

Lecture Notes on Predictive Parsing

15-411: Compiler Design
Frank Pfenning

Lecture 8
September 18, 2009

1 Introduction

In this lecture we discuss two parsing algorithms, both of which traverse the input string from left to right. The first, LL(1), makes a decision on which grammar production to use based on the first character of the input string. If that were ambiguous, the grammar would have to be rewritten to fall into this class, which is not always possible. The second, LR(1), can postpone the decision at first by pushing input characters onto a stack and then deciding on the production later, taking into account both the first input character and the stack. It is variations on the latter which are typically used in modern parser generation tools.

Alternative presentations of the material in this lecture can be found in the textbook [[App98](#), Chapter 3] and a paper by Shieber et al. [[SSP95](#)].

2 LL(1) Parsing

We have seen in the previous section, that the general idea of recursive descent parsing without restrictions forces us to non-deterministically choose between several productions which might be applied and potentially backtrack if parsing gets stuck after a choice, or even loop (if the grammar is left-recursive). Backtracking is not only potentially very inefficient, but it makes it difficult to produce good error messages in case the string is not grammatically well-formed. Say we try three different ways to parse a given term and all fail. How could we say which of these is the source of

the error? This is compounded because nested choices multiply the number of possibilities. We therefore have to look for ways to disambiguate the choices.

One way is to *require* of the grammar that at each potential choice point we can look at the next input token and based on that token decide which production to take. This is called *1 token lookahead*, and grammars that satisfy this restriction are called *LL(1)* grammars. Here, the first *L* stands for *Left-to-right*; the second *L* stands for *Leftmost parse* (which a recursive descent parser generates) and *1* stands for *1 token lookahead*. Potentially, we could also define *LL(2)*, *LL(3)*, etc., but these are of limited practical utility.

Since we are restricting ourselves to parsing by a left-to-right traversal of the input string, we will consider only tails, or postfixes of the input strings, and also of the strings in the grammar, when we restrict our inference rules. For short, we will say γ is a *postfix substring* of the grammar, or w is a postfix substring of the input string w_0 . For example, in the grammar

$$\begin{array}{l} [\text{emp}] \quad S \longrightarrow \\ [\text{pars}] \quad S \longrightarrow [S] \\ [\text{dup}] \quad S \longrightarrow S S \end{array}$$

the only postfix substrings are ϵ , $[S]$, $S]$, $]$, S , and $S S$, but not $[S$.

We begin by defining two kinds of predicates (later we will have occasion to add a third), where β is either a non-terminal or postfix substring of the grammar.

$$\begin{array}{ll} \text{first}(\beta, a) & \text{Token } a \text{ can be first in string } \beta \\ \text{null}(\beta) & \text{String } \beta \text{ can produce the empty string } \epsilon \end{array}$$

These predicates must be computed entirely statically, by an analysis of the grammar before any concrete string is ever parsed. This is because we want to be able to tell if the parser can do its work properly with 1 token look-ahead regardless of the string it has to parse.

We define the relation $\text{first}(\beta, a)$ by the following rules.

$$\frac{}{\text{first}(a \beta, a)} F_1$$

This rule *seeds* the first predicate. Then is it propagated to other strings appearing in the grammar by the following three rules.

$$\frac{\text{first}(X, a)}{\text{first}(X \beta, a)} F_2 \quad \frac{\text{null}(X) \quad \text{first}(\beta, a)}{\text{first}(X \beta, a)} F_3 \quad \frac{[r]X \longrightarrow \gamma \quad \text{first}(\gamma, a)}{\text{first}(X, a)} F_4(r)$$

Even though ϵ may be technically a postfix substring of every grammar, it can never arise in the first argument of the first predicate. The auxiliary predicate `null` is also easily defined.

$$\frac{}{\text{null}(\epsilon)} N_1 \quad \frac{\text{null}(X) \quad \text{null}(\gamma)}{\text{null}(X \gamma)} N_2 \quad \frac{[r]X \longrightarrow \gamma \quad \text{null}(\gamma)}{\text{null}(X)} N_3$$

