

# Context-free Languages, Type Theoretically

Jaro Reinders<sup>1</sup>[0000–0002–6837–3757] and Casper Bach<sup>2</sup>[0000–0003–0622–7639]

<sup>1</sup> Delft University of Technology, Delft, The Netherlands

<sup>2</sup> University of Southern Denmark, Odense, Denmark

**Abstract.** Parsing is the process of recovering structure from strings, an essential part of implementing programming languages. Previous work has shown that formalizing languages and parsers using an idiomatic type theoretic approach can be simple and enlightening. Unfortunately, this approach has only been applied to regular languages, which are not expressive enough for many practical applications. We are working on extending the type theoretic formalization to context-free languages which substantially more expressive. We hope our formalization can serve as a foundation for reasoning about new disambiguation techniques and even more expressive formalisms such as data-dependent grammars.

**Keywords:** Language · Parsing · Type Theory

## 1 Introduction

Parsing remains an open problem

E.g., dealing with ambiguities, but also dealing with programming languages beyond context-free languages

While parsing serves a practical purpose, it is important to also have a deep theoretical understanding of the semantics of parsing

While parsing is often studied using automata theory [CITE], there is also value in studying more denotational approaches to parsing.

Indeed, a main purpose of denotational semantics is to abstract away operational concerns, as such concerns tends to be a hindrance for equational reasoning.

A denotational approach could thus provide a framework for studying and proving the correctness of, e.g., parser optimizations and disambiguation techniques.

It could also provide a building block for obtaining correct parsers of expressive grammar formalisms, such as data-dependent grammars [CITE].

Recent work by Elliott demonstrated that parsers for regular grammars can be given a simple and direct semantics

In turn, we can obtain parsers for such languages that are (practically) correct by construction, by taking their *derivatives*.

While it was well-known [2] that we can parse regular grammars using Brzozowski derivatives, Elliot provides a simple and direct mechanization in Agda of the denotational semantics of these derivatives.

Elliott leaves open the question of how the approach scales to more expressive grammar formalisms, such as context-free languages and beyond.

This paper addresses that question.

Specifically, we study the problem of mechanizing, also in Agda, (1) the denotational semantics of context-free grammars; and (2) a simple and direct denotational semantics of the derivative of context-free grammars.

A main challenge for this mechanization is dealing with the recursive nature of context-free languages.

## 44 1.1 The Challenge with Automated Differentiation of Context-Free 45 Grammars

46 Derivatives provide an automated procedure for differentiation of languages.

47 In this subsection we consider the problem of automatically taking the derivative of  
48 a context-free grammar w.r.t. a given symbol.

49 To illustrate, let us consider the following context-free grammar of palindromic bit  
50 strings:

$$\langle pal \rangle ::= 0 \mid 1 \mid 0 \langle pal \rangle 0 \mid 1 \langle pal \rangle 1$$

51 Say we want to use this grammar to parse the string 0110.

52 The idea of automatic differentiation is this:

53 We first compute the derivative of the grammar w.r.t. the first bit (0) of our bit  
54 string (0110).

55 Let us call this grammar  $\langle pal_0 \rangle$ .

56 Then, we take the derivative of  $\langle pal_0 \rangle$  w.r.t. the next bit (1).

57 We continue this procedure until we either (1) get stuck because the derivative is  
58 invalid, in which case the bit string is not well-formed w.r.t. our grammar, or (2) the  
59 derivative grammar contains the empty production (we will use the symbol  $\epsilon$  to denote  
60 the empty grammar), in which case the bit string is well-formed w.r.t. our grammar.

61 Taking the derivative of the  $\langle pal \rangle$  grammar w.r.t. the bit 0 yields the following derived  
62 grammar:

$$\langle pal_0 \rangle ::= \epsilon \mid \langle pal \rangle 0$$

63 This grammar essentially represents the residual parsing obligations after parsing a  
64 0 bit.

65 The derived grammar contains fewer productions than the original grammar because  
66 we have pruned those productions that started with the bit 1 (because the derivative of  
67 the bit 0 w.r.t. the terminal symbol 1 is invalid).

68 Now, how do we take the derivative of the grammar  $\langle pal_0 \rangle$  w.r.t. the next bit (1) in  
69 our string?

70 A simple solution is to recursively unfold the  $\langle pal \rangle$  non-terminal.

71 Doing so for the  $\langle pal \rangle$  grammar yields the following derived grammar:

$$\langle pal'_0 \rangle ::= \epsilon \mid 00 \mid 10 \mid 0 \langle pal \rangle 00 \mid 1 \langle pal \rangle 10$$

72 By continuing this procedure, with additional recursive unfolding where needed, we  
73 eventually yield a grammar that contains the the empty production  $\epsilon$ , whereby we can  
74 conclude that 0110 is, in fact, a palindromic bit string.

75 However, the recursive unfolding we performed above is not safe to do for all gram-  
76 mars.

77 Consider, for example, the infinitely recursive grammar:

$$\langle rec \rangle ::= \langle rec \rangle$$

78 We cannot ever unfold this grammar to expose a terminal symbol to derive w.r.t.,  
79 akin to the informal procedure we applied above.

