

# Verified parsers using the refinement calculus and algebraic effects

Tim Baanen and Wouter Swierstra

Vrije Universiteit Amsterdam, Utrecht University  
{t.baanen@vu.nl,w.s.swierstra@uu.nl}

## 1 Introduction

Parsing is hard, despite a great deal of existing tooling and libraries. There are various ways to write a parser in functional languages—famously using parser combinators [others?; Hut92; SD96; Wad85].

How do we ensure that these parser combinators are correct? There are several different papers that attempt to answer this question [Dan10; Fir16]. In this paper, we take a different approach, drawing inspiration from recent work on algebraic effects [BP15; WSH14; McB15]. In particular, we demonstrate how an algebraic treatment of general recursion lets us separate the (partial) correctness of the combinators from their termination in a clean fashion. Most existing proofs require combinators to show that the string being parsed decreases, conflating termination and correctness. We aim to split apart various aspects of parsers and their correctness into a library of individually useful components. With this library, we can easily introduce concepts only when they are needed, and do this without having to rework the previous definitions.

Finally, this paper serves an extended example of the recent work on using predicate transformers to reason about (combinations of) algebraic effects [SB19].

To these ends, this paper makes the following novel contributions:

- The non-recursive fragment of regular expressions can be correctly parsed using non-determinism (Section 3); by combining non-determinism with general recursion (Section 4), support for the Kleene star can be added without compromising our previous definitions (Section 5). Although the resulting parser is not guaranteed to terminate, it is *refined* by another implementation using Brzozowski derivatives that does terminate (Section 6).
- Next, we show how this approach may be extended to handle context-free languages. To do so, we show how to write parsers using algebraic effects (Section 7), and map grammars to parsers (Section 8). Once again, we can cleanly separate the proofs of partial correctness (Section 9) and termination (Section 10).

We write our programs and proofs in the dependently typed language Agda [Nor07], ensuring our work is formally correct.

## 2 Recap: algebraic effects and predicate transformers

Algebraic effects were introduced to allow for incorporating side effects in functional languages. For example, the effect *ENondet* allows for nondeterministic programs:

```

record Effect : Set where
  constructor eff
  field
    C : Set
    R : C → Set
data CNondet : Set where
  Fail : CNondet
  Choice : CNondet
  RNondet : CNondet → Set
  RNondet Fail = ⊥
  RNondet Choice = Bool
  ENondet = eff CNondet RNondet

```

We represent effectful programs using the *Free* datatype.

```

data Free (e : Effect) (a : Set) : Set where
  Pure : a → Free e a
  Step : (c : C e) → (R e c → Free e a) → Free e a

```

This gives a monad, with the bind operator defined as follows:

```

_>>_ : Free e a → (a → Free e b) → Free e b
Pure x >> f = f x
Step c k >> f = Step c (λ x → k x >> f)

```

The easiest way to use effects is with smart constructors:

```

fail : Free ENondet a
fail = Step Fail λ ()
choice : Free ENondet a → Free ENondet a → Free ENondet a
choice S1 S2 = Step Choice λ b → if b then S1 else S2

```

To give specifications of programs that incorporate effects, we can use predicate transformers.

```

wp : {C : Set} {R : C → Set} → ((c : C) → (R c → Set) → Set) →
  {a : Set} → Free (eff C R) a → (a → Set) → Set
wp alg (Pure x) P = P x
wp alg (Step c k) P = alg c λ x → wp alg (k x) P

```

Interestingly, these predicate transformers are exactly the catamorphisms from *Free* to *Set*.

$$\begin{aligned}
ptAll &: (c : CNondet) \rightarrow (RNondet\ c \rightarrow Set) \rightarrow Set \\
ptAll\ Fail\ P &= \top \\
ptAll\ Choice\ P &= P\ True \wedge P\ False
\end{aligned}$$

$$\begin{aligned}
wpNondetAll &: Free\ ENondet\ a \rightarrow (a \rightarrow Set) \rightarrow Set \\
wpNondetAll\ S\ P &= wp\ ptAll\ S\ P
\end{aligned}$$

We use pre- and postconditions to give a specification for a program. If the precondition holds on the input, the program needs to ensure the postcondition holds on the output.

```

module Spec where
  record Spec (a : Set) : Set where
    constructor [_, _]
    field
      pre : Set
      post : a → Set
  wpSpec : Spec a → (a → Set) → Set
  wpSpec [pre, post] P = pre ∧ (∀ o → post o → P o)

```

The refinement relation expresses when one program is “better” than another. We need to take into account the semantics we want to impose on the program, so we define it in terms of the predicate transformer associated with the program.

$$\begin{aligned}
\_ \sqsubseteq \_ &: (pt_1\ pt_2 : (a \rightarrow Set) \rightarrow Set) \rightarrow Set \\
pt_1 \sqsubseteq pt_2 &= \forall P \rightarrow pt_1\ P \rightarrow pt_2\ P
\end{aligned}$$

### 3 Almost parsing regular languages

To see how we can use the *Free* monad for writing and verifying a parser, and more specifically how we use the *ENondet* effect for writing and the *wpNondetAll* semantics for verifying a parser, we will look at parsing a given regular language. Our approach is first to define the specification of a parser, then inspect this specification to write the first implementation and prove (partial) correctness of this implementation. We will later improve this implementation by refining it.

**Definition 1** ([AU77]) *The class of regular languages is the smallest class such that:*

- the empty language is regular,
- the language containing only the empty string is regular,
- for each character  $\mathfrak{x}$ , the language containing only the string “ $\mathfrak{x}$ ” is regular,
- the union and concatenation of regular languages are regular, and
- the repetition of a regular language is regular.

A regular language can be defined using a regular expression, which we will represent as an element of the *Regex* datatype. An element of this type represents the syntax of a regular language, and we will generally identify a regular expression with the language it denotes.

```
data Regex : Set where
  Empty      : Regex
  Epsilon    : Regex
  Singleton  : Char → Regex
  _|_         : Regex → Regex → Regex
  _ · _       : Regex → Regex → Regex
  _ ★        : Regex → Regex
```

Here, *Empty* is an expression for empty language (which matches no strings at all), while *Epsilon* is an expression for the language of the empty string (which matches exactly one string: "").

What should a parser for regular languages output? If we only want to know whether a string matches a regular expression, we can return a *Bool*. If we want to know more, we could annotate the regular expression with capture groups, and say that the output of the parser maps each capture group to the substring that the capture group matches. We can also return a full parse tree, mirroring the structure of the expression.

```
Tree : Regex → Set
Tree Empty      = ⊥
Tree Epsilon    = ⊤
Tree (Singleton _) = Char
Tree (l | r)     = Either (Tree l) (Tree r)
Tree (l · r)     = Pair (Tree l) (Tree r)
Tree (r ★)       = List (Tree r)
```

In Agda, we can represent the semantics of the *Regex* type by giving a relation between a *Regex* and a *String* on the one hand (the input of the matcher), and a parse tree on the other hand (the output of the parser). Note that the *Tree* type itself is not sufficient to represent the semantics, since it does not say which strings result in any given parse tree. If the *Regex* and *String* do not match, there should be no output, otherwise the output consists of all relevant parse trees. We give the relation using the following inductive definition:

```
data Match : (r : Regex) → String → Tree r → Set where
  Epsilon      : Match Epsilon Nil tt
  Singleton    : Match (Singleton x) (x :: Nil) x
  OrLeft       : Match l xs x → Match (l | r) xs (Inl x)
  OrRight      : Match r xs x → Match (l | r) xs (Inr x)
  Concat       : Match l ys y → Match r zs z →
                  Match (l · r) (ys ++ zs) (y , z)
  StarNil      : Match (r ★) Nil Nil
  StarConcat   : Match (r · (r ★)) xs (y , ys) → Match (r ★) xs (y :: ys)
```

Note that there is no constructor for *Match Empty xs ms* for any *xs* or *ms*, which we interpret as that there is no way to match the *Empty* language with a string *xs*. Similarly, the only constructor for *Match Epsilon xs ms* is where *xs* is the empty string *Nil*.

Since the definition of *Match* allows for multiple ways that a given *Regex* and *String* may match, such as in the trivial case where the *Regex* is of the form  $r \mid r$ , and it also has cases where there is no way to match a *Regex* and a *String*, such as where the *Regex* is *Empty*, we can immediately predict some parts of the implementation of the *match* function. Whenever we encounter an expression of the form  $l \mid r$ , we make a nondeterministic *Choice* between either *l* or *r*. Similarly, whenever we encounter the *Empty* expression, we immediately *fail*. In the previous analysis steps, we have already assumed that we implement the parser by structural recursion on the *Regex*, so let us consider other cases.

