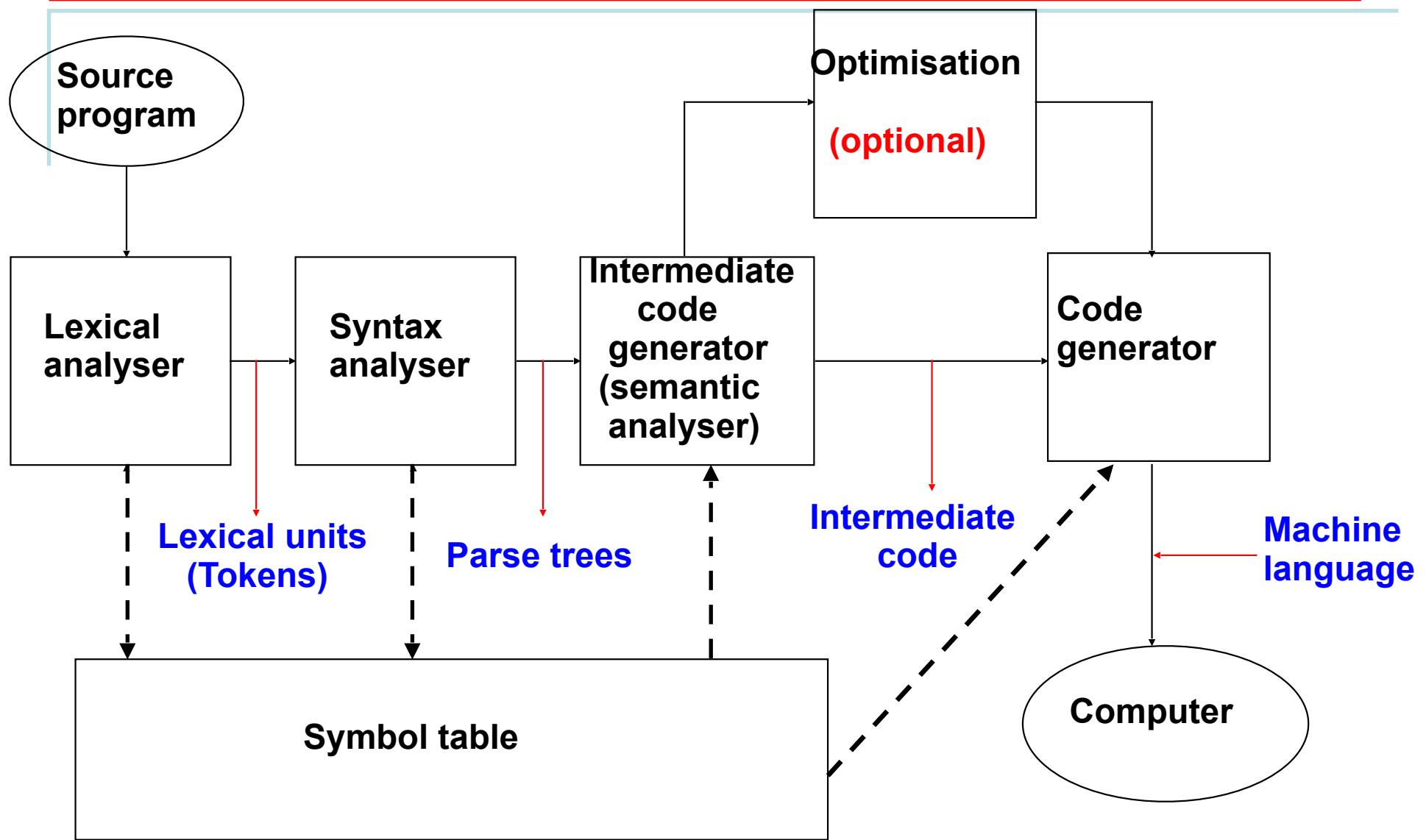


CC410: System Programming

Dr. Manal Helal – Spring 2015 – Lecture 14–Compilers 1

Chapter 5 - Compilers



Terminology

- » Statement
 - Declaration, assignment containing expression
- » Grammar
 - A set of rules specify the form of legal statements
- » Syntax vs. Semantics
 - Example: assuming I, J, K: integer and X,Y: float
 - I:=J+K vs. X:= Y+I
- » Compilation
 - Matching statements written by the programmer to structures defined by the grammar and generating the appropriate object code.

