

Chapter 4

Lexical and Syntax Analysis

長庚大學資訊工程學系 陳仁暉 助理教授

Tel: (03) 211-8800 Ext: 5990

E-mail: jhchen@mail.cgu.edu.tw

URL: <http://www.csie.cgu.edu.tw/jhchen>

- © All rights reserved. No part of this publication and file may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without prior written permission of Professor Jenhui Chen (E-mail: jhchen@mail.cgu.edu.tw).

Chapter 4 Topics

- Introduction
- Lexical Analysis
- The Parsing Problem
- Recursive-Descent Parsing
- Bottom-Up Parsing

Introduction

- Language implementation systems must analyze source code, regardless of the specific implementation approach
- Nearly all syntax analysis is based on a formal description of the syntax of the source language (BNF)

Syntax Analysis

- The syntax analysis portion of a language processor nearly always consists of two parts:
 - A low-level part called a *lexical analyzer* (mathematically, a finite automaton based on a **regular grammar**)
 - A high-level part called a *syntax analyzer*, or parser (mathematically, a push-down automaton based on a context-free grammar, or BNF)

Using BNF to Describe Syntax

- Provides a clear and concise syntax description
- The parser can be based directly on the BNF
- Parsers based on BNF are easy to maintain

Reasons to Separate Lexical and Syntax Analysis

- *Simplicity* – less complex approaches can be used for lexical analysis; separating them simplifies the parser
- *Efficiency* – separation allows optimization of the lexical analyzer
- *Portability* – parts of the lexical analyzer may not be portable, but the parser always is portable

Lexical Analysis

- A lexical analyzer is a pattern matcher for character strings
- A lexical analyzer is a “front-end” for the parser
- Identifies substrings of the source program that belong together – *lexemes*
 - Lexemes match a character pattern, which is associated with a lexical category called a *token*
 - `sum` is a lexeme; its token may be `IDENT`

Lexical Analysis (continued)

- The lexical analyzer is usually a function that is called by the parser when it needs the next token
- Three approaches to building a lexical analyzer:
 - Write a formal description of the tokens and use a software tool that constructs table-driven lexical analyzers given such a description
 - Design a state diagram that describes the tokens and write a program that implements the state diagram
 - Design a state diagram that describes the tokens and hand-construct a table-driven implementation of the state diagram

State Diagram Design

- A naïve state diagram would have a transition from every state on every character in the source language – such a diagram would be very large!

Lexical Analysis (cont.)

- In many cases, transitions can be combined to simplify the state diagram
 - When recognizing an identifier, all uppercase and lowercase letters are equivalent
 - Use a character class that includes all letters
 - When recognizing an integer literal, all digits are equivalent – use a digit class

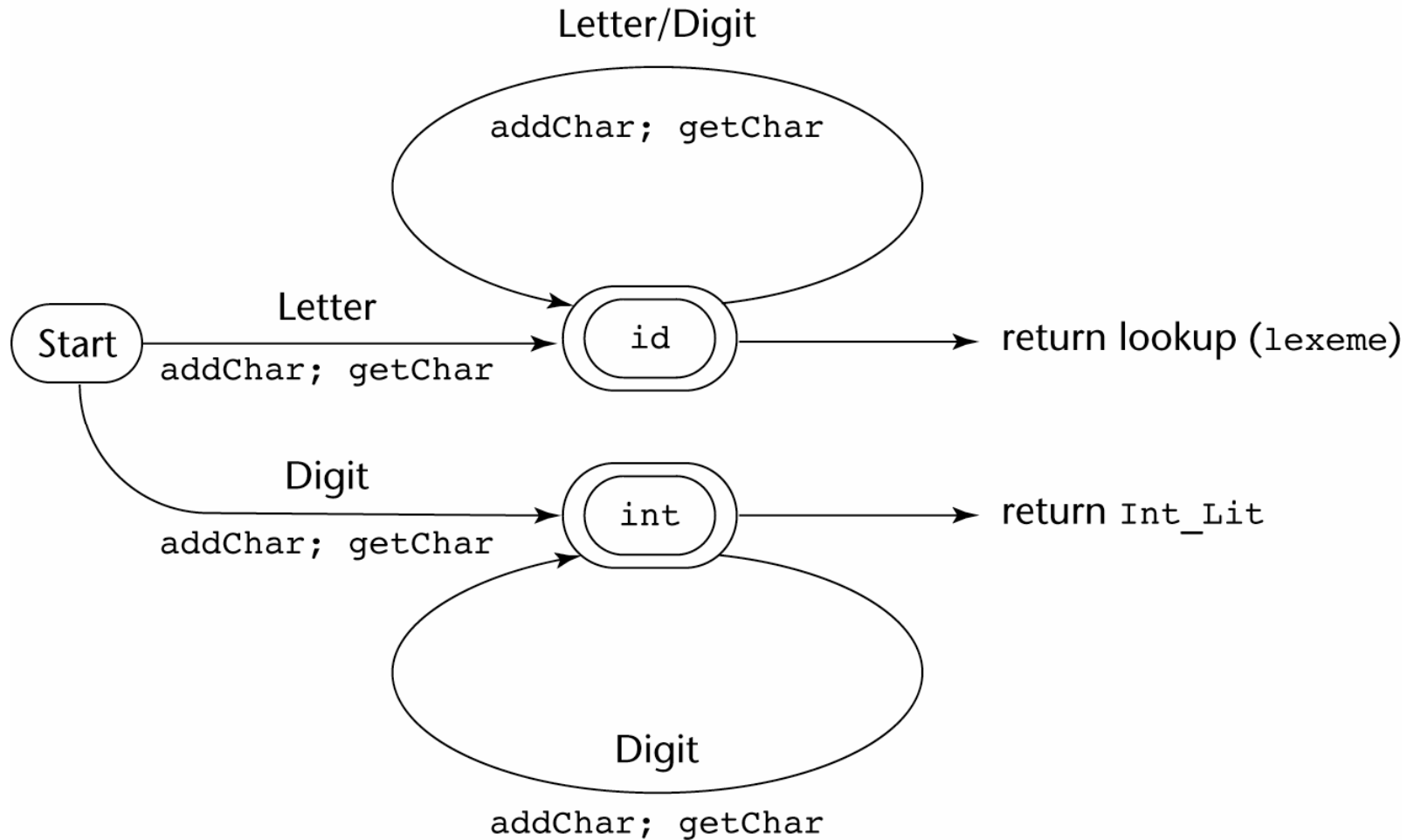
Lexical Analysis (cont.)

- Reserved words and identifiers can be recognized together (rather than having a part of the diagram for each reserved word)
 - Use a table lookup to determine whether a possible identifier is in fact a reserved word

Lexical Analysis (cont.)

- Convenient utility subprograms:
 - `getChar` – gets the next character of input, puts it in `nextChar`, determines its class and puts the class in `charClass`
 - `addChar` – puts the character from `nextChar` into the place the lexeme is being accumulated, `lexeme`
 - `lookup` – determines whether the string in `lexeme` is a reserved word (returns a code)

State Diagram



Lexical Analysis (cont.)

Implementation (assume initialization):

```
int lex() {
    getChar();
    switch (charClass) {
        case LETTER:
            addChar();
            getChar();
            while (charClass == LETTER || charClass == DIGIT)
            {
                addChar();
                getChar();
            }
            return lookup(lexeme);
            break;

        ...
    }
}
```

Lexical Analysis (cont.)

```
...
case DIGIT:
    addChar();
    getChar();
    while (charClass == DIGIT) {
        addChar();
        getChar();
    }
    return INT_LIT;
    break;
} /* End of switch */
} /* End of function lex */
```

The Parsing Problem

- Goals of the parser, given an input program:
 - Find all syntax errors; for each, produce an appropriate diagnostic message, and recover quickly
 - Produce the parse tree, or at least a trace of the parse tree, for the program

The Parsing Problem (cont.)

- Two categories of parsers
 - *Top down* – produce the parse tree, beginning at the root
 - Order is that of a leftmost derivation
 - Traces or builds the parse tree in preorder
 - *Bottom up* – produce the parse tree, beginning at the leaves
 - Order is that of the reverse of a rightmost derivation
- Parsers look only one token ahead in the input

The Parsing Problem (cont.)

- Top-down Parsers
 - Given a sentential form, $xA\alpha$, the parser must choose the correct A-rule to get the next sentential form in the leftmost derivation, using only the first token produced by A
- The most common top-down parsing algorithms:
 - Recursive descent – a coded implementation
 - LL parsers – table driven implementation
 - LL means ‘Left-to-right Leftmost derivation’

The Parsing Problem (cont.)

- Bottom-up parsers
 - Given a right sentential form, α , determine what substring of α is the right-hand side of the rule in the grammar that must be reduced to produce the previous sentential form in the right derivation
 - The most common bottom-up parsing algorithms are in the LR family
 - LR stands for ‘Left-to-right Rightmost derivation’

The Parsing Problem (cont.)

