarrow
       match C (r <sup>+s</sup>) s' (\lambda {p'} {s''} eq' inL' \rightarrow
           k (replace-right p s' p' s'' s eq' eq) (C+ refl inL inL'))
          (f_{-}(suffix-after-\in L^{s} eq inL))) (CanRec f)
```

Figure 2: Complete definition of the **match** function for the higher-order intrinsic matcher.

$$\begin{split} & \mathsf{match-completeness} \ : \ (\mathsf{C} \ : \ \mathsf{Set}) \ (\mathsf{r} \ : \ \mathsf{StdRegExp}) \ (\mathsf{s} \ : \ \mathsf{List} \ \mathsf{Char}) \\ & \to (\mathsf{k} \ : \ \forall \ \{\mathsf{p} \ \mathsf{s}'\} \to \mathsf{p} \ + \ \mathsf{s}' \equiv \mathsf{s} \to \mathsf{p} \in \mathsf{L}^{\mathsf{s}} \ \mathsf{r} \to \mathsf{Maybe} \ \mathsf{C}) \\ & \to (\mathsf{perm} \ : \ \mathsf{RecursionPermission} \ \mathsf{s}) \\ & \to \Sigma_{-} (\lambda \ \{(\mathsf{p},\mathsf{s}') \to \Sigma_{-} (\lambda \ \mathsf{eq} \to \Sigma_{-} (\lambda \ \mathsf{inL} \to \mathsf{isJust} \ (\mathsf{k} \ \{\mathsf{p}\} \ \{\mathsf{s}'\} \ \mathsf{eq} \ \mathsf{inL}))) \}) \\ & \to \mathsf{isJust} \ (\mathsf{match} \ \mathsf{C} \ \mathsf{r} \ \mathsf{s} \ \mathsf{perm}) \end{split}$$

The type can be read as follows: Suppose we have C, r, s, k and perm. Suppose there exists a split of  $p + s' \equiv s$  such that there exists a derivation of type  $p \in L^s r$  such that the continuation called with those arguments does not return nothing. Then we have to show that the match function does not fail. Like above, we cannot make a stronger claim and say that the calls to the continuation and the match function return the same derivations, because there can be more than one derivation of a string matching a regexp. The proof is in the companion code, and follows the same pattern as the proof of match-completeness for the defunctionalized version.

## 5 Matching non-standard regular expressions

Both of the above matchers work on syntactically standard regular expressions, which is used to show termination. We now show that this suffices to define a matcher for non-standard regular expressions, which we represent by a type RegExp:

data RegExp : Set where  $\emptyset$  : RegExp  $\varepsilon$  : RegExp Lit : Char  $\rightarrow$  RegExp  $\_\cdot\_$  : RegExp  $\rightarrow$  RegExp  $\rightarrow$  RegExp  $\_\oplus\_$  : RegExp  $\rightarrow$  RegExp  $\rightarrow$  RegExp  $\_*$  : RegExp  $\rightarrow$  RegExp G : RegExp  $\rightarrow$  RegExp

 $\varnothing$  matches the empty language,  $\varepsilon$  matches the empty string, Lit is character literals, \_\_\_\_ is concatenation, \_ $\oplus$ \_ is alternation, \_\*\_ is Kleene star, and G is for reporting matching strings.

```
\begin{array}{l} \mbox{mutual} \\ \_ \in L\_: \mbox{ List Char} \to \mbox{RegExp} \to \mbox{Set} \\ \_ \in L \varnothing = \bot \\ s \in L \ \epsilon = s \equiv [] \\ s \in L \ (\mbox{Lit c}) = s \equiv c :: [] \end{array}
```

$$\begin{array}{l} s \in L \ (r_1 \oplus r_2) \ = \ (s \in L \ r_1) \uplus (s \in L \ r_2) \\ s \in L \ (r_1 \cdot r_2) \ = \\ \Sigma \ (List \ Char \times List \ Char) \ (\lambda \ \{(p,q) \rightarrow (p \ + \ q \equiv s) \times (p \in L \ r_1) \times (q \in L \ r_2)\}) \\ s \in L \ (r \ *) \ = \ s \in L^x \ r \\ s \in L \ (G \ r) \ = \ s \in L \ r \\ \hline data \ \_ \in L^x \_ : \ List \ Char \rightarrow RegExp \rightarrow Set \ where \\ Ex \ : \ \forall \ \{s \ r\} \rightarrow s \equiv [] \rightarrow s \in L^x \ r \\ Cx \ : \ \forall \ \{s \ s_1 \ s_2 \ r\} \rightarrow s_1 \ + \ s_2 \equiv s \rightarrow s_1 \in L \ r \rightarrow s_2 \in L^x \ r \rightarrow s \in L^x \ r \end{array}$$

The definition of the language of these regular expressions is similar to above, but with an  $\epsilon$  case that requires an empty string, and a base case for  $s \in L^{\times} r$  that allows the empty string; note that G does not change the language.

Next, we define a translation from regexps to syntactically standard regexps. The translation uses a helper function that checks if a regular expression accepts the empty string. Instead of giving this function the type  $\text{RegExp} \rightarrow \text{Bool}$ , we give it a more informative type stating that it decides whether the empty string is in the language of its input:

$$\delta'$$
: (r : RegExp)  $\rightarrow$  ([]  $\in$ L r)  $\uplus$  ( $\neg$  ([]  $\in$ L r))

Using  $\delta',$  we can easily define  $\delta$  :  $\mathsf{RegExp} \to \mathsf{Bool}$  by forgetting the extra information.

The specification for standardization, which we prove below, is that

$$(\forall s) | s \in L(r) \iff [(\delta(r) = true \land s = []) \lor s \in L(\mathsf{standardize}(r))] |$$

That is, s is in the language of r if and only if either and r accepts the empty string and the string is empty, or s is in the language of the standardized version of r.

We standardize as follows:

