

**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?



# A Formally Verified Parser for a JavaScript Subset in a New Extensible Framework for Synthesizing Efficient Verified Parsers

by

Jason S. Gross

Submitted to the Department of Electrical Engineering and Computer  
Science

in partial fulfillment of the requirements for the degree of

Master of Science in Computer Science and Engineering

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 2015

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

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

Author .....  
Department of Electrical Engineering and Computer Science  
July 9, 2015

Certified by .....  
Adam Chlipala  
Associate Professor  
Thesis Supervisor

Accepted by .....  
Leslie A. Kolodziejcki  
Chair, Department Committee on Graduate Students



# **A Formally Verified Parser for a JavaScript Subset in a New 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

<b>1</b>	<b>Introduction</b>	<b>15</b>
1.1	Related Work . . . . .	15
1.2	What's New . . . . .	16
<b>2</b>	<b>Parsing Context Free Grammars</b>	<b>17</b>
2.1	Parsing . . . . .	17
2.1.1	Infinite regress . . . . .	18
2.1.2	Aborting early . . . . .	19
2.2	Standard Formal Definitions . . . . .	21
2.2.1	Context Free Grammar . . . . .	21
2.2.2	Parse Trees . . . . .	21
<b>3</b>	<b>Completeness and Soundness</b>	<b>23</b>
<b>4</b>	<b>Completeness, Soundness, and Parsing Parse Trees</b>	<b>25</b>
4.1	Proving Completeness: Conceptual Approach . . . . .	25
4.2	Minimal Parse Trees: Formal Definition . . . . .	26
4.3	Parser Interface . . . . .	27

4.3.1	Parsing Parses . . . . .	31
4.3.2	Example . . . . .	33
4.3.3	Parametricity . . . . .	36
4.3.4	Putting it all together . . . . .	37
4.4	Missteps, Insights, and Dependently Typed Lessons . . . . .	38
4.4.1	The trouble of choosing the right types . . . . .	38
4.4.2	Misordered splitters . . . . .	38
4.4.3	Minimal Parse Trees vs. Parallel Traces . . . . .	39
<b>5</b>	<b>Refining Splitters by Fiat</b>	<b>41</b>
5.1	Splitters at a Glance . . . . .	41
5.2	What counts as efficient? . . . . .	41
5.3	Introducing Fiat . . . . .	42
5.3.1	Incremental Construction by Refinement . . . . .	42
5.3.2	The Fiat Mindset . . . . .	42
5.4	Optimizations . . . . .	44
5.4.1	An Easy First Optimization: Indexed Representation of Strings	44
5.4.2	Upcoming Optimizations . . . . .	44
<b>6</b>	<b>fixed length nonterminals, parsing (ab)*; parsing #s; parsing #, ()</b>	<b>45</b>
<b>7</b>	<b>disjoint items, parsing #, +</b>	<b>47</b>
<b>8</b>	<b>Parsing well-parenthesized expressions</b>	<b>49</b>
8.1	At a Glance . . . . .	49

8.2	Grammars we can parse . . . . .	49
8.3	The Splitting Strategy . . . . .	51
8.3.1	The Main Idea . . . . .	51
8.3.2	Building the Lookup Table . . . . .	51
8.3.3	The Code . . . . .	51
8.3.4	The Correctness Proof . . . . .	51
<b>9</b>	<b>Future work</b>	<b>53</b>
9.1	Future work with dependent types . . . . .	53



# List of Figures

4-1	The dependently typed interface of our parser . . . . .	28
4-2	Pseudo-Implementation of our parser. We take the convention that dependent indices to functions (e.g., <b>unseen</b> ) are implicit. . . . .	34



# List of Tables





# Chapter 1

## Introduction

**TODO:** Very basic Background on parsers, grammars

### 1.1 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**

## 1.2 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 2

## Parsing Context Free Grammars

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

### 2.1 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.

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

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

oracle must give at least one splitting location that also yields a valid parse. Different splitters yield different simple recursive-descent parsers.