- The Complexity of Parsing
 - Parsers that work for any unambiguous grammar are complex and inefficient ($O(n^3)$, where n is the length of the input)
 - Compilers use parsers that only work for a subset of all unambiguous grammars, but do it in linear time ($O(n)$, where n is the length of the input)

Recursive–Descent Parsing

- There is a subprogram for each nonterminal in the grammar, which can parse sentences that can be generated by that nonterminal
- EBNF is ideally suited for being the basis for a recursive–descent parser, because EBNF minimizes the number of nonterminals

Recursive-Descent Parsing (cont.)

- A grammar for simple expressions:

`<expr> → <term> { (+ | -) <term> }`

`<term> → <factor> { (* | /) <factor> }`

`<factor> → id | (<expr>)`

Recursive-Descent Parsing (cont.)

- Assume we have a lexical analyzer named `lex`, which puts the next token code in `nextToken`
- The coding process when there is only one RHS:
 - For each terminal symbol in the RHS, compare it with the next input token; if they match, continue, else there is an error
 - For each nonterminal symbol in the RHS, call its associated parsing subprogram

Recursive-Descent Parsing (cont.)

```
/* Function expr
   Parses strings in the language
   generated by the rule:
   <expr> → <term> { (+ | -) <term> }
*/

void expr() {

/* Parse the first term */

    term();
    ...
}
```


Recursive-Descent Parsing (cont.)

```
/* As long as the next token is + or -, call  
lex to get the next token, and parse the  
next term */
```

```
while (nextToken == PLUS_CODE ||  
       nextToken == MINUS_CODE) {  
    lex();  
    term();  
}  
}
```

- This particular routine does not detect errors
- Convention: Every parsing routine leaves the next token in `nextToken`

Recursive–Descent Parsing (cont.)

- A nonterminal that has more than one RHS requires an initial process to determine which RHS it is to parse
 - The correct RHS is chosen on the basis of the next token of input (the lookahead)
 - The next token is compared with the first token that can be generated by each RHS until a match is found
 - If no match is found, it is a syntax error

Recursive-Descent Parsing (cont.)

```
/* Function factor
   Parses strings in the language
   generated by the rule:
   <factor> -> id | (<expr>) */

void factor() {

    /* Determine which RHS */

    if (nextToken) == ID_CODE)

    /* For the RHS id, just call lex */

    lex();
```

Recursive-Descent Parsing (cont.)

```
/* If the RHS is (<expr>) - call lex to pass
   over the left parenthesis, call expr, and
   check for the right parenthesis */

else if (nextToken == LEFT_PAREN_CODE) {
    lex();
    expr();
    if (nextToken == RIGHT_PAREN_CODE)
        lex();
    else
        error();
} /* End of else if (nextToken == ... */

else error(); /* Neither RHS matches */
}
```

Recursive–Descent Parsing (cont.)

- The LL Grammar Class
 - The Left Recursion Problem
 - If a grammar has left recursion, either direct or indirect, it cannot be the basis for a top–down parser
 - A grammar can be modified to remove left recursion

Recursive-Descent Parsing (cont.)

- The other characteristic of grammars that disallows top-down parsing is the lack of pairwise disjointness
 - The inability to determine the correct RHS on the basis of one token of lookahead
 - Def: $\text{FIRST}(\alpha) = \{a \mid \alpha \Rightarrow^* a\beta\}$
(If $\alpha \Rightarrow^* \varepsilon$, ε is in $\text{FIRST}(\alpha)$)

Recursive–Descent Parsing (cont.)

- Pairwise Disjointness Test:
 - For each nonterminal, A , in the grammar that has more than one RHS, for each pair of rules, $A \rightarrow \alpha_i$ and $A \rightarrow \alpha_j$, it must be true that

$$\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \phi$$

- Examples:

$A \rightarrow a \mid bB \mid cAb$

$A \rightarrow a \mid aB$

Recursive-Descent Parsing (cont.)

- Left factoring can resolve the problem

Replace

$\langle \text{variable} \rangle \rightarrow \text{identifier} \mid \text{identifier} [\langle \text{expression} \rangle]$

with

$\langle \text{variable} \rangle \rightarrow \text{identifier} \langle \text{new} \rangle$

$\langle \text{new} \rangle \rightarrow \varepsilon \mid [\langle \text{expression} \rangle]$

or

$\langle \text{variable} \rangle \rightarrow \text{identifier} [[\langle \text{expression} \rangle]]$

(the outer brackets are metasymbols of EBNF)

FIRST Sets

- $\text{FIRST}(\alpha)$ is the set of all terminal symbols that can begin some sentential form that starts with α
- $\text{FIRST}(\alpha) = \{a \text{ in } V_t \mid \alpha \rightarrow^* a\beta\} \cup \{\varepsilon\}$ if $\alpha \rightarrow^* \varepsilon$
- Example:
 $\langle \text{stmt} \rangle \rightarrow \text{simple} \mid \text{begin } \langle \text{stmts} \rangle \text{ end}$
 $\text{FIRST}(\langle \text{stmt} \rangle) = \{\text{simple}, \text{begin}\}$

Computing FIRST sets

Initially $\text{FIRST}(A)$ is empty

1. For productions $A \rightarrow a \beta$, where $a \in V_t$
Add $\{ a \}$ to $\text{FIRST}(A)$
2. For productions $A \rightarrow \varepsilon$
Add $\{ \varepsilon \}$ to $\text{FIRST}(A)$
3. For productions $A \rightarrow \alpha B \beta$, where $\alpha \xrightarrow{*} \varepsilon$ and
NOT $(B \rightarrow \varepsilon)$
Add $\text{FIRST}(\alpha B)$ to $\text{FIRST}(A)$
4. For productions $A \rightarrow \alpha$, where $\alpha \xrightarrow{*} \varepsilon$
Add $\text{FIRST}(\alpha)$ and $\{ \varepsilon \}$ to $\text{FIRST}(A)$

To compute FIRST across strings of terminals and non-terminals:

$$\text{FIRST}(\varepsilon) = \{ \varepsilon \}$$

$$\begin{aligned} \text{FIRST}(A\alpha) &= A \quad \text{if } A \text{ is a terminal} \\ &= \text{FIRST}(A) \cup \text{FIRST}(\alpha) \\ &\quad \text{if } A \rightarrow \varepsilon \\ &= \text{FIRST}(A) \text{ otherwise} \end{aligned}$$

Example 1

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

Example 1

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

Example 2

- $P \rightarrow i \mid c \mid n T S$
- $Q \rightarrow P \mid a S \mid b S c S T$
- $R \rightarrow b \mid \varepsilon$
- $S \rightarrow c \mid R n \mid \varepsilon$
- $T \rightarrow R S q$
- $\text{FIRST}(P) =$
- $\text{FIRST}(Q) =$
- $\text{FIRST}(R) =$
- $\text{FIRST}(S) =$
- $\text{FIRST}(T) =$

Example 2

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

Example 3

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

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.
- $\text{FOLLOW}(A) = \{a \text{ in } V_t \mid S \rightarrow^+ \dots Aa\dots\} \cup \{\$$
(end of file) $\}$ if $S \rightarrow^+ \dots A$
- For example, consider $L \rightarrow^+ (())(L)L$ --
Both ')' and end of file can follow L

Computing FOLLOW(A)

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

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

$S \xrightarrow{+} \alpha B \beta \xrightarrow{+} \alpha A X \beta \xrightarrow{+} \alpha A c \delta \beta$

4. Productions of the form $B \rightarrow \alpha A$

or $B \rightarrow \alpha A \beta$ where $\beta \rightarrow^* \varepsilon$

Add FOLLOW(B) to FOLLOW(A)

INTUITION:

– Suppose $B \rightarrow Y A$

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

– Suppose $B \rightarrow A X$ and $X \rightarrow \varepsilon$

$S \rightarrow^+ \alpha B \beta \rightarrow \alpha A X \beta \rightarrow \alpha A \beta$

NOTE: ε *never* in FOLLOW sets

Example 4

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

Example 4

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

Example 5

- $S \rightarrow (A) \mid \varepsilon$
- $A \rightarrow T E$
- $E \rightarrow , T E \mid \varepsilon$
- $T \rightarrow (A) \mid a \mid b \mid c$

- $\text{FIRST}(T) = \{(, a, b, c\}$
- $\text{FIRST}(E) = \{', \varepsilon\}$
- $\text{FIRST}(A) = \{(, a, b, c\}$
- $\text{FIRST}(S) = \{(, \varepsilon\}$

- $\text{FOLLOW}(S) =$
- $\text{FOLLOW}(A) =$
- $\text{FOLLOW}(E) =$
- $\text{FOLLOW}(T) =$

