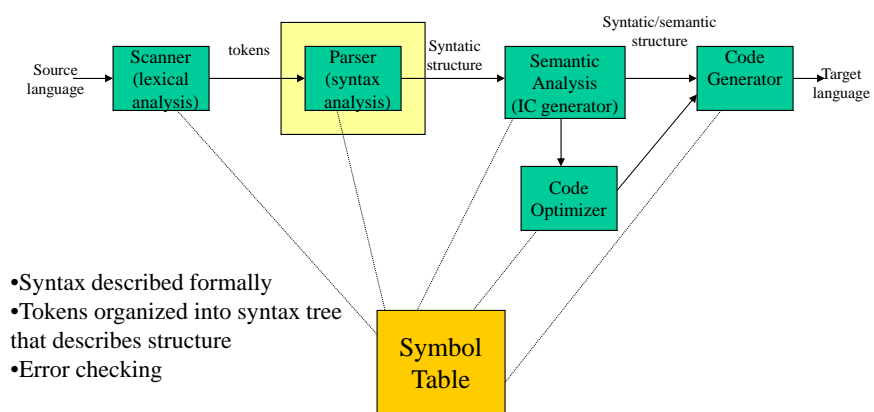


# Lecture 4a: Parsing

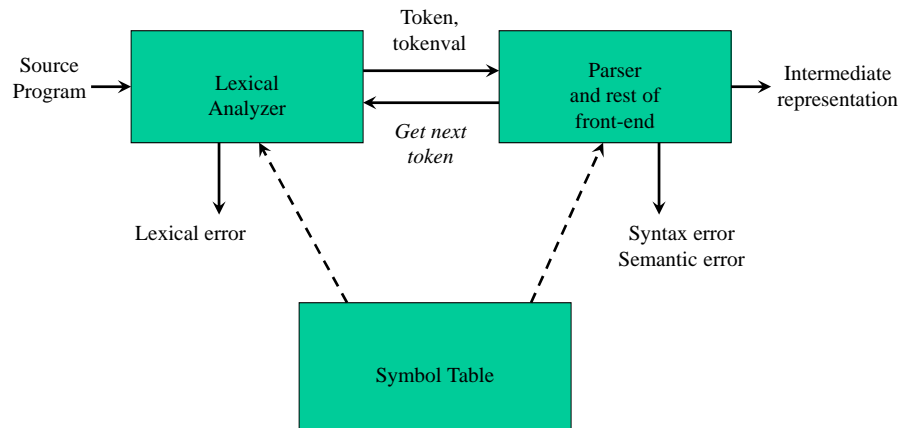
COSC 4316

(grateful acknowledgement to Robert van Engelen and Elizabeth White for some of the material from which these slides have been adapted)

## Syntax Analysis - Parsing



## Position of a Parser in the Compiler Model



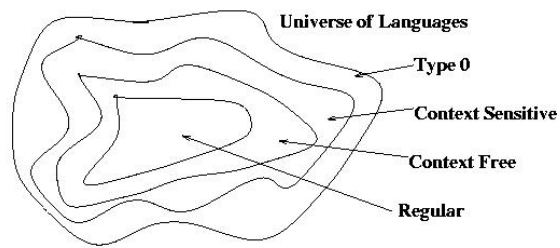
3

## Syntax Described by CF Grammars (BNF)

- Advantages of CFGs
  - Precise, easily understood syntactic specification of a language
  - Automatic construction of parsers
  - Makes syntactic ambiguities more obvious
  - Allows extension of language to be done more readily

# Static Analysis - Parsing

We can use context free grammars to specify the syntax of programming languages.



COSC 4316, Timothy J. McGuire

5

## Role of the Parser

- Obtain a string of tokens from the lexical analyzer
- Verify that the string can be generated by the grammar for the source language
- Report syntax errors (intelligibly)
- Recover from some syntax errors

COSC 4316, Timothy J. McGuire

6

## 3 Types of Parsers

- Universal parsing algorithms
  - Examples: CYK algorithm or Earley's algorithm
    - May be used for any CF grammar
    - Too inefficient for practical use
- Top down parsers (often constructed by hand)
  - Build parse trees from the root down to the leaves
  - Works with certain classes of grammars (e.g. LL grammars – *later* )
- Bottom up parsers (often build with automated tools)
  - Build parsers from the leaves up to the root
  - Work with a broader class of grammars (LR)

COSC 4316, Timothy J. McGuire

7

## Error Handling

- Errors **will** occur
- A good compiler should assist in identifying and locating errors
  - *Lexical errors*: important, compiler can easily recover and continue
  - *Syntax errors*: most important for compiler, can almost always recover
  - *Static semantic errors*: important, can sometimes recover
  - Goal: detection of an error *as soon as possible* without further consuming unnecessary input
  - How: detect an error as soon as the prefix of the input does not match a prefix of any string in the language (the *viable prefix property*)

8

## Error Recovery Strategies

- *Bail out* (sudden death)
  - Stop when first error found
- *Panic mode*
  - Discard input until a token in a set of designated *synchronizing tokens* is found (e.g., a semicolon or a `}`)
- *Phrase-level recovery*
  - Perform local correction on the input to repair the error
  - Easy to do in predictive parsing because you know what is expected (**match**)
- *Error productions*
  - Augment grammar with productions for erroneous constructs
- *Global correction*
  - Choose a minimal sequence of changes to obtain a global least-cost correction

9

## Context-Free Grammars (Recap)

- Context-free grammar is a 4-tuple  $G = (N, T, P, S)$  where
  - $T$  is a finite set of tokens (*terminal* symbols)
  - $N$  is a finite set of *nonterminals*
  - $P$  is a finite set of *productions* of the form
$$A \rightarrow \beta$$
where  $A \in N$  and  $\beta \in (N \cup T)^*$
  - $S \in N$  is a designated *start symbol*

10

## Notational Conventions Used

- Terminals  
 $a, b, c, \dots \in T$  (lowercase letters early in the alphabet)  
specific terminals: **0**, **1**, **id**, **+**
- Nonterminals (uppercase letters early in the alphabet)  
 $A, B, C, \dots \in N$   
specific nonterminals: *expr*, *term*, *stmt*
- Grammar symbols (uppercase letters late in the alphabet)  
 $X, Y, Z \in (N \cup T)$
- Strings of terminals (lowercase letters late in the alphabet)  
 $u, v, w, x, y, z \in T^*$
- Strings of grammar symbols (Greek letters)  
 $\alpha, \beta, \gamma \in (N \cup T)^*$

11

## Derivations

- The *one-step derivation* is defined by  
 $\alpha A \beta \Rightarrow \alpha \gamma \beta$   
where  $A \rightarrow \gamma$  is a production in the grammar
- In addition, we define
  - $\Rightarrow$  is *leftmost*  $\Rightarrow_{lm}$  if  $\alpha$  does not contain a nonterminal
  - $\Rightarrow$  is *rightmost*  $\Rightarrow_{rm}$  if  $\beta$  does not contain a nonterminal
  - Transitive closure  $\Rightarrow^*$  (zero or more steps)
  - Positive closure  $\Rightarrow^+$  (one or more steps)

12

## Derivations

- A derivation is an alternative to constructing a parse tree.
- We view a production as a *rewriting rule*.
- The sequence of replacements is called a *derivation*.
- If  $S \Rightarrow^* \alpha$ , then  $\alpha$  is said to be a *sentential form* of the CFG,  $G$
- The *language generated by  $G$*  is defined by
$$L(G) = \{w \in T^* \mid S \Rightarrow^+ w\}$$
and  $w$  is called a sentence in  $L(G)$   
(a sentential form with no non-terminals)

