



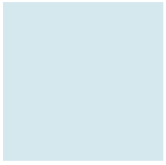
# COMPILER CONSTRUCTION

## Grammars and Parsing

Chia-Heng Tu

Dept. of Computer Science and Information  
Engineering

National Cheng Kung University  
Spring 2017



# Chapter 4

## Grammars and Parsing



# Why Grammar?

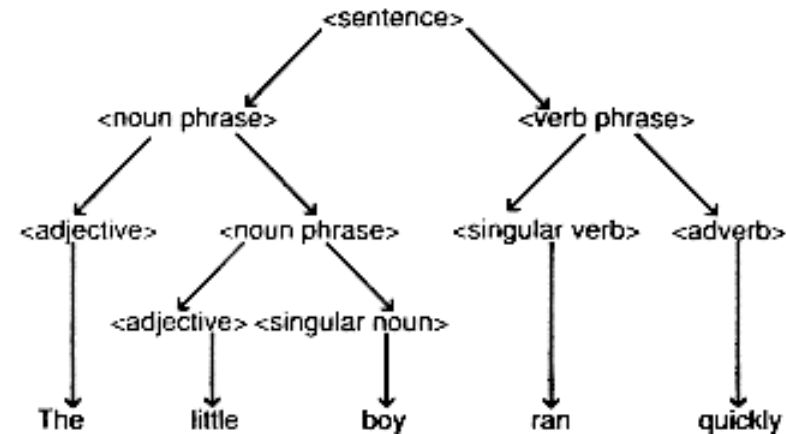
- Grammar rules of **natural languages**, such as English or Chinese
  - **Define proper sentence structure**, e.g., defining phrases in terms of subjects, verbs, and objects; and phrases and conjunctions
  - Served as a tool for **diagnosing malformed sentences** (validity check)
- It is possible to construct sentences in a natural language
  - that are grammatically correct
  - but still make no sense



## Why Grammar? (Cont'd)

- A structured sentence can be diagrammed
  - to show how its components conform to a language's grammar
  - Grammars can also explain what is absent or superfluous in a malformed sentence
- **Ambiguity** of a sentence
  - Can often be explained by providing **multiple diagrams** for the same sentence

“The little boy ran quickly”  
can be diagrammed as:





# Grammars for Programming Languages

- Modern programming languages
  - contain a **grammar in their specification** as a guide to those who teach, study, or use the language
- A compiler front-end for the language
  - Scans for tokens in input stream based on the **regular sets**
  - Parses the structures formed by the tokens using the **grammar** that specifies a programming language's syntax

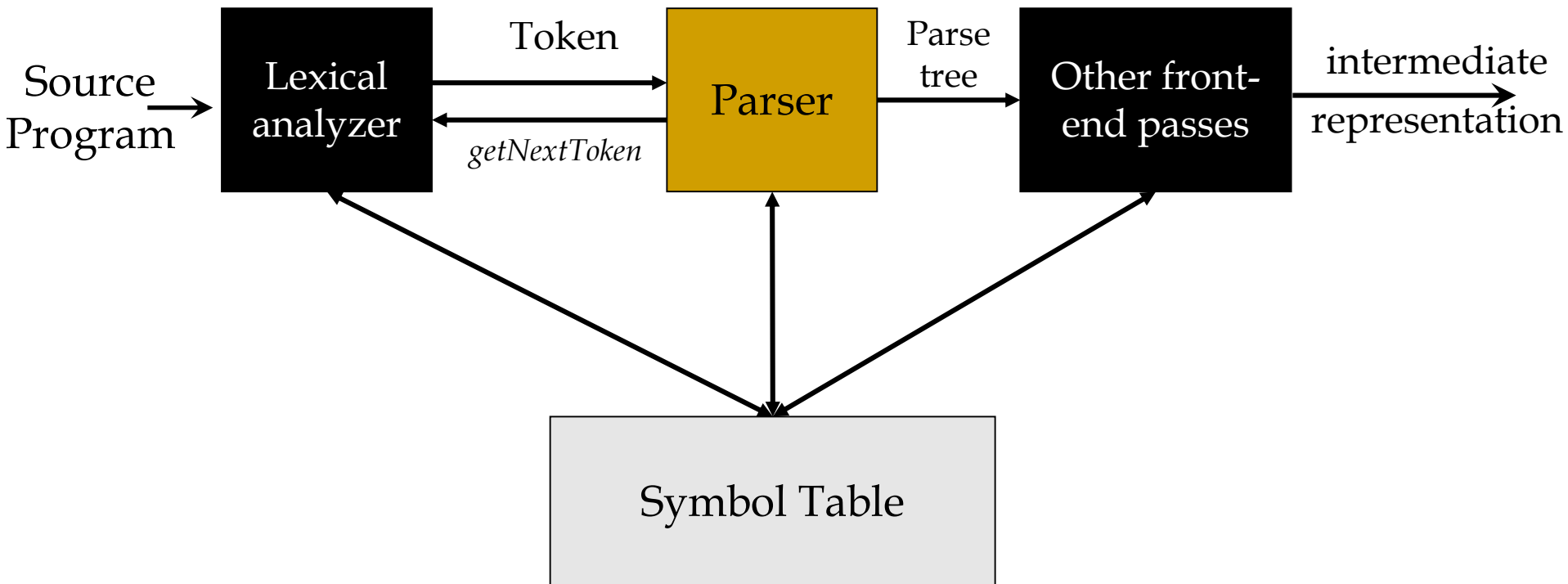


# This Chapter is for CFGs

- We discuss the basics of **context-free grammars** (CFGs) in Ch. 2
- In Ch. 4 we
  - **formalize the definition and notation for CFGs** and
  - present algorithms that analyze such grammars in preparation for the parsing techniques covered in Ch. 5 and 6



# The Role of the Parser





# Context-Free Grammar

- Context-free grammar is a 4-tuple  $G = \langle \Sigma, N, P, S \rangle$  where
  - $\Sigma$  is a finite set of **terminal alphabet**, which is the set of tokens produced by the scanner
  - $N$  is a finite set of **nonterminal alphabet**
  - $P$  is a finite set of **productions** of the form  $A \rightarrow \beta$   
where  $A \in N$  and  $\beta \in (N \cup \Sigma)^*$
  - $S \in N$  is a designated **start symbol**, which initiates all derivations





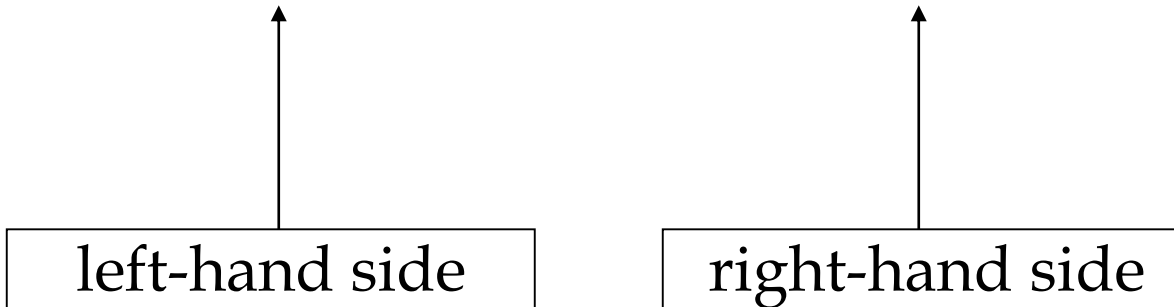
# Conventions

- Terminals
  - $a, b, c, \dots \in \Sigma$
  - More example: **0, 1, +, \*, id, if**
- Nonterminals
  - $A, B, C, \dots \in N$
  - More example: *expr, term, stmt*
- Grammar symbols
  - $X, Y, Z \in (N \cup \Sigma)$
- Strings of grammar symbols (sentential form)
  - $\alpha, \beta, \gamma \in (N \cup \Sigma)^*$
- The head of the first production is often the start symbol



# Production Rules

- A **grammar** consists of a set of **production rules** and **a start symbol** (left symbol of first rule)
- A **production rule** consists of two parts:  
a left-hand side and a right-hand side
  - ex:  $\text{expression} \rightarrow \text{expression} + \text{term}$





## Production Rules (Cont.)

- The **left-hand side** (LHS) is the **name** of the syntactic construct
- The **right-hand side** (RHS) shows a **possible form** of the syntactic construct
- E.g., there are **two possible forms** (rules) derived by **the name** “expression”:
  - expression  $\rightarrow$  expression ‘+’ term (rule 1)
  - expression  $\rightarrow$  expression ‘-’ term (rule 2)



# Production Rules (Cont.)

- LHS must be a single **non-terminal** symbol (or *non-terminal*)
- RHS of a production rule can contain zero or more **terminals** and **non-terminals**
- A **terminal symbol** (or *terminal*) is a grammar symbol
  - that cannot be rewritten
  - Is also an end point of the production process, also called **token**
  - Use lower-case letters such as a and b
- A **non-terminal symbol** (or *non-terminal*) is able to be rewritten
  - Use upper-case letters such as A, B, and S
- Non-terminal and terminal together are called **grammar symbols (or vocabulary)**



# Different Forms of Production Rules

- There is often more than one way to rewrite a given nonterminal
- For example, when multiple productions share the same LHS symbol, it could be presented in one of the two forms

$$\begin{array}{l} A \rightarrow \alpha \\ A \rightarrow \beta \\ \cdot \quad \cdot \quad \cdot \\ A \rightarrow \zeta \text{ (Zeta)} \end{array}$$

$$\begin{array}{l} A \rightarrow \alpha \\ | \beta \\ \cdot \quad \cdot \quad \cdot \\ | \zeta \end{array}$$



