[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.

# 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 framework for synthesizing efficient, verified parsers by refinement, and a demonstration of this framework by synthesizing a JavaScript parser competitive with existing open-source parsers. 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

## Acknowledgments

This is the acknowledgements section. You should replace this with your own acknowledgements.

# Contents

1	Par	sing C	ontext Free Grammars	15
	1.1	Parsing		
		1.1.1	Infinite regress	16
		1.1.2	Aborting early	17
	1.2 Standard Formal Definitions			19
		1.2.1	Context Free Grammar	19
		1.2.2	Parse Trees	19
2	Rela	ated W	Vork and Other Approaches to Parsing	21
	2.1	Relate	d Work	21
	2.2	What'	s New	22
3	Con	mpleteness and Soundness		
4	Con	ompleteness, Soundness, and Parsing Parse Trees		
	4.1	Provin	ng Completeness: Conceptual Approach	27
	4.2	Minim	al Parse Trees: Formal Definition	28
	4.3	Parser	Interface	29

		4.3.1	Parsing Parses	33
		4.3.2	Example	35
		4.3.3	Parametricity	38
		4.3.4	Putting it all together	39
	4.4	Misste	ps, Insights, and Dependently Typed Lessons	40
		4.4.1	The trouble of choosing the right types	40
		4.4.2	Misordered splitters	40
		4.4.3	Minimal Parse Trees vs. Parallel Traces	41
-	D «£	: C	mlittang by Diet	49
5	кеп		plitters by Fiat	43
	5.1	Splitte	rs at a Glance	43
	5.2	What	counts as efficient?	43
	5.3	Introd	ucing Fiat	44
		5.3.1	Incremental Construction by Refinement	44
		5.3.2	The Fiat Mindset	44
	5.4	Optimizations		46
		5.4.1	An Easy First Optimization: Indexed Representation of Strings	46
		5.4.2	Upcoming Optimizations	46
6	fixed	d lengt	th nonterminals, parsing (ab)*; parsing #s; parsing #, ()	47
7	disjoint items, parsing $\#$ , $+$			
8	Pars	sing we	ell-parenthesized expressions	61
	8.1	At a C	lance	61

	8.2	Gram	mars we can parse	61
	8.3	3 The Splitting Strategy		
		8.3.1	The Main Idea	62
		8.3.2	Building the Lookup Table	62
		8.3.3	Table Correctness	64
		8.3.4	The Code	65
		8.3.5	The Correctness Proof	65
9	Future work			
	9.1	Future	e work with dependent types	73

# List of Figures

4-1	The dependently typed interface of our parser	30
4-2	Pseudo-Implementation of our parser. We take the convention that dependent indices to functions (e.g., unseen) are implicit	36

# List of Tables

## 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.

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

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  $\c^*c_1' \cdot \c^*c_2' \cdot \s$  as  $ab(ab)^*$ , splits the string into  $\c^*c_1', \c^*c_2', \s$ ; and when trying to parse  $\s^*s$  as  $\s^*s$ , does 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:

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

$$\frac{s_0 \in (\mathsf{ab})^* \quad s_1 \in (\mathsf{ab})^*}{s_0 s_1 \in (\mathsf{ab})^*} ((\texttt{ab})^*(\mathsf{ab})^*)$$

Now, let us try to parse the string "ab" as (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.

```
\begin{array}{l} \text{any\_parses} : [\text{String} \times \text{String}] \to \text{Bool} \\ \\ \text{any\_parses} \ [] \coloneqq \text{false} \\ \\ \text{any\_parses} \ (("a","b") :: \ ) \coloneqq \text{true} \\ \\ \text{any\_parses} \ ((s_1,s_2) :: \text{rest\_splits}) \\ \\ \quad \coloneqq (\text{parses} \ s_1 \ \&\& \ \text{parses} \ s_2) \ | \ | \ \text{any\_parses} \ \text{rest\_splits} \\ \\ \text{parses} : \text{String} \to \text{Bool} \\ \\ \text{parses} \ \text{str} \coloneqq \text{true} \\ \\ \text{parses} \ \text{str} \coloneqq \text{any\_parses} \ (\text{split} \ \text{str}) \\ \end{array}
```

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:

$$\boxed{"a" \in 'a'} \qquad \boxed{"b" \in 'b'}$$

Productions and nonterminals:

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

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

#### 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:

```
(rule): ParseItemsTreeOf its s \rightarrow ParseTreeOf nt 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.

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 fits in

We begin with an overview of the general theory and setting of parsing, with an eye towards verified parsers.

[TODO: Very basic Background on parsers, grammars]

#### 2.1 Related Work

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. 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]

#### 2.2 What's New

#### [TODO: Fill out what's new]

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

Start with naive recursive-descent parsing. We ensure termination via memoization, a la [11]. We parameterize the parser on a "splitting oracle", which describes how to recurse. <section citation> As far as we can tell, the idea of factoring the algorithmic complexity like this is new. <section citation>

We use fiat to incrementally build efficient parsers by refinement. < section citation >

The idea of reusing the parsing algorithm to prove its own completeness, by parsing parse trees rather than strings, is not found in the literature, to the authors'

 $knowledge.\ < section\ citation >$ 

## 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} {\sf has\_parse\_sound} : ({\sf nt} : {\sf Nonterminal}) \to ({\sf s} : {\sf String}) \\ & \to {\sf has\_parse} \ {\sf nt} \ {\sf s} = {\sf true} \\ & \to {\sf ParseTreeOf} \ {\sf nt} \ {\sf s} \\ {\sf has\_parse\_complete} : ({\sf nt} : {\sf Nonterminal}) \to ({\sf s} : {\sf String}) \\ & \to {\sf ParseTreeOf} \ {\sf nt} \ {\sf s} \\ & \to {\sf has\_parse} \ {\sf nt} \ {\sf s} = {\sf 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} & parse: (nt: Nonterminal) \\ & \rightarrow (s: String) \\ & \rightarrow ParseTreeOf \ nt \ s + (ParseTreeOf \ nt \ s \rightarrow \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 mini-

mal 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<sub>'ch'</sub>: MinParseTreeOf s<sub>0</sub> 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 \leq s_0 \\ & \rightarrow \texttt{MinParseTreeOf} \ \ s_0 \ \ unseen \ \ it \ \ s_1 \\ & \rightarrow \texttt{MinItemsTreeOf} \ \ s_0 \ \ unseen \ \ its \ \ s_2 \\ & \rightarrow \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  $\overline{s} \in X^{<(s_0,v)}$  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 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].

```
\mathtt{T_{success}}, \ \mathtt{T_{failure}} : \mathtt{String} 
ightarrow \mathtt{[Nonterminal]} 
ightarrow \mathtt{ParseQuery} 
ightarrow \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 : Nonterminal) \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

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

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

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  $\top$ :

```
T_{\text{success}}
             \coloneqq \top
T_{\text{failure}}
terminal_success
                       := ()
terminal_failure
                           := ()
nil_success
                 = ()
                   := ()
nil_failure
                               := ()
cons success
                           := ()
cons failure
                                     := ()
production_success
                                       := ()
production_success_
production_failure
                                := ()
production_failure_
                                 := ()
production_failure∉
                                 := ()
```

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.

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

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).

$$\begin{array}{ll} \mathtt{T}_{\mathtt{success}} & := \top \\ \\ \mathtt{T}_{\mathtt{failure}} & \mathtt{s}_0 \text{ unseen } (\overline{\mathtt{s} \in \mathtt{X}}) \coloneqq \left(\overline{\mathtt{s} \in \mathtt{X}} \stackrel{<(\mathtt{s}_0,\mathtt{unseen})}{}\right) \to \bot \end{array}$$

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

```
terminal_success
                       := ()
terminal_failure
nil success
               := ()
nil_failure
                                  := ()
cons_success
                              := !
cons_failure
production_success
                                        := ()
production_success=
                                           := ()
production_failure
production_failure_
                                    := \mathcal{I}
production_failure<sub>∉</sub>
```

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:

$$\begin{array}{lll} \mathtt{split} & (\mathtt{s} \in \mathtt{it} :: \mathtt{its}) \coloneqq \\ & \mathbf{case} & \mathtt{parse\_tree\_data} & \mathtt{s} & \mathbf{of} \\ & \left| \frac{\frac{p_1}{s_1 \in \mathtt{it}} & \frac{p_2}{s_2 \in \mathtt{its}}}{s_1 s_2 \in \mathtt{it} :: \mathtt{its}} \right. \to & \left[ (\mathtt{deloop} \ p_1, \mathtt{deloop} \ p_2) \right] \\ & \left| \begin{array}{c} \to & \boldsymbol{\ell} \\ \mathtt{split} & \coloneqq [] \end{array} \right. \end{array}$$

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:

$$\begin{array}{ll} T_{\texttt{success}} & s_0 & \texttt{unseen} & (\texttt{d}: \overline{s \in X}) \coloneqq \overline{s \in X} < (s_0, \texttt{unseen}) \\ T_{\texttt{failure}} & \coloneqq \bot \end{array}$$

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

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

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

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".

We first ensure that we are not in an infinite loop. We check if s < s<sub>0</sub> (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

```
parse nt s := parse' (s_0 := s) (unseen := unseen_0) (s \in nt)
parse' ("ch" \in 'ch') := inl terminal_success
parse' ( \in \cite{ch'}) := inr (terminal_failure 1)
parse' ( ( ( ( ) ) ) = inl nil_success 
parse' (\overline{\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)
            else inr (production_failure
    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 []). f)
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 ret_1 ret_2)
              , , inl ret' 
ightarrow inl ret' , ret_2 , inr ret' 
ightarrow inr
where the hole on the last line constructs a proof of
\forall x' \in (x :: xs), T_{failure} (take_{x'} s \in it) + T_{failure} (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.

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 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

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

- (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

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

- (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  $\overline{\ }^{\parallel}a^{\parallel} \in \overline{\ }^{\perp}a^{\parallel}$ , and the second is the parse tree  $\overline{\ }^{\parallel}b^{\parallel} \in \overline{\ }^{\perp}b^{\parallel}$ .
- (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" and "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

