

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 [2]. 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 [2]. 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 derivatives.

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.

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

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)

```

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.

```
permissions-parse = read-parse *-parse (write-parse *-parse execute-parse)
read-parse       = ('-parse '-' ) ∪ -parse ('-parse 'r' )
write-parse      = ('-parse '-' ) ∪ -parse ('-parse 'w' )
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 [2] defines a Kleene star combinator which enables him to specify regular languages such as file paths.

[JR] does this need citation?

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:

```
<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.

[JR] todo: flesh out this outline

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).

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

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

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

- μ , with that exact type, cannot be implemented
- 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?

3 Context-free Languages

3.1 Fixed Points

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

- If $F : \text{Type} \rightarrow \text{Type}$ is a strictly positive functor, then we know its fixed point is well-defined.
- So we could restrict the argument of our fixed point combinator to only accept strictly positive functors.
- 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.

- 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 Desc if we want to be able to write parsers
206   – but for specifiction 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 = ◇.[]
213 [ ε ]o = ◇.ε
214 [ ' c ]o = ◇.' 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   roll : [ D ]o [ D ] w → [ D ] w

```

```

222 unroll : [ D ] w → [ D ]o [ D ] w
223 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.⁵

```

225 bracketsD = ε ∪ ' ' [ ' * var * ' ' ] ∪ var * var
226 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

⁴ There is work on implementing positivity annotations.

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

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, ϵ , 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 μ , which has the meaning of locally taking a fixed point of a subexpression.

```
data Desc : Type1 where
  - ...
  μ : Desc → Desc

[ ]o : Desc → ◇.Lang → ◇.Lang
- ...
[ μ D ]o _ = [ D ]
```

[JR] How is this used in our example?

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

```
νbrackets : Dec (◇.ν brackets)
νbrackets = yes (roll (inj1 refl))
```

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

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

```
bracketsD' = μ bracketsD * '[' ] ∪ νbrackets · var ∪ var * μ bracketsD
brackets'  = [ bracketsD' ]
```

Conclusion:

- We can reuse many of the results of finite languages (Section 2).

- We need a new μ combinator to nest fixed points in descriptions. This is necessary to refer back to the original language before reduction.
- Reducing a self-reference simply results in a self-reference again, because self-references in the reduct refer to the reduct.

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

3.3 Parsing in General

Our goal is to define:

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

We approach this by decomposing parsing into ν and δ .

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

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

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

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

The actual parsing can now be done character by character:

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

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

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

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 μ case we use the nullability of the fixed point, which we will implement shortly.

[JR] Reiterate that the cases for the basic combinators are the same as in Figure 2.

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

$\nu_o \ \nu P \ \emptyset = \text{no } \lambda \ ()$

$\nu_o \ \nu P \ \epsilon = \text{yes refl}$

$\nu_o \ \nu P \ (' \ c) = \text{no } \lambda \ ()$

$\nu_o \ \nu P \ (D \cup D_1) = \nu_o \ \nu P \ D \uplus\text{-dec } \nu_o \ \nu P \ D_1$

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

$\nu_o \ \nu P \ (x \cdot D) = x \times\text{-dec } \nu_o \ \nu P \ D$

$\nu_o \ \nu P \ \text{var} = \nu P$

$\nu_o \ \nu P \ (\mu \ D) = \nu D \ D$

- 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

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

$$319 \quad \nu D \emptyset \Leftrightarrow \nu D : \diamond.\nu \left(\llbracket D \rrbracket_o \diamond \emptyset \right) \Leftrightarrow \diamond.\nu \llbracket D \rrbracket$$

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

$$323 \quad \nu D \emptyset \rightarrow \nu D : \forall D \rightarrow \diamond.\nu \left(\llbracket D \rrbracket_o \diamond \emptyset \right) \rightarrow \diamond.\nu \left(\llbracket D \rrbracket_o \llbracket D_0 \rrbracket \right)$$

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

$$330 \quad \nu D \emptyset \leftarrow \nu D : \forall D \rightarrow \diamond.\nu \left(\llbracket D \rrbracket_o \llbracket D_0 \rrbracket \right) \rightarrow \diamond.\nu \left(\llbracket D \rrbracket_o \diamond \emptyset \right) \uplus \diamond.\nu \left(\llbracket D_0 \rrbracket_o \diamond \emptyset \right)$$