# Derivations of Production Rules

- Derivation
  - The rewriting step replaces a **nonterminal** by the RHS of one of **its** production rules
- Example:
  - If  $A \rightarrow \zeta$  is a production, then  $\alpha A \beta \Rightarrow \alpha \zeta \beta$  denotes one step of a derivation using this production
  - The one-step derivation can be denoted as  $\alpha A \beta \Rightarrow \alpha \zeta \beta$
  - The sequence of replacement a derivation of  $\zeta$  from  $A$



# Notation for Derivations

- $\Rightarrow$  derives in one step
- $\Rightarrow^+$  derives in one or more steps
- $\Rightarrow^*$  derives in zero or more steps
- $\Rightarrow_{lm}$  refers to **leftmost derivation**, which expands nonterminals left to right
- $\Rightarrow_{rm}$  refers to **right most derivation**, which expands nonterminals right to left
- Example:
  - $\alpha A \beta \Rightarrow \alpha \zeta \beta$
  - is leftmost derivation, if  $\alpha$  does not contain a nonterminal
  - is rightmost derivation, if  $\beta$  does not contain a nonterminal



# More about the Grammar

- We say that  $\alpha$  is a **sentential form** of  $G$ 
  - if  $S \Rightarrow^* \alpha$ , where  $S$  is the start symbol of a grammar  $G$
  - Note that a sentential form may contain both terminals and nonterminals, and may be empty
- A **sentence** of  $G$  is a sentential form with no nonterminals
- The language generated by a grammar is its set of sentences
  - Hence, a string of terminals  $w$  is in  $L(G)$ , the language generated by  $G$ , if and only if  $w$  is a sentence of  $G$  (or  $S \Rightarrow^* w$ )





# Example of Rule Derivations

- Given the production rules:
  - $\text{expr} \rightarrow '(\text{expr op expr})'$
  - $\text{expr} \rightarrow '1'$
  - $\text{op} \rightarrow '+'$
  - $\text{op} \rightarrow '*'$



# Example of Rule Derivations (Cont'd)

- Derivation of the string  $(1*(1+1))$

- $\text{expr}$
- $(\text{expr op expr})$
- $((1 \text{ op expr}))$
- $((1 * \text{expr}))$
- $((1 * (\text{expr op expr})))$
- $((1 * ((1 \text{ op expr})))$
- $((1 * ((1 + \text{expr})))$
- $((1 * ((1 + 1)))$

$\text{expr} \rightarrow '(\text{expr op expr})'$

$\text{expr} \rightarrow '1'$

$\text{op} \rightarrow '+'$

$\text{op} \rightarrow '*'$

- Each of the strings is a **sentential form**
- The *strings* refer to  $\text{expr}$ ,  $(\text{expr op expr})$ ,  $((1 \text{ op expr}))$ ,  $((1 * \text{expr}))$ , etc.

- It forms a **leftmost derivation**, in which the leftmost nonterminal is always rewritten in each sentential form



# A CFG Example

- $G = \langle \Sigma, N, P, S \rangle$ 
  - $\Sigma = \{ +, *, (, ), -, \text{id} \}$
  - $N = \{ E \}$
  - $P = \begin{aligned} E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow ( E ) \\ E &\rightarrow - E \\ E &\rightarrow \text{id} \end{aligned}$
  - $S = E$

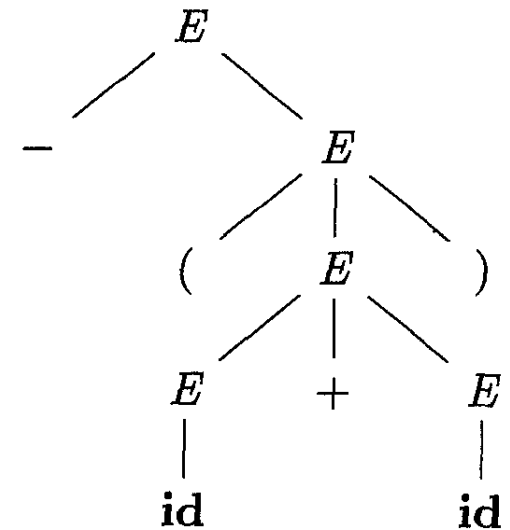
- Is the string  $-(\text{id} + \text{id})$  a sentence of  $G$ ?

Yes, there is a derivation of the given string

$\Rightarrow E \Rightarrow_{\text{lm}} - E \Rightarrow_{\text{lm}} -(E) \Rightarrow_{\text{lm}} -(E+E) \Rightarrow_{\text{lm}} -(\text{id} + E) \Rightarrow_{\text{lm}} -(\text{id} + \text{id})$

- The strings  $E, - E, - (E), \dots, - (\text{id} + \text{id})$  are all sentential forms of this grammar
  - We write  $E \Rightarrow^* - (\text{id} + \text{id})$  to indicate that  $-(\text{id} + \text{id})$  can be derived from  $E$

Parse tree for  $-(\text{id} + \text{id})$





# You would ...

- Refer to Section 4.1 for more information
  - Examples of leftmost and right most derivations in Section 4.1.1 and 4.1.2
  - Parse Trees in Section 4.1.3



# Properties of CFGs

- Some grammars have one or more of the following problems that preclude their use
  - **(Reduced Grammar)** The grammar may include useless symbols
  - **(Ambiguity)** The grammar may allow multiple, distinct derivations (parse trees) for some input string
  - **(Faulty Language Definition)** The grammar may include strings that do not belong in the language, or
    - the grammar may exclude strings that are in the language



## We are ...

- going to know more about the ambiguity next
- You would ...
  - read Section 4.2.1 and 4.2.3 by yourself



# Ambiguity

- A grammar that produces more than one parse tree for some sentence is said to be *ambiguous*
- Hence, an ambiguous grammar is one that produces
  - more than one leftmost derivation or
  - more than one rightmost derivation for the same sentence



# Example of Ambiguous Grammar

- Given the sentence: **id + id \* id**
- Production rules:

$$E \rightarrow E + E \mid E * E \mid ( E ) \mid \text{id}$$

- Two possible derivations:

$$\begin{aligned} E &\rightarrow E + E \\ &\rightarrow \text{id} + E \\ &\rightarrow \text{id} + E * E \\ &\rightarrow \text{id} + \text{id} * E \\ &\rightarrow \text{id} + \text{id} * \text{id} \end{aligned}$$

$$\begin{aligned} E &\rightarrow E * E \\ &\rightarrow E + E * E \\ &\rightarrow \text{id} + E * E \\ &\rightarrow \text{id} + \text{id} * E \\ &\rightarrow \text{id} + \text{id} * \text{id} \end{aligned}$$





## Another Example

- Given the sentence: **9-5+2**
- Production rules:

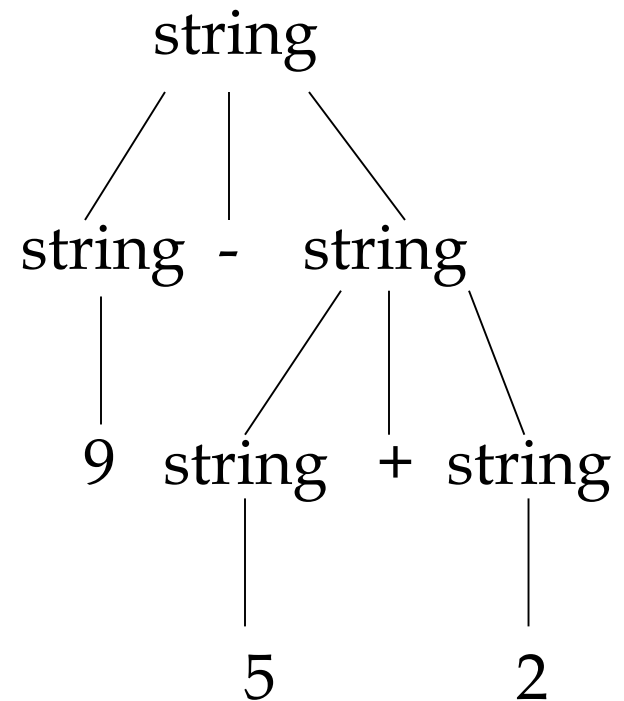
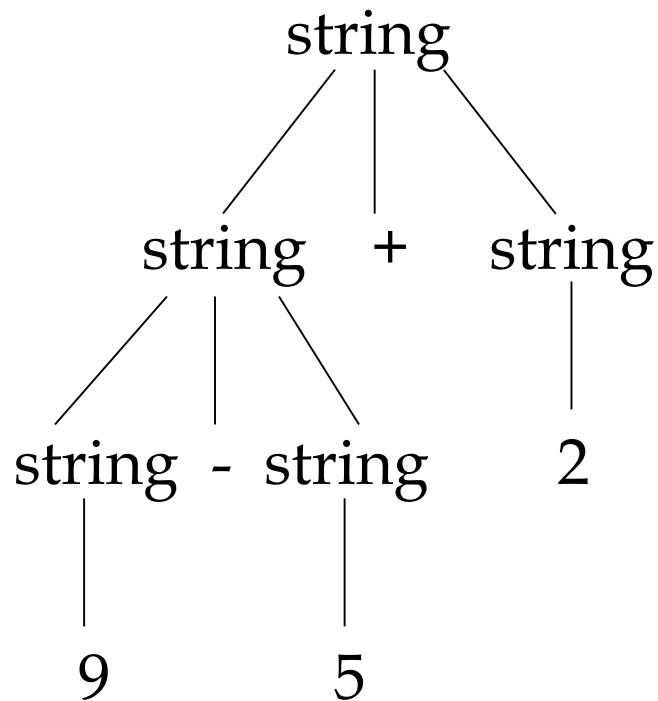
string  $\rightarrow$  string + string

| string - string

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9



# Another Example (Con'td)





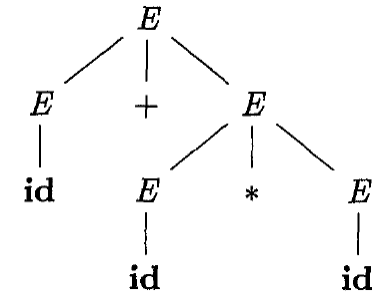
# To Deal with Ambiguity

- Disambiguating rules are used to throw away undesirable parse trees
- Two options to deal with ambiguity:
  - Enforce precedence and associativity of the existing rules
  - Rewrite the grammar (rules)

