

## **Lecture Overview**

- **Context-Free Grammars**
- **Derivations and Parse Trees**
- **Parsing**
  - **Left-factoring**
  - **Predictive Parsing Example**
- **Making a Predictive Grammar**
  - **Removing Direct Left-recursion**

Chapter 3 in text

## Context-Free Grammars

$$G = (N, T, S, P)$$

N: Set of nonterminals (often capital letters)

T: Set of terminals (often lower-case letters and symbols)

S: Start symbol (element of N)

P: Productions of the form

$$A \rightarrow \alpha$$

where  $\alpha$  is a string on  $N \cup T \cup \{\epsilon\}$ .

This is the BNF form of productions (Backus-Naur Form). The milestones use the EBNF (Extended BNF) form, where  $\alpha$  can be a *regular expression* on  $N \cup T \cup \{\epsilon\}$ . This is just a shorthand and does not change which languages have a grammar.

There are standard conventions for giving a grammar. We only list the productions. N is taken to be the set of all symbols on the left-hand side of any production. T is taken to be the set of all other symbols appearing in productions. S is taken to be the left-hand symbol of the first production.

For instance, in the grammar:

$$E \rightarrow T + E$$

$$E \rightarrow T$$

$$T \rightarrow F * T$$

$$T \rightarrow F$$

$$F \rightarrow ( E )$$

$$F \rightarrow c$$

$$F \rightarrow i$$

The nonterminals are the set  $\{E, T, F\}$ , the terminals are the set  $\{+, *, (, ), c, i\}$ , and the start symbol is  $E$ .

This grammar models an arithmetic expression;  $E$  represents “expression”,  $T$  represents “term”, and  $F$  represents “factor”.  $c$  is for “constant”, and  $i$  for “identifier.”

We can also write this grammar as follows:

$$E \rightarrow T + E$$

$T$

$$T \rightarrow F * T$$

$F$

$$F \rightarrow ( E )$$

$c$

$i$

In this, the lines without a production LHS have the LHS of the most recent production with a LHS.

We can even write the grammar as

$$E \rightarrow T + E \mid T$$

$$T \rightarrow F * T \mid F$$

$$F \rightarrow ( E ) \mid c \mid i$$

where the vertical bar is read as “or” (as in regular expressions).

## Derivations and Parse Trees

A production  $A \rightarrow \alpha$  is called a *production for A*.

A **derivation** in a grammar starts with  $\alpha_0$  being the start symbol, and in each step  $k$  it derives  $\alpha_k$  from  $\alpha_{k-1}$  by finding a nonterminal in  $\alpha_{k-1}$  and replacing it with the RHS of a production for that nonterminal. The derivation halts if at some step  $k$  there are no nonterminals remaining in  $\alpha_{k-1}$ .

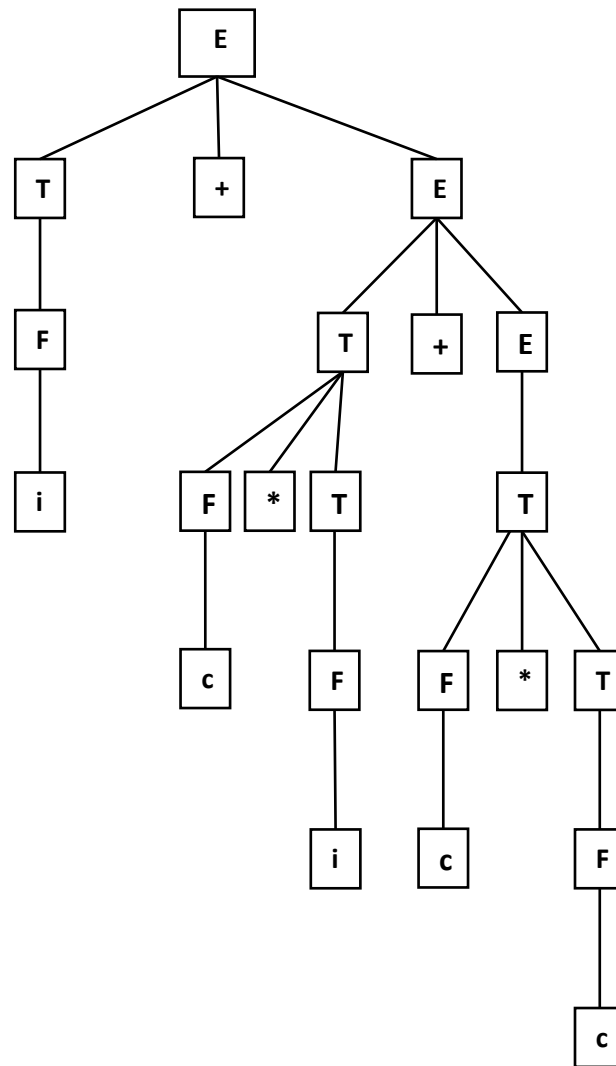
$$\begin{aligned} \underline{E} &\Rightarrow T + \underline{E} \Rightarrow T + T + \underline{E} \Rightarrow \underline{I} + T + T \Rightarrow \underline{E} + T + T \Rightarrow \\ i + \underline{I} + T &\Rightarrow i + \underline{E} * T + T \Rightarrow i + c * \underline{I} + T \Rightarrow \\ i + c * \underline{E} + T &\Rightarrow i + c * i + \underline{I} \Rightarrow i + c * i + \underline{E} * T \Rightarrow \\ i + c * i + c * \underline{I} &\Rightarrow i + c * i + c * \underline{E} \Rightarrow i + c * i + c * c \end{aligned}$$