13

## Derivations

- When deriving a token sequence, if more than one nonterminal is present, we have a choice of which to replace next.
- One convention:
  - Leftmost derivation –
    - Choose the leftmost possible nonterminal at each step.
    - $\Rightarrow_{lm} \quad \Rightarrow_{lm}^* \quad \Rightarrow_{lm}^+$
- A sentential form produced via a leftmost derivation is called a ***left sentential form***

14

## Derivation (Example)

$S \rightarrow P ( S ) \mid \underline{\text{var}} R$

$P \rightarrow \underline{\text{func}} \mid \varepsilon$

$R \rightarrow + S \mid \varepsilon$

Expressions of variables and functions

- A leftmost derivation of  $\underline{\text{func}} ( \underline{\text{var}} + \underline{\text{var}} )$  is

$$\begin{aligned} S &\Rightarrow_{lm} P ( S ) \Rightarrow_{lm} \underline{\text{func}} ( S ) \Rightarrow_{lm} \underline{\text{func}} ( \underline{\text{var}} R ) \\ &\Rightarrow_{lm} \underline{\text{func}} ( \underline{\text{var}} + S ) \Rightarrow_{lm} \underline{\text{func}} ( \underline{\text{var}} + \underline{\text{var}} R ) \\ &\Rightarrow_{lm} \underline{\text{func}} ( \underline{\text{var}} + \underline{\text{var}} ) \end{aligned}$$

15

## Derivations

- Analogous to leftmost derivations, we have *rightmost derivations*.
- These seem less intuitive, but they correspond to a large class of parsers (the *bottom-up parsers*.)
- Leftmost derivations are usually associated with top-down parsing.
- Rightmost derivations are sometimes called *canonical derivations*.

16



## CFG Examples

Indicates a production

$T = \{+, -, 0..9\}, N = \{L, D\}, S = L$

$L \rightarrow L + D \mid L - D \mid D$

$D \rightarrow 0 \mid \dots \mid 9$

Shorthand for multiple productions

$T = \{ (, ) \}, N = \{ L \}, S = L$

$L \rightarrow ( L ) L$

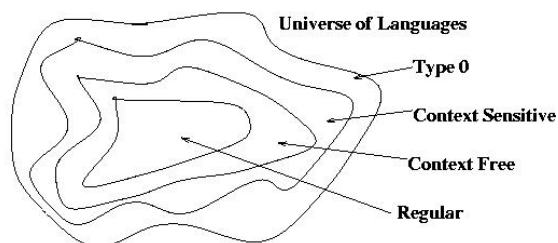
$L \rightarrow \varepsilon$

COSC 4316, Timothy J. McGuire

17

## Languages

Regular	$A \rightarrow a B, C \rightarrow \varepsilon$
Context free	$A \rightarrow \alpha$
Context sensitive	$\alpha A \beta \rightarrow \alpha \gamma \beta$
Type 0	$\alpha \rightarrow \beta$



COSC 4316, Timothy J. McGuire

18

## Any regular language can be expressed using a CFG

Starting with a NFA:

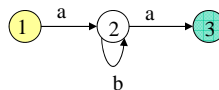
- For each state  $S_i$  in the NFA
  - Create non-terminal  $A_i$
  - If transition  $(S_i, a) = S_k$ , create production  $A_i \rightarrow a A_k$
  - If transition  $(S_i, \epsilon) = S_k$ , create production  $A_i \rightarrow A_k$
  - If  $S_i$  is a final state, create production  $A_i \rightarrow \epsilon$
  - If  $S_i$  is the NFA start state,  $s = A_i$
- *What does the existence of this algorithm tell us about the relationship between regular and context free languages?*

COSC 4316, Timothy J. McGuire

19

## NFA to CFG Example

$ab^*a$



$A_1 \rightarrow a A_2$

$A_2 \rightarrow b A_2$

$A_2 \rightarrow a A_3$

$A_3 \rightarrow \epsilon$

COSC 4316, Timothy J. McGuire

20

## Writing Grammars

When writing a grammar (or RE) for some language, the following must be true:

1. All strings generated are in the language.
2. Your grammar produces all strings in the language.

## Try these:

- Integers divisible by 2
- Legal postfix expressions
- Floating point numbers with no extra zeros
- Strings of 0,1 where there are more 0 than 1 (hard)

## Regular Expressions vs. CFGs

- A grammar in which every production is of the form:

$$A \rightarrow w B \quad (w, x \in T^* \text{ and } A, B \in N)$$

or  $A \rightarrow x$

is called a ***right-linear grammar***. (*Left-linear grammar* defined analogously.)

- It can be proven that any right-linear grammar generates a regular language, and *vice versa*.

## Regular Expressions vs. CFGs

- Why use regular expressions to denote the lexical syntax of a language if we could use CFGs instead?
  - Lexical rules are simple – don't use a chainsaw to prune a rose
  - Regular expressions are more concise and easier to understand than CFGs
  - Easier to generate a lexical analyzer from a regular expression than an arbitrary grammar.
  - Promotes modularity of the front end.

# Parsing

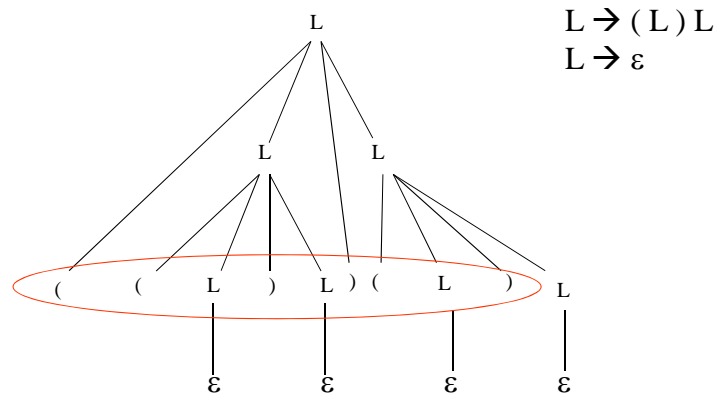
- The task of parsing is figuring out what the **parse tree** looks like for a given input and language.
- If a string is in the given language, a parse tree must exist.
- However, just because a parse tree exists for some string in a given language doesn't mean a given algorithm can find it.

# Parse Trees

The parse tree for some string in a language that is defined by the grammar  $G$  as follows:

- The root is the start symbol of  $G$
- The leaves are terminals or  $\epsilon$ . When visited from left to right, the leaves form the input string
- The interior nodes are non-terminals of  $G$
- For every non-terminal  $A$  in the tree with children  $B_1 \dots B_k$ , there is some production  $A \rightarrow B_1 \dots B_k$

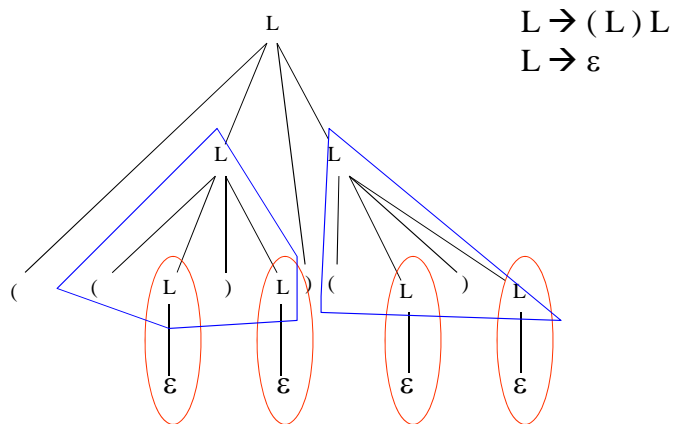
## Parse Tree for $((\ ))(\ ))$



COSC 4316, Timothy J. McGuire

27

## Parse Tree for $((\ ))(\ ))$



COSC 4316, Timothy J. McGuire

28

## Parse Trees & Derivations

- A derivation is a *linear representation* of a parse tree.
- Equivalently, a parse tree is a *graphical representation* of a derivation.