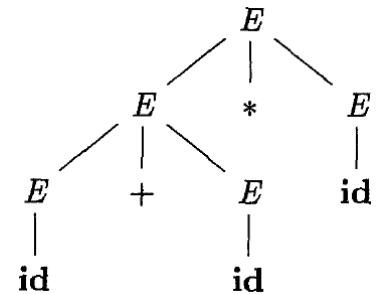


# Precedence and Associativity to Resolve Conflicts

- The grammar  $G$  with the rules  
 $E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$
- $G$  is **ambiguous** because
  - it does not specify the associativity or precedence of the operators  $+$  and  $*$



- The following grammar  $G'$  with the rules  
 $E \rightarrow E + T \mid T$   
 $T \rightarrow T * F \mid F$   
 $F \rightarrow (E) \mid \text{id}$



- $G'$  is **unambiguous** grammar that generates the same language
  - but gives  $+$  lower precedence than  $*$  (as the top figure), and makes both operators left associative
- There are parsers for both handling unambiguous and ambiguous grammars respectively



# Example

- Given the sentence: **id+id\*id**
- The following grammar  $G'$  with the rules

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

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

$$F \rightarrow ( E ) \mid \mathbf{id}$$

- The corresponding derivations:  
$$\begin{aligned} E &\rightarrow E+T \\ &\rightarrow T+T \\ &\rightarrow F+T \\ &\rightarrow \mathbf{id}+T*F \\ &\rightarrow \mathbf{id}+F*F \\ &\rightarrow \mathbf{id}+\mathbf{id}*F \\ &\rightarrow \mathbf{id}+\mathbf{id}*\mathbf{id} \end{aligned}$$



# Extended Grammars

- Backus-Naur Form (BNF)
  - is a **formal grammar** for expressing context-free grammars
- The single grammar rule format:
  - Non-terminal  $\rightarrow$  zero or more grammar symbols
- Actually, BNF **extends** the grammar notation defined above, with syntax for defining **optional** and **repeated symbols**



# Optional Symbols

- **Optional** symbols are enclosed in square brackets
- In the production rule
$$A \rightarrow \alpha [ X1 \dots Xn ] \beta$$
  - the symbols  $X1 \dots Xn$  are **entirely present or absent** between the symbols of  $\alpha$  and  $\beta$
  - Refer to Fig. 4.4 for the algorithm to transform a BNF grammar into standard form



# Repeated Symbols

- **Repeated** symbols are enclosed in braces. In the production
- In the production rule
$$B \rightarrow \gamma \{ X_1 \dots X_m \} \delta$$
  - the entire sequence of symbols  $X_1 \dots X_m$  can be **repeated zero or more times**
  - Refer to Fig. 4.4 for the algorithm to transform a BNF grammar into standard form, which is accepted by parsers





# Example of BNF

- The extensions are useful in representing many programming language constructs
- For example, in Java
  - Declarations can optionally include modifiers, such as **final**, **static**, and **const**, and
  - each declaration can include a *list* of identifiers
  - A production specifying a Java-like declaration could be as follows:  
Declaration → [ **final** ] [ **static** ] [ **const** ] Type identifier { , **identifier** }
  - Possible declarations:  
int a  
int a,b,c  
static int a
- This declaration insists that the modifiers be ordered as shown



# Parsers and Recognizers

- Compilers are expected to
  - verify the **syntactic validity** of their **inputs** with respect to a **grammar** that
  - defines the programming language's syntax
- A **recognizer** is an algorithm determines if  $x \in L(G)$ , given a grammar  $G$  and an input string  $x$
- A **parser** is a module that determines the string's validity and its structure (or parse tree)
  - The process of finding the structure in the flat stream of tokens is called **parsing**



# Two Parsing Approaches

- Top-down parsers
  - Left-scan, Leftmost derivation
  - Best-known parser in this category, called **LL** parsers
- Bottom-up parsers
  - Left-scan, Rightmost derivation in reverse
  - Best-known parser in this category, called **LR** parsers



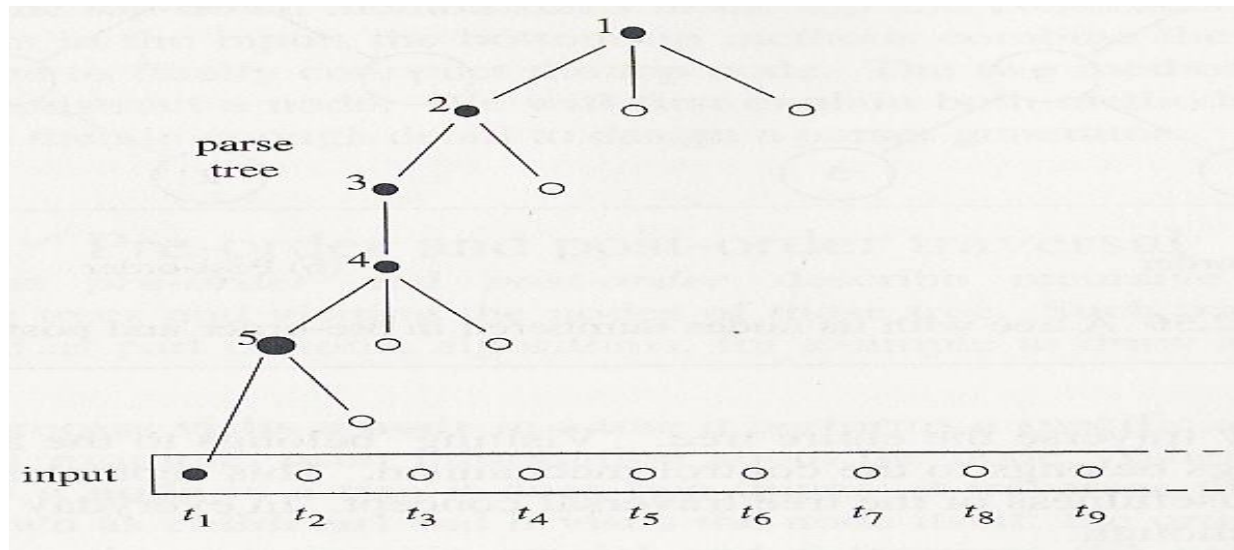
# Top-Down Parsers

- A parser is considered **top-down**
  - if it generates a parse tree by starting at the root of the tree (the start symbol),
  - expanding the tree by applying productions in a **depth-first** manner
  - It corresponds to a **preorder traversal** of the parse tree
- Top-down parsing techniques are **predictive**
  - because they always predict the production that is to be matched before matching actually begins



# Illustration of Top-Down Parsing

- A top-down parser begins by constructing the top node of the parse tree, which is the start symbol





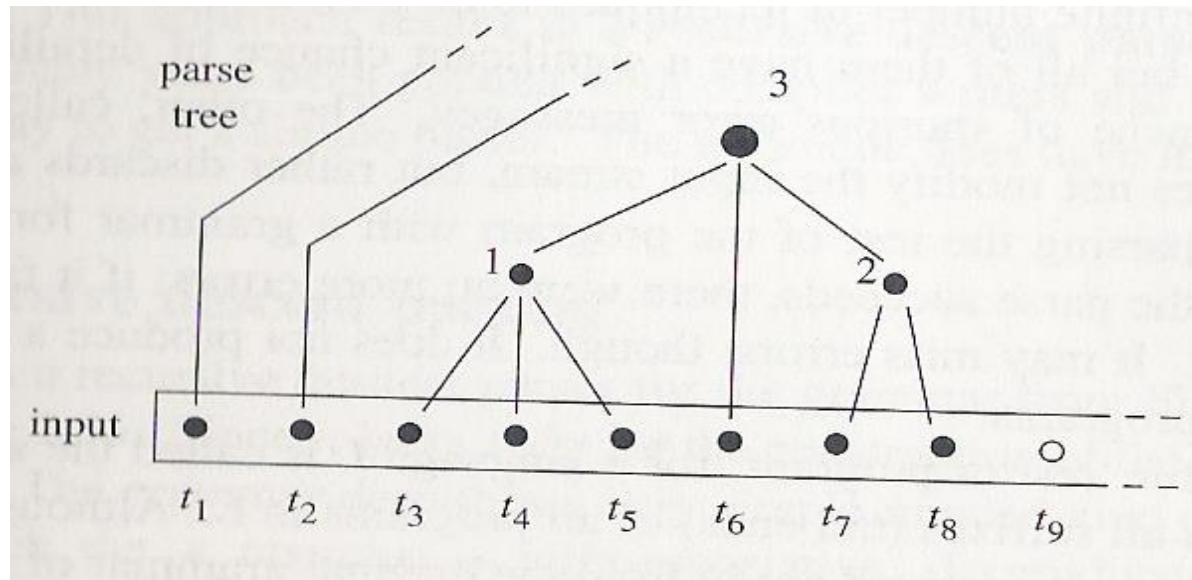
# Bottom-Up Parsers

- The **bottom-up** parsers generate a parse tree by
  - starting at the tree's leaves and working toward its root
  - A node is inserted in the tree only after its children have been inserted
- A bottom-up parse corresponds to a **postorder traversal** of the parse tree