The last string in a terminating derivation is called a **sentence** or a **word** in the language of the grammar.

For instance,  $i + c * i + c * c$  is a word in the language of our grammar; i.e. it is a valid arithmetic expression.

The **language of a grammar** is all possible sentences that you can derive from the grammar. In most grammars we will be interested in, there are an infinite number of different derivations and therefore the language of the grammar has an infinite number of members.

A **parse tree** of a derivation is the tree one gets by placing the RHS of a production used as a step in the derivation below the LHS and as children of it. The parse tree corresponding to the previous derivation is:



Different derivations can lead to the same parse tree, but different parse trees must lead to different derivations.

## Parsing

It is the job of the parser to look at an input stream (of tokens) and decide which derivation or parse tree of the language that input stream belongs to, if any. Languages are designed to make this task possible and efficient.

For instance, in many languages, each type of statement starts with its own keyword (including assignment statements).

```
stmt → while ( cond ) stmt  
      for ( forControl ) stmt  
      let target = expr ;  
      return expr ;  
      call ident ( args ) ;  
      ...
```

This allows the parser to decide which alternative to use by examining the currently-being-processed token. Parsing becomes more difficult if two



alternatives for a nonterminal start with the same token.

In our sample grammar, we have two different alternatives for E that start with the same symbol. (We also have two different alternatives for T that start with the same symbol.) Here I let the bold parentheses stand for literal parentheses.

$$E \rightarrow \mathbf{T} + E \mid \mathbf{T}$$

$$T \rightarrow \mathbf{F} * T \mid \mathbf{F}$$

$$F \rightarrow (\mathbf{E}) \mid c \mid i$$

We can eliminate this problem through a process known as **left-factoring**. In EBNF, with the nonbold parentheses standing for EBNF grouping:

$$E \rightarrow \mathbf{T} (+ E \mid \epsilon)$$

$$T \rightarrow \mathbf{F} (* T \mid \epsilon)$$

$$F \rightarrow (\mathbf{E}) \mid c \mid i$$

Or, in BNF:

$$E \rightarrow T Y$$

$$Y \rightarrow + E \mid \varepsilon$$

$$T \rightarrow F Z$$

$$Z \rightarrow * T \mid \varepsilon$$

$$F \rightarrow ( E ) \mid c \mid i$$

This grammar is now **predictive**—in parsing a sentence, you can always tell which production to use by what the next token is.

Here's another example. In this grammar, nonterminals are words/identifiers and terminals are written in bold.

doStatement  $\rightarrow$  **do** statement **while** ( condition )  
doStatementVariant

doStatementVariant  $\rightarrow$

**do** statement **until** ( condition )

This is **not** left-factored, as the two alternatives for `doStatement` can both start with the token **do**. (This holds even though the second option doesn't explicitly have the **do** in the RHS.) So first we have to substitute up until both options do explicitly have the **do**. This is simple in this case.

**doStatement** → **do statement while ( condition )  
do statement until ( condition )**

Now we factor out the common prefix:

doStatement  $\rightarrow$  **do** statement doStatementEnd  
doStatementEnd  $\rightarrow$  **while** ( condition )  
**until** ( condition )

## Predictive Parsing Example

For instance, suppose we have the sentence  $3 + 4 * p$  and the expression grammar shown above.