Example 5

- $S \rightarrow (A) \mid \varepsilon$
- $A \rightarrow T E$
- $E \rightarrow , T E \mid \varepsilon$
- $T \rightarrow (A) \mid a \mid b \mid c$
- $\text{FIRST}(T) = \{(, a, b, c\}$
- $\text{FIRST}(E) = \{', \varepsilon\}$
- $\text{FIRST}(A) = \{(, a, b, c\}$
- $\text{FIRST}(S) = \{(, \varepsilon\}$
- $\text{FOLLOW}(S) = \{\$\}$
- $\text{FOLLOW}(A) = \{) \}$
- $\text{FOLLOW}(E) = \{) \}$
- $\text{FOLLOW}(T) = \{ ', ',) \}$

Example 6

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

Example 6

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

Example 7

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

Example 7

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

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

Writing an LL(1) Grammar

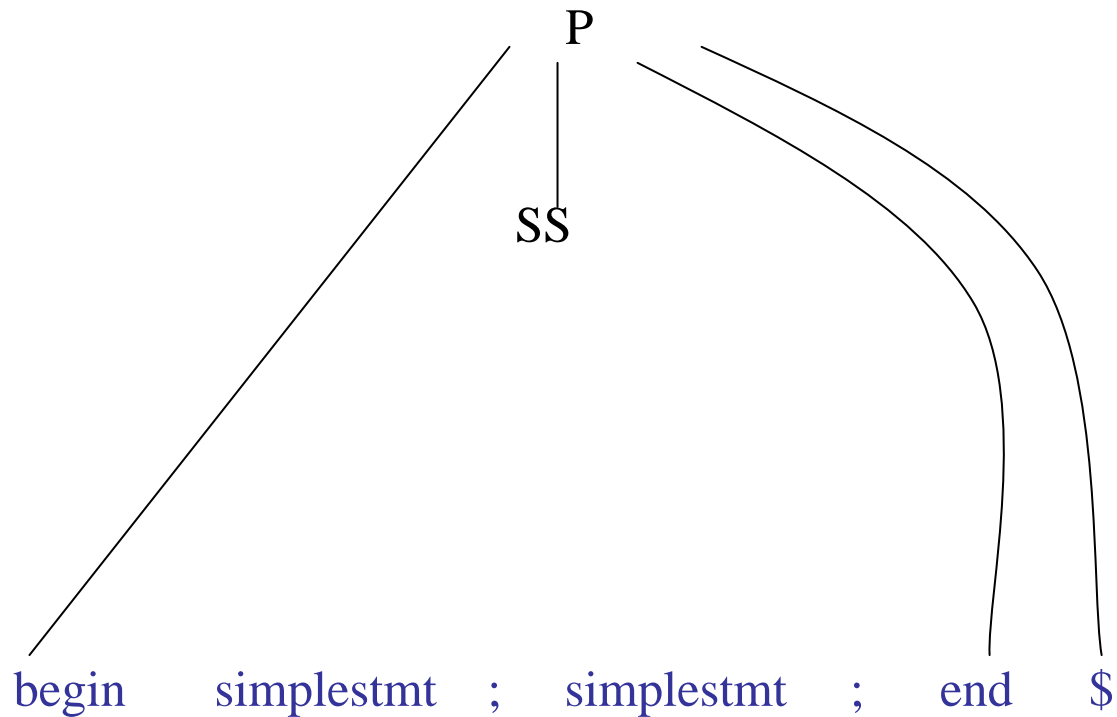
- The two most common obstacles to “LL(1)-ness” are
 - Left recursion
 - Common prefixes