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

$$335 \quad \nu D \emptyset \Leftrightarrow \nu D \{D\} = \text{mk} \Leftrightarrow (\text{roll} \circ \nu D \emptyset \rightarrow \nu D D) (\text{reduce} \circ \nu D \emptyset \leftarrow \nu D \{D_0 = D\} D \circ \text{unroll})$$

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

$$337 \quad \nu D = \text{Dec.map } \nu D \emptyset \Leftrightarrow \nu D \circ \nu_o (\text{no } \lambda ())$$

338 *Remark 4.* Lemma 1 does not define an isomorphism on types. In particular, the back-
 339 wards direction is not injective. Consider the brackets language. It has the following null
 340 element, where we first choose the third disjunct, `var * var`, and then the first disjunct
 341 `ε` for both branches.

$$342 \quad \text{brackets}_0 : \diamond.\nu \text{ brackets}$$

$$343 \quad \text{brackets}_0 = \text{roll} (\text{inj}_2 (\text{inj}_2 (\llbracket \rrbracket, \llbracket \rrbracket, \text{refl}, \text{roll} (\text{inj}_1 \text{ refl}), \text{roll} (\text{inj}_1 \text{ refl}))))$$

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

$$345 \quad \text{brackets}_0' : \nu D \emptyset \Leftrightarrow \nu D \{\text{bracketsD}\} .\text{to} (\nu D \emptyset \Leftrightarrow \nu D \{\text{bracketsD}\} .\text{from } \text{brackets}_0)$$

$$346 \quad \quad \equiv \text{roll} (\text{inj}_1 \text{ refl})$$

$$347 \quad \text{brackets}_0' = \text{refl}$$

348 It now directly takes the first disjunct, `ε`.

349 In practice, such problems should be avoided by using unambiguous languages, en-
 350 suring 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.

351 3.5 Reduction

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

354 In Section 3.2, we established that the meaning of self-references changes and thus
 355 they need to be replaced by local fixed points of the original language. We define a
 356 function σD to perform this substitution. It is a simple recursive function which replaces
 357 the **var** constructor with a given D' description.

```

358  $\sigma : \text{Desc} \rightarrow \text{Desc} \rightarrow \text{Desc}$ 
359  $\sigma \emptyset \quad D' = \emptyset$ 
360  $\sigma \epsilon \quad D' = \epsilon$ 
361  $\sigma (' c) \quad D' = ' c$ 
362  $\sigma (D \cup D_1) \quad D' = \sigma D D' \cup \sigma D_1 D'$ 
363  $\sigma (D * D_1) \quad D' = \sigma D D' * \sigma D_1 D'$ 
364  $\sigma (x \cdot D) \quad D' = x \cdot \sigma D D'$ 
365  $\sigma \text{var} \quad D' = D'$ 
366  $\sigma (\mu D) \quad D' = \mu D$ 

```

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

```

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

```

370 The proof follows directly by induction and computation.

```

371  $\delta_o : \text{Desc} \rightarrow \text{Char} \rightarrow \text{Desc} \rightarrow \text{Desc}$ 
372  $\delta_o D_0 c \emptyset = \emptyset$ 
373  $\delta_o D_0 c \epsilon = \emptyset$ 
374  $\delta_o D_0 c (' c') = (c \stackrel{?}{=} c') \cdot \epsilon$ 
375  $\delta_o D_0 c (D \cup D_1) = \delta_o D_0 c D \cup \delta_o D_0 c D_1$ 
376  $\delta_o D_0 c (D * D_1) = \nu_o (\nu D D_0) D \cdot \delta_o D_0 c D_1 \cup \delta_o D_0 c D * \sigma D_1 (\mu D_0)$ 
377  $\delta_o D_0 c (x \cdot D) = x \cdot \delta_o D_0 c D$ 
378  $\delta_o D_0 c \text{var} = \text{var}$ 
379  $\delta_o D_0 c (\mu D) = \mu (\delta D c D)$ 

```

380 $\delta D c D = \delta_o D c D$

```

381  $\delta D\text{-to} : \forall D \rightarrow \llbracket \delta_o D_0 c D \rrbracket_o \llbracket \delta D c D_0 \rrbracket w \rightarrow \diamond \delta c (\llbracket D \rrbracket_o \llbracket D_0 \rrbracket) w$ 