# Illustration of Bottom-Up Parsing

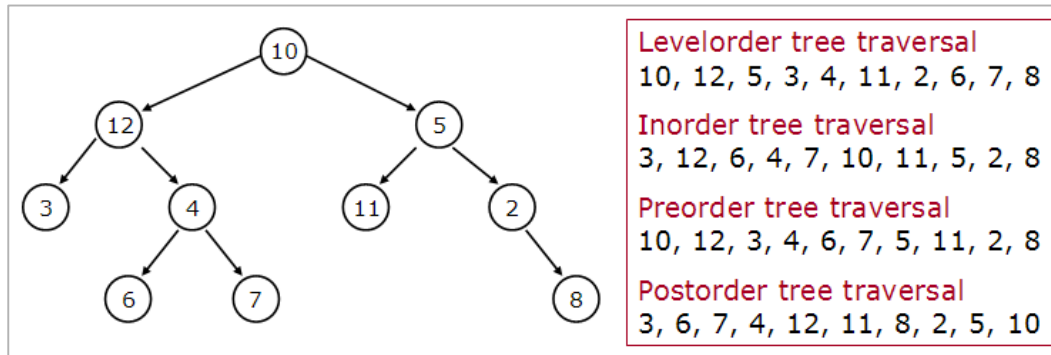
- The bottom-up parsing method constructs the nodes in the parse tree in post-order





# To Refresh Your (and My) Memory

- Several ways for tree traversals



← Top-down parsers

← Bottom-up parsers

- When traversing a node N in pre-order
  - the process first visits the node N and then traverses N's subtrees in left-to-right order
- When traversing a node N in post-order
  - the process first traverses N's subtrees in left-to-right order and then visits the node N

Quick note

Preorder: Root, Left, Right

Postorder: Left, Right, Root





# Example of the Parsers

- An example grammar generates the skeletal block structure of a hypothetical programming language

1 **Program**  $\rightarrow$  **begin Stmts end \$**

2 **Stmts**  $\rightarrow$  **Stmt ; Stmts**

3  $\quad \quad \quad | \Lambda$

4 **Stmt**  $\rightarrow$  **simplestmt**

5  $\quad \quad \quad | \text{begin Stmts end}$

- The Fig. 4.5 and 4.6 illustrate a top-down and bottom-up parse of the given string:

**begin simplestmt ; simplestmt ; end \$**

# Top-Down Parsing Example

### Legend:

1. Each box shows one step of the parse, with the particular rule denoted by bold lines between a parent (the rule's LHS) and its children (the rule's RHS)
2. Solid, non-bold lines indicate rules that have already been applied
3. Dashed lines indicate rules that have not yet been applied

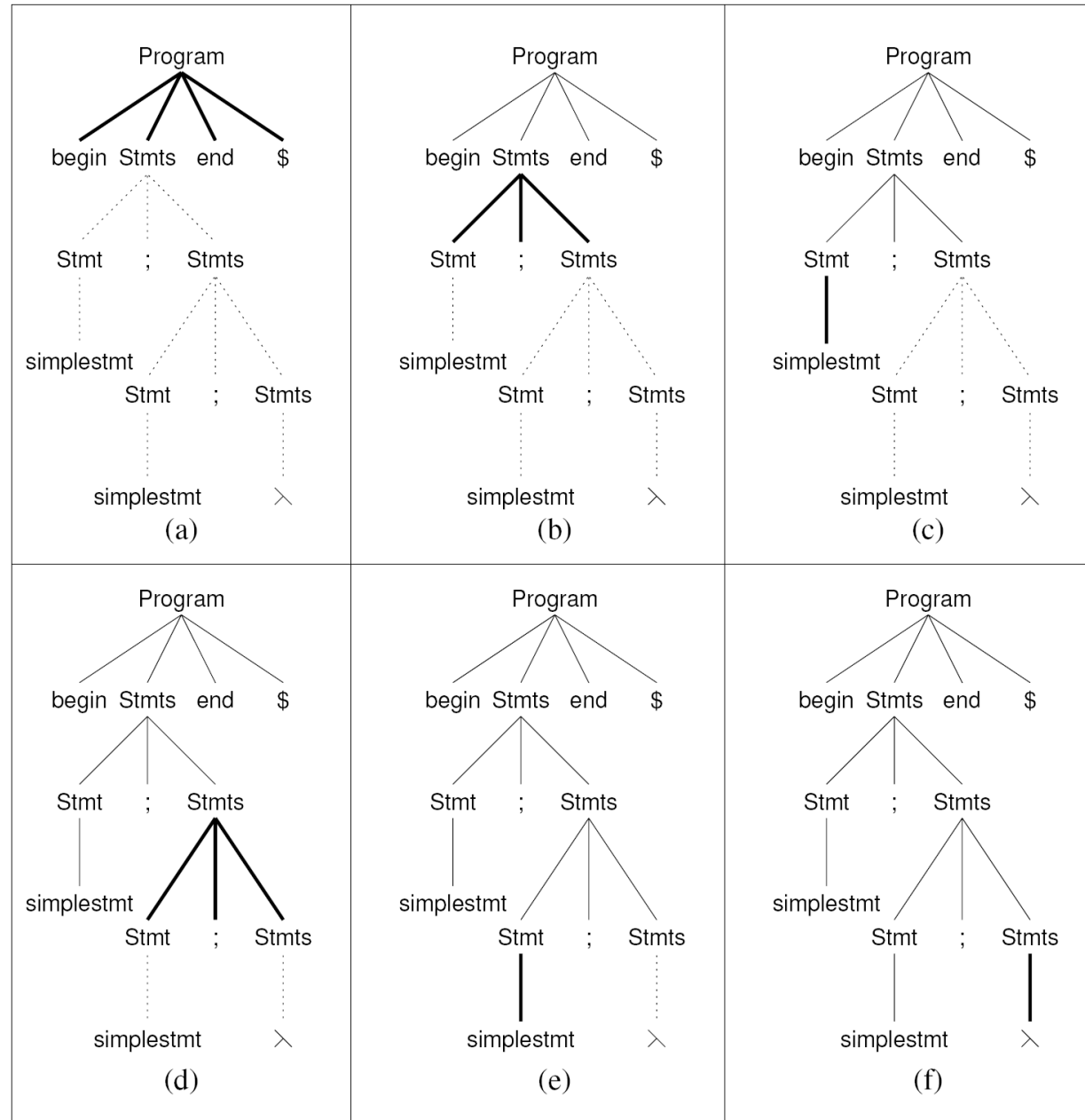


Figure 4.5: Parse of “begin simplestmt ; simplestmt ; end \$” using the top-down technique. Legend explained on page 126.

# Top-Down Parsing Example'

- Fig. 4.5(a) shows the rule **Program**→**begin Stmts end \$** applied as the first step of a top-down parse

## Legend:

- Each box shows one step of the parse, with the particular rule denoted by bold lines between a parent (the rule's LHS) and its children (the rule's RHS)
- Solid, non-bold lines indicate rules that have already been applied
- Dashed lines indicate rules that have not yet been applied