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

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

Although this simple splitter is sufficient for proving the algorithm correct, it is horribly inefficient, running in time  $\mathcal{O}(n!)$ , where  $n$  is the length of the string. We synthesize more efficient splitters in later chapters; we believe that parameterizing the parser over a splitter gives us enough expressiveness to implement essentially all optimizations of interest, while being a sufficiently simple language to make proofs relatively straightforward. For example, to achieve linear parse time on the  $(ab)^*$  grammar, we could have a splitter that, when trying to parse  $'c_1' \cdot 'c_2' \cdot s$  as  $ab(ab)^*$ , splits the string into  $('c_1', 'c_2', s)$ ; and when trying to parse  $s$  as  $\epsilon$ , 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.

### 2.1.1 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:

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

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

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

[illegible]

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 → [String × String]` which is our splitting oracle, we may write out a potentially divergent parser specialized to this grammar.

```
any_pareses : [String × String] → Bool
any_pareses [] := false
any_pareses (("a", "b") :: _) := true
any_pareses ((s1, s2) :: rest_splits)
    := (pareses s1 && pareses s2) || any_pareses rest_splits

pareses : String → Bool
pareses "" := true
pareses str := any_pareses (split str)
```

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

### 2.1.2 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`.

```

any_pares : String → [String × String] → Bool
any_pares _ [] := false
any_pares _ (("a", "b") :: _) := true
any_pares prev_s ((s1, s2) :: rest_splits)
  := (s1 < prev_s && s2 < prev_s
      && parses s1 && parses s2)
    || any_pares prev_s rest_splits

parses : String → Bool
parses "" := true
parses str := any_pares str (split str)

```

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.

## 2.2 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.

### 2.2.1 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:

$$\frac{}{''a'' \in 'a'} \quad \frac{}{''b'' \in 'b'}$$

Productions and nonterminals:

$$\frac{s \in \epsilon}{s \in (ab)^*} \quad \frac{}{'''' \in \epsilon}$$

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

### 2.2.2 Parse Trees

A string  $s$  *parses* as:

- a given terminal  $ch$  iff  $s = 'ch'$ .

- a given sequence of items  $\mathbf{x}_i$  iff  $\mathbf{s}$  splits into a sequence of strings  $\mathbf{s}_i$ , each of which parses as the corresponding item  $\mathbf{x}_i$ .
- a given nonterminal  $\mathbf{nt}$  iff  $\mathbf{s}$  parses as one of the item sequences that  $\mathbf{nt}$  maps to under the set of productions.

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

$$\begin{aligned} \text{ParseTreeOf} &: \text{Item} \rightarrow \text{String} \rightarrow \mathbf{Type} \\ \text{ParseItemsTreeOf} &: [\text{Item}] \rightarrow \text{String} \rightarrow \mathbf{Type} \end{aligned}$$

For any terminal character  $\mathbf{ch}$ , we have the constructor

$$(\mathbf{'ch'}) : \text{ParseTreeOf } \mathbf{'ch'} \text{ "ch"}$$

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

$$(\mathbf{rule}) : \text{ParseItemsTreeOf } \mathbf{its} \ \mathbf{s} \rightarrow \text{ParseTreeOf } \mathbf{nt} \ \mathbf{s}$$

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

$$\begin{array}{c} \overline{\text{""} \in \epsilon : \text{ParseItemsTreeOf } [] \text{ ""}} \\ \frac{s_1 \in \mathbf{it} \quad s_2 \in \mathbf{its}}{s_1 s_2 \in \mathbf{it} :: \mathbf{its}} : \text{ParseTreeOf } \mathbf{it} \ s_1 \\ \rightarrow \text{ParseItemsTreeOf } \mathbf{its} \ s_2 \\ \rightarrow \text{ParseItemsTreeOf } (\mathbf{it} :: \mathbf{its}) \ s_1 s_2 \end{array}$$