80 Another challenge with context-free grammars is how to encode their recursive na-  
81 ture in a proof assistant such as Agda in a way that our encoding of grammars is  
82 strictly positive, and in a way that ensures that automated differentiation—that is, con-  
83 tinuously applying the method we informally illustrated above for taking the derivative  
84 of a grammar w.r.t. a symbol—is guaranteed to terminate.

## 1.2 Contributions

This paper tackles the challenges discussed in the previous section by providing a mechanization in Agda of automated differentiation of a subset of context-free grammars.

The subset we consider corresponds to context-free grammars with non-terminal symbols without mutual recursion.

We conjecture that our approach is compatible with all context-free grammars, at the cost of some additional book-keeping while taking the derivative.

We leave verifying this conjecture as a challenge for future work.

We make the following technical contributions:

1. We provide a semantics in Agda of context-free grammars without mutual recursion.

2. We provide a semantics of automated differentiation for this class of grammars, along with its simple and direct correctness proof.

The rest of this paper is structured as follows: [...]

We assume basic familiarity with Agda.

## 2 Finite Languages

In this section, we introduce background information, namely how we define languages, basic language combinators, and parsers. Our exposition follows Elliott [3]. In Section 3, we extend these concepts to context free languages.

### 2.1 Languages

We define languages as being functions from strings to types.<sup>3</sup>

`Lang = String → Type`

The result type can be thought of as the type of proofs that the string is in the language.

*Remark 1.* Note that a language may admit multiple different proofs for the same string. That is an important difference between the type theoretic approach and the more common set theoretic approach, which models languages as sets of strings.

This is a broad definition of what a language is; it includes languages that are outside the class of context-free languages.

*Example 1.* The language  $a^n b^n c^n$  can be specified as follows:

```
abc : Lang
abc w = Σ[ n ∈ ℕ ] w ≡ repeat n 'a' ++ repeat n 'b' ++ repeat n 'c'
```

We can show that the string `aabbcc` is in this language by choosing  $n$  to be 2, from which the required equality follows by reflexivity after normalization:

```
aabbcc : abc "aabbcc"
aabbcc = 2 , refl
```

Example 1 shows that it is possible to specify languages and prove that certain strings are in those languages, but for practical applications we do not want to be burdened with writing such proofs ourselves. The compiler should be able to decide whether or not your program is valid by itself.

- Agda is too powerful: it can specify undecidable languages
- So, we need to define a simpler language which still supports all the features we need.

[JR] do I need to give an example?

<sup>3</sup> We use `Type` as a synonym for Agda's `Set` to avoid confusion.

## 128 2.2 Basic Language Combinators

129 Let's start with a simple example: POSIX file system permissions. These are usually  
 130 summarized using the characters 'r', 'w', and 'x' if the permissions are granted, or '-' in  
 131 place of the corresponding character if the permission is denied. For example the string  
 132 "r-x" indicates that read and execute permissions are granted, but the write permission  
 133 is denied. The full language can be expressed using the following BNF grammar:

```

134 <permissions> ::= <read> <write> <execute>
135 <read>         ::= '-' | 'r'
136 <write>        ::= '-' | 'w'
137 <execute>      ::= '-' | 'x'

```

|   |   |
|---|---|
| $\_ : \text{Char} \rightarrow \text{Lang}$                                  | $\emptyset : \text{Lang}$   |
| $(' c) w = w \equiv c :: []$  | $\emptyset \_ = \perp$  |
| $\_ \cup \_ : \text{Lang} \rightarrow \text{Lang} \rightarrow \text{Lang}$  | $\epsilon : \text{Lang}$  |
| $(P \cup Q) w = P w \uplus Q w$   | $\epsilon w = w \equiv []$  |
| $\_ * \_ : \text{Lang} \rightarrow \text{Lang} \rightarrow \text{Lang}$     | $\_ \cdot \_ : \text{Type} \rightarrow \text{Lang} \rightarrow \text{Lang}$ |
| $(P * Q) w = \exists [u] \exists [v] w \equiv u ++ v \times P u \times Q v$ | $(A \cdot P) w = A \times P w$  |

Fig. 1. Basic language combinators.

142 This grammar uses three important features: sequencing, choice, and matching char-  
 143 acter literals. We can define these features as combinators in Agda as shown in Figure 1  
 144 and use them to write our permissions grammar as follows:

```

145 permissions = read * write * execute
146 read        = ' '-' | 'r'
147 write       = ' '-' | 'w'
148 execute     = ' '-' | 'x'

```

## 149 2.3 Parsers

150 We want to write a program which can prove for us that a given string is in the language.  
 151 What should this program return for strings that are not in the language? We want to  
 152 make sure our program does find a proof if it exists, so if it does not exist then we want  
 153 a proof that the string is not in the language. We can capture this using a type called  
 154 **Dec** from the Agda standard library. It can be defined as follows:

```

155 data Dec (A : Type) : Type where
156   yes : A → Dec A
157   no  : ¬ A → Dec A

```

158 A parser for a language, then, is a program which can tell us whether any given  
 159 string is in the language or not.

```

160 Parser : Lang → Set
161 Parser P = (w : String) → Dec (P w)

```

162 *Remark 2.* Readers familiar with Haskell might see similarity between this type and the  
163 type `String -> Maybe a`, which is one way to implement parser combinators (although  
164 usually the return type is `Maybe (a, String)` giving parsers the freedom to consume  
165 only a prefix of the input string and return the rest). The differences are that the result  
166 of our `Parser` type depends on the language specification and input string, and that a  
167 failure carries with it a proof that the string cannot be part of the language. This allows  
168 us to separate the specification of our language from the implementation while ensuring  
169 correctness.