```
\begin{array}{lll} standardize \,:\, \mathsf{RegExp} \to \mathsf{StdRegExp} \\ standardize \, \varnothing &= \, \varnothing^s \\ standardize \, (\mathsf{Lit} \, x) \,=\, \mathsf{Lit}^s \, x \\ standardize \, (\mathsf{Lit} \, x) \,=\, \mathsf{Lit}^s \, x \\ standardize \, (\mathsf{r}_1 \cdot \mathsf{r}_2) \, \text{with} \, standardize \, \mathsf{r}_1 \,\mid\, standardize \, \mathsf{r}_2 \,\mid\, \delta \, \mathsf{r}_1 \,\mid\, \delta \, \mathsf{r}_2 \\ \ldots \,\mid\, \mathsf{r}_1' \mid\, \mathsf{r}_2' \mid\, \mathsf{false} \mid\, \mathsf{false} \,=\, \mathsf{r}_1' \cdot^s \, \mathsf{r}_2' \\ \ldots \,\mid\, \mathsf{r}_1' \mid\, \mathsf{r}_2' \mid\, \mathsf{false} \mid\, \mathsf{true} \,=\, \mathsf{r}_1' \oplus^s (\mathsf{r}_1' \cdot^s \, \mathsf{r}_2') \\ \ldots \,\mid\, \mathsf{r}_1' \mid\, \mathsf{r}_2' \mid\, \mathsf{true} \mid\, \mathsf{false} \,=\, \mathsf{r}_2' \oplus^s (\mathsf{r}_1' \cdot^s \, \mathsf{r}_2') \\ \ldots \,\mid\, \mathsf{r}_1' \mid\, \mathsf{r}_2' \mid\, \mathsf{true} \mid\, \mathsf{true} \,=\, \mathsf{r}_1' \oplus^s \mathsf{r}_2' \oplus^s (\mathsf{r}_1' \cdot^s \, \mathsf{r}_2') \\ standardize \, (\mathsf{r}_1 \oplus \mathsf{r}_2) \,=\, \mathsf{standardize} \, \mathsf{r}_1 \oplus^s \, \mathsf{standardize} \, \mathsf{r}_2 \end{array}
```

standardize (r \*) = (standardize r)  $^{+s}$ standardize (G r) = standardize r

The empty string language  $\varepsilon$  becomes the empty set  $\emptyset^s$  and Kleene star becomes Kleene plus, because the emptiness checking is performed outside matching the standardized regexp. For the concatenation case, we write  $r_1$ ' and  $r_2$ ' for the standardizations of  $r_1$  and  $r_2$ . Because **standardize** r will not accept the empty string even when r does, it is necessary to check  $r_1$ ' and  $r_2$ ' by themselves in the case where the other one accepts the empty string, because otherwise we would miss strings that rely on one component but not the other being empty.

Our definition of the concatenation case is a bit different than Harper's, where  $\delta$  returns not a boolean, but a regexp  $\emptyset$  (if r does not accept the empty string) or  $\varepsilon$  (if it does), and the clause is as follows:

 $\begin{array}{ll} \mathsf{standardize}\;(\mathsf{r}_1\cdot\mathsf{r}_2)\;=\;((\delta\;\mathsf{r}_1)\cdot^{\mathsf{s}}\;(\mathsf{standardize}\;\mathsf{r}_2))\oplus^{\mathsf{s}}\\ &\quad((\mathsf{standardize}\;\mathsf{r}_1)\cdot^{\mathsf{s}}\;(\delta\;\mathsf{r}_2))\oplus^{\mathsf{s}}\\ &\quad((\mathsf{standardize}\;\mathsf{r}_1)\cdot^{\mathsf{s}}\;(\mathsf{standardize}\;\mathsf{r}_2))\end{array}$ 