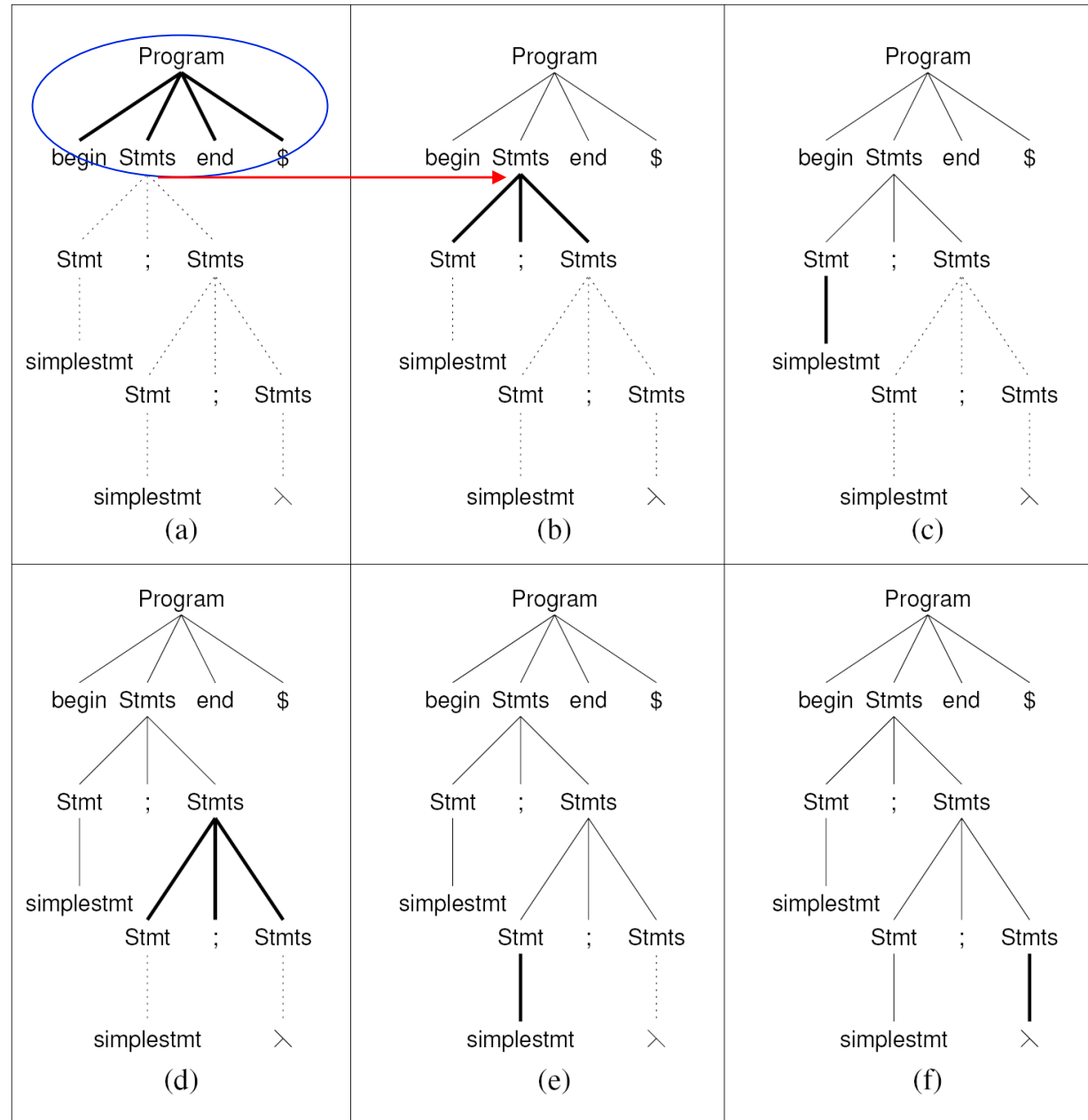


Figure 4.5: Parse of "begin simplestmt ; simplestmt ; end \$" using the top-down technique. Legend explained on page 126.



# Top-Down Parsing Example'

- 1 Program  $\rightarrow$  begin Stmts end \$
- 2 Stmts  $\rightarrow$  Stmt ; Stmts
- 3 |  $\lambda$
- 4 Stmt  $\rightarrow$  simplestmt
- 5 | begin Stmts end

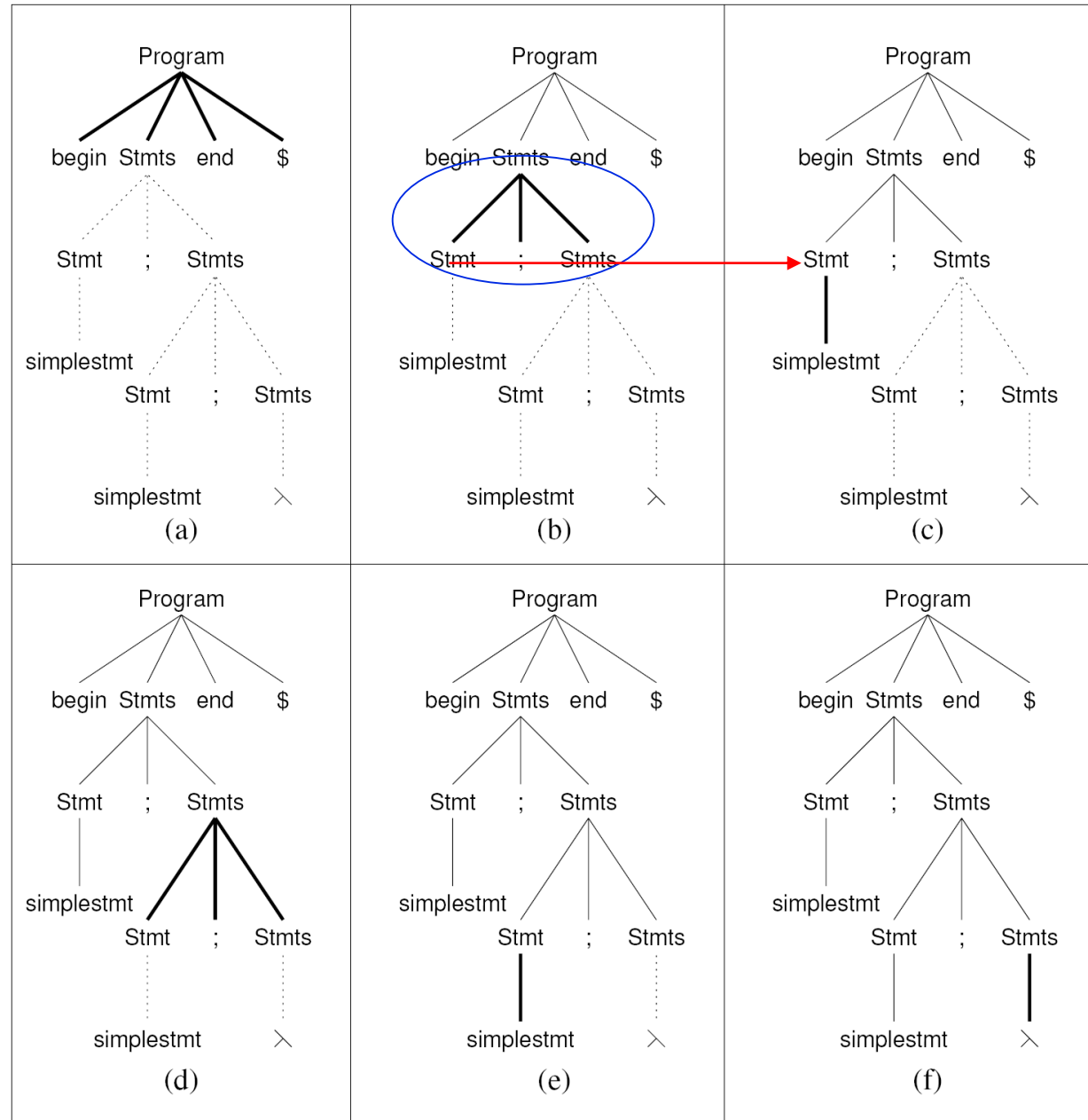
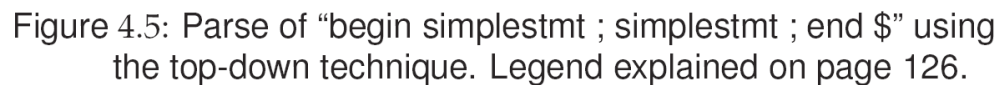


Figure 4.5: Parse of “begin simplestmt ; simplestmt ; end \$” using the top-down technique. Legend explained on page 126.

- 1 Program  $\rightarrow$  begin Stmts end \$
- 2 Stmts  $\rightarrow$  Stmt ; Stmts
- 3       |  $\lambda$
- 4 Stmt  $\rightarrow$  simplestmt
- 5       | begin Stmts end



# Top-Down Parsing Example'

- 1 Program  $\rightarrow$  begin Stmts end \$
- 2 Stmts  $\rightarrow$  Stmt ; Stmts
- 3 |  $\lambda$
- 4 Stmt  $\rightarrow$  simplestmt
- 5 | begin Stmts end

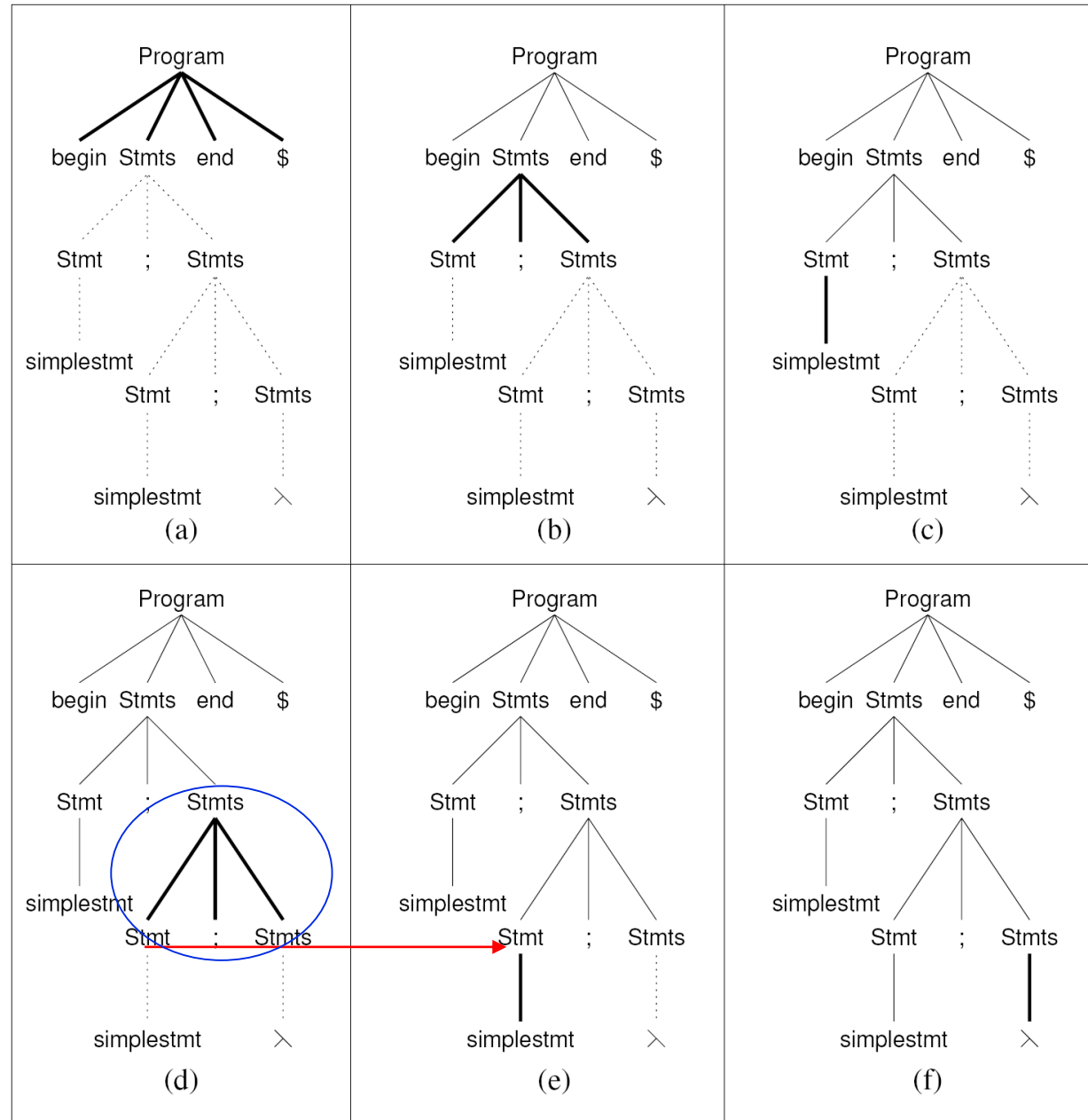


Figure 4.5: Parse of "begin simplestmt ; simplestmt ; end \$" using the top-down technique. Legend explained on page 126.