Top Down (LL) Parsing

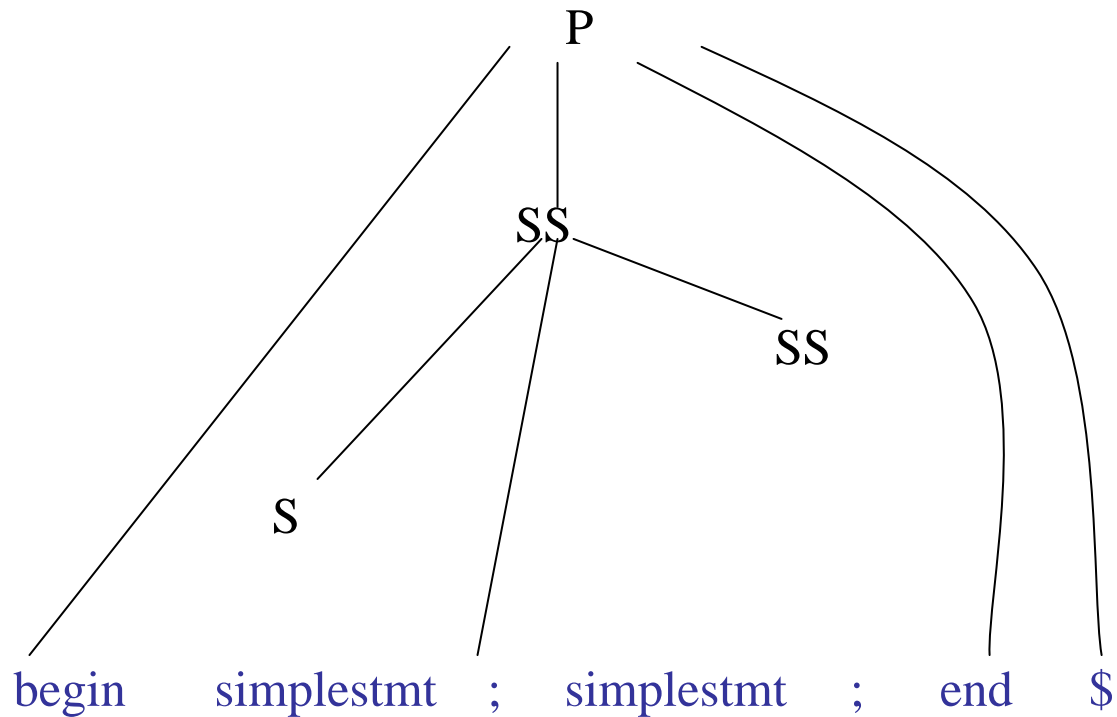
P

begin simplestmt ; simplestmt ; end \$

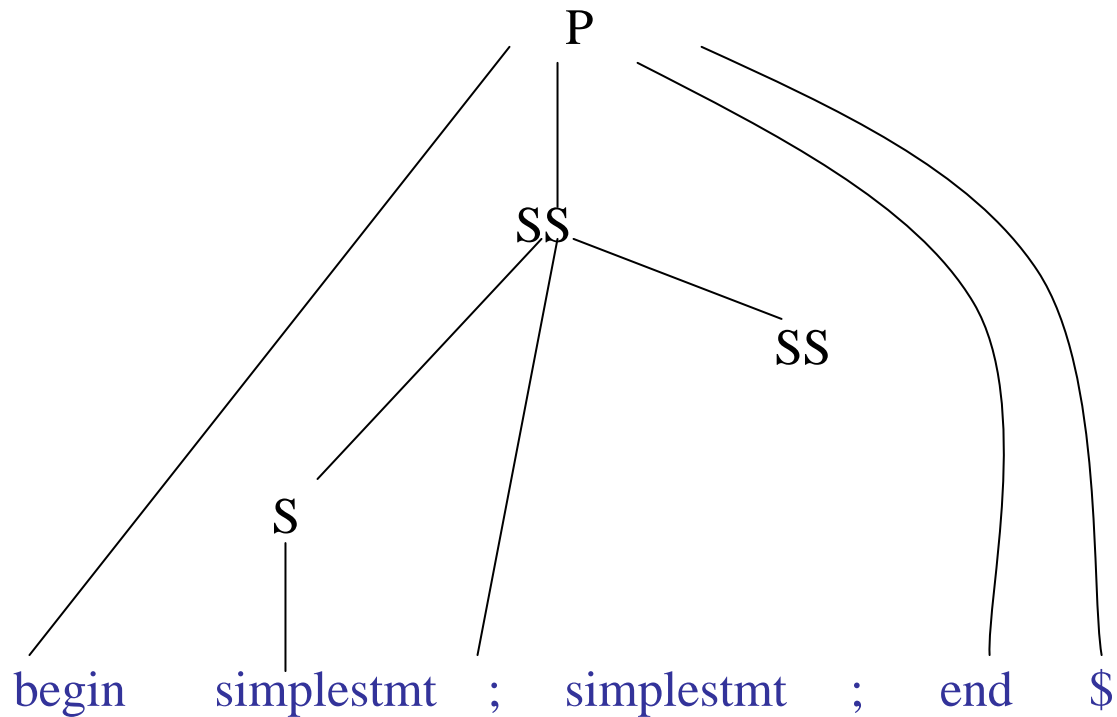
Top Down (LL) Parsing



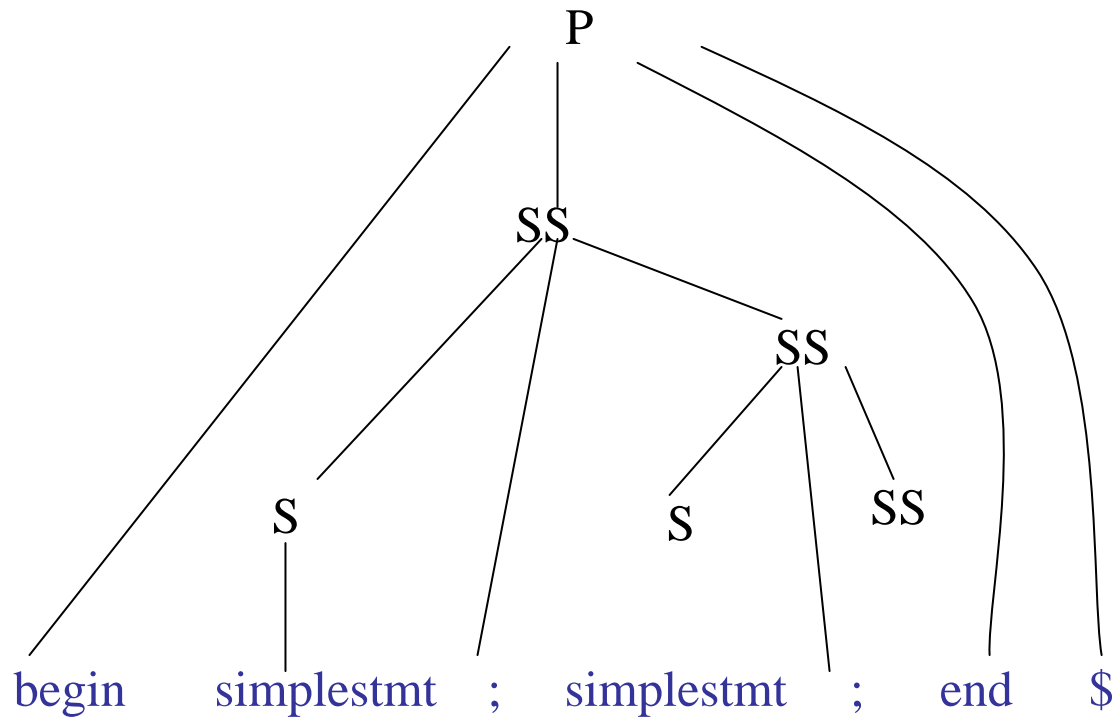
Top Down (LL) Parsing



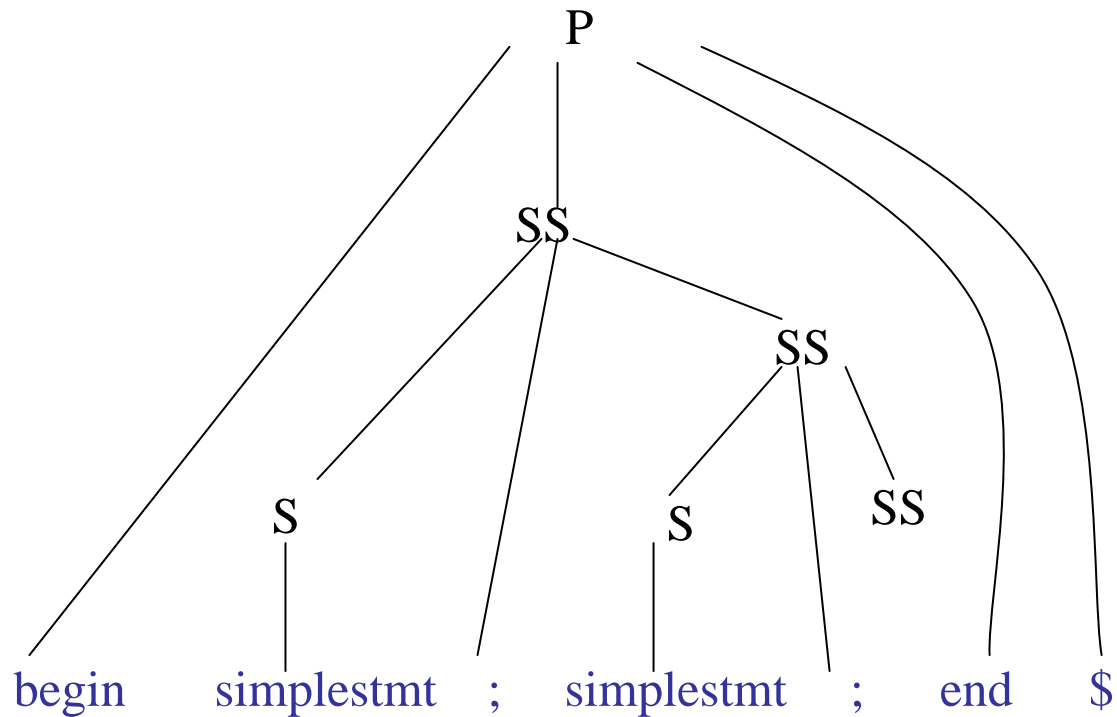
Top Down (LL) Parsing



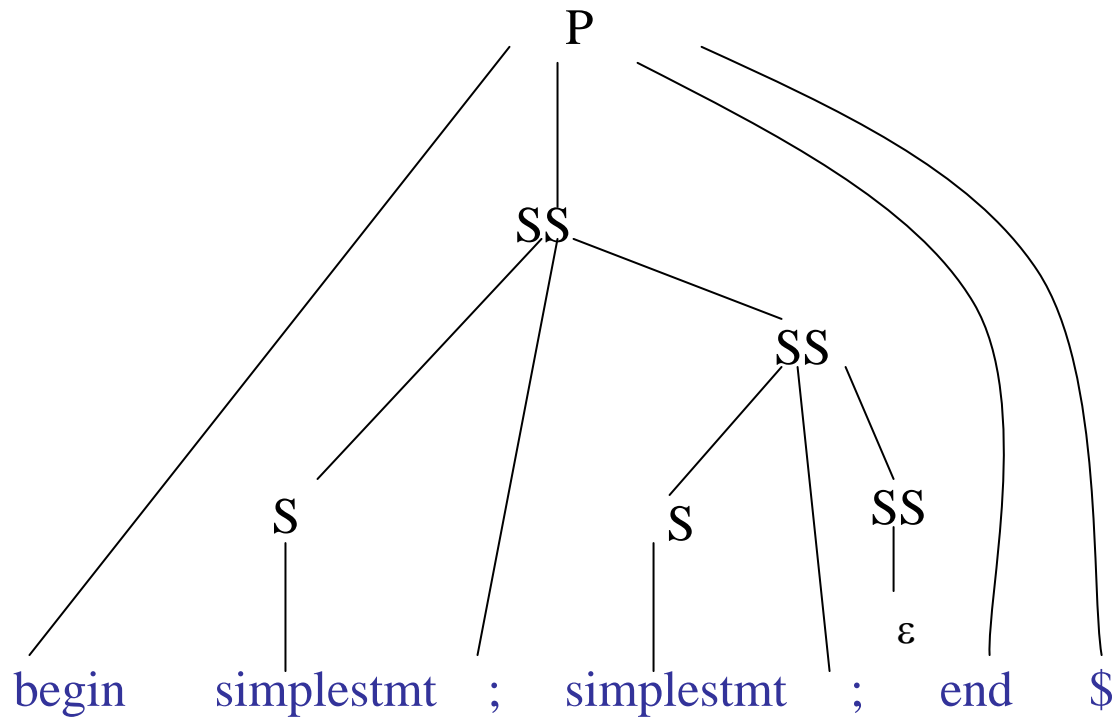
Top Down (LL) Parsing



Top Down (LL) Parsing



Top Down (LL) Parsing



Grammar

$S \rightarrow a B$

$\quad | b C$

$B \rightarrow b b C$

$C \rightarrow c c$

Two strings in the language: `abbcc` and `bcc`
Can choose between them based on the first
character of the input.

LL(k) parsing

- Process input k symbols at a time.
- Initially, current non-terminal is start symbol.
- Algorithm
 - Given next k input tokens and current non-terminal T , choose a rule $R (T \rightarrow \dots)$
 - For each element X in rule R from left to right,
if X is a non-terminal, call function for X
else if symbol X is a terminal, see if next input symbol matches X ; if so, update from the input
- Typically, we consider LL(1)

Two Approaches

- Recursive Descent parsing
 - Code tailored to the grammar
- Table Driven – predictive parsing
 - Table tailored to the grammar
 - General Algorithm

Writing a Recursive Descent Parser

- Procedure for each non-terminal.

Use next token (lookahead) to choose which production to mimic.

- for non-terminal X , call procedure $X()$
- for terminals X , call 'match(X)'

- ```
match(symbol) {
 if (symbol = lookahead)
 lookahead = yylex()
 else error() }
```

- Call `yylex()` before the first call to get first lookahead.



