# Lecture Notes on
# Top-Down Predictive LL Parsing

15-411: Compiler Design
Frank Pfenning[*]

Lecture 8

## 1 Introduction

In this lecture we discuss a parsing algorithm that traverses the input string from left to right giving a left-most derivation and makes a decision on which grammar production to use based on the first character/token of the input string. This parsing algorithm is called LL(1). If that were ambiguous, the grammar would have to be rewritten to fall into this class, which is not always possible. Most hand-written parsers are recursive descent parsers, which follow the LL(1) principles.

Alternative presentations of the material in this lecture can be found in a paper by Shieber et al. [SSP95]. The textbook [App98, Chapter 3] covers parsing. For more detailed background on parsing, also see [WM95].

## 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 input and all fail. How could we say which of these is the source

---

[*]With edits by André Platzer

of the error? How could we report an error location? This is compounded because nested choices multiply the number of possibilities. We therefore have to look for ways to disambiguate the choices, making parsing more efficient and syntax errors easier to locate.

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 *L*eft-to-right reading of the input; the second *L* stands for *L*eftmost parse (which a recursive descent parser generates) and *1* stands for *1 token lookahead*. We can also define LL(2), LL(3), etc. But for higher $k$, LL($k$) parsers become fairly inefficient and are used less frequently. The LL parser generators (or SLL parser generators which are a simplification of LL) themselves, however, are much more efficient than the LR parser generators.

Since we are restricting ourselves to parsing by a left-to-right traversal of the input string, we will consider only tails, or suffixes of the input strings, and also of the strings in the grammar, when we restrict our inference rules. Those suffixes capture what still has to be parsed or worked on. For short, we will say $\gamma$ is a *suffix substring* of the grammar, or $w$ is a suffix substring of the input string $w_0$. For example, in the grammar

$$
\begin{array}{lll}
[\mathsf{emp}] & S & \longrightarrow \\
[\mathsf{pars}] & S & \longrightarrow \quad [\, S \,] \\
[\mathsf{dup}] & S & \longrightarrow \quad S\, S
\end{array}
$$

the only suffix substrings are $\epsilon$, $[\, S \,]$, $S \,]$, $]$, $S$, and $S\, S$, but not the prefix $[\, S$.

**Firsts & Nulls**    The main idea behind LL(1)-parsing is that we restrict our attention to grammars where we can use a 1 token lookahead to disambiguate which grammar production to use. That is, by peeking at the next input token, we want to know which grammar production helps us. The first thing we want to know for this is what the tokens are that could be generated from a sequence $\beta$ of terminals and non-terminals at all. So imagine $\beta$ describes a shape of how far we have come with parsing and what we expect for the rest of the input. It could be something like $S]$. We begin by defining two kinds of predicates (later we will have occasion to add a third), where $\beta$ is either a non-terminal or suffix substring of the grammar. The predicate first($\beta$, a) captures if token a can be the first token occurring in a word that matches the expression $\beta$. What we also need to know is if

$\beta$ could possibly match the empty word $\epsilon$, because then the first token of $\beta\gamma$ could actually come from the first token of $\gamma$. This is what we use the predicate null($\beta$) for.

| | |
|---|---|
| first($\beta$, a) | Token a can be first in any string produced by $\beta$ |
| null($\beta$) | String $\beta$ can produce the empty string $\epsilon$ |

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 actual input strings it will later have to parse. We want the parser generator to tell us right away if it will work on all input. We do not want to wait till runtime to tell us that it doesn't know what to do with some particular input.

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

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

This rule *seeds* the first predicate with the knowledge that parse strings starting with a token a always start with that token a, no matter what $\beta$ yields. Then is it propagated to other strings appearing in the grammar by the following three rules.

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

Rule $F_2$ says that if $X$ can start with a, then so can $X\beta$. Rule $F_3$ says that if $X$ can produce $\epsilon$ and $\beta$ can start with a, then $X\beta$ can start with a as well. Rule $F_4(r)$ captures that $X$ can start with whatever any of the right-hand sides $\gamma$ of its productions $[r]X \to \gamma$ can start with. Even though $\epsilon$ may be technically a suffix substring of every grammar, it can never arise in the first argument of the first predicate, because it is not a proper token, so we do not need any information about it. The auxiliary predicate null is also easily defined.

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

$N_1$ expresses that $\epsilon$ can produce $\epsilon$ – surprise. That $X\beta$ can produce $\epsilon$ if both $X$ and $\beta$ can ($N_2$). And that $X$ can produce $\epsilon$ if one of its productions has a right-hand side $\gamma$ hat can ($N_3$).

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 $\approx O(|G|)$ and so is reasonably efficient. Moreover, it only happens once, before any parsing takes place.