# Top-Down Parsing Example'

- 1 Program → begin Stmts end \$
- 2 Stmts → Stmt ; Stmts
- 3 | λ
- 4 Stmt → simplestmt
- 5 | begin Stmts end

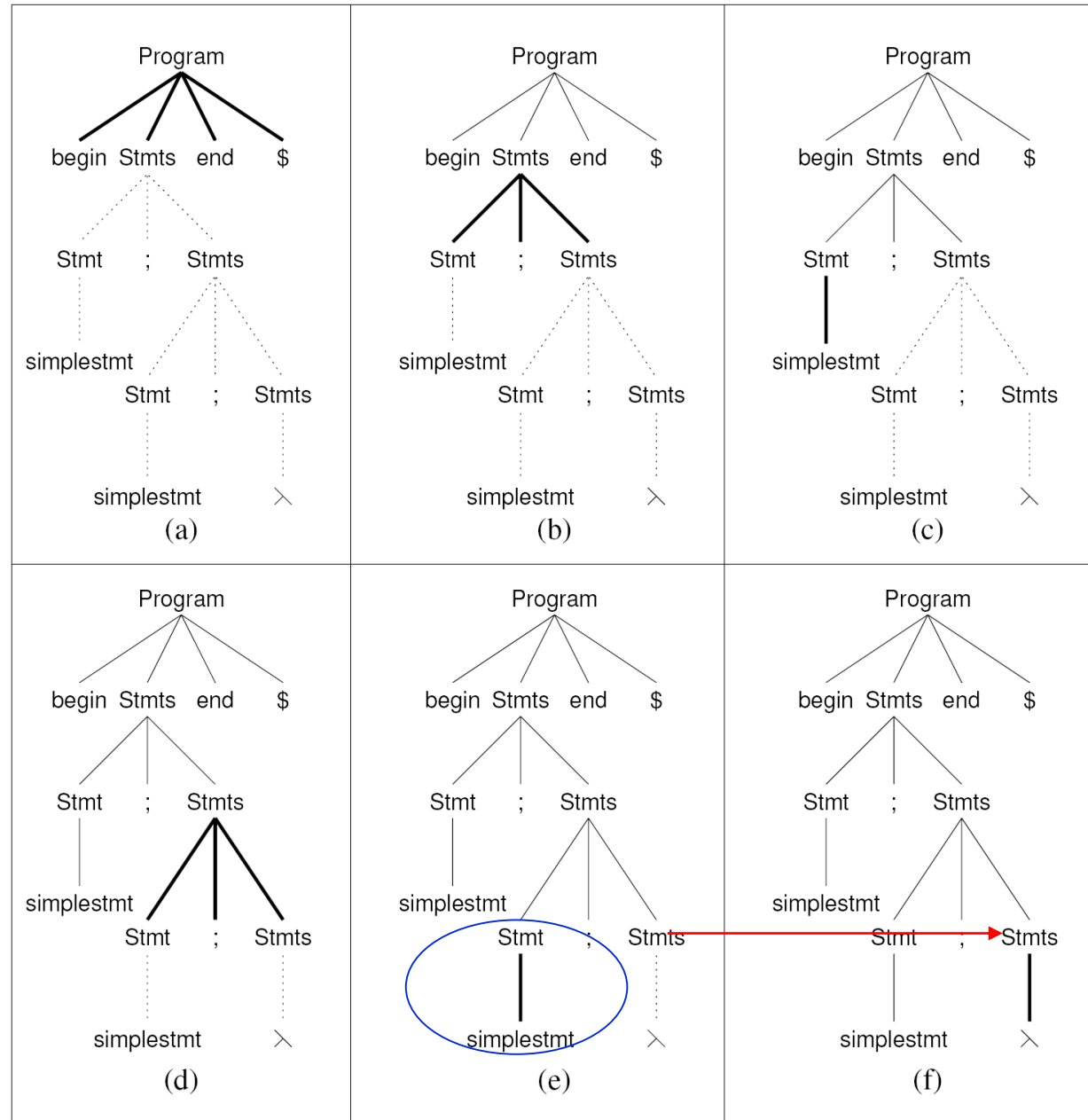


Figure 4.5: Parse of "begin simplestmt ; simplestmt ; end \$" using the top-down technique. Legend explained on page 126.

# Top-Down Parsing Example'

- 1 Program  $\rightarrow$  begin Stmts end \$
- 2 **Stmts**  $\rightarrow$  Stmt ; Stmts
- 3     |  $\lambda$
- 4 Stmt  $\rightarrow$  simplestmt
- 5     | begin Stmts end

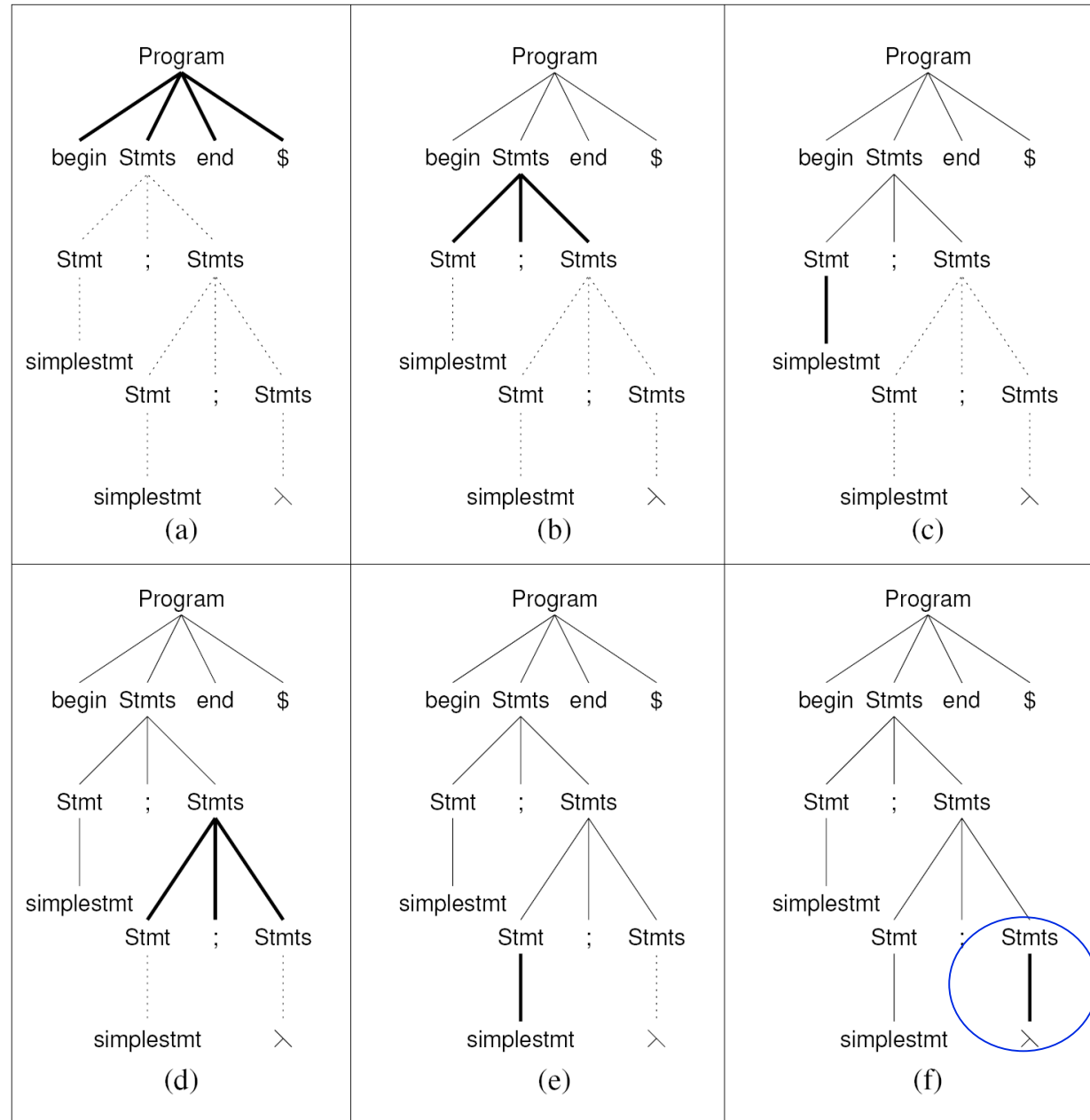


Figure 4.5: Parse of “begin simplestmt ; simplestmt ; end \$” using the top-down technique. Legend explained on page 126.





