[TODO: Fill in acknowledgements in cover.tex]

[TODO: Move related work to end, so it flows better as explaining various other ways of parsing? Or introduce CFGs earlier.]

[TODO: Fill in related work with detailed explanations of various ways of writing parsers]

[TODO: Splitter now returns numbers, not strings]

[TODO: Explain how soundness can be done without parser extensionality, at the cost of algorithmic complexity elsewhere.]

[TODO: Find a citation for Fiat, test with [?]]

[TODO: (Optional) Section on showing that parser has "reasonable" performance on grammars with non-brute-force splitter (by using arrays and native strings)]

[TODO: (Really Optional) Section on building parse trees, not just recognizers)]

[QUESTION FOR ADAM: Should I include an appendix of all of the code of Fiat and parsers that is used, rendered verbatim?]

# An Extensible Framework for Synthesizing Efficient, Verified Parsers

by

## Jason S. Gross

Submitted to the Department of Electrical Engineering and Computer Science

in partial fulfillment of the requirements for the degree of

Master of Science in Computer Science and Engineering

at the

#### MASSACHUSETTS INSTITUTE OF TECHNOLOGY

## September 2015

© Jason S. Gross, MMXV. All rights reserved.

The author hereby grants to MIT permission to reproduce and to distribute publicly paper and electronic copies of this thesis document in whole or in part in any medium now known or hereafter created.

Author	
	Department of Electrical Engineering and Computer Science
	$\mathrm{July}\ 9,\ 2015$
Certified	by
	Adam Chlipala
	Associate Professor without Tenure of Computer Science
	Thesis Supervisor
Accepted	by
	Leslie A. Kolodziejski
	Chair, Department Committee on Graduate Students

# An Extensible Framework for Synthesizing Efficient, Verified Parsers

by Jason S. Gross

Submitted to the Department of Electrical Engineering and Computer Science on July 9, 2015, in partial fulfillment of the requirements for the degree of Master of Science in Computer Science and Engineering

#### Abstract

Parsers have a long history in computer science. This thesis proposes a novel approach to synthesizing efficient, verified parsers by refinement, and presents a demonstration of this approach in the Fiat framework by synthesizing a parser for arithmetic expressions. The benefits of this framework may include more flexibility in the parsers that can be described, more control over the low-level details when necessary for performance, and automatic or mostly automatic correctness proofs.

Thesis Supervisor: Adam Chlipala

Title: Associate Professor without Tenure of Computer Science

# Acknowledgments

This is the acknowledgements section. You should replace this with your own acknowledgements. [TODO: Fill this in]

# Contents

1	Parsing Context-Free Grammars					
	1.1	Parsing	9			
		1.1.1 Infinite Regress	10			
		1.1.2 Aborting Early	11			
	1.2	Standard Formal Definitions	12			
		1.2.1 Context-Free Grammar	13			
		1.2.2 Parse Trees	13			
<b>2</b>	Related Work and Other Approaches to Parsing 1					
	2.1	Recursive Descent Parsing	15			
		2.1.1 Parser Combinators	15			
	2.2	What's New and What's Old	19			
3	Con	mpleteness and Soundness 21				
4	Con	npleteness, Soundness, and Parsing Parse Trees	23			
	4.1	Proving Completeness: Conceptual Approach	23			
	4.2	Minimal Parse Trees: Formal Definition	24			
	4.3	Parser Interface	25			
		4.3.1 Parsing Parses	29			
		4.3.2 Example	31			
		4.3.3 Parametricity	34			
		4.3.4 Putting It All Together	35			
	4.4	Missteps, Insights, and Dependently Typed Lessons	35			
		4.4.1 The Trouble of Choosing the Right Types	36			
		4.4.2 Misordered Splitters	36			
		4.4.3 Minimal Parse Trees vs. Parallel Traces	37			
5	Ref	ining Splitters by Fiat	39			
	5.1	Splitters at a Glance	39			
	5.2	What counts as efficient?	39			
	5.3	Introducing Fiat	40			
		5.3.1 Incremental Construction by Refinement	40			
		5.3.2 The Fiat Mindset	40			
	5.4	Optimizations	43			

		5.4.1	An Easy First Optimization: Indexed Representation of Strings	43		
		5.4.2	Upcoming Optimizations	44		
6	Fixe	ed Len	gth Items, Parsing (ab)*; Parsing #s; Parsing #, ()	47		
	6.1	Parsin	g (ab)*: At Most One Nonterminal	47		
	6.2	Parsin	g Parenthesized Numbers: Fixed Lengths	48		
	6.3		g It Together	49		
7	Disj	isjoint Items, Parsing $\#$ , + 51				
8 Parsing well-parenthesized expressions				55		
	8.1	At a C	${f Glance}$	55		
	8.2	Gramı	mars we can parse	55		
		plitting Strategy	56			
		8.3.1		56		
		8.3.2	Building the Lookup Table	56		
		8.3.3	Table Correctness	57		
		8.3.4	Diving into Refinement Code	58		
9	Fut	Future work 6				
	9.1	Future	e work with dependent types	63		

# Chapter 1

# Parsing Context-Free Grammars

We begin with an overview of the general setting and a description of our approach to parsing.

## 1.1 Parsing

The job of a parser is to decompose a flat list of characters, called a *string*, into a structured tree, called a *parse tree*, on which further operations can be performed. As a simple example, we can parse "ab" as an instance of the regular expression  $(ab)^*$ , giving this parse tree, where we write  $\cdot$  for string concatenation.

$$\frac{\|\mathbf{a}\| \in \mathbf{a}\|}{\|\mathbf{a}\| \in \mathbf{b}\|} \frac{\|\mathbf{a}\| \in \epsilon}{\|\mathbf{a}\|} \frac{\|\mathbf{a}\| \in \epsilon}{\|\mathbf{a}\|}$$

$$\frac{\|\mathbf{a}\| \cdot \|\mathbf{b}\| \cdot \|\| \in \mathbf{ab}(\mathbf{ab})^*}{\|\mathbf{ab}\| \in (\mathbf{ab})^*}$$

Our parse tree is implicitly constructed from a set of general inference rules for parsing. There is a naive approach to parsing a string s: run the inference rules as a logic program. Several execution orders work: we may proceed bottom-up, by generating all of the strings that are in the language and not longer than s, checking each one for equality with s; or top-down, by splitting s into smaller parts in a way that mirrors the inference rules. In this paper, we present an implementation based on the second strategy, parameterizing over a "splitting oracle" that provides a list of candidate locations at which to split the string, based on the available inference rules. Soundness of the algorithm is independent of the splitting oracle; each location in the list is attempted. To be complete, if any split of the string yields a valid parse, the oracle must give at least one splitting location that also yields a valid parse. Different splitters yield different simple recursive-descent parsers.

Note that, for soundness and completeness, there is a trivial splitter: it returns a list of all numbers between 0 and the length of the string.

There is a trivial, brute-force splitter that suffices for proving correctness: simply return the list of all locations in the string, the list of all numbers between 0 and the length of the string. Because we construct a parser that terminates no matter what list it is given, and all valid splits are trivially in this list, this splitting "oracle" is enough to fill the oracle-shaped-hole in the correctness proofs. Thus, we can largely separate concerns about correctness and concerns about efficiency. In sections [TODO: secref], we focus only on correctness, we set up the framework we use to achieve efficiency in section [TODO: secref], and we demonstrate the use of the framework in sections [TODO: secref].

Although this simple splitter is sufficient for proving the algorithm correct, it is horribly inefficient, running in time  $\mathcal{O}(n!)$ , where n is the length of the string. We synthesize more efficient splitters in later chapters; we believe that parameterizing the parser over a splitter gives us enough expressiveness to implement essentially all optimizations of interest, while being a sufficiently simple language to make proofs relatively straightforward. For example, to achieve linear parse time on the (ab)\* grammar, we could have a splitter that, when trying to parse  $c_1 c_2 c_3 c_4 c_5 c_5 c_5 c_6 c_6$  as ab(ab)\*, splits the string into  $c_1 c_2 c_3 c_5 c_5 c_6 c_5 c_6 c_6$  and when trying to parse s as  $c_6 c_6 c_6 c_6 c_6 c_6$  not split the string at all.

Parameterizing over a splitting oracle allows us to largely separate correctness concerns from efficiency concerns.

Proving completeness—that our parser succeeds whenever there is a valid parse tree—is conceptually straightforward: trace the algorithm, showing that if the parser returns false at a given point, then assuming a corresponding parse tree exists yields a contradiction. The one wrinkle in this approach is that the algorithm, the logic program, is not guaranteed to terminate.

## 1.1.1 Infinite Regress

Since we have programmed our parser in Coq, our program must be terminating by construction. However, naive recursive-descent parsers do not always terminate!

To see how such parsers can diverge, consider the following example. When defining the grammar (ab)\*, perhaps we give the following production rules:

$$\frac{s \in \epsilon}{s \in (\mathsf{ab})^*} \stackrel{(\epsilon)}{=} \frac{s_0 \in \texttt{'a'} \quad s_1 \in \texttt{'b'}}{s_0 s_1 \in (\mathsf{ab})^*} \stackrel{(\texttt{"ab"})}{=} \frac{s_0 \in (\mathsf{ab})^* \quad s_1 \in (\mathsf{ab})^*}{s_0 s_1 \in (\mathsf{ab})^*} \stackrel{(\texttt{(ab)}^*(\mathsf{ab})^*)}{=}$$

Now, let us try to parse the string "ab" as (ab)\*:

$$\frac{\frac{}{\text{"""} \in \epsilon}}{\text{"""} \in (ab)^*} \frac{\frac{}{\text{"ab"} \in (ab)^*}}{\text{"ab"} \in (ab)^*}$$

$$\frac{\text{"""} \cdot \text{"ab"} \in (ab)^*}{\text{"ab"} \in (ab)^*}$$

$$\frac{\text{"""} \cdot \text{"ab"} \in (ab)^*}{\text{"ab"} \in (ab)^*}$$

$$\frac{\text{"ab"} \in (ab)^*}{\text{"ab"} \in (ab)^*}$$

Thus, by making a poor choice in how we split strings and choose productions, we can quickly hit an infinite regress.

Assuming we have a function split: String  $\rightarrow$  [String  $\times$  String] which is our splitting oracle, we may write out a potentially divergent parser specialized to this grammar.

```
any_parses: [String \times String] \rightarrow Bool any_parses [] := false any_parses (("a", "b") :: _) := true any_parses ((s<sub>1</sub>, s<sub>2</sub>) :: rest_splits) := (parses s<sub>1</sub> && parses s<sub>2</sub>) || any_parses rest_splits parses: String \rightarrow Bool parses "" := true parses str := any_parses (split str)
```

If split returns ("", "ab") as the first item in its list when given "ab", then the code given above will diverge in the way demonstrated above with the infinite derivation tree.

## 1.1.2 Aborting Early

To work around this wrinkle, we keep track of what nonterminals we have not yet tried to parse the current string as, and we abort early if we see a repeat. Note that this strategy only works for grammars with finite sets of nonterminals, in line with most formalizations of context-free grammars. For our example grammar, since there is only one nonterminal, we only need to keep track of the current string. We refactor the above code to introduce a new parameter prev\_s, recording the previous string we were parsing. We use s < prev\_s to denote the test that s is strictly shorter than

prev\_s.

We can convince Coq that this definition is total via well-founded recursion on the length of the string passed to parses. For a more-complicated grammar, we'd need to use a well-founded relation that also included the number of nonterminals not yet tried for this string; we do this in Figure 4-2 in Subsection 4.3.2.

With this refactoring, however, completeness is no longer straightforward. We must show that aborting early does not eliminate good parse trees.

We devote the rest of this paper to describing an elegant approach to proving completeness. Ridge [11] carried out a proof about essentially the same algorithm in HOL4, a proof assistant that does not support dependent types. We instead refine our parser to have a more general polymorphic type signature that takes advantage of dependent types, supporting a proof strategy with a different kind of aesthetic appeal. Relational parametricity frees us from worrying about different control flows with different instantiations of the arguments: when care is taken to ensure that the execution of the algorithm does not depend on the values of the arguments, we are guaranteed that all instantiations succeed or fail together. Freed from this worry, we convince our parser to prove its own soundness and completeness by instantiating its arguments correctly.

## 1.2 Standard Formal Definitions

Before proceeding, we pause to standardize on terminology and notation for context-free grammars and parsers. In service of clarity for some of our later explanations, we formalize grammars via natural-deduction inference rules, a slightly nonstandard choice.

#### 1.2.1 Context-Free Grammar

A context-free grammar consists of items, which may be either terminals (characters) or nonterminals; plus a set of productions, each mapping a nonterminal to a sequence of items.

#### Example: (ab)\*

The inference rules of the regular-expression grammar (ab)\* are:

Terminals:

Productions and nonterminals:

$$\frac{s \in \epsilon}{s \in (ab)^*} \qquad \boxed{"" \in \epsilon}$$

$$\frac{s_0 \in \texttt{'a'} \quad s_1 \in \texttt{'b'} \quad s_2 \in (\texttt{ab})^*}{s_0 s_1 s_2 \in (\texttt{ab})^*}$$

An alternative, more standard, more compact, notation for specifying context free grammars would present this grammar as:

$$(ab)^* := \epsilon \mid 'a' \mid b' (ab)^*$$

This can be read as saying that there is a single nonterminal (ab)\*, which parses empty strings, as well as parsing strings which are an 'a', followed by a 'b', followed by a string which parses as the nonterminal (ab)\*.