# Back to grammar

---

```
S() {
 if (lookahead==a) { match(a);B(); }
 else if (lookahead == b) { match(b);
 C(); }
 else error("expecting a or b");
}
B() {match(b); match(b); C();}
C() { match(c) ; match(c) ;}
```

```
main() {
 lookahead==yylex();
 S();
}
```

$$S \rightarrow a B$$
$$| b C$$
$$B \rightarrow b b C$$
$$C \rightarrow c c$$

# Parsing abbcc

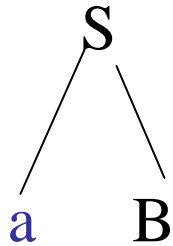
---

S

Remaining input: abbcc

# Parsing abbcc

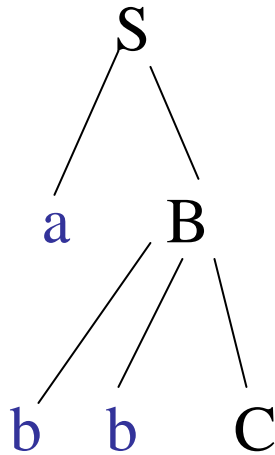
---



Remaining input: bbcc

# Parsing abbcc

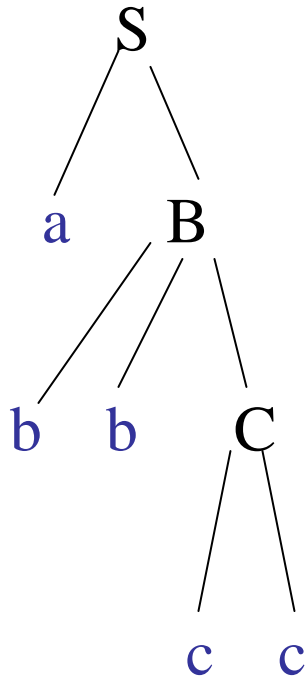
---



Remaining input: cc

# Parsing abbcc

---



Remaining input:

# How do we find the lookaheads?

---

- Can compute PREDICT sets from FIRST and FOLLOW
- $\text{PREDICT}(A \rightarrow \alpha) =$   
     $\text{FIRST}(\alpha) - \{\varepsilon\} \cup \text{FOLLOW}(A)$  if  $\varepsilon$  in  $\text{FIRST}(\alpha)$   
     $\text{FIRST}(\alpha)$  if  $\varepsilon$  not in  $\text{FIRST}(\alpha)$

**NOTE:  $\varepsilon$  never in PREDICT sets**

For  $\text{LL}(k)$  grammars, the PREDICT sets for a given non-terminal will be disjoint.