## You would ...

- Practice the bottom-up parsing in Fig. 4.6
- If you get stuck
  - I assume you forget to read Section 4.1.1 and 4.1.2 ☺
  - Please carefully read the sections first and try it again



# LL and LR Parsers

- **LL** and **LR** reflect
  - **how the input is processed** and **which kind of parse is produced**
  - The first character (L) states that the token sequence is processed from left to right
  - The second letter (L or R) indicates whether a leftmost or rightmost parse is produced
- The parsing technique can be further characterized by the **number of lookahead symbols**
  - i.e., symbols beyond the current token that the parser may consult to make
  - parsing choices
  - LL(1) and LR(1) parsers are the most common, requiring only **one symbol of lookahead**



# Summary

- There is a many-to-one relationship between derivations and parse trees
  - A parse tree ignores variations in the order in which symbols in sentential forms are replaced
- While the parsing sequences of top-down and bottom-up parsing are different, two parsing techniques construct the same parse tree, as shown in Fig. 4.5 and 4.6



# Why FIRST and FOLLOW?

- FIRST and FOLLOW
  - are the two important functions for the construction of top-down and bottom-up parsers
- FIRST and FOLLOW allow the parsers to choose which production to apply
- During panic-mode error recovery, sets of tokens produced by FOLLOW can be used as synchronizing tokens



# FIRST

- FIRST( $\alpha$ )

- refers to the set of **terminals** that begin strings derived from  $\alpha$
- where  $\alpha$  is any string of grammar symbols

Example:

- If  $A \Rightarrow^* cy$ , then  $c$  is in FIRST( $A$ )
- If  $\alpha \Rightarrow^* \lambda$ , then  $\lambda$  is also in FIRST( $\alpha$ )

- Given  $A \Rightarrow a \mid b$

- FIRST( $A$ ) is the union of FIRST( $a$ ) and FIRST( $b$ ),
- where FIRST( $a$ ) and FIRST( $b$ ) are disjoint sets



## FIRST (Cont'd)

- \*Compute FIRST( $X$ )
  - $X$  is grammar symbol
  - We apply the following rules until no more terminals or  $\lambda$  can be added to it
- 1. If  $X$  is a terminal (i.e.,  $X \in \Sigma$ ), then  $\text{FIRST}(X) = \{X\}$
- 2. If  $X \Rightarrow \lambda$  is a production (i.e.,  $X \in N$  and  $X \Rightarrow \lambda$ ), then add  $\lambda$  to  $\text{FIRST}(X)$
- 3. If  $X$  is a non-terminal (i.e.,  $X \in N$ ) and  $X \Rightarrow Y_1 Y_2 \dots Y_k$  is a production for some  $k \geq 1$ , then place  $\alpha$  in  $\text{FIRST}(X)$  if for some  $i$ ,  $\alpha$  is in  $\text{FIRST}(Y_i)$ , and  $\lambda$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$

If  $\lambda$  is in  $\text{FIRST}(Y_j)$  for all  $j=1, 2, \dots, k$ , then add  $\lambda$  to  $\text{FIRST}(X)$



# FIRST( $\alpha$ ) Function

```

function FIRST( $\alpha$ ) returns Set
    foreach  $A \in \text{NonTerminals}()$  do VisitedFirst( $A$ )  $\leftarrow$  false
    ans  $\leftarrow$  INTERNALFIRST( $\alpha$ )
    return (ans)
end
function INTERNALFIRST( $X\beta$ ) returns Set
    if  $X\beta = \perp$  → Rule 2
    then return ( $\emptyset$ )
    if  $X \in \Sigma$  → Rule 1
    then return ( $\{X\}$ )
    /*  $X$  is a nonterminal. → Rule 3
    ans  $\leftarrow \emptyset$ 
    if not VisitedFirst( $X$ )
    then
        VisitedFirst( $X$ )  $\leftarrow$  true
        foreach rhs  $\in$  ProductionsFor( $X$ ) do
            ans  $\leftarrow$  ans  $\cup$  INTERNALFIRST(rhs)
        if SymbolDerivesEmpty( $X$ )
        then ans  $\leftarrow$  ans  $\cup$  INTERNALFIRST( $\beta$ )
        return (ans)
    end

```

⑨  
 ⑩  
 ⑪  
 ★/ ⑫  
 ⑬  
 ⑭  
 ⑮  
 ⑯

- **SymbolDerivesEmpty( $A$ )** indicates whether or not the nonterminal  $A$  can derive  $\lambda$
- **VisitedFirst( $X$ )** is to indicate that the productions of  $X$  is already participate in the computation of FIRST( $\alpha$ )

Figure 4.8: Algorithm for computing First( $\alpha$ ).



# Handling Endless Recursion

$A \Rightarrow B$   
 $\cdot \quad \cdot \quad \cdot$   
 $B \Rightarrow C$   
 $\cdot \quad \cdot \quad \cdot$   
 $C \Rightarrow A$

```

function FIRST( $\alpha$ ) returns Set
  foreach  $A \in \text{NonTerminals}()$  do  $\text{VisitedFirst}(A) \leftarrow \text{false}$ 
   $\text{ans} \leftarrow \text{INTERNALFIRST}(\alpha)$ 
  return ( $\text{ans}$ )
end
function INTERNALFIRST( $X\beta$ ) returns Set
  if  $X\beta = \perp$ 
  then return ( $\emptyset$ )
  if  $X \in \Sigma$ 
  then return ( $\{X\}$ )
  /*  $X$  is a nonterminal.
   $\text{ans} \leftarrow \emptyset$ 
  if not  $\text{VisitedFirst}(X)$ 
  then
     $\text{VisitedFirst}(X) \leftarrow \text{true}$ 
    foreach  $\text{rhs} \in \text{ProductionsFor}(X)$  do
       $\text{ans} \leftarrow \text{ans} \cup \text{INTERNALFIRST}(\text{rhs})$ 
  if  $\text{SymbolDerivesEmpty}(X)$ 
  then  $\text{ans} \leftarrow \text{ans} \cup \text{INTERNALFIRST}(\beta)$ 
  return ( $\text{ans}$ )
end
  
```

Figure 4.8: Algorithm for computing  $\text{FIRST}(\alpha)$ .

- Termination of  $\text{FIRST}(A)$  must be handled properly in grammars
  - where **the computation of  $\text{FIRST}(A)$  appears to depend on  $\text{FIRST}(A)$**
- VisitedFirst(X)** is to indicate that the productions of  $X$  is already participate in the computation of  $\text{FIRST}(\alpha)$





## FIRST(Cont'd)

- More about the Rule 3 in the previous page
- 3. If  $X$  is a non-terminal (i.e.,  $X \in N$ ) and  $X \Rightarrow Y_1 Y_2 \dots Y_k$  is a production for some  $k \geq 1$ , then place  $\alpha$  in  $\text{FIRST}(X)$  if for some  $i$ ,  $\alpha$  is in  $\text{FIRST}(Y_i)$ , and  $\lambda$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$
- Examples:
  - Everything in  $\text{FIRST}(Y_1)$  is surely in  $\text{FIRST}(X)$
  - If  $Y_1$  does not derive  $\lambda$ , then we add nothing more to  $\text{FIRST}(X)$
  - If  $Y_1 \Rightarrow^* \lambda$ , then we add  $\text{FIRST}(Y_2)$ , and so on



# Example

Given the grammar,

- $E \Rightarrow TE'$
- $E' \Rightarrow +TE' \mid \lambda$
- $T \Rightarrow FT'$
- $T' \Rightarrow *FT' \mid \lambda$
- $F \Rightarrow (E) \mid id$

where E is for expression, T is for term, and F is for factor.  
Please find the FIRST set of each symbol

- FIRST sets for E, T, F, E', and T':
  - $\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ (, id \}$
  - $\text{FIRST}(E') = \{ +, \lambda \}$
  - $\text{FIRST}(T') = \{ *, \lambda \}$



# Another Example for FIRST

Given the grammar,

- $input \Rightarrow expression$
- $expression \Rightarrow term \ rest\_expression$
- $term \Rightarrow ID \mid parenthesized\_expression$
- $parenthesized\_expression \Rightarrow '(' \ expression \ ')'$
- $rest\_expression \Rightarrow '+' \ expression \mid \lambda$

FIRST sets for input, expression, term, parenthesized\_expression, and rest\_expression:

- $FIRST(input) = FIRST(expression) = FIRST(term) = \{ ID, ( \}$
- $FIRST(parenthesized\_expression) = \{ ( \}$
- $FIRST(rest\_expression) = \{ +, \lambda \}$



## You Should ...

- Refer to Section 4.5.3 to do the exercise in Fig. 4.9 and 4.10

# Exercise

- Note the *ans* does not contain  $\lambda$ , which may cause trouble for you to derive the results

```

1 E      → Prefix ( E )
2       | v Tail
3 Prefix → f
4       | λ
5 Tail   → + E
6       | λ

```

Figure 4.1: A simple expression grammar.

```

function FIRST( $\alpha$ ) returns Set
  foreach  $A \in \text{NonTerminals}()$  do VisitedFirst( $A$ )  $\leftarrow$  false
   $ans \leftarrow \text{INTERNALFIRST}(\alpha)$ 
  return ( $ans$ )
end
function INTERNALFIRST( $X\beta$ ) returns Set
  if  $X\beta = \perp$ 
  then return ( $\emptyset$ )
  if  $X \in \Sigma$ 
  then return ( $\{X\}$ )
  /*  $X$  is a nonterminal.
   $ans \leftarrow \emptyset$ 
  if not VisitedFirst( $X$ )
  then
    VisitedFirst( $X$ )  $\leftarrow$  true
    foreach  $rhs \in \text{ProductionsFor}(X)$  do
       $ans \leftarrow ans \cup \text{INTERNALFIRST}(rhs)$ 
  if SymbolDerivesEmpty( $X$ )
  then  $ans \leftarrow ans \cup \text{INTERNALFIRST}(\beta)$ 
  return ( $ans$ )
end

```

Figure 4.8: Algorithm for computing First( $\alpha$ ).