**Constructing LL(1) Parsing Rules**   Next, we modify the rules for recursive descent parsing from the last lecture to take these restrictions into account. The first two compare rules stay the same.

$$\frac{}{\epsilon : \epsilon} \, L_1 \qquad \frac{w : \gamma}{\mathsf{a}\, w : \mathsf{a}\, \gamma} \, L_2$$

The third generate rule,

$$\frac{\begin{array}{c}[r]X \longrightarrow \beta \\ w : \beta\, \gamma\end{array}}{w : X\, \gamma} \, L_3(r)$$

is split into two, each of which has an additional precondition on when to use it:

$$\frac{\begin{array}{c}[r]X \longrightarrow \beta \\ \mathsf{first}(\beta, \mathsf{a}) \\ \mathsf{a}\, w : \beta\, \gamma\end{array}}{\mathsf{a}\, w : X\, \gamma} \, L_3' \qquad \frac{\begin{array}{c}[r]X \longrightarrow \beta \\ \mathsf{null}(\beta) \\ w : \beta\, \gamma\end{array}}{w : X\, \gamma} \, L_3''?$$

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 $\mathsf{a}$. Unfortunately, this does not quite work yet in the presence non-terminals that can rewrite to $\epsilon$, because the second rule above does not even look at the input string. This nondeterministic rule cannot possibly ensure the appropriate parse without looking at the input. To disambiguate, we need to know what tokens could follow $X$. Only if the first token of $w$ could follow $X$ would we want to use a production for $X$ that would make it vanish.

**Follows**   We thus need to know which tokens could follow $X$ in order to know if it would make sense to skip over it by using one of its productions that produces $\epsilon$. For that, we define one additional predicate, again on suffix strings in the grammar and non-terminals.

$\quad$ follow$(\beta, \texttt{a})$   Token $\texttt{a}$ can follow string $\beta$ in a valid string

See Figure 1 for an illustration of where first and follow come from in a parse tree. The set of all $\texttt{a}$ for which first$(X, \texttt{a})$ characterizes what tokens $X$ itself can start with. The set of all $\texttt{a}$ for which follow$(X, \texttt{a})$ characterizes what tokens can follow $X$.



nonterminals:

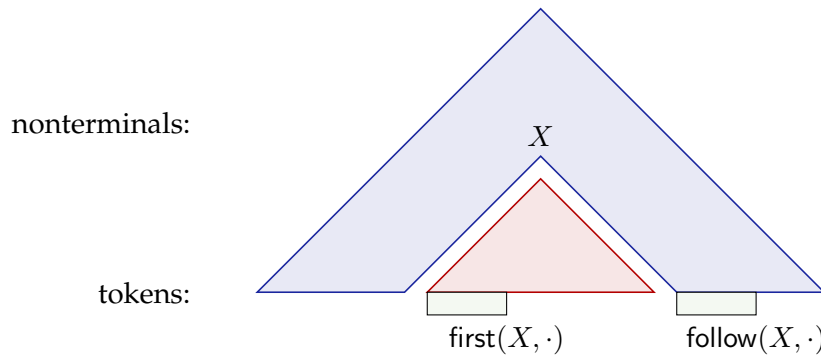tokens:

first$(X, \cdot)$        follow$(X, \cdot)$

Figure 1: Illustration of first and follow sets of non-terminal $X$ in a parse tree in which $X$ occurs

We seed this follows relation with the rules

$$\frac{\begin{array}{c} X\gamma \ \textit{suffix} \\ \mathsf{first}(\gamma, \texttt{a}) \end{array}}{\mathsf{follow}(X, \texttt{a})} \ W_1$$

Here, $X\gamma$ *suffix* means that the string $X\gamma$ appears as a suffix substring on the right-hand side of some production, because then we would want to know if we should choose the productions for $X$ that ultimately produce $\epsilon$, because the lookahead token comes from $\gamma$. Whatever can come first in $\gamma$ can follow $X$ if $X\gamma$ occurs in the right-hand sides of some production. We then propagate this information applying the following rules from premises to conclusion until saturation is reached.

$$\frac{\mathsf{follow}(\texttt{b}\,\gamma, \texttt{a})}{\mathsf{follow}(\gamma, \texttt{a})} \ W_2 \qquad \frac{\mathsf{follow}(X\,\gamma, \texttt{a})}{\mathsf{follow}(\gamma, \texttt{a})} \ W_3 \qquad \frac{\begin{array}{c} \mathsf{follow}(X\,\gamma, \texttt{a}) \\ \mathsf{null}(\gamma) \end{array}}{\mathsf{follow}(X, \texttt{a})} \ W_4 \qquad \frac{\begin{array}{c} [r]X \longrightarrow \gamma \\ \mathsf{follow}(X, \texttt{a}) \end{array}}{\mathsf{follow}(\gamma, \texttt{a})} \ W_5$$

Rule $W_2$ says that whatever can follow $\mathsf{b}\gamma$ can also follow the last part $\gamma$. Similarly, rule $W_3$ says that whatever can follow $X\gamma$ can also follow the last part $\gamma$. Rule $W_4$ says that, in addition, everything that can follow $X\gamma$ can follow $X$ itself if $\gamma$ can produce $\epsilon$ (otherwise the follows of $X\gamma$ are no follows of $X$, because there is always at least one token from $\gamma$ in between). Rule $W_5$ says that whatever follows a nonterminal $X$ can also follow the right-hand side $\gamma$ of each production $[r]X \to \gamma$.

