

Lexical and Syntax Analysis

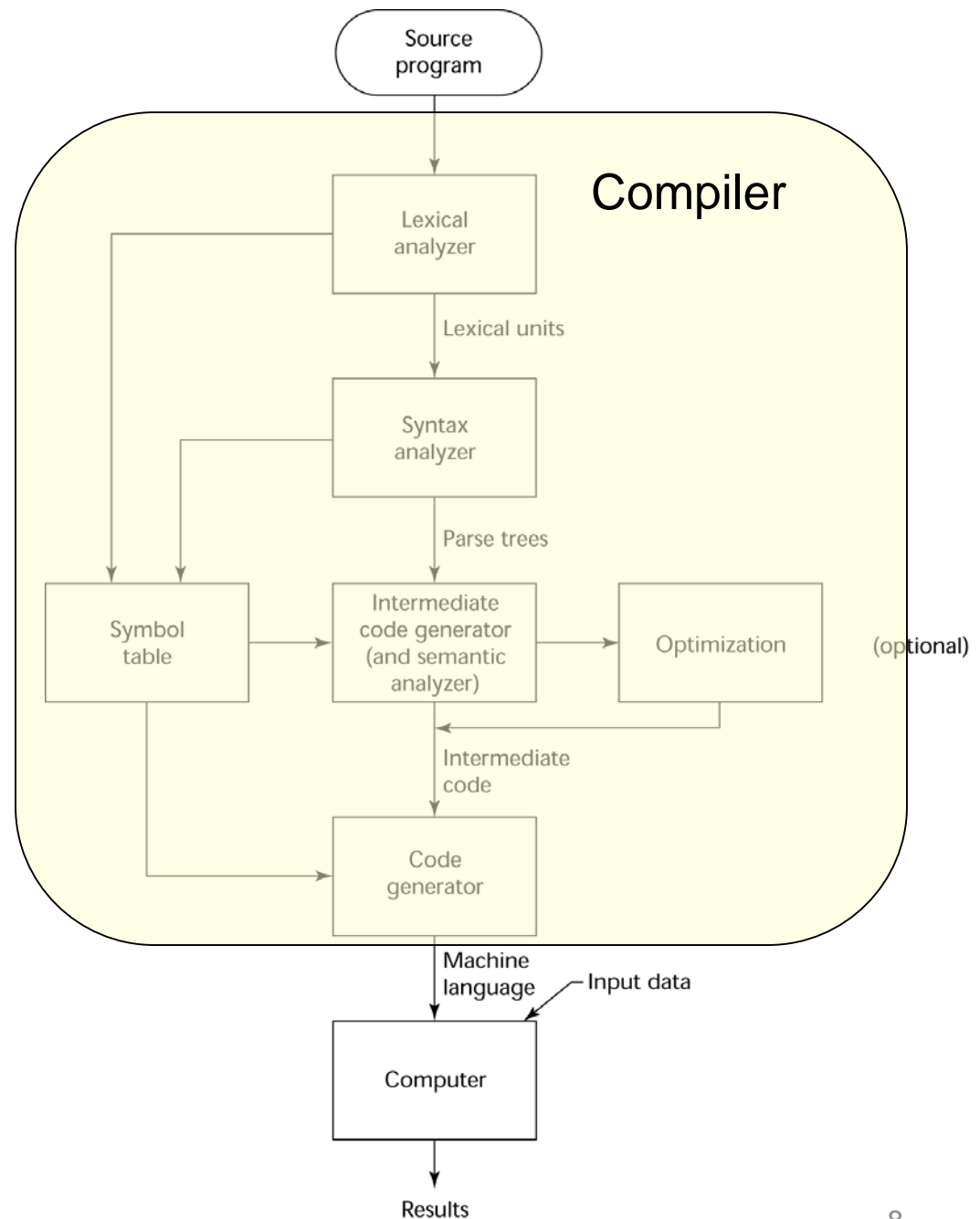
Chapter 4

Chapter 4 Topics

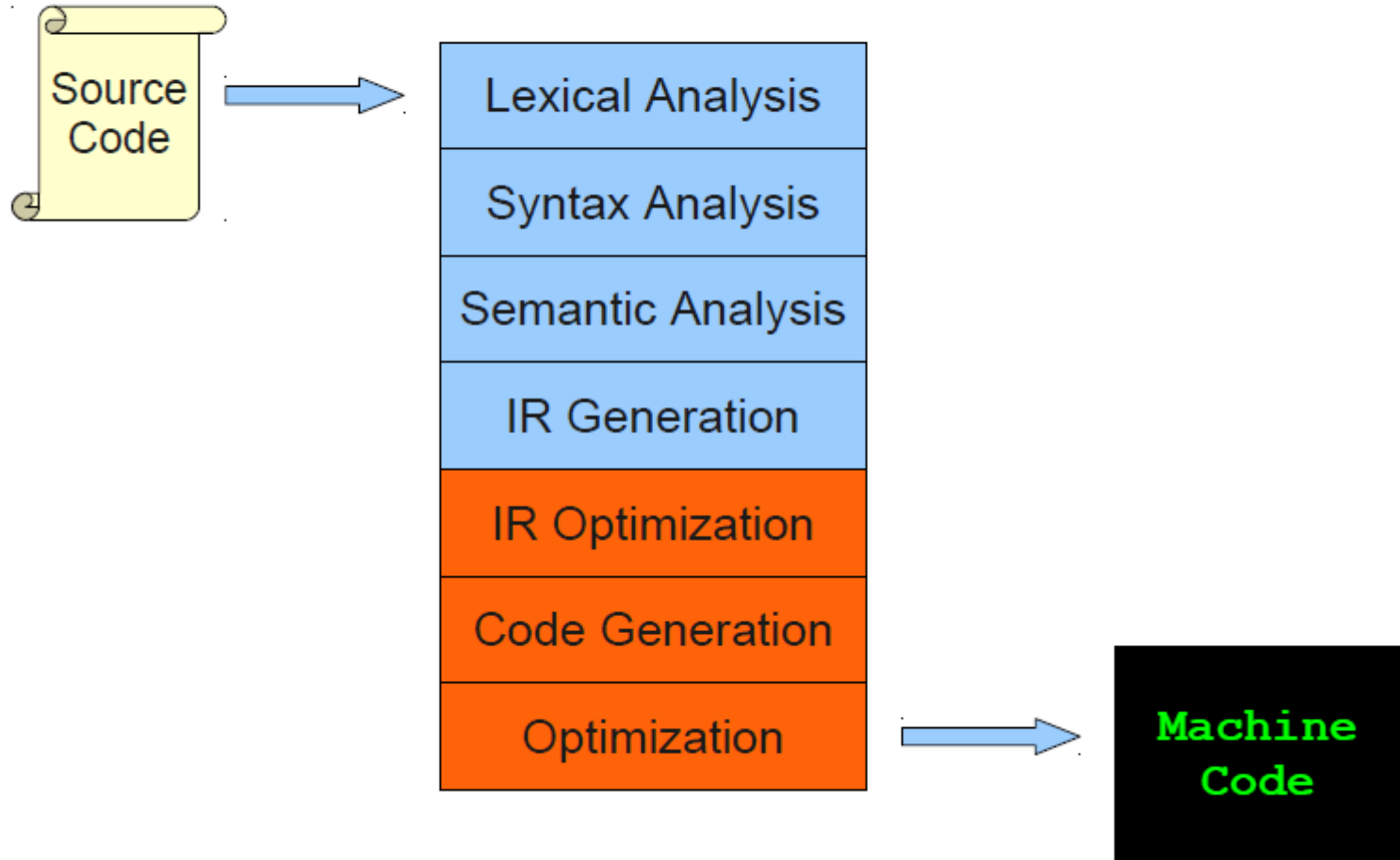
- Introduction
- The Structure of a Modern Compiler
- Lexical Analysis
- The Parsing Problem