For brevity, we will sometimes use the notation  $\mathbf{s} \in \mathbf{X}$  to denote both `ParseTreeOf  $\mathbf{X} \ \mathbf{s}$`  and `ParseItemsTreeOf  $\mathbf{X} \ \mathbf{s}$` , relying on context to disambiguate based on the type of  $\mathbf{X}$ . Additionally, we will sometimes fold the constructors of `ParseItemsTreeOf` into the  $(\mathbf{rule})$  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

$$\{(\mathbf{nt}, \mathbf{s}) : \text{Nonterminal} \times \text{String} \mid \text{ParseTreeOf } \mathbf{nt} \ \mathbf{s}\}$$



# 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 → String → 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:

```
has_parse_sound : (nt : Nonterminal) → (s : String)
  → has_parse nt s = true
  → ParseTreeOf nt s
has_parse_complete : (nt : Nonterminal) → (s : String)
  → ParseTreeOf nt s
  → has_parse nt s = true
```

For any parser

```
parse : Nonterminal → String → option ParseTree,
```

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

following type signatures, using  $p_1$  to denote the first projection of  $p$ :

```

parse_sound : (nt : Nonterminal)
  → (s : String)
  → (p : ParseTree)
  → parse nt s = Some p
  → p1 = (nt, s)
parse_complete : (nt : Nonterminal)
  → (s : String)
  → ParseTreeOf nt s
  → parse nt s ≠ None

```

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):

```

parse : (nt : Nonterminal)
  → (s : String)
  → ParseTreeOf nt s + (ParseTreeOf nt s → ⊥)

```

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  $\perp$ , 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 2.1.2 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 2.2.2, 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 `unseen0`.

```
MinParseTreeOf : String → [Nonterminal]
                → Item → String → Type
MinItemsTreeOf : String → [Nonterminal]
                → [Item] → String → Type
```

Much as in the case of parse trees, for any terminal character `ch`, any string `s0`, and any list of nonterminals `unseen`, we have the constructor

```
min_parsech : MinParseTreeOf s0 unseen 'ch' "ch"
```

For any production `rule` mapping a nonterminal `nt` to a sequence of items `its`, any string `s0`, 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.

```
(rule)< : s < s0
        → MinItemsTreeOf s unseen0 its s
        → MinParseTreeOf s0 unseen nt s
(rule)= : s = s0
        → nt ∈ unseen
        → MinItemsTreeOf s0 unseen' its s
        → MinParseTreeOf s0 unseen nt s
```

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`.

```
min_parse[] : MinItemsTreeOf s0 unseen [] ""
min_parse:: : s1s2 ≤ s0
              → MinParseTreeOf s0 unseen it s1
              → MinItemsTreeOf s0 unseen its s2
              → MinItemsTreeOf s0 unseen (it :: its) s1s2
```

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 s0 v X s` and `MinItemsTreeOf s0 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 `optionParseTree`,