The implementation for concatenation is not as immediately obvious. One way that we can deal with it is to not only return a *Tree* from the parser. Instead, the parser also returns the unmatched portion of the string, and when we have to match a regular expression of the form  $l \cdot r$  with a string *xs*, we match *l* with *xs* giving a left over string *ys*, then match *r* with *ys*. We can also do without changing the return values of the parser, by nondeterministically splitting the string *xs* into  $ys \mathrel{+} zs$ . That is what we do in a helper function *allSplits*, which nondeterministically chooses such *ys* and *zs* and returns them as a pair.

```
allSplits : (xs : List a) → Free ENondet (List a × List a)
allSplits Nil = Pure (Nil , Nil)
allSplits (x :: xs) = choice
  (Pure (Nil , (x :: xs)))
  (allSplits xs ≫= λ { (ys , zs) → Pure ((x :: ys) , zs) })
```

Armed with this helper function, we can write the first part of a nondeterministic regular expression matcher, that does a case distinction on the expression and then checks that the string has the correct format.

```
match : (r : Regex) (xs : String) → Free ENondet (Tree r)
match Empty xs = fail
match Epsilon Nil = Pure tt
match Epsilon ( _ :: _ ) = fail
match (Singleton c) Nil = fail
match (Singleton c) (x :: Nil) with c  $\stackrel{?}{=}$  x
match (Singleton c) (.c :: Nil) | yes refl = Pure c
match (Singleton c) (x :: Nil) | no ¬p = fail
match (Singleton c) ( _ :: _ :: _ ) = fail
match (l · r) xs = do
  (ys , zs) ← allSplits xs
  y ← match l ys
  z ← match r zs
  Pure (y , z)
match (l | r) xs = choice (Inl ⟨$⟩ match l xs) (Inr ⟨$⟩ match r xs)
```

Unfortunately, we get stuck in the case of  $\_ \star$ . We could do a similar construction to  $l \cdot r$ , where we split the string into two parts and match the first part with  $r$  and the second part with  $r \star$ , but this definition will be rejected by Agda, since it does not terminate. Since there is no easy way to handle this case for now, we just *fail* when we encounter a regex  $r \star$ .

$$\text{match } (r \star) \text{ } xs = \text{fail}$$

Still, we can prove that this matcher works, as long as the regular expression does not contain  $\_ \star$ . In other words, we can prove that the *match* function satisfies the postcondition given by the type *Match*, as long as the precondition *hasNo\** holds:

$$\begin{aligned} \text{hasNo*} &: \text{Regex} \rightarrow \text{Set} \\ \text{hasNo* Empty} &= \top \\ \text{hasNo* Epsilon} &= \top \\ \text{hasNo* (Singleton } x) &= \top \\ \text{hasNo* } (l \cdot r) &= \text{hasNo* } l \wedge \text{hasNo* } r \\ \text{hasNo* } (l \mid r) &= \text{hasNo* } l \wedge \text{hasNo* } r \\ \text{hasNo* } (r \star) &= \perp \\ \text{pre} &: (r : \text{Regex}) (xs : \text{String}) \rightarrow \text{Set} \\ \text{pre } r \text{ } xs &= \text{hasNo* } r \\ \text{post} &: (r : \text{Regex}) (xs : \text{String}) \rightarrow \text{Tree } r \rightarrow \text{Set} \\ \text{post} &= \text{Match} \end{aligned}$$

In order to state that *match* works correctly, we need to determine its semantics: is the nondeterminism angelic or demonic? Since the use of nondeterminism in *match* is to find all correct matches, we want that all values potentially returned are correct, as specified by the *ptAll* semantics used in *wpNondetAll*.

If we now try to give a correctness proof with respect to this pre- and postcondition, we run into an issue in cases where the definition makes use of the  $\_ \gg \_$  operator. The *wp*-based semantics completely unfolds the left hand side, before it can talk about the right hand side. Whenever our matcher makes use of structural recursion on the left hand side of a  $\_ \gg \_$  (more specifically, in the definition of *allSplits* and in the cases of  $l \cdot r$  and  $l \mid r$ ), we cannot make progress in our proof without reducing this left hand side to a recursion-less expression. We need a lemma relating the semantics of program composition to the semantics of individual programs, which is also known as the law of consequence for traditional predicate transformer semantics.[cite?](#)

$$\begin{aligned} \text{consequence} &: \forall pt (mx : \text{Free es } a) (f : a \rightarrow \text{Free es } b) \rightarrow \\ &\quad \text{wp } pt \text{ } mx (\lambda x \rightarrow \text{wp } pt (f x) P) == \text{wp } pt (mx \gg f) P \\ \text{consequence } pt (Pure x) f &= \text{refl} \\ \text{consequence } pt (Step c k) f &= \text{cong } (pt c) \\ &\quad (\text{extensionality } \lambda x \rightarrow \text{consequence } pt (k x) f) \\ \text{wpToBind} &: (mx : \text{Free es } a) (f : a \rightarrow \text{Free es } b) \rightarrow \end{aligned}$$

$$\begin{aligned}
& wp \text{ pt } mx (\lambda x \rightarrow wp \text{ pt } (f \ x) \ P) \rightarrow wp \text{ pt } (mx \gg= f) \ P \\
& wpToBind \ mx \ f \ H = \text{subst id } (\text{consequence pt } mx \ f) \ H \\
& wpFromBind : (mx : Free \ es \ a) (f : a \rightarrow Free \ es \ b) \rightarrow \\
& \quad wp \text{ pt } (mx \gg= f) \ P \rightarrow wp \text{ pt } mx (\lambda x \rightarrow wp \text{ pt } (f \ x) \ P) \\
& wpFromBind \ mx \ f \ H = \text{subst id } (\text{sym } (\text{consequence pt } mx \ f)) \ H
\end{aligned}$$

The correctness proof for *match* closely matches the structure of *match* (and by extension *allSplits*). It uses the same recursion on *Regex* as in the definition of *match*. Since we make use of *allSplits* in the definition, we first give its correctness proof.

$$\begin{aligned}
& allSplitsPost : String \rightarrow String \times String \rightarrow Set \\
& allSplitsPost \ xs \ (ys, zs) = xs == ys \# zs \\
& allSplitsSound : \forall \ xs \rightarrow \\
& \quad wpSpec [\top, allSplitsPost \ xs] \sqsubseteq wpNondetAll (allSplits \ xs) \\
& allSplitsSound \ Nil \quad P \ (preH, postH) = postH \_ refl \\
& allSplitsSound \ (x :: xs) \ P \ (preH, postH) = postH \_ refl, \\
& \quad wpToBind (allSplits \ xs) \_ (allSplitsSound \ xs \_ (tt, \\
& \quad \lambda \_ \ H \rightarrow postH \_ (cong \ (x :: \_) \ H)))
\end{aligned}$$

Then, using *wpToBind*, we incorporate this correctness proof in the correctness proof of *match*. Apart from having to introduce *wpToBind*, the proof essentially follows automatically from the definitions.

$$\begin{aligned}
& matchSound : \forall \ r \ xs \rightarrow \\
& \quad wpSpec [\text{pre } r \ xs, \text{post } r \ xs] \sqsubseteq wpNondetAll (match \ r \ xs) \\
& matchSound \ Empty \ xs \quad P \ (preH, postH) = tt \\
& matchSound \ Epsilon \ Nil \quad P \ (preH, postH) = postH \_ Epsilon \\
& matchSound \ Epsilon \ (x :: xs) \ P \ (preH, postH) = tt \\
& matchSound \ (Singleton \ x) \ Nil \ P \ (preH, postH) = tt \\
& matchSound \ (Singleton \ x) \ (c :: Nil) \ P \ (preH, postH) \textbf{ with } x \stackrel{?}{=} c \\
& \dots \mid \text{yes refl} = postH \_ Singleton \\
& \dots \mid \text{no } \neg p = tt \\
& matchSound \ (Singleton \ x) \ (\_ :: \_ :: \_) \ P \ (preH, postH) = tt \\
& matchSound \ (l \cdot r) \ xs \ P \ ((preL, preR), postH) = \\
& \quad wpToBind (allSplits \ xs) \_ (allSplitsSound \ xs \_ (tt, \\
& \quad \lambda \{ (ys, zs) \} \text{ splitH} \rightarrow wpToBind (match \ l \ ys) \_ (matchSound \ l \ ys \_ (preL, \\
& \quad \lambda \ y \ lH \rightarrow wpToBind (match \ r \ zs) \_ ((matchSound \ r \ zs \_ (preR, \\
& \quad \lambda \ z \ rH \rightarrow postH \ (y, z) \ (\text{subst } (\lambda \ xs \rightarrow Match \_ \ xs \_) \ (\text{sym splitH}) \\
& \quad (\text{Concat } lH \ rH)))))) \}))) \\
& matchSound \ (l \mid r) \ xs \ P \ ((preL, preR), postH) = \\
& \quad wpToBind (match \ l \ xs) \_ (matchSound \ l \ xs \_ (preL, \\
& \quad \lambda \_ \ lH \rightarrow postH \_ (OrLeft \ lH))), \\
& \quad wpToBind (match \ r \ xs) \_ (matchSound \ r \ xs \_ (preR, \\
& \quad \lambda \_ \ rH \rightarrow postH \_ (OrRight \ rH))) \\
& matchSound \ (r \star) \ xs \ P \ ((), postH)
\end{aligned}$$

## 4 Combining nondeterminism and general recursion

The matcher we have defined in the previous section is unfinished, since it is not able to handle regular expressions that incorporate the Kleene star. The fundamental issue is that the Kleene star allows for arbitrarily many distinct matchings in certain cases. For example, matching *Epsilon*  $\star$  with the empty string "" will allow for repeating the *Epsilon* arbitrarily often, since *Epsilon*  $\cdot$  (*Epsilon*  $\star$ ) is equivalent to both *Epsilon* and *Epsilon*  $\star$ . Thus, we cannot implement *match* on the  $\_ \star$  operator by helping Agda's termination checker.

