

Context-free Languages, Type Theoretically

Jaro Reinders¹[0000-0002-6837-3757] and Casper Bach²[0000-0003-0622-7639]

¹ Delft University of Technology, Delft, The Netherlands

² 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 is the conversion of flat, human-readable text into a tree structure that is easier for computers to manipulate. As one of the central pillars of compiler tooling since the 1960s, today almost every automated transformation of computer programs requires a form of parsing. Though it is a mature research subject, it is still actively studied, for example the question of how to resolve ambiguities in context-free grammars [1].

Most parsing works mix the essence of the parsing technique with operational details. Our understanding and ability to improve upon these parsing techniques is hindered by the additional complexity of these inessential practical concerns. To address this issue, we are developing natural denotational semantics for traditional parsing techniques.

Recent work by Elliott uses interactive theorem provers to state simple specifications of languages and that proofs of desirable properties of these language specifications transfer easily to their parsers [3]. Unfortunately, this work only considers regular languages which are not powerful enough to describe practical programming languages.

In this paper, we formalize context-free languages and show how to parse them, extending Elliott's type theoretic approach to language specification. One of the main challenges is that the recursive nature of context-free languages does not map directly onto interactive theorem provers as they do not support general recursion (for good reasons). We encode context-free languages as fixed points of functors (initial algebras).

We make the following concrete contributions:

- We extend Elliott's type theoretic formalization of regular languages to context-free languages.

For this paper we have chosen Agda as our type theory and interactive theorem prover. We believe our definitions should transfer easily to other theories and tools. This paper itself is a literate Agda file; all highlighted Agda code has been accepted by Agda's type checker, giving us a high confidence of correctness.

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.

[JR] such as... state machines, continuations, memoization?

[JR] Elliott has kicked off this effort...

[JR] Make the problem clear through an example: if we have a left-recursive grammar then naively unfolding it gets us into an infinite loop.

[JR] Say something about the limitation that we only study acyclic grammars: there must be a total order on nonterminals and a nonterminal is not allowed to refer to nonterminals that come before it. We wanted to start by limiting ourselves to grammars with only one nonterminal, but those are not closed under derivations.

2.1 Languages

We define languages as being functions from strings to types.³

```
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.

[JR] do I need to give an example?

- 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.

2.2 Basic Language Combinators

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

[JR] cite: BNF

```
⟨permissions⟩ ::= ⟨read⟩ ⟨write⟩ ⟨execute⟩
⟨read⟩        ::= '-' | 'r'
⟨write⟩       ::= '-' | 'w'
⟨execute⟩     ::= '-' | 'x'
```

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

```
permissions = read * write * execute
read        = ' '-' ∪ ' 'r'
write       = ' '-' ∪ ' 'w'
execute     = ' '-' ∪ ' 'x'
```

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

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

Fig. 1. Basic language combinators.

2.3 Parsers

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

```
data Dec (A : Type) : Type where
  yes : A → Dec A
  no  : ¬ A → Dec A
```

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

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

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

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

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

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

To construct a parser for our permissions language, we start by defining parsers for each of the language combinators. Let us start by considering the character combinator. If the given string is empty or has more than one character, it can never be in a language

[JR] cite: monadic parser combinators

[JR] This should be explained in more detail

124 formed by one character. If the string does consist of only one character, then it is in
 125 the language if that character is the same as from the language specification. In Agda,
 126 we can write such a parser for characters as follows:

```

127 'parse_ : (x : Char) → Parser (' x)
128 ('parse _) [] = no λ ()
129 ('parse x) (c :: []) = Dec.map (mk⇔ (λ { refl → refl }) (λ { refl → refl }))) (c ? x)
130 ('parse _) (_ :: _ :: _) = no λ ()

```

131 This is a correct implementation of a parser for languages that consist of a single
 132 character, but the implementation is hard to read and does not give much insight.
 133 Instead, we can factor this parser into two cases: the empty string case and the case
 134 where the string has at least one character. We call the former nullability and use
 135 the greek character ν to signify it, and we call the latter derivative and use the greek
 136 character δ to signify it. Figure 2 shows how these cases can be defined and how they
 137 relate to the basic combinators. These properties motivate the introduction of three new
 138 basic combinators: guards $_ \cdot _$, the language consisting of only the empty string ϵ , and
 139 the empty language \emptyset .

[JR] This does not motivate the split into ν and δ well enough. Also, the new combinators can be motivated more clearly.

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

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

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

