

## Context-Free Grammars

## Context-Free Grammars

- Adds recursion/allows non-terminals to be expressed in terms of themselves
- Can be used to count/impart structure – e.g., nested parentheses
- Notation – grammar  $G(S,N,T,P)$ 
  - S is the start symbol
  - N is a set of non-terminal symbols (LHS)
  - T is a set of terminal symbols (tokens)
  - P is a set of productions or rewrite rules  
( $P : N \rightarrow NUT$ )

## Derivations

- A sequence of application of the rewrite rules is a derivation or a parse (e.g., deriving the string  $x + 2 - y$ )
- The process of discovering a derivation is called parsing

## Parse Trees

A parse tree for a grammar G is a tree where

- The root is the start symbol for G
- The interior nodes are non-terminals of G
- The leaf nodes are terminal symbols of G
- The children of a node T (from left to right) correspond to the symbols on the right hand side of some production for T in G

Every terminal string generated by a grammar has at least one corresponding parse tree; every valid parse tree represents a string generated by the grammar (yield of the parse tree)

## Advantages of CFGs

- Precise syntactic specification of a programming language
- Easy to understand, avoids ad hoc definition
- Easier to maintain, add new language features
- Can automatically construct efficient parser
- Parser construction reveals ambiguity, other difficulties
- Imparts structure to language
- Supports syntax-directed translation

## Calculator Example

- All variables are integers
- There are no declarations
- The only statements are assignments, input, and output
- Expressions use one of four arithmetic operators and parentheses
- Operators are left associative, with the usual precedence
- There are no unary operators

## Regular Expressions

id  $\rightarrow$  letter ( letter | digit)\*  
literal  $\rightarrow$  digit digit\*  
read, write, ":", "+", "-", "\*", "/", "(", ")"  
\$\$ /\* end of input \*/

## Grammar for Calculator

$\langle \text{pgm} \rangle \rightarrow \langle \text{stmtlist} \rangle \$\$$   
 $\langle \text{stmtlist} \rangle \rightarrow \langle \text{stmtlist} \rangle \langle \text{stmt} \rangle \mid \epsilon$   
 $\langle \text{stmt} \rangle \rightarrow \text{id} := \langle \text{expr} \rangle \mid \text{read} \langle \text{id} \rangle \mid \text{write} \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle \langle \text{add op} \rangle \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \mid \langle \text{term} \rangle \langle \text{multop} \rangle \langle \text{factor} \rangle$   
 $\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle) \mid \text{id} \mid \text{literal}$   
 $\langle \text{addop} \rangle \rightarrow + \mid -$   
 $\langle \text{multop} \rangle \rightarrow * \mid /$

## Types of derivations

- Leftmost derivation
  - The leftmost non-terminal is replaced at each step
- Rightmost derivation
  - The rightmost non-terminal is replaced at each step
- Ambiguous grammar – one with multiple leftmost (or multiple rightmost) derivations for a single sentential form

## Types of parsers

- Top-down (LL) parsers
  - Left to right, leftmost derivation
  - Starts at the root of the derivation tree and fills in
  - Predicts next state with n lookahead
- Bottom-up (LR) parsers
  - Left to right, rightmost derivation
  - Starts at the leaves and fills in
  - Start with input string, end with start symbol
  - Starts in a state valid for legal first tokens
  - Change state to encode possibilities as input is consumed
  - Use a stack to store both state and sentential form

## Top-Down Parsing

- At a node labeled A, select a production with A on its LHS and for each symbol on its RHS, construct the appropriate child
- When a terminal is added that does not match the input, backtrack
- Find the next node to be expanded (must have a label in NT)

## Eliminating Left Recursion

- A grammar is left recursive if there exists A in NT such that  $A \rightarrow A\delta$  for some string  $\delta$
- Transform the grammar to remove left recursion  
 $\langle \text{foo} \rangle \rightarrow \langle \text{foo} \rangle \delta$   
 $\quad \quad \quad \mid \mu$   
 $\rightarrow$   
 $\langle \text{foo} \rangle \rightarrow \mu \langle \text{bar} \rangle$   
 $\langle \text{bar} \rangle \rightarrow \delta \langle \text{bar} \rangle \mid \epsilon$   
where  $\langle \text{bar} \rangle$  is a new non-terminal

## Eliminating Common Prefixes

foo  $\rightarrow$  bar  $\delta$   
 $\rightarrow$  bar (  $\mu$  )  
 **$\rightarrow$**   
foo  $\rightarrow$  bar footail  
footail  $\rightarrow$   $\delta$  | (  $\mu$  )