What we will do instead is to deal with the recursion as an effect. A recursively defined (dependent) function of type  $(i : I) \rightarrow O\ i$  can instead be given as an element of the type  $(i : I) \rightarrow \text{Free } (ERec\ I\ O)\ (O\ i)$ , where *ERec* *I* *O* is the effect of *general recursion* [McB15]:

$$\begin{aligned} ERec &: (I : \text{Set}) (O : I \rightarrow \text{Set}) \rightarrow \text{Effect} \\ ERec\ I\ O &= \text{eff}\ I\ O \end{aligned}$$

Defining *match* with the *ERec* effect is not sufficient to implement it fully either, since replacing the effect *ENondet* with *ERec* does not allow for non-determinism anymore, so while the Kleene star might work, the other parts of *match* do not work anymore. We need a way to combine effects.

We can combine two effects in a straightforward way: given *eff* *C*<sub>1</sub> *R*<sub>1</sub> and *eff* *C*<sub>2</sub> *R*<sub>2</sub>, we can define a new effect by taking the disjoint union of the commands and responses, resulting in *eff* (*Either* *C*<sub>1</sub> *C*<sub>2</sub>) [*R*<sub>1</sub> , *R*<sub>2</sub> ], where [*R*<sub>1</sub> , *R*<sub>2</sub> ] is the unique map given by applying *R*<sub>1</sub> to values in *C*<sub>1</sub> and *R*<sub>2</sub> to *C*<sub>2</sub> [WSH14]. If we want to support more effects, we can repeat this process of disjoint unions, but this quickly becomes somewhat cumbersome. For example, the disjoint union construction is associative semantically, but not syntactically. If two programs have the same set of effects that is associated differently, we cannot directly compose them.

Instead of building a new effect type, we modify the *Free* monad to take a list of effects instead of a single effect. The *Pure* constructor remains as it is, while the *Step* constructor takes an index into the list of effects and the command and continuation for the effect with this index.

$$\begin{aligned} \text{data } Free\ (es : List\ Effect)\ (a : Set) &: Set \text{ where} \\ \text{Pure} &: a \rightarrow Free\ es\ a \\ \text{Step} &: (i : e \in es)\ (c : C\ e)\ (k : R\ e\ c \rightarrow Free\ es\ a) \rightarrow Free\ es\ a \end{aligned}$$

By using a list of effects instead of allowing arbitrary disjoint unions, we have effectively chosen that the disjoint unions canonically associate to the right. Since the disjoint union is also commutative, it would be cleaner to have the collection of effects be unordered as well. Unfortunately, Agda does not provide a multiset type that is easy to work with.

We choose to use the same names and almost the same syntax for this new definition of *Free*, since all definitions that use the old version can be ported



over with almost no change. Thus, we will not repeat definitions such as  $\_\gg\_\$  and *consequence* for the new *Free* type.

Most of this bookkeeping can be inferred by Agda's typeclass inference, so we make the indices instance arguments, indicated by the double curly braces  $\{\!\!\{\}$  surrounding the arguments.

```
fail : {iND : ENondet ∈ es} → Free es a
fail {iND} = Step iND Fail λ ()
choice : {iND : ENondet ∈ es} → Free es a → Free es a → Free es a
choice {iND} S1 S2 = Step iND Choice λ b → if b then S1 else S2
call : {iRec : ERec I O ∈ es} → (i : I) → Free es (O i)
call {iRec} i = Step iRec i Pure
```

For convenience of notation, we introduce the  $\_ \overset{es}{\dashv} \_$  notation for general recursion, i.e. Kleisli arrows into *Free* (*ERec*  $\_ \_ :: es$ ).

```
_ \overset{es}{\dashv} _ : (C : Set) (es : List Effect) (R : C → Set) → Set
C \overset{es}{\dashv} R = (c : C) → Free (eff C R :: es) (R c)
```

With the syntax for combinations of effects defined, let us turn to semantics. Since the weakest precondition predicate transformer for a single effect is given as a fold over the effect's predicate transformer, the semantics for a combination of effects can be given as a fold over a (dependent) list of predicate transformers.

```
record PT (e : Effect) : Set where
  constructor mkPT
  field
    pt : (c : C e) → (R e c → Set) → Set
    mono : ∀ c P P' → P ⊆ P' → pt c P → pt c P'
data PTs : List Effect → Set where
  Nil : PTs Nil
  _ :: _ : ∀ {e es} → PT e → PTs es → PTs (e :: es)
```

The record type *PT* not only contains a predicate transformer *pt*, but also a proof that *pt* is monotone in its predicate. The requirement of monotonicity is needed to prove some lemmas later on **which exactly?**, and makes intuitive sense: if the precondition holds for a certain postcondition, a weaker postcondition should also have its precondition hold.

Given a such a list of predicate transformers, defining the semantics of an effectful program is a straightforward generalization of *wp*. The *Pure* case is identical, and in the *Step* case we find the predicate transformer at the corresponding index to the effect index  $i : e \in es$  using the *lookupPT* helper function.

```
lookupPT : (pts : PTs es) (i : eff C R ∈ es) →
  (c : C) → (R c → Set) → Set
```

$$\begin{aligned} \text{lookupPT } (pt :: pts) \in \text{Head} &= PT.pt \ pt \\ \text{lookupPT } (pt :: pts) \in \text{Tail } i &= \text{lookupPT } pts \ i \end{aligned}$$

This results in the following definition of  $wp$  for combinations of effects.

$$\begin{aligned} wp &: (pts : PTs \ es) \rightarrow Free \ es \ a \rightarrow (a \rightarrow Set) \rightarrow Set \\ wp \ pts \ (Pure \ x) \ P &= P \ x \\ wp \ pts \ (Step \ i \ c \ k) \ P &= \text{lookupPT } pts \ i \ c \ \lambda x \rightarrow wp \ pts \ (k \ x) \ P \end{aligned}$$

The effects we are planning to use for *match* are a combination of nondeterminism and general recursion. We re-use the *ptAll* semantics of nondeterminism, packaging them in a *PT* record. However, it is not as easy to give a predicate transformer for general recursion, since the intended semantics of a recursive call depend on the function that is being called, i.e. the function that is being defined.

However, if we have a specification of a function of type  $(i : I) \rightarrow O \ i$ , for example in terms of a relation of type  $(i : I) \rightarrow O \ i \rightarrow Set$ , we are able to define a predicate transformer:

$$\begin{aligned} ptRec &: ((i : I) \rightarrow O \ i \rightarrow Set) \rightarrow PT \ (ERec \ I \ O) \\ PT.pt \ (ptRec \ R) \ i \ P &= \forall o \rightarrow R \ i \ o \rightarrow P \ o \\ PT.mono \ (ptRec \ R) \ c \ P \ P' \ imp \ asm \ o \ h &= imp \ _ \ (asm \ _ \ h) \end{aligned}$$

In the case of verifying the *match* function, the *Match* relation will play the role of  $R$ . If we use  $ptRec \ R$  as a predicate transformer to check that a recursive function satisfies the relation  $R$ , then we are proving *partial correctness*, since we assume each recursive call terminates according to the relation  $R$ .

## 5 Recursively parsing every regular expression

To deal with the Kleene star, we rewrite *match* as a generally recursive function using a combination of effects. Since *match* makes use of *allSplits*, we also rewrite that function to use a combination of effects. The types become:

$$\begin{aligned} allSplits &: \{iND : ENondet \in es\} \rightarrow List \ a \rightarrow Free \ es \ (List \ a \times List \ a) \\ match &: \{iND : ENondet \in es\} \rightarrow Regex \times String \overset{es}{\rightsquigarrow} Tree \circ Pair.fst \end{aligned}$$

Since the index argument to the smart constructor is inferred by Agda, the only change in the definition of *match* and *allSplits* will be that *match* now implements the Kleene star:

$$\begin{aligned} match \ ((r \star), Nil) &= Pure \ Nil \\ match \ ((r \star), xs@(\_ :: \_)) &= \mathbf{do} \\ &\quad (y, ys) \leftarrow call \ ((r \cdot (r \star)), xs) \\ &\quad Pure \ (y :: ys) \end{aligned}$$

The effects we need to use for running *match* are a combination of nondeterminism and general recursion. As discussed, we first need to give the specification for *match* before we can verify a program that performs a recursive *call* to *match*.

$$\begin{aligned} \text{matchSpec} &: (r, xs : \text{Pair Regex String}) \rightarrow \text{Tree} (\text{Pair.fst } r, xs) \rightarrow \text{Set} \\ \text{matchSpec } (r, xs) \text{ ms} &= \text{Match } r \text{ xs ms} \\ \text{wpMatch} &: \text{Free} (\text{ERec} (\text{Pair Regex String}) (\text{Tree} \circ \text{Pair.fst}) :: \text{ENondet} :: \text{Nil}) a \rightarrow \\ &\quad (a \rightarrow \text{Set}) \rightarrow \text{Set} \\ \text{wpMatch} &= \text{wp} (\text{ptRec matchSpec} :: \text{ptAll} :: \text{Nil}) \end{aligned}$$

We can reuse exactly the same proof to show *allSplits* is correct, since we use the same semantics for the effects in *allSplits*. Similarly, the correctness proof of *match* will be the same on all cases except the Kleene star. Now we are able to prove correctness of *match* on a Kleene star.

$$\begin{aligned} \text{matchSound } ((r \star), \text{Nil}) \quad P (\text{preH}, \text{postH}) &= \\ \text{postH} - \text{StarNil} & \\ \text{matchSound } ((r \star), (x :: xs)) P (\text{preH}, \text{postH}) \circ H &= \\ \text{postH} - (\text{StarConcat } H) & \end{aligned}$$

At this point, we have defined a parser for regular languages and formally proved that its output is always correct. However, *match* does not necessarily terminate: if *r* is a regular expression that accepts the empty string, then calling *match* on *r*  $\star$  and a string *xs* results in the first nondeterministic alternative being an infinitely deep recursion.