The first argument of follow should remain a non-empty suffix or a nonterminal here, because we are not interested in what could follow $\epsilon$.

**Final LL(1) Parser Rules**    Now we can refine the proposed $L_3''$ rule from above into one which is no longer ambiguous (for LL(1) grammars).

$$\frac{\begin{array}{l}[r]X \longrightarrow \beta \\ \mathsf{first}(\beta, \mathsf{a}) \\ \mathsf{a}\,w : \beta\,\gamma\end{array}}{\mathsf{a}\,w : X\,\gamma}\,L_3' \qquad \frac{\begin{array}{l}[r]X \longrightarrow \beta \\ \mathsf{null}(\beta) \\ \mathsf{follow}(X, \mathsf{a}) \\ \mathsf{a}\,w : \beta\,\gamma\end{array}}{\mathsf{a}\,w : X\,\gamma}\,L_3''$$

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

$$\frac{}{\epsilon : \epsilon}\,L_1 \qquad \frac{w : \gamma}{\mathsf{a}\,w : \mathsf{a}\,\gamma}\,L_2$$

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.

1. If $X$ cannot derive $\epsilon$, this amounts to checking that there is at most one production $X \longrightarrow \beta$ such that $\mathsf{first}(\beta, \mathsf{a})$.

2. Otherwise there is a *first/first conflict* between the two ambiguous productions $X \longrightarrow \beta$ and $X \longrightarrow \gamma$ that share a $\mathsf{first}(\beta, \mathsf{a})$ and $\mathsf{first}(\gamma, \mathsf{a})$.

For nullable non-terminals there are more complicated extra conditions, because it also depends on the follow sets. The conflicts can still easily be read off from the rules. First/first conflicts stay the same.
In addition,

3. There is a *first/follow conflict* if we cannot always decide between $L'_3$ and $L''_3$. This happens if there is a token $a$ with $\text{first}(\beta, a)$ for a production $X \longrightarrow \beta$ that would trigger $L'_3$, but that token also satisfies $\text{follow}(X, a)$ with a nonterminal $X$ that is nullable and trigger $L''_3$.

**Example** 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.

$$
\begin{array}{llcl}
[\text{start}] & S & \longrightarrow & S'\, \$ \\
[\text{emp}] & S' & \longrightarrow & \epsilon \\
[\text{pars}] & S' & \longrightarrow & [\, S'\, ]
\end{array}
$$

This grammar generates all strings starting with an arbitrary number of opening parentheses followed by the same number of closing parentheses and an end-of-string marker, i.e., the language $[^n\, ]^n \$$.
    We have:

$$
\begin{array}{ll}
\text{null}(\epsilon) & N_1 \\
\text{null}(S') & N_3 \\[4pt]
\text{first}(\,[\, S'\, ]\,, [\,) & F_1 \\
\text{first}(\,]\,, ]\,) & F_1 \\
\text{first}(S'\, ]\,, ]\,) & F_3 \\
\text{first}(S', [\,) & F_4\ [\text{pars}] \\
\text{first}(S'\, ]\,, [\,) & F_2 \\[4pt]
\text{first}(\$, \$) & F_1 \\
\text{first}(S'\, \$, \$) & F_3 \\
\text{first}(S'\, \$, [\,) & F_2 \\
\text{first}(S, \$) & F_4\ [\text{start}] \\
\text{first}(S, [\,) & F_4\ [\text{start}] \\[4pt]
\text{follow}(S', \$) & W_1 \\
\text{follow}(S', ]\,) & W_1 \\
\text{follow}(\,[\, S'\, ]\,, \$) & W_5 \\
\text{follow}(\,[\, S'\, ]\,, ]\,) & W_5 \\
\text{follow}(S'\, ]\,, \$) & W_3 \\
\text{follow}(S'\, ]\,, ]\,) & W_3 \\
\text{follow}(\,]\,, \$) & W_4 \\
\text{follow}(\,]\,, ]\,) & W_4
\end{array}
$$

## 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 $L_3'(r)$ and $L_3''(r)$ by matching the first two and three premises against grammar productions and saturated null, first, and follow database. This is essentially partial evaluation on inference rules. In this case, this leads to the following specialized rules (repeating once again the two initial rules).