For the grammar and the string

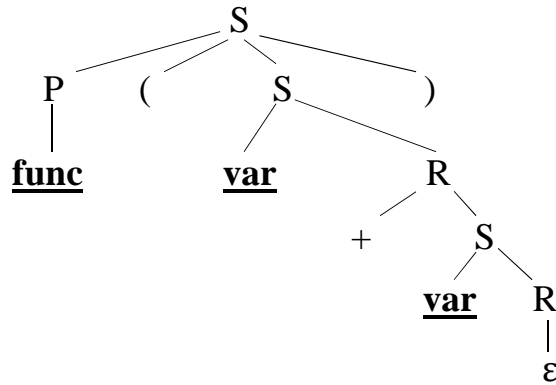
func ( var + var ):

$S \rightarrow P(S) \mid \text{var } R$

$P \rightarrow \text{func} \mid \varepsilon$

$R \rightarrow + S \mid \varepsilon$

$S \Rightarrow_{lm} P(S) \Rightarrow_{lm} \text{func}(S) \Rightarrow_{lm} \text{func}(\text{var } R)$   
 $\Rightarrow_{lm} \text{func}(\text{var} + S)$   
 $\Rightarrow_{lm} \text{func}(\text{var} + \text{var } R)$   
 $\Rightarrow_{lm} \text{func}(\text{var} + \text{var})$



COSC 4316, Timothy J. McGuire

29

## Single Step Derivation

**Definition:** Given  $\alpha A \beta$  (with  $\alpha, \beta$  in  $(V_n \cup V_t)^*$ ) and a production  $A \rightarrow \gamma$ ,

$\alpha A \beta \Rightarrow \alpha \gamma \beta$  is a **single step derivation**.

Examples:

$L + D \Rightarrow L - D + D$

$L \rightarrow L - D$

$(L)(L) \Rightarrow ((L)L)(L)$

$L \rightarrow (L)L$

Greek letters ( $\alpha, \beta, \gamma, \dots$ ) denote a (possibly empty) sequence of terminals and non-terminals.

COSC 4316, Timothy J. McGuire

30

# Derivations

**Definition:** A sequence of the form:

$$w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_n$$

is a **derivation** of  $w_n$  from  $w_0$  ( $w_0 \Rightarrow^* w_n$ )

$L$	production $L \rightarrow (L)L$
$\Rightarrow (L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow ()L$	production $L \rightarrow \varepsilon$
$\Rightarrow ()$	

$$L \Rightarrow^* ()$$

If  $w_1$  has non-terminal symbols, it is referred to as *sentential form*.

$$L \Rightarrow^* (( )) ( )$$

$L$	production $L \rightarrow (L)L$
$\Rightarrow (L)L$	production $L \rightarrow (L)L$
$\Rightarrow (L)(L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow (L)(L)$	production $L \rightarrow (L)L$
$\Rightarrow ((L)L)(L)$	production $L \rightarrow \varepsilon$
$\Rightarrow ((L)L)(L)$	production $L \rightarrow \varepsilon$
$\Rightarrow ((L)L)()$	production $L \rightarrow \varepsilon$
$\Rightarrow (( )) ( )$	



- $L(G)$ , the language generated by grammar  $G$  is  
 $\{w \text{ in } T^*: S \Rightarrow^* w, \text{ for start symbol } S\}$
- Both  $()$  and  $(( ))()$  are in  $L(G)$  for the following grammar.
  - $L \rightarrow ( L ) L$
  - $L \rightarrow \varepsilon$

## Leftmost Derivations

- Recall that a leftmost derivation is one where the leftmost nonterminal is always chosen
- If a string is in a given language (i.e. a derivation exists), then a leftmost derivation *must* exist
- Rightmost derivation defined as you would expect

## Leftmost Derivation for $((()))$

<b>L</b>	production $L \rightarrow (L)L$
$\Rightarrow (L)L$	production $L \rightarrow (L)L$
$\Rightarrow ((L)L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow ((L)L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow (((L)L)L)L$	production $L \rightarrow (L)L$
$\Rightarrow (((L)L)L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow (((L)L)L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow (((L)L)L)L$	

$L \rightarrow (L)L$   
 $L \rightarrow \varepsilon$

COSC 4316, Timothy J. McGuire

35

## Rightmost Derivation for $((()))$

<b>L</b>	production $L \rightarrow (L)L$
$\Rightarrow (L)L$	production $L \rightarrow (L)L$
$\Rightarrow (L)(L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow (L)(L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow (L)(L)L$	production $L \rightarrow (L)L$
$\Rightarrow ((L)L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow ((L)L)L$	production $L \rightarrow \varepsilon$
$\Rightarrow (((L)L)L)L$	

$L \rightarrow (L)L$   
 $L \rightarrow \varepsilon$

COSC 4316, Timothy J. McGuire

36

# Ambiguity

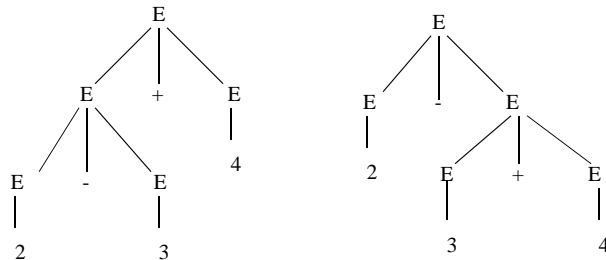
- An **ambiguous** grammar is one in which two (or more) parse trees or leftmost derivations exist for *some string in the language*

$E \rightarrow E + E$

$E \rightarrow E - E$

$E \rightarrow 0 \mid \dots \mid 9$

$2 - 3 + 4$



COSC 4316, Timothy J. McGuire

37

- Two leftmost derivations

$E \Rightarrow E + E$   
 $\Rightarrow E - E + E$   
 $\Rightarrow 2 - E + E$   
 $\Rightarrow 2 - 3 + E$   
 $\Rightarrow 2 - 3 + 4$

$E \Rightarrow E - E$   
 $\Rightarrow 2 - E$   
 $\Rightarrow 2 - E + E$   
 $\Rightarrow 2 - 3 + E$   
 $\Rightarrow 2 - 3 + 4$

- We must either write unambiguous grammars **or** have *disambiguating* rules.