The parser starts with the start symbol  $E$ . The current token is the 3, an instance of  $c$ . There's only one production to take, so it takes the production (does a step in the derivation):

$$E \Rightarrow T Y$$

The first nonterminal in the current string is  $T$ , and the current token is still  $c$ . There's again only one production to take, so we do the step:

$$T Y \Rightarrow F Z Y$$

The first nonterminal in the current string is now  $F$ , and the current token still  $c$ . The three alternatives for  $F$  start with  $($ ,  $c$ , and  $i$ —we take the one that matches the current token.

$$F Z Y \Rightarrow c Z Y$$

This match consumes the current token and so the current token becomes  $+$ . The first nonterminal in the string is now  $Z$ , and no nonempty alternative for  $Z$  starts with  $+$ , so we take the empty alternative.

$$cZY \Rightarrow cY \quad (= c\epsilon Y)$$

Now the first nonterminal is  $Y$ , and the current token is  $+$ , so we take the alternative for  $Y$  that starts with  $+$ .

$$cY \Rightarrow c + E$$

This matches the plus, so we read the next token, which is  $4$ , and instance of  $c$ . The first nonterminal is  $E$ , and it has only one RHS.  $T$  also has only one RHS:

$$c + E \Rightarrow c + T Y \Rightarrow c + F Z Y$$

The first nonterminal in the current string is now  $F$ , and the current token still  $c$ . We take the alternative for  $F$  that matches the current token.

$$c + F Z Y \Rightarrow c + c Z Y$$

Our **c** is now matched, so we get the next token, which is **\***. The first nonterminal is **Z**, so we take the alternative for **Z** that starts with **\***:

$$c + c Z Y \Rightarrow c + c * T Y$$

This matches the **\***, so we get the next token, which is **p**, an instance of **i**. The first nonterminal is **T**, for which there is only one alternative:

$$c + c * T Y \Rightarrow c + c * F Z Y$$

The leftmost nonterminal is **F**, so we take the alternative for **F** that matches the current token, **i**.

$$c + c * F Z Y \Rightarrow c + c * i Z Y$$

We advance the current token to **end-of-input**. **Z** is the first nonterminal, and we do not match the first alternative so we take the empty alternative.

$$c + c * i Z Y \Rightarrow c + c * i Y$$

Now **Y** is the first nonterminal, and again we do not match its first alternative, so we take the empty alternative.

$$c + c * i Y \Rightarrow c + c * i$$

Since we now have all terminals, we are done; we have successfully parsed  $3 + 4 * p$ .

## Making a Predictive Grammar

A predictive grammar has the following two properties. We assume the grammar is in BNF.

- 1) No two alternatives for any production can possibly start with the same terminal. (If this occurs, then **left-factor** the grammar.)
- 2) There is no **left-recursion**. This is a situation where some

$$A \alpha \Rightarrow \dots \Rightarrow A \beta$$

when always expanding the leftmost nonterminal.

For example,  $E \rightarrow E + E$  is a single production that has left recursion (this is called **direct** left-recursion). Another example (this time of **indirect** left-recursion) is:

$$E \rightarrow T ** E$$

$$T \rightarrow E + E$$

because

$$E \Rightarrow T ** E \Rightarrow E + E ** E$$

## Removing direct left-recursion

To remove **direct left-recursion**, write all of the productions as two (left-factored) productions of the form:

$$A \rightarrow A \alpha$$

$$A \rightarrow \beta$$

where  $\alpha$  and  $\beta$  can be regular expressions. Clearly,  $A$  will derive a bunch of  $\alpha$ 's at the end of its string (0 or more) and finally a  $\beta$  at the beginning. So, in regular expressions,  $A$  derives  $\beta \alpha^*$ . We can write this in BNF as:

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A'$$

$$A' \rightarrow \varepsilon$$



For example, if you had:

$$E \rightarrow E + E$$

$$E * E$$

$$( E )$$

$$i$$

$$c$$

You could make this the two productions

$$E \rightarrow E (+ E \mid * E)$$

$$E \rightarrow ( E ) \mid i \mid c$$

Here  $\alpha = (+ E \mid * E)$  and  $\beta = ( E ) \mid i \mid c$

So we could write

$$E \rightarrow ( ( E ) \mid i \mid c ) E'$$

$$E' \rightarrow (+ E \mid * E) E' \mid \varepsilon$$

and this is the equivalent grammar without left-recursion.

For eliminating **indirect left-recursion**, use the following algorithm:

---

Assume nonterminals are  $A_1, A_2, A_3, \dots, A_n$ .

for  $i = 1$  to  $n$

    for  $j = 1$  to  $i-1$

        if  $\exists$  a production  $A_i \rightarrow A_j \gamma$

            replace that production by expanding  $A_j$   
            with all alternatives for  $A_j$

    rewrite any direct left recursion on  $A_i$ .