We can run these rules to saturation because there are only $O(|G|)$ possible strings in the first argument to both of these predicates, and at most the number of possible terminal symbols in the grammar, $O(|\Sigma|)$, in the second argument. Naive counting the number of prefix firings (see [GM02]) gives a complexity bound of $O(|G| \times |\Xi| \times |\Sigma|)$ where $|\Xi|$ is the number of non-terminals in the grammar. Since usually the number of symbols is a small constant, this is roughly equivalent to $O(|G|)$ and so is reasonably efficient. Moreover, it only happens once, before any parsing takes place.

Next, we modify the rules for recursive descent parsing from the last lecture to take these restrictions into account. The first two stay the same.

$$\frac{}{\epsilon : \epsilon} R_1 \quad \frac{w : \gamma}{a w : a \gamma} R_2$$

The third,

$$\frac{[r]X \longrightarrow \beta \quad w : \beta \gamma}{w : X \gamma} R_3(r)$$

is split into two, each of which has an additional precondition.

$$\frac{[r]X \longrightarrow \beta \quad \text{first}(\beta, a) \quad a w : \beta \gamma}{a w : X \gamma} R'_3 \quad \frac{[r]X \longrightarrow \beta \quad \text{null}(\beta) \quad w : \beta \gamma}{w : X \gamma} R''_3?$$

We would like to say that a grammar is LL(1) if the additional preconditions in these last two rules make all choices unambiguous when an arbitrary non-terminal X is matched against a string starting with an arbitrary terminal a . Unfortunately, this does not quite work yet in the presence non-terminals that can rewrite to ϵ , because the second rule above does not even look at the input string.

To further refine this we need one additional predicate, again on postfix string in the grammar and non-terminals.

$\text{follow}(\beta, a)$ Token a can follow string β in a valid string

We seed this relation with the rules

$$\frac{X \gamma \text{ postfix} \quad \text{first}(\gamma, a)}{\text{follow}(X, a)} W_1$$

Here, $X \gamma \text{ postfix}$ means that the string $X \gamma$ appears as a postfix substring on the right-hand side of a production. We then propagate this information applying the following rules from premises to conclusion until saturation is reached.

$$\frac{\text{follow}(b \gamma, a)}{\text{follow}(\gamma, a)} W_2 \quad \frac{\text{follow}(X \gamma, a)}{\text{follow}(\gamma, a)} W_3 \quad \frac{\text{follow}(X \gamma, a) \quad \text{null}(\gamma)}{\text{follow}(X, a)} W_4 \quad \frac{[r]X \rightarrow \gamma \quad \text{follow}(X, a)}{\text{follow}(\gamma, a)} W_5$$

The first argument here should remain a non-empty postfix or a non-terminal. Now we can refine the proposed R_3'' rule from above into one which is much less likely to be ambiguous.

$$\frac{[r]X \rightarrow \beta \quad \text{first}(\beta, a) \quad a w : \beta \gamma}{a w : X \gamma} R_3' \quad \frac{[r]X \rightarrow \beta \quad \text{null}(\beta) \quad \text{follow}(X, a) \quad a w : \beta \gamma}{a w : X \gamma} R_3''$$

We avoid creating an explicit rule to treat the empty input string by appending a special $\$$ symbol at the end before starting the parsing process. We repeat the remaining rules for completeness.

$$\frac{}{\epsilon : \epsilon} R_1 \quad \frac{w : \gamma}{a w : a \gamma} R_2$$

These rules are interpreted as a parser by proof search, applying them from the conclusion to the premise. We say the grammar is LL(1) if for any goal $w : \gamma$ at most one rule applies. If X cannot derive ϵ , this amounts to checking that there is at most one production $X \rightarrow \beta$ such that $\text{first}(\beta, a)$.

For nullable non-terminals the condition is slightly more complicated, but can still easily be read off from the rules.