```
∀ (nt: Nonterminal) (s: String),
bool_of_sum (parse₁ nt s) = bool_of_sum (parse₂ nt s)
```

where bool\_of\_sum is, for any types A and B, the function of type  $A + B \to Bool$  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} \\ \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}) \to ({\tt s}: {\tt String}) \\ & \to {\tt ParseTreeOf} \ {\tt nt} \ {\tt s} + ({\tt ParseTreeOf} \ {\tt nt} \ {\tt s} \to \bot) \\ {\tt parse\_full} \ {\tt nt} \ {\tt s} \coloneqq \\ {\tt case} \ \ {\tt parse} \ {\tt nt} \ {\tt s}, \ {\tt has\_no\_parse} \ {\tt nt} \ {\tt s} \ \ {\tt of} \\ & \mid {\tt Some} \ {\tt d}, \qquad \to \ \ {\tt inl} \ {\tt d} \\ & \mid & , \ \ {\tt inr} \ {\tt nd} \ \to \ \ {\tt inr} \ ({\tt nd} \circ {\tt min\_parse}) \\ & \mid & , \qquad \to \ \rlap{\it 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.¹ 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  $\mathcal{O}(\text{length of string} \times (\text{product over all nonterminals of the number of possible rules for that nonterminal}).$ 

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.
```

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. 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, if we are looking to split the string at the location of the first '+' character, we might first pick the list of "valid locations to split at", and then refine that into a computation that picks the location of the first '+', and returns the singleton list containing that, and then replace the remaining bit of nondeterminism with a computation of the location of the first '+'.

The key 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.

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.

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 (ab)\* grammar looks like this: [TODO: <INSERT GOAL>]. To get to this goal, we write this code: [TODO: <COQ CODE HERE, with comments describing what each line does>] We thus have to describe how to split a string for the rules [TODO: <RULES HERE>], and provide proofs that these splitting strategies are complete. We begin the next section with this splitting strategy.

# Chapter 6

# fixed length nonterminals, parsing (ab)\*; parsing #s; parsing #, ()

#### • Goals

- Explore the framework
- Demonstrate that we can implement the "obvious" rules to handle a large swath of CFG rules

#### • At a Glance

- There are some strategies that are obvious enough that it's easy for the computer to decide whether or not it's good to apply them. For example, if a rule starts with a terminal, then we should split off one character, because terminals are always single characters. These rules are enough to parse some simple grammars, such as the regular expression grammar (ab)\*; the grammar accepting numbers; and the grammar accepting parenthesized numbers. [GIVE GRAMMARS HERE] In this section, we explain how we parse these grammars.
- The splitting strategy: if all strings parsed by a given item are the same length, then we can always split the string when faced with that nonterminal.
  - Walk through an example of parsing "abab" as "(ab)\*", describing the splitting at each point.
  - Another example: "((123))"
  - For (ab)\*, the item "a", and the item "b" only parse strings of length 1. Thus when asked for the split for "a", then "b(ab)\*", we can split after one character, and similarly after "b".
  - For numbers with parentheses (of which numbers are a subgrammar), digits always have length 1, so we can split after the first character.