The next step is then to write a parser that always terminates and show that *match* is refined by it. Our approach is to do recursion on the input string instead of on the regular expression.

## 6 Termination, using derivatives

Since recursion on the structure of a regular expression does not guarantee termination of the parser, we can instead perform recursion on the string to be parsed. To do this, we make use of an operation on languages called the Brzozowski derivative.

**Definition 2 ([Brz64])** *The Brzozowski derivative of a formal language  $L$  with respect to a character  $x$  consists of all strings  $xs$  such that  $x :: xs \in L$ .*

Importantly, if  $L$  is regular, so are all its derivatives. Thus, let  $r$  be a regular expression, and  $d \ r / d \ x$  an expression for the derivative with respect to  $x$ , then  $r$  matches a string  $x :: xs$  if and only if  $d \ r / d \ x$  matches  $xs$ . This suggests the following implementation of matching an expression  $r$  with a string  $xs$ : if  $xs$  is empty, check whether  $r$  matches the empty string; otherwise let  $x$  be the head of the string and  $xs'$  the tail and go in recursion on matching  $d \ r / d \ x$  with  $xs'$ .

The first step in implementing a parser using the Brzozowski derivative is to compute the derivative for a given regular expression. Following Brzozowski

[Brz64], we use a helper function  $\varepsilon?$  that decides whether an expression matches the empty string.

$$\varepsilon? : (r : \text{Regex}) \rightarrow \text{Dec} (\sum (Tree\ r) (Match\ r\ Nil))$$

The definitions of  $\varepsilon?$  is given by structural recursion on the regular expression, just as the derivative operator is:

$$\begin{aligned} d\_ / d\_ : \text{Regex} &\rightarrow \text{Char} \rightarrow \text{Regex} \\ d\ \text{Empty} / d\ c &= \text{Empty} \\ d\ \text{Epsilon} / d\ c &= \text{Empty} \\ d\ \text{Singleton}\ x / d\ c &\mathbf{with}\ c \stackrel{?}{=} x \\ \dots \mid \text{yes}\ p &= \text{Epsilon} \\ \dots \mid \text{no}\ \neg p &= \text{Empty} \\ d\ l \cdot r / d\ c &\mathbf{with}\ \varepsilon?\ l \\ \dots \mid \text{yes}\ p &= ((d\ l / d\ c) \cdot r) \mid (d\ r / d\ c) \\ \dots \mid \text{no}\ \neg p &= (d\ l / d\ c) \cdot r \\ d\ l \mid r / d\ c &= (d\ l / d\ c) \mid (d\ r / d\ c) \\ d\ r \star / d\ c &= (d\ r / d\ c) \cdot (r \star) \end{aligned}$$

In order to use the derivative of  $r$  to compute a parse tree for  $r$ , we need to be able to convert a tree for  $d\ r / d\ x$  to a tree for  $r$ . We do this with the function *integralTree*:

$$\text{integralTree} : (r : \text{Regex}) \rightarrow \text{Tree}\ (d\ r / d\ x) \rightarrow \text{Tree}\ r$$

We can also define it with exactly the same case distinction as we used to define  $d\_ / d\_$ .

The code for the parser, *dmatch*, itself is very short. As we sketched, for an empty string we check that the expression matches the empty string, while for a non-empty string we use the derivative to perform a recursive call.

$$\begin{aligned} \text{dmatch} : \{\{ iND : ENondet \in es \} \} &\rightarrow \text{Regex} \times \text{String} \xrightarrow{es} \text{Tree} \circ \text{Pair.fst} \\ \text{dmatch}\ (r, Nil) &\mathbf{with}\ \varepsilon?\ r \\ \dots \mid \text{yes}\ (ms, \_) &= \text{Pure}\ ms \\ \dots \mid \text{no}\ \neg p &= \text{fail} \\ \text{dmatch}\ (r, (x :: xs)) &= \text{integralTree}\ r\ \langle \$ \rangle\ \text{call}\ ((d\ r / d\ x), xs) \end{aligned}$$

Since *dmatch* always consumes a character before going in recursion, we can easily prove that each recursive call only leads to finitely many other calls. This means that for each input value we can unfold the recursive step in the definition a bounded number of times and get a computation with no recursion. Intuitively, this means that *dmatch* terminates on all input. If we are going to give a formal proof of termination, we should first determine the correct formalization of this notion. For that, we need to consider what it means to have no recursion in the unfolded computation. A definition for the *while* loop using general recursion looks as follows:

$$\begin{aligned}
& \text{while} : \{ iRec : ERec\ a\ (K\ a) \in es \} \rightarrow \\
& \quad (a \rightarrow Bool) \rightarrow (a \rightarrow a) \rightarrow (a \rightarrow Free\ es\ a) \\
& \text{while}\ cond\ body\ i = \text{if}\ cond\ i\ \text{then}\ Pure\ i\ \text{else}\ (call\ (body\ i))
\end{aligned}$$

We would like to say that some *while* loops terminate, yet the definition of *while* always contains a *call* in it. Thus, the requirement should not be that there are no more calls left, but that these calls are irrelevant.

Intuitively, we could say that a definition *S* calling *f* terminates if we make the unfolded definition into a *Partial* computation by replacing *call* with *fail*, the definition terminates if the *Partial* computation still works the same, i.e. it refines *S*. However, this mixes the concepts of correctness and termination. We want to see that the *Partial* computation gives some output, without caring about which output this is. Thus, we should only have a trivial postcondition. We formalize this idea in the *terminates-in* predicate.

$$\begin{aligned}
& \text{terminates-in} : (pts : PTs\ es) \\
& \quad (f : C \overset{es}{\rightharpoonup} R) (S : Free\ (eff\ C\ R :: es)\ a) \rightarrow \mathbb{N} \rightarrow Set \\
& \text{terminates-in}\ pts\ f\ (Pure\ x)\ n = \top \\
& \text{terminates-in}\ pts\ f\ (Step \in Head\ c\ k)\ Zero = \perp \\
& \text{terminates-in}\ pts\ f\ (Step \in Head\ c\ k)\ (Succ\ n) = \\
& \quad \text{terminates-in}\ pts\ f\ (f\ c \gg k)\ n \\
& \text{terminates-in}\ pts\ f\ (Step \in Tail\ i)\ c\ k\ n = \\
& \quad \text{lookupPT}\ pts\ i\ c\ (\lambda x \rightarrow \text{terminates-in}\ pts\ f\ (k\ x)\ n)
\end{aligned}$$

Since *dmatch* always consumes a character before going in recursion, we can bound the number of recursive calls with the length of the input string. The proof goes by induction on this string. Unfolding the recursive *call* gives *integralTree*  $\langle \$ \rangle$  *dmatch* (*d r* / *d x* , *xs*), which we rewrite using the associativity monad law in a lemma called *terminates-fmap*.

$$\begin{aligned}
& \text{dmatchTerminates} : (r : Regex) (xs : String) \rightarrow \\
& \quad \text{terminates-in}\ (ptAll :: Nil)\ (\text{dmatch})\ (\text{dmatch}\ (r\ ,\ xs))\ (\text{length}\ xs) \\
& \text{dmatchTerminates}\ r\ Nil\ \text{with}\ \varepsilon? r \\
& \text{dmatchTerminates}\ r\ Nil\ | \text{yes}\ p = tt \\
& \text{dmatchTerminates}\ r\ Nil\ | \text{no}\ \neg p = tt \\
& \text{dmatchTerminates}\ r\ (x :: xs) = \text{terminates-fmap}\ (\text{length}\ xs) \\
& \quad (\text{dmatch}\ ((d\ r\ /\ d\ x)\ ,\ xs)) \\
& \quad (\text{dmatchTerminates}\ (d\ r\ /\ d\ x)\ xs)
\end{aligned}$$

To show partial correctness of *dmatch*, we can use the transitivity of the refinement relation. If we apply transitivity, it suffices to show that *dmatch* is a refinement of *match*. Our first step is to show that the derivative operator is correct, i.e. *d r* / *d x* matches those strings *xs* such that *r* matches *x :: xs*.

$$\begin{aligned}
& \text{derivativeCorrect} : \forall r \rightarrow \\
& \quad \text{Match}\ (d\ r\ /\ d\ x)\ xs\ y \rightarrow \text{Match}\ r\ (x :: xs)\ (\text{integralTree}\ r\ y)
\end{aligned}$$

The proof mirrors the definitions of these functions, being structured as a case distinction on the regular expression.