#### 1.2.2 Parse Trees

A string s parses as:

- a given terminal ch iff s = 'ch'.
- a given sequence of items  $x_i$  iff s splits into a sequence of strings  $s_i$ , each of which parses as the corresponding item  $x_i$ .
- a given nonterminal nt iff s parses as one of the item sequences that nt maps to under the set of productions.

We may define mutually inductive dependent type families of ParseTreeOfs and ParseItemsTreeOfs for a given grammar:

```
\label{eq:parseTreeOf} {\tt ParseTreeOf}: {\tt Item} \to {\tt String} \to {\tt Type} \\ {\tt ParseItemsTreeOf}: [{\tt Item}] \to {\tt String} \to {\tt Type} \\
```

For any terminal character ch, we have the constructor

```
('ch'): ParseTreeOf 'ch' "ch"
```

For any production rule mapping a nonterminal nt to a sequence of items its, and any string s, we have this constructor:

```
(\mathtt{rule}): \mathtt{ParseItemsTreeOf} its \mathtt{s} 	o \mathtt{ParseTreeOf} nt \mathtt{s}
```

We have the following two constructors of ParseItemsTree. In writing the type of the latter constructor, we adopt a common space-saving convention where we assume that all free variables are quantified implicitly with dependent function ( $\Pi$ ) types. We also write constructors in the form of schematic natural-deduction rules, since that notation will be convenient to use later on.

```
\begin{tabular}{ll} \hline & $"" \in \epsilon$ : ParseItemsTreeOf [] & "" \\ \hline & s_1 \in it & s_2 \in its \\ \hline & s_1s_2 \in it::its \\ \hline & \rightarrow ParseItemsTreeOf & its & s_2 \\ \hline & \rightarrow ParseItemsTreeOf & (it::its) & s_1s_2 \\ \hline \end{tabular}
```

For brevity, we will sometimes use the notation  $s \in X$  to denote both ParseTreeOf X s and ParseItemsTreeOf X s, relying on context to disambiguate based on the type of X. Additionally, we will sometimes fold the constructors of ParseItemsTreeOf into the (rule) constructors of ParseTreeOf, to mimic the natural-deduction trees.

We also define a type of all parse trees, independent of the string and item, as this dependent-pair  $(\Sigma)$  type, using set-builder notation; we use ParseTree to denote the type

```
\{(\mathtt{nt},\mathtt{s}) : \mathtt{Nonterminal} \times \mathtt{String} \mid \mathtt{ParseTreeOf} \ \mathtt{nt} \ \mathtt{s}\}
```

# Chapter 2

# Related Work and Other Approaches to Parsing

Stepping back a bit, we describe how our approach to parsing relates to existing work.

[TODO: Mention LR (and LL(1)?) parsers: "The field of parsing is one of the most venerable in computer-science. Still with us are a variety of parsing approaches born in times of much more severe constraints on memory and processor speed, including various flavors of LR parsers, which apply only to strict subsets of the context-free grammars, to guarantee ability to predict which production applies based on finite look-ahead into a string."

## 2.1 Recursive Descent Parsing

The most conceptually straightforward approaches to parsing fall into the class called recursive descent parsing, where, to parse a string s as a given production p, you attempt to parse various parts of s as each of the items in the list p. The control flow of the code mirrors the structure of the grammar, as well as the structure of the eventual parse tree, descending down the branches of the parse tree, recursively calling itself at each step.

#### 2.1.1 Parser Combinators

A popular approach, called *combinator parsing* [5], to implementing recursive descent parsing involves writing a small set of typed combinators, or higher-order functions, which are then applied to each other in various combinations to write a parser that mimics closely the structure of the grammar.

Essentially, parsers defined via parser combinators answer the question "what prefixes

of a given string be parsed as a given item?" Each function returns a list of postfixes of the string they are passed, indicating all of the strings that might remain for the other items in a given rule. [QUESTION FOR ADAM: Should I mention semantic actions?]

#### **Basic Combinators**

We now define the basic combinators. In the simplest form, each combinator takes in a string, and returns a list of strings (the postfixes); we can define the type parser as String  $\rightarrow$  [String]. We can define the empty-string parser, as well as the parser for a nonterminal with no production rules, which always fails:

```
\epsilon : parser \epsilon str = [str] fail : parser fail _ = []
```

Failure is indicated by returning the empty list; success at parsing the entire string is indicated by returning a list containing the empty string.

The parser for a given terminal fails if the string does not start with that character, and returns all but the first character if it does:

```
terminal : Char \rightarrow parser
terminal ch (ch :: str) = [str]
terminal _ "" = []
```

We now define combinators for sequencing and alternatives:

```
(>>>) : parser \rightarrow parser \rightarrow parser (p_0 >>> p_1) str = flat_map p_1 (p_0 \text{ str}) (|||) : parser \rightarrow parser \rightarrow parser (p_0 |||| p_1) str = p_0 \text{ str} ++ p_1 \text{ str}
```

where ++ is list concatenation, and flat\_map, which concatenates the lists returned by mapping its first argument over each of the elements in its second argument, has type ( $A \rightarrow [B]$ )  $\rightarrow [A] \rightarrow [B]$ , where [A] is the type denoting a list of elements, all of type A.[TODO: Check if I've defined this notation already.]

#### An Example

We can now easily define a parser for the grammar (ab)\*:

```
parse-(ab)* : parser parse-(ab)* = (terminal 'a' >>> terminal 'b' >>> parse-(ab)*) ||| \epsilon
```

Note that, by putting  $\epsilon$  last, we ensure that this parser returns the list in order of longest parse (shortest postfix) to shortest parse (longest postfix).

#### Proving Correctness and Dealing with Nontermination

Although parser combinators are straightforward, it is easy to make them loop forever. It is well-known that parsers defined naively using parser combinators don't handle grammars with *left recursion*, where the first item in a given production rule is the nonterminal currently being defined. For example, if we have the nonterminal expr ::= number | expr '+' expr, then the parser for expr '+' expr will call the parser for expr, which will call the parse for expr, which will quickly loop forever.

The algorithm we presented in Subsection 1.1.2 is essentially the same as the algorithm Ridge presents in [11] to deal with this problem. By wrapping the calls to the parsers, in each combinator, with a function that prunes duplicative calls, Ridge provides a way to ensure that parsers terminate. Also included in [11] are proofs in HOL4 that such wrapped parsers are both sound (and therefore terminating) and complete. Furthermore, Ridge's parser has worst-case  $O(n^5)$  running time in the input-string length.

## [TODO: Fill in content] [QUESTION FOR ADAM: Am I missing any branches?]

- Say something about recursive descent parsing seeming obvious and trivial only after writing out the types inductively?
- Say something about LR parsers and processor/memory constraints?
- Parser combinators
  - parse by consuming characters from the string.
  - Allow incremental parsing
  - Can be extended to return lists of possible parses
  - Ridge describes a way to wrap parser combinators to guarantee termination and ensure completeness; the algorithm we presented above uses essentially the same strategy. He also proves correctness.
- Parsing with derivatives

- a kind-of conceptual dual to standard parser-combinators; rather than mutating the string as we go along, we mutate the language
- With parser combinators, we drop bits of the string handled by each rule as we go along. With derivatives, we drop bits of each rule (or entire rules) as we go down the string.
- [TODO: Look into whether its been verified]
- Look into what [4]s and [12]s are.
- Approaches to verifying parsers: correct-by-construction, induction. Other way? [TODO: Read paper]

#### [TODO: Rewrite the rest of this part so it's not just Adam's words]

However, despite rumors to the contrary, the field of parsing is far from dead. In the twentieth century, the functional-programming world experimented with a variety of approaches to parser combinators [5], where parsers are higher-order functions built from a small set of typed combinators. In the twenty-first century alone, a number of new parsing approaches have been proposed or popularized, including parsing expression grammars (PEGs) [4], derivative-based parsing [8], and GLL parsers [12].

However, our approach is essentially the same, algorithmically, as the one that Ridge demonstrated with a verified parser-combinator system [11], taking naive recursive-descent parsing and adding a layer to prune duplicative calls to the parser. His proof was carried out in HOL4, necessarily without using dependent types. Our new work may be interesting for the aesthetic appeal of our unusual application of dependent types to get the parser to generate some of its own soundness proof. Ridge's parser also has worst-case  $O(n^5)$  running time in the input-string length. In the context of our verified implementation, we plan to explore a variety of optimizations based on clever, grammar-specific choices of string-splitter functions, which should have a substantial impact on the run-time cost of parsing some relevant grammars, and which we conjecture will not require any changes to the development presented in this paper.

A few other past projects have verified parsers with proof assistants, applying to derivative-based parsing [1] and SLR [2] and LR(1) [6] parsers. Several projects have used proof assistants to apply verified parsers within larger programming-language tools. RockSalt [9] does run-time memory-safety enforcement for x86 binaries, relying on a verified machine-code parser that applies derivative-based parsing for regular expressions. The verified Jitawa [10] and CakeML [7] language implementations include verified parsers, handling Lisp and ML languages, respectively.

Our final parser derivation relies on a relational parametricity property for polymorphic functions in Coq's type theory Gallina. With Coq as it is today, we need to prove

this property manually for each eligible function, even though we can prove metatheoretically that it holds for them all. Bernardy and Guilhem [3] have shown how to extend type theories with support for materializing "free theorem" parametricity facts internally, and we might be able to simplify our implementation using such a feature. [TODO: Flesh this out, rewrite it so it's not all Adam's words, read the papers]

#### 2.2 What's New and What's Old

The goal of this project is to demonstrate a new approach to generating parsers: incrementally building efficient parsers by refinement.

We begin with naive recursive-descent parsing. [TODO: Citation?] [TODO: Say something about viewing recursive-descent as an instance of general inhabitation-decision not being anywhere in the literature?] We ensure termination via memoization, a la [11]. We parameterize the parser on a "splitting oracle", which describes how to recurse. [TODO: section citation] As far as we can tell, the idea of factoring the algorithmic complexity like this is new.

We use Fiat to incrementally build efficient parsers by refinement. [TODO: section citation]

Additionally, we take a digression in Chapter 4 to describe how our parser can be used to prove its own completeness; the idea of reusing the parsing algorithm to generate proofs, parsing parse trees rather than strings, is not found in the literature, to the authors' knowledge.

# Chapter 3

# Completeness and Soundness

Parsers come in a number of flavors. The simplest flavor is the *recognizer*, which simply says whether or not there exists a parse tree of a given string for a given nonterminal; it returns Booleans. There is also a richer flavor of parser that returns inhabitants of option ParseTree.

For any recognizer has\_parse: Nonterminal  $\rightarrow$  String  $\rightarrow$  Bool, we may ask whether it is *sound*, meaning that when it returns true, there is always a parse tree; and *complete*, meaning that when there is a parse tree, it always returns true. We may express these properties as theorems (alternatively, dependently typed functions) with the following type signatures:

```
\begin{array}{c} {\tt has\_parse\_sound}: ({\tt nt}: {\tt Nonterminal}) \to ({\tt s}: {\tt String}) \\ & \to {\tt has\_parse} \ {\tt nt} \ {\tt s} = {\tt true} \\ & \to {\tt ParseTreeOf} \ {\tt nt} \ {\tt s} \\ {\tt has\_parse\_complete}: ({\tt nt}: {\tt Nonterminal}) \to ({\tt s}: {\tt String}) \\ & \to {\tt ParseTreeOf} \ {\tt nt} \ {\tt s} \\ & \to {\tt has\_parse} \ {\tt nt} \ {\tt s} = {\tt true} \end{array}
```

For any parser

```
parse: Nonterminal \rightarrow String \rightarrow option ParseTree,
```

we may also ask whether it is sound and complete, leading to theorems with the

following type signatures, using p<sub>1</sub> to denote the first projection of p:

```
\begin{array}{l} \texttt{parse\_sound} : (\texttt{nt} : \texttt{Nonterminal}) \\ & \to (\texttt{s} : \texttt{String}) \\ & \to (\texttt{p} : \texttt{ParseTree}) \\ & \to \texttt{parse} \ \texttt{nt} \ \texttt{s} = \texttt{Some} \ \texttt{p} \\ & \to \texttt{p_1} = (\texttt{nt}, \texttt{s}) \\ \texttt{parse\_complete} : (\texttt{nt} : \texttt{Nonterminal}) \\ & \to (\texttt{s} : \texttt{String}) \\ & \to \texttt{ParseTreeOf} \ \texttt{nt} \ \texttt{s} \\ & \to \texttt{parse} \ \texttt{nt} \ \texttt{s} \neq \texttt{None} \end{array}
```

Since we are programming in Coq, this separation into code and proof actually makes for more awkward type assignments. We also have the option of folding the soundness and completeness conditions into the types of the code. For instance, the following type captures the idea of a sound and complete parser returning parse trees, using the type constructor + for disjoint union (i.e., sum or variant type):

```
\begin{aligned} & \texttt{parse}: (\texttt{nt}: \texttt{Nonterminal}) \\ & \to (\texttt{s}: \texttt{String}) \\ & \to \texttt{ParseTreeOf} \ \texttt{nt} \ \texttt{s} + (\texttt{ParseTreeOf} \ \texttt{nt} \ \texttt{s} \to \bot) \end{aligned}
```

That is, given a nonterminal and a string, parse either returns a valid parse tree, or returns a proof that the existence of any parse tree is contradictory (i.e., implies  $\bot$ , the empty type). Our implementation follows this dependently typed style. Our main initial goal in the project was to arrive at a parse function of just this type, generic in an arbitrary choice of context-free grammar, implemented and proven correct in an elegant way.