170 *Remark 3.* Note that the `Dec` type only requires our parsers to produce a single result; it  
171 does not have to exhaustively list all possible ways to parse the input string. In Haskell,  
172 one might write `String -> [(a, String)]`, which allows a parser to return multiple  
173 results but still does not enforce exhaustiveness. Instead, we could use:

[JR] cite: monadic parser combinators

[JR] This should be explained in more detail

- 174 – completely unique account of enumeration.
- 175 – bijection with `Fin n` for some `n` or `Nat`.

176 In this paper, however, we use `Dec` to keep the presentation simple.

177 To construct a parser for our permissions language, we start by defining parsers for  
178 each of the language combinators. Let us start by considering the character combinator.  
179 If the given string is empty or has more than one character, it can never be in a language  
180 formed by one character. If the string does consist of only one character, then it is in  
181 the language if that character is the same as from the language specification. In Agda,  
182 we can write such a parser for characters as follows:

```

183 '-parse_ : (x : Char) → Parser (' x)
184 ('-parse _) [] = no λ ()
185 ('-parse x) (c :: []) = Dec.map (mk⇔ (λ { refl → refl }) (λ { refl → refl })) (c ≐? x)
186 ('-parse _) (_ :: _ :: _) = no λ ()

```

187 This is a correct implementation of a parser for languages that consist of a single  
188 character, but the implementation is hard to read and does not give much insight.  
189 Instead, we can factor this parser into two cases: the empty string case and the case  
190 where the string has at least one character. We call the former nullability and use  
191 the greek character  $\nu$  to signify it, and we call the latter derivative and use the greek  
192 character  $\delta$  to signify it. Figure 2 shows how these cases can be defined and how they  
193 relate to the basic combinators. These properties motivate the introduction of three new  
194 basic combinators: guards `_ · _`, the language consisting of only the empty string  $\epsilon$ , and  
195 the empty language  $\emptyset$ .

[JR] This does not motivate the split into  $\nu$  and  $\delta$  well enough. Also, the new combinators can be motivated more clearly.

196 Now the implementation of parsers for languages consisting of a single character  
197 follows completely from the decomposition into nullability and derivatives.

```

198 '-parse_ : (c' : Char) → Parser (' c')
199 ('-parse _) [] = Dec.map ν' ⊥-dec
200 ('-parse c') (c :: w) = Dec.map δ' (((c ≐? c') · -parse ε-parse) w)

```

201 The implementation of `· -parse`, `ε-parse`, and `∅-parse` are straightforward and can be  
202 found in our source code artifact.

[JR] todo: reference this nicely

```

203 _∪-parse_ : Parser P → Parser Q → Parser (P ∪ Q)
204 (φ ∪-parse ψ) [] = Dec.map ν∪ (ν φ ⊔-dec ν ψ)
205 (φ ∪-parse ψ) (c :: w) = Dec.map δ∪ ((δ c φ ∪-parse δ c ψ) w)

```

$$\begin{array}{ll}
\nu P = P \square & (\delta c P) w = P (c :: w) \\
\\
A \Leftrightarrow B = (A \rightarrow B) \times (B \rightarrow A) & P \Leftrightarrow Q = \forall \{w\} \rightarrow P w \Leftrightarrow Q w \\
\\
\begin{array}{llll}
\nu \emptyset : \perp & \Leftrightarrow \nu \emptyset & \delta \emptyset : \emptyset & \Leftrightarrow \delta c \emptyset \\
\nu \epsilon : \top & \Leftrightarrow \nu \epsilon & \delta \epsilon : \emptyset & \Leftrightarrow \delta c \epsilon \\
\nu \cdot : (A \times \nu P) & \Leftrightarrow \nu (A \cdot P) & \delta \cdot : (A \cdot \delta c P) & \Leftrightarrow \delta c (A \cdot P) \\
\nu ' : \perp & \Leftrightarrow \nu (' c') & \delta ' : ((c \equiv c') \cdot \epsilon) & \Leftrightarrow \delta c (' c') \\
\nu \cup : (\nu P \uplus \nu Q) & \Leftrightarrow \nu (P \cup Q) & \delta \cup : (\delta c P \cup \delta c Q) & \Leftrightarrow \delta c (P \cup Q) \\
\nu * : (\nu P \times \nu Q) & \Leftrightarrow \nu (P * Q) & \delta * : (\nu P \cdot \delta c Q \cup \delta c P * Q) & \Leftrightarrow \delta c (P * Q)
\end{array}
\end{array}$$

**Fig. 2.** Nullability, derivatives, and how they relate to the basic combinators.

```

206  _*-parse_ : Parser P → Parser Q → Parser (P * Q)
207  (φ *-parse ψ) [] = Dec.map ν* (ν φ ×-dec ν ψ)
208  (φ *-parse ψ) (c :: w) = Dec.map δ* ((ν φ ·-parse δ c ψ ∪-parse δ c φ *-parse ψ) w)

```

Using these combinators we can define a parser for the permissions language by simply mapping each of the language combinators onto their respective parser combinators.

```

211  permissions-parse = read-parse *-parse (write-parse *-parse execute-parse)
212  read-parse       = ('-parse '-'') ∪-parse ('-parse 'r')
213  write-parse      = ('-parse '-'') ∪-parse ('-parse 'w')
214  execute-parse    = ('-parse '-'') ∪-parse ('-parse 'x')

```

## 2.4 Infinite Languages

This permissions language is very simple. In particular, it is finite. In practice, many languages are infinite, for which the basic combinators will not suffice. For example, file paths can be arbitrarily long on most systems. Elliott [3] defines a Kleene star combinator which enables him to specify regular languages such as file paths.