---

The first iteration of the  $i$  loop eliminates direct left-recursion from  $A_1$ . The second iteration first replaces any  $A_2 \rightarrow A_1 \gamma$  productions and then eliminates direct left-recursion from  $A_2$ . Now  $A_1$  and  $A_2$  have no direct left recursion, and  $A_2$  has no production starting with  $A_1$ .

The third iteration of the  $i$  loop replaces  $A_3 \rightarrow A_1 \gamma$  productions (which could introduce  $A_3 \rightarrow A_2 \gamma$  productions) then replaces  $A_3 \rightarrow A_2 \gamma$  productions

(which can't introduce any  $A_3 \rightarrow A_1 \gamma$  productions), and then eliminates left recursion from  $A_3$ . Now  $A_3$  has no direct left recursion and no production starting with  $A_1$  or  $A_2$ .

And so on. Iteration  $i$  makes it so that  $A_i$  has no production starting with  $A_k$  for any  $k < i$ , and clears left-recursion from  $A_i$ .

At the end, all productions are left-recursion-free.

$$A \rightarrow Ax \mid ABy \mid Bz \mid x$$

$$B \rightarrow Ax \mid Cw \mid Cx$$

$$C \rightarrow Az \mid Bz \mid Cy \mid v$$

$i=1$ : eliminate direct left-recursion from  $A$ .

$$A \rightarrow A(x \mid By) \mid (Bz \mid x) \qquad A \rightarrow A \alpha \mid \beta$$

$$A \rightarrow (Bz \mid x) A'$$

$$A' \rightarrow (x \mid By) A' \mid \varepsilon$$

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \varepsilon$$

This leaves us with the grammar:

$$A \rightarrow (Bz \mid x) A'$$

$$A' \rightarrow (x \mid By) A' \mid \varepsilon$$

$$B \rightarrow Ax \mid Cw \mid Cx$$

$$C \rightarrow Az \mid Bz \mid Cy \mid v$$

$i=2: j=1$ : replace productions for second symbol (B) that start with the first (A):

$$B \rightarrow Ax \mid Cw \mid Cx$$

$$B \rightarrow (Bz \mid x) A' x \mid Cw \mid Cx$$

We now have:

$$A \rightarrow (Bz \mid x) A'$$

$$A' \rightarrow (x \mid By) A' \mid \varepsilon$$

$$B \rightarrow (Bz \mid x) A' x \mid Cw \mid Cx$$

$$C \rightarrow Az \mid Bz \mid Cy \mid v$$

i=2: eliminate direct left-recursion from B.

$$B \rightarrow Bz \mid A'x \mid (xA'x \mid Cw \mid Cx) \quad B \rightarrow B\alpha \mid \beta$$

$$\begin{aligned} B &\rightarrow (xA'x \mid Cw \mid Cx) B' & B &\rightarrow \beta B' \\ B' &\rightarrow zA'xB' \mid \varepsilon & B' &\rightarrow \alpha B' \mid \varepsilon \end{aligned}$$

So we now have:

$$\begin{aligned} A &\rightarrow (Bz \mid x) A' \\ A' &\rightarrow (x \mid By) A' \mid \varepsilon \\ B &\rightarrow (xA'x \mid Cw \mid Cx) B' \\ B' &\rightarrow zA'xB' \mid \varepsilon \\ C &\rightarrow Az \mid Bz \mid Cy \mid v \end{aligned}$$

i=3: j=1: replace productions for third symbol (C) that start with the first (A):

$$C \rightarrow (Bz \mid x) A' z \mid Bz \mid Cy \mid v$$

or

$$C \rightarrow Bz A' z \mid x A' z \mid Bz \mid Cy \mid v$$

i=3: j=2: replace productions for third symbol (C) that start with the second (B):

$$C \rightarrow (x A' x \mid Cw \mid Cx) B' z A' z \mid x A' z \mid \\ (x A' x \mid Cw \mid Cx) B' z \mid Cy \mid v$$

we can rewrite this as:

$$C \rightarrow xA'xB'zA'z \mid CwB'zA'z \mid CxB'zA'z \mid xA'z \mid \\ xA'xB'z \mid CwB'z \mid CxB'z \mid Cy \mid v$$

and change the order of the terms and regroup to get:

$$C \rightarrow C(wB'zA'z \mid xB'zA'z \mid wB'z \mid xB'z \mid y) \mid \\ (A'xB'zA'z \mid xA'z \mid xA'xB'z \mid v)$$

$$C \rightarrow C \alpha \mid \beta$$

i=3: eliminate direct left-recursion from C.