COSC 4316, Timothy J. McGuire


38

- An ambiguous grammar can sometimes be made unambiguous:

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

$$T \rightarrow 0 \mid \dots \mid 9$$

enforces the correct associativity



- Precedence can be specified as well:

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

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

$$F \rightarrow ( E ) \mid 0 \mid \dots \mid 9$$

## Another example of ambiguity

- **if-else** in Java

```
if (expr)
  if (expr)
    stmt;
  else
    stmt;
```

- Which **if** is the **else** associated with?
- The last one, but we can't specify that via the definition.
- Could fix this by requiring the use of an **endif** keyword

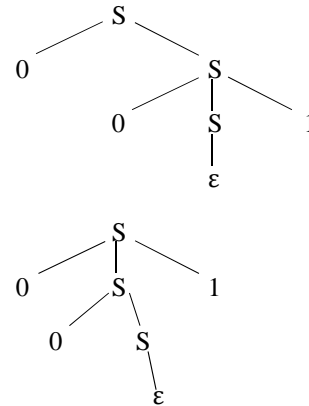
## Yet another example of ambiguity

- $L = \{ 0^i 1^j \mid i \geq j \geq 0 \}$  ( $0 \equiv \text{if}, 1 \equiv \text{else}$ )
- Expressed by the grammar

$$S \rightarrow 0 S \mid 0 S 1 \mid \epsilon$$

Parse trees for **001**

Disambiguating rule: match each 1 with the closest unmatched 0



COSC 4316, Timothy J. McGuire

41

## Disambiguating rule incorporated into the grammar

- $S \rightarrow 0 S \mid A$
- $A \rightarrow 0 A 1 \mid \epsilon$

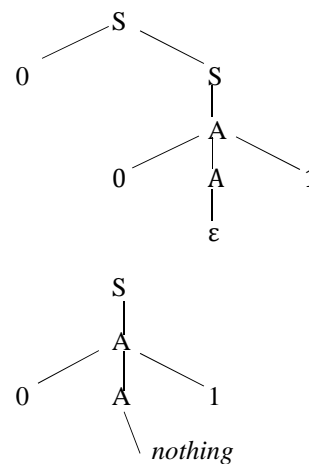
Parse trees for **001**

Disambiguating rule: match each 1 with the closest unmatched 0

Grammar still has issues, because we can't use a predictive parser – can't predict whether the 0 is matched or unmatched.

The dangling else is really a language design issue.

**Conclusion: If you ever design a programming language, you need to know the issues involved in parsing that language!**

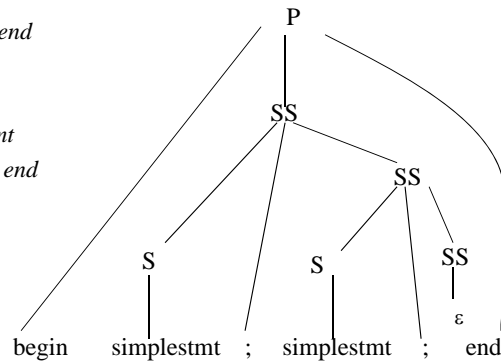


COSC 4316, Timothy J. McGuire

42

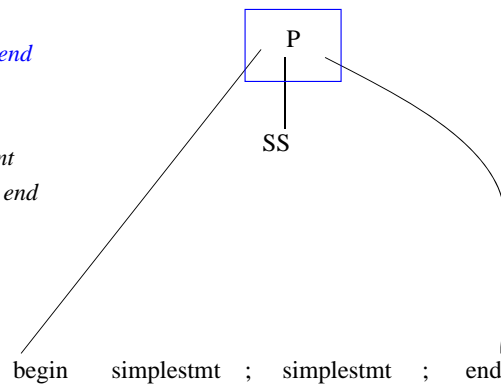
Input: **begin simplestmt;  
simplestmt; end**

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$



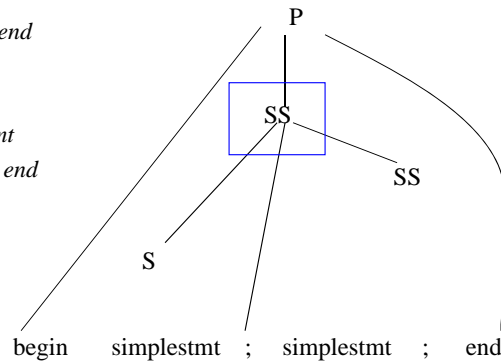
## Top Down (LL) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$



# Top Down (LL) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$

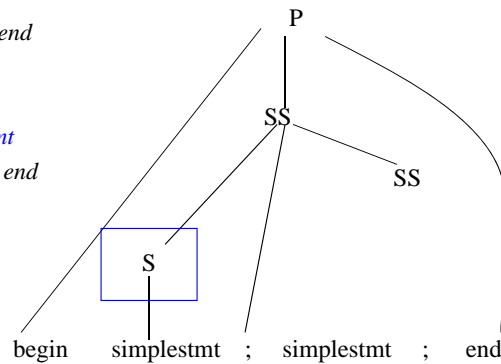


COSC 4316, Timothy J. McGuire

45

# Top Down (LL) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$

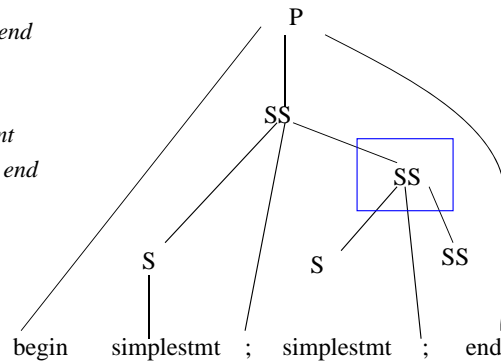


COSC 4316, Timothy J. McGuire

46

# Top Down (LL) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$

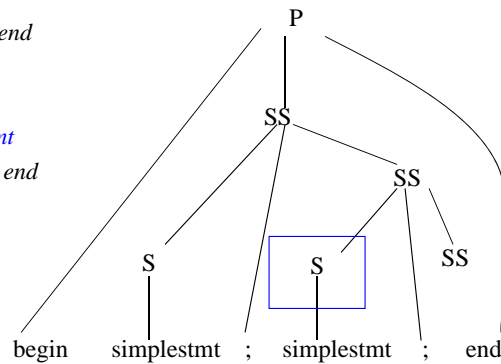


COSC 4316, Timothy J. McGuire

47

# Top Down (LL) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$



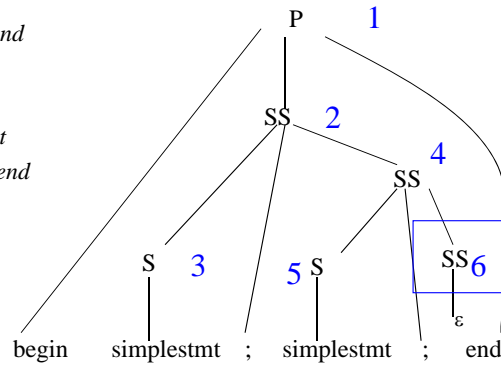
COSC 4316, Timothy J. McGuire

48



# Top Down (LL) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$

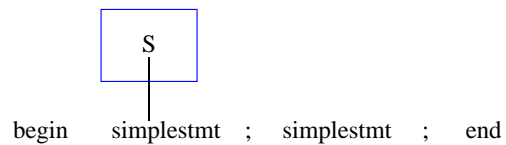


COSC 4316, Timothy J. McGuire

49

# Bottomup (LR) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$



COSC 4316, Timothy J. McGuire

50

## Bottomup (LR) Parsing

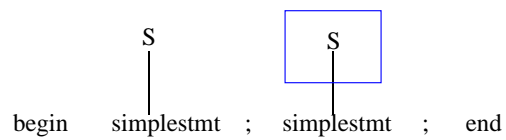
$P \rightarrow \text{begin } SS \text{ end}$