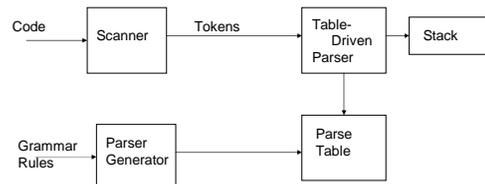
## LL Grammar for Calculator

$\langle \text{pgm} \rangle \rightarrow \langle \text{stmtlist} \rangle \$\$$   
 $\langle \text{stmtlist} \rangle \rightarrow \langle \text{stmt} \rangle \langle \text{stmtlist} \rangle \mid \epsilon$   
 $\langle \text{stmt} \rangle \rightarrow \text{id} := \langle \text{expr} \rangle \mid \text{read } \langle \text{id} \rangle \mid \text{write } \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \langle \text{termtail} \rangle$   
 $\langle \text{termtail} \rangle \rightarrow \langle \text{addop} \rangle \langle \text{term} \rangle \langle \text{termtail} \rangle \mid \epsilon$   
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \langle \text{factortail} \rangle$   
 $\langle \text{factortail} \rangle \rightarrow \langle \text{multop} \rangle \langle \text{factor} \rangle \langle \text{factortail} \rangle \mid \epsilon$   
 $\langle \text{factor} \rangle \rightarrow \langle \text{expr} \rangle \mid \text{id} \mid \text{literal}$   
 $\langle \text{addop} \rangle \rightarrow + \mid -$   
 $\langle \text{multop} \rangle \rightarrow * \mid /$

## Parser Construction

- Recursive descent parsing
  - Top-down parsing algorithm
  - Built on procedure calls (may be recursive)
  - Write procedure for each non-terminal, turning each production into clause
  - Insert call to procedure A() for non-terminal A and to match(x) for terminal x
  - Start by invoking procedure for start symbol S

## Predictive (Table-Driven) Parser



## Predictive (Table-Driven) Parsing

- Actions –
  - Match a terminal
  - Predict a production
  - Announce a syntax error
- Push as yet unseen portions of productions onto a stack
- Use –
  - FIRST(A)
  - FOLLOW(A)

## The FIRST Set

- FIRST( $\alpha$ ) is the set of terminal symbols that begin strings derived from  $\alpha$
- To build FIRST(X):
  - If X is a terminal, FIRST(X) is {X}
  - If  $X \leftarrow \epsilon$ , then  $\epsilon \in \text{FIRST}(X)$
  - If  $X \leftarrow Y_1 Y_2 \dots Y_k$  then put FIRST( $Y_1$ ) in FIRST(X)
  - If X is a non-terminal and  $X \leftarrow Y_1 Y_2 \dots Y_k$ , then a  $\epsilon \in \text{FIRST}(X)$  if a  $\epsilon \in \text{FIRST}(Y_i)$  and  $\epsilon \in \text{FIRST}(Y_j)$ , for all  $1 \leq j < i$

## The Follow Set

- For a non-terminal A, FOLLOW(A) is the set of terminals that can appear immediately to the right of A in some sentential form
- To build FOLLOW(B) for all B -
  - Starting with goal, place eof in FOLLOW(<goal>)
  - If  $A \leftarrow \alpha B \beta$ , then put  $\{FIRST(\beta) - \epsilon\}$  in FOLLOW(B)
  - If  $A \leftarrow \alpha B$ , then put FOLLOW(A) in FOLLOW(B)
  - If  $A \leftarrow \alpha B \beta$  and  $\epsilon \in FIRST(\beta)$ , then put FOLLOW(A) in FOLLOW(B)

## Using FIRST and FOLLOW

- For each production  $A \leftarrow \alpha$  and lookahead token
  - Expand A using the production if token  $\epsilon \in FIRST(\alpha)$
  - If  $\epsilon \in FIRST(\alpha)$ , expand A using the production if token  $\epsilon \in FOLLOW(A)$
  - All other tokens return error
- If there are multiple choices, the grammar is not LL(1) (predictive)

## LL(1) Grammars

A Grammar G is LL(1) if and only if, for all non-terminals A, each distinct pair of productions  $A \leftarrow \alpha$  and  $A \leftarrow \beta$  satisfy the condition  $FIRST(\alpha) \cap FIRST(\beta) = \Phi$ , i.e.,

For each set of productions  $A \leftarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$

- $FIRST(\alpha_1), FIRST(\alpha_2), \dots, FIRST(\alpha_n)$  are all pairwise disjoint
- If  $\alpha_i \leftarrow \epsilon$  for any i, then  $FIRST(\alpha_j) \cap FOLLOW(A) = \Phi$ , for all  $j \neq i$

## The Complexity of LL(1) Parsing

- Inside main loop – bounded by constant (function of symbols on RHS)
- How many times does the main loop execute?
  - Number of iterations is the number of nodes in the parse tree, which is  $N * P$  (N is the number of tokens in the input, P is the number of productions)
  - P is a constant, therefore running time is  $O(N)$

## CFGs versus Regular Expressions

- CFGs strictly more powerful than REs
  - Any language that can be generated using an RE can be generated by a CFG (proof by induction)
  - There are languages that can be generated by a CFG that cannot be generated by an RE (proof by contradiction)

## Example non-LL Grammar Construct

```
stmt → if condition then_clause else_clause
      | other_stmt
then_clause → then stmt
else_clause → else stmt | ε
→ Ambiguous – allows dangling else to be
paired with either then in
if A then if B then C else D
```

## Fix – Bottom-up parsing (LR)

stmt  $\rightarrow$  balanced\_stmt | unbalanced\_stmt  
balanced\_stmt  $\rightarrow$  if condition then balanced\_stmt  
                  else balanced\_stmt  
                          | other\_stmt  
unbalanced\_stmt  $\rightarrow$  if condition then stmt  
                          | if condition then balanced\_stmt  
                          else unbalanced\_stmt

-----  
OR

Use special disambiguating rules  $\rightarrow$  use production  
that occurs first in case of conflict