The big picture

Compilation Process



The Structure of a Modern Compiler



```

while (y < z) {
    int x = a + b;
    y += x;
}

```

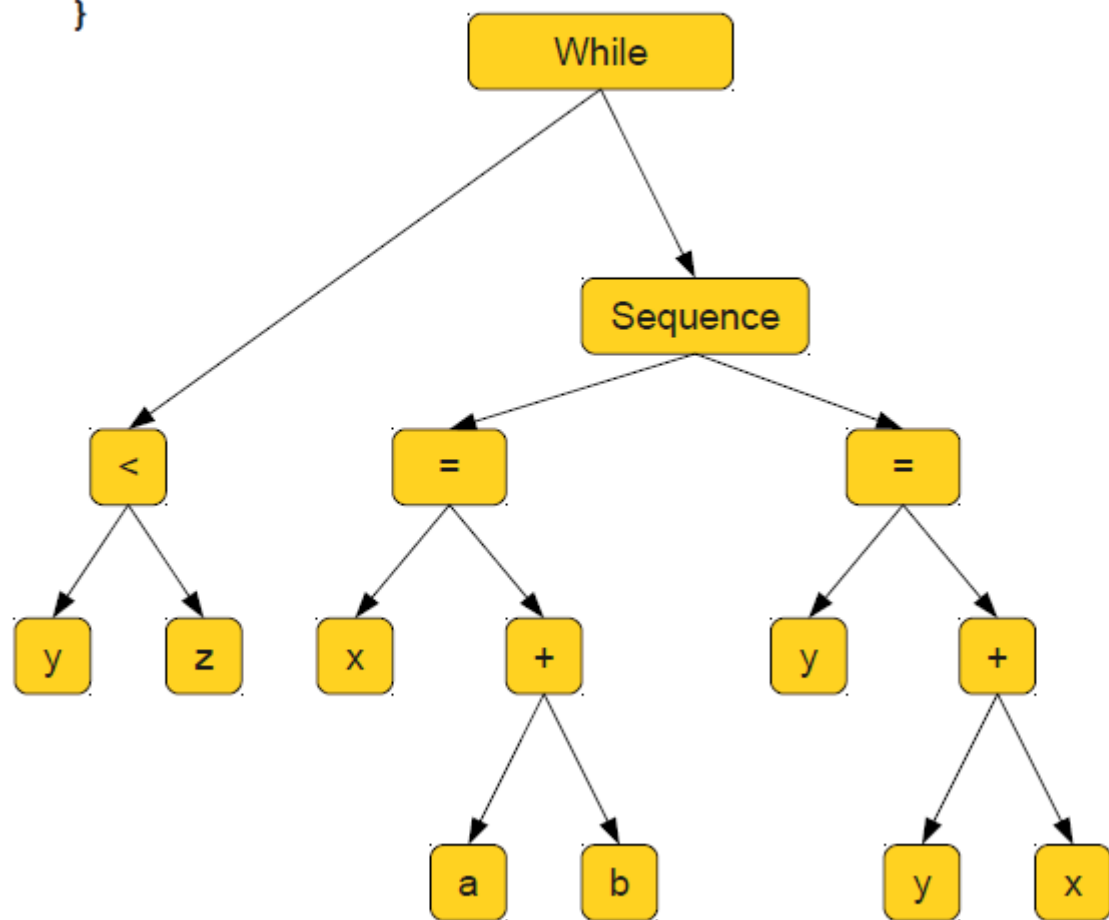
```

T_While
T_LeftParen
T_Identifier y
T_Less
T_Identifier z
T_RightParen
T_OpenBrace
T_Int
T_Identifier x
T_Assign
T_Identifier a
T_Plus
T_Identifier b
T_Semicolon
T_Identifier y
T_PlusAssign
T_Identifier x
T_Semicolon
T_CloseBrace

```

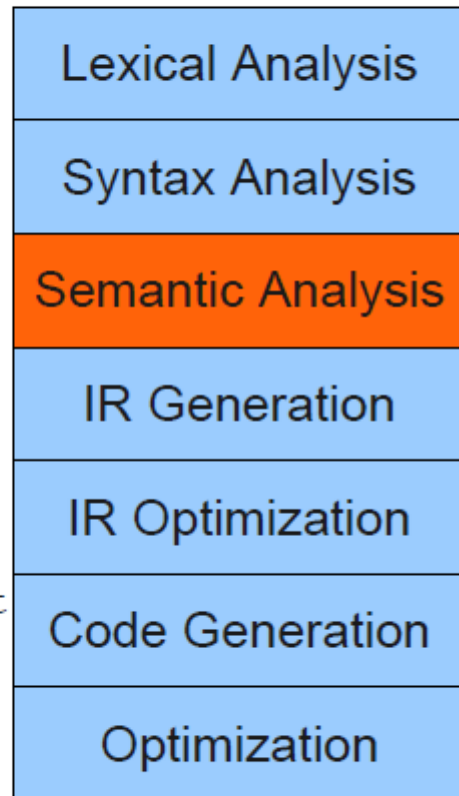
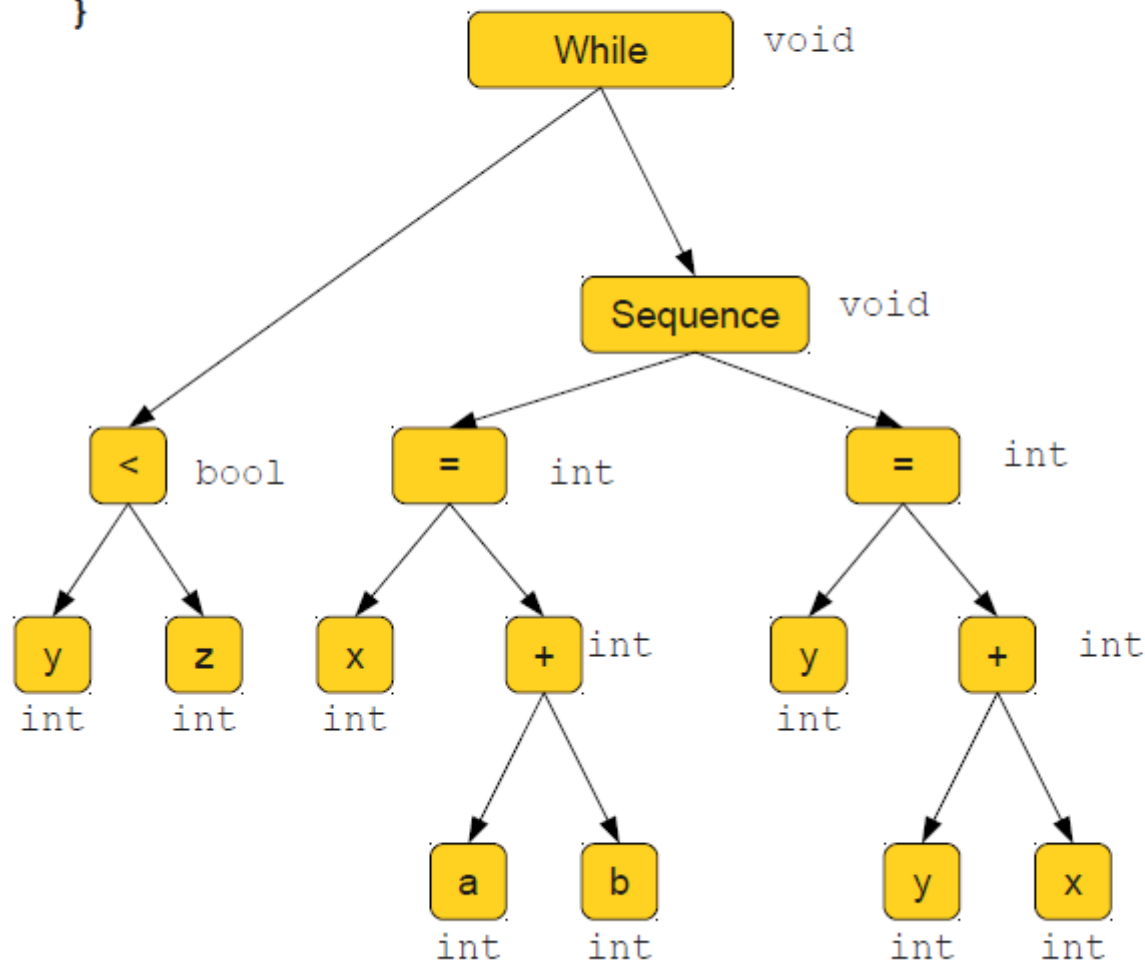
Lexical Analysis
Syntax Analysis
Semantic Analysis
IR Generation
IR Optimization
Code Generation
Optimization

```
while (y < z) {  
    int x = a + b;  
    y += x;  
}
```