$SS \rightarrow S ; SS$

$SS \rightarrow \epsilon$

$S \rightarrow \text{simplestmt}$

$S \rightarrow \text{begin } SS \text{ end}$



COSC 4316, Timothy J. McGuire

51

## Bottomup (LR) Parsing

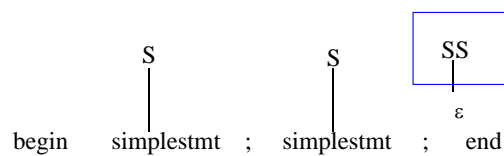
$P \rightarrow \text{begin } SS \text{ end}$

$SS \rightarrow S ; SS$

$SS \rightarrow \epsilon$

$S \rightarrow \text{simplestmt}$

$S \rightarrow \text{begin } SS \text{ end}$

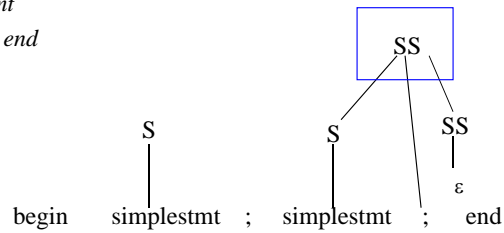


COSC 4316, Timothy J. McGuire

52

# Bottomup (LR) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$

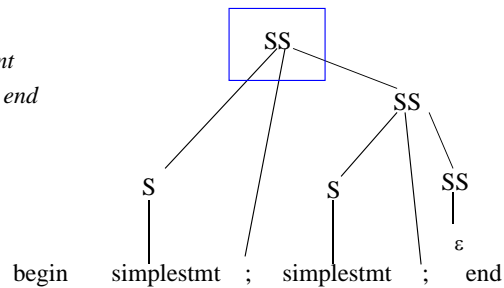


COSC 4316, Timothy J. McGuire

53

# Bottomup (LR) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$

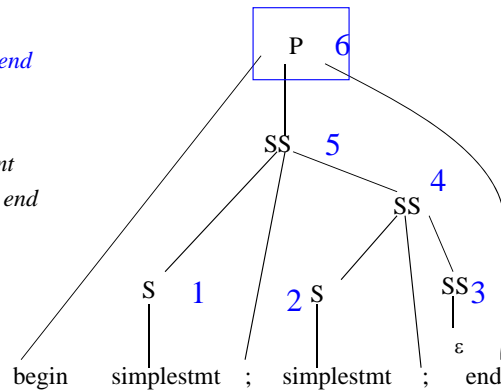


COSC 4316, Timothy J. McGuire

54

## Bottomup (LR) Parsing

$P \rightarrow \text{begin } SS \text{ end}$   
 $SS \rightarrow S ; SS$   
 $SS \rightarrow \epsilon$   
 $S \rightarrow \text{simplestmt}$   
 $S \rightarrow \text{begin } SS \text{ end}$



COSC 4316, Timothy J. McGuire

55

## Left Recursion (Recap)

- Productions of the form
 
$$A \rightarrow A \alpha$$

$$/ \beta$$
 are left recursive
- When one of the productions in a grammar is left recursive then a predictive parser loops forever on certain inputs

56

## Left Recursion (Recap)

- Replace

$$A \rightarrow A \alpha$$
$$/ \beta$$

with

$$A \rightarrow \beta R$$
$$R \rightarrow \alpha R / \varepsilon$$

57

## Indirect Left Recursion

- What about indirect left recursion?

$$B \rightarrow D \alpha$$

$$D \rightarrow B \beta$$

58

## General Left Recursion Elimination Method

Arrange the nonterminals in some order  $A_1, A_2, \dots, A_n$

**for**  $i = 1, \dots, n$  **do**

**for**  $j = 1, \dots, i-1$  **do**

**for** each production  $A_i \rightarrow A_j \alpha$  and  $A_j \rightarrow \beta$  **do**

            Replace  $A_i \rightarrow A_j \alpha$  with  $A_i \rightarrow \beta \alpha$

**endfor**

**endfor**

    eliminate any direct *left recursion* among  $A_i$

**endfor**

59

## Example Left Recursion Elimination

$$\left. \begin{array}{l} A \rightarrow B \mathbf{a} \mid \mathbf{b} C \\ B \rightarrow B \mathbf{c} \mid A \mathbf{d} \\ C \rightarrow \mathbf{e} \mid \mathbf{f} \end{array} \right\} \text{Choose arrangement: } A = A_1, B = A_2, C = A_3$$

A unchanged

$B \rightarrow B \mathbf{c} \mid B \mathbf{a} \mathbf{d} \mid \mathbf{b} C \mathbf{d}$

$B \rightarrow \mathbf{b} C \mathbf{d} D$  (*eliminating direct left recursion*)

$D \rightarrow \mathbf{c} D \mid \mathbf{a} \mathbf{d} D \mid \varepsilon$

C unchanged

Result:

$A \rightarrow B \mathbf{a} \mid \mathbf{b} C$

$B \rightarrow \mathbf{b} C \mathbf{d} D$

$C \rightarrow \mathbf{e} \mid \mathbf{f}$

$D \rightarrow \mathbf{c} D \mid \mathbf{a} \mathbf{d} D \mid \varepsilon$

60

## Left Factoring

- Most problems with predictive parsing are either (a) left recursion (which we have already dealt with) or (b) common prefixes
- Example:
  - $stmt \rightarrow \text{if } expr \text{ then } seq\text{-of-}stmts \text{ endif}$
  - $stmt \rightarrow \text{if } expr \text{ then } seq\text{-of-}stmts \text{ else } seq\text{-of-}stmts \text{ endif}$
- Solution: Rewrite the production to defer the decision until we have enough information to make the right choice.

61

## Left Factoring

- Replace productions
$$A \rightarrow \alpha \beta \mid \alpha \gamma$$
with
$$A \rightarrow \alpha A_R$$
$$A_R \rightarrow \beta \mid \gamma$$
- e.g.
  - $stmt \rightarrow \text{if } expr \text{ then } seq\text{-of-}stmts \text{ opt\_end}$
  - $opt\_end \rightarrow \text{endif} \mid \text{else } seq\text{-of-}stmts \text{ endif}$

62

## Non-context-free Language Constructs

- Not all syntactic rules are expressible using CFGs.
  - *e.g.*, “variables must be declared before they are used” cannot be expressed in a CFG.
  - *or*, “the number of formal parameters for a function must equal the number of actual parameters.”
- In practice, syntactic details that cannot be represented in a CFG are considered part of the *static semantics* and deferred to the semantic analysis phase.

COSC 4316, Timothy J. McGuire

63

## Top-Down Parsing

- General Algorithms:
  - LL (top down)
    - (We looked at elementary recursive descent parsing in module 2)
  - LR (bottom up)
- Both algorithms are driven by the input grammar and the input to be parsed
- LL(1) grammars are those suitable for predictive parsing
- “LL(1)”  $\equiv$  scans input from Left to right, Leftmost derivation, 1 token lookahead.

An LL(1) grammar can always be parsed top-down without backtracking.