# Example

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

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

# Parsing $a + b * c$

---

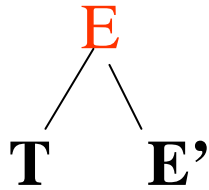
E

Remaining input:  $a+b*c$



# Parsing $a + b * c$

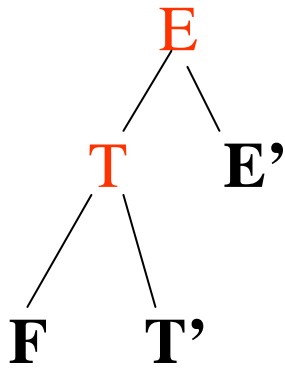
---



Remaining input:  $a+b*c$

# Parsing $a + b * c$

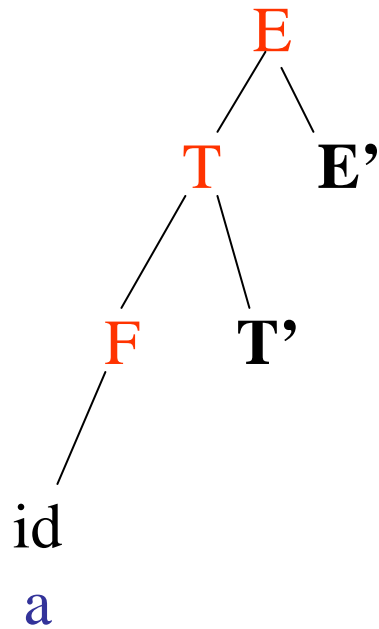
---



Remaining input:  $a+b*c$

# Parsing $a + b * c$

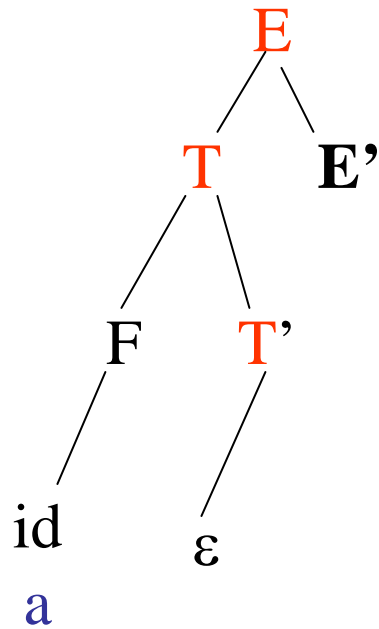
---



Remaining input:  $+b*c$

# Parsing $a + b * c$

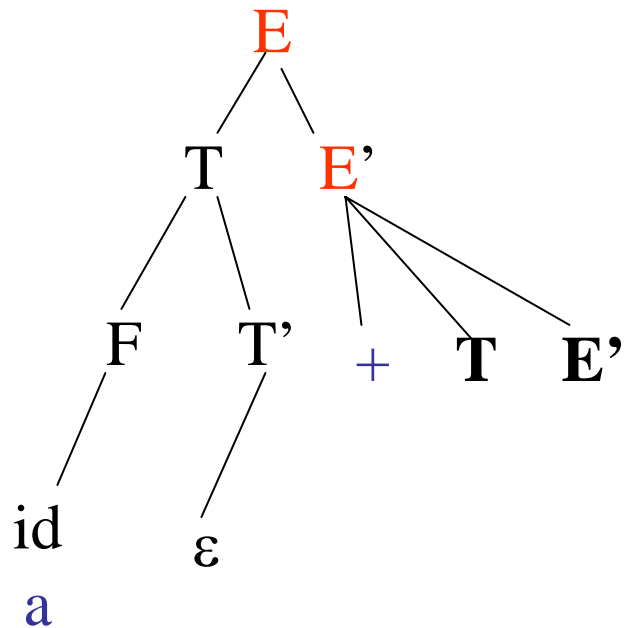
---



Remaining input:  $+b*c$

# Parsing $a + b * c$

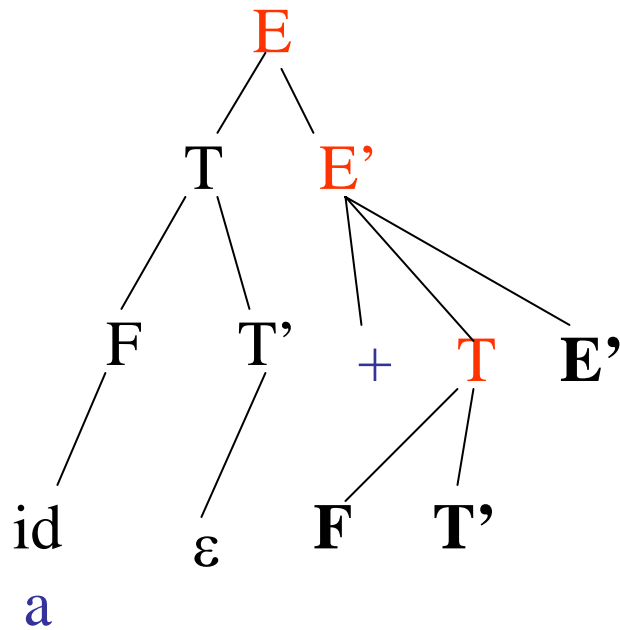
---



Remaining input:  $b * c$

# Parsing $a + b * c$

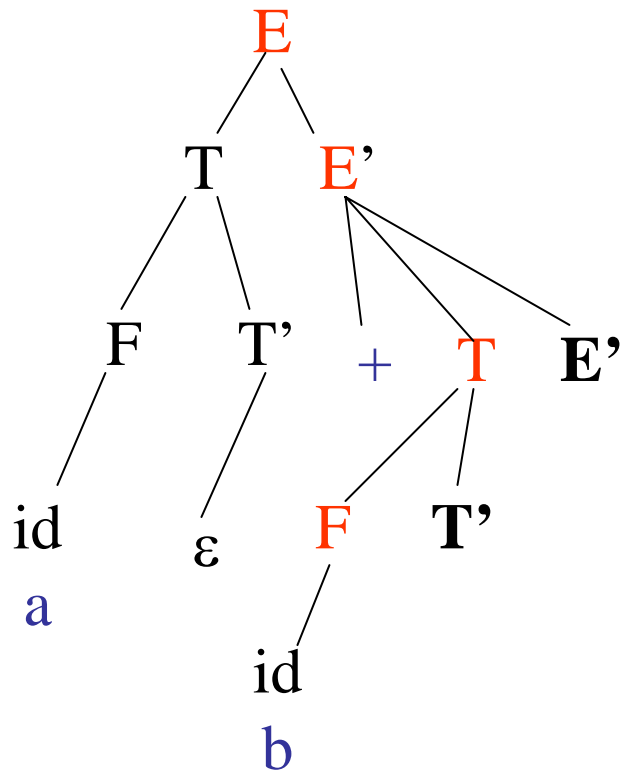
---



Remaining input:  $b * c$

# Parsing $a + b * c$

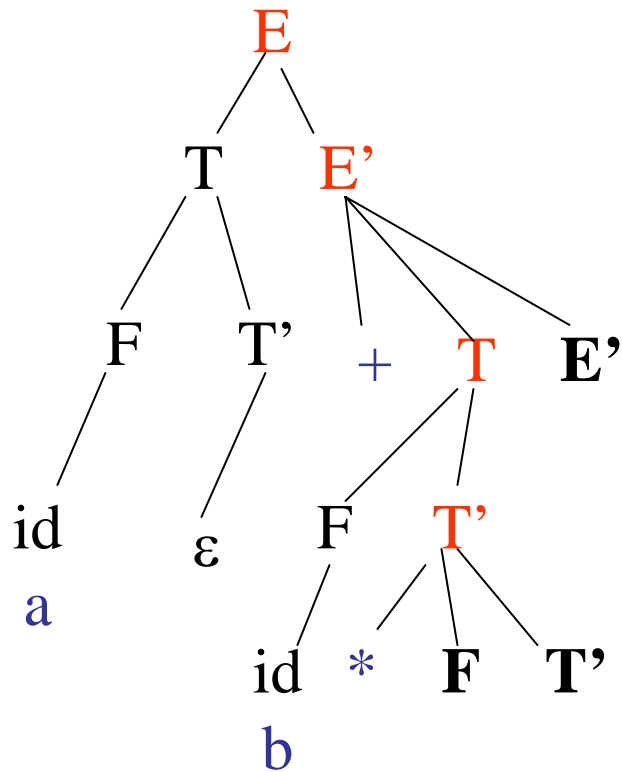
---



Remaining input:  $*c$

# Parsing $a + b * c$

---

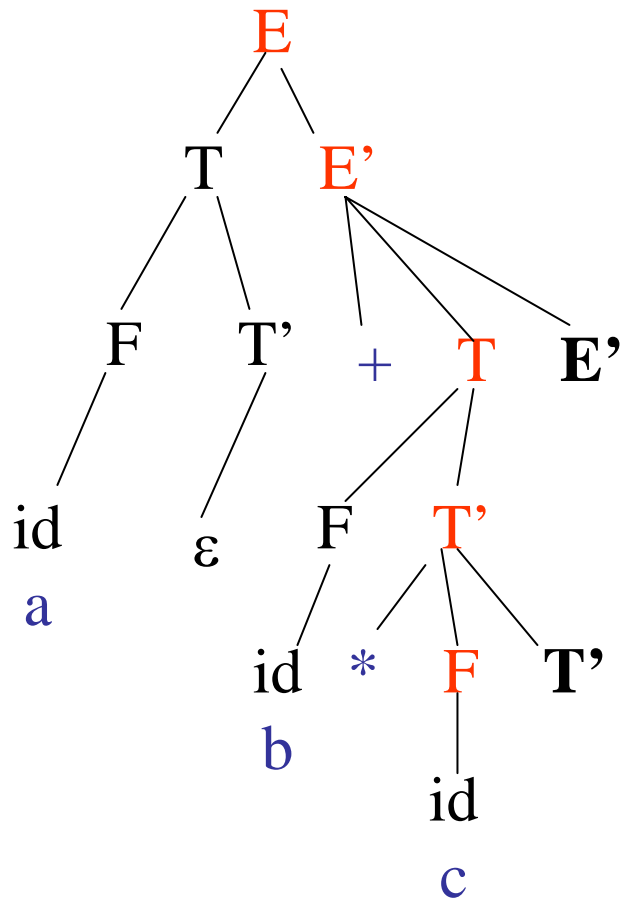


Remaining input:  $c$



# Parsing $a + b * c$

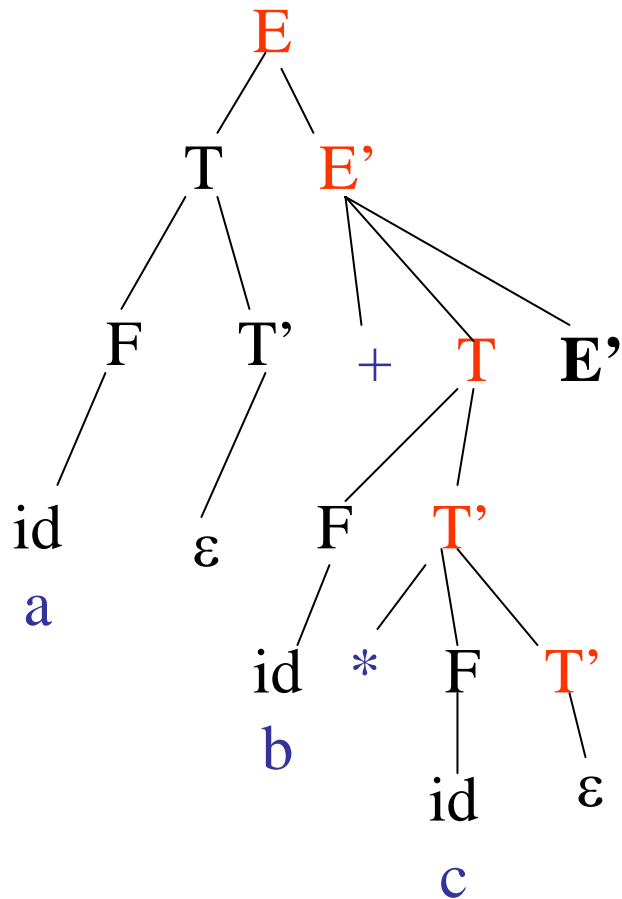
---



Remaining input:

# Parsing $a + b * c$

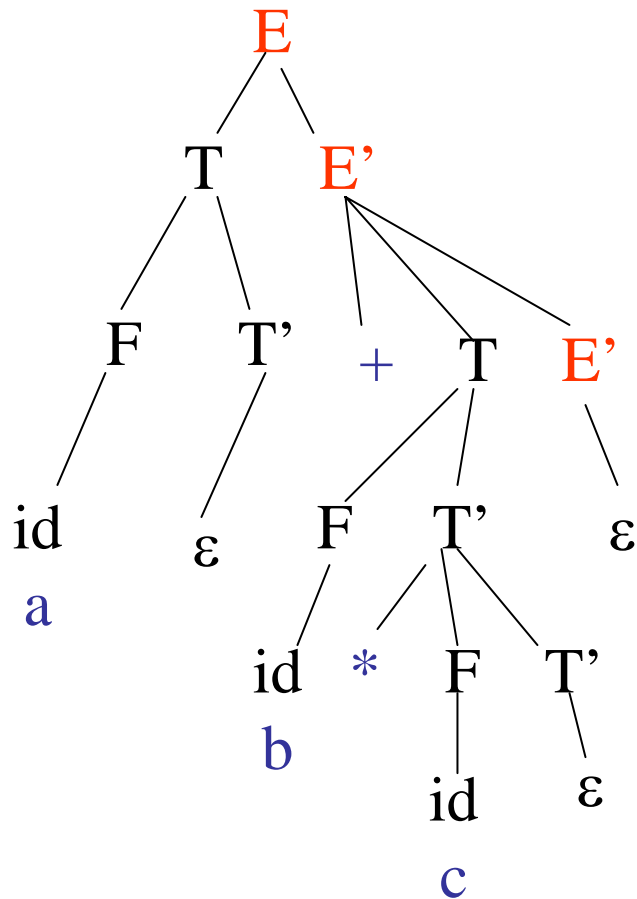
---



Remaining input:

# Parsing $a + b * c$

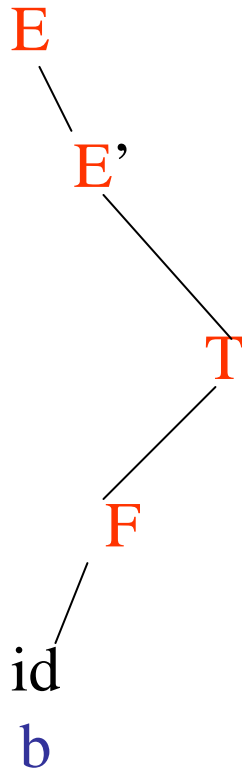
---



Remaining input:

# Stacks in Recursive Descent Parsing

---



- Runtime stack
- Procedure activations correspond to a path in parse tree from root to some interior node

# LL(1) Predictive Parse Tables

---

An LL(1) Parse table is a mapping  $T: V_n \times V_t \rightarrow$   
production  $P$  or error

1. For all productions  $A \rightarrow \alpha$  do
  - For each terminal  $a$  in  $\text{Predict}(A \rightarrow \alpha)$ ,  
 $T[A][a] = A \rightarrow \alpha$
2. Every undefined table entry is an error.