Lexical Analysis
Syntax Analysis
Semantic Analysis
IR Generation
IR Optimization
Code Generation
Optimization

```
while (y < z) {
    int x = a + b;
    y += x;
}
```



```
while (y < z) {  
    int x = a + b;  
    y += x;  
}
```

```
Loop: x    = a + b  
      y    = x + y  
      _t1  = y < z  
      if _t1 goto Loop
```

Lexical Analysis
Syntax Analysis
Semantic Analysis
IR Generation
IR Optimization
Code Generation
Optimization


```
while (y < z) {  
    int x = a + b;  
    y += x;  
}
```

```
        x      = a + b  
Loop:   y      = x + y  
        _t1    = y < z  
        if _t1 goto Loop
```

Lexical Analysis

Syntax Analysis

Semantic Analysis

IR Generation

IR Optimization

Code Generation

Optimization

```
while (y < z) {  
    int x = a + b;  
    y += x;  
}
```

```
      add $1, $2, $3  
Loop: add $4, $1, $4  
      slt $6, $4, $5  
      beq $6, loop
```

Lexical Analysis

Syntax Analysis

Semantic Analysis

IR Generation

IR Optimization

Code Generation

Optimization

```
while (y < z) {  
    int x = a + b;  
    y += x;  
}
```

```
          add $1, $2, $3  
Loop:    add $4, $1, $4  
          blt $4, $5, loop
```

Lexical Analysis

Syntax Analysis

Semantic Analysis

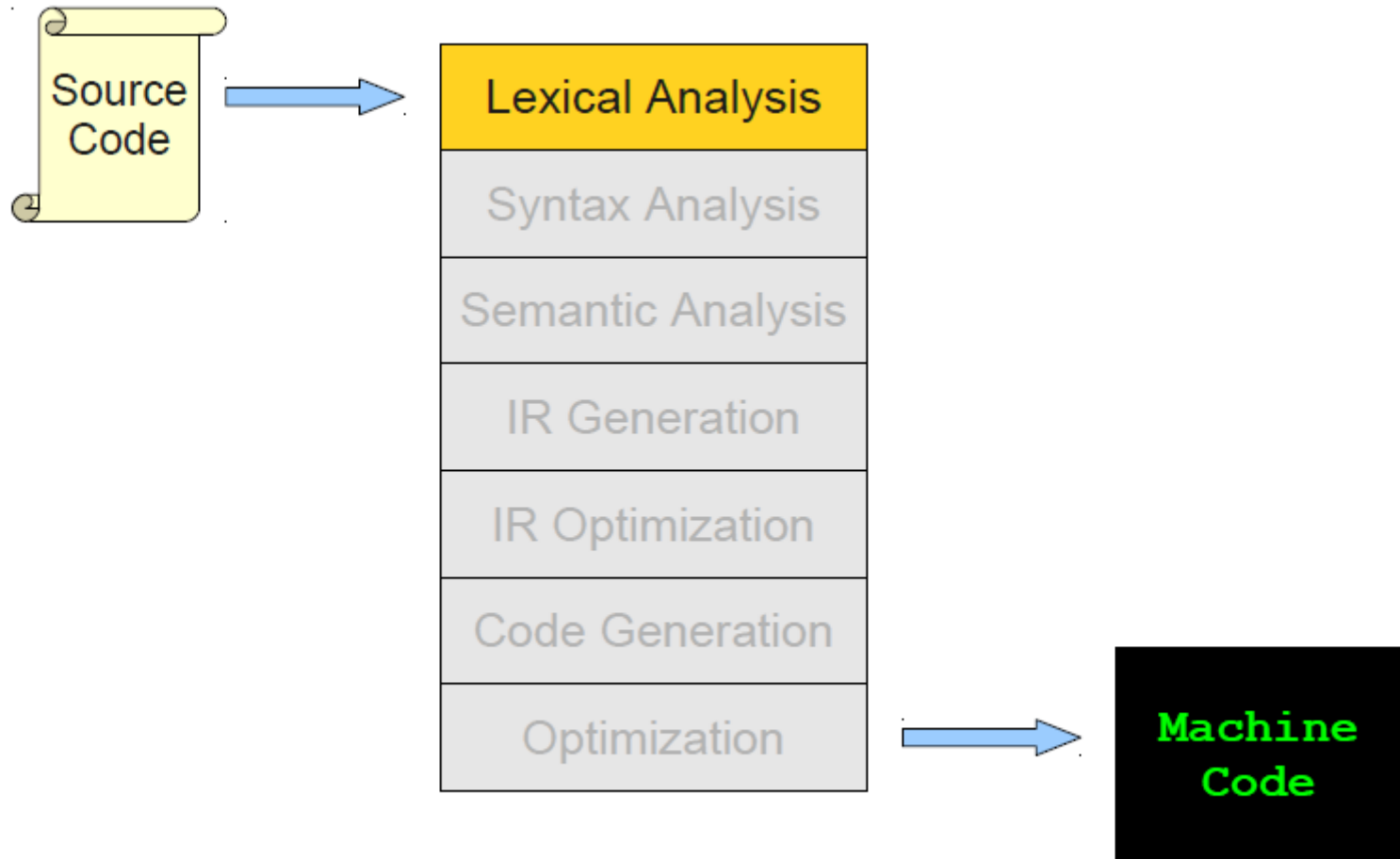
IR Generation

IR Optimization

Code Generation

Optimization

Next →



Definitions

- A *lexeme* is the lowest level syntactic unit of a language (e.g. the, boy)
- A *token* is a category of lexemes (e.g., identifier)

```
result = oldsum - value / 100;
```

Following are the tokens and lexemes of this statement:

<i>Token</i>	<i>Lexeme</i>
IDENT	result
ASSIGN_OP	=
IDENT	oldsum
SUB_OP	-
IDENT	value
DIV_OP	/
INT_LIT	100
SEMICOLON	;

Introduction

For all practical purposes, we can think of the source code as a very long string of characters from some alphabet.