$$C \rightarrow (A'xB'zA'z \mid xA'z \mid xA'xB'z \mid v)C'$$

$$C' \rightarrow (wB'zA'z \mid xB'zA'z \mid wB'z \mid xB'z \mid y)C' \mid \varepsilon$$

$$C \rightarrow \beta C'$$

$$C' \rightarrow \alpha C' \mid \varepsilon$$

We end up with the grammar:

$$A \rightarrow (Bz \mid x) A'$$

$$A' \rightarrow (x \mid By) A' \mid \varepsilon$$

$$B \rightarrow (x A' x \mid Cw \mid Cx) B'$$

$$B' \rightarrow z A' x B' \mid \varepsilon$$

$$C \rightarrow (A'xB'zA'z \mid xA'z \mid xA'xB'z \mid v)C'$$

$$C' \rightarrow (wB'zA'z \mid xB'zA'z \mid wB'z \mid xB'z \mid y)C' \mid \varepsilon$$

This has no left-recursion, which was the goal.

So, as you can see, this algorithmic process can get quite messy. Hopefully one doesn't have a grammar that gives this much mess, but often this process is done by machine and not by hand.

---

Suppose we had:

$E \rightarrow E + E$	low precedence
$E - E$	also low precedence
$E * E$	medium precedence
$-E$	high precedence
$( E )$	<b>highest precedence</b>
$i$	
$c$	

To parse this, we will construct a grammar from high to low precedence. Consider  $i$  and  $c$  as super-high precedence. So we'll call them literals, and have a nonterminal LE (Literal expression) for them:

$$LE \rightarrow i \mid c$$



The next highest precedence level is parentheses; we'll introduce a nonterminal PE (Parenthesized expression) for them. The kind of expression we want in a parenthesized expression is a full-blown expression (E). We also allow a literal to be a parenthesized expression:

$$PE \rightarrow ( E ) \mid LE$$

$$LE \rightarrow i \mid c$$

The next highest level is unary minus; we'll introduce UE (Unary expression) for them. We allow a parenthesized expression to be a unary expression. If we allow just one unary operator to occur before the expression, we get:

$$UE \rightarrow -PE \mid PE$$

Note that we're using the next highest level for the operand of the unary and for skipping the unary. If instead we want more than one unary operator to be allowed (e.g. ---14), then we should use UE itself for the operand.

$$UE \rightarrow -UE \mid PE$$

$$PE \rightarrow ( E ) \mid LE$$

$$LE \rightarrow i \mid c$$

The next highest level is multiplication. Here we use UE as the first argument and the skip option.

$$ME \rightarrow UE * ME \mid UE$$

ME is used as the second argument so we can have expressions with three or more things multiplied together (e.g.  $3 * 4 * n$ ).

The lowest level is addition and subtraction. We'll call the nonterminal AE (additive expression). We use the next level up (ME) as the first argument and the skip option.

$$AE \rightarrow ME + AE \mid ME - AE \mid ME$$

$$ME \rightarrow UE * ME \mid UE$$

$$UE \rightarrow -UE \mid PE$$

$$PE \rightarrow ( E ) \mid LE$$

$$LE \rightarrow i \mid c$$

Let's left-factor the productions:

$$AE \rightarrow ME (+ AE \mid - AE \mid \varepsilon)$$

$$ME \rightarrow UE (* ME \mid \varepsilon)$$

$$UE \rightarrow -UE \mid PE$$

$$PE \rightarrow ( E ) \mid LE$$

$$LE \rightarrow i \mid c$$

Let's work on the first production.

$$AE \rightarrow ME ((+ \mid -) AE \mid \varepsilon) \quad \text{remove recursion}$$

$$AE \rightarrow ME ((+ \mid -) ME)^*$$

Now the second:

$$ME \rightarrow UE (* UE)^*$$

$$ME \rightarrow UE ((* \mid / \mid \%) UE)^*$$

The third:

$$UE \rightarrow (-)^* PE$$

$$UE \rightarrow (- \mid +)^* PE$$

The fourth and fifth are about what we want, so we get:

$$E \rightarrow AE$$

$$AE \rightarrow ME ((+|-) ME)^*$$

$$ME \rightarrow UE ((*|/) UE)^*$$

$$UE \rightarrow (-|+)^* PE$$

$$PE \rightarrow (E) | LE$$

$$LE \rightarrow i | c$$

This grammar is much like what you see in the Tan parser.

Suppose I put in binary operators @, \$ that are even lower precedence than the additive operators. Then we'd get:

$$E \rightarrow LPE$$

$$LPE = AE ((@|\$) AE)^*$$

...