Basic Compiler

- » Lexical analysis - scanner
 - Scanning the source statement, recognising and classifying the various tokens
- » Syntactic analysis - parser
 - Recognising the statement as some language construct.
 - Construct a parser tree (syntax tree)
- » Code generation – code generator
 - Generate assembly language codes that require an assembler
 - Generate machine codes (Object codes) directly

Scanner

```
SUM:=0;  
SUMSQ:= 0;
```

```
PROGRAM  
STATS  
VAR  
SUM,SUMSQ,I
```

FIGURE 5.1 Example of a Pascal program.

```
1 PROGRAM STATS  
2 VAR  
3     SUM, SUMSQ, I, VALUE, MEAN, VARIANCE : INTEGER  
4 BEGIN  
5     SUM := 0;  
6     SUMSQ := 0;  
7     FOR I := 1 TO 100 DO  
8         BEGIN  
9             READ(VALUE);  
10            SUM := SUM + VALUE;  
11            SUMSQ := SUMSQ + VALUE * VALUE  
12        END;  
13        MEAN := SUM DIV 100;  
14        VARIANCE := SUMSQ DIV 100 - MEAN * MEAN;  
15        WRITE(MEAN, VARIANCE)  
16    END.
```

```
READ (VALUE);
```

Grammar

- » Grammar: a set of rules
 - Backus-Naur Form (BNF)
 - Ex: Figure 5.2

```
<read> ::= READ ( <id-list> )
<id-list>    ::= id | <id-list> , id
```

- » Terminology
 - Define symbol **::=**
 - Nonterminal symbols **<>**
 - Alternative symbols **|**
 - Terminal symbols

Simplified Pascal Grammar

```
1 <prog>           ::= PROGRAM <prog-name> VAR <dec-list> BEGIN <stmt-list> END.
2 <prog-name>       ::= id
3 <dec-list>        ::= <dec> | <dec-list> ; <dec>
4 <dec>             ::= <id-list> : <type>
5 <type>            ::= INTEGER
6 <id-list>         ::= id | <id-list> , id
7 <stmt-list>       ::= <stmt> | <stmt-list> ; <stmt>
8 <stmt>             ::= <assign> | <read> | <write> | <for>
9 <assign>          ::= id := <exp>
10 <exp>             ::= <term> | <exp> + <term> | <exp> - <term>
11 <term>            ::= <factor> | <term> * <factor> | <term> DIV <factor>
12 <factor>          ::= id | int | ( <exp> )
13 <read>            ::= READ ( <id-list> )
14 <write>           ::= WRITE ( <id-list> )
15 <for>             ::= FOR <index-exp> DO <body>
16 <index-exp>       ::= id := <exp> TO <exp>
17 <body>            ::= <stmt> | BEGIN <stmt-list> END
```

Figure 5.2 Simplified Pascal grammar.