$$\frac{}{\epsilon : \epsilon} L_1 \qquad \frac{w : \gamma}{\mathsf{a}\, w : \mathsf{a}\, \gamma} L_2$$

$$\frac{[\, w : S'\, \$\, \gamma}{[\, w : S\, \gamma} L_3'(\mathsf{start}) \qquad \frac{\$\, w : S'\, \$\, \gamma}{\$\, w : S\, \gamma} L_3'(\mathsf{start})$$

$$\frac{[\, w : [\, S'\, ]\, \gamma}{[\, w : S'\, \gamma} L_3'(\mathsf{pars}) \qquad \frac{]\, w : \gamma}{]\, w : S'\, \gamma} L_3''(\mathsf{emp}) \qquad \frac{\$\, w : \gamma}{\$\, w : S'\, \gamma} L_3''(\mathsf{emp})$$

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 rules applies, none of the other rules apply. Moreover, each rule except for $L_1$ (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 LL 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 report: *Found '('while expecting ')'* along with a report of the error location. Similarly for other cases that match none of the conclusions of the rules.

## 4   Grammar Transformations

This predictive parsing or LL(1) parsing works quite well. But now suppose we have a grammar

$$\begin{aligned} X &\longrightarrow Y\,b \\ X &\longrightarrow Y\,c \end{aligned}$$

for a nonterminal $Y$ that may produce token streams of unbounded length. This grammar is left-recursive, which is a problem, because we cannot know which production to use in predictive parsing without unbounded token lookahead. Yet we can *left-factorize* the grammar in order to capture the common prefix of the two productions and only distinguish their suffix in a new common nonterminal $X'$. This gives

$$
\begin{aligned}
X &\longrightarrow Y X' \\
X' &\longrightarrow b \\
X' &\longrightarrow c
\end{aligned}
$$

This left-factorization is also useful to remove a First/First conflict from grammar

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

turning it into a form where common prefixes have been factored out

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

More generally, consider the grammar

$$
\begin{aligned}
X &\longrightarrow \alpha \gamma \\
X &\longrightarrow \alpha \delta
\end{aligned}
$$

where we have difficulties deciding between the two productions based on token lookahead, because both start with $\alpha$. We can left-factorize it into the following grammar where common prefixes have been factored out and the distinction between $\gamma$ and $\delta$ happens at a new nonterminal $X'$

$$
\begin{aligned}
X &\longrightarrow \alpha X' \\
X' &\longrightarrow \gamma \\
X' &\longrightarrow \delta
\end{aligned}
$$

The same phenomenon can also happen indirectly, though, even if the grammar rules do not literally start with the same expression $Y$. Suppose that we have a left-recursive grammar (snippet)

$$
\begin{aligned}
E &\longrightarrow E + T \\
E &\longrightarrow T \\
T &\longrightarrow \dots..
\end{aligned}
$$

We cannot choose which grammar rule to use because both $E$ productions can start with terms $(T)$, which may be arbitrarily long. We can change its left recursion into a right recursion by pulling commonalities up. Then we start with what the productions of the nonterminal have in common and postpone the distinction to a new subsequent nonterminal $E'$:

$$
\begin{aligned}
E &\longrightarrow TE' \\
E' &\longrightarrow +TE' \\
E' &\longrightarrow \epsilon
\end{aligned}
$$

This process eliminates left recursion using right recursion. These changes of the grammar have also changed the structure of the parse trees, which needs to be taken into account, e.g., through postprocessing transformations.

## Quiz

1. Is there a language that CYK can parse but LL cannot parse? Give example or explain why none exists.

2. Is there a language that LL can parse but CYK cannot parse? Give example or explain why none exists.

3. Is there a language that recursive descent can parse but LL cannot parse? Give example or explain why none exists.

4. Is there a language that LL can parse but recursive descent cannot parse? Give example or explain why none exists.

5. Why do we limit the token lookahead to 1? Does it make a big difference if we would instead limit it to 5?

6. Should we allow token lookahead $\infty$?

7. Suppose you only have an LL parser generator. Can you still write a compiler for any programming language?

8. Which features of programming languages are difficult or inconvenient for LL parsers?

9. Can $follow(X, \cdot)$ be computed from the grammar productions for $X$ alone? Can $first(X, \cdot)$ be computed from only $X$ productions?

10. Why should parser combinator libraries worry about whether a grammar is LL? Isn't that irrelevant?

11. What is very likely wrong with a compiler when you see a grammar production

$$S ::= E" = "" > "E";"$$

12. Your combinator parser is beautiful, but for some input files it just produces a stack overflow without giving you any other information. What is wrong?

# 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. *J. Log. Program.*, 24(1&2):3–36, 1995.

[WM95] Reinhard Wilhelm and Dieter Maurer. *Compiler Design*. Addison-Wesley, 1995.