```

142 'parse_ : (c' : Char) → Parser (' c')
143 ('parse _) [] = Dec.map \nu' \_ -dec
144 ('parse c') (c :: w) = Dec.map \delta' (((c ? c') \cdot -parse \epsilon -parse) w)

```

145 The implementation of \cdot -parse, ϵ -parse, and \emptyset -parse are straightforward and can be
 146 found in our source code artifact.

[JR] todo: reference this nicely

```

147 \_ \cup -parse_ : Parser P → Parser Q → Parser (P \cup Q)
148 (\phi \cup -parse \psi) [] = Dec.map \nu \cup (\nu \phi \uplus -dec \nu \psi)
149 (\phi \cup -parse \psi) (c :: w) = Dec.map \delta \cup ((\delta c \phi \cup -parse \delta c \psi) w)
150 \_ * -parse_ : Parser P → Parser Q → Parser (P * Q)
151 (\phi * -parse \psi) [] = Dec.map \nu * (\nu \phi \times -dec \nu \psi)
152 (\phi * -parse \psi) (c :: w) = Dec.map \delta * ((\nu \phi \cdot -parse \delta c \psi \cup -parse \delta c \phi * -parse \psi) w)

```

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

```
155 permissions-parse = read-parse *-parse (write-parse *-parse execute-parse)
156 read-parse       = ('-parse '-' ) ∪ -parse ('-parse 'r' )
157 write-parse      = ('-parse '-' ) ∪ -parse ('-parse 'w' )
158 execute-parse    = ('-parse '-' ) ∪ -parse ('-parse 'x' )
```

159 2.4 Infinite Languages

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

[JR] does this need citation?

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

```
168 ⟨brackets⟩ ::= ε | '[' ⟨brackets⟩ ']' | ⟨brackets⟩ ⟨brackets⟩
```

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

[JR] todo: flesh out this outline

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

```
178 {-# NON_TERMINATING #-}
179 brackets = ε ∪ '[' * brackets * ']' ∪ brackets * brackets
```

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

```
181 postulate μ : (Lang → Lang) → Lang
182 brackets_μ = μ (λ P → ε ∪ '[' * P * ']' ∪ P * P)
```

- 183 – μ , with that exact type, cannot be implemented
- 184 – 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?

186 3 Context-free Languages

187 3.1 Fixed Points

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

- 189 – If $F : \text{Type} \rightarrow \text{Type}$ is a strictly positive functor, then we know its fixed point is
 190 well-defined.
- 191 – So we could restrict the argument of our fixed point combinator to only accept
 192 strictly positive functors.
- 193 – We are dealing with languages and not types directly, but luckily our definition of
 194 language is based on types and our basic combinators are strictly positive.

- One catch is that Agda currently has no way of directly expressing that a functor is strictly positive.⁴
- 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:

```

199 data Desc : Type1 where
200   [] : Desc
201   ε : Desc
202   ' _ : Char → Desc
203   _ ∪ _ : Desc → Desc → Desc
204   _ * _ : Desc → Desc → Desc
205   – We need Dec if we want to be able to write parsers
206   – but for specification it is not really needed
207   _ · _ : {A : Type} → Dec A → Desc → Desc
208   var : Desc

```

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

[JR] todo: proper ref

```

211 [ ]o : Desc → ◇.Lang → ◇.Lang
212 [ [] ]o = const ◇.[]
213 [ ε ]o = const ◇.ε
214 [ ' c ]o = const (◇.' c)
215 [ D1 ∪ D2 ]o P = [ D1 ]o P ◇ ∪ [ D2 ]o P
216 [ D1 * D2 ]o P = [ D1 ]o P ◇ * [ D2 ]o P
217 [ _ · {A} _ D ]o P = A ◇ · [ D ]o P
218 [ var ]o P = P

```

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

```

220 data [ ] (D : Desc) : ◇.Lang where
221   step : [ D ]o [ D ] w → [ D ] w

```

So we can finally define the brackets language.⁵

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