# Chapter 4

# Completeness, Soundness, and Parsing Parse Trees

## 4.1 Proving Completeness: Conceptual Approach

Recall from Subsection 1.1.2 that the essential difficulty with proving completeness is dealing with the cases where our parser aborts early; we must show that doing so does not eliminate good parse trees.

The key is to define an intermediate type, that of "minimal parse trees." A "minimal" parse tree is simply a parse tree in which the same (string, nonterminal) pair does not appear more than once in any path of the tree. Defining this type allows us to split the completeness problem in two; we can show separately that every parse tree gives rise to a minimal parse tree, and that having a minimal parse tree in hand implies that our parser succeeds (returns true or Some \_).

Our dependently typed parsing algorithm subsumes the soundness theorem, the minimization of parse trees, and the proof that having a minimal parse tree implies that our parser succeeds. We write one parametrically polymorphic parsing function that supports all three modes, plus the several different sorts of parsers (recognizers, generating parse trees, running semantic actions). That level of genericity requires us to be flexible in which type represents "strings," or inputs to parsers. We introduce a parameter that is often just the normal String type, but which needs to be instantiated as the type of parse trees themselves to get a proof of parse tree minimizability. That is, we "parse" parse trees to minimize them, reusing the same logic that works for the normal parsing problem.

Before presenting our algorithm's interface, we will formally define and explain minimal parse trees, which will provide motivation for the type signatures of our parser's arguments.

#### 4.2 Minimal Parse Trees: Formal Definition

In order to make tractable the second half of the completeness theorem, that having a minimal parse tree implies that parsing succeeds, it is essential to make the inductive structure of minimal parse trees mimic precisely the structure of the parsing algorithm. A minimal parse tree thus might better be thought of as a parallel trace of parser execution.

As in Subsection 1.2.2, we define mutually inductive type families of MinParseTreeOfs and MinItemsTreeOfs for a given grammar. Because our parser proceeds by well-founded recursion on the length of the string and the list of nonterminals not yet attempted for that string, we must include both of these in the types. Let us call the initial list of all nonterminals unseen<sub>0</sub>.

```
\label{eq:minParseTreeOf} \begin{split} & \operatorname{MinParseTreeOf}: \operatorname{String} \to [\operatorname{Nonterminal}] \\ & \to \operatorname{Item} \to \operatorname{String} \to \operatorname{Type} \\ & \operatorname{MinItemsTreeOf}: \operatorname{String} \to [\operatorname{Nonterminal}] \\ & \to [\operatorname{Item}] \to \operatorname{String} \to \operatorname{Type} \end{split}
```

Much as in the case of parse trees, for any terminal character ch, any string  $s_0$ , and any list of nonterminals unseen, we have the constructor

```
min_parse'ch': MinParseTreeOf so unseen 'ch' "ch"
```

For any production rule mapping a nonterminal nt to a sequence of items its, any string  $s_0$ , any list of nonterminals unseen, and any string s, we have two constructors, corresponding to the two ways of progressing with respect to the well-founded relation. Letting unseen' := unseen - {nt}, we have the following, where we interpret the < relation on strings in terms of lengths.

```
\begin{array}{l} (\text{rule})_<: \texttt{s} < \texttt{s}_0 \\ & \to \texttt{MinItemsTreeOf} \ \texttt{s} \ \texttt{unseen}_0 \ \texttt{its} \ \texttt{s} \\ & \to \texttt{MinParseTreeOf} \ \texttt{s}_0 \ \texttt{unseen} \ \texttt{nt} \ \texttt{s} \\ (\text{rule})_=: \texttt{s} = \texttt{s}_0 \\ & \to \texttt{nt} \in \texttt{unseen} \\ & \to \texttt{MinItemsTreeOf} \ \texttt{s}_0 \ \texttt{unseen}' \ \texttt{its} \ \texttt{s} \\ & \to \texttt{MinParseTreeOf} \ \texttt{s}_0 \ \texttt{unseen} \ \texttt{nt} \ \texttt{s} \end{array}
```

In the first case, the length of the string has decreased, so we may reset the list of not-yet-seen nonterminals, as long as we reset the base of well-founded recursion  $s_0$  at the same time. In the second case, the length of the string has not decreased, so we require that we have not yet seen this nonterminal, and we then remove it from the list of not-yet-seen nonterminals.

Finally, for any string  $s_0$  and any list of nonterminals unseen, we have the following two constructors of MinItemsTreeOf.

```
\begin{split} & \min\_parse_{[]} : \texttt{MinItemsTreeOf} \ \ s_0 \ \ unseen \ \ [] \ \ "" \\ & \min\_parse_{::} : s_1s_2 \le s_0 \\ & \longrightarrow \texttt{MinParseTreeOf} \ \ s_0 \ \ unseen \ \ it \ \ s_1 \\ & \longrightarrow \texttt{MinItemsTreeOf} \ \ s_0 \ \ unseen \ \ its \ \ s_2 \\ & \longrightarrow \texttt{MinItemsTreeOf} \ \ s_0 \ \ unseen \ \ (it::its) \ \ s_1s_2 \end{split}
```

The requirement that  $s_1s_2 \leq s_0$  in the second case ensures that we are only making well-founded recursive calls.

Once again, for brevity, we will sometimes use the notation  $\overline{s} \in \overline{X}^{<(s_0,v)}$  to denote both MinParseTreeOf  $s_0$  v X s and MinItemsTreeOf  $s_0$  v X s, relying on context to disambiguate based on the type of X. Additionally, we will sometimes fold the constructors of MinItemsTreeOf into the two (rule) constructors of MinParseTreeOf, to mimic the natural-deduction trees.

## 4.3 Parser Interface

Roughly speaking, we read the interface of our general parser off from the types of the constructors for minimal parse trees. Every constructor leads to one parameter passed to the parser, much as one derives the types of general "fold" functions for arbitrary inductive datatypes. For instance, lists have constructors nil and cons, so a fold function for lists has arguments corresponding to nil (initial accumulator) and cons (step function). The situation for the type of our parser is similar, though we need parallel success (managed to parse the string) and failure (could prove that no parse is possible) parameters for each constructor of minimal parse trees.

The type signatures in the interface are presented in Figure 4-1. We explain each type one by one, presenting various instantiations as examples. Note that the interface we actually implemented is also parameterized over a type of Strings, which we will instantiate with parse trees later in this paper. The interface we present here fixes String, for conciseness.

Since we want to be able to specialize our parser to return either Bool or option ParseTree, we want to be able to reuse our soundness and completeness proofs for both. Our strategy for generalization is to parameterize on dependent type families for "success" and "failure", so we can use relational parametricity to ensure that all instantiations of the parser succeed or fail together. The parser has the rough type signature

```
\texttt{parse}: \texttt{Nonterminal} \ \to \ \texttt{String} \ \to \ \texttt{T}_{\texttt{success}} + \texttt{T}_{\texttt{failure}}.
```

We use ParseQuery to denote the type of all propositions like ""a"  $\in$  'a'"; a query consists of a string and a grammar rule the string might be parsed into. We use the same notation for ParseQuery and ParseTree inhabitants. All \*\_success and \*\_failure type signatures are implicitly parameterized over a string  $s_0$  and a list of nonterminals unseen. We assume we are given unseen<sub>0</sub>: [Nonterminal].

```
nonterminals unseen. We assume we are given unseen_0: [Nonterminal].
           \mathtt{T}_{\mathtt{success}}, \ \mathtt{T}_{\mathtt{failure}} : \mathtt{String} 	o [\mathtt{Nonterminal}] 	o \mathtt{ParseQuery} 	o \mathtt{Type}
                             \mathtt{split}:\mathtt{String} 	o [\mathtt{Nonterminal}] 	o \mathtt{ParseQuery} 	o [\mathbb{N}]
      \texttt{terminal\_success}: (\texttt{ch}: \texttt{Char}) \to \texttt{T}_{\texttt{success}} \ \texttt{s}_0 \ \texttt{unseen} \ (\texttt{"ch"} \in \texttt{'ch'})
      \texttt{terminal\_failure}: (\texttt{ch}: \texttt{Char}) \rightarrow (\texttt{s}: \texttt{String}) \rightarrow \texttt{s} \neq \texttt{"ch"} \rightarrow \texttt{T}_{\texttt{failure}} \ \texttt{s}_0 \ \texttt{unseen} \ (\overline{\texttt{s} \in \texttt{'ch'}})
                nil\_success: T_{success} s_0 unseen ( ( \in \epsilon ) )
                nil_failure: (s:String) \rightarrow s \neq "" \rightarrow T_{failure} s_0 unseen (\overline{s \in \epsilon})
               cons_success: (it: Item) \rightarrow (its: [Item]) \rightarrow (s<sub>1</sub>: String) \rightarrow (s<sub>2</sub>: String)
                                        \rightarrow s_1 s_2 < s_0
                                         \rightarrow T_{\text{success}} s_0 \text{ unseen } (s_1 \in it)
                                         \rightarrow T_{\text{success}} s_0 \text{ unseen } (s_2 \in \text{its})
                                         \rightarrow T_{\text{success}} s<sub>0</sub> unseen (s_1s_2 \in \text{it} :: \text{its})
              \texttt{cons\_failure}: (\texttt{it}: \texttt{Item}) \rightarrow (\texttt{its}: \texttt{[Item]}) \rightarrow (\texttt{s}: \texttt{String})
                                        \rightarrow s < s_0
                                        \rightarrow (\forall (s_1, s_2) \in split s_0 unseen (\overline{s \in it :: its}),
                                                    T_{\text{failure}} s<sub>0</sub> unseen (s_1 \in it) + T_{\text{failure}} s<sub>0</sub> unseen (s_2 \in its))
                                         \rightarrow T_{\text{failure}} s<sub>0</sub> unseen (s \in it::its)
production\_success_{<}:(its:[Item]) \rightarrow (nt:Nonterminal) \rightarrow (s:String)
                                         \rightarrow s < s_0
                                         \rightarrow (p: a production mapping nt to its)
                                         \rightarrow T_{\text{success}} s unseen<sub>0</sub> (s \in its)
                                         \rightarrow T_{\text{success}} s<sub>0</sub> unseen (\overline{s \in nt})
\mathtt{production\_success} = : (\mathtt{its} : [\mathtt{Item}]) \rightarrow (\mathtt{nt} : \mathtt{Nonterminal}) \rightarrow (\mathtt{s} : \mathtt{String})
                                         \rightarrow nt \in unseen
                                         \rightarrow (p: a production mapping nt to its)
                                         \rightarrow T_{\text{success}} \ s_0 \ (\text{unseen} - \{\text{nt}\}) \ (\overline{s \in \text{its}})
                                         \rightarrow T_{\text{success}} s_0 \text{ unseen } (\overline{s \in nt})
production\_failure_{<}:(nt:Nonterminal) \rightarrow (s:String)
                                         \rightarrow s < s_0
                                        \rightarrow (\forall (its:[Item]) (p:a production mapping nt to its), T_{\texttt{failure}} s unseen
                                         \rightarrow T_{\text{failure}} s<sub>0</sub> unseen (\overline{s \in nt})
production\_failure_= : (nt : Nontermina_k) \rightarrow (s : String)
                                         \rightarrow s = s_0
```

 $\rightarrow$  ( $\forall$  (its: [Item]) (p:a production mapping nt to its),  $T_{\text{failure}}$  s<sub>0</sub> (unse

To instantiate the parser as a Boolean recognizer, we instantiate everything trivially; we use the fact that  $\top + \top \cong Bool$ . Just to show how trivial everything is, here is a precise instantiation of the parser, still parameterized over the initial list of nonterminals and the splitter, where  $\top$  is the one constructor of the one-element type  $\top$ :

```
\begin{array}{lll} T_{\text{success}} & \_ & \_ & = \top \\ T_{\text{failure}} & \_ & = \top \\ \end{array} \begin{array}{lll} \text{terminal\_success} & \_ & \_ & = () \\ \text{terminal\_failure} & \_ & \_ & = () \\ \text{nil\_success} & \_ & = & = () \\ \text{nil\_failure} & \_ & \_ & = & () \\ \text{cons\_success} & \_ & \_ & = & = () \\ \text{cons\_failure} & \_ & \_ & = & = & () \\ \text{production\_success} & \_ & \_ & = & = & () \\ \text{production\_success} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & = & () \\ \text{production\_failure} & \_ & \_ & = & () \\ \text{production\_failure} & \_ & \_ & = & () \\ \text{pro
```