- Implementation as a refinement rule:
  - Describe the obligation Fiat presents us with for the splitter for (#) [CODE HERE]
  - Describe how each rule is handled
  - For the relevant rules, we can compute the length at compile time. Here is the algorithm. <Coq code here>
  - To actually make use of this, we must satisfy the correctness criterion.
     This is what refinement means. We relate the length to the parse trees by a few correctness criteria.
    - \* Note that we need to use only well-founded recursion.
  - We provide a decision procedure for the validity of this rule.
- In practice, we don't actually need to rewrite it; because it's never suboptimal to apply this rule (returning a single split location is just about the best we can do (TODO: handle invalid parses and backtracking better)), so we do it automatically, baking it into the initial goal, along with the indexed representation change (explain)

```
Lemma ComputationalSplitter'
: FullySharpened (string_spec ab_star_grammar).
Proof.
  start honing parser using indexed representation.
 hone method "splits".
  {
    simplify parser splitter.
    finish honing parser method.
  }
  Fully Sharp en Each {\tt MethodWithoutDelegation}.
  extract delegate-free implementation.
  simpl; higher_order_reflexivityT.
Defined.
3 subgoals, subgoal 1 (ID 17)
  Sharpened
     (ADTRep (string * (nat * nat))
     { Def Constructor "new" (s: string) : rep :=
         ret (s, (0, length s)),
```

```
ret (s, string_of_indexed s) ,
        Def Method "is_char" (s : rep, ch : Ascii.ascii) : bool :=
          ret (s, Equality.string_beq (string_of_indexed s) (String ch "")) ,
        Def Method "get" (s : rep, n : nat) : (option Ascii.ascii) :=
          ret (s, get n (string_of_indexed s)) ,
        Def Method "length" (s : rep, _ : ()) : nat :=
          ret (s, ilength s),
        Def Method "take" (s : rep, n : nat) : () :=
          ret (fst s, (fst (snd s), min (snd (snd s)) n), ()),
        Def Method "drop" (s : rep, n : nat) : () :=
          ret (fst s, (n + fst (snd s), NatFacts.minusr (snd (snd s)) n), ()),
        Def Method "splits" (s : rep, p : (item Ascii.ascii *
                                            production Ascii.ascii)) :
        (list nat) :=
          dummy <- {_ : list nat | True};</pre>
          ls <- If [Terminal "a"%char; Terminal "b"%char;</pre>
                   NonTerminal "(ab)*"] =p fst p :: snd p Then
                ret [1]
                Else (If [Terminal "b"%char; NonTerminal "(ab)*"] =p
                         fst p :: snd p Then ret [1]
                      Else (If [NonTerminal "(ab)*"] =p fst p :: snd p
                            Then ret [ilength s] Else
                            ret dummy));
          ret (s, ls),
        Def Method "rules" (s : rep, p : (productions Ascii.ascii)) :
        (list nat) :=
          dummy <- {_ : list nat | True};</pre>
          ls <- If ContextFreeGrammarEquality.productions_beq</pre>
                     Equality.ascii_beq
                     ΓΠ;
                     [Terminal "a"%char; Terminal "b"%char;
                     NonTerminal "(ab)*"]] p
                Then option_rect
                       (fun _ : option Ascii.ascii => Comp (list nat))
                       (fun ch : Ascii.ascii =>
                        If Equality.ascii_beq "a"%char ch Then
                        ret [1] Else ret dummy) (ret [0])
                       (get 0 (string_of_indexed s)) Else
                ret dummy;
          ret (s, ls) })%ADT
subgoal 2 (ID 9) is:
forall idx: StringBound.BoundedString,
refineADT (Sharpened_DelegateSpecs ?4 idx)
```

Def Method "to\_string" (s : rep, \_ : ()) : string :=

```
(Computational ADT. LiftcADT
     (existT (ComputationalADT.pcADT (Sharpened_DelegateSigs ?4 idx))
         (?6 idx) (?7 idx)))
subgoal 3 (ID 10) is:
refineADT (ComputationalADT.LiftcADT (Sharpened_Implementation ?4 ?6 ?7))
   (Computational ADT. LiftcADT ?5)
(dependent evars: ?3 using , ?4 open, ?5 open, ?6 open, ?7 open, ?11 using ,)
(* hone method "splits". *)
4 subgoals, subgoal 1 (ID 60)
 r_n : string * (nat * nat)
 n : item Ascii.ascii * production Ascii.ascii
  H := ?50 : hiddenT
  refine
     (dummy <- {_ : list nat | True};</pre>
     ls <- If [Terminal "a"%char; Terminal "b"%char; NonTerminal "(ab)*"] =p
              fst n :: snd n Then ret [1]
           Else (If [Terminal "b"%char; NonTerminal "(ab)*"] =p
                    fst n :: snd n Then ret [1]
                 Else (If [NonTerminal "(ab)*"] =p fst n :: snd n
                       Then ret [ilength r_n] Else
                       ret dummy));
     ret (r_n, ls)) (H r_n n)
(* simplify parser splitter. *)
1 focused subgoals (unfocused: 3)
, subgoal 1 (ID 362)
 r_n : string * (nat * nat)
 n : item Ascii.ascii * production Ascii.ascii
  H := ?50 : hiddenT
  refine
     (ret
        (r_n,
       [If ([Terminal "a"%char; Terminal "b"%char; NonTerminal "(ab)*"] =p
            fst n :: snd n)
           || ([Terminal "b"%char; NonTerminal "(ab)*"] =p fst n :: snd n)
        Then 1
        Else (If [NonTerminal "(ab)*"] =p fst n :: snd n Then
```