Before we can prove the correctness of *dmatch* in terms of *match*, it turns out that we also need to describe *match* itself better. The meaning of our goal, to show that *match* is refined by *dmatch*, is to prove that the output of *dmatch* is a subset of that of *match*. Since *match* makes use of *allSplits*, we first prove that *allSplits* returns all possible splittings of a string.

$$\text{allSplitsComplete} : (xs\ ys\ zs : \text{String}) (P : \text{String} \times \text{String} \rightarrow \text{Set}) \rightarrow \\ \text{wpMatch} (\text{allSplits}\ xs) P \rightarrow (xs == ys \uplus zs) \rightarrow P (ys, zs)$$

The proof mirrors *allSplits*, performing induction on *xs*.

Using the preceding lemmas, we can prove the partial correctness of *dmatch* by showing it refines *match*:

$$\text{dmatchSound} : \forall r\ xs \rightarrow \\ \text{wpMatch} (\text{match}\ (r, xs)) \sqsubseteq \text{wpMatch} (\text{dmatch}\ (r, xs))$$

Since we need to perform the case distinctions of *match* and of *dmatch*, the proof is longer than that of *matchSoundness*. Despite the length, most of it consists of performing the case distinction, then giving a simple argument for each case. Therefore, we omit the proof.

With the proof of *dmatchSound* finished, we can conclude that *dmatch* always returns a correct parse tree, i.e. that *dmatch* is sound. However, *dmatch* is *not* complete with respect to the *Match* relation: since *dmatch* never makes a nondeterministic choice, it will not return all possible parse trees as specified by *Match*, only the first tree that it encounters. Still, we can express the property that *dmatch* finds a parse tree if it exists. In other words, we will show that if there is a valid parse tree, *dmatch* returns any parse tree (and this is a valid tree by *dmatchSound*). To express that *dmatch* returns something, we use a trivially true postcondition, and replace the demonic choice of the *ptAll* semantics with the angelic choice of *ptAny*:

$$\text{dmatchComplete} : \forall r\ xs\ y \rightarrow \text{Match}\ r\ xs\ y \rightarrow \\ \text{wp} (\text{ptRec}\ \text{matchSpec} :: \text{ptAny} :: \text{Nil}) (\text{dmatch}\ (r, xs)) (\lambda \_ \rightarrow \top)$$

The proof is short, since *dmatch* can only *fail* when it encounters an empty string and a regex that does not match the empty string, contradicting the assumption immediately:

$$\text{dmatchComplete}\ r\ \text{Nil}\ y\ H \text{ with } \varepsilon? \ r \\ \dots \mid \text{yes } p = tt \\ \dots \mid \text{no } \neg p = \neg p (\_, H) \\ \text{dmatchComplete}\ r\ (x :: xs)\ y\ H\ y'\ H' = tt$$

Note that *dmatchComplete* does not show that *dmatch* terminates: the semantics for the recursive case assume that *dmatch* always returns some value *y'*.

In the proofs of *dmatchSound* and *dmatchComplete*, we demonstrate the power of predicate transformer semantics for effects: by separating syntax and semantics, we can easily describe different aspects (soundness and completeness) of the one definition of *dmatch*. Since the soundness and completeness result we have proved imply partial correctness, and partial correctness and termination imply total correctness, we can conclude that *dmatch* is a totally correct parser for regular languages.

Note the correspondences of this section with a Functional Pearl by Harper [Har99], which also uses the parsing of regular languages as an example of principles of functional software development. Starting out with defining regular expressions as a data type and the language associated with each expression as an inductive relation, both use the relation to implement essentially the same *match* function, which does not terminate. In both, the partial correctness proof of *match* uses a specification expressed as a postcondition, based on the inductive relation representing the language of a given regular expression. Where we use nondeterminism to handle the concatenation operator, Harper uses a continuation-passing parser for control flow. Since the continuations take the unparsed remainder of the string, they correspond almost directly to the *EParser* effect of the following section. Another main difference between our implementation and Harper’s is in the way the non-termination of *match* is resolved. Harper uses the derivative operator to rewrite the expression in a standard form which ensures that the *match* function terminates. We use the derivative operator to implement a different matcher *dmatch* which is easily proved to be terminating, then show that *match*, which we have already proven partially correct, is refined by *dmatch*. The final major difference is that Harper uses manual verification of the program and our work is formally computer-verified. Although our development takes more work, the correctness proofs give more certainty than the informal arguments made by Harper. In general, choosing between informal reasoning and formal verification will always be a trade-off between speed and accuracy.

## 7 Parsing as effect

In the previous sections, we wrote parsers as nondeterministic functions. For more complicated classes of languages than regular expressions, explicitly passing around the string to be parsed becomes cumbersome quickly. The traditional solution is to switch from nondeterminism to stateful nondeterminism, where the state contains the unparsed portion of the string [Hut92]. The combination of nondeterminism and state can be represented by the *Parser* monad:

$$\begin{aligned} \text{Parser} &: \text{Set} \rightarrow \text{Set} \\ \text{Parser } a &= \text{String} \rightarrow \text{List } (a \times \text{String}) \end{aligned}$$

Since our development makes use of algebraic effects, we can introduce the effect of mutable state without having to change existing definitions. We introduce this using the *EParser* effect, which has one command *Symbol*. Calling

*Symbol* will return the current symbol in the state (advancing the state by one) or fail if all symbols have been consumed.

```

data CParser : Set where
  Symbol : CParser
  RParser : CParser → Set
  RParser Symbol = Char
  EParser = eff CParser RParser
  symbol : {iP : EParser ∈ es} → Free es Char
  symbol {iP} = Step iP Symbol Pure

```

We could add more commands such as *EOF* for detecting the end of the input, but we do not need them in the current development. In the semantics we will define that parsing was successful if the input string has been completely consumed.

Note that *EParser* is not sufficient by itself to implement even simple parsers such as *dmatch*: we need to be able to choose between parsing the next character or returning a value for the empty string. This is why we usually combine *EParser* with the nondeterminism effect *ENondet*, and the general recursion effect *ERec*.

The denotational semantics of a parser in the *Free* monad take the form of a fold, handling each command in the *Parser* monad.

```

toParser : Free (ENondet :: EParser :: Nil) a → Parser a
toParser (Pure x) Nil = (x , Nil) :: Nil
toParser (Pure x) (– :: –) = Nil
toParser (Step ∈ Head Fail k) xs = Nil
toParser (Step ∈ Head Choice k) xs =
  toParser (k True) xs ++ toParser (k False) xs
toParser (Step (∈ Tail ∈ Head) Symbol k) Nil = Nil
toParser (Step (∈ Tail ∈ Head) Symbol k) (x :: xs) = toParser (k x) xs

```

In this article, we are more interested in the predicate transformer semantics of *EParser*. Since the semantics of the *EParser* effect refer to a state, the predicates depend on this state. We can incorporate a mutable state of type *s* in predicate transformer semantics by replacing the propositions in *Set* with predicates over the state in  $s \rightarrow \text{Set}$ . We define the resulting type of stateful predicate transformers for an effect *e* to be  $PT^S s e$ , as follows:

```

record PTS (s : Set) (e : Effect) : Set where
  constructor mkPTS
  field
    pt : (c : C e) → (R e c → s → Set) → s → Set
    mono : ∀ c P P' → (∀ x t → P x t → P' x t) → pt c P ⊆ pt c P'

```

If we define  $PTs^S$  and *lookupPTS* analogously to *PTs* and *lookupPT*, we can find a weakest precondition that incorporates the current state:



$$\begin{aligned}
wp^S &: (pts : PTS^S s es) \rightarrow Free es a \rightarrow (a \rightarrow s \rightarrow Set) \rightarrow s \rightarrow Set \\
wp^S pts (Pure x) P &= P x \\
wp^S pts (Step i c k) P &= lookupPTS pts i c \lambda x \rightarrow wp^S pts (k x) P
\end{aligned}$$

In this definition for  $wp^S$ , we assume that all effects share access to one mutable variable of type  $s$ . We can allow for more variables by setting  $s$  to be a product type over the effects. With a suitable modification of the predicate transformers, we could set it up so that each effect can only modify its own associated variable. Thus, the previous definition is not limited in generality by writing it only for one variable.