To instantiate our parser so that it returns option ParseTree (rather, the dependently typed flavor, ParseTreeOf), we take advantage of the isomorphism  $T + T \cong$  option T. We show only the success instantiations, as the failure ones are identical with the Boolean recognizer. For readability of the code, we write schematic natural-deduction proof trees inline.

$$\begin{array}{l} T_{\texttt{success}} \ \_ \ \_ \ (s \in X) \coloneqq s \in X \\ \\ \text{terminal\_success} \ \_ \ \_ \ ch \coloneqq (\texttt{'ch'}) \\ \\ \text{nil\_success} \ \_ \ \_ \coloneqq \overline{\texttt{"""} \in \epsilon} \\ \\ \text{cons\_success} \ \_ \ \_ \ \text{it its} \ s_1 \ s_2 \ \_ \ d_1 \ d_2 \coloneqq \frac{\frac{d_1}{s_1 \in \text{it}} \quad \frac{d_2}{s_2 \in \text{its}}}{s_1 s_2 \in \text{it} :: \text{its}} \\ \\ \text{production\_success}_{<} \ \_ \ \_ \ \text{it nt s} \ \_ \ p \ d \coloneqq \frac{\frac{d}{s \in \text{its}}}{s \in \text{nt}} \ _{(p)} \\ \\ \text{production\_success}_{=} \ \_ \ \ \text{it nt s} \ \_ \ p \ d \coloneqq \frac{\frac{d}{s \in \text{its}}}{s \in \text{nt}} \ _{(p)} \\ \\ \text{production\_success}_{=} \ \_ \ \ \text{it nt s} \ \_ \ p \ d \coloneqq \frac{\frac{d}{s \in \text{its}}}{s \in \text{nt}} \ _{(p)} \\ \\ \end{array}$$

What remains is to instantiate the parser in such a way that proving completeness is trivial. The simpler of our two tasks is to show that when the parser fails, no minimal parse tree exists. Hence we instantiate the types as follows, where  $\bot$  is the empty type (equivalently, the false proposition).

$$\begin{split} & \textbf{T}_{\texttt{success}} \text{ \_ \_ } := \top \\ & \textbf{T}_{\texttt{failure}} \text{ } \textbf{s}_0 \text{ } \textbf{unseen} \text{ } (\overline{\textbf{s} \in \textbf{X}}) \coloneqq \left(\overline{\textbf{s} \in \textbf{X}} \overset{<(\textbf{s}_0, \textbf{unseen})}{}\right) \to \bot \end{split}$$

Using  $\mathcal{I}$  to denote deriving a contradiction, we can unenlighteningly instantiate the arguments as

```
terminal_success _ _ _ := ()

terminal_failure _ _ _ _ := \mathcal{\epsilon}

nil_success _ _ := ()

nil_failure _ _ _ := \mathcal{\epsilon}

cons_success _ _ _ _ = _ _ := \mathcal{\epsilon}

cons_failure _ _ _ _ := \mathcal{\epsilon}

production_success < _ _ _ _ _ := \mathcal{\epsilon}

production_failure < _ _ _ := \mathcal{\epsilon}

production_failure = _ _ _ := \mathcal{\epsilon}

production_failure = _ _ _ := \mathcal{\epsilon}

production_failure = _ _ _ := \mathcal{\epsilon}
```

A careful inspection of the proofy arguments to each failure case will reveal that

there is enough evidence to derive the appropriate contradiction. For example, the  $s \neq$  "" hypothesis of nil\_failure contradicts the equalities implied by the type signature of min\_parse[], and the use of [] contradicts the equality implied by the use of it::its in the type signature of min\_parse[]. Similarly, the  $s \neq$  "ch" hypothesis of terminal\_failure contradicts the equality implied by the usage of the single identifier ch in two different places in the type signature of min\_parse\_ch'.

#### 4.3.1 Parsing Parses

We finally come to the most twisty part of the parser: parsing parse trees. Recall that our parser definition is polymorphic in a choice of String type. We proceed with the straw-man solution of literally passing in parse trees as strings to be parsed, such that parsing generates minimal parse trees, as introduced in Section 4.1 and defined formally in Section 4.2. Intuitively, we run a top-down traversal of the tree, pausing at each node before descending to its children. During that pause, we eliminate one level of wastefulness: if the parse tree is proving  $s \in X$ , we look for any subtrees also proving  $s \in X$ . If we find any, we replace the original tree with the smallest duplicative subtree. If we do not find any, we leave the tree unchanged. In either case, we then descend into "parsing" each subtree.

We define a function deloop to perform the one step of eliminating waste:

```
deloop : ParseTreeOf nt s \rightarrow ParseTreeOf nt s
```

This transformation is straightforward to define by structural recursion.

To implement all of the generic parameters of the parser, we must actually augment the result type of deloop with stronger types. Define the predicate Unloopy(t) on parse trees t to mean that, where the root node of t proves  $s \in nt$ , for every subtree proving  $s \in nt'$  (same string, possibly different nonterminal), (1) nt' is in the set of allowed nonterminals, unseen, associated to the overall tree with dependent types, and (2) if this is not the root node, then  $nt' \neq nt$ .

We augment the return type of deloop, writing:

```
\{t : ParseTreeOf nt s \mid Unloopy(t)\}.
```

We instantiate the generic "string" type parameter of the general parser with this type family, so that, in implementing the different parameters to pass to the parser, we have the property available to us.

Another key ingredient is the "string" splitter, which naturally breaks a parse tree

into its child trees. We define it like so:

$$\begin{array}{l} \mathtt{split} \ \_ \ ( \mathtt{s} \in \mathtt{it} :: \mathtt{its} ) \coloneqq \\ \\ \mathbf{case} \ \ \mathtt{parse\_tree\_data} \ \ \mathtt{s} \ \ \mathbf{of} \\ \big| \ \frac{\frac{p_1}{s_1 \in \mathtt{it}} \quad \frac{p_2}{s_2 \in \mathtt{its}}}{s_1 s_2 \in \mathtt{it} :: \mathtt{its}} \ \ \to \ \big[ \, (\mathtt{deloop} \ p_1, \mathtt{deloop} \ p_2) \, \big] \\ \big| \ \_ \ \ \to \ \ \not\{ \\ \\ \mathtt{split} \ \_ \ \_ \coloneqq [] \end{array}$$

Note that we use it and its nonlinearly; the pattern only binds if its it and its match those passed as arguments to split. We thus return a nonempty list only if the query is about a nonempty sequence of items. Because we use dependent types to enforce the requirement that the parse tree associated with a string match the query we are considering, we can derive contradictions in the non-matching cases.

This splitter satisfies two important properties. First, it never returns the empty list on a parse tree whose list of productions is nonempty; call this property nonempty preservation. Second, it preserves Unloopy. We use both facts in the other parameters to the generic parser (and we leave their proofs as exercises for the reader—Coq solutions may be found in our source code).

Now recall that our general parser always returns a type of the form  $T_{\text{success}} + T_{\text{failure}}$ , for some  $T_{\text{success}}$  and  $T_{\text{failure}}$ . We want our tree minimizer to return just the type of minimal trees. However, we can take advantage of the type isomorphism  $T + \bot \cong T$  and instantiate  $T_{\text{failure}}$  with  $\bot$ , the uninhabited type; and then apply a simple fix-up wrapper on top. Thus, we instantiate the general parser like so:

$$\begin{array}{ll} T_{\texttt{success}} & \texttt{s}_0 & \texttt{unseen} & (\texttt{d}: \texttt{s} \in \texttt{X}) \coloneqq \texttt{s} \in \texttt{X} & <(\texttt{s}_0, \texttt{unseen}) \\ \\ T_{\texttt{failure}} & \_ \ \_ \ \coloneqq \bot & \end{array}$$

The success cases are instantiated in an essentially identical way to the instantiation we used to get option ParseTree. The terminal\_failure and nil\_failure cases provide enough information ( $s \neq "ch"$  and  $s \neq ""$ , respectively) to derive  $\bot$  from the existence of the appropriately typed parse tree. In the cons\_failure case, we make use of the splitter's nonempty preservation behavior, after which all that remains is  $\bot + \bot \to \bot$ , which is trivial. In the production\_failure and production\_failure cases, it is sufficient to note that every nonterminal is mapped by some production to some sequence of items. Finally, to instantiate the production\_failure case, we need to appeal to the Unloopy-ness of the tree to deduce that  $nt \in unseen$ . Then we can derive  $\bot$  from the hypothesis that  $nt \notin unseen$ , and we are done.

We instantiate the general parser with an input type that requires Unloopy, so our

final tree minimizer is really the composition of the instantiated parser with deloop, ensuring that invariant as we kick off the recursion.

## 4.3.2 Example

In Subsection 1.1.1, we defined an ambiguous grammar for (ab)\* which led our naive parser to diverge. We will walk through the minimization of the following parse tree of "abab" into this grammar. For reference, Figure 4-2 contains the fully general implementation of our parser, modulo type signatures.

For reasons of space, define  $\overline{T}$  to be the parse tree

$$\frac{\frac{}{\text{"""} \in \epsilon}}{\text{"""} \in (ab)^*} \frac{\frac{\text{"a"} \in \text{'a'}}{\text{"a"} \cdot \text{"b"} \in (ab)^*}}{\frac{\text{"ab"} \in (ab)^*}{\text{"ab"} \in (ab)^*}}{\text{((ab)*(ab)*)}}$$

Then we consider minimizing the parse tree:

$$\frac{\overline{T}}{\text{"ab"} \in (ab)^*} \frac{\overline{T}}{\text{"ab"} \in (ab)^*}$$

$$\frac{\text{"ab"} \cdot \text{"ab"} \in (ab)^*}{\text{"abab"} \in (ab)^*}$$

Letting  $\overline{T'_m}$  denote the same tree as  $\overline{T'}$ , but constructed as a MinParseTree rather than a ParseTree, the tree we will end up with is:

$$\frac{\overline{T'_m}}{\text{"ab"} \in (ab)^*} < (\text{"ab"}, [(ab)^*]) \qquad \frac{\overline{T'_m}}{\text{"ab"} \in (ab)^*} < (\text{"ab"}, [(ab)^*]) \\ \frac{\text{"ab"} \cdot \text{"ab"} \in (ab)^*}{\text{"abab"} \in (ab)^*} < (\text{"abab"}, [(ab)^*])$$

To begin, we call parse, passing in the entire tree as the string, and (ab)\* as the nonterminal. To transform the tree into one that satisfies Unloopy, the first thing parse does is call deloop on our tree. In this case, deloop is a no-op; it promotes the deepest non-root nodes labeled with ("abab"  $\in$  (ab)\*), of which there are none.

We then take the following execution steps, starting with unseen := unseen<sub>0</sub> := [(ab)\*], the singleton list containing the only nonterminal, and  $s_0 :=$  "abab".

- 1. We first ensure that we are not in an infinite loop. We check if  $s < s_0$  (it is not, for they are both equal to "abab"), and then check if our current nonterminal, (ab)\*, is in unseen. Since the second check succeeds, we remove (ab)\* from unseen; calls made by this stack frame will pass [] for unseen.
- 2. We may consider only the productions for which the parse tree associated to

```
parse nt s := parse' (s_0 := s) (unseen := unseen_0) (s \in nt)
parse' ("ch" \in 'ch') := inl terminal_success
parse' (\underline{\ } \in \underline{\ } ch') := inr (terminal_failure 1)
parse' ( (iiii \in \epsilon) := inl nil_success)
parse' (\underline{\phantom{a}} \in \epsilon) := inr (nil_failure f)
parse' (s \in it :: its) :=
     case any_parse s it its (split (s \in it :: its)) of
       inl ret \rightarrow inl ret
       inr ret \rightarrow inr (cons\_failure \_ ret)
parse' (\overline{s} \in nt) :=
     if s < s_0
    then if (parse' (s_0 := s) (unseen := unseen<sub>0</sub>) (s \in its)) succeeds returning d
                                 for any production p mapping nt to its
              then inl (production_success< _ p d)</pre>
              else inr (production_failure< _ _)</pre>
     else if nt \in unseen
              then if (parse' (unseen = unseen - {nt}) (s \in its)) succeeds returning d
                                          for any production p mapping nt to its
                       then inl (production_success_ _ p d)
                       else inr (production_failure= _ _)
              else inr (production_failure<sub>∉</sub> _ _)
any_parse s it its [] := inr(\lambda_{-}: (_{-} \in []). \mathcal{I})
any_parse s it its (x :: xs) :=
     case parse' (take s \in it), parse' (drop s \in its), any parse s it its xs of
       inl ret_1, inl ret_2, \_ \rightarrow inl (cons_success <math>\_ ret_1 ret_2)
                      , inl ret' 
ightarrow inl ret'
                , \operatorname{ret}_2 , \operatorname{inr}\operatorname{ret}' 	o \operatorname{inr} _
where the hole on the last line constructs a proof of
\forall \ x' \in (x :: xs), \ T_{\text{failure}} \ \_ \ \_ \ (\overline{\text{take}_{x'}} \ s \in it) + T_{\text{failure}} \ \_ \ \_ \ (\overline{\text{drop}_{x'}} \ s \in its)
by using ret' directly when x' \in xs, and using whichever one of ret<sub>1</sub> and ret<sub>2</sub> is on
the right when x' = x. While straightforward, the use of sum types makes it painfully
verbose without actually adding any insight; we prefer to elide the actual term.
```

Figure 4-2: Pseudo-Implementation of our parser. We take the convention that dependent indices to functions (e.g., unseen) are implicit.

the string is well-typed; we will describe the headaches this seemingly innocuous simplification caused us in Subsection 4.4.2. The only such production in this case is the one that lines up with the production used in the parse tree, labeled (ab)\*(ab)\*.