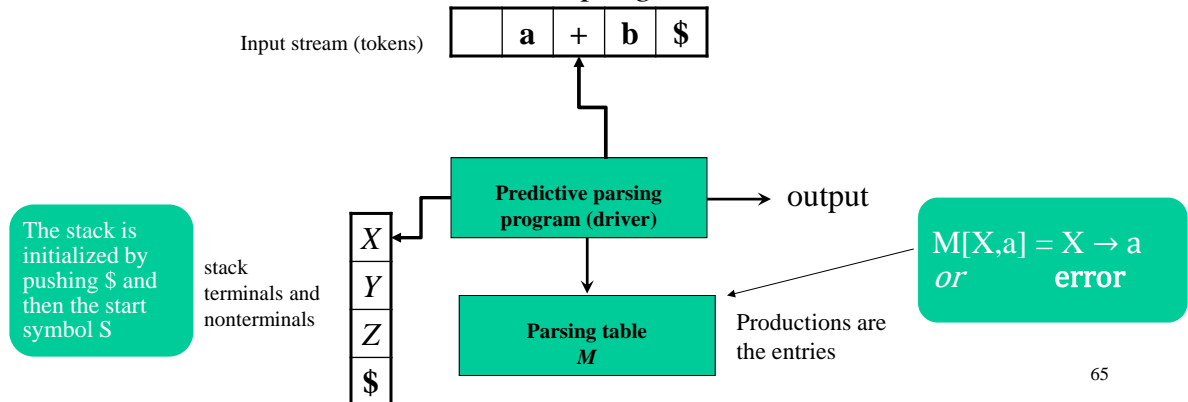
COSC 4316, Timothy J. McGuire

64



# Non-Recursive Predictive Parsing: Table-Driven Parsing

- Given an LL(1) grammar  $G = (N, T, P, S)$  construct a table  $M[A, a]$  for  $A \in N, a \in T$  and use a *driver program* with a *stack*



65

## Predictive Parsing Algorithm

- If  $X = a = \$$   $\cdots$  parser halts. Success!
- If  $X = a \neq \$$ 
  - Pop  $X$  off stack
  - Advance input pointer
- If  $X$  is a non-terminal, look up  $M[X, a]$ 
  - If, for example  $M[X, a] = X \rightarrow UVW$ , push  $WVU$  on the stack ( $U$  on top)

## Predictive Parsing Program (Driver)

```

push($)
push(S)
a := lookahead
repeat
    X := pop()
    if X is a terminal or X = $ then
        match(X)           // moves to next token and a := lookahead
    else if M[X,a] = X → Y1Y2...Yk then
        push(Yk, Yk-1, ..., Y2, Y1)    // such that Y1 is on top
        ... invoke actions and/or produce IR output ...
    else
        error()
    endif
until X = $
    
```

67


## Example Table

$E \rightarrow E + T / T$   
 $T \rightarrow T * F / F$   
 $F \rightarrow ( E ) | \text{id}$

Eliminate left recursion

$E \rightarrow T E_R$   
 $E_R \rightarrow + T E_R | \epsilon$   
 $T \rightarrow F T_R$   
 $T_R \rightarrow * F T_R | \epsilon$   
 $F \rightarrow ( E ) | \text{id}$

Nevermind how the  
parse table is built just  
yet



	id	+	*	(	)	\$
E	$E \rightarrow T E_R$			$E \rightarrow T E_R$		
E <sub>R</sub>		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
T	$T \rightarrow F T_R$			$T \rightarrow F T_R$		
T <sub>R</sub>		$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow ( E )$		

68

# Parsing Example

	id	+	*	(	)	\$
$E$	$E \rightarrow T E_R$			$E \rightarrow T$ $E_R$		
$E_R$		$E_R \rightarrow + T E_R$			$E_R \rightarrow \epsilon$	$E_R \rightarrow \epsilon$
$T$	$T \rightarrow F T_R$			$T \rightarrow F T_R$		
$T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \epsilon$	$T_R \rightarrow \epsilon$
$F$	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

Stack	Input	Production applied
$\$E$	<u>id</u> +id*id\$	$E \rightarrow T E_R$
$\$E_R T$	<u>id</u> +id*id\$	$T \rightarrow F T_R$
$\$E_R T_R F$	<u>id</u> +id*id\$	$F \rightarrow \text{id}$
$\$E_R T_R \text{id}$	<u>id</u> +id*id\$	
$\$E_R T_R$	+ <u>id</u> *id\$	$T_R \rightarrow \epsilon$
$\$E_R$	+ <u>id</u> *id\$	$E_R \rightarrow + T E_R$
$\$E_R T$	+ <u>id</u> *id\$	
$\$E_R T$	<u>id</u> *id\$	$T \rightarrow F T_R$
$\$E_R T_R F$	<u>id</u> *id\$	$F \rightarrow \text{id}$
$\$E_R T_R \text{id}$	<u>id</u> *id\$	
$\$E_R T_R$	* <u>id</u> \$	$T_R \rightarrow * F T_R$
$\$E_R T_R F$	* <u>id</u> \$	
$\$E_R T_R F$	<u>id</u> \$	$F \rightarrow \text{id}$
$\$E_R T_R \text{id}$	<u>id</u> \$	
$\$E_R T_R$	\$	$T_R \rightarrow \epsilon$
$\$E_R$	\$	$E_R \rightarrow \epsilon$
$\$$	\$	

- So, the big question:

**How do we compute the parse tables?**

- We're going to need some more theory, folks, so hang on to your hats.

## FIRST Sets

$\text{FIRST}(\alpha)$  is the set of all terminal symbols that can begin some sentential form in a derivation that starts with  $\alpha$

$\alpha \Rightarrow \dots \Rightarrow a \beta$   
(also include  $\varepsilon$  if  $\alpha \Rightarrow^* \varepsilon$ )

- $\text{FIRST}(\alpha) = \{ \mathbf{a} \in T \mid \alpha \Rightarrow^* \mathbf{a} \beta \} \cup \{ \varepsilon \text{ if } \alpha \Rightarrow^* \varepsilon \}$
- Example:

*simple*  $\rightarrow$  **integer** | **char** | **num dotdot num**

$\text{FIRST}(\textit{simple}) = \{ \mathbf{integer}, \mathbf{char}, \mathbf{num} \}$

## Computing FIRST

- To compute  $\text{FIRST}(X)$  for any single grammar symbol  $X$ :
  1. If  $X$  is a terminal,  $\text{FIRST}(X) = \{X\}$
  2. If  $X \rightarrow \varepsilon$  is a production, add  $\varepsilon$  to  $\text{FIRST}(X)$
  3. If  $X$  is a nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_n$  is a production, add **a** to  $\text{FIRST}(X)$  **if**  $a \in \text{FIRST}(Y_i)$  **and**  
 $\varepsilon \in \text{FIRST}(Y_1) \cap \text{FIRST}(Y_2) \cap \dots \cap \text{FIRST}(Y_{i-1})$   
(i.e,  $Y_1 Y_2 \dots Y_{i-1} \Rightarrow^* \varepsilon$ )

# Computing FIRST

- To compute  $\text{FIRST}(\alpha)$  where  $\alpha = X_1 X_2 \dots X_n$  :
  1. Add non- $\epsilon$  symbols of  $\text{FIRST}(X_1)$
  2. If  $\epsilon \in \text{FIRST}(X_1)$  then add non- $\epsilon$  symbols of  $\text{FIRST}(X_2)$ .
  3. As long as  $\epsilon \in \text{FIRST}(X_{i-1})$  then add non- $\epsilon$  symbols of  $\text{FIRST}(X_i)$ .
  4. If  $\epsilon$  is a member of **all** the  $\text{FIRST}(X_i)$  sets then add  $\epsilon$  to  $\text{FIRST}(\alpha)$

## Example 1

- $S \rightarrow a S e$
  - $S \rightarrow B$
  - $B \rightarrow b B e$
  - $B \rightarrow C$
  - $C \rightarrow \underline{c} C e$
  - $C \rightarrow \underline{d}$
- $\text{FIRST}(C) = \{c, d\}$
  - $\text{FIRST}(B) =$
  - $\text{FIRST}(S) =$
- Start with the 'simplest' non-terminal
-