#### ilength r\_n Else 0)])) (H r\_n n)

```
Definition indexed_spec : ADT (string_rep Ascii.ascii) := ADTRep T {
 Def Constructor "new"(s : String.string) : rep :=
   ret (s, (0, String.length s)),
 Def Method "to_string"(s : rep, x : unit) : String.string :=
   ret (s, string_of_indexed s),
 Def Method "is_char"(s : rep, ch : Ascii.ascii) : bool :=
   ret (s, string_beq (string_of_indexed s) (String.String ch "")),
 Def Method "get"(s : rep, n : nat) : option Ascii.ascii :=
   ret (s, iget n s),
 Def Method "length"(s : rep, x : unit) : nat :=
   ret (s, ilength s),
 Def Method "take"(s : rep, n : nat) : unit :=
   ret ((fst s, (fst (snd s), min (snd (snd s)) n)), tt),
 Def Method "drop"(s : rep, n : nat) : unit :=
   ret ((fst s, (fst (snd s) + n, snd (snd s) - n)), tt),
 Def Method "splits"(s : rep, p : item Ascii.ascii * production Ascii.ascii) : lis
   fallback_ls <- { ls : list nat
                   | match fst p with
                       | Terminal _
                         => True
                       | NonTerminal _
                         => if has_only_terminals (snd p)
                            else split_list_is_complete G (string_of_indexed s) (fs
                     end };
   let ls := (match snd p, fst p with
                 | nil, _
                   => [ilength s]
                 | _::_, Terminal _
                   => [1]
                 | _::_, NonTerminal _
                   => if has_only_terminals (snd p)
                      then [ilength s - List.length (snd p)]
                      else fallback_ls
               end) in
```

```
ret (s, ls)
  }.
 Definition split_list_is_complete '{HSL : StringLike Char} (str : String) (it : ite
            (splits : list nat)
  : Prop
    := forall n,
        n <= length str
         -> production_is_reachable (it::its)
         -> forall (pit : parse_of_item G (take n str) it)
                   (pits : parse_of_production G (drop n str) its),
             Forall_parse_of_item (fun _ nt => List.In nt (Valid_nonterminals G)) pi
             -> Forall_parse_of_production (fun _ nt => List.In nt (Valid_nontermina
             -> List. In n splits.
Inductive length_result := same_length (n : nat) | different_lengths | cyclic_length
Coercion collapse_length_result (1 : length_result) : option nat
  := match l with
      | same_length n => Some n
      | _ => None
    end.
Fixpoint length_of_any_production' {Char} (length_of_any_nt : String.string -> length
         (its : production Char) : length_result
  := match its with
      | nil => same_length 0
       | (Terminal _)::xs => match length_of_any_production' length_of_any_nt xs with
                              | same_length n => same_length (S n)
                              | different_lengths => different_lengths
                              | cyclic_length => cyclic_length
                              | not_yet_handled_empty_rule => not_yet_handled_empty_
                            end
       | (NonTerminal nt)::xs => match length_of_any_nt nt, length_of_any_production'
                                  | same_length n1, same_length n2 => same_length (n
                                  | cyclic_length, _ => cyclic_length
                                  | _, cyclic_length => cyclic_length
                                  | different_lengths, _ => different_lengths
                                  | _, different_lengths => different_lengths
```