- 3. We invoke split on our parse tree.
  - (a) The split that we defined then invokes deloop on the two copies of the parse tree

$$\frac{\overline{T}}{\text{"ab"} \in (ab)^*}$$

Since there are non-root nodes labeled with ("ab"  $\in$  (ab)\*), the label of the root node, we promote the deepest one. Letting T' denote the tree

$$\frac{\boxed{"a" \in 'a'} \qquad \boxed{"b" \in 'b'}}{"a" \cdot "b" \in (ab)^*} ("ab")$$

the result of calling deloop is the tree

$$\frac{T'}{\text{"ab"} \in (ab)^*}$$

- (b) The return of split is thus the singleton list containing a single pair of two parse trees; each element of the pair is the parse tree for "ab" ∈ (ab)\* that was returned by deloop.
- 4. We invoke parse on each of the items in the sequence of items associated to (ab)\* via the rule ((ab)\*(ab)\*). The two items are identical, and their associated elements of the pair returned by split are identical, so we only describe the execution once, on

$$\frac{T'}{\text{"ab"} \in (ab)^*}$$

- (a) We first ensure that we are not in an infinite loop. We check if  $s < s_0$ . This check succeeds, for "ab" is shorter than "abab". We thus reset unseen and  $s_0$ ; calls made by this stack frame will pass  $unseen_0 \equiv [(ab)^*]$  for unseen, and  $s \equiv "ab"$  for  $s_0$ .
- (b) We may again consider only the productions for which the parse tree associated to the string is well-typed. The only such production in this case is the one that lines up with the production used in the parse tree T', labeled ("ab").
- (c) We invoke split on our parse tree.
  - i. The split that we defined then invokes deloop on the trees  $\overline{\ } a'' \in \ 'a''$  and  $\overline{\ } b'' \in \ 'b'$ . Since these trees have no non-root nodes (let alone non-root nodes sharing a label with the root), deloop is a no-op.

- ii. The return of split is thus the singleton list containing a single pair of two parse trees; the first is the parse tree  $\overline{\ "a" \in \ "a"}$ , and the second is the parse tree  $\overline{\ "b" \in \ "b"}$ .
- (d) We invoke parse on each of the items in the sequence of items associated to (ab)\* via the rule ("ab"). Since both of these items are terminals, and the relevant equality check (that "a" is equal to "a", and similarly for "b") succeeds, parse returns terminal\_success. We thus have the two MinParseTrees: "a" \in 'a" \in 'a" \in 'b" \in 'b".
- (e) We combine these using cons\_success (and nil\_success, to tie up the base case of the list). We thus have the tree  $\overline{T'_m}$ .
- (f) We apply production\_success< to this tree, and return the tree

$$\frac{T'_m}{\text{"ab"} \in (\text{ab})^*} < (\text{"ab"}, [(\text{ab})^*])$$

5. We now combine the two identical trees returned by parse using cons\_success (and nil\_success, to tie up the base case of the list). We thus have the tree

$$\frac{\overline{T'_m}}{\text{"ab"} \in (ab)^*} < (\text{"ab"}, [(ab)^*]) \qquad \frac{\overline{T'_m}}{\text{"ab"} \in (ab)^*} < (\text{"ab"}, [(ab)^*])$$

$$\frac{\text{"ab"} \cdot \text{"ab"} \in (ab)^*}{\text{"ab"} \cdot \text{"ab"} \in (ab)^*} < (\text{"abab"}, [])$$

6. We apply production\_success<sub>=</sub> to this tree, and return the tree we claimed we would end up with,

$$\frac{\overline{T_m'}}{\text{"ab"} \in (ab)^*} < (\text{"ab"}, [(ab)^*]) \qquad \frac{\overline{T_m'}}{\text{"ab"} \in (ab)^*} < (\text{"ab"}, [(ab)^*]) \\ \frac{\text{"ab"} \cdot \text{"ab"} \in (ab)^*}{\text{"abab"} \in (ab)^*} < (\text{"abab"}, [(ab)^*])$$

## 4.3.3 Parametricity

Before we can combine different instantiations of this interface, we need to know that they behave similarly. Inspection of the code, together with relational parametricity, validates assuming the following axiom, which should also be internally provable by straightforward induction (though we have not bothered to prove it).

The parser extensionality axiom states that, for any fixed instantiation of split, and any arbitrary instantiations of the rest of the interface, giving rise to two different functions parse<sub>1</sub> and parse<sub>2</sub>, we have

where bool\_of\_sum is, for any types A and B, the function of type  $A + B \rightarrow Bool$ 

obtained by sending everything in the left component to true, and everything in the right component to false.

### 4.3.4 Putting It All Together

Now we have parsers returning the following types:

```
\begin{array}{c} {\tt has\_parse} : {\tt Nonterminal} \to {\tt String} \to {\tt Bool} \\ {\tt parse} : ({\tt nt} : {\tt Nonterminal}) \to ({\tt s} : {\tt String}) \\ \to {\tt option} \; ({\tt ParseTreeOf} \; {\tt nt} \; {\tt s}) \\ {\tt has\_no\_parse} : ({\tt nt} : {\tt Nonterminal}) \to ({\tt s} : {\tt String}) \\ \to \top + ({\tt MinParseTreeOf} \; {\tt nt} \; {\tt s} \to \bot) \\ {\tt min\_parse} : ({\tt nt} : {\tt Nonterminal}) \to ({\tt s} : {\tt String}) \\ \to {\tt ParseTreeOf} \; {\tt nt} \; {\tt s} \\ \to {\tt MinParseTreeOf} \; {\tt nt} \; {\tt s} \\ \end{array}
```

Note that we have taken advantage of the isomorphism  $\top + \top \cong \texttt{Bool}$  for has\_parse, the isomorphism  $A + \top \cong \texttt{option}$  A for parse, and the isomorphism  $A + \bot \cong A$  for min\_parse.

We can compose these functions to obtain our desired correct-by-construction parser:

```
\begin{array}{c} {\tt parse\_full}: ({\tt nt}: {\tt Nonterminal}) \rightarrow ({\tt s}: {\tt String}) \\ \qquad \rightarrow {\tt ParseTreeOf} \ {\tt nt} \ {\tt s} + ({\tt ParseTreeOf} \ {\tt nt} \ {\tt s} \rightarrow \bot) \\ {\tt parse\_full} \ {\tt nt} \ {\tt s} \coloneqq \\ \qquad {\tt case} \ \ {\tt parse} \ {\tt nt} \ {\tt s}, \ {\tt has\_no\_parse} \ {\tt nt} \ {\tt s} \ \ {\tt of} \\ \qquad & | \ {\tt Some} \ {\tt d}, \ \_ \qquad \rightarrow \ \ {\tt inl} \ {\tt d} \\ \qquad & | \ \_ \qquad , \ \ {\tt inr} \ {\tt nd} \ \rightarrow \ \ {\tt inr} \ ({\tt nd} \circ {\tt min\_parse}) \\ \qquad & | \ \_ \qquad , \ \_ \qquad \rightarrow \ \ \not {\tt f} \\ \end{array}
```

In the final case, we derive a contradiction by applying the parser extensionality axiom, which says that parse and has\_no\_parse must agree on whether or not s parses as nt.

### 4.4 Missteps, Insights, and Dependently Typed Lessons

We will now take a step back from the parser itself, and briefly talk about the process of coding it. We encountered a few pitfalls that we think highlight some key aspects of dependently typed programming, and our successes suggest benefits to be reaped from using dependent types.

### 4.4.1 The Trouble of Choosing the Right Types

Although we began by attempting to write the type-signature of our parser, we found that trying to write down the correct interface, without any code to implement it, was essentially intractable. Giving your functions dependent types requires performing a nimble balancing act between being uselessly general on the one hand, and too overly specific on the other, all without falling from the highropes of well-typedness onto the unforgiving floor of type errors.

We have found what we believe to be the worst sin the typechecker will let you get away with: having different levels of generality in different parts of your code base, which are supposed to interface with each other without a thoroughly vetted abstraction barrier between them. Like setting your highropes at different tensions, every trip across the interface will be costly, and if the abstraction levels get too far away, recovering your balance will require Herculean effort.

We eventually gave up on writing a dependently typed interface from the start, and decided instead to implement a simply typed Boolean recognizer, together with proofs of soundness and completeness. Once we had in hand these proofs, and the data types required to carry them out, we found that it was mostly straightforward to write down the interface and refine our parser to inhabit its newly generalized type.

### 4.4.2 Misordered Splitters

One of our goals in this presentation was to hide most of the abstraction-level mismatch that ended up in our actual implementation, often through clever use of notation overloading. One of the most significant mismatches we managed to overcome was the way to represent the set of productions. In this paper, we left the type as an abstract mathematical set, allowing us to forgo concerns about ordering, quantification, and occasionally well-typedness.

In our Coq implementation, we fixed the type of productions to be a list very early on, and paid the price when we implemented our parse-tree parser. As mentioned in the execution of the example in Subsection 4.3.2, we wanted to restrict our attention to certain productions, and rule out the other ones using dependent types. This should be possible if we parameterize over not just a splitter, but a production-selector, and only require that our string type be well-typed for productions given by the production-selector. However, the implementation that we currently have requires a well-typed string type for all productions; furthermore, it does not allow the order in which productions are considered to depend on the augmented string data. We paid for this with the extra 300 lines of code we had to write to interleave two different splitters, so that we could handle the cases that we dismissed above as being ill-typed and therefore not necessary to consider. That is, because our types were not formulated in a way that actually made these cases ill-typed, we had to deal with them, much to our displeasure.

#### 4.4.3 Minimal Parse Trees vs. Parallel Traces

Taking another step back, our biggest misstep actually came before we finished the completeness proof for our simply typed Boolean recognizer.

When first constructing the type MinParseTree, we thought of them genuinely as minimal parse trees (ones without a duplicate label in any single path). After much head-banging, of knowledge that a theorem was obviously true, against proof goals that were obviously impossible, we discovered the single biggest insight—albeit a technical one—of the project. The type of "minimal parse trees" we had originally formulated did not match the parse trees produced by our algorithm. A careful examination of the algorithm execution in Subsection 4.3.2 should reveal the difference. Our insight, thus, was to conceptualize the data type as the type of traces of parallel executions of our particular parser, rather than as truly minimal parse trees.

This may be an instance of a more general phenomenon present when programming with dependent types: subtle friction between what you think you are doing and what you are actually doing often manifests as impossible proof goals.

<sup>&</sup>lt;sup>1</sup>For readers wanting to skip that examination: the algorithm we described allows a label  $(s \in nt)$  to appear one extra time along a path if, the first time it appears, its parent node's label,  $(s' \in nt')$ , satisfies s < s'. That is, whenever the string being parsed shrinks, the first nonterminal the shrunken string is parsed as may be duplicated once before shrinking the string again.

### Chapter 5

### Refining Splitters by Fiat

### 5.1 Splitters at a Glance

We have now finished describing the general parsing algorithm, as well as its correctness proofs; we have an algorithm, that decides whether or not a given structure can be imposed on any block of unstructured text. The algorithm is parametrized on an "oracle" that describes how to split the string for each rule; essentially all of the algorithmically interesting content is in the splitters. For the remainder of this paper, we will focus on how to implement the splitting oracle. Correctness is not enough, in general; algorithms also need to be fast to use. We thus focus primarily on efficiency when designing splitting algorithms, and work in a framework that guarantees correctness.

The goals of this work, as mentioned in Section 2.2, are to present a framework for constructing proven-correct parsers incrementally, and argue for its eventual feasibility. To this end, we build on the previous work of Fiat [?], to allow us to build programs incrementally while maintaining correctness guarantees. This section will describe Fiat, and how it is used in this project. The following sections will focus more on the details of the splitting algorithms, and less on Fiat itself.

### 5.2 What counts as efficient?

To guide our implementations, we characterize efficient splitters informally, as follows. Although our eventual concrete efficiency target is to be competitive with extant open source JavaScript parsers, when designing algorithms, we aim at the asymptotic efficiency target of linearity in the length of the string. In practice, the dominating concern is that doubling the length of the string should only double the duration of the parse, and not quadruple it (or more!). [TODO: CITATION NEEDED] To be efficient, it suffices to have the splitter return at most one index. In this case, the parsing time is  $\mathcal{O}(\text{length of string} \times (\text{product over all nonterminals of the number of the string})$ 

possible rules for that nonterminal)).

Here is an example of hitting the worst-case scenario. [TODO: Is this actually possible?]

To avoid hitting this worst-case scenario, we can use a nonterminal-picker, which returns the list of possible production rules for a given string and nonterminal. As long as it returns at most one possible rule in most cases, in constant time, the parsing time will be  $\mathcal{O}(\text{length of string})$ ; backtracking will never happen. This is future work.