Level	First $X$	$\beta$	<i>ans</i>	Marker	Done? ( $\star$ =Yes)	Comment
FIRST(Tail)						
0	Tail	$\perp$	{ }	(12)		
1	+	E	{+}	(11)	$\star$	Tail $\rightarrow$ +E
1	$\perp$	$\perp$	{ }	(10)	$\star$	Tail $\rightarrow$ $\lambda$
0			{+}	(14)		After all rules for Tail
1	$\perp$	$\perp$	{ }	(10)	$\star$	Since $\beta = \perp$
0			{+}	(15)	$\star$	Final answer

FIRST(Prefix)						
0	Prefix	$\perp$	{ }	(12)		
1	f	$\perp$	{f}	(11)	$\star$	Prefix $\rightarrow$ f
1	$\perp$	$\perp$	{ }	(10)	$\star$	Prefix $\rightarrow$ $\lambda$
0			{f}	(14)		After all rules for Prefix
1	$\perp$	$\perp$	{ }	(10)	$\star$	Since $\beta = \perp$
0			{f}	(15)	$\star$	Final answer

FIRST(E)						
0	E	$\perp$	{ }	(12)		
1	Prefix	( E )	{ }	(12)		E $\rightarrow$ Prefix ( E )
1			{f}	(16)		Computation shown above
2	(	E )	{ ( }	(11)	$\star$	Since Prefix $\Rightarrow^* \lambda$
1			{f, ( }	(15)	$\star$	Results due to E $\rightarrow$ Prefix ( E )
1	v	Tail	{v}	(11)	$\star$	E $\rightarrow$ v Tail
1	$\perp$	$\perp$	{ }	(10)		Since $\beta = \perp$
0			{f, (, v}	(15)	$\star$	Final answer

Figure 4.9: First sets for the nonterminals of Figure 4.1.



# FOLLOW

- FOLLOW( $\beta$ )

- refers to the set of terminals  $a$  that can appear immediately to the right of non-terminal  $\beta$  in some sentential form

Example:

- If  $S \Rightarrow \alpha A a \beta$ , then  $a$  is in the set of FOLLOW( $A$ )



# FOLLOW (Cont'd)

- \*FOLLOW(B)
  - where B is non-terminal ,
  - S is the start symbol for the grammar, and
  - \$ is the input right end-marker
  - We apply the following rules for all nonterminals B until nothing can be added
- 1. Place \$ in FOLLOW(S)
- 2. If there is a production  $A \Rightarrow \alpha B \beta$ ,  
then everything in **FIRST**( $\beta$ ) except  $\lambda$  is added to FOLLOW(B)
- 3. If there is a production
  - (a)  $A \Rightarrow \alpha B$ , or
  - (b)  $A \Rightarrow \alpha B \beta$ , where FIRST( $\beta$ ) contains  $\lambda$ ,
 then everything in **FOLLOW**(A) is added to FOLLOW(B)



# FOLLOW (A) Function

```

function FOLLOW(A) returns Set
  foreach A ∈ NonTERMINALS( ) do
    VisitedFollow(A) ← false
    ans ← INTERNALFOLLOW(A)
    return (ans)
end
function INTERNALFOLLOW(A) returns Set
  ans ← ∅
  if not VisitedFollow(A)
  then
    VisitedFollow(A) ← true
    foreach a ∈ OCCURRENCES(A) do
      ans ← ans ∪ FIRST(TAIL(a))
      if ALLDERIVEEMPTY(TAIL(a))
      then
        targ ← LHS(PRODUCTION(a))
        ans ← ans ∪ INTERNALFOLLOW(targ)
    return (ans)
  end
function ALLDERIVEEMPTY( $\gamma$ ) returns Boolean
  foreach  $X \in \gamma$  do
    if not SymbolDerivesEmpty( $X$ ) or  $X \in \Sigma$ 
    then return (false)
  return (true)
end

```

(17)   
 (18)   
 (19)   
 (20)   
 (21)   
 (22)   
 (23)   
 (24)

Rule 2   
 Rule 3   
 Rule 3

(when argument  $\gamma$  is empty set (3a) or the following non-terminals in  $\gamma$  all contain  $\lambda$  (3b))

- **SymbolDerivesEmpty( $\gamma$ (A))** indicates whether or not the nonterminal A can derive  $\lambda$
- **VisitedFollow(X)** is to indicate that the productions of X is already participate in the computation of FOLLOW(A)
- What happens to '\$'? It is fine. Rule 1 is not defined in the primary textbook

Figure 4.11: Algorithm for computing Follow(A).





# Example

Given the grammar,

- $E \Rightarrow TE'$
- $E' \Rightarrow +TE' \mid \lambda$
- $T \Rightarrow FT'$
- $T' \Rightarrow *FT' \mid \lambda$
- $F \Rightarrow (E) \mid id$

where E is for expression, T is for term, and F is for factor

Please find the FOLLOW set of each symbol



# Example (Cont'd)

Given the grammar,

- $E \Rightarrow TE'$
  - $E' \Rightarrow +TE' \mid \lambda$
  - $T \Rightarrow FT'$
  - $T' \Rightarrow *FT' \mid \lambda$
  - $F \Rightarrow (E) \mid id$
- E is start symbol*

- FIRST sets for E, T, F, E', and T':

- $\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ (, id \}$
- $\text{FIRST}(E') = \{ +, \lambda \}$
- $\text{FIRST}(T') = \{ *, \lambda \}$

- FOLLOW sets for E, T, F, E', and T':

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

## FOLLOW(B)

1. Place \$ in FOLLOW(S)
2. If there is a production  $A \Rightarrow \alpha B \beta$ , then everything in  $\text{FIRST}(\beta)$  except  $\lambda$  is added to FOLLOW(B)
3. If there is a production  $A \Rightarrow \alpha B$ , or  $A \Rightarrow \alpha B \beta$ , where  $\text{FIRST}(\beta)$  contains  $\lambda$ , then everything in FOLLOW(A) is added to FOLLOW(B)



# Another Example for FOLLOW

Given the grammar,

- $input \Rightarrow expression$
- $expression \Rightarrow term \ rest\_expression$
- $term \Rightarrow ID \mid parenthesized\_expression$
- $parenthesized\_expression \Rightarrow '(' \ expression \ ')'$
- $rest\_expression \Rightarrow '+' \ expression \mid \lambda$

FOLLOW(B)

1. Place \$ in FOLLOW(S)
2. If there is a production  $A \Rightarrow \alpha B \beta$ , then everything in  $FIRST(\beta)$  except  $\lambda$  is added to FOLLOW(B)
3. If there is a production (a)  $A \Rightarrow \alpha B$ , or (b)  $A \Rightarrow \alpha B \beta$ , where  $FIRST(\beta)$  contains  $\lambda$ , then everything in FOLLOW(A) is added to FOLLOW(B)

FOLLOW sets for input, expression, term, parenthesized\_expression, and rest\_expression

- FOLLOW (input) = { \$ } Rule 1
- FOLLOW (expression) = { \$, ) } Rule 3(a) got \$; Rule 2 got )
- FOLLOW (term) = FOLLOW (parenthesized\_expression) Rule 3(a)  
= { +, \$, ) } Rule 2 got +; Rule 3(b) got \$ )
- FOLLOW (rest\_expression) = { \$, ) } Rule 3(a)



## You Should ...

- Refer to Section 4.5.4 to do the exercise in Fig. 4.12 and 4.13



# Just Another Example for FOLLOW

- Fig. 4.10 grammar
- Can you derive the FOLLOW sets for A, B?
- Because of different definition used in the book, S does not contain \$

```

1 S → A B c
2 A → a
3   | λ
4 B → b
5   | λ

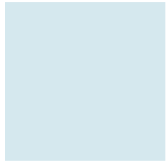
function FOLLOW(A) returns Set
  foreach A ∈ NonTERMINALS() do
    VisitedFollow(A) ← false
    ans ← INTERNALFOLLOW(A)
  return (ans)
end
function INTERNALFOLLOW(A) returns Set
  ans ← ∅
  if not VisitedFollow(A)
  then
    VisitedFollow(A) ← true
    foreach a ∈ OCCURRENCES(A) do
      ans ← ans ∪ FIRST(TAIL(a))
      if ALLDERIVEEMPTY(TAIL(a))
      then
        targ ← LHS(PRODUCTION(a))
        ans ← ans ∪ INTERNALFOLLOW(targ)
    return (ans)
  end
function ALLDERIVEEMPTY(γ) returns Boolean
  foreach X ∈ γ do
    if not SymbolDerivesEmpty(X) or X ∈ Σ
    then return (false)
  return (true)
end

```

Figure 4.11: Algorithm for computing Follow(A).

Level	Rule	Marker	Result	Comment
FOLLOW(B)				
0			FOLLOW(B)	
0	S → A <u>B</u> c	(21)	{ c }	
0		(24)	{ c }	Returns
FOLLOW(A)				
0			FOLLOW(A)	
0	S → <u>A</u> B c	(21)	{ b,c }	
0		(24)	{ b,c }	Returns
FOLLOW(S)				
0			FOLLOW(S)	
0		(24)	{ }	Returns

Figure 4.12: Follow sets for the grammar in Figure 4.10. Note that Follow(S) = { } because S does not appear on the RHS of any production.



# QUESTIONS?