We use `ParseQuery` to denote the type of all propositions like “`"a" ∈ '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 `s0` and a list of nonterminals `unseen`. We assume we are given `unseen0 : [Nonterminal]`.

```

Tsuccess, Tfailure : String → [Nonterminal] → ParseQuery → Type
split : String → [Nonterminal] → ParseQuery → [N]

terminal_success : (ch : Char) → Tsuccess s0 unseen ("ch" ∈ 'ch')
terminal_failure : (ch : Char) → (s : String) → s ≠ "ch" → Tfailure s0 unseen (s ∈ 'ch')
nil_success : Tsuccess s0 unseen ("" ∈ ε)
nil_failure : (s : String) → s ≠ "" → Tfailure s0 unseen (s ∈ ε)

cons_success : (it : Item) → (its : [Item]) → (s1 : String) → (s2 : String)
  → s1s2 ≤ s0
  → Tsuccess s0 unseen (s1 ∈ it)
  → Tsuccess s0 unseen (s2 ∈ its)
  → Tsuccess s0 unseen (s1s2 ∈ it :: its)
cons_failure : (it : Item) → (its : [Item]) → (s : String)
  → s ≤ s0
  → (∀ (s1, s2) ∈ split s0 unseen (s ∈ it :: its),
    Tfailure s0 unseen (s1 ∈ it) + Tfailure s0 unseen (s2 ∈ its))
  → Tfailure s0 unseen (s ∈ it :: its)

production_success< : (its : [Item]) → (nt : Nonterminal) → (s : String)
  → s < s0
  → (p : a production mapping nt to its)
  → Tsuccess s unseen0 (s ∈ its)
  → Tsuccess s0 unseen (s ∈ nt)
production_success= : (its : [Item]) → (nt : Nonterminal) → (s : String)
  → nt ∈ unseen
  → (p : a production mapping nt to its)
  → Tsuccess s0 (unseen - {nt}) (s ∈ its)
  → Tsuccess s0 unseen (s ∈ nt)
production_failure< : (nt : Nonterminal) → (s : String)
  → s < s0
  → (∀ (its : [Item]) (p : a production mapping nt to its), Tfailure s unseen
  → Tfailure s0 unseen (s ∈ nt)
production_failure= : (nt : Nonterminal) → (s : String)
  → s = s0
  → (∀ (its : [Item]) (p : a production mapping nt to its), Tfailure s0 (unseen

```

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

$$\text{parse} : \text{Nonterminal} \rightarrow \text{String} \rightarrow T_{\text{success}} + T_{\text{failure}}.$$

To instantiate the parser as a Boolean recognizer, we instantiate everything trivially; we use the fact that  $\top + \top \cong \text{Bool}$ . 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  $\top$ :

```

Tsuccess _ _ _ :=  $\top$ 
Tfailure _ _ _ :=  $\top$ 

terminal_success _ _ _ := ()
terminal_failure _ _ _ _ := ()
nil_success _ _ := ()
nil_failure _ _ _ _ := ()
cons_success _ _ _ _ _ _ _ := ()
cons_failure _ _ _ _ _ _ _ := ()
production_success< _ _ _ _ _ _ _ := ()
production_success= _ _ _ _ _ _ _ := ()
production_failure< _ _ _ _ _ _ _ := ()
production_failure= _ _ _ _ _ _ _ := ()
production_failure $\notin$  _ _ _ _ _ _ _ := ()

```

To instantiate our parser so that it returns `option ParseTree` (rather, the dependently typed flavor, `ParseTreeOf`), we take advantage of the isomorphism  $T + \top \cong \text{option } T$ . 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.

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

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  $\perp$  is the empty type (equivalently, the false proposition).

$$\begin{aligned}
T_{\text{success}} \_ \_ \_ &:= \top \\
T_{\text{failure}} \text{ } s_0 \text{ } \text{unseen} \_ \_ (\overline{s \in X}) &:= \left( \overline{s \in X} <_{(s_0, \text{unseen})} \right) \rightarrow \perp
\end{aligned}$$

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

$$\begin{aligned}
\text{terminal\_success} \_ \_ \_ &:= () \\
\text{terminal\_failure} \_ \_ \_ \_ &:= \text{!} \\
\text{nil\_success} \_ \_ &:= () \\
\text{nil\_failure} \_ \_ \_ \_ &:= \text{!} \\
\text{cons\_success} \_ \_ \_ \_ \_ \_ \_ &:= () \\
\text{cons\_failure} \_ \_ \_ \_ \_ \_ \_ &:= \text{!} \\
\text{production\_success}_{<} \_ \_ \_ \_ \_ \_ \_ &:= () \\
\text{production\_success}_{=} \_ \_ \_ \_ \_ \_ \_ &:= () \\
\text{production\_failure}_{<} \_ \_ \_ \_ \_ \_ &:= \text{!} \\
\text{production\_failure}_{=} \_ \_ \_ \_ \_ \_ &:= \text{!} \\
\text{production\_failure}_{\neq} \_ \_ \_ \_ \_ \_ &:= \text{!}
\end{aligned}$$