| \_, not\_yet\_handled\_empty\_rule => not\_yet\_handled

end.

```
Lemma length_of_any_production',ext {Char}
      (ext : forall b, f b = g b)
: @length_of_any_production' Char f b = length_of_any_production' g b.
Proof.
  induction b as [ | x ]; try reflexivity; simpl.
  destruct x; rewrite ?IHb, ?ext; reflexivity.
Qed.
Definition length_of_any_productions'_f
  := (fun x1 x2)
      => match x1, x2 with
           | same_length n1, same_length n2 => if Nat.eq_dec n1 n2 then same_length n
           cyclic_length, _ => cyclic_length
           | _, cyclic_length => cyclic_length
           | _, different_lengths => different_lengths
           | different_lengths, _ => different_lengths
           | not_yet_handled_empty_rule, _ => not_yet_handled_empty_rule
           | _, not_yet_handled_empty_rule => not_yet_handled_empty_rule
         end).
Arguments length_of_any_productions'_f !_ !_ / .
Lemma length_of_any_productions'_f_same_length {n x1 x2}
: length_of_any_productions'_f x1 x2 = same_length n
  <-> (x1 = same\_length n / x2 = same\_length n).
Proof.
  destruct x1, x2; simpl in *;
  repeat match goal with
           | _ => reflexivity
           | [ H : context[Nat.eq_dec ?x ?y] |- _ ] => destruct (Nat.eq_dec x y)
           | [ |- context[Nat.eq_dec ?x ?y] ] => destruct (Nat.eq_dec x y)
           | _ => progress subst
           | [ H : same_length _ = same_length _ |- _ ] => inversion H; clear H
           _ => intro
           | [ H : _ /\ _ |- _ ] => destruct H
           | [ |- _ /\ _ ] => split
           | [ |- _ <-> _ ] => split
           | _ => congruence
           | _ => tauto
         end.
```

```
Qed.
```

```
Lemma length_of_any_productions'_f_same_length_fold_right {n x1 x2}
: fold_right length_of_any_productions'_f x1 x2 = same_length n
  <-> (x1 = same_length n /\ fold_right and True (map (fun k => k = same_length n) x2
Proof.
  induction x2; simpl in *; eauto; try tauto.
  rewrite length_of_any_productions'_f_same_length.
  rewrite IHx2.
  tauto.
Qed.
Definition length_of_any_productions' {Char} (length_of_any_nt : String.string -> len
           (prods : productions Char)
: length_result
  := match prods with
       | nil => not_yet_handled_empty_rule
       | p::ps => fold_right
                    length_of_any_productions'_f
                    (length_of_any_production', length_of_any_nt p)
                    (map (length_of_any_production, length_of_any_nt) ps)
     end.
Lemma length_of_any_productions'_ext {Char}
      (ext : forall b, f b = g b)
: @length_of_any_productions' Char f b = length_of_any_productions' g b.
Proof.
  unfold length_of_any_productions'.
  destruct b as [ | ? b]; trivial; [].
  induction b; try reflexivity; simpl;
  erewrite length_of_any_production'_ext by eassumption; trivial; [].
  edestruct (length_of_any_production', g);
    rewrite IHb; reflexivity.
Qed.
Definition length_of_any_nt_step {Char} (G : grammar Char) (predata := @rdp_list_pred
           (valid0 : nonterminals_listT)
           (length_of_any_nt : forall valid, nonterminals_listT_R valid valid)
                                              -> String.string -> length_result)
           (nt : String.string)
: length_result.
Proof.
  refine (if Sumbool.sumbool_of_bool (is_valid_nonterminal valid0 nt)
```

```
then length_of_any_productions'
                 (@length_of_any_nt (remove_nonterminal valid0 nt) (remove_nontermina
                 (Lookup G nt)
          else different_lengths);
  assumption.
Defined.
Lemma length_of_any_nt_step_ext {Char G}
      x0 f g
      (ext : forall y p b, f y p b = g y p b)
: @length_of_any_nt_step Char G x0 f b = length_of_any_nt_step g b.
Proof.
 unfold length_of_any_nt_step.
  edestruct Sumbool.sumbool_of_bool; trivial.
  apply length_of_any_productions'_ext; eauto.
Qed.
Definition length_of_any_nt {Char} (G : grammar Char) initial : String.string -> leng
  := let predata := @rdp_list_predata _ G in
     @Fix _ _ ntl_wf _
          (@length_of_any_nt_step _ G)
          initial.
Definition length_of_any {Char} (G : grammar Char) : String.string -> length_result
  := @length_of_any_nt Char G initial_nonterminals_data.
Definition length_of_any_productions {Char} G := @length_of_any_productions' Char (@l
Global Arguments ComputationalSplitter / .
Require Import Fiat.Parsers.ParserFromParserADT.
Require Import Fiat.Parsers.ExtrOcamlParsers.
Import Fiat.Parsers.ExtrOcamlParsers.HideProofs.
Time Definition ab_star_parser (str : String.string) : bool
  := Eval simpl in has_parse (parser ComputationalSplitter) str.
Print ab_star_parser.
Recursive Extraction ab_star_parser.
```

Recursive Extraction test0 test1 test2.

# Chapter 7

# disjoint items, parsing #, +

- Goals
  - More exploration
  - Demonstrate a slightly less obvious strategy that handles even more rules
- [Intro]
- The splitting strategy: if the set of all characters in one item are disjoint from the set of possible first characters of the next item, then we can split at either the first character not in the first set, or at the first character that is in the second set.
  - For example, if the nonterminal "number" accepts 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, and the nonterminal "binop" accepts +, -, then we can either look for the first non-digit and split there, or we can look for the first + or -, and split there.
- Write down the grammars we want to handle.
- We compute the sets reflectively. Here is the algorithm. Again, note the well-founded recursion. <Coq code here>
- We relate the computation to the parse trees by a few correctness criteria, similar to above.
  - Describe the structures and lemmas that go into it.
- This rule is not applied automatically, because we have this choice about what sets to look for. In the future, we might pick whichever set is smaller, and do that (but perhaps we think one is more likely than the other?) Instead we use setoid\_rewrite, with [reflexivity] to solve the side-conditions.
- Example: Parse "1+2+3"

```
Section IndexedImpl.
Typeclasses Opaque If_Then_Else.
  Lemma ComputationalSplitter'
  : FullySharpened (string_spec plus_expr_grammar).
  Proof.
    start honing parser using indexed representation.
    hone method "splits".
      simplify parser splitter.
      setoid_rewrite refine_disjoint_search_for; [ | reflexivity.. ].
      finish honing parser method.
    Fully Sharpen Each {\tt MethodWithoutDelegation}\,.
    extract delegate-free implementation.
    simpl; higher_order_reflexivityT.
  Defined.
  Lemma ComputationalSplitter
  : FullySharpened (string_spec plus_expr_grammar).
  Proof.
    let impl := (eval simpl in (projT1 ComputationalSplitter')) in
    refine (existT _ impl _).
    abstract (exact (projT2 ComputationalSplitter')).
  Defined.
End IndexedImpl.
Lemma refine_disjoint_search_for {G : grammar Ascii.ascii} {str nt its}
      (Hvalid : grammar_rvalid G)
      (H_disjoint : disjoint ascii_beq
                              (possible_terminals_of G nt)
                              (possible_first_terminals_of_production G its))
: refine {splits : list nat
         | split_list_is_complete
             Gstr
             (NonTerminal nt)
             its splits}
         (ret [find_first_char_such_that str (fun ch => list_bin ascii_beq ch (possib
Lemma refine_disjoint_search_for_not {G : grammar Ascii.ascii} {str nt its}
      (Hvalid : grammar_rvalid G)
```