### 5.3 Introducing Fiat

### 5.3.1 Incremental Construction by Refinement

Efficiency targets in hand, we move on to incremental construction. The key idea is that parsing rules tend to fall into clumps that are similar between grammars. For example, many grammars use delimiters (such as whitespace, commas, or binary operation symbols) as splitting points, but only between well-balanced brackets (such as double quotes, parentheses, or comment markers). We can take advantage of these similarities by baking the relevant algorithms into basic building blocks, which can then be reused across different grammars. To allow this reuse, we construct the splitters incrementally, allowing us to deal with different rules in different ways.

The Fiat framework [?] is the scaffolding of our splitter implementations. As a framework, the goal of Fiat is to enable library-writers to construct algorithmic building blocks packaged with correctness guarantees, in such a way that users can easily and mostly-automatically make use of these building blocks when they apply.

#### 5.3.2 The Fiat Mindset

The correctness guarantees of Fiat are based on specifications in the form of propositions in Gallina, the mathematical language used by Coq. For example, the specification of a valid has\_parse method is that has\_parse nt str = true  $\longleftrightarrow$  inhabited (ParseTreeOf nt s). Fiat allows incremental construction of algorithms by providing a language for seamlessly mixing specifications and code. The language is a light-weight monadic syntax with one extra operator: a non-deterministic choice operator; we define the following combinators:

```
x \leftarrow c; c' Run c and store the result in x; continue with c', which may mention x
```

c;; c' Run c. If it terminates, throw away the result, and run c'

ret x A computation that immediately returns the value x

 $\{x \mid P(x)\}$  Nondeterministically choose a value of x satisfying P.

If none exists, the program is considered to not terminate.

An algorithm starts out as a nondeterministic choice of a value satisfying the specification. Coding then proceeds by refinement. Formally, we say that a computation c' refines a computation c, written  $c' \subseteq c$ , if every value that c' can compute to, c' can also compute to. We freely generate the relation "the computation c' can compute to the value c'", written c' v, by the rules:

In our use case, we express the specification of the splitter as a nondeterministic choice of a list of split locations, such that any splitting location that results in a valid parse tree is contained in the list. More formally, we can define the proposition

```
\begin{array}{l} {\rm split\_list\_is\_complete} : {\rm Grammar} \ \to \ {\rm String} \ \to \ [{\rm Item}] \ \to \ {\rm Prop} \\ {\rm split\_list\_is\_complete} \ {\rm G} \ {\rm str} \ [] \ {\rm splits} \ = \ \bot \\ {\rm split\_list\_is\_complete} \ {\rm G} \ {\rm str} \ ({\rm it} :: {\rm its}) \ {\rm splits} \\ {\rm =} \ \forall \ {\rm n}, \ {\rm n} < {\rm length} \ {\rm str} \\ {\rm \to} \ ({\rm has\_parse} \ {\rm it} \ ({\rm take} \ {\rm n} \ {\rm str}) \ \land \ {\rm has\_parse} \ {\rm its} \ ({\rm drop} \ {\rm n} \ {\rm str})) \\ {\rm \to} \ {\rm n} \in {\rm splits} \end{array}
```

where we overload has\_parse to apply to items and productions alike. In practice, we pass the first item, and the rest of the items, as separate arguments, so that we don't have to deal with the empty list case.

Let production\_is\_reachable G p to be the proposition that p could show up during parsing, i.e., that p is a tail of a rule associated to some valid nonterminal in the grammar; we define this by folding over the list of valid nonterminals. The specification of the splitter, as a nondeterministic computation, for a given grammar G, a given string str, and a given rule it::its, is then:

```
{ splits : list N
| production_is_reachable G (it :: its)
  → split_list_is_complete G str it its splits }
```

We then refine this into a choice of a splitting location for each rule actually in the grammar (checking for equality with the given rule), and then can refine (implement) the splitter for each rule separately. For example, for the grammar (ab)\*, defined to have a single nonterminal (ab)\* which can either be empty, or be mapped to

'a' 'b' (ab)\*, we would refine this to the computation:

```
If [(ab)^*] =_p it :: its Then { splits : list \mathbb{N} | split_list_is_complete G str (ab)^* [] splits } Else If <math>['b', (ab)^*] =_p it :: its Then { splits : list \mathbb{N} | split_list_is_complete G str 'b' [(ab)^*] splits } Else If <math>['a', 'b', (ab)^*] =_p it :: its Then { splits : list \mathbb{N} | split_list_is_complete G str 'a' ['b', (ab)^*] splits } Else { <math>dummy\_splits : list \mathbb{N} \mid T }
```

where  $=_p$  refers to a boolean equality test for productions. Note that in the final case, we permit any list to be picked, because whenever the production we are handling is reachable, it will never be used.

We can now refine each of these cases separately, using setoid\_rewrite. The key to doing this, to making Fiat work, is that the refinement rules package their correctness properties, so users don't have to worry about correctness when programming by refinement. We use Coq's setoid rewriting machinery to automatically glue together the various correctness proofs when refining only a part of a program.

For example, we might have a lemma singleton which says that returning the length of the string is a valid refinement for any rule that has only one nonterminal; it's type, for a particular grammar G, a particular string str, and a particular nonterminal nt, would be

```
Then setoid_rewrite (singleton \_ _ (ab)*) would refine
         If [(ab)^*] =_p it :: its Then
               \{ \text{ splits } : \text{ list } \mathbb{N} 
               | split_list_is_complete G str (ab)* [] splits }
         Else If ['b', (ab)*] =<sub>p</sub> it :: its Then
               \{ \text{ splits } : \text{ list } \mathbb{N} 
               | split_list_is_complete G str 'b' [(ab)*] splits }
         Else If ['a', 'b', (ab)*] =_p it :: its Then
               \{ \text{ splits } : \text{ list } \mathbb{N} 
               | split_list_is_complete G str 'a' ['b', (ab)*] splits }
         Else
               { dummy_splits : list \mathbb{N} \mid \top }
into
          If [(ab)^*] =_p it :: its Then
               ret [length str]
         Else If ['b', (ab)*] =_p it :: its Then
               \{ \text{ splits } : \text{ list } \mathbb{N} 
               | split_list_is_complete G str 'b' [(ab)*] splits }
         Else If ['a', 'b', (ab)*] =_p it :: its Then
               \{ \text{ splits : list } \mathbb{N} \}
               | split_list_is_complete G str 'a' ['b', (ab)*] splits }
         Else
               \{ dummy\_splits : list \mathbb{N} \mid \top \}
```

We now describe the refinements that we do within this framework, to implement efficient splitters.

### 5.4 Optimizations

## 5.4.1 An Easy First Optimization: Indexed Representation of Strings

One optimization that is always possible is to represent the current string being parsed in this recursive call as a pair of indices into the original string. This allows us to optimize the code doing string manipulation, as it will no longer need to copy strings around, only do index arithmetic.

### 5.4.2 Upcoming Optimizations

In the next few sections, we build up various strategies for splitters. Although our eventual target is JavaScript, we cover only a more modest target of very simple arithmetical expressions in this paper. We begin by tying up the (ab)\* grammar, and then moving on to parse numbers, parenthesized numbers, expressions with only numbers and '+', and then expressions with numbers, '+' and parentheses.

[TODO: The following is mostly redundant with what was said above; remove most of it.] [QUESTION FOR ADAM: Is there any part here worth keeping, that isn't taken care of by the above.] For each grammar, the Fiat framework presents us with goals describing the unimplemented portion of the splitter for this particular grammar. For example, the goal for the grammar that parses arithmetic expressions involving plusses and parentheses, after taking care of the trivial obligations that we describe in the next chapter, looks like this:

```
1 focused subgoals (unfocused: 3)
, subgoal 1 (ID 3491)
r_n : string * (nat * nat)
n : item ascii * production ascii
H := ?88 : hiddenT
refine
   (ls <- If ([NonTerminal "pexpr"] =p fst n :: snd n)</pre>
             || (([NonTerminal "expr"] =p fst n :: snd n)
             || ([NonTerminal "number"] =p fst n :: snd n)
             || (([Terminal ")"] =p fst n :: snd n)
             || ([Terminal "0"] =p fst n :: snd n)
             || ([Terminal "1"] =p fst n :: snd n)
             || ([Terminal "2"] =p fst n :: snd n)
             || ([Terminal "3"] =p fst n :: snd n)
             || ([Terminal "4"] =p fst n :: snd n)
             || ([Terminal "5"] =p fst n :: snd n)
             || ([Terminal "6"] =p fst n :: snd n)
             || ([Terminal "7"] =p fst n :: snd n)
             || ([Terminal "8"] =p fst n :: snd n)
             || ([Terminal "9"] =p fst n :: snd n)
             || ([NonTerminal "digit"] =p fst n :: snd n)))
          Then ret [ilength r_n]
          Else (If [NonTerminal "pexpr"; Terminal "+";
                   NonTerminal "expr"] =p fst n :: snd n
          Then {splits : list nat
               | ParserInterface.split_list_is_complete
                 plus_expr_grammar
```

```
(string_of_indexed r_n)
                  (NonTerminal "pexpr")
                  [Terminal "+"; NonTerminal "expr"]
                  splits}
          Else
            ret [If ([Terminal "+"; NonTerminal "expr"]
                         =p fst n :: snd n)
                     || ([Terminal "("; NonTerminal "expr"; Terminal ")"]
                         =p fst n :: snd n)
                     || ([NonTerminal "digit"; NonTerminal "number"]
                         =p fst n :: snd n)
                  Then 1
                  Else (If [NonTerminal "expr"; Terminal ")"]
                               =p fst n :: snd n
                  Then pred (ilength r_n)
                  Else (If [NonTerminal "number"]
                              =p fst n :: snd n
                  Then ilength r_n
                  Else 0))]);
  ret (r_n, ls)) (H r_n n)
The important part of this goal is:
{splits : list nat
| ParserInterface.split_list_is_complete
    plus_expr_grammar
    (string_of_indexed r_n)
    (NonTerminal "pexpr") [Terminal "+"; NonTerminal "expr"]
    splits}
It says that we have to find a refinement of non-deterministically picking a complete
list of splits for the rule "a pexpr, followed by a '+', followed by a expr."
To get to this goal, we write this code:
Lemma ComputationalSplitter'
: FullySharpened (string_spec plus_paren_expr_grammar).
Proof.
  start honing parser using indexed representation.
 hone method "splits".
  {
    simplify parser splitter.
```

The first line of the proof takes care of the trivial rules, and hone method "splits"

says that we want to refine the splitter. The tactic simplify parser splitter performs a number of straightforward and simple optimizations, such as combining nested if statements which return the same value.

We begin the next section by describing the strategies for the "trivial" rules.

### Chapter 6

# Fixed Length Items, Parsing (ab)\*; Parsing #s; Parsing #, ()

In this chapter, we explore the Fiat framework with a few example grammars, which we describe how to parse. Because these rules are so straightforward, they can be handled automatically, in the very first step of the derivation; we will explain how this works, too.

Recall the grammar for the regular expression (ab)\*:

$$(ab)^* := \epsilon \mid 'a' \mid b' \mid (ab)^*$$

In addition to parsing this grammar, we will also be able to parse the grammar for non-negative parenthesized integers:

```
\begin{array}{l} {\tt pexpr} ::= \ '\ (\ '\ {\tt pexpr}\ '\ )\ '\ |\ {\tt number} \\ \\ {\tt number} ::= \ {\tt digit}\ {\tt number} \\ \\ {\tt number} ::= \ {\tt \epsilon}\ |\ {\tt number} \\ \\ {\tt digit} ::= \ '0'\ |\ '1'\ |\ '2'\ |\ '3'\ |\ '4'\ |\ '5'\ |\ '6'\ |\ '7'\ |\ '8'\ |\ '9' \end{array}
```

### 6.1 Parsing (ab)\*: At Most One Nonterminal

The simpler of these grammars is the one for (ab)\*. The idea is that if any rule has at most one nonterminal, then there is only one possible split: we assign one character to each terminal, and the remaining characters to the single nonterminal.

For any given rule, we can compute straightforwardly whether or not this is the case; Haskell-like pseudocode for doing so is:

```
has-at-most-one-nt : [Item] -> Bool
has-at-most-one-nt [] = true
has-at-most-one-nt (Terminal _)::xs = has-at-most-one-nt xs
has-at-most-one-nt (NonTerminal _)::xs = has-only-terminals xs
has-only-terminals : [Item] -> Bool
has-only-terminals [] = true
has-only-terminals (Terminal _)::xs = has-only-terminals xs
has-only-terminals (NonTerminal _)::xs = false
```

The code for determining the split location is even easier: if the first item of the rule is a terminal, then split at character 1; if the first item of the rule is a nonterminal, and there are n remaining items in the rule, then split n characters before the end of the string.

### 6.2 Parsing Parenthesized Numbers: Fixed Lengths

The grammar for parenthsized numbers has only one rule with multiple nonterminals: the rule for number := digit number?. The strategy here is also simple: because digit only accepts strings of length exactly 1, we always want to split after the first character.