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  $s \in \bar{X}$ , we look for any subtrees also proving  $s \in \bar{X}$ . 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:

$$\text{deloop} : \text{ParseTreeOf } \text{nt } s \rightarrow \text{ParseTreeOf } \text{nt } s$$

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  $\text{Unloopy}(t)$  on parse trees  $t$  to mean that, where the root node of  $t$  proves  $s \in \bar{nt}$ , for every subtree proving  $s \in \bar{nt}'$  (same string, possibly different nonterminal), (1)  $\bar{nt}'$  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  $\bar{nt}' \neq \bar{nt}$ .

We augment the return type of `deloop`, writing:

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

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:

```

split _ _ (s ∈ it :: its) :=
  case parse_tree_data s of
  |  $\frac{p_1}{s_1 \in \text{it}} \quad \frac{p_2}{s_2 \in \text{its}}$  → [(deloop p1, deloop p2)]
  | _ → []
split _ _ _ := []

```

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  $T + \perp \cong T$  and instantiate  $T_{\text{failure}}$  with  $\perp$ , the uninhabited type; and then apply a simple fix-up wrapper on top. Thus, we instantiate the general parser like so:

```

Tsuccess s0 unseen (d : s ∈ X̄) := s ∈ X̄ <^(s0, unseen)
Tfailure _ _ _ := ⊥

```

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 ( $s \neq \text{"ch"}$  and  $s \neq \text{" "}$ , respectively) to derive  $\perp$  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  $\perp + \perp \rightarrow \perp$ , 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` ∈ `unseen`. Then we can derive  $\perp$  from the hypothesis that `nt` ∉ `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 2.1.1, 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

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

Then we consider minimizing the parse tree:

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

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

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:

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

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

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  $(ab)^*$ , of which there are none.

We then take the following execution steps, starting with `unseen` := `unseen0` := `[(ab)*]`, the singleton list containing the only nonterminal, and `s0` := "abab".

1. We first ensure that we are not in an infinite loop. We check if `s` < `s0` (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' (s0 := s) (unseen := unseen0) (s ∈ nt)

parse' ("ch" ∈ "ch") := inl terminal_success
parse' (_ ∈ "ch") := inr (terminal_failure ℓ)
parse' (" " ∈ ε) := inl nil_success
parse' (_ ∈ ε) := inr (nil_failure ℓ)
parse' (s ∈ it :: its) :=
  case any_parse s it its (split (s ∈ it :: its)) of
  | inl ret → inl ret
  | inr ret → inr (cons_failure _ ret)
parse' (s ∈ nt) :=
  if s < s0
  then if (parse' (s0 := s) (unseen := unseen0) (s ∈ 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 ∈ unseen
        then if (parse' (unseen := unseen - {nt}) (s ∈ 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≠ _ _)