Lexical analyzers collect characters into logical groupings (lexemes) and assign internal codes (tokens) to the grouping.

As an analogy, we can think of the lexical analyzer as a machine that chops up source code into the individual “words”.

What lexical analyzers do

- The “front-end” for the parser
- A **pattern matcher** for character strings
- Identifies substrings of the source program that belong together => lexemes

– Example: sum = B -5;

<u>Lexeme</u>	<u>Token</u>
sum	ID (identifier)
=	ASSIGN_OP
B	ID
-	SUBTRACT_OP
5	INT_LIT (integer literal)
;	SEMICOLON

What lexical analyzers do/2

- **Functions:**

- ⇒ **Extract lexemes** from a given input string and produce the corresponding tokens, while skipping comments and blanks
- ⇒ **Insert lexemes** for user-defined names into symbol table, which is used by later phases of the compiler
- ⇒ **Detect syntactic errors in tokens** and report such errors to user

How to build a lexical analyzer?

- Three ways (p175)
- We will take the second approach:
- Formal regular grammar => software package that generates the lexical analyzer.
- **State transition diagrams that describe token types => write a program to implement the diagram**
- State transition diagrams that describe token types => table driven approach to implement the diagram.

State Transition Diagram: Example

$\langle id \rangle \rightarrow \langle letter \rangle \{ \langle letter \rangle | \langle int \rangle \}$ **Tokens in CFG**

$\langle int \rangle \rightarrow \langle digit \rangle \{ \langle digit \rangle \}$