However, we want to go one step further, specifying and parsing context-free languages. Most practical programming languages are at least context-free, if not more complicated. An essential feature of many languages is the ability to recognize balanced brackets. A minimal example language with balanced brackets is the following:

```

225  ⟨brackets⟩ ::= ε | '[' ⟨brackets⟩ ']' | ⟨brackets⟩ ⟨brackets⟩

```

This is the language of all strings which consist of balanced square brackets. Many practical programming languages include some form of balanced brackets. Furthermore, this language is well known to be context-free and not regular. Thus, we need more powerful combinators.

We could try to naively transcribe the brackets grammar using our basic combinators, but Agda will justifiably complain that it is not terminating (here we have added a `NON_TERMINATING` pragma to make Agda to accept it any way).

```

234  {-# NON_TERMINATING #-}
235  brackets = ε ∪ '[' * brackets * ']' ∪ brackets * brackets

```

236 We need to find a different way to encode this recursive relation.

```
237 postulate  $\mu : (\text{Lang} \rightarrow \text{Lang}) \rightarrow \text{Lang}$ 
238 brackets  $\mu = \mu (\lambda P \rightarrow \epsilon \cup ' [' * P * ' ' ]' \cup P * P)$ 
```

- 239 –  $\mu$ , with that exact type, cannot be implemented
- 240 – The  $\text{Lang} \rightarrow \text{Lang}$  function needs to be restricted

[JR] Can we give a concrete example of how  $\text{Lang} \rightarrow \text{Lang}$  is too general?

## 242 3 Context-free Languages

### 243 3.1 Fixed Points

[JR] Make it clear that we depart from Elliott's work at this point.

- 245 – If  $F : \text{Type} \rightarrow \text{Type}$  is a strictly positive functor, then we know its fixed point is well-defined.
- 246 – So we could restrict the argument of our fixed point combinator to only accept strictly positive functors.
- 247 – We are dealing with languages and not types directly, but luckily our definition of language is based on types and our basic combinators are strictly positive.
- 248 – One catch is that Agda currently has no way of directly expressing that a functor is strictly positive.<sup>4</sup>
- 249 – We can still make this evident to Agda by defining a data type of descriptions such as those used in the paper "gentle art of levitation".

[JR] todo: cite this

```
255 data Desc : Type1 where
256   [] : Desc
257   [ε] : Desc
258   [' c] : Char → Desc
259   [ _ ∪ _ ] : Desc → Desc → Desc
260   [ _ * _ ] : Desc → Desc → Desc
261   – We need Dec if we want to be able to write parsers
262   – but for specification it is not really needed
263   [ _ · _ ] : {A : Type} → Dec A → Desc → Desc
264   var : Desc
```

265 We can give semantics to our descriptions in terms of languages that we defined in the previous section.

[JR] todo: proper ref

```
267 [ ]o : Desc → ◇.Lang → ◇.Lang
268 [ [] ]o      = ◇.[]
269 [ [ε] ]o     = ◇.ε
270 [ [' c] ]o   = ◇.' c
271 [ [D1 ∪ D2] ]o = P = [ [D1] ]o P ◇.∪ [ [D2] ]o P
272 [ [D1 * D2] ]o = P = [ [D1] ]o P ◇.* [ [D2] ]o P
273 [ [ _ · _ ] {A} _ D ]o = P = A ◇.· [ [D] ]o P
274 [ [ var ] ]o      = P = P
```

275 Using these descriptions, we can define a fixed point as follows:

```
276 data [ ] (D : Desc) : ◇.Lang where
277   roll : [ D ]o [ D ] w → [ D ] w
```

<sup>4</sup> There is work on implementing positivity annotations.

```

278 unroll : [ D ] w → [ D ]o [ D ] w
279 unroll (roll x) = x

```

[JR] Brackets is one example, but can we characterise the whole class of languages we can define using these descriptions?

So we can finally define the brackets language.<sup>5</sup>

```

281 bracketsD = ε ∪ ' '[' * var * ' ' ∪ var * var
282 brackets = [ bracketsD ]

```

[JR] This modularity and nesting is not clear enough.

This representation is not modular, however. We cannot nest fixed points in descriptions. This problem comes up naturally when considering reduction, which we discuss next.

### 3.2 Reduction by Example

As we have seen with finite languages in Section 2, when writing parsers it is useful to consider how a language changes after one character has been parsed. We will call this *reduction*. For example, we could consider what happens to our brackets languages after one opening brackets has been parsed:  $\delta$  ' '[' brackets. In this section, we search for a description of this reduced language (the *reduct*).