The following pseudocode computes the determines whether or not all strings parsed by a given item are a fixed length, and, if so, what that length is: [TODO: Colorize code; use macros to ensure it matches the conventions of the rest of the document]

We have proven that for any nonterminal for which this method returns just k, the only valid split of the any string for this rule is at location k. This is the correctness

obligation that Fiat demands of us to be able to use this rule.

### 6.3 Putting It Together

Both of these rules are simple and complete for the rules they handle; if a rule has at most one nonterminal, or if the first element of a rule has a fixed length, then we can't do any better than these rules. Therefore, we incorperate them into the initial invocation of start honing parser using indexed representation. To do this, we express the splitter by folding if statements over all of the rules of the grammar that are reachable from valid nonterminals. The if statements check equality of the rule against the one we were given, and, if they match, looks to see if either of these strategies applies. If either does, than we return the appropriate singleton value. If neither applies, then we leave over a non-deterministic pick of a list containing all possible valid splits. The results of applying this procedure without treating any rules specially was shown in Subsection 5.3.2. The results of applying this procedure, including the rules of this chapter, was shown in Subsection 5.4.2.

[QUESTION FOR ADAM: Is it worth including more code? I don't actually write out the refinement rules in the code I use, since I never need to setoid\_rewrite with them; it's all baked into the initial first-step refinement rule. The general code for defining the computation, in case it's worth including is:

```
Definition expanded_fallback_list'
           (P : String -> item Ascii.ascii -> production Ascii.ascii -> item Ascii.as
           (it : item Ascii.ascii) (its : production Ascii.ascii)
           (dummy : list nat)
: Comp (T * list nat)
  := (ls <- (forall_reachable_productions
               (fun p else_case
                => If production_beq ascii_beq p (it::its)
                   Then (match p return Comp (list nat) with
                            | nil => ret dummy
                            | _::nil => ret [ilength s]
                            | (Terminal _):: _ :: _ => ret [1]
                            | (NonTerminal nt):: p'
                             => If has_only_terminals p'
                                   ret [(ilength s - Datatypes.length p')%natr]
                                Else
                                 (option_rect
                                    (fun _ => Comp (list nat))
                                    (fun (n : nat) => ret [n])
```

### Chapter 7

### Disjoint Items, Parsing #, +

Consider now the following grammar for arithmetic expressions involving '+' and numbers:

```
\begin{array}{l} {\rm expr} \coloneqq {\rm number} + {\rm expr}? \\ {\rm +expr}? \coloneqq \epsilon \mid \ '+\ ' \ {\rm expr} \\ {\rm number} \coloneqq {\rm digit} \ {\rm number}? \\ {\rm number}? \coloneqq \epsilon \mid {\rm number} \\ {\rm digit} \coloneqq \ '0' \mid \ '1' \mid \ '2' \mid \ '3' \mid \ '4' \mid \ '5' \mid \ '6' \mid \ '7' \mid \ '8' \mid \ '9' \end{array}
```

The only rule not handled by the strategies of the previous chapter is the rule expr ::= number +expr?. We can handle this rule by noticing that the set of characters in the strings accepted by number is disjoint from the set of possible first characters of the strings accepted by +expr?. Namely, all characters in strings accepted by number are digits, while the first character of a string accepted by +expr? can only be '+'.

The following code computes the set of possible characters of a rule:

```
else possible-terminals valid_nonterminals xs
```

```
possible-terminals : Grammar -> [Item] -> [Char]
possible-terminals G its = possible-terminals, (valid_nonterminals_of G) its
```

In the case where the nonterminal is not in the list of valid nonterminals, we assume that we have already seen that nonterminal higher up the chain of recursion (which we will have, if it is valid according to the initial list), and thus don't have to recompute its possible terminals.

The following code computes the set of possible first characters of a rule:

```
possible-first-terminals' : [String] -> [Item] -> [Char]
possible-first-terminals' _ [] = []
possible-first-terminals' valid_nonterminals (Terminal ch :: xs)
  = [ch]
possible-first-terminals' valid_nonterminals (NonTerminal nt :: xs)
  = (if nt 'in' valid_nonterminals
     then fold
            union
            П
            (map (possible-first-terminals' (remove nt valid_nonterminals))
                 (Lookup nt))
     else [])
    'union'
    (if has_parse nt ""
     then possible-first-terminals' valid_nonterminals xs
     else [])
possible-first-terminals : Grammar -> [Item] -> [Char]
possible-first-terminals G its = possible-first-terminals' (valid_nonterminals_of G)
```

We can decide has\_parse at refinement-time with the brute-force parser, which tries every split; when the string we're parsing is empty,  $\mathcal{O}(\text{length!})$  is not that long. The idea is that we recursively look at the first element of each production; if it is a terminal, then that is the only possible first terminal of that production. If it's a nonterminal, then we have to fold the recursive call over the alternatives. Additionally, if the nonterminal might end up parsing the empty string, then we have to also move on to the next item in the production, and see what its first characters might be.

By computing whether or not these two lists are disjoint, we can decide whether or not this rule applies. When it applies, we can either look for the first character not in the first list (in this example, the list of digits), or we can look for the first character which is in the second list (in this case, the '+'). Since there are two alternatives,

we leave it up to the user to decide which one to use.

For this grammar, we choose the shorter list. We define a function:

```
index\_of\_first\_character\_in : String \to [Char] \to \mathbb{N} by folding over the string. We can then phrase the refinement rule as having type is\_disjoint \ (possible\_terminals \ G \ [it]) \ (possible\_first\_terminals \ G \ its) = true \\ \to ret \ [index\_of\_first\_character\_in \ str \ (possible\_first\_terminals \ G \ its)] \\ \subseteq \\ \{ \ splits : list \ \mathbb{N} \\ | \ split\_list\_is\_complete \ G \ str \ it \ its \ splits \ \}
```

### Chapter 8

## Parsing well-parenthesized expressions

### 8.1 At a Glance

We finally get to a grammar that requires a non-trivial splitting strategy. In this section, we describe how to parse strings for a grammar that accepts arithmetical expressions involving numbers, pluses, and well-balanced parentheses. More generally, this strategy handles any binary operation with guarded brackets.

### 8.2 Grammars we can parse

Consider the following two grammars, with digit denoting the nonterminal that accepts any single decimal digit.

Parenthesized addition:

```
\begin{array}{l} \exp r ::= \operatorname{pexpr} + \exp r \\ + \exp r ::= \epsilon \mid '+' \exp r \\ \operatorname{pexpr} ::= \operatorname{number} \mid '(' \exp r ')' \\ \operatorname{number} ::= \operatorname{digit} \operatorname{number} ? \\ \operatorname{number} ::= \epsilon \mid \operatorname{number} \\ \operatorname{digit} ::= '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9' \end{array}
```

We have carefully constructed this grammar so that the first character of the string suffices to uniquely determine which rule of any given nonterminal to apply.

S-expressions are a notation for nested space-separated lists. By replacing digit with

a nonterminal that accepts any symbol in a given set, which must not contain either of the brackets, nor whitespace, and replacing '+' with a space character ' ', we get a grammar for S-expressions:

```
\begin{array}{c} \operatorname{expr} ::= \operatorname{pexpr} \operatorname{sexpr} \\ \operatorname{sexpr} ::= \epsilon \mid \operatorname{whitespace} \operatorname{expr} \\ \operatorname{pexpr} ::= \operatorname{atom} \mid '(' \operatorname{expr}')' \\ \operatorname{atom} ::= \operatorname{symbol} \operatorname{atom} ? \\ \operatorname{atom} ::= \epsilon \mid \operatorname{atom} \\ \operatorname{whitespace} ::= \operatorname{whitespace-char} \operatorname{whitespace} ? \\ \operatorname{whitespace} ::= \epsilon \mid \operatorname{whitespace} \\ \operatorname{whitespace-char} ::= ' \cdot ' \mid ' \setminus \operatorname{n'} \mid ' \setminus \operatorname{t'} \mid ' \setminus \operatorname{r'} \end{cases}
```

### 8.3 The Splitting Strategy

#### 8.3.1 The Main Idea

The only rule not already handled is the rule that says that a pexpr +expr is an expr. The key insight here is that, to know where to split, we need to know where the next '+' at the current level of parenthesization is. If we can compute an appropriate lookup table in time linear in the length of the string, then our splitter overall with be linear.

### 8.3.2 Building the Lookup Table

We build the table by reading the string from right to left, storing for each character the location of the next '+' at the current level of parenthesization. To compute this location we keep a list of the location of next '+' at every level of parenthesization.

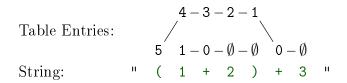
[QUESTION FOR ADAM: Should this example be set off with a \subsubsection or a bold Example. or something?] Let's start with a very simple example, before moving to a more interesting one. To parse "4+5", we are primarily interested in the case where we are parsing something that is a number, or parenthesized on the left, followed by a '+', followed by any expression. For this item, we want to split the string right before the '+', and say that the "4" can be parsed as a number (or parenthesized expression), and that the "+5" can be parsed as a '+' followed by an expression.

To do this, we build a table that keeps track of the location of the next '+', starting from the right of the string. We will end up with the table:

```
Table Entries: 1 0 \emptyset
String: " 4 + 5 "
```

At the '5', there is no next '+', and we are not parenthesized at all, so we record this as  $\emptyset$ . At the '+', we record that there is a '+' at the current level of parenthesization, with 0. Then, since the '4' is not a '+', we increment the previous number, and store 1. This is the correct location to split the string: we parse 1 character as a number, and the rest as +expr.

Now let's try a more complicated example: "(1+2)+3". We want to split this string into "(1+2)" and "+3". The above strategy is insufficient to do this; we need to keep track of the next '+' at all levels of parenthesization at each point. We will end up with the following table, where the bottom-most element is the current level, and the ones above that are higher levels. We use lines to indicate chains of numbers at the same level of parenthesization.



We again start from the right. Since there are no '+'s that we have seen, we store the singleton list  $[\emptyset]$ , indicating that we know about only the current level of parenthesization, and that there is no '+' to the right. At the '+' before the "3", we store the singleton list [0], indicating that the current character is a '+', and we only know about one level of parenthesization. At the ')', we increment the counter for '+', but we also now know about a new level of parenthesization. So we store the two element list  $[\emptyset, 1]$ . At the '3', we increment all numbers, storing  $[\emptyset, 2]$ . At the '+' before the "2", we store 0 at the current level, and increment higher levels, storing [0, 3]. At the '1', we simply increment all numbers, storing [1, 4]. Finally, at the '(', we pop a level of parenthesization, and increment the remaining numbers, storing [5]. This is correct; we have 5 characters in the first string, and when we go to split "1+2" into "1" and "+2", we have the list [1, 4], and the first string does indeed contain 1 character.

As an optimization, we can drop all but the first element of each list once we're done computing, and, in fact, can do this as we construct the table. However, for correctness, it is easier to reason about a list located at each location.

#### 8.3.3 Table Correctness

What is the correctness condition on this table? The correctness condition Fiat gives us is that the splits we compute must be the only ones that give rise to valid parses. This is hard to reason about directly, so we use an intermediate correctness condition: for any cell of the table [QUESTION FOR ADAM: IS MEANING OF CELL OBVIOUS?], if it is empty (is  $\emptyset$ , or does not exist at all), then there must not be a well-parenthesized fragment of the string starting at that point and ending with a '+', which closes the appropriate number of parentheses for this level of

parenthesization. [TODO: FIXME, THIS PREVIOUS SENTENCE SEEMS LIKE A CONFUSING EXPLANATION] [QUESTION FOR ADAM: Any suggestions?] If the cell points to a given location, then that location must contain a '+', and the fragment of the string starting at the current location and going up to but not including the '+', must not contain any '+'s which are "exposed". [QUESTION FOR ADAM: THIS ALSO SEEMS CONFUSING. Ideas? Should I put code/math first?]

More formally, we can define a notation of paren-balanced and paren-balanced-hiding'+'. Say that a string is paren-balanced at level n if it closes exactly n more parentheses than it opens, and there is no point at which it has closed more than n more
parentheses than it has opened. So the string "1+2)" is paren-balanced at level 1
(because it closes 1 more parenthesis than it opens), and the string "1+2)+(3" is
not paren balanced at any level (though the string "1+2)+(3)" is paren-balanced
at level 1). A string is paren-balanced-hiding-'+' at level n if it is paren-balanced
at level n, and, at any point at which there is a '+', at most n-1 more parentheses have been closed than opened. So "(1+2)" is paren-balanced-hiding-'+' at
level 0, and "(1+2))" is paren-balanced-hiding-'+' at level 1, and "(1+2)+3" is
not paren-balanced-hiding-'+' at any level, though it is paren-balanced at level 0.

[QUESTION FOR ADAM: THIS IS VERBOSE. MAYBE I SHOULD START
WITH CODE?]

Then, the formal correctness condition is that if a cell at parenthesis level n points to a location  $\ell$ , then the string from the cell up to but not including  $\ell$  must be parenbalanced-hiding-'+' at level n, and the character at location  $\ell$  must be a '+'. If the cell is empty, then the string up to but not including any subsequent '+' must not be paren-balanced at level n.

The table computed by the algorithm given above satisfies this correctness condition, and this correctness condition implies that the splitting locations given by the table are the only ones that produce valid parse trees; there is a unique table satisfying this correctness condition (because it picks out the *first* '+' at the relevant level), and any split location which is not appropriately paren-balanced/paren-balanced-hiding results in no parse tree. [QUESTION FOR ADAM: THIS SEEMS REPETITIVE?]