## Example 1

- $S \rightarrow a S e$
- $S \rightarrow B$
- $B \rightarrow \underline{b} B e$
- $B \rightarrow \underline{C}$
- $C \rightarrow c C e$
- $C \rightarrow d$

- $\text{FIRST}(C) = \{c, d\}$
- $\text{FIRST}(B) = \{b, c, d\}$
- $\text{FIRST}(S) =$

Now that we know  $\text{FIRST}(C) \dots$



## Example 1

- $S \rightarrow \underline{a} S e$
- $S \rightarrow \underline{B}$
- $B \rightarrow b B e$
- $B \rightarrow C$
- $C \rightarrow c C e$
- $C \rightarrow d$

- $\text{FIRST}(C) = \{c, d\}$
- $\text{FIRST}(B) = \{b, c, d\}$
- $\text{FIRST}(S) = \{a, b, c, d\}$

## Example 2

- $P \rightarrow \underline{i} \mid \underline{c} \mid \underline{n} T S$
- $Q \rightarrow P \mid a S \mid d S c S T$
- $R \rightarrow \underline{b} \mid \underline{\epsilon}$
- $S \rightarrow e \mid R n \mid \epsilon$
- $T \rightarrow R S q$
- $\text{FIRST}(P) = \{i, c, n\}$
- $\text{FIRST}(Q) =$
- $\text{FIRST}(R) = \{b, \epsilon\}$
- $\text{FIRST}(S) =$
- $\text{FIRST}(T) =$

## Example 2

- $P \rightarrow i \mid c \mid n T S$
- $Q \rightarrow \underline{P} \mid \underline{a} S \mid \underline{d} S c S T$
- $R \rightarrow b \mid \epsilon$
- $S \rightarrow e \mid R n \mid \epsilon$
- $T \rightarrow R S q$
- $\text{FIRST}(P) = \{i, c, n\}$
- $\text{FIRST}(Q) = \{i, c, n, a, d\}$
- $\text{FIRST}(R) = \{b, \epsilon\}$
- $\text{FIRST}(S) =$
- $\text{FIRST}(T) =$

## Example 2

- $P \rightarrow i \mid c \mid n \ T \ S$
  - $Q \rightarrow P \mid a \ S \mid d \ S \ c \ S \ T$
  - $R \rightarrow b \mid \epsilon$
  - $S \rightarrow \underline{e} \mid \underline{R \ n} \mid \underline{\epsilon}$
  - $T \rightarrow R \ S \ q$
- $FIRST(P) = \{i, c, n\}$
  - $FIRST(Q) = \{i, c, n, a, d\}$
  - $FIRST(R) = \{b, \epsilon\}$
  - $FIRST(S) = \{e, b, n, \epsilon\}$
  - $FIRST(T) = \{b, c, n, q\}$

Note:

$S \Rightarrow R \ n \Rightarrow n$  because  $R \Rightarrow^* \epsilon$

COSC 4316, Timothy J. McGuire

79

## Example 2

- $P \rightarrow i \mid c \mid n \ T \ S$
  - $Q \rightarrow P \mid a \ S \mid d \ S \ c \ S \ T$
  - $R \rightarrow b \mid \epsilon$
  - $S \rightarrow e \mid R \ n \mid \epsilon$
  - $T \rightarrow \underline{R \ S} \ q$
- $FIRST(P) = \{i, c, n\}$
  - $FIRST(Q) = \{i, c, n, a, d\}$
  - $FIRST(R) = \{b, \epsilon\}$
  - $FIRST(S) = \{e, b, n, \epsilon\}$
  - $FIRST(T) = \{b, c, n, q\}$

Note:

$T \Rightarrow R \ S \ q \Rightarrow S \ q \Rightarrow q$   
because both  $R$  and  $S \Rightarrow^* \epsilon$

COSC 4316, Timothy J. McGuire

80



## Example 3

- $S \rightarrow a S e \mid S T S$
  - $T \rightarrow R S e \mid Q$
  - $R \rightarrow r S r \mid \varepsilon$
  - $Q \rightarrow S T \mid \varepsilon$
- $\text{FIRST}(S) =$
  - $\text{FIRST}(R) =$
  - $\text{FIRST}(T) =$
  - $\text{FIRST}(Q) =$

## Example 3

- $S \rightarrow a S e \mid S T S$
  - $T \rightarrow R S e \mid Q$
  - $R \rightarrow r S r \mid \varepsilon$
  - $Q \rightarrow S T \mid \varepsilon$
- $\text{FIRST}(S) = \{a\}$
  - $\text{FIRST}(R) = \{r, \varepsilon\}$
  - $\text{FIRST}(T) = \{r, a, \varepsilon\}$
  - $\text{FIRST}(Q) = \{a, \varepsilon\}$

## FOLLOW Sets

- FOLLOW(A) is the set of terminals (including end of file - \$) that may follow non-terminal A in some sentential form.
- $\text{FOLLOW}(A) = \{a \in T \mid S \Rightarrow^+ \alpha A a \beta\} \cup \{\$\}$  if  $S \Rightarrow^+ \gamma A$
- For example, consider  $L \Rightarrow^+ (())(L)L$   
Both '(' and end of file can follow L
- NOTE:  $\epsilon$  is *never* in FOLLOW sets

## Computing FOLLOW(A)

1. If S is the start symbol, put \$ in FOLLOW(S)
2. Productions of the form  $B \rightarrow \alpha A \beta$ ,  
Add  $\text{FIRST}(\beta) - \{\epsilon\}$  to FOLLOW(A)

INTUITION: Suppose  $B \rightarrow AX$  and  $\text{FIRST}(X) = \{c\}$

$$S \Rightarrow^+ \alpha B \beta \Rightarrow^+ \alpha A X \beta \Rightarrow^+ \alpha A c \delta \beta$$

= FIRST(X)



3. Productions of the form  $B \rightarrow \alpha A$  or

$B \rightarrow \alpha A \beta$  where  $\beta \Rightarrow^* \epsilon$  (i.e.,  $\epsilon \in \text{FIRST}(\beta)$ )

Add everything in  $\text{FOLLOW}(B)$  to  $\text{FOLLOW}(A)$

---

**INTUITION:**

- Suppose  $B \rightarrow Y A$

$S \Rightarrow^+ \alpha B \beta \Rightarrow \alpha Y A \beta$

$\text{FOLLOW}(B)$

- Suppose  $B \rightarrow A X$  and  $X \Rightarrow^* \epsilon$

$S \Rightarrow^+ \alpha B \beta \Rightarrow \alpha A X \beta \Rightarrow^* \alpha A \beta$

$\text{FOLLOW}(B)$

COSC 4316, Timothy J. McGuire

85

Assume the first non-terminal is  
the start symbol

## Example 4

- $S \rightarrow a S e \mid B$
- $B \rightarrow b B C f \mid C$
- $C \rightarrow c C g \mid d \mid \epsilon$
- $\text{FOLLOW}(C) =$
- $\text{FOLLOW}(B) =$
- $\text{FIRST}(C) = \{c, d, \epsilon\}$
- $\text{FOLLOW}(S) = \{\$ \}$
- $\text{FIRST}(B) = \{b, c, d, \epsilon\}$
- $\text{FIRST}(S) = \{a, b, c, d, \epsilon\}$