To give the predicate transformer semantics of the *EParser* effect, we need to choose the meaning of failure, for the case where the next character is needed and all characters have already been consumed. Since we want all results returned by the parser to be correct, we use demonic choice and the *ptAll* predicate transformer as the semantics for *ENondet*. Using *ptAll*'s semantics for the *Fail* command gives the following semantics for the *EParser* effect.

$$\begin{aligned}
ptParse &: PT^S String EParser \\
PTS.pt ptParse Symbol P Nil &= \top \\
PTS.pt ptParse Symbol P (x :: xs) &= P x xs
\end{aligned}$$

With the predicate transformer semantics of *EParser*, we can define the language accepted by a parser in the *Free* monad as a predicate over strings: a string  $xs$  is in the language of a parser  $S$  if the postcondition “all characters have been consumed” is satisfied.

$$\begin{aligned}
empty? &: List a \rightarrow Set \\
empty? Nil &= \top \\
empty? (_ :: _) &= \perp \\
\_ \in [\_] &: String \rightarrow Free (ENondet :: EParser :: Nil) a \rightarrow Set \\
xs \in [S] &= wp^S (ptAll :: ptParse :: Nil) S (\lambda _ \rightarrow empty?) xs
\end{aligned}$$

## 8 Parsing context-free languages

In Section 6, we developed and formally verified a parser for regular languages. The class of regular languages is small, and does not include most programming languages. A class of languages that is more expressive than the regular languages, while remaining tractable in parsing is that of the *context-free language*. The expressiveness of context-free languages is enough to cover most programming languages used in practice [AU77]. We will represent context-free languages in Agda by giving a grammar in the style of Brink, Holdermans, and Löh [BHL10], in a similar way as we represent a regular language using an element of the *Regex* type. Following their development, we parametrize our definitions over a collection of nonterminal symbols.

```

record GrammarSymbols : Set where
  field
    Nonterm : Set
     $\llbracket \_ \rrbracket : \text{Nonterm} \rightarrow \text{Set}$ 
     $\_ \stackrel{?}{=} \_ : \text{Decidable } \{A = \text{Nonterm}\} \_ == \_$ 

```

The elements of the type *Char* are the *terminal* symbols. The elements of the type *Nonterm* are the *nonterminal* symbols, representing the language constructs. As for *Char*, we also need to be able to decide the equality of nonterminals. The (disjoint) union of *Char* and *Nonterm* gives all the symbols that we can use in defining the grammar.

```

Symbol = Either Char Nonterm
Symbols = List Symbol

```

For each nonterminal *A*, our goal is to parse a string into a value of type  $\llbracket A \rrbracket$ , based on a set of production rules. A production rule  $A \rightarrow xs$  gives a way to expand the nonterminal *A* into a list of symbols *xs*, such that successfully matching each symbol of *xs* with parts of a string gives a match of the string with *A*. Since matching a nonterminal symbol *B* with a (part of a) string results in a value of type  $\llbracket B \rrbracket$ , a production rule for *A* is associated with a *semantic function* that takes all values arising from submatches and returns a value of type  $\llbracket A \rrbracket$ , as expressed by the following type:

```

 $\llbracket \_ \rrbracket : \text{Symbols} \rightarrow \text{Nonterm} \rightarrow \text{Set}$ 
 $\llbracket Nil \parallel A \rrbracket = \llbracket A \rrbracket$ 
 $\llbracket Inl\ x \ ::\ xs \parallel A \rrbracket = \llbracket xs \parallel A \rrbracket$ 
 $\llbracket Inr\ B \ ::\ xs \parallel A \rrbracket = \llbracket B \rrbracket \rightarrow \llbracket xs \parallel A \rrbracket$ 

```

Now we can define the type of production rules. A rule of the form  $A \rightarrow BcD$  is represented as  $prod\ A\ (Inr\ B \ ::\ Inl\ c \ ::\ Inr\ D \ ::\ Nil)\ f$  for some *f*.

```

record Prod : Set where
  constructor prod
  field
    lhs : Nonterm
    rhs : Symbols
    sem :  $\llbracket rhs \parallel lhs \rrbracket$ 

```

We use the abbreviation *Prods* to represent a list of productions, and a grammar will consist of the list of all relevant productions.

We want to show that a generally recursive function making use of the effects *EParser* and *ENondet* can parse any context-free grammar. To show this claim, we implement a function *fromProds* that constructs a parser for any context-free grammar given as a list of *Prods*, then formally verify the correctness of *fromProds*. Our implementation mirrors the definition of the *generateParser* function by Brink, Holdermans, and Löh, differing in the naming and in the

system that the parser is written in: our implementation uses the *Free* monad and algebraic effects, while Brink, Holdermans, and Löh use a monad *Parser* that is based on parser combinators.

We start by defining two auxiliary types, used as abbreviations in our code.

```
FreeParser = Free (eff Nonterm  $\llbracket \_ \rrbracket$  :: ENondet :: EParser :: Nil)
record ProdRHS (A : Nonterm) : Set where
  constructor prodrhs
  field
    rhs : Symbols
    sem :  $\llbracket$  rhs  $\parallel$  A  $\rrbracket$ 
```

The core algorithm for parsing a context-free grammar consists of the following functions, calling each other in mutual recursion:

```
fromProds  : (A : Nonterm) → FreeParser  $\llbracket$  A  $\rrbracket$ 
filterLHS  : (A : Nonterm) → Prods → List (ProdRHS A)
fromProd   : ProdRHS A → FreeParser  $\llbracket$  A  $\rrbracket$ 
buildParser : (xs : Symbols) → FreeParser ( $\llbracket$  xs  $\parallel$  A  $\rrbracket$  →  $\llbracket$  A  $\rrbracket$ )
exact      : a → Char → FreeParser a
```

The main function is *fromProds*: given a nonterminal, it selects the productions with this nonterminal on the left hand side using *filterLHS*, and makes a non-deterministic choice between the productions.

```
filterLHS A Nil = Nil
filterLHS A (prod lhs rhs sem :: ps) with A  $\stackrel{?}{=}$  lhs
... | yes refl = prodrhs rhs sem :: filterLHS A ps
... | no _    = filterLHS A ps
fromProds A = foldr (choice) (fail) (map fromProd (filterLHS A prods))
```

The function *fromProd* takes a single production and tries to parse the input string using this production. It then uses the semantic function of the production to give the resulting value.

```
fromProd (prodrhs rhs sem) = buildParser rhs  $\ggg$   $\lambda$  f → Pure (f sem)
```

The function *buildParser* iterates over the *Symbols*, calling *exact* for each literal character symbol, and making a recursive *call* to *fromProds* for each nonterminal symbol.

```
buildParser Nil = Pure id
buildParser (Inl x :: xs) = exact tt x  $\ggg$   $\lambda$  _ → buildParser xs
buildParser (Inr B :: xs) = do
  x ← call B
  o ← buildParser xs
  Pure  $\lambda$  f → o (f x)
```

Finally, *exact* uses the *symbol* command to check that the next character in the string is as expected, and *fails* if this is not the case.

$$\text{exact } x \ t = \text{symbol} \gg \lambda t' \rightarrow \text{if } t \stackrel{?}{=} t' \text{ then Pure } x \text{ else fail}$$

## 9 Partial correctness of the parser

Partial correctness of the parser is relatively simple to show, as soon as we have a specification. Since we want to prove that *fromProds* correctly parses any given context free grammar given as an element of *Prods*, the specification consists of a relation between many sets: the production rules, an input string, a nonterminal, the output of the parser, and the remaining unparsed string. Due to the many arguments, the notation is unfortunately somewhat unwieldy. To make it a bit easier to read, we define two relations in mutual recursion, one for all productions of a nonterminal, and for matching a string with a single production rule.

```

data  $\_ \vdash \_ \in \llbracket \_ \rrbracket \Rightarrow \_, \_ \text{ prods}$  where
  Produce :  $\text{prod lhs rhs sem} \in \text{prods} \rightarrow$ 
     $\text{prods} \vdash xs \sim \text{rhs} \Rightarrow f, ys \rightarrow$ 
     $\text{prods} \vdash xs \in \llbracket \text{lhs} \rrbracket \Rightarrow f \text{ sem}, ys$ 
data  $\_ \vdash \_ \sim \_ \Rightarrow \_, \_ \text{ prods}$  where
  Done :  $\text{prods} \vdash xs \sim \text{Nil} \Rightarrow \text{id}, xs$ 
  Next :  $\text{prods} \vdash xs \sim ps \Rightarrow o, ys \rightarrow$ 
     $\text{prods} \vdash (x :: xs) \sim (\text{Inl } x :: ps) \Rightarrow o, ys$ 
  Call :  $\text{prods} \vdash xs \in \llbracket A \rrbracket \Rightarrow o, ys \rightarrow$ 
     $\text{prods} \vdash ys \sim ps \Rightarrow f, zs \rightarrow$ 
     $\text{prods} \vdash xs \sim (\text{Inr } A :: ps) \Rightarrow (\lambda g \rightarrow f (g o)), zs$ 

```

With these relations, we can define the specification *parserSpec* to be equal to  $\_ \vdash \_ \in \llbracket \_ \rrbracket \Rightarrow \_, \_$  (up to reordering some arguments), and show that *fromProds* refines this specification. To state that the refinement relation holds, we first need to determine the semantics of the effects. We choose *ptAll* as the semantics of nondeterminism, since we want to ensure all output of the parser is correct.

```

pts prods = ptRec (parserSpec prods) :: ptAll :: ptParse :: Nil
wpFromProd prods = wpS (pts prods)
partialCorrectness : (prods : Prods) (A : Nonterm) →
  wpSpec [⊤, (parserSpec prods A)] ⊆
  wpFromProd prods (fromProds prods A)

```

Let us fix the production rules *prods*. How do we prove the partial correctness of a parser for *prods*? Since the structure of *fromProds* is of a nondeterministic choice between productions to be parsed, and we want to show that all alternatives for a choice result in success, we will first give a lemma expressing the

correctness of each alternative. Correctness in this case is expressed by the semantics of a single production rule, i.e. the  $\_ \vdash \_ \sim \_ \Rightarrow \_$  relation. Thus, we want to prove the following lemma:

$$\begin{aligned} \text{parseStep} &: \forall A \, xs \, P \, str \rightarrow \\ &(\forall o \, str' \rightarrow \text{prods} \vdash str \sim xs \Rightarrow o, str' \rightarrow P \, o \, str') \rightarrow \\ &\text{wpFromProd} \, \text{prods} \, (\text{buildParser} \, \text{prods} \, xs) \, P \, str \end{aligned}$$

