

A context-free grammar (CFG) is defined as:

- A finite terminal set V<sub>t</sub>; these are the tokens produced by the scanner.
- A set of intermediate symbols, called non-terminals, V<sub>n</sub>.
- A start symbol, a designated nonterminal, that starts all derivations.
- A set of productions (sometimes called rewriting rules) of the form

 $A \rightarrow X_1 \dots X_m$   $X_1$  to  $X_m$  may be any combination of terminals and non-terminals. If m =0 we have  $A \rightarrow \lambda$ which is a valid production.

CS 536 Spring 2006

## Example

 $\begin{array}{ll} \mbox{Prog} \rightarrow & \{ \mbox{Stmts} \} \\ \mbox{Stmts} \rightarrow & \mbox{Stmts} ; \mbox{Stmt} \\ \mbox{Stmts} \rightarrow & \mbox{Stmt} \\ \mbox{Stmt} \rightarrow & \mbox{id} = \mbox{Expr} \\ \mbox{Expr} \rightarrow & \mbox{id} \\ \mbox{Expr} \rightarrow & \mbox{Expr} + \mbox{id} \\ \end{array}$ 

CS 536 Spring 2006

194

Often more than one production shares the same left-hand side. Rather than repeat the left hand side, an "or notation" is used:

## Derivations

Starting with the start symbol, nonterminals are rewritten using productions until only terminals remain.

Any terminal sequence that can be generated in this manner is syntactically valid.

If a terminal sequence can't be generated using the productions of the grammar it is invalid (has syntax errors).

The set of strings derivable from the start symbol is the *language* of the grammar (sometimes denoted L(G)).

195



200

CS 536 Spring 2006®



 $\{ id = id ; id = Expr \} \Rightarrow_{L}$   $\{ id = id ; id = Expr + id \} \Rightarrow_{L}$   $\{ id = id ; id = id + id \}$   $Prog \Rightarrow_{L}^{+} \{ id = id ; id = id + id \}$ 

**Rightmost Derivations** 

An alternative to a leftmost derivation is a rightmost derivation, in which the rightmost possible nonterminal is always expanded.

This derivation sequence may seem less intuitive given our normal leftto-right bias, but it corresponds to an important class of parsers (the bottom-up parsers, including CUP).

As a bottom-up parser discovers the productions used to derive a token sequence, it discovers a rightmost derivation, but in *reverse order*.

The last production applied in a rightmost derivation is the first that is discovered, while the first production used, involving the start symbol, is the last to be discovered.

204

CS 536 Spring 2006

203



 $\{ id = id ; id = id + id \}$  $\mathsf{Prog} \Rightarrow^+ \{ \mathsf{id} = \mathsf{id} ; \mathsf{id} = \mathsf{id} + \mathsf{id} \}$ 

You can derive the same set of tokens using leftmost and rightmost derivations; the only difference is the order in which productions are used.

CS 536 Spring 2006

206

## **Ambiguous Grammars**

Some grammars allow more than one parse tree for the same token sequence. Such grammars are ambiguous. Because compilers use syntactic structure to drive translation, ambiguity is undesirableit may lead to an unexpected translation.

Consider

$$\textbf{E} \rightarrow \textbf{E}$$
 -  $\textbf{E}$ 

When parsing the input a-b-c (where a, b and c are scanned as identifiers)

we can build the following two parse trees:



The effect is to parse a-b-c as either (a-b)-c or a-(b-c). These two groupings are certainly not equivalent.

Ambiguous grammars are usually voided in building compilers; the tools we use, like Yacc and CUP, strongly prefer unambiguous grammars.

208

207



## **Operator Precedence**

Most programming languages have operator precedence rules that state the order in which operators are applied (in the absence of explicit parentheses). Thus in C and Java and CSX, a+b\*c means compute b\*c, then add in a.

These operators precedence rules can be incorporated directly into a CFG.

Consider

$$\begin{array}{c} \mathsf{E} \rightarrow \mathsf{E} + \mathsf{T} \\ | & \mathsf{T} \\ \mathsf{T} \rightarrow \mathsf{T} * \mathsf{P} \\ | & \mathsf{P} \\ \mathsf{P} \rightarrow \mathsf{id} \\ | & \mathsf{(E)} \end{array}$$

CS 536 Spring 2006®

211



PP I I id id id

The other grouping can't be obtained unless explicit parentheses are used. (Why?)