This definition is equivalent to above, using the equivalences that for any r,  $\varnothing \cdot^{s} r = \varnothing = r \cdot^{s} \varnothing$  and  $\epsilon \cdot^{s} r = r = r \cdot^{s} \epsilon$ . For example, when  $\delta r_{1}$  is true, Harper's translation gives a  $\epsilon \cdot^{s}$  (standardize  $r_{2}$ ) summand, which is standard but *not* syntactically standard, but we can simplify it to standardize  $r_{2}$ . When  $\delta r_{1}$  is false, Harper's translation gives an  $\varnothing \cdot^{s}$  (standardize  $r_{2}$ ) summand, which drops out.

This definition of standardization satisfies the above correctness theorem; we have proved

$$\begin{split} \in &\mathsf{L}\text{-soundness} : (\mathsf{s} : \mathsf{List} \mathsf{Char}) \, (\mathsf{r} : \mathsf{RegExp}) \\ & \to ((\delta \ \mathsf{r} \equiv \mathsf{true}) \times (\mathsf{s} \equiv [])) \uplus (\mathsf{s} \in \mathsf{L}^{\mathsf{s}} \ (\mathsf{standardize} \ \mathsf{r})) \\ & \to \mathsf{s} \in \mathsf{L} \ \mathsf{r} \\ \in &\mathsf{L}\text{-completeness} : (\mathsf{s} : \mathsf{List} \mathsf{Char}) \ (\mathsf{r} : \mathsf{RegExp}) \\ & \to \mathsf{s} \in \mathsf{L} \ \mathsf{r} \\ & \to ((\delta \ \mathsf{r} \equiv \mathsf{true}) \times (\mathsf{s} \equiv [])) \uplus (\mathsf{s} \in \mathsf{L}^{\mathsf{s}} \ (\mathsf{standardize} \ \mathsf{r})) \end{split}$$

Now that we have a verified standardize function, we can define \_accepts\_ as follows, where accepts<sup>s</sup>-intrinsic can be either of the above matchers:

accepts-intrinsic : (r : RegExp)  $\rightarrow$  (s : List Char)  $\rightarrow$  Maybe (s  $\in$  L r) accepts-intrinsic r s with  $\delta$ ' r accepts-intrinsic r [] | inj<sub>1</sub> x = just x accepts-intrinsic r s | \_ = map ( $\in$ L-soundness s r  $\circ$  inj<sub>2</sub>) (accepts<sup>s</sup>-intrinsic (standardize r) s) If r accepts the empty string, we return true if xs is empty or the standardization of r accepts xs. If r does not accept the empty string, then we only have the latter option. In that case, we call accepts<sup>s</sup>-intrinsic to get an optional derivation of the type  $s \in L^s$  (standardize r) and use that on  $\in$ L-soundness to get an optional derivation of the type  $s \in L^r$ .

As usual, we have proved completeness extrinsically:

```
\begin{array}{l} \mbox{correct-completeness} \ : \ (r \ : \ \mbox{RegExp}) \ (s \ : \ \mbox{List Char}) \\ \rightarrow s \in L \ r \\ \rightarrow \ \mbox{isJust} \ (r \ \mbox{accepts} \ s) \end{array}
```

Finally, we have proved decidability of matching as a corollary of soundness and completeness:

decidability : (r : RegExp) (s : List Char)  $\rightarrow$  (s  $\in$ L r)  $\uplus$  ( $\neg$  (s  $\in$ L r))

#### 5.1 Capturing groups

The "capturing group" constructor G is intended to allow the user to specify parts of a regular expression whose matching strings should be extracted and reported. For example, if our regular expression checks if a string is a valid e-mail address, we might use this to parse the username and domains. If we have a regular expression alphanumeric : RegExp that accepts a single alphanumeric character (this can be defined in the above language as a big  $\oplus$  of character literals), we can define a (naïve) regular expression for e-mail addresses such as

```
e-mail : RegExp
e-mail = G (alphanumeric *) · Lit '0' · G (alphanumeric *) · Lit '. ' · G (alphanumeric *)
```

Now, if we match the string "jdoe@wesleyan.edu" with e-mail, we should extract and report "jdoe", "wesleyan" and "edu", because each of those substrings matched a sub-regexp that was marked with G.

This extraction can be computed from the derivation of  $s \in L r$ , which provides a parse tree that says which substring of xs is matched by which part of r. Thus, to report the groups, we do an in-order traversal of the regexp and derivation tree, and collect the strings matching a capturing group to a list. The function to do this can be defined as follows:

```
extract : {r : RegExp} \rightarrow {xs : List Char} \rightarrow xs \inL r \rightarrow List (List Char) extract {\emptyset} ()
```

```
\begin{array}{l} \text{extract } \{\epsilon\} \ \text{refl} \ = \ [] \\ \text{extract } \{\text{Lit } x\} \ \text{refl} \ = \ [] \\ \text{extract } \{\text{Lit } x\} \ \text{refl} \ = \ [] \\ \text{extract } \{r_1 \cdot r_2\} \ ((\text{as, bs}), \text{eq, a, b}) \ = \ \text{extract } \{r_1\} \ \{\text{as}\} \ a \ + \ \text{extract } \{r_2\} \ \{\text{bs}\} \ b \\ \text{extract } \{r_1 \oplus r_2\} \ (\text{inj}_1 \ x) \ = \ \text{extract } \{r_1\} \ x \\ \text{extract } \{r_1 \oplus r_2\} \ (\text{inj}_2 \ y) \ = \ \text{extract } \{r_2\} \ y \\ \text{extract } \{r^*\} \ (\text{Ex refl}) \ = \ [] \\ \text{extract } \{r^*\} \ (\text{Ex refl}) \ = \ [] \\ \text{extract } \{r^*\} \ (\text{Cx } \{s\} \ \{s_1\} \ \{s_2\} \ x \ x_1 \ \text{inL}) \ = \ \text{extract } \{r\} \ x_1 \ + \ \text{extract } \{r^*\} \ \text{inL} \\ \text{extract } \{G\ r\} \ \{x_S\} \ \text{inL} \ = \ xs \ :: \ \text{extract } \{r\} \ \text{inL} \end{array}
```

Base cases  $\emptyset$ ,  $\varepsilon$ , Lit will return an empty list because if they are captured by a group, the substring is already added to the list in the previous calls to extract. In concatenation, we make two recursive calls and append the results because  $r_1$  and  $r_2$  match different substrings and they may have different capturing groups inside them. In alternation, the entire string matches either  $r_1$  or  $r_2$ , so we make one recursive call to the one it matches. The Kleene star case follows the same principles.

Combining this with our intrinsic matcher, we can define an overall function

```
groups : (r : RegExp) (s : List Char) \rightarrow Maybe (List (List Char))
groups r s = map extract (accepts-intrinsic r s)
```

# 6 Conclusion

We have studied three variations on Harper's algorithm for regular expression matching, which were inspired by programming and verifying this algorithm using dependent types: defunctionalizing the matcher allows Agda to see termination without an explicit metric, and provides an alternative to Yi's first-order matcher based on state machines; intrinsically verifying soundness allows extracting matching strings; and a syntactic definition of standard regular expressions simplifies the staging of the development. We believe that these variations provide a nice illustration of the benefits of thinking in a dependently typed language, and that they have some pedagogical value for teaching this material in courses on dependently typed programming—or, by porting the observations back to simply-typed languages, on introductory programming.

# References

- Danielsson, Nils Anders. (2010). Total parser combinators. Pages 285– 296 of: Proceedings of the 15th acm sigplan international conference on functional programming. ICFP '10. New York, NY, USA: ACM.
- Firsov, Denis, & Uustalu, Tarmo. (2013). Certified parsing of regular languages. Pages 98–113 of: Proceedings of the third international conference on certified programs and proofs - volume 8307. New York, NY, USA: Springer-Verlag New York, Inc.
- Harper, Robert. (1999). Functional pearl. proof-directed debugging corrigendum. Journal of functional programming, **9**(4), 463–469.
- Norell, Ulf. (2007). Towards a practical programming language based on dependent type theory. Ph.D. thesis, Chalmers University of Technology.
- Reynolds, John C. (1972). Definitional interpreters for higher-order programming languages. Pages 717–740 of: Proceedings of the acm annual conference - volume 2. ACM '72. ACM.
- Ridge, Tom. (2011). Simple, functional, sound and complete parsing for all context-free grammars. Pages 103–118 of: Proceedings of the first international conference on certified programs and proofs. CPP'11. Berlin, Heidelberg: Springer-Verlag.
- Yi, Kwangkeun. (2006). Educational pearl: "proof-directed debugging" revisited for a first-order version. Journal of functional programming, 16(6), 663–670.