$\langle letter \rangle \rightarrow A | B | C | \dots | Z | a | b | \dots | z$

$\langle digit \rangle \rightarrow 0 | 1 | 2 | \dots | 9$

$\langle special \rangle \rightarrow * | / | + | - | (|)$

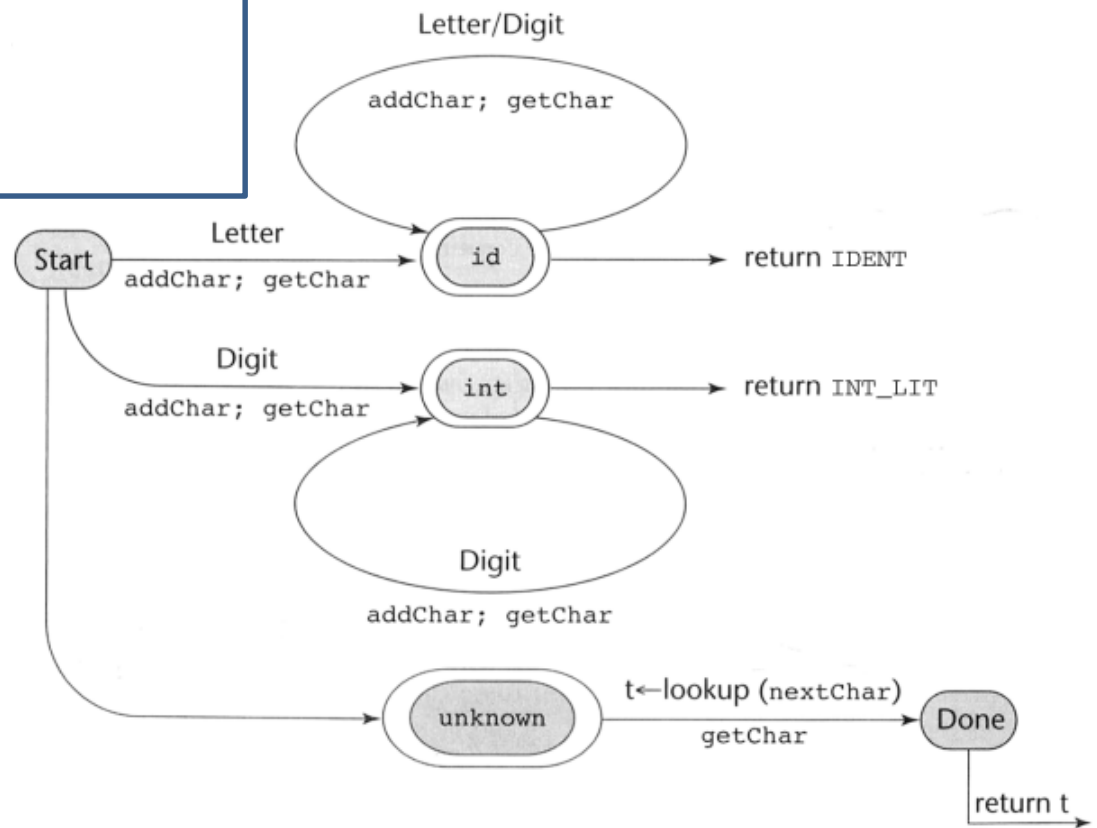
$Id \rightarrow letter(letter|digit)^*$

$int \rightarrow digit\ digit^*$

$digit \rightarrow [0-9]$

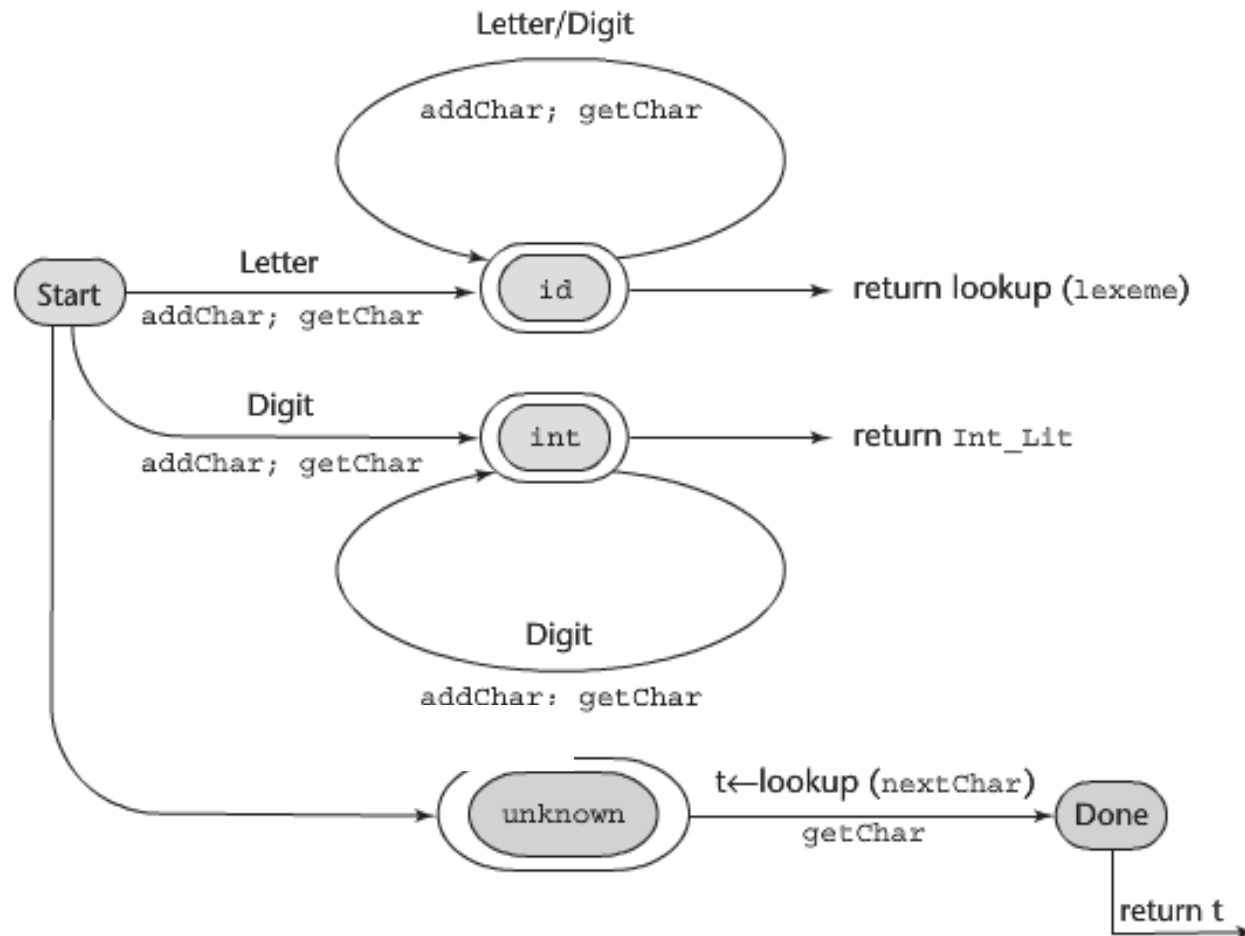
$letter \rightarrow [a-zA-Z]$ **Tokens in**

$special \rightarrow [* / + - ()]$ **Regular grammar**



Suppose we need a lexer that recognizes only names, integer literals, parentheses, and arithmetic operators.

State Transition Diagram: Example



Suppose we need a lexer that recognizes only names, reserved words, integer literals, parentheses, and arithmetic operators.

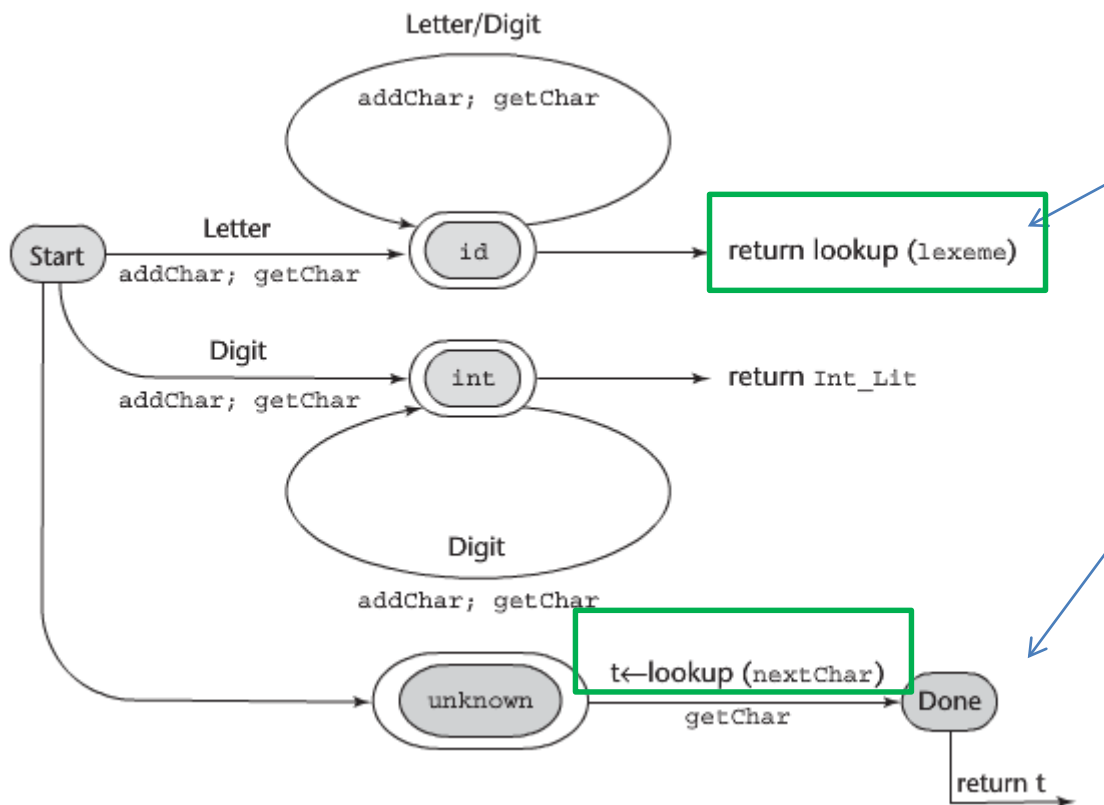
Example/3

Need to distinguish reserved words from identifiers

e.g., reserved words: main and int

identifiers: sum and B

Use a table lookup to determine whether a possible identifier is in fact a reserved word



To determine whether name is a reserved word (not in the book Implementation.)

To determine whether a character is a parenthesis, or arithmetic operator

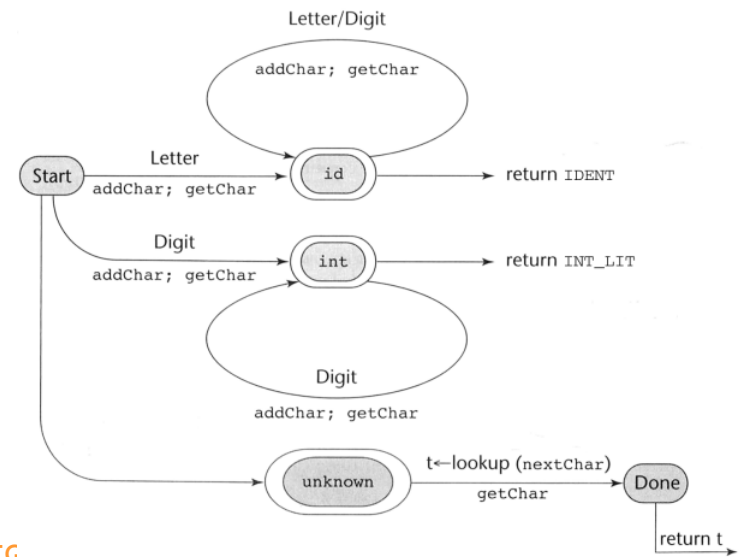
Example :Implementation/3

```
int lex() {
    lexLen = 0;
    getNonBlank();
    switch (charClass) {
    case LETTER:
        addChar();
        getChar();
        while (charClass == LETTER || charClass == DIGIT) {
            addChar();
            getChar();
        }
        nextToken = IDENT;
        break;
    case DIGIT:
        addChar();
        getChar();
        while (charClass == DIGIT) {
            addChar();
            getChar();
        }
        nextToken = INT_LIT;
        break;
    case UNKNOWN:
        lookup(nextChar);
        getChar();
        break
    }
```