### 8.3.4 Diving into Refinement Code

**TODO:** Check the following assumptions: that the reader has seen:

- what Fiat goals look like
- what the initial representation change to the indexed representation looks like
- the general structure of remaining splitter goals

[TODO: Remove me. GOALS: rule is non-trivial, but use should be trivial. pieces to making it trivial: no representation change, automatic discovery of binary operation and brackets, reflective side-conditions]

Although the rule itself is non-trivial, our goal is to make using this rule as trivial as possible; we now explain how this refinement rule requires only one line of code to use (modulo long tactic names):

```
setoid_rewrite refine_binop_table;
[ presimpl_after_refine_binop_table | reflexivity.. ].
```

[QUESTION FOR ADAM: How much explanation at this point? Mention that presimpl... only unfolds things?]

There are three components to making a rule that can be used with a one-liner: not requiring a change of representation; reflective computation of the side-conditions, allowing them all to be discharged with reflexivity; and automatic discovery of any arguments to the rule. We cover each of these separately.

#### Doing Without a New Representation

Recall from [TODO: secref] that the first step in any parser/splitter [TODO: decide on which one to use] is to use an indexed representation of strings, where splitting a string only involves updating the indices into the original string, as well as taking care of all of the straightforward grammar rules. Naively, we'd need to store a fresh table every time we split the string.

How do we get around computing a new table on splits? We pull the same trick here that we pulled on strings; we refer only to the table that is built from the initial string, and describe the correctness condition on the table in terms of arbitrary substrings of that string.

Fiat presents us with a goal involving a statement of the form "non-deterministically pick a list of splitting locations that is complete for the substring of str starting at n and ending at m, for the rule pexpr+expr." In code form, this is:

```
{splits : list nat
| split_list_is_complete
    G
    (substring n m str)
    (NonTerminal nt)
    (Terminal ch :: its)
    splits}
```

[QUESTION FOR ADAM: How much should I say about this code? For example, the current Coq code doesn't line up perfectly with the explanation; the code only handles the case where the rule is "a nonterminal, followed by a terminal, followed by other things", not the form presented above, where the nonterminal is hidden by one layer of indirection (to prevent backtracking).]

The final refinement rule, which we use with setoid\_rewrite, says that this is refined by a lookup into a suitably defined table: [TODO: Standardize notation for goals, or drop code]

By phrasing the rule in terms of substring n m str, rather than in terms of an arbitrary string, the computation of the table is the same in every call to the splitter. All that remains is to lift the computation outside of the recursive call to the parsing function during the extraction phase, which we plan to do soon. [TODO: Section reference future work. Add bit to future work.]

Before moving on to the other components of making usage of this rule require only one line of code, we note that we make use of the essential property that removing characters from the end of the string doesn't add new locations where splitting could yield a valid parse; if a given location is the first '+' at the current level of parenthesization, this does not change when we remove characters from the end of the string.

[QUESTION FOR ADAM: Should there be more content here?]

#### Discharging the Side-Conditions Trivially

To prove the correctness condition on the table, we need to know some things about the grammar that we are handling. In particular, we need to know that if we are trying to parse a string as a rule analogous to pexpr, then there can be no exposed '+' characters, and, furthermore, that every such parseable string has well-balanced parentheses. To allow discharging the side conditions trivially, we define a function that computes whether or not this is the case for a given nonterminal in a given grammar. We then prove that, whenever this function returns true, valid tables yield complete lists of splitting locations.

To make things simple, we approximate which grammars are valid; we require that every open parenthesis be closed in the same rule (rather than deeply nested in further nonterminals). In Haskell-like pseudocode, the function we use to check validity of a grammar can be written as:

```
grammar-and-nonterminal-is-valid : Grammar -> NonTerminal -> Bool
grammar-and-nonterminal-is-valid G nt
  = fold (&&) paren-balanced-hiding G (G.(Lookup) nt)
pb': Grammar -> Char -> Nat -> [Item] -> Bool
pb'Gn[]
                   = (n == 0)
pb' G n (NonTerminal nt :: s)
  = fold (&&) (pb' G 0) (G.(Lookup) nt) && pb' G n s
pb' G n ('(' :: s) = pb' G (n + 1) s
pb' G n (')' :: s) = n > 0 \&\& pb' G (n - 1) s
pb' G n (\_ :: s) = pb' G n s
paren-balanced G = pb' G 0
pbh': Grammar -> Char -> Nat -> [Item] -> Bool
                    = (n == 0)
pbh' G n []
pbh' G n (NonTerminal nt :: s)
  = fold (pbh' G 0) (G.(Lookup) nt) && pb' G n s
pbh' G n ('+' :: s) = n > 0 \&\& pbh' G n s
pbh' G n ('(' :: s) = pbh' G (n + 1) s
pbh' G n (')' :: s) = n > 0 \&\& pbh' G (n - 1) s
pbh' G n (\_ :: s) = pbh' G n s
paren-balanced-hiding G = pbh, G 0
```

#### **QUESTION FOR ADAM:** Does this code need more explanation?

There is one subtlety here, that was swept under the rug in the above code: this computation might not terminate! We could deal with this by memoizing this computation in much the same way that we memoized the parser to deal with potentially infinite parse trees. Rather than dealing with the subtleties of figuring out what to do when we hit repeated nonterminals, we perform the computation in two steps. First, we trace the algorithm above, building up a list of which nonterminals need to be paren-balanced, and which ones need to be paren-balanced-hiding. Second, we fold the computation over these lists, replacing the recursive calls for nonterminals with list membership tests. [QUESTION FOR ADAM: Does this need code?]

#### Automatic Discovery of Arguments

Throughout this chapter, we've been focusing on the arithmetic expression example. However, the exact same rule can handle S-expressions, with just a bit of generalization. There are two things to be computed: the binary operation, and the parenthesis characters.<sup>1</sup>

We require that the first character of any string parsed by the nonterminal analogous to +expr be uniquely determined; that character will be the binary operator; we can reuse the code from Chapter 7 to compute this character and ensure that it is unique.

To find the parenthesis characters, we first compute a list of candidate character pairs: for each rule associated to the nonterminal analogous to pexpr, we consider the pair of the first character and the last character (assuming both are terminals) to be a candidate pair.<sup>2</sup> We then filter the list for characters which yield the conclusion that this rule is applicable to the grammar, according to the algorithm of the previous subsubsection. We then require, as a side-condition, that the length of this list be positive.

<sup>&</sup>lt;sup>1</sup>Currently, our code requires the binary operation to be exposed as a terminal in the rule we are handling. We plan on generalizing this to the grammars described in this chapter shortly.

<sup>&</sup>lt;sup>2</sup>Again, generalizing this to characters hidden by nested nonterminals should be straightforward.

### Chapter 9

### Future work

- Grammars and efficiency: The eventual target for this demonstration of the framework is the JavaScript grammar, and we aim to be competitive, performancewise, with popular open-source JavaScript implementations. <lookup list of JS implementations> We plan to profile our parser against these on <lookup test suite>
  - Description of anticipated challenges, based on the JavaScript grammar
     clookup JS grammar>
- Generating Parse Trees
  - We plan to eventually generate parse trees, and error messages, rather than just booleans, in the complete pipeline. We have already demonstrated that this requires only small adjustments to the algorithm in the section on the dependently typed parser.
- Validating extraction
  - By adapting <Clement's work>, our parsers will be able to be compiled to verified bedrock/assembly, within Coq

### 9.1 Future work with dependent types

Recall from Chapter 4 that dependent types have allowed us to refine our parsing algorithm to prove its own soundness and completeness.

However, we still have some work left to do to clean up the implementation of the dependently typed version of the parser.

Formal extensionality/parametricity proof To completely finish the formal proof of completeness, as described in this paper, we need to prove the parser extensionality axiom from Subsection 4.3.3. We need to prove that the parser does not make any decisions based on any arguments to its interface other than split, internalizing the obvious parametricity proof. (Alternatively, as mentioned above, we could hope to use an extension of Coq with internalized parametricity [3].)

Even more self-reference We might also consider reusing the same generic parser to generate the extensionality proofs, by instantiating the type families for success and failure with families of propositions saying that all instantiations of the parser, when called with the same parsing problem, always return values that are equivalent when converted to Booleans. A more specialized approach could show just that has\_parse agrees with parse on successes and with has\_no\_parse on failures:

```
\begin{split} &T_{\texttt{success} \ \_ \ \_} \ (s \in nt) \\ &\coloneqq \texttt{has\_parse} \ \texttt{nt} \ s = \texttt{true} \land \texttt{parse} \ \texttt{nt} \ s \neq \texttt{None} \\ &T_{\texttt{failure} \ \_ \ \_} \ (s \in nt) \\ &\coloneqq \texttt{has\_parse} \ \texttt{nt} \ s = \texttt{false} \land \texttt{has\_no\_parse} \neq \texttt{inl} \ () \end{split}
```

Synthesizing dependent types automatically? Although finding sufficiently general (dependent) type signatures was a Herculean task before we finished the completeness proof and discovered the idea of using parallel parse traces, it was mostly straightforward once we had proofs of soundness and completeness of the simply typed parser in hand; most of the issues we faced involving having to figure out how to thread additional hypotheses, which showed up primarily at the very end of the proof, through the entire parsing process. Subsequently instantiating the types was also mostly straightforward, with most of our time and effort being spent writing transformations between nearly identical types that had slightly different hypotheses, e.g., converting a Foo involving strings shorter than  $s_1$  into another analogous Foo, but allowing strings shorter than  $s_2$ , where  $s_1$  is not longer than  $s_2$ . Our experience raises the question of whether it might be possible to automatically infer dependently typed generalizations of an algorithm, which subsume already-completed proofs about it, and perhaps allow additional proofs to be written more easily.

Further generalization Finally, we believe our parser could be generalized even further; the algorithm we have implemented is essentially an algorithm for inhabiting arbitrary inductive type families, subject to some well-foundedness, enumerability, and finiteness restrictions on the arguments to the type family. The interface we described is, conceptually, a composition of this inhabitation algorithm with recursion and inversion principles for the type family we are inhabiting (ParseTreeOf in this paper). Our techniques for refining this algorithm so that it could prove itself sound and complete should therefore generalize to this viewpoint.

### Bibliography

- [1] José Bacelar Almeida, Nelma Moreira, David Pereira, and Simão Melo de Sousa. Partial derivative automata formalized in Coq. In *Proceedings of the 15th International Conference on Implementation and Application of Automata*, CIAA'10, pages 59–68, Berlin, Heidelberg, 2011. Springer-Verlag.
- [2] Aditi Barthwal and Michael Norrish. Verified, executable parsing. In *Proceedings* of the 18th European Symposium on Programming Languages and Systems: Held As Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2009, ESOP '09, pages 160–174, Berlin, Heidelberg, 2009. Springer-Verlag.
- [3] Jean-Philippe Bernardy and Moulin Guilhem. Type-theory in color. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming*, ICFP '13, pages 61–72, New York, NY, USA, 2013. ACM.
- [4] Bryan Ford. Parsing expression grammars: A recognition-based syntactic foundation. In *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '04, pages 111–122, New York, NY, USA, 2004. ACM.
- [5] Graham Hutton. Higher-order functions for parsing. *Journal of Functional Programming*, 2(3):323–343, July 1992.
- [6] Jacques-Henri Jourdan, François Pottier, and Xavier Leroy. Validating LR(1) parsers. In *Proceedings of the 21st European Conference on Programming Languages and Systems*, ESOP'12, pages 397–416, Berlin, Heidelberg, 2012. Springer-Verlag.
- [7] Ramana Kumar, Magnus O. Myreen, Michael Norrish, and Scott Owens. CakeML: A verified implementation of ML. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '14, pages 179–191, New York, NY, USA, 2014. ACM.
- [8] Matthew Might, David Darais, and Daniel Spiewak. Parsing with derivatives: A functional pearl. In *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming*, ICFP '11, pages 189–195, New York, NY, USA, 2011. ACM.

- [9] Greg Morrisett, Gang Tan, Joseph Tassarotti, Jean-Baptiste Tristan, and Edward Gan. RockSalt: better, faster, stronger SFI for the x86. In *Proceedings* of the 33rd ACM SIGPLAN conference on Programming Language Design and Implementation, PLDI '12, pages 395–404, New York, NY, USA, 2012. ACM.
- [10] Magnus O. Myreen and Jared Davis. A verified runtime for a verified theorem prover. In *Proceedings of the Second International Conference on Interactive Theorem Proving*, ITP'11, pages 265–280, Berlin, Heidelberg, 2011. Springer-Verlag.
- [11] Tom Ridge. Simple, functional, sound and complete parsing for all context-free grammars. In *Proceedings of the First International Conference on Certified Programs and Proofs*, CPP'11, pages 103–118, Berlin, Heidelberg, 2011. Springer-Verlag.
- [12] Elizabeth Scott and Adrian Johnstone. GLL parsing. In *Proceedings of the Ninth Workshop on Language Descriptions Tools and Applications*, LDTA '09, pages 177–189, 2009.