We now use a very simple grammar to illustrate these rules. We have transformed it in the way indicated above, by assuming a special token $\$$ to indicate the end of the input string.

$$\begin{aligned} [\text{start}] \quad S &\longrightarrow S' \$ \\ [\text{emp}] \quad S' &\longrightarrow \epsilon \\ [\text{pars}] \quad S' &\longrightarrow [S'] \end{aligned}$$

This generates all string starting with an arbitrary number of opening parentheses followed by the same number of closing parentheses and an end-of-string marker.

We have:

$\text{null}(\epsilon)$	N_1
$\text{null}(S')$	N_3
$\text{first}([S'], [)$	F_1
$\text{first}(,)$	F_1
$\text{first}(S'],)$	F_3
$\text{first}(S', [)$	F_4 [pars]
$\text{first}(S'], [)$	F_2
$\text{first}(\$, \$)$	F_1
$\text{first}(S' \$, \$)$	F_3
$\text{first}(S' \$, [)$	F_2
$\text{first}(S , \$)$	F_4 [start]
$\text{first}(S , [)$	F_4 [start]
$\text{follow}(S' , \$)$	W_1
$\text{follow}(S' ,])$	W_1
$\text{follow}([S'], \$)$	W_5
$\text{follow}([S'],])$	W_5
$\text{follow}(S'], \$)$	W_3
$\text{follow}(S'],])$	W_3
$\text{follow}(, \$)$	W_4
$\text{follow}(,])$	W_4

3 Parser Generation

Parser generation is now a very simple process. Once we have computed the null, first, and follow predicates by saturation from a given grammar,

we specialize the inference rules $R'_3(r)$ and $R''_3(r)$ by matching the first two and three premises against grammar productions and saturated database. In this case, this leads to the following specialized rules (repeating once again the two initial rules).

$$\begin{array}{c} \frac{}{\epsilon : \epsilon} R_1 \qquad \frac{w : \gamma}{a w : a \gamma} R_2 \\ \\ \frac{[w : S' \ \$ \ \gamma]}{[w : S \ \gamma]} R'_3(\text{start}) \qquad \frac{\$ w : S' \ \$ \ \gamma}{\$ w : S \ \gamma} R'_3(\text{start}) \\ \\ \frac{[w : [S'] \ \gamma]}{[w : S' \ \gamma]} R'_3(\text{pars}) \qquad \frac{] w : \gamma]}{] w : S' \ \gamma]} R''_3(\text{emp}) \qquad \frac{\$ w : \gamma}{\$ w : S' \ \gamma} R''_3(\text{emp}) \end{array}$$

Recall that these rules are applied from the bottom-up, starting with the goal $w_0 \$: S$, where w_0 is the input string. It is easy to observe by pattern matching that each of these rules are mutually exclusive: if one of the applies, none of the other rules applies. Moreover, each rule except for R_1 (which accepts) has exactly one premise, so the input string is traversed linearly from left-to-right, without backtracking. When none of the rules applies, then the input string is not in the language defined by the grammar. This proves that our simple language $(^n)^n$ is LL(1).

Besides efficiency, an effect of this approach to parser generation is that it supports good error messages in the case of failure. For example, if we see the parsing goal $(w :) \ \gamma$ we can state: *Found ' (' while expecting ') '*, and similarly for other cases that match none of the conclusions of the rules.

4 LR(1) Parsing

One difficulty with LL(1) parsing is that it is often difficult or impossible to rewrite a grammar so that 1 token look-ahead during a left-to-right traversal becomes unambiguous. To illustrate this, we return to the earlier example of simple arithmetic expressions.

$$\begin{array}{l} [\text{plus}] \quad E \longrightarrow E + E \\ [\text{times}] \quad E \longrightarrow E * E \\ [\text{ident}] \quad E \longrightarrow id \\ [\text{number}] \quad E \longrightarrow num \\ [\text{parens}] \quad E \longrightarrow (E) \end{array}$$

If we see a simple expression such as $3 + 4 * 5$ (which becomes the token stream $num + num * num$), we cannot predict when we see the $+$ symbol which production to use because of the inherent ambiguity of the grammar.