# Using LL(1) Parse Tables

---

## ALGORITHM

INPUT: token sequence to be parsed, followed by '\$' (end of file)

## DATA STRUCTURES:

- Parse stack: Initialized by pushing '\$' and then pushing the start symbol
- Parse table T

---

```
push($); push(start_symbol); lookahead = yylex()
repeat
 X = pop(stack)
 if X is a terminal symbol or $ then
 if X = lookahead then
 lookahead = yylex()
 else error()
 else /* X is non-terminal */
 if T[X][lookahead] = X → Y1 Y2 ...Ym
 push(Ym) ... push (Y1)
 else error()
until X = $ token
```

# Expression Grammar

---

| NT/T | +                         | *                    | (                   | )                         | ID                 | \$                        |
|------|---------------------------|----------------------|---------------------|---------------------------|--------------------|---------------------------|
| E    |                           |                      | $\rightarrow T E'$  |                           | $\rightarrow T E'$ |                           |
| E'   | $\rightarrow + T E'$      |                      |                     | $\rightarrow \varepsilon$ |                    | $\rightarrow \varepsilon$ |
| T    |                           |                      | $\rightarrow F T'$  |                           | $\rightarrow F T'$ |                           |
| T'   | $\rightarrow \varepsilon$ | $\rightarrow * F T'$ |                     | $\rightarrow \varepsilon$ |                    | $\rightarrow \varepsilon$ |
| F    |                           |                      | $\rightarrow ( E )$ |                           | $\rightarrow ID$   |                           |



# Parsing $a + b * c$

| Stack    | Input   | Action                       |
|----------|---------|------------------------------|
| \$E      | a+b*c\$ | $E \rightarrow T E'$         |
| \$E'T    | a+b*c\$ | $T \rightarrow F T'$         |
| \$E'T'F  | a+b*c\$ | $F \rightarrow id$           |
| \$E'T'id | a+b*c\$ | match                        |
| \$E'T'   | +b*c\$  | $T' \rightarrow \varepsilon$ |
| \$E'     | +b*c\$  | $E' \rightarrow + T E'$      |
| \$E'T+   | +b*c\$  | match                        |
| \$E'T    | b*c\$   | $T \rightarrow F T'$         |

| Stack    | Input | Action                       |
|----------|-------|------------------------------|
| \$E'T'F  | b*c\$ | $F \rightarrow id$           |
| \$E'T'id | b*c\$ | match                        |
| \$E'T'   | *c\$  | $T' \rightarrow * F T'$      |
| \$E'T'F* | *c\$  | match                        |
| \$E'T'F  | c\$   | $F \rightarrow id$           |
| \$E'T'id | c\$   | match                        |
| \$E'T'   | \$    | $T' \rightarrow \varepsilon$ |
| \$E'     | \$    | $E' \rightarrow \varepsilon$ |
| \$       | \$    | accept                       |

# Stack in Predictive Parsing

---

- Algorithm data structure
- Holds terminals and non-terminals from the grammar
  - terminals – still need to be matched from the input
  - non-terminals – still need to be expanded

# Making a grammar LL(1)

---

- Not all context free languages have LL(1) grammars
- Can show a grammar is not LL(1) by looking at the predict sets
  - For LL(a) grammars, the PREDICT sets for a given non-terminal will be disjoint.

# Example

| Production                | Predict                                |
|---------------------------|----------------------------------------|
| $E \rightarrow E + T$     | $= \text{FIRST}(E) = \{(, \text{id}\}$ |
| $E \rightarrow T$         | $= \text{FIRST}(T) = \{(, \text{id}\}$ |
| $T \rightarrow T * F$     | $= \text{FIRST}(T) = \{(, \text{id}\}$ |
| $T \rightarrow F$         | $= \text{FIRST}(F) = \{(, \text{id}\}$ |
| $F \rightarrow \text{id}$ | $= \{\text{id}\}$                      |
| $F \rightarrow ( E )$     | $= \{()\}$                             |

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

Two problems: E and T

# Making a non-LL(1) grammar LL(1)

---

- Eliminate common prefixes

Ex:  $A \rightarrow B a C D \mid B a C E$

- Transform left recursion to right recursion

Ex:  $E \rightarrow E + T \mid T$

# Eliminate Common Prefixes

---

- $A \rightarrow \alpha \beta \mid \alpha \delta$

Can become:

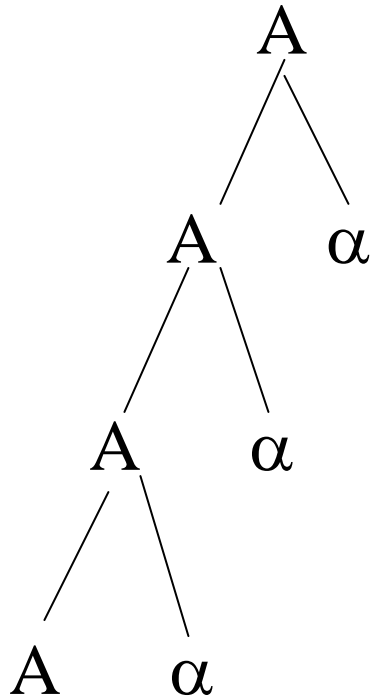
$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta \mid \delta$$

Doesn't always remove the problem. *Why?*

# Why is left recursion a problem?

---



# Remove Left Recursion

---

$A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \dots \mid \beta_1 \mid \beta_2 \mid \dots$

becomes

$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots$

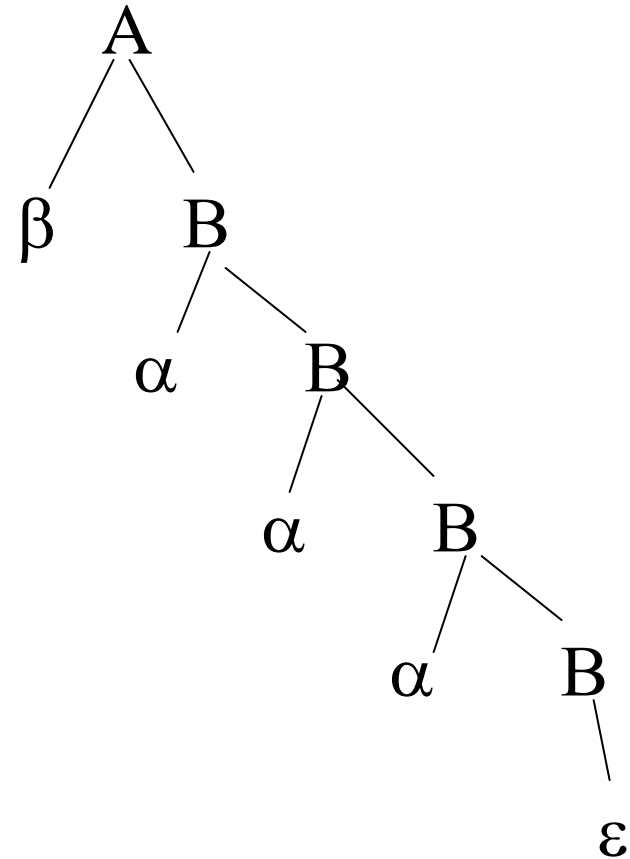
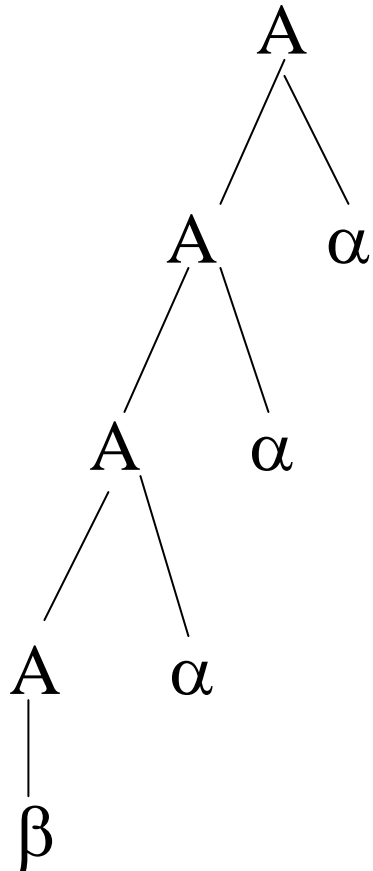
$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \varepsilon$

The left recursion becomes right recursion



$A \rightarrow A \alpha \mid \beta$  becomes  $A \rightarrow \beta B, B \rightarrow \alpha B \mid \varepsilon$

---



# Bottom-up Parsing

---

- The parsing problem is finding the correct RHS in a right-sentential form to reduce to get the previous right-sentential form in the derivation