Using rule #1

COSC 4316, Timothy J. McGuire

86

## Example 4

- $S \rightarrow a \underline{S} e \mid B$
- $B \rightarrow b \underline{B} \underline{C} f \mid C$
- $C \rightarrow c \underline{C} g \mid d \mid \varepsilon$
- $\text{FOLLOW}(C) = \{f, g\}$
- $\text{FOLLOW}(B) = \{c, d, f\}$
- $\text{FIRST}(C) = \{c, d, \varepsilon\}$
- $\text{FOLLOW}(S) = \{\$, e\}$
- $\text{FIRST}(B) = \{b, c, d, \varepsilon\}$
- $\text{FIRST}(S) = \{a, b, c, d, \varepsilon\}$

Using rule #2

## Example 4

- $S \rightarrow a S e \mid \underline{B}$
- $B \rightarrow b B C f \mid \underline{C}$
- $C \rightarrow c C g \mid d \mid \varepsilon$
- $\text{FOLLOW}(C) = \{f, g\} \cup \text{FOLLOW}(B) = \{c, d, e, f, g, \$\}$
- $\text{FOLLOW}(B) = \{c, d, f\} \cup \text{FOLLOW}(S) = \{c, d, e, f, \$\}$
- $\text{FOLLOW}(S) = \{\$, e\}$
- $\text{FIRST}(C) = \{c, d, \varepsilon\}$
- $\text{FIRST}(B) = \{b, c, d, \varepsilon\}$
- $\text{FIRST}(S) = \{a, b, c, d, \varepsilon\}$

Using rule #3

## Example 5

- $S \rightarrow A B C \mid A D$
- $A \rightarrow \varepsilon \mid a A$
- $B \rightarrow b \mid c \mid \varepsilon$
- $C \rightarrow D d C$
- $D \rightarrow e b \mid f c$
- $FOLLOW(S) =$
- $FOLLOW(A) =$
- $FOLLOW(B) =$
- $FOLLOW(C) =$
- $FOLLOW(D) =$
- $FIRST(D) = \{e, f\}$
- $FIRST(C) = \{e, f\}$
- $FIRST(B) = \{b, c, \varepsilon\}$
- $FIRST(A) = \{a, \varepsilon\}$
- $FIRST(S) = \{a, b, c, e, f\}$

COSC 4316, Timothy J. McGuire

89

## Example 5

- $S \rightarrow A B C \mid A D$
- $A \rightarrow \varepsilon \mid a A$
- $B \rightarrow b \mid c \mid \varepsilon$
- $C \rightarrow D d C$
- $D \rightarrow e b \mid f c$
- $FOLLOW(S) = \{\$ \}$
- $FOLLOW(A) = \{b, c, e, f\}$
- $FOLLOW(B) = \{e, f\}$
- $FOLLOW(C) = \{\$ \}$
- $FOLLOW(D) = \{\$ \}$
- $FIRST(D) = \{e, f\}$
- $FIRST(C) = \{e, f\}$
- $FIRST(B) = \{b, c, \varepsilon\}$
- $FIRST(A) = \{a, \varepsilon\}$
- $FIRST(S) = \{a, b, c, e, f\}$

COSC 4316, Timothy J. McGuire

90

## Example 6

- $S \rightarrow ( A ) \mid \varepsilon$
- $A \rightarrow T E$
- $E \rightarrow \& T E \mid \varepsilon$
- $T \rightarrow ( A ) \mid a \mid b \mid c$
- $\text{FOLLOW}(S) =$
- $\text{FOLLOW}(A) =$
- $\text{FOLLOW}(E) =$
- $\text{FOLLOW}(T) =$
- $\text{FIRST}(T) = \{(, a, b, c\}$
- $\text{FIRST}(E) = \{\&, \varepsilon\}$
- $\text{FIRST}(A) = \{(, a, b, c\}$
- $\text{FIRST}(S) = \{(, \varepsilon\}$

COSC 4316, Timothy J. McGuire

91

## Example 6

- $S \rightarrow ( A ) \mid \varepsilon$
- $A \rightarrow T E$
- $E \rightarrow \& T E \mid \varepsilon$
- $T \rightarrow ( A ) \mid a \mid b \mid c$
- $\text{FOLLOW}(S) = \{\$ \}$
- $\text{FOLLOW}(A) = \{ ) \}$
- $\text{FOLLOW}(E) =$   
 $\text{FOLLOW}(A) = \{ ) \}$
- $\text{FOLLOW}(T) =$   
 $\text{FIRST}(E) \cup \text{FOLLOW}(A) \cup$   
 $\text{FOLLOW}(E) = \{\&, )\}$
- $\text{FIRST}(T) = \{(, a, b, c\}$
- $\text{FIRST}(E) = \{\&, \varepsilon\}$
- $\text{FIRST}(A) = \{(, a, b, c\}$
- $\text{FIRST}(S) = \{(, \varepsilon\}$

COSC 4316, Timothy J. McGuire

92

## Example 7

- $E \rightarrow T E'$
- $E' \rightarrow + T E' \mid \varepsilon$
- $T \rightarrow F T'$
- $T' \rightarrow * F T' \mid \varepsilon$
- $F \rightarrow ( E ) \mid \text{id}$
- $\text{FOLLOW}(E) =$
- $\text{FOLLOW}(E') =$
- $\text{FOLLOW}(T) =$
- $\text{FOLLOW}(T') =$
- $\text{FOLLOW}(F) =$
- $\text{FIRST}(F) = \text{FIRST}(T) = \text{FIRST}(E) = \{ (, \text{id} \}$
- $\text{FIRST}(T') = \{ *, \varepsilon \}$
- $\text{FIRST}(E') = \{ +, \varepsilon \}$

COSC 4316, Timothy J. McGuire

93

## Example 7

- $E \rightarrow T E'$
- $E' \rightarrow + T E' \mid \varepsilon$
- $T \rightarrow F T'$
- $T' \rightarrow * F T' \mid \varepsilon$
- $F \rightarrow ( E ) \mid \text{id}$
- $\text{FOLLOW}(E) = \{ \$, ) \}$
- $\text{FOLLOW}(E') = \text{FOLLOW}(E) = \{ \$, ) \}$
- $\text{FOLLOW}(T) = \text{FIRST}(E') \cup \text{FOLLOW}(E) \cup \text{FOLLOW}(E') = \{ +, \$, ) \}$
- $\text{FOLLOW}(T') = \text{FOLLOW}(T) = \{ +, \$, ) \}$
- $\text{FOLLOW}(F) = \text{FIRST}(T') \cup \text{FOLLOW}(T) \cup \text{FOLLOW}(T') = \{ *, +, \$, ) \}$
- $\text{FIRST}(F) = \text{FIRST}(T) = \text{FIRST}(E) = \{ (, \text{id} \}$
- $\text{FIRST}(T') = \{ *, \varepsilon \}$
- $\text{FIRST}(E') = \{ +, \varepsilon \}$

COSC 4316, Timothy J. McGuire

94

## Using FIRST and FOLLOW to Write a Recursive Descent Parser

$expr \rightarrow term\ rest$   
 $rest \rightarrow +\ term\ rest$   
          |  $- term\ rest$   
          |  $\epsilon$   
 $term \rightarrow id$

$FIRST(+\ term\ rest) = \{ + \}$   
 $FIRST(-\ term\ rest) = \{ - \}$   
 $FOLLOW(rest) = \{ \$ \}$

```
procedure rest();  
begin  
  if lookahead in FIRST(+ term rest) then  
    match('+'); term(); rest()  
  else if lookahead in FIRST(- term rest) then  
    match('-'); term(); rest()  
  else if lookahead in FOLLOW(rest) then  
    return  
  else error()  
end;
```