We can mechanically derive this new language from the brackets definition by going over each of the disjuncts. The first disjunct,  $\epsilon$ , does not play a role because we know the string contains at least the opening bracket. The second disjunct, brackets surrounding a self-reference, is trickier. The opening bracket clearly matches, but it would be a mistake to say the new disjunct should be a self-reference followed by a closing bracket: `var * ' ' ]`.

Note that the self-reference in the new language would refer to the derivative of the old language, not to the old language itself. We would like to refer to the original bracket language: `brackets * ' ' ]`, but we cannot nest the brackets language into another description.

There are cases where we do want to use self-reference in the new language. Consider the third disjunct, `var * var`. It is a sequence so we expect from the finite case of Section 2 that matching one character results in two new disjuncts: one where the first sequent matches the empty string and the second is reduced and one where the first is reduced and the second is unchanged. In this case both sequents are self-references, so we need to know three things:

[JR] Why? That is what we saw in Section 2

1. Does the original language match the empty string?
2. What is the reduct of the language? (With reduct I mean the new language that results after one character is matched.)
3. What does it mean for the language to remain the same?

At first glance, the last point seems obvious, but remember that we are reducing the language, so self-references will change meaning even if they remain unchanged. Similarly to the previous disjunct, we want to refer to the original brackets in this case. To resolve this issue of referring to the original brackets expression, we introduce a new combinator  $\mu$ , which has the meaning of locally taking a fixed point of a subexpression.

```

317 data Desc : Type1 where
318   - ...
319   μ : Desc → Desc
320

```

<sup>5</sup> We split this definition into two because we want to separately reuse the description later.



```

321   $\llbracket \_ \rrbracket_o : \text{Desc} \rightarrow \diamond.\text{Lang} \rightarrow \diamond.\text{Lang}$ 
322  - ...
323   $\llbracket \mu D \rrbracket_o \_ = \llbracket D \rrbracket$ 
324

```

[JR] How is this used in our example?

325 The first question is easy to answer: yes, the first disjunct of brackets is epsilon which  
326 matches the empty string.

```

327   $\nu\text{brackets} : \text{Dec } (\diamond.\nu \text{ brackets})$ 
328   $\nu\text{brackets} = \text{yes } (\text{roll } (\text{inj}_1 \text{ refl}))$ 

```

329 The second question is where having a self-reference in the new language is useful.  
330 We can refer to the reduct of brackets by using self-reference.

331 This enables us to write the reduct of brackets with respect to the opening bracket.

```

332   $\text{bracketsD}' = \mu \text{ bracketsD} * ']' \cup \nu\text{brackets} \cdot \text{var} \cup \text{var} * \mu \text{ bracketsD}$ 
333   $\text{brackets}' = \llbracket \text{bracketsD}' \rrbracket$ 

```

334 Conclusion:

- 335 - We can reuse many of the results of finite languages (Section 2).
- 336 - We need a new  $\mu$  combinator to nest fixed points in descriptions. This is necessary  
337 to refer back to the original language before reduction.
- 338 - Reducing a self-reference simply results in a self-reference again, because self-references  
339 in the reduct refer to the reduct.

340 Again, we do not want to have to do this reduction manually. Instead, we show how to  
341 do it in general for any description in the next section.

### 342 3.3 Parsing in General

343 Our goal is to define:

```

344   $\text{parse} : \forall D \rightarrow \diamond.\text{Parser } \llbracket D \rrbracket$ 

```

345 We approach this by decomposing parsing into  $\nu$  and  $\delta$ .

```

346   $\nu D : \forall D \rightarrow \text{Dec } (\diamond.\nu \llbracket D \rrbracket)$ 
347   $\delta D : \text{Char} \rightarrow \text{Desc} \rightarrow \text{Desc}$ 

```

348 The  $\nu D$  function can easily be written to be correct by construction, however  $\delta D$   
349 must be proven correct separately as follows:

```

350   $\delta D\text{-correct} : \llbracket \delta D \ c \ D \rrbracket \diamond.\iff \diamond.\delta \ c \llbracket D \rrbracket$ 

```

351 The actual parsing can now be done character by character:

```

352   $\text{parse } D \llbracket \rrbracket = \nu D \ D$ 
353   $\text{parse } D \ (c :: w) = \text{Dec.map } \delta D\text{-correct} \ (\text{parse } (\delta D \ c \ D) \ w)$ 

```

354 That is the main result of this paper. The remainder of the paper concerns the  
355 implementation of  $\nu D$ ,  $\delta D$ ,  $\delta D\text{-correct}$ .

### 3.4 Nullability

If we know the nullability of a language,  $P$ , then the nullability of a description functor applied to  $P$  is the same as the empty string parsers for our finite languages, but with the nullability of the variables given by the nullability of  $P$ . For the  $\mu$  case we use the nullability of the fixed point, which we will implement shortly.