# 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} ::= ' \mid | \operatorname{'} \operatorname{'n} \mid | \operatorname{'} \operatorname{'} \operatorname{'} | | \operatorname{'} \operatorname{'r} | \\ \end{array}
```

# 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

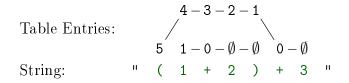
the string right before the '+', and say that the "4" can be parsed as a number (or parenthesized expression), and that the "+5" can be parsed as a '+' followed by an expression.

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?]

Although proving and formalizing this rule took some doing, once it was formalized, using it for any particular grammar is a one-liner with setoid\_rewrite:

setoid\_rewrite refine\_binop\_table; [ presimpl\_after\_refine\_binop\_table | reflexivit

[QUESTION FOR ADAM: How much explanation does this code deserve? Should it be in here at all?]

In order to make it this easy, we had to formalize a correctness condition on the grammar: any use of the '+' character must be hidden by parentheses when on the left of the pexpr+expr rule, and pexpr must only generate well-parenthesized strings. By defining a function that computes this as a boolean we can use reflexivity to prove that any particular grammar is valid, when it is.

For example, for the string "((1+2)+3)+4", we have the following lists at each point in the string: [TODO: Insert computation here]

• The generated table is "5, 1, 0, -, -, 0 . , 4, 3, 2, 1"

#### 8.3.4 The Code

[TODO: Insert Haskell-like code here]

Optimization: Compute based on original string, lookup based on indices, don't need to compute substrings.

#### 8.3.5 The Correctness Proof

To use this as a refinement rule, we need to prove it complete. To do this, we prove [TODO: explain correctness criterion and insert relation of paren-balanced-hiding on a string to properties of grammar.]

Section IndexedImpl.

```
Lemma ComputationalSplitter'
  : FullySharpened (string_spec plus_expr_grammar).
    start honing parser using indexed representation.
    hone method "splits".
      simplify parser splitter.
      setoid_rewrite refine_binop_table; [ presimpl_after_refine_binop_table | reflex
      finish honing parser method.
    }
    Fully Sharpen Each {\tt MethodWithoutDelegation}.
    extract delegate-free implementation.
    simpl; higher_order_reflexivityT.
  Defined.
  Lemma ComputationalSplitter
  : FullySharpened (string_spec plus_expr_grammar).
  Proof.
    let impl := (eval simpl in (projT1 ComputationalSplitter')) in
    refine (existT _ impl _).
    abstract (exact (projT2 ComputationalSplitter')).
  Defined.
End IndexedImpl.
  Local Notation retT table
    := (refine {splits : list nat
               | split_list_is_complete
                   G (substring n m str)
                    (NonTerminal nt)
                    (Terminal ch :: its) splits}
               (ret [match List.nth n table None with
                       | Some idx => idx
                       | None => dummy_value
                     end])).
    Lemma refine_binop_table
          (predata := rdp_list_predata (G := G))
          (pdata := correct_open_close)
```

```
(H_nt_hiding
       : match possible_valid_open_closes with
           | nil => false
           | => true
         end)
: retT (list_of_next_bin_ops_opt str).
 unfold correct_open_close, possible_valid_open_closes in *.
  subst pdata.
  revert H_nt_hiding.
  generalize possible_open_closes.
  intro ls.
  induction ls; simpl.
  { intro; congruence. }
  { match goal with
      | [ |- context[if ?e then _ else nil] ] => destruct e eqn:?
    end.
    { simpl; intro.
      apply refine_binop_table''; try assumption.
      apply ascii_lb; reflexivity. }
    { simpl; assumption. } }
Qed.
Definition bin_op_data_of (open close : Ascii.ascii)
: paren_balanced_hiding_dataT Ascii.ascii
  := {| is_bin_op := ascii_beq ch;
        is_open := ascii_beq open;
        is_close := ascii_beq close |}.
Definition maybe_open_closes {Char} (p : production Char)
: list (Char * Char)
  := match hd None (map Some p), hd None (map Some (rev p)) with
       | Some (Terminal open), Some (Terminal close)
         => [(open, close)]
       | _, _ => nil
     end.
Definition possible_open_closes
: list (Ascii.ascii * Ascii.ascii)
  := fold_right
       (@app _)
       nil
       (map maybe_open_closes (Lookup G nt)).
Definition possible_valid_open_closes
```