* This does not have logic for checking reserved words

```
case EOF:
    nextToken = EOF;
    lexeme[0] = 'E';
    lexeme[1] = 'O';
    lexeme[2] = 'F';
    lexeme[3] = '\0';
    break;
}
```

```
printf("Next token is: %d, Next lexeme is %s\n",
       nextToken, lexeme);
return nextToken;
```



```

void getChar()
{
    nextChar = getc(in_fp);
    // printf("next char %c\t", nextChar);
    if (nextChar != EOF)
    {
        if (isalpha(nextChar))
            charClass = LETTER;

        else if (isdigit(nextChar))
            charClass = NUMBER;

        else
            charClass = UNKNOWN;
    }
    else
        charClass = EOF;
    //printf("char class %d\t", charClass);
}

```

```

int lookup(char ch)
{
    switch (ch)
    {
        case '(':
            addChar();
            nextToken = LEFT_PAREN;
            break;

        case ')':
            addChar();
            nextToken = RIGHT_PAREN;
            break;

        case '+':
            addChar();
            nextToken = ADD_OP;
            break;

        case '-':
            addChar();
            nextToken = SUB_OP;
            break;

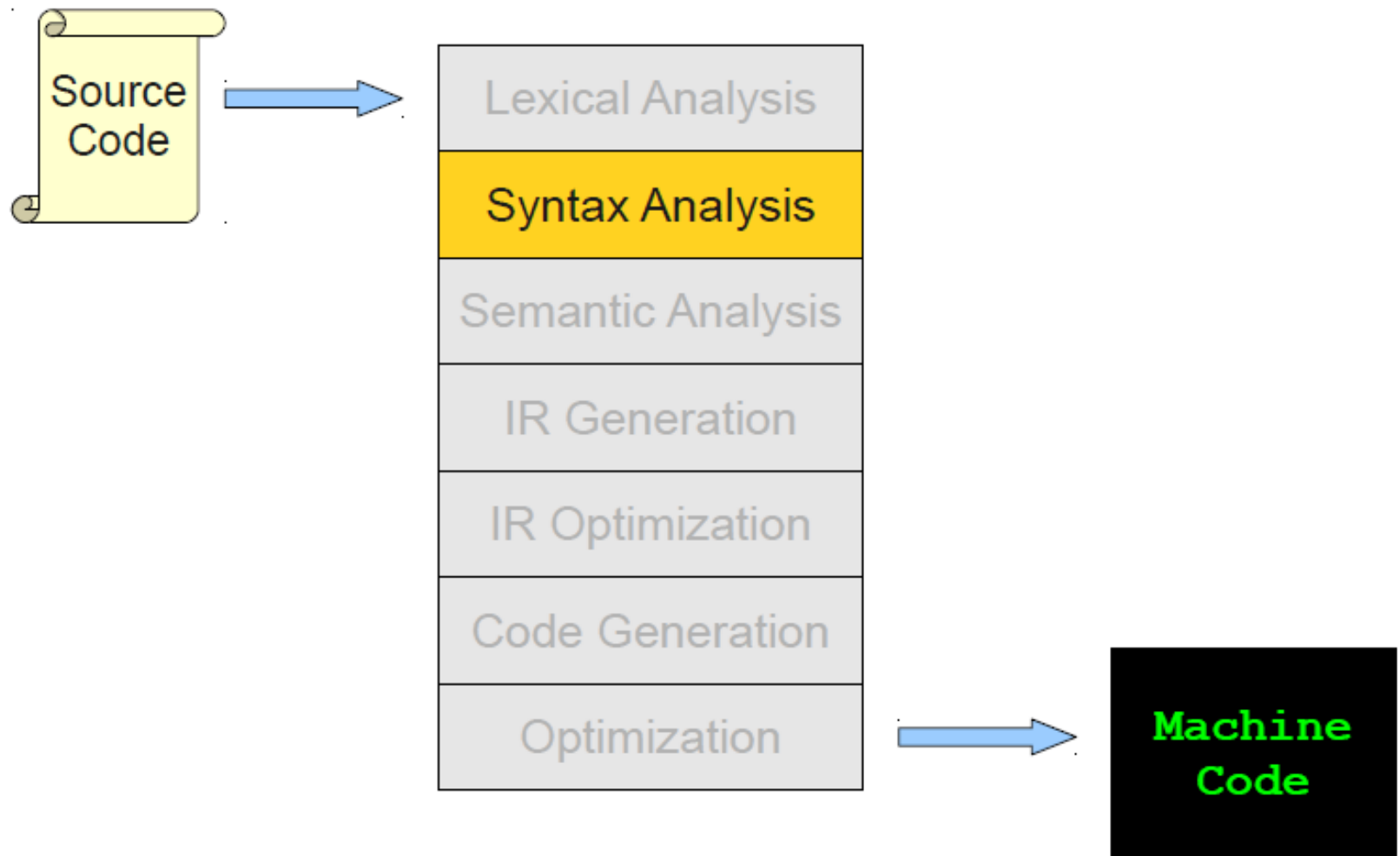
        case '*':
            addChar();
            nextToken = MULT_OP;
            break;

        case '/':
            addChar();
            nextToken = DIV_OP;
            break;

        default:
            addChar();
            nextToken = EOF;
    }
}

```

Where We Are



Introduction

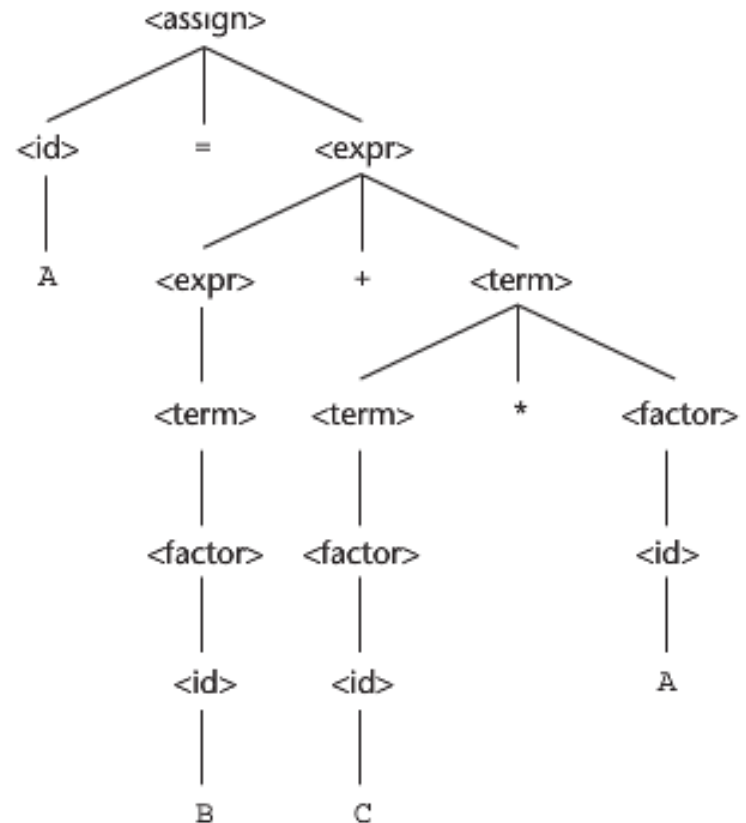
- The syntax analyzer receives the source code in the form of tokens from the lexical analyzer and performs syntax analysis, which determines **the structure** of the program.
- Typically, the result generated by the syntax analyzer is a parse tree.