# Bottom-up Parsing (cont.)

---

- Intuition about handles:
  - Def:  $\beta$  is the *handle* of the right sentential form  $\gamma = \alpha\beta w$  if and only if  $S \Rightarrow^* \alpha A w \Rightarrow \alpha \beta w$
  - Def:  $\beta$  is a *phrase* of the right sentential form  $\gamma$  if and only if  $S \Rightarrow^* \gamma = \alpha_1 A \alpha_2 \Rightarrow + \alpha_1 \beta \alpha_2$
  - Def:  $\beta$  is a *simple phrase* of the right sentential form  $\gamma$  if and only if  $S \Rightarrow^* \gamma = \alpha_1 A \alpha_2 \Rightarrow \alpha_1 \beta \alpha_2$

# Bottom-up Parsing (cont.)

---

- Intuition about handles:
  - The handle of a right sentential form is its leftmost simple phrase
  - Given a parse tree, it is now easy to find the handle
  - Parsing can be thought of as handle pruning

# Bottom-up Parsing (cont.)

---

- Shift-Reduce Algorithms
  - Reduce is the action of replacing the handle on the top of the parse stack with its corresponding LHS
  - Shift is the action of moving the next token to the top of the parse stack

# Bottom-up Parsing (cont.)

---

- Advantages of LR parsers:
  - They will work for nearly all grammars that describe programming languages.
  - They work on a larger class of grammars than other bottom-up algorithms, but are as efficient as any other bottom-up parser.
  - They can detect syntax errors as soon as it is possible.
  - The LR class of grammars is a superset of the class parsable by LL parsers.

# Bottom-up Parsing (cont.)

---

- LR parsers must be constructed with a tool
- Knuth's insight: A bottom-up parser could use the entire history of the parse, up to the current point, to make parsing decisions
  - There were only a finite and relatively small number of different parse situations that could have occurred, so the history could be stored in a parser state, on the parse stack

# Bottom-up Parsing (cont.)

---

- An LR configuration stores the state of an LR parser

$(S_0 X_1 S_1 X_2 S_2 \dots X_m S_m, a_i a_{i+1} \dots a_n \$)$



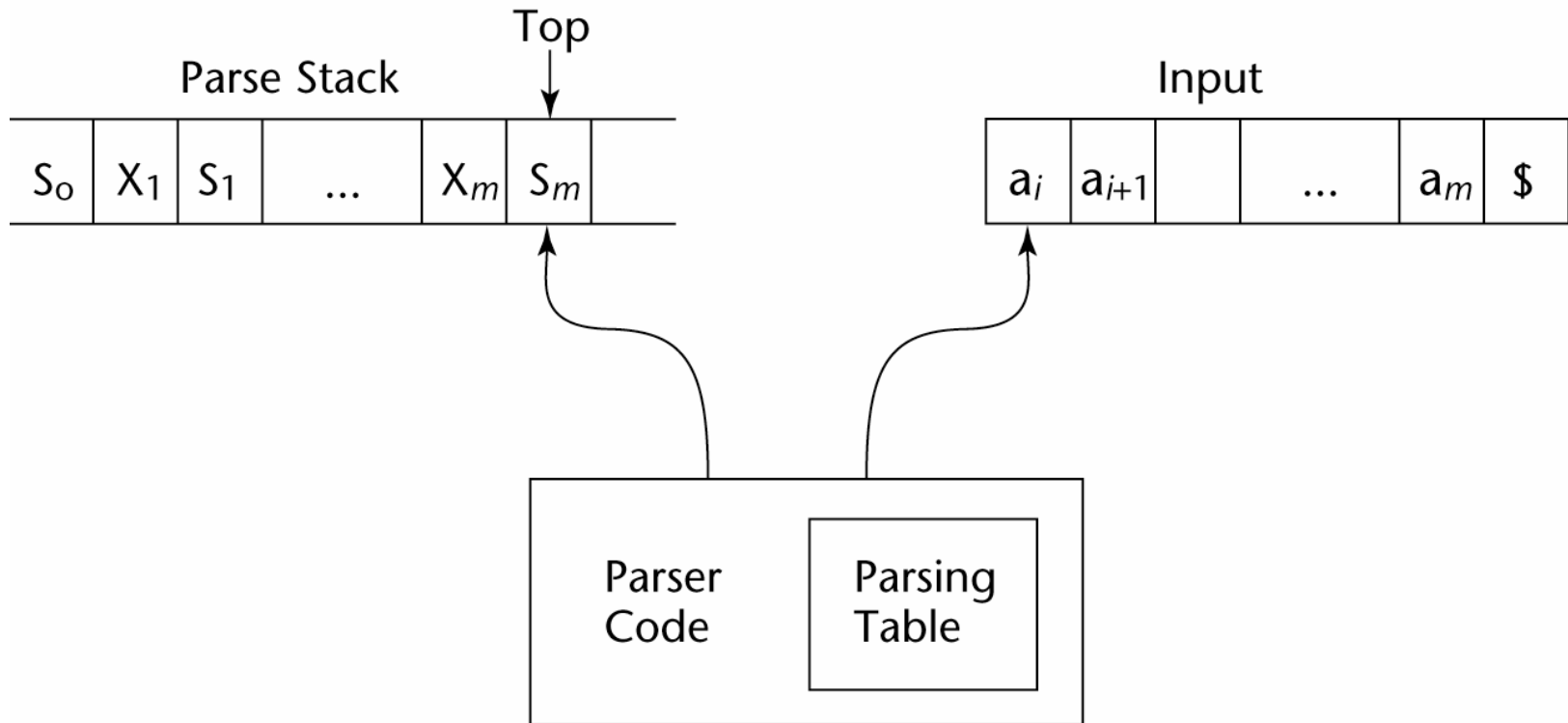
# Bottom-up Parsing (cont.)

---

- LR parsers are table driven, where the table has two components, an ACTION table and a GOTO table
  - The ACTION table specifies the action of the parser, given the parser state and the next token
    - Rows are state names; columns are terminals
  - The GOTO table specifies which state to put on top of the parse stack after a reduction action is done
    - Rows are state names; columns are nonterminals

# Structure of An LR Parser

---



# Bottom-up Parsing (cont.)

---

- Initial configuration:  $(S_0, a_1 \dots a_n \$)$
- Parser actions:
  - If  $\text{ACTION}[S_m, a_i] = \text{Shift } S$ , the next configuration is:  
 $(S_0 X_1 S_1 X_2 S_2 \dots X_m S_m a_i S, a_{i+1} \dots a_n \$)$
  - If  $\text{ACTION}[S_m, a_i] = \text{Reduce } A \rightarrow \beta$  and  $S = \text{GOTO}[S_{m-r}, A]$ , where  $r = \text{length of } \beta$ , the next configuration is  
 $(S_0 X_1 S_1 X_2 S_2 \dots X_{m-r} S_{m-r} A S, a_i a_{i+1} \dots a_n \$)$

# Bottom-up Parsing (cont.)

---

- Parser actions (continued):
  - If  $\text{ACTION}[S_m, a_i] = \text{Accept}$ , the parse is complete and no errors were found.
  - If  $\text{ACTION}[S_m, a_i] = \text{Error}$ , the parser calls an error-handling routine.

# LR Parsing Table

| State | Action |    |    |    |     |        | Goto |   |    |
|-------|--------|----|----|----|-----|--------|------|---|----|
|       | id     | +  | *  | (  | )   | \$     | E    | T | F  |
| 0     | S5     |    | S4 |    |     |        | 1    | 2 | 3  |
| 1     |        | S6 |    |    |     | accept |      |   |    |
| 2     |        | R2 | S7 |    | R2  | R2     |      |   |    |
| 3     |        | R4 | R4 |    | R4  | R4     |      |   |    |
| 4     | S5     |    |    | S4 |     |        | 8    | 2 | 3  |
| 5     |        | R6 | R6 |    | R6  | R6     |      |   |    |
| 6     | S5     |    |    | S4 |     |        |      | 9 | 3  |
| 7     | S5     |    |    | S4 |     |        |      |   | 10 |
| 8     |        | S6 |    |    | S11 |        |      |   |    |
| 9     |        | R1 | S7 |    | R1  | R1     |      |   |    |
| 10    |        | R3 | R3 |    | R3  | R3     |      |   |    |
| 11    |        | R5 | R5 |    | R5  | R5     |      |   |    |

# Bottom-up Parsing (cont.)

---

- A parser table can be generated from a given grammar with a tool, e.g., `yacc`

# Summary

---

- Syntax analysis is a common part of language implementation
- A lexical analyzer is a pattern matcher that isolates small-scale parts of a program
  - Detects syntax errors
  - Produces a parse tree
- A recursive-descent parser is an LL parser
  - EBNF
- Parsing problem for bottom-up parsers: find the substring of current sentential form
- The LR family of shift-reduce parsers is the most common bottom-up parsing approach