```
: list (Ascii.ascii * Ascii.ascii)
  := fold_right
       (@app _)
       nil
       (map
          (fun oc
           => if paren_balanced_hiding_correctness_type (pdata := bin_op_data_of
              then [oc]
              else nil)
          possible_open_closes).
Definition bin_op_data_of_maybe (oc : option (Ascii.ascii * Ascii.ascii))
: paren_balanced_hiding_dataT Ascii.ascii
  := {| is_bin_op := ascii_beq ch;
        is_open ch' := match oc with
                          | Some oc' => ascii_beq (fst oc') ch'
                         | None => false
                       end;
        is_close ch' := match oc with
                          | Some oc' => ascii_beq (snd oc') ch'
                          | None => false
                        end | }.
Definition correct_open_close
: paren_balanced_hiding_dataT Ascii.ascii
  := bin_op_data_of_maybe
       (hd None (map Some possible_valid_open_closes)).
Definition paren_balanced'_step (ch : Char) (pbh_rest : nat -> bool) (start_level
: bool
  := if is_bin_op ch
     then pbh_rest start_level
     else if is_open ch
          then pbh_rest (S start_level)
          else if is_close ch
               then ((Compare_dec.gt_dec start_level 0)
                       && pbh_rest (pred start_level))%bool
               else pbh_rest start_level.
Global Instance paren_balanced'_step_Proper {ch}
: Proper ((eq ==> eq) ==> eq ==> eq) (paren_balanced'_step ch).
```

```
Proof.
  unfold paren_balanced'_step.
  repeat intro; subst.
  unfold respectful in *.
  edestruct Compare_dec.gt_dec; simpl;
  repeat match goal with
           | _ => reflexivity
           | [ |- context[if ?e then _ else _] ] => destruct e
           | [ H : _ |- _ ] => apply H
         end.
Qed.
Definition paren_balanced' (str : String) (start_level : nat)
: bool
  := fold
       paren_balanced'_step
       Compare_dec.zerop
       str
       start level.
Definition paren_balanced (str : String) := paren_balanced' str 0.
Definition paren_balanced_hiding'_step (ch : Char) (pbh_rest : nat -> bool) (star
: bool
  := if is_bin_op ch
     then ((Compare_dec.gt_dec start_level 0)
             && pbh_rest start_level)%bool
     else paren_balanced'_step ch pbh_rest start_level.
Global Instance paren_balanced_hiding'_step_Proper {ch}
: Proper ((eq ==> eq) ==> eq ==> eq) (paren_balanced_hiding'_step ch).
Proof.
  unfold paren_balanced_hiding'_step.
  repeat intro; subst.
  edestruct Compare_dec.gt_dec; simpl;
  repeat match goal with
           | _ => reflexivity
           | [ H : _ |- _ ] => erewrite !H; reflexivity
         end.
Qed.
Definition paren_balanced_hiding' (str : String) (start_level : nat)
: bool
  := fold
       paren_balanced_hiding',_step
       (Compare_dec.zerop)
```

```
Definition paren_balanced_hiding (str : String) := paren_balanced_hiding' str 0.
  Definition cell_of_next_bin_ops_spec'' (level : nat) (cell : option nat) (str : Str
    := (cell = Some idx
        -> index_points_to_binop offset idx str
           /\ paren_balanced_hiding' (take idx (drop offset str)) level)
       /\ (cell = None
           -> paren_balanced' (take idx (drop offset str)) level
           -> index_not_points_to_binop offset idx str).
  Definition list_of_next_bin_ops_spec'' (level : nat) (table : list (option nat)) (s
    := cell_of_next_bin_ops_spec', level (nth offset table None) str offset idx.
  Definition list_of_next_bin_ops_spec' (level : nat) (table : list (option nat)) (st
    := forall offset idx, list_of_next_bin_ops_spec', level table str offset idx.
  Definition list_of_next_bin_ops_spec
    := list_of_next_bin_ops_spec' 0.
  (** We build a version of paren-balanced-hiding to compute each cell
      of the table. *)
  (**
pb' ch n "" = (n == 0)
pb' ch n (ch :: s) = n > 0 \&\& pb' ch n s
pb' ch n ('(' :: s) = pb' ch (n + 1) s
pb' ch n (')' :: s) = n > 0 && pb' ch (n - 1) s
pb, ch n (_ :: s) = pb, ch n s
pb = pb', +, 0
>>
*)
  Definition compute_next_bin_op'_step
    := (fun ch next level
        => if is_bin_op ch
           then if Compare_dec.gt_dec level 0
```

str

start\_level.

```
then option_map S (next level)
              else Some 0
         else if is_open ch
              then option_map S (next (S level))
              else if is_close ch
                   then if Compare_dec.gt_dec level 0
                        then option_map S (next (pred level))
                        else None
                   else option_map S (next level)).
Definition compute_next_bin_op' (str : String) (level : nat)
: option nat
  := fold
       compute_next_bin_op'_step
       (fun _ => None)
       str
       level.
```

# 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, performance-wise, 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}} & (\overline{s \in \texttt{nt}}) \\ & \coloneqq \texttt{has\_parse nt } s = \texttt{true} \land \texttt{parse nt } s \neq \texttt{None} \\ &T_{\texttt{failure}} & (\overline{s \in \texttt{nt}}) \\ & \coloneqq \texttt{has\_parse 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.