any_parse s it its [] := inr (λ _ : (_ ∈ []). ℓ)
any_parse s it its (x :: xs) :=
  case parse' (takex s ∈ it), parse' (dropx s ∈ its), any_parse s it its xs of
  | inl ret1, inl ret2, _ → inl (cons_success _ ret1 ret2)
  | _ , _ , inl ret' → inl ret'
  | ret1 , ret2 , inr ret' → inr _

```

where the hole on the last line constructs a proof of

$$\forall x' \in (x :: xs), T_{\text{failure}} \_ \_ (\text{take}_{x'} s \in \text{it}) + T_{\text{failure}} \_ \_ (\text{drop}_{x'} s \in \text{its})$$

by using  $\text{ret}'$  directly when  $x' \in xs$ , and using whichever one of  $\text{ret}_1$  and  $\text{ret}_2$  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.,  $\text{unseen}$ ) 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

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

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

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

the result of calling **deloop** is the tree

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

- (b) The return of **split** is thus the singleton list containing a single pair of two parse trees; each element of the pair is the parse tree for  $"ab" \in (ab)^*$  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

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

- (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  $\text{unseen}_0 \equiv [(ab)^*]$  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  $\overline{\text{"a"} \in \text{'a'}}$  and  $\overline{\text{"b"} \in \text{'b'}}$ . 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  $\overline{\text{"a"} \in \text{'a'}}$ , and the second is the parse tree  $\overline{\text{"b"} \in \text{'b'}}$ .
- (d) We invoke **parse** on each of the items in the sequence of items associated to  $(ab)^*$  via the rule  $(\text{"ab"})$ . Since both of these items are terminals, and the relevant equality check (that  $\text{"a"}$  is equal to  $\text{"a"}$ , and similarly for  $\text{"b"}$ ) succeeds, **parse** returns **terminal\_success**. We thus have the two **MinParseTrees**:  $\overline{\text{"a"} \in \text{'a'}}$  and  $\overline{\text{"b"} \in \text{'b'}}$ .
- (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**<sub><</sub> to this tree, and return the tree

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

- 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

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

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

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

### 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 `parse1` and `parse2`, we have

$$\forall (\text{nt} : \text{Nonterminal}) (\text{s} : \text{String}), \\ \text{bool\_of\_sum} (\text{parse}_1 \text{ nt s}) = \text{bool\_of\_sum} (\text{parse}_2 \text{ nt s})$$

where `bool_of_sum` is, for any types  $A$  and  $B$ , the function of type  $A + B \rightarrow \text{Bool}$  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:

```
has_parse : Nonterminal → String → Bool
parse : (nt : Nonterminal) → (s : String)
      → option (ParseTreeOf nt s)
has_no_parse : (nt : Nonterminal) → (s : String)
      →  $\top$  + (MinParseTreeOf nt s →  $\perp$ )
min_parse : (nt : Nonterminal) → (s : String)
      → ParseTreeOf nt s
      → MinParseTreeOf nt s
```

Note that we have taken advantage of the isomorphism  $\top + \top \cong \text{Bool}$  for `has_parse`, the isomorphism  $A + \top \cong \text{option } A$  for `parse`, and the isomorphism  $A + \perp \cong A$  for `min_parse`.

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