[JR] Reiter  
are the same

```

361  $\nu_o : \text{Dec } (\diamond.\nu P) \rightarrow \forall D \rightarrow \text{Dec } (\diamond.\nu (\llbracket D \rrbracket_o P))$ 
362  $\nu_o \_ \emptyset = \text{no } \lambda ()$ 
363  $\nu_o \_ \epsilon = \text{yes refl}$ 
364  $\nu_o \_ (' c) = \text{no } \lambda ()$ 
365  $\nu_o z (D \cup D_1) = \nu_o z D \uplus\text{-dec } \nu_o z D_1$ 
366  $\nu_o z (D * D_1) = \text{Dec.map } \diamond.\nu * (\nu_o z D \times\text{-dec } \nu_o z D_1)$ 
367  $\nu_o z (x \cdot D) = x \times\text{-dec } \nu_o z D$ 
368  $\nu_o z \text{var} = z$ 
369  $\nu_o \_ (\mu D) = \nu D D$ 

```

- Naively we might try  $\nu D D = \nu_o (\nu D D) D$
- But that obviously will not terminate (consider the language  $\llbracket \text{var} \rrbracket$ ).
- Instead we use Lemma 1

**Lemma 1.** *The nullability of a fixed point is determined completely by a single application of the underlying functor to the empty language.*

$$\nu D \emptyset \Leftrightarrow \nu D : \diamond.\nu (\llbracket D \rrbracket_o \diamond.\emptyset) \Leftrightarrow \diamond.\nu (\llbracket D \rrbracket)$$

*Proof.* The forward direction is easily proven by noting that nullability and the semantics of a description are functors and that the empty language is initial. It is also straightforward to write the proof directly.

$$\nu D \emptyset \rightarrow \nu D : \forall D \rightarrow \diamond.\nu (\llbracket D \rrbracket_o \diamond.\emptyset) \rightarrow \diamond.\nu (\llbracket D \rrbracket_o \llbracket D_0 \rrbracket)$$

The backwards direction is more difficult. We prove a more general lemma from which our desired result follows. The generalized lemma states that, if the application of a descriptor functor to a fixed point of another descriptor is nullable, then either the fixed point plays no role and the descriptor functor is also nullable if applied to the empty language, or the other descriptor (that we took the fixed point of) is nullable when applied to the empty language.

$$\nu D \emptyset \leftarrow \nu D : \forall D \rightarrow \diamond.\nu (\llbracket D \rrbracket_o \llbracket D_0 \rrbracket) \rightarrow \diamond.\nu (\llbracket D \rrbracket_o \diamond.\emptyset) \uplus \diamond.\nu (\llbracket D_0 \rrbracket_o \diamond.\emptyset)$$

If we choose  $D_0 = D$  then both cases of the resulting disjoint union have the same type, so we can just pick whichever of the two we get as a result using the **reduce** :  $A \uplus A \rightarrow A$  function. Modulo wrapping and unwrapping of the fixed point (using the **roll** constructor), we now have the two functions which prove the lemma:

$$\nu D \emptyset \Leftrightarrow \nu D \{D\} = \text{mk} \Leftrightarrow (\text{roll} \circ \nu D \emptyset \rightarrow \nu D D) (\text{reduce} \circ \nu D \emptyset \leftarrow \nu D \{D_0 = D\} D \circ \text{unroll})$$

Using Lemma 1, we can easily define nullability for our description functors.

$$\nu D = \text{Dec.map } \nu D \emptyset \Leftrightarrow \nu D \circ \nu_o (\text{no } \lambda ())$$

394 *Remark 4.* Lemma 1 does not define an isomorphism on types. In particular, the back-  
 395 wards direction is not injective. Consider the brackets language. It has the following null  
 396 element, where we first choose the third disjunct,  $\text{var} * \text{var}$ , and then the first disjunct  
 397  $\epsilon$  for both branches.

```
398 brackets0 :  $\diamond.\nu$  brackets
399 brackets0 = roll (inj2 (inj2 ([], [], refl, roll (inj1 refl), roll (inj1 refl))))
```

400 When we round-trip this through our lemma, we get a different result:

```
401 brackets0' :  $\nu D \emptyset \Leftrightarrow \nu D \{ \text{bracketsD} \}$  .to ( $\nu D \emptyset \Leftrightarrow \nu D \{ \text{bracketsD} \}$  .from brackets0)
402            $\equiv$  roll (inj1 refl)
403 brackets0' = refl
```

404 It now directly takes the first disjunct,  $\epsilon$ .

405 In practice, such problems should be avoided by using unambiguous languages, ensu-  
 406 ring that there is only one valid parse result for each string.

[JR] todo: give recommendations for future work, for example to use data-dependent grammars.

### 407 3.5 Reduction

408 The final piece of the puzzle is reduction. This tells us how the language descriptions  
 409 change after parsing each input character.