We can rewrite the grammar, at significant expense of readability, or we could just specify that multiplication has higher precedence than addition, $+ < *$. Obviously, the latter is more convenient, but how can we make it work?

The idea is to put off the decision on which productions to use and just *shift* the input symbols onto a stack until we can make the decision! We write

$$\gamma \mid w \quad \text{parse input } w \text{ under stack } \gamma$$

where, as generally in predictive parsing, the rules are interpreted as transitions from the conclusion to the premises. The parsing attempt succeeds if we can consume all of w and produce the start symbol S on the left-hand side. That is, the deduction representing a successful parse of terminal string w_0 has the form

$$\frac{}{S \mid \epsilon} R_1$$

$$\vdots$$

$$\epsilon \mid w_0$$

Parsing is defined by the following rules:

$$\frac{}{S \mid \epsilon} R_1$$

$$\frac{\gamma a \mid w}{\gamma \mid a w} R_2 (= \text{shift}) \quad \frac{\begin{array}{l} [r]X \longrightarrow \beta \\ \gamma X \mid w \end{array}}{\gamma \beta \mid w} R_3(r) (= \text{reduce}(r))$$

We resume the example above, parsing $num + num * num$. After one step (reading this bottom-up)

$$\begin{array}{l|l} num & + num * num \quad ? \\ \epsilon & num + num * num \quad \text{shift} \end{array}$$

we already have to make a decision: should we shift + or should we reduce num using rule [number]. In this case the action to reduce is forced, because we will never get another chance to see this num as an E .

E	$+ num * num$?
num	$+ num * num$	reduce(number)
ϵ	$num + num * num$	shift

At this point we need to shift +; no other action is possible. We take a few steps and arrive at

$E + E$	$* num$	
$E + num$	$* num$	reduce(number)
$E +$	$num * num$	shift
E	$+ num * num$	shift
num	$+ num * num$	reduce(number)
ϵ	$num + num * num$	shift

At this point, we have a real conflict. We can either reduce, viewing $E + E$ as a subexpression, or shift and later consider $E * E$ as a subexpression. Since the $*$ has higher precedence than +, we need to shift.

E	ϵ	R_1
$E + E$	ϵ	reduce(plus)
$E + E * E$	ϵ	reduce(times)
$E + E * num$	ϵ	reduce(number)
$E + E *$	num	shift
$E + E$	$* num$	shift
$E + num$	$* num$	reduce(number)
$E +$	$num * num$	shift
E	$+ num * num$	shift
num	$+ num * num$	reduce(number)
ϵ	$num + num * num$	shift

Since E was the start symbol in this example, this concludes the deduction. If we now read the lines from the top to the bottom, ignoring the separator, we see that it represents a *rightmost* derivation of the input string. So we have parsed analyzing the string from left to right, constructing a rightmost derivation. This type of parsing algorithms is called LR-parsing, where the L stands for left-to-right and the R stands for rightmost.

The decisions above are based on the postfix of the stack on the left-hand side and the first token on the right-hand side. Here, the postfix of

the stack on the left-hand side must be a *prefix substring* of a grammar production. If not, it would be impossible to complete it in such a way that a future grammar production can be applied in a reduction step: the parse attempt is doomed to failure.

5 LR(1) Parsing Tables

We could now define again a slightly different version of $\text{follow}(\gamma, a)$, where γ is a prefix substring of the grammar or a non-terminal, and then specialize the rules. An alternative, often used to describe parser generators, is to construct a *parsing table*. For an LR(1) grammar, this table contains an entry every for every prefix substring of the grammar and token seen on the input. An entry describes whether to shift, reduce (and by which rule), or to signal an error. If the action is ambiguous, the given grammar is not LR(1), and either an error message is issued, or some default rule comes into effect that chooses between the option.