```
parse_full : (nt : Nonterminal) → (s : String)
      → ParseTreeOf nt s + (ParseTreeOf nt s →  $\perp$ )
parse_full nt s :=
  case parse nt s, has_no_parse nt s of
  | Some d, _      → inl d
  | _ , inr nd     → inr (nd ∘ min_parse)
  | _ , _         →  $\text{!}$ 
```

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.<sup>1</sup> 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>1</sup>For readers wanting to skip that examination: the algorithm we described allows a label  $(\underline{s} \in \underline{nt})$  to appear one extra time along a path if, the first time it appears, its parent node’s label,  $(\underline{s}' \in \underline{nt}')$ , satisfies  $\underline{s} < \underline{s}'$ . That is, whenever the string being parsed shrinks, the first nonterminal the shrunk 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 1.2, 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  $\iff$  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:

$\mathbf{x} \leftarrow \mathbf{c}; \mathbf{c}'$  Run  $\mathbf{c}$  and store the result in  $\mathbf{x}$ ; continue with  $\mathbf{c}'$ , which may mention  $\mathbf{x}$   
 $\mathbf{c};; \mathbf{c}'$  Run  $\mathbf{c}$ . If it terminates, throw away the result, and run  $\mathbf{c}'$   
 $\mathbf{ret} \ \mathbf{x}$  A computation that immediately returns the value  $\mathbf{x}$   
 $\{\mathbf{x} \mid \mathbf{P}(\mathbf{x})\}$  Nondeterministically choose a value of  $\mathbf{x}$  satisfying  $\mathbf{P}$ .

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

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

$$\begin{aligned} \mathbf{ret} \ \mathbf{v} &\rightsquigarrow \mathbf{v} \\ \{\mathbf{x} \mid \mathbf{P}(\mathbf{x})\} &\rightsquigarrow \mathbf{v} \text{ iff } \mathbf{x} \text{ satisfies } \mathbf{P} \\ (\mathbf{c};; \mathbf{c}') &\rightsquigarrow \mathbf{v} \text{ iff there is a } \mathbf{v}' \text{ such that } \mathbf{c} \rightsquigarrow \mathbf{v}' \text{ and } \mathbf{c}' \rightsquigarrow \mathbf{v} \\ (\mathbf{x} \leftarrow \mathbf{c}; \mathbf{c}'(\mathbf{x})) &\rightsquigarrow \mathbf{v} \text{ iff there is a } \mathbf{v}' \text{ such that } \mathbf{c} \rightsquigarrow \mathbf{v}' \text{ and } \mathbf{c}'(\mathbf{v}') \rightsquigarrow \mathbf{v} \end{aligned}$$

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)



# 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"



# 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: **QUESTION FOR ADAM:** Should I write these as inference rules, or as standard CFG presentations?

$$\begin{array}{c} \frac{s \in \epsilon}{s \in \text{number?}} \text{ (number?-\epsilon)} \qquad \frac{s \in \text{number}}{s \in \text{number?}} \text{ (number?)} \\[10pt] \frac{s_0 \in \text{digit} \quad s_1 \in \text{number?}}{s_0 s_1 \in \text{number}} \text{ (number)} \\[10pt] \frac{s_0 \in '(' \quad s_1 \in \text{expr} \quad s_2 \in ')'}{s \in \text{pexpr}} \text{ (pexpr)} \end{array}$$

$$\frac{s \in \text{number}}{s \in \text{expr}} \text{ (expr-number)} \qquad \frac{s_0 \in \text{pexpr} \quad s_1 \in +\text{expr}}{s \in \text{expr}} \text{ (expr-pexpr)}$$

$$\frac{s \in \epsilon}{s \in +\text{expr}} \text{ (\epsilon+expr)} \qquad \frac{s_0 \in '+' \quad s_1 \in \text{expr}}{s_0 s_1 \in +\text{expr}} \text{ (+expr)}$$

**TODO:** Pick one of these

Or, in the standard presentation:

```

expr ::= number | pexpr +expr
+expr ::=  $\epsilon$  | '+' expr
pexpr ::= '(' expr ')'
number ::= digit number?
number? ::=  $\epsilon$  | number
digit ::= '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'

```

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:

```

expr ::= atom | pexpr sexpr
sexpr ::=  $\epsilon$  | whitespace expr
pexpr ::= '(' expr ')'
atom ::= symbol atom?
atom? ::=  $\epsilon$  | atom
whitespace ::= whitespace-char whitespace?
whitespace? ::=  $\epsilon$  | whitespace
whitespace-char ::= ' ' | '\n' | '\t' | '\r'

```

## 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 parenthetization is. If we can compute an appropriate lookup table in time linear in the length of the string, then our splitter overall will 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 parenthetization. To compute this location we keep a list of the location of next '+' at every level of parenthetization. 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 "

### 8.3.3 The Code

**TODO: Insert Haskell-like code here**

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

### 8.3.4 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.**



# 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 <lookup 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:

$$\begin{aligned}
T_{\text{success}} &= \lambda s. (s \in \text{nt}) \\
&:= \text{has\_parse nt } s = \text{true} \wedge \text{parse nt } s \neq \text{None} \\
T_{\text{failure}} &= \lambda s. (s \notin \text{nt}) \\
&:= \text{has\_parse nt } s = \text{false} \wedge \text{has\_no\_parse } s \neq \text{inl } ()
\end{aligned}$$

**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.