The lemma can be proved by reproducing the case distinctions used to define *buildParser*; there is no complication apart from having to use the *wpToBind* lemma to deal with the  $\_ \gg \_$  operator in a few places.

$$\begin{aligned} \text{parseStep} \, A \, \text{Nil} \, P \, t \, H &= H \, \text{id} \, t \, \text{Done} \\ \text{parseStep} \, A \, (\text{Inl} \, x :: xs) \, P \, \text{Nil} \, H &= tt \\ \text{parseStep} \, A \, (\text{Inl} \, x :: xs) \, P \, (x' :: t) \, H &\textbf{with} \, x \stackrel{?}{=} x' \\ \dots \mid \text{yes refl} &= \text{parseStep} \, A \, xs \, P \, t \, \lambda o \, t' \, H' \rightarrow H \, o \, t' \, (\text{Next} \, H') \\ \dots \mid \text{no } \neg p &= tt \\ \text{parseStep} \, A \, (\text{Inr} \, B :: xs) \, P \, t \, H \, o \, t' \, Ho &= \\ &\text{wpToBind} \, (\text{buildParser} \, \text{prods} \, xs) \, - \, - \\ &(\text{parseStep} \, A \, xs \, - \, t' \, \lambda o' \, str' \, Ho' \rightarrow H \, - \, - \, (\text{Call} \, Ho \, Ho')) \end{aligned}$$

To combine the *parseStep* for each of the productions that the nondeterministic choice is made between, it is tempting to define another lemma *filterStep* by induction on the list of productions. But we must be careful that the productions that are used in the *parseStep* are the full list *prods*, not the sublist *prods'* used in the induction step. Additionally, we must also make sure that *prods'* is indeed a sublist, since using an incorrect production rule in the *parseStep* will result in an invalid result. Thus, we parametrise *filterStep* by a list *prods'* and a proof that it is a sublist of *prods*. Again, the proof uses the same distinction as *fromProds* does, and uses the *wpToBind* lemma to deal with the  $\_ \gg \_$  operator.

$$\begin{aligned} \text{filterStep} &: \forall \text{prods}' \rightarrow (p \in \text{prods}' \rightarrow p \in \text{prods}) \rightarrow \\ &\forall A \rightarrow \text{wpSpec} \, [\top, \text{parserSpec} \, \text{prods} \, A] \sqsubseteq \text{wpFromProd} \, \text{prods} \\ &(\text{foldr} \, (\text{choice}) \, (\text{fail}) \, (\text{map} \, (\text{fromProd} \, \text{prods}) \, (\text{filterLHS} \, \text{prods} \, A \, \text{prods}')))) \\ \text{filterStep} \, \text{Nil} \, \text{subset} \, A \, P \, xs \, H &= tt \\ \text{filterStep} \, (\text{prod} \, lhs \, rhs \, sem :: \text{prods}') \, \text{subset} \, A \, P \, xs \, H &\textbf{with} \, A \stackrel{?}{=} lhs \\ \text{filterStep} \, (\text{prod} \, A \, rhs \, sem :: \text{prods}') \, \text{subset} \, A \, P \, xs \, (-, H) \mid \text{yes refl} &= \\ &\text{wpToBind} \, (\text{buildParser} \, \text{prods} \, rhs) \, - \, - \\ &(\text{parseStep} \, A \, rhs \, - \, xs \, \lambda o \, t' \, H' \rightarrow H \, - \, - \, (\text{Produce} \, (\text{subset} \in \text{Head}) \, H')) \\ &, \text{filterStep} \, \text{prods}' \, (\text{subset} \circ \in \text{Tail}) \, A \, P \, xs \, (-, H) \\ \dots \mid \text{no } \neg p &= \text{filterStep} \, \text{prods}' \, (\text{subset} \circ \in \text{Tail}) \, A \, P \, xs \, H \end{aligned}$$

With these lemmas, *partialCorrectness* just consists of applying *filterStep* to the subset of *prods* consisting of *prods* itself.

## 10 Termination of the parser

To show termination we need a somewhat more subtle argument: since we are able to call the same nonterminal repeatedly, termination cannot be shown sim-

ply by inspecting each alternative in the definition. Consider the grammar given by  $E \rightarrow aE; E \rightarrow b$ , where we see that the string that matches  $E$  in the recursive case is shorter than the original string, but the definition itself can be expanded to unbounded length. By taking into account the current state, i.e. the string to be parsed, in the variant, we can show that a decreasing string length leads to termination.

But not all grammars feature this decreasing string length in the recursive case, with the most pathological case being those of the form  $E \rightarrow E$ . The issues do not only occur in edge cases: the grammar  $E \rightarrow E + E; E \rightarrow 1$  representing very simple expressions will already result in non-termination for *fromProds* as it will go in recursion on the first non-terminal without advancing the input string. Since the position in the string and current nonterminal together fully determine the state of *fromParsers*, it will not terminate. We need to ensure that the grammars passed to the parser do not allow for such loops.

Intuitively, the condition on the grammars should be that they are not *left-recursive*, since in that case, the parser should always advance its position in the string before it encounters the same nonterminal. This means that the number of recursive calls to *fromProds* is bounded by the length of the string times the number of different nonterminals occurring in the production rules. The type we will use to describe the predicate “there is no left recursion” is constructively somewhat stronger: we define a left-recursion chain from  $A$  to  $B$  to be a sequence of nonterminals  $A, \dots, A_i, A_{i+1}, \dots, B$ , such that for each adjacent pair  $A_i, A_{i+1}$  in the chain, there is a production of the form  $A_{i+1} \rightarrow B_1 B_2 \dots B_n A_i \dots$ , where  $B_1 \dots B_n$  are all nonterminals. In other words, we can advance the parser to  $A$  starting in  $B$  without consuming a character. Disallowing (unbounded) left recursion is not a limitation for our parsers: Brink, Holdermans, and Löh [BHL10] have shown that the *left-corner transform* can transform left-recursive grammars into an equivalent grammar without left recursion. Moreover, they have implemented this transform, including formal verification, in Agda. In this work, we assume that the left-corner transform has already been applied if needed, so that there is an upper bound on the length of left-recursive chains in the grammar.

We formalize one link of this left-recursive chain in the type *LRec*, while a list of such links forms the *Chain* data type.

```
record LRec (prods : Prods) (A B : Nonterm) : Set where
  field
```

```
    rec : prod A (map Inr xs ++ (Inr B :: ys)) sem ∈ prods
```

(We leave *xs*, *ys* and *sem* as implicit fields of *LRec*, since they are fixed by the type of *rec*.)

```
data Chain (prods : Prods) : Nonterm → Nonterm → Set where
```

```
  Nil : Chain prods A A
```

```
  _ :: _ : LRec prods B A → Chain prods A C → Chain prods B C
```

Now we say that a set of productions has no left recursion if all such chains have an upper bound on their length.

$$\begin{aligned}
& \text{chainLength} : \text{Chain prods } A \ B \rightarrow \mathbb{N} \\
& \text{chainLength } \text{Nil} = 0 \\
& \text{chainLength } (c :: cs) = \text{Succ } (\text{chainLength } cs) \\
& \text{leftRecBound} : \text{Prods} \rightarrow \mathbb{N} \rightarrow \text{Set} \\
& \text{leftRecBound prods } n = (cs : \text{Chain prods } A \ B) \rightarrow \text{chainLength } cs < n
\end{aligned}$$

If we have this bound on left recursion, we are able to prove termination, since each call to *fromProds* will be made either after we have consumed an extra character, or it is a left-recursive step, of which there is an upper bound on the sequence.

This informal proof fits better with a different notion of termination than in the petrol-driven semantics. The petrol-driven semantics are based on a syntactic argument: we know a computation terminates because expanding the call tree will eventually result in no more *calls*. Here, we want to capture the notion that a recursive definition terminates if all recursive calls are made to a smaller argument, according to a well-founded relation.

**Definition 3 ([Acz77])** *In intuitionistic type theory, we say that a relation  $\_ \prec \_$  on a type  $a$  is well-founded if all elements  $x : a$  are accessible, which is defined by (well-founded) recursion to be the case if all elements in the downset of  $x$  are accessible.*

$$\begin{aligned}
& \text{data Acc } (\_ \prec \_ : a \rightarrow a \rightarrow \text{Set}) : a \rightarrow \text{Set} \text{ where} \\
& \text{acc} : (\forall y \rightarrow y \prec x \rightarrow \text{Acc } \_ \prec \_ y) \rightarrow \text{Acc } \_ \prec \_ x
\end{aligned}$$

To see that this is equivalent to the definition of well-foundedness in set theory, recall that a relation  $\_ \prec \_$  on a set  $a$  is well-founded if and only if there is a monotone function from  $a$  to a well-founded order. Since all inductive data types are well-founded, and the termination checker ensures that the argument to *acc* is a monotone function, there is a function from  $x : a$  to  $\text{Acc } \_ \prec \_ x$  if and only if  $\_ \prec \_$  is a well-founded relation in the set-theoretic sense.

The condition that all calls are made to a smaller argument is related to the notion of a loop *variant* in imperative languages. While an invariant is a predicate that is true at the start and end of each looping step, the variant is a relation that holds between successive looping steps.