410 In Section 3.2, we established that the meaning of self-references changes and thus  
 411 they need to be replaced by local fixed points of the original language. We define a  
 412 function  $\sigma D$  to perform this substitution. It is a simple recursive function which replaces  
 413 the  $\text{var}$  constructor with a given  $D'$  description.

```
414  $\sigma : \text{Desc} \rightarrow \text{Desc} \rightarrow \text{Desc}$ 
415  $\sigma \emptyset \quad D' = \emptyset$ 
416  $\sigma \epsilon \quad D' = \epsilon$ 
417  $\sigma (' c) \quad D' = ' c$ 
418  $\sigma (D \cup D_1) \quad D' = \sigma D D' \cup \sigma D_1 D'$ 
419  $\sigma (D * D_1) \quad D' = \sigma D D' * \sigma D_1 D'$ 
420  $\sigma (x \cdot D) \quad D' = x \cdot \sigma D D'$ 
421  $\sigma \text{var} \quad D' = D'$ 
422  $\sigma (\mu D) \quad D' = \mu D$ 
```

423 It turns out that the only the sequencing case,  $*$  leaves the variables untouched, thus  
 424 we only need to apply the substitution there. This substitution does mean we need to  
 425 keep track of the original description,  $D_0$ , through the recursion. Most other cases follow  
 426 the structure we uncovered in Figure 2. For the self-reference case,  $\text{var}$ , we produce a  
 427 self-reference again, which works because it now refers to the reduct. Finally, for the  
 428 internal fixed point,  $\mu$ , we can simply recursively call the reduction function. Thus, our  
 429 reduction helper function is defined as follows:

```
430  $\delta_o : \text{Desc} \rightarrow \text{Char} \rightarrow \text{Desc} \rightarrow \text{Desc}$ 
431  $\delta_o D_0 c \emptyset = \emptyset$ 
432  $\delta_o D_0 c \epsilon = \emptyset$ 
433  $\delta_o D_0 c (' c') = (c \stackrel{?}{=} c') \cdot \epsilon$ 
434  $\delta_o D_0 c (D \cup D_1) = \delta_o D_0 c D \cup \delta_o D_0 c D_1$ 
435  $\delta_o D_0 c (D * D_1) = \nu_o (\nu D D_0) D \cdot \delta_o D_0 c D_1 \cup \delta_o D_0 c D * \sigma D_1 (\mu D_0)$ 
436  $\delta_o D_0 c (x \cdot D) = x \cdot \delta_o D_0 c D$ 
437  $\delta_o D_0 c \text{var} = \text{var}$ 
438  $\delta_o D_0 c (\mu D) = \mu (\delta D c D)$ 
```

439 At the top level, we simply delegate to the helper by passing  $D_0 = D$ .

$$440 \quad \delta D \ c \ D = \delta_o \ D \ c \ D$$

441 **Lemma 2.** *Substitution of a local fixed point into a description is the same as applying*  
 442 *the corresponding functor to the semantic fixed point.*

$$443 \quad \sigma\mu : \forall \ D \rightarrow \llbracket \sigma \ D \ (\mu \ D_0) \rrbracket_o \ P \ w \equiv \llbracket D \rrbracket_o \llbracket D_0 \rrbracket \ w$$

444 The proof follows directly by induction and computation.

$$\begin{aligned} 445 \quad & \delta D\text{-to} : \forall \ D \rightarrow \llbracket \delta_o \ D_0 \ c \ D \rrbracket_o \llbracket \delta D \ c \ D_0 \rrbracket \ w \rightarrow \diamond.\delta \ c \ (\llbracket D \rrbracket_o \llbracket D_0 \rrbracket) \ w \\ 446 \quad & \delta D\text{-to} \ (' \ c') \ (\text{refl} \ , \ \text{refl}) = \text{refl} \\ 447 \quad & \delta D\text{-to} \ (D \cup D_1) \ (\text{inj}_1 \ x) = \text{inj}_1 \ (\delta D\text{-to} \ D \ x) \\ 448 \quad & \delta D\text{-to} \ (D \cup D_1) \ (\text{inj}_2 \ y) = \text{inj}_2 \ (\delta D\text{-to} \ D_1 \ y) \\ 449 \quad & \delta D\text{-to} \ (D * D_1) \ (\text{inj}_1 \ (x, y)) = \llbracket \ , \ \_ , \ \text{refl} \ , \ x , \ \delta D\text{-to} \ D_1 \ y \rrbracket \\ 450 \quad & \delta D\text{-to} \ (D * D_1) \ (\text{inj}_2 \ (\_ :: \_ , \_ , \text{refl} \ , \ x, y)) = \_ :: \_ , \_ , \text{refl} \ , \ \delta D\text{-to} \ D \ x , \ \text{subst id} \ (\sigma\mu \ D_1) \ y \\ 451 \quad & \delta D\text{-to} \ (A \cdot D) \ (x, y) = x , \ \delta D\text{-to} \ D \ y \\ 452 \quad & \delta D\text{-to} \ \{D_0 = D\} \ \text{var} \ (\text{roll} \ x) = \text{roll} \ (\delta D\text{-to} \ D \ x) \\ 453 \quad & \delta D\text{-to} \ (\mu \ D) \ (\text{roll} \ x) = \text{roll} \ (\delta D\text{-to} \ D \ x) \end{aligned}$$