```

```

382  $\delta D\text{-to} (' c') = \diamond \delta' . \text{to}$ 
383  $\delta D\text{-to} (D \cup D_1) (\text{inj}_1 x) = \text{inj}_1 (\delta D\text{-to} D x)$ 
384  $\delta D\text{-to} (D \cup D_1) (\text{inj}_2 y) = \text{inj}_2 (\delta D\text{-to} D_1 y)$ 
385  $\delta D\text{-to} (D * D_1) (\text{inj}_1 (x, y)) = \llbracket , \_ , \text{refl} , x , \delta D\text{-to} D_1 y$ 
386  $\delta D\text{-to} (D * D_1) (\text{inj}_2 (u, v, \text{refl}, x, y)) = (\_ :: \_) , \_ , \text{refl} , \delta D\text{-to} D x , \text{subst id } (\sigma \mu D_1) y$ 
387  $\delta D\text{-to} (A \cdot D) (x, y) = x , \delta D\text{-to} D y$ 
388  $\delta D\text{-to} \{D_0 = D\} \text{var} (\text{roll } x) = \text{roll } (\delta D\text{-to} D x)$ 
389  $\delta D\text{-to} (\mu D) (\text{roll } x) = \text{roll } (\delta D\text{-to} D x)$ 

```

```

390  $\delta D\text{-from} : \forall D \rightarrow \diamond \delta c (\llbracket D \rrbracket_o \llbracket D_0 \rrbracket) w \rightarrow \llbracket \delta_o D_0 c D \rrbracket_o \llbracket \delta D c D_0 \rrbracket w$ 

```

```

391  $\delta D\text{-from} (' c') = \diamond \delta' . \text{from}$ 
392  $\delta D\text{-from} (D \cup D_1) (\text{inj}_1 x) = \text{inj}_1 (\delta D\text{-from} D x)$ 
393  $\delta D\text{-from} (D \cup D_1) (\text{inj}_2 y) = \text{inj}_2 (\delta D\text{-from} D_1 y)$ 
394  $\delta D\text{-from} (D * D_1) (\llbracket , w , \text{refl} , x , y \rrbracket) = \text{inj}_1 (x , \delta D\text{-from} D_1 y)$ 
395  $\delta D\text{-from} (D * D_1) (c :: u , v , \text{refl} , x , y) = \text{inj}_2 (u , v , \text{refl} , \delta D\text{-from} D x , \text{subst id } (\text{sym } (\sigma \mu D_1) y))$ 
396  $\delta D\text{-from} (A \cdot D) (x, y) = x , \delta D\text{-from} D y$ 
397  $\delta D\text{-from} \{D_0 = D\} \text{var} (\text{roll } x) = \text{roll } (\delta D\text{-from} D x)$ 
398  $\delta D\text{-from} (\mu D) (\text{roll } x) = \text{roll } (\delta D\text{-from} D x)$ 

```

```

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

```

400 4 Discussion

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

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

404 *Expressiveness* We conjecture that our grammars which include variables and fixed
405 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 μ -regular languages

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

413 *Performance* For a parser to practically useful, it must at least have linear asymptotic
414 complexity for practical grammars. Might et al. [3] show that naively parsing using
415 derivatives does not achieve that bound, but optimizations might make it possible. In
416 particular, they argue that we could achieve $O(n|G|)$ time complexity (where $|G|$ is the
417 grammar size) if the grammar size stays approximately constant after every derivative.
418 By compacting the grammar, they conjecture it is possible to achieve this bound for
419 any unambiguous grammar. We want to investigate if similar optimizations could be
420 applied to our parser and if we can prove that we achieve this bound.

[JR] cite Jeremy Yallop's work

421 *Simplicity* One of the main contributions of Elliott's type theoretic formalization of
422 languages [2] is its simplicity of implementation and proof. To be able to extend his
423 approach to context-free languages we have had to introduce some complications.

[JR] TODO: finish this paragraph

424 5 Related Work

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

439 6 Conclusion

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

References

1. Afroozeh, 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. Elliott, C.: Symbolic and automatic differentiation of languages. *Proc. ACM Program. Lang.* **5**(ICFP) (Aug 2021). <https://doi.org/10.1145/3473583>
3. 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>