An illustration

Figure 3.3

The unique parse tree for $A = B + C * A$ using an unambiguous grammar

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$
 $\langle \text{id} \rangle \rightarrow A \mid B \mid C$
 $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle$
 $\mid \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle$
 $\mid \langle \text{factor} \rangle$
 $\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle)$
 $\mid \langle \text{id} \rangle$



Lexer vs. Parser

<i>Phase</i>	<i>Input</i>	<i>Output</i>
Lexer	String of characters	String of tokens
Parser	String of tokens	Parse tree

What do parsers do?

- Goals of a parser:
 - Find all syntax errors
 - Produce parse trees for input program
- Two categories of parsers:
 - **Top down**
 - produces the parse tree, beginning at the root
 - **Bottom up**
 - produces the parse tree, beginning at the leaves

Top-down parsers

- Launch the process with the start symbol and apply the productions until you arrive at the desired string.
 - Builds a parse tree in **preorder**
 - Corresponds **leftmost derivation**
 - **LL parser**

$$\begin{array}{lcl} S & \rightarrow & AB \\ A & \rightarrow & aA \mid \epsilon \\ B & \rightarrow & b \mid bB \end{array}$$
$$S \rightarrow aAc$$
$$A \rightarrow aA \mid b$$
$$S \Rightarrow aAc \Rightarrow aaAc \Rightarrow aabc$$

With this grammar

Here is a top-down parse of `aaab`

S	
\Rightarrow AB	$S \rightarrow AB$
\Rightarrow aAB	$A \rightarrow aA$
\Rightarrow aaAB	$A \rightarrow aA$
\Rightarrow aaaAB	$A \rightarrow aA$
\Rightarrow aaa ϵ B	$A \rightarrow \epsilon$
\Rightarrow aaab	$B \rightarrow b$

Top-down Parsers Again

From deviation aspect

- 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
 - x : all terminals
 - A : the leftmost nonterminal in the sentential form
 - α : combination of terminal and non-terminal

(Predicative) Top-down Parsers

From deviation aspect

- A nonterminal has more than one RHS
- 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

$$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ ('+' \mid '-') \langle \text{term} \rangle \}$$
$$A \rightarrow bB \mid cBb \mid a$$

(Predicative)Recursive Descent Parser

- A top-down parser implementation
 - Based on directly on BNF/EBNF
- Consists of a collection of subprograms
 - A recursive descent parser has a subprogram for each non-terminal symbol
 - The responsibility of the subprogram: parse sentences that can be generated by that nonterminal

Example

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ ('+' | '-') \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ ('*' | '/') \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle \rightarrow \text{id} | '(' \langle \text{expr} \rangle ')'$

```
void expr()  
{...}
```

```
void term()  
{...}
```

```
void factor()  
{...}
```

```
void error()  
{...}
```

-- "id" here is not a terminal, but it stands
for the identifier token

$\langle \text{factor} \rangle \rightarrow \langle \text{id} \rangle | '(' \langle \text{expr} \rangle ')'$

-- language literals are enclosed with ' '

- For each **terminal** symbol/token 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

Example/2

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ ('+' | '-') \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ ('*' | '/') \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle \rightarrow \text{id} \mid '(\langle \text{expr} \rangle)'$

```
void expr() {  
    printf("Enter <expr>\n");  
    term();  
    while (nextToken == PLUS_CODE  
        || nextToken == MINUS_CODE)  
    {  
        lex();  
        term();  
    }  
    printf("Exit <expr>\n");  
}
```

- ❖ lex() is the lexical analyzer function. It gets the next lexeme and puts its token code in the global variable nextToken
- ❖ All subprograms are written with the convention that each one leaves the next token of input in nextToken

Example/3

```
void term() {  
    printf("Enter <term>\n");  
    /* parse the first factor */  
    factor();  
    while (nextToken == MUL_CODE ||  
           nextToken == DIV_CODE)  
    {  
        lex();  
        factor();  
    }  
    printf("Exit <term>\n");  
}
```

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ ('+' | '-') \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ ('*' | '/') \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle \rightarrow \text{id} | '(' \langle \text{expr} \rangle ')'$

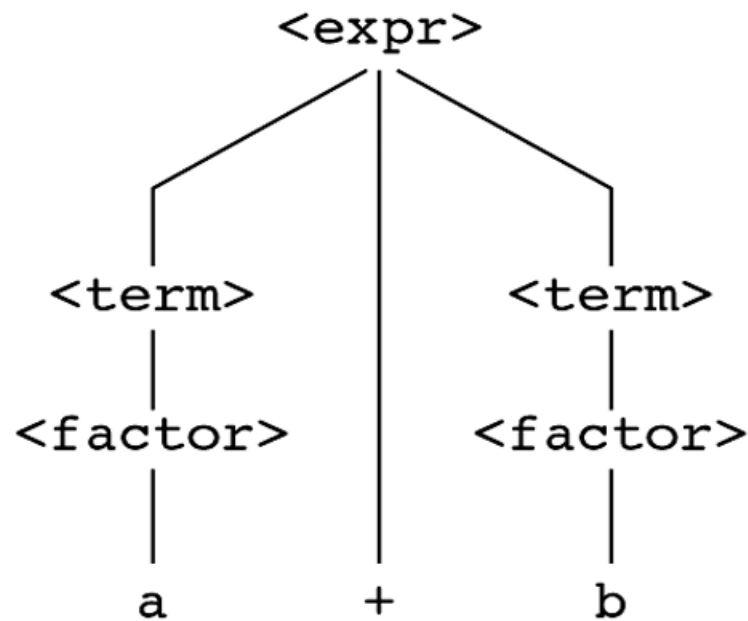
$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ ('+' | '-') \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ ('*' | '/') \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle \rightarrow \text{id} | '(\langle \text{expr} \rangle)'$

Example/4

```
void factor() {
    printf("Enter <factor>\n");
    /* Determine which RHS */
    if (nextToken == ID_CODE)
        lex();
    else if (nextToken == LEFT_PAREN_CODE) {
        lex();
        expr();
        if (nextToken == RIGHT_PAREN_CODE)
            lex();
        else
            error();
    }
    else {
        error(); /* Neither RHS matches */
    }
    printf("Exit <factor>\n");
}
```

Figure 4.2

Parse tree for $a + b$



Calling sequences:

Call lex //a
Enter $\langle \text{expr} \rangle$
Enter $\langle \text{term} \rangle$
Enter $\langle \text{factor} \rangle$
Call lex //+
Exit $\langle \text{factor} \rangle$
Exit $\langle \text{term} \rangle$
Call lex //b
Enter $\langle \text{term} \rangle$
Enter $\langle \text{factor} \rangle$
Call lex // end of input
Exit $\langle \text{factor} \rangle$
Exit $\langle \text{term} \rangle$
Exit $\langle \text{expr} \rangle$

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ ('+' | '-') \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ ('*' | '/') \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle \rightarrow \text{id} | '(' \langle \text{expr} \rangle ')'$

Calling sequences:

Next Token is : 25 Next lexeme is (

Enter <expr>

Enter <term>

Enter <factor>

Next Token is : 11 Next lexeme is sum

Enter <expr>

Enter <term>

Enter <factor>

Next Token is : 21 Next lexeme is +

Exit <factor>

Exit <term>

Next Token is : 10 Next lexeme is 47

Enter <term>

Enter <factor>

Next Token is : 26 Next lexeme is)

Exit <factor>

Exit <term>

Exit <expr>

Next Token is : 24 Next lexeme is /

Exit <factor>

Next Token is : 11 Next lexeme is total

Enter <factor>

Next Token is : -1 Next lexeme is EOF

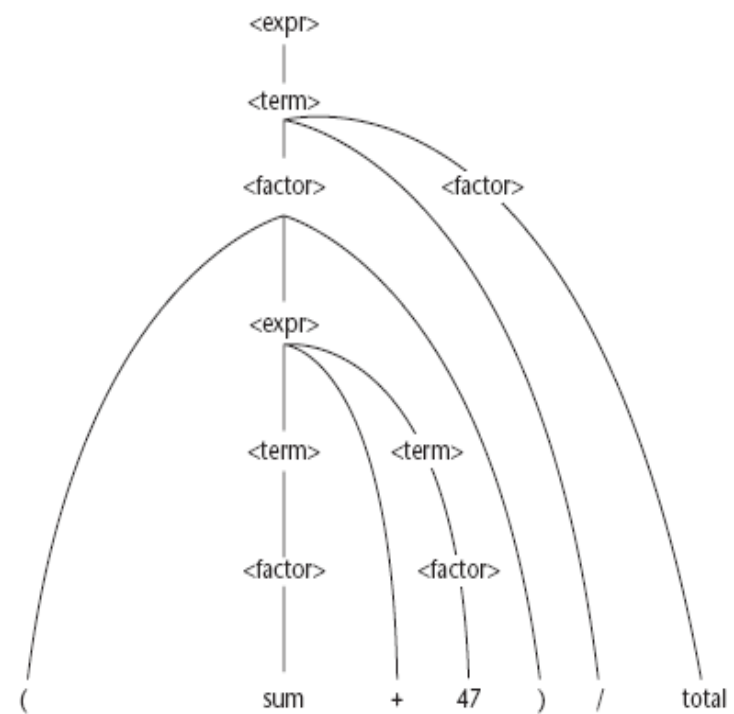
Exit <factor>

Exit <term>

Exit <expr>

Figure 4.2

Parse tree for (sum + 47) / total



$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ ('+' | '-') \langle \text{term} \rangle \}$

$\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ ('*' | '/') \langle \text{factor} \rangle \}$

$\langle \text{factor} \rangle \rightarrow \text{id} \mid '(\langle \text{expr} \rangle)'$

Recursive Descent Parser: drawback

- Problem with **left recursion**
 - $\langle A \rangle \rightarrow \langle A \rangle + \langle B \rangle$ (direct left recursion)
 - $\langle A \rangle \rightarrow \langle B \rangle c \langle D \rangle$ (indirect left recursion)
 $\langle B \rangle \rightarrow \langle A \rangle b$
 - A grammar can be modified to remove left recursion (p181)
- Inability to **determine the correct RHS** on the basis of one token of lookahead
 - **Yes** Example: $A \rightarrow aB \mid bAb \mid Bb$
 $B \rightarrow cB \mid d$
 - **No** Example: $A \rightarrow aC \mid Bd$
 $B \rightarrow ac$
 $C \rightarrow c$
 - Pairwise disjoint test and left factoring
 - EBNF does not have this problem

Bottom-up parsers

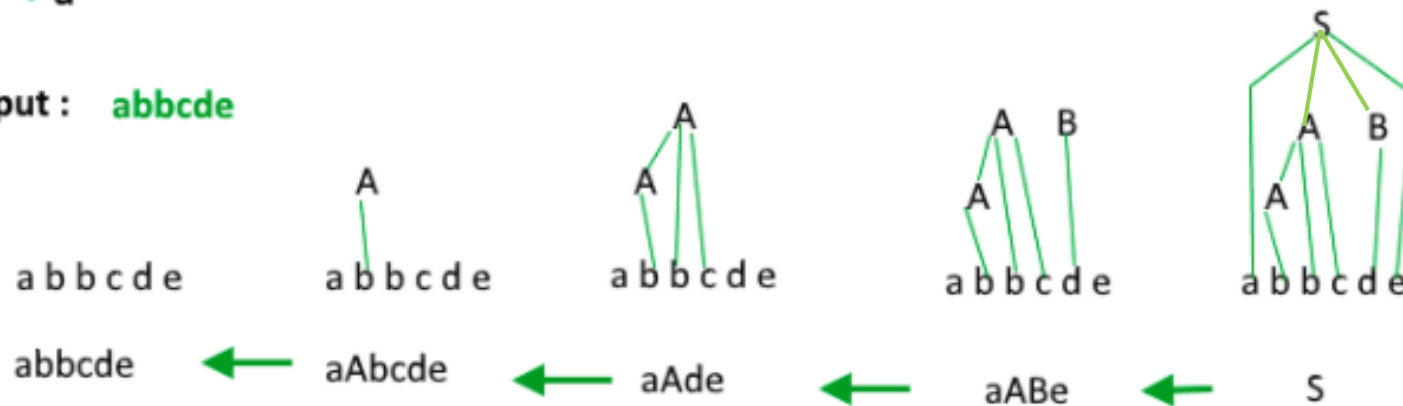
- Begin at the leaves and progressing toward the root.
 - Corresponds to the reverse of **rightmost derivation**
 - are more powerful than top-down methods
 - **LR**

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

Input : **abbcd e**



Bottom-up parsers Again

from the derivation aspect

- 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 rightmost derivation
 - The most common bottom-up parsing algorithms are in the LR family

Every string of symbols (terminal and nonterminal) in a derivation is a *sentential form*
A right sentential form is a sentential form appears in a rightmost derivation.

Shift-Reduce Algorithms

- Implementation of bottom-up parsing
 - Uses stack (parse stack)
- Reduce is the action of replacing substring (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

An shift-reduce example

from the book

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

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

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

$$E \Rightarrow E + T$$

$$\Rightarrow E + \underline{T} * F$$

$$\Rightarrow E + T * \underline{\text{id}}$$

$$\Rightarrow E + \underline{F} * \text{id}$$

$$\Rightarrow E + \underline{\text{id}} * \text{id}$$

$$\Rightarrow \underline{T} + \text{id} * \text{id}$$

$$\Rightarrow \underline{F} + \text{id} * \text{id}$$

$$\Rightarrow \underline{\text{id}} + \text{id} * \text{id}$$

Corresponds to the
reverse of
rightmost derivation

An shift-reduce example

from the book

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow (E) \mid id$$

. id + id * id (shift)
⇒ Id . + id * id (reduce)
⇒ F. + id * id (reduce)
⇒ T. + id * id (reduce)
⇒ E. + id * id (shift)
⇒ E + . Id * id (shift)
⇒ E + id. * id (reduce)
⇒ E + F. * id (reduce)
⇒ E + T. * id (pick shift over reduce)
⇒ E + T *. Id (shift)
⇒ E + T * F (reduce)
⇒ E + T
⇒ E

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
- A recursive-descent parser is an LL parser
- Parsing problem for bottom-up parsers: find the substring of current sentential form that can be reduced to some left-hand side.