```

223 bracketsD = ε ∪ ' '[' * var * ' ' ∪ var * var
224 brackets = [ bracketsD ]

```

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.

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

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: δ '[' **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

238 would be a mistake to say the new disjunct should be a self-reference followed by a
239 closing bracket.

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

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

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

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

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

```
259 data Desc : Type1 where
260   - ...
261    $\mu$  : Desc → Desc
262
263    $\llbracket \_ \rrbracket_o$  : Desc →  $\diamond$ .Lang →  $\diamond$ .Lang
264   - ...
265    $\llbracket \mu D \rrbracket_o \_ = \llbracket D \rrbracket$ 
```

[JR] How is this used in our example?

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

```
269  $\nu$ brackets : Dec ( $\diamond$ . $\nu$  brackets)
270  $\nu$ brackets = yes (step (inj1 refl))
```

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

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

```
274 bracketsD' =  $\mu$  bracketsD * ' ' ]  $\cup$   $\nu$ brackets · var  $\cup$  var *  $\mu$  bracketsD
275 brackets'  =  $\llbracket$  bracketsD'  $\rrbracket$ 
```

276 Conclusion:

- 277 - We can reuse many of the results of finite languages (Section 2).
- 278 - We need a new μ combinator to nest fixed points in descriptions. This is necessary
279 to refer back to the original language before reduction.
- 280 - Reducing a self-reference simply results in a self-reference again, because self-references
281 in the reduct refer to the reduct.

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

⁴ There is work on implementing positivity annotations.

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

284 3.3 Parsing in General

285 Our goal is to define:

286 $\text{parse} : \forall D \rightarrow \diamond.\text{Parser } \llbracket D \rrbracket$

287 We approach this by decomposing parsing into ν and δ .

288 $\nu D : \forall D \rightarrow \text{Dec } (\diamond.\nu \llbracket D \rrbracket)$

289 $\delta D : \text{Char} \rightarrow \text{Desc} \rightarrow \text{Desc}$

290 The νD function can easily be written to be correct by construction, however δD
291 must be proven correct separately as follows:

292 $\delta D\text{-correct} : \llbracket \delta D \ c \ D \rrbracket \diamond.\iff \diamond.\delta \ c \llbracket D \rrbracket$

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

294 $\text{parse } D \llbracket \rrbracket = \nu D \ D$

295 $\text{parse } D \ (c :: w) = \text{Dec.map } \delta D\text{-correct} \ (\text{parse } (\delta D \ c \ D) \ w)$

296 That is the main result of this paper. The remainder of the paper concerns the
297 implementation of νD , δD , $\delta D\text{-correct}$.

298 3.4 Implementation and Proof

299 **Lemma 1.** *The nullability of a fixed point is determined completely by a single appli-*
300 *cation of the underlying functor to the empty language.*

301 $\nu D\emptyset \iff \nu D : \diamond.\nu (\llbracket D \rrbracket_o \diamond.\emptyset) \iff \diamond.\nu \llbracket D \rrbracket$

302 $\nu D\emptyset \iff \nu D = \{\text{!} \ \text{!}\}$

303 $\nu D\emptyset : \forall D \rightarrow \text{Dec } (\diamond.\nu (\llbracket D \rrbracket_o \diamond.\emptyset))$

304 $\nu D\emptyset \ \emptyset = \text{no } \lambda \ ()$

305 $\nu D\emptyset \ \epsilon = \text{yes refl}$

306 $\nu D\emptyset \ (' \ x) = \text{no } \lambda \ ()$

307 $\nu D\emptyset \ (D \cup D_1) = \nu D\emptyset \ D \uplus\text{-dec } \nu D\emptyset \ D_1$

308 $\nu D\emptyset \ (D * D_1) = \text{Dec.map } \diamond.\nu * (\nu D\emptyset \ D \times\text{-dec } \nu D\emptyset \ D_1)$

309 $\nu D\emptyset \ (x \cdot D) = x \times\text{-dec } \nu D\emptyset \ D$