$$\begin{aligned} 454 \quad & \delta D\text{-from} : \forall \ D \rightarrow \diamond.\delta \ c \ (\llbracket D \rrbracket_o \llbracket D_0 \rrbracket) \ w \rightarrow \llbracket \delta_o \ D_0 \ c \ D \rrbracket_o \llbracket \delta D \ c \ D_0 \rrbracket \ w \\ 455 \quad & \delta D\text{-from} \ (' \ c') \ \text{refl} = \text{refl} \ , \ \text{refl} \\ 456 \quad & \delta D\text{-from} \ (D \cup D_1) \ (\text{inj}_1 \ x) = \text{inj}_1 \ (\delta D\text{-from} \ D \ x) \\ 457 \quad & \delta D\text{-from} \ (D \cup D_1) \ (\text{inj}_2 \ y) = \text{inj}_2 \ (\delta D\text{-from} \ D_1 \ y) \\ 458 \quad & \delta D\text{-from} \ (D * D_1) \ (\llbracket \ , \ \_ , \ \text{refl} \ , \ x , \ y \rrbracket) = \text{inj}_1 \ (x , \ \delta D\text{-from} \ D_1 \ y) \\ 459 \quad & \delta D\text{-from} \ (D * D_1) \ (\_ :: \_ , \_ , \text{refl} \ , \ x, y) = \text{inj}_2 \ (\_ , \_ , \text{refl} \ , \ \delta D\text{-from} \ D \ x , \ \text{subst id} \ (\text{sym} \ (\sigma\mu \ D_1))) \\ 460 \quad & \delta D\text{-from} \ (A \cdot D) \ (x, y) = x , \ \delta D\text{-from} \ D \ y \\ 461 \quad & \delta D\text{-from} \ \{D_0 = D\} \ \text{var} \ (\text{roll} \ x) = \text{roll} \ (\delta D\text{-from} \ D \ x) \\ 462 \quad & \delta D\text{-from} \ (\mu \ D) \ (\text{roll} \ x) = \text{roll} \ (\delta D\text{-from} \ D \ x) \end{aligned}$$

$$463 \quad \delta D\text{-correct} \ \{D = D\} = \text{mk} \Leftrightarrow (\text{roll} \circ \delta D\text{-to} \ D \circ \text{unroll}) \ (\text{roll} \circ \delta D\text{-from} \ D \circ \text{unroll})$$

## 464 4 Discussion

465 Finally, we want to discuss three aspects of our work: expressiveness, performance, and  
 466 simplicity.

[JR] TODO:  $\mu$ -regular expressions have been studied before, cite

468 *Expressiveness* We conjecture that our grammars which include variables and fixed  
 469 points can describe any context-free language. .

[JR] mention that we only support context-free languages without mutual recursion and how we use a subset of  $\mu$ -regular languages

470 Going beyond context-free languages, many practical programming languages cannot  
 471 be adequately described as context-free languages. For example, features such as  
 472 associativity, precedence, and indentation sensitivity cannot be expressed directly using  
 473 context-free grammars. Recent work by Afrozeh and Izmaylova [1] shows that all these  
 474 advanced features can be supported if we extend our grammars with data-dependencies.  
 475 Our framework can form a foundation for such extensions and we consider formalizing  
 476 it as future work.

*Performance* For a parser to be practically useful, it must at least have linear asymptotic complexity for practical grammars. Might et al. [4] show that naively parsing using derivatives does not achieve that bound, but optimizations might make it possible. In particular, they argue that we could achieve  $O(n|G|)$  time complexity (where  $|G|$  is the grammar size) if the grammar size stays approximately constant after every derivative. By compacting the grammar, they conjecture it is possible to achieve this bound for any unambiguous grammar. We want to investigate if similar optimizations could be applied to our parser and if we can prove that we achieve this bound.

*Simplicity* One of the main contributions of Elliott’s type theoretic formalization of languages [3] is its simplicity of implementation and proof. To be able to extend his approach to context-free languages we have had to introduce some complications.

[JR] TODO: finish this paragraph

## 5 Related Work

- Jeremy Yallop - performance
- Peter Thiemann - derivatives of  $\mu$ -regular expressions. This is the closest to our work, we have a mechanized proof and use type theory instead of set theory.
- Guillaume Allais’ Agdarsec
- Danielsson’s coinductive parser combinators
- ”Certified Parsing of Dependent Regular Grammars” John Sarracino; Gang Tan; Greg Morrisett
- Brink et al. (MPC 2010), they formalize the left-corner transformation
- Jean-Philippe Bernardy and Patrik Jansson, ”Certified Context-Free Parsing: A formalisation of Valiant’s Algorithm in Agda.” This is a formalization of a performant matrix-based parsing algorithm.
- Conal Elliott, of course
- Introduction to Automata Theory, Languages, and Computation - Hopcroft, Motswani, Ullman

## 6 Conclusion

In conclusion, we have formalized (acyclic) context-free grammars using a type theoretic approach to provide fertile ground for further formalizations of disambiguation strategies and parsers that are both correct and performant.

## References

1. Afrozeh, A., Izmaylova, A.: One parser to rule them all. In: 2015 ACM International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward!). p. 151–170. Onward! 2015, Association for Computing Machinery, New York, NY, USA (2015). <https://doi.org/10.1145/2814228.2814242>
2. Brzozowski, J.A.: Derivatives of regular expressions. J. ACM **11**(4), 481–494 (Oct 1964). <https://doi.org/10.1145/321239.321249>
3. Elliott, C.: Symbolic and automatic differentiation of languages. Proc. ACM Program. Lang. **5**(ICFP) (Aug 2021). <https://doi.org/10.1145/3473583>
4. Might, M., Darais, D., Spiewak, D.: Parsing with derivatives: a functional pearl. In: Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming. p. 189–195. ICFP ’11, Association for Computing Machinery, New York, NY, USA (2011). <https://doi.org/10.1145/2034773.2034801>