We now construct the parsing table, assuming $+ < *$, that is, multiplication binds more tightly than addition. Moreover, we specify that both addition and multiplication are *left associative* so that, for example, $3 + 4 + 5$ should be parsed as $(3 + 4) + 5$. We have removed *id* since it behaves identically to *num*.

[plus]	E	\longrightarrow	$E + E$
[times]	E	\longrightarrow	$E * E$
[number]	E	\longrightarrow	<i>num</i>
[parens]	E	\longrightarrow	(E)

As before, we assume that a special end-of-file token $\$$ has been added to the end of the input string. When the parsing goal has the form $\gamma\beta \mid a w$ where β is a prefix substring of the grammar, we look up β in the left-most column and a in the top row to find the action to take. The non-terminal ϵE in the last line is a special case in that E must be the only thing on the stack. In that case we can accept if the next token is $\$$ because we know that $\$$ can only be the last token of the input string.

	+	*	<i>num</i>	()	\$
$E + E$	reduce(plus) (+ left assoc.)	shift (+ < *)	error	error	reduce(plus)	reduce(plus)
$E * E$	reduce(times) (+ < *)	reduce(times) (* left assoc.)	error	error	reduce(times)	reduce(times)
<i>num</i>	reduce(number)	reduce(number)	error	error	reduce(number)	reduce(number)
(<i>E</i>)	reduce(parens)	reduce(parens)	error	error	reduce(parens)	reduce(parens)
$E +$	error	error	shift	shift	error	error
$E *$	error	error	shift	shift	error	error
(E	shift	shift	error	error	shift	error
(error	error	shift	shift	error	error
ϵE	shift	shift	error	error	error	accept(E)

We can see that the bare grammar has four shift/reduce conflicts, while all other actions (including errors) are uniquely determined. These conflicts arise when $E + E$ or $E * E$ is on the stack and either + or * is the first character in the remaining input string. It is called a shift/reduce conflict, because either a shift action or a reduce action could lead to a valid parse. Here, we have decided to resolve the conflicts by giving a precedence to the operators and declaring both of them to be left-associative.

It is also possible to have reduce/reduce conflicts, if more than one reduction could be applied in a given situation, but it does not happen in this grammar.

Parser generators will generally issue an error or warning when they detect a shift/reduce or reduce/reduce conflict. For many parser generators, the default behavior of a shift/reduce conflict is to shift, and for a reduce/reduce conflict to apply the textually first production in the grammar. Particularly the latter is rarely what is desired, so we strongly recommend rewriting the grammar to eliminate any conflicts in an LR(1) parser.

One interesting special case is the situation in a language where the else-clause of a conditional is optional. For example, one might write (among other productions)

$$\begin{aligned} S &\longrightarrow \text{if } E \text{ then } E \text{ else } E \\ S &\longrightarrow \text{if } E \text{ then } E \end{aligned}$$

Now an statement

if (b) then if (c) then x else y

is ambiguous because it would be read as

```
if (b) then (if (c) then x) else y
```

or

```
if (b) then (if (c) then x else y)
```

In a shift/reduce parser, typically the default action for a shift/reduce conflict is to shift. This means that the above grammar in a tool such as ML-Lex will parse the ambiguous statement into the second form, that is, the `else` is match with the most recent unmatched `if`. This is consistent with language such as C (or the language used in this course), so we can tolerate the above shift/reduce conflict, if you wish, instead of rewriting the grammar to make it unambiguous.

References

- [App98] Andrew W. Appel. *Modern Compiler Implementation in ML*. Cambridge University Press, Cambridge, England, 1998.
- [GM02] Harald Ganzinger and David A. McAllester. Logical algorithms. In P. Stuckey, editor, *Proceedings of the 18th International Conference on Logic Programming*, pages 209–223, Copenhagen, Denmark, July 2002. Springer-Verlag LNCS 2401.
- [SSP95] Stuart M. Shieber, Yves Schabes, and Fernando C. N. Pereira. Principles and implementation of deductive parsing. *Journal of Logic Programming*, 24(1–2):3–36, 1995.