310 $\nu D\emptyset \ \text{var} = \text{no } \lambda \ ()$

311 $\nu D\emptyset \ (\mu \ D) = \text{Dec.map } \nu D\emptyset \iff \nu D \ (\nu D\emptyset \ D)$

312 $\nu D \ D = \text{Dec.map } \nu D\emptyset \iff \nu D \ (\nu D\emptyset \ D)$

313 $\sigma D : \text{Desc} \rightarrow \text{Desc} \rightarrow \text{Desc}$

314 $\sigma D \ \emptyset \quad D' = \emptyset$

315 $\sigma D \ \epsilon \quad D' = \epsilon$

316 $\sigma D \ (' \ c) \quad D' = ' \ c$

317 $\sigma D \ (D \cup D_1) \ D' = \sigma D \ D \ D' \cup \sigma D \ D_1 \ D'$

318 $\sigma D \ (D * D_1) \ D' = \sigma D \ D \ D' * \sigma D \ D_1 \ D'$

319 $\sigma D \ (x \cdot D) \ D' = x \cdot \sigma D \ D \ D'$

$\sigma D \text{ var} \quad D' = D'$
 $\sigma D (\mu D) \quad D' = \mu D$

$\delta D_o : \text{Desc} \rightarrow \text{Char} \rightarrow \text{Desc} \rightarrow \text{Desc}$
 $\delta D_o D_0 c \emptyset = \emptyset$
 $\delta D_o D_0 c \epsilon = \emptyset$
 $\delta D_o D_0 c (' c') = (c \stackrel{?}{=} c') \cdot \epsilon$
 $\delta D_o D_0 c (D \cup D_1) = \delta D_o D_0 c D \cup \delta D_o D_0 c D_1$
 $\delta D_o D_0 c (D * D_1) = \nu D D \cdot \delta D_o D_0 c D_1 \cup \delta D_o D_0 c D * \sigma D D_1 D_0$
 $\delta D_o D_0 c (x \cdot D) = x \cdot \delta D_o D_0 c D$
 $\delta D_o D_0 c \text{ var} = \text{var}$
 $\delta D_o D_0 c (\mu D) = \mu (\delta D_o D c D)$

$\delta D c D = \delta D_o D c D$

$\delta D\text{-correct} = \{! \ !\}$

4 Discussion

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

[JR] TODO: μ -regular expressions have been studied before, cite

Expressiveness We conjecture that our grammars which include variables and fixed points can describe any context-free language. We have shown the example of balanced the bracket language which is known to be context-free. Furthermore, Grenrus shows that any context-free grammar can be converted to his grammars [4], which are similar to our grammars. The main problem is showing that mutually recursive nonterminals can be expressed using our simple fixed points, which requires Bekić's bisection lemma [2]. Formalizing this in our framework is future work.

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

Performance For a parser to practically useful, it must at least have linear asymptotic complexity for practical grammars. Might et al. [5] 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.

[JR] cite Jeremy Yallop's work

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

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.

365 **References**

- 366 1. Afroozeh, A., Izmaylova, A.: One parser to rule them all. In: 2015 ACM International
367 Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software
368 (Onward!). p. 151–170. Onward! 2015, Association for Computing Machinery, New York,
369 NY, USA (2015). <https://doi.org/10.1145/2814228.2814242>
- 370 2. Bekić, H.: Definable operations in general algebras, and the theory of automata and
371 flowcharts, pp. 30–55. Springer Berlin Heidelberg, Berlin, Heidelberg (1984). [https://doi.](https://doi.org/10.1007/BFb0048939)
372 [org/10.1007/BFb0048939](https://doi.org/10.1007/BFb0048939), <https://doi.org/10.1007/BFb0048939>
- 373 3. Elliott, C.: Symbolic and automatic differentiation of languages. *Proc. ACM Program. Lang.*
374 **5**(ICFP) (Aug 2021). <https://doi.org/10.1145/3473583>
- 375 4. Grenrus, O.: Fix-ing regular expressions (2020), [https://well-typed.com/blog/2020/06/](https://well-typed.com/blog/2020/06/fix-ing-regular-expressions/)
376 [fix-ing-regular-expressions/](https://well-typed.com/blog/2020/06/fix-ing-regular-expressions/), accessed: 2024-12-12
- 377 5. Might, M., Darais, D., Spiewak, D.: Parsing with derivatives: a functional pearl. In: Pro-
378 ceedings of the 16th ACM SIGPLAN International Conference on Functional Programming.
379 p. 189–195. ICFP '11, Association for Computing Machinery, New York, NY, USA (2011).
380 <https://doi.org/10.1145/2034773.2034801>