**Definition 4** *Given a recursive definition  $f : I \overset{es}{\rightsquigarrow} O$ , a relation  $\_ \prec \_$  on  $C$  is a recursive variant if for each argument  $c$ , and each recursive call made to  $c'$  in the evaluation of  $f \ c$ , we have  $c' \prec c$ . Formally:*

$$\begin{aligned}
& \text{variant}' : (\text{pts} : \text{PTS}^S \ s \ (\text{eff } C \ R :: es)) (f : C \overset{es}{\rightsquigarrow} R) \\
& \quad (\_ \prec \_ : (C \times s) \rightarrow (C \times s) \rightarrow \text{Set}) \\
& \quad (c : C) (t : s) (S : \text{Free } (\text{eff } C \ R :: es) \ a) \rightarrow s \rightarrow \text{Set} \\
& \text{variant}' \text{ pts } f \ \_ \prec \_ \ c \ t \ (\text{Pure } x) \ t' = \top \\
& \text{variant}' \text{ pts } f \ \_ \prec \_ \ c \ t \ (\text{Step } \in \text{Head } c' \ k) \ t' \\
& \quad = ((c', t') \prec (c, t)) \times \text{lookupPTS } \text{pts} \in \text{Head } c'
\end{aligned}$$

$$\begin{aligned}
& (\lambda x \rightarrow \text{variant}' \text{pts} f \_ \prec \_ c t (k x)) t' \\
\text{variant}' \text{pts} f \_ \prec \_ c t (\text{Step } (\in \text{Tail } i) c' k) t' \\
& = \text{lookupPTS } \text{pts} (\in \text{Tail } i) c' (\lambda x \rightarrow \text{variant}' \text{pts} f \_ \prec \_ c t (k x)) t' \\
\text{variant} : (\text{pts} : \text{PTS}^S s (\text{eff } C R :: \text{es})) (f : C \overset{\text{es}}{\rightsquigarrow} R) \rightarrow \\
& (\_ \prec \_ : (C \times s) \rightarrow (C \times s) \rightarrow \text{Set}) \rightarrow \text{Set} \\
\text{variantpts} f \_ \prec \_ = \forall c t \rightarrow \text{variant}' \text{pts} f \_ \prec \_ c t (f c) t
\end{aligned}$$

Note that *variant* depends on the semantics *pts* we give to the recursive function *f*. We cannot derive the semantics in *variant* from the structure of *f* as we do for the petrol-driven semantics, since we do not yet know whether *f* terminates. Using *variant*, we can define another termination condition on *f*: there is a well-founded variant for *f*.

$$\begin{aligned}
& \mathbf{record} \text{Termination} (\text{pts} : \text{PTS}^S s (\text{eff } C R :: \text{es})) (f : C \overset{\text{es}}{\rightsquigarrow} R) : \text{Set} \mathbf{where} \\
& \mathbf{field} \\
& \_ \prec \_ : (C \times s) \rightarrow (C \times s) \rightarrow \text{Set} \\
& w - f : \forall c t \rightarrow \text{Acc } \_ \prec \_ (c, t) \\
& \text{var} : \text{variant } \text{pts} f \_ \prec \_
\end{aligned}$$

A generally recursive function that terminates in the petrol-driven semantics also has a well-founded variant, given by the well-order  $\_ \prec \_$  on the amount of fuel consumed by each call. The converse also holds: if we have a descending chain of calls *cs* after calling *f* with argument *c*, we can use induction on the type  $\text{Acc } \_ \prec \_ c$  to bound the length of *cs*. This bound gives the amount of fuel consumed by evaluating a call to *f* on *c*.

In our case, the relation *RecOrder* will work as a recursive variant for *fromProds*:

$$\begin{aligned}
& \mathbf{data} \text{RecOrder} (\text{prods} : \text{Prods}) : (x y : \text{Nonterm} \times \text{String}) \rightarrow \text{Set} \mathbf{where} \\
& \text{Adv} : \text{length } \text{str} < \text{length } \text{str}' \rightarrow \\
& \quad \text{RecOrder } \text{prods} (A, \text{str}) (B, \text{str}') \\
& \text{Rec} : \text{length } \text{str} \leq \text{length } \text{str}' \rightarrow \\
& \quad \text{LRec } \text{prods } A B \rightarrow \text{RecOrder } \text{prods} (A, \text{str}) (B, \text{str}')
\end{aligned}$$

With the definition of *RecOrder*, we can complete the correctness proof of *fromProds*, by giving an element of the corresponding *Termination* type. We assume that the length of recursion is bounded by *bound* :  $\mathbb{N}$ .

$$\begin{aligned}
& \text{fromProdsTerminates} : \forall \text{prods } \text{bound} \rightarrow \text{leftRecBound } \text{prods } \text{bound} \rightarrow \\
& \quad \text{Termination } (\text{pts } \text{prods}) (\text{fromProds } \text{prods}) \\
& \text{Termination.} \_ \prec \_ (\text{fromProdsTerminates } \text{prods } \text{bound } H) = \text{RecOrder } \text{prods}
\end{aligned}$$

To show that the relation *RecOrder* is well-founded, we need to show that there is no infinite descending chain starting from some nonterminal *A* and string *str*. The proof is based on iteration on two natural numbers *n* and *k*, which form an upper bound on the number of allowed left-recursive calls in sequence and unconsumed characters in the string respectively. Note that the number *bound*



is an upper bound for  $n$  and the length of the input string is an upper bound for  $k$ . Since each nonterminal in the production will decrease  $n$  and each terminal will decrease  $k$ , we eventually reach the base case  $0$  for either. If  $n$  is zero, we have made more than  $bound$  left-recursive calls, contradicting the assumption that we have bounded left recursion. If  $k$  is zero, we have consumed more than  $length\ str$  characters of  $str$ , also a contradiction.

*Termination*.  $w - f\ (fromProdsTerminates\ prods\ bound\ H)\ A\ str$   
 $=\ acc\ (go\ A\ str\ (length\ str) \leq\text{-refl}\ bound\ Nil \leq\text{-refl})$   
**where**  
 $go : \forall A\ str \rightarrow$   
 $(k : \mathbb{N}) \rightarrow length\ str \leq k \rightarrow$   
 $(n : \mathbb{N}) (cs : Chain\ prods\ A\ B) \rightarrow bound \leq chainLength\ cs + n \rightarrow$   
 $\forall y \rightarrow RecOrder\ prods\ y\ (A , str) \rightarrow Acc\ (RecOrder\ prods)\ y$

Our next goal is that *RecOrder* is a variant for *fromProds*, as abbreviated by the *prodsVariant* type. We cannot follow the definitions of *fromProds* as closely as we did for the partial correctness proof; instead we need a complicated case distinction to keep track of the left-recursive chain we have followed in the proof. For this reason, we split the *parseStep* apart into two lemmas *parseStepAdv* and *parseStepRec*, both showing that *buildParser* maintains the variant. We also use a *filterStep* lemma that calls the correct *parseStep* for each production in the nondeterministic choice.

$prodsVariant = variant'\ (pts\ prods)\ (fromProds\ prods)\ (RecOrder\ prods)$   
 $parseStepAdv : \forall A\ xs\ str\ str' \rightarrow length\ str' < length\ str \rightarrow$   
 $prodsVariant\ A\ str\ (buildParser\ xs)\ str'$   
 $parseStepRec : \forall A\ xs\ str\ str' \rightarrow length\ str' \leq length\ str \rightarrow$   
 $\forall ys \rightarrow prod\ A\ (map\ Inr\ ys \uplus xs)\ sem \in prods \rightarrow$   
 $prodsVariant\ A\ str\ (buildParser\ xs)\ str'$   
 $filterStep : \forall prods' \rightarrow (x \in prods' \rightarrow x \in prods) \rightarrow$   
 $\forall A\ str\ str' \rightarrow length\ str' \leq length\ str \rightarrow$   
 $prodsVariant\ A\ str$   
 $(foldr\ (choice)\ (fail)\ (map\ fromProd\ (filterLHS\ A\ prods')))$   
 $str'$

In the *parseStepAdv*, we deal with the situation that the parser has already consumed at least one character since it was called. This means we can repeatedly use the *Adv* constructor of *RecOrder* to show the variant holds.

In the *parseStepRec*, we deal with the situation that the parser has only encountered nonterminals in the current production. This means that we can use the *Rec* constructor of *RecOrder* to show the variant holds until we consume a character, after which we call *parseStepAdv* to finish the proof.

The lemma *filterStep* shows that the variant holds on all subsets of the production rules, analogously to the *filterStep* of the partial correctness proof. It calls *parseStepRec* since the parser only starts consuming characters after it selects a production rule.

```

filterStep Nil A str str' lt subset = tt
filterStep (prod lhs rhs sem :: prods') subset A str str' lt with A  $\stackrel{?}{=}$  lhs
... | yes refl
    = variant - fmap (pts prods) (fromProds prods) (buildParser rhs)
      (parseStepRec A rhs str str' lt Nil (subset ∈ Head))
      , filterStep prods' (subset ∘ ∈ Tail) A str str' lt
... | no ¬p = filterStep prods' (subset ∘ ∈ Tail) A str str' lt

```

As for partial correctness, we obtain the proof of termination by applying *filterStep* to the subset of *prods* consisting of *prods* itself.

## 11 Conclusions and discussion

Fill this!