Parser

» **READ(VALUE)**

» **SUM := 0**

» **SUM := SUM + VALUE**

» **MEAN := SUM DIV 100**

» **<read> ::= READ (<id-list>)**

» **<id-list> ::= id | <id-list>, id**

» **<assign> ::= id := <exp>**

» **<exp> ::= <term> |
<exp>+<term> |
<exp>-<term>**

» **<term> ::= <factor> |
<term>*<factor> | <term> DIV
<factor>**

» **<factor> ::= id | int | (<exp>)**

Syntax Tree

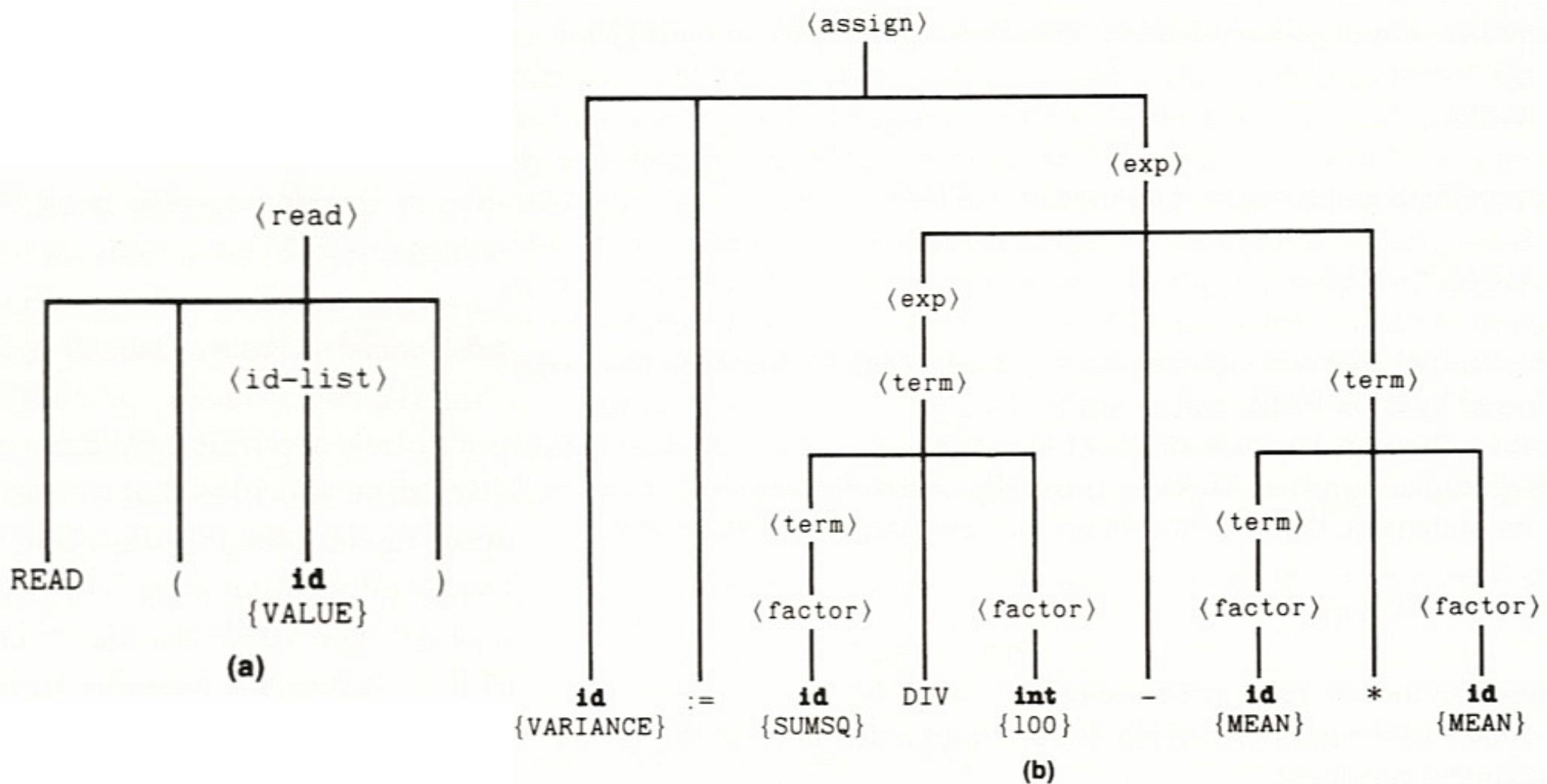
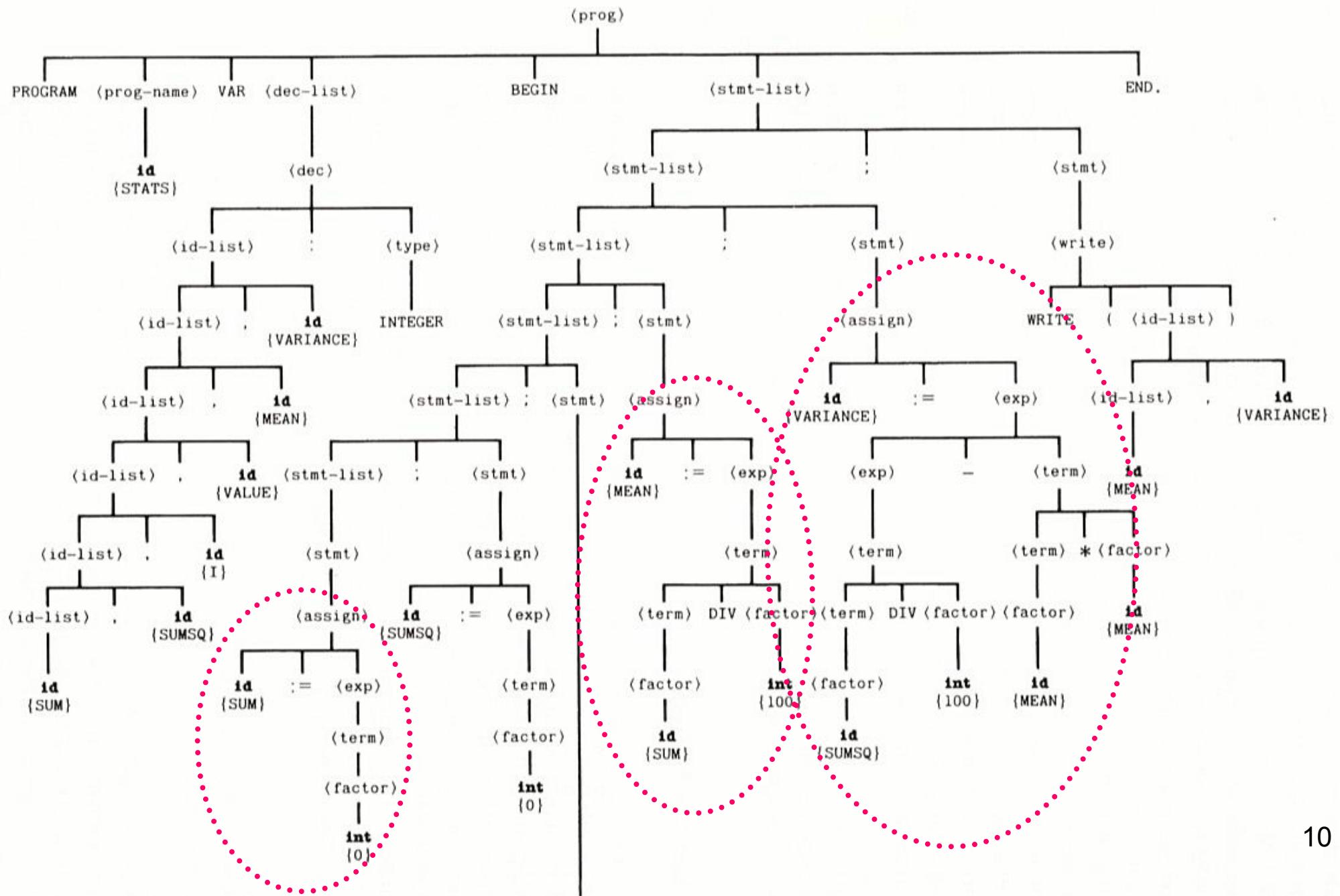


Figure 5.3 Parse trees for two statements from Fig. 5.1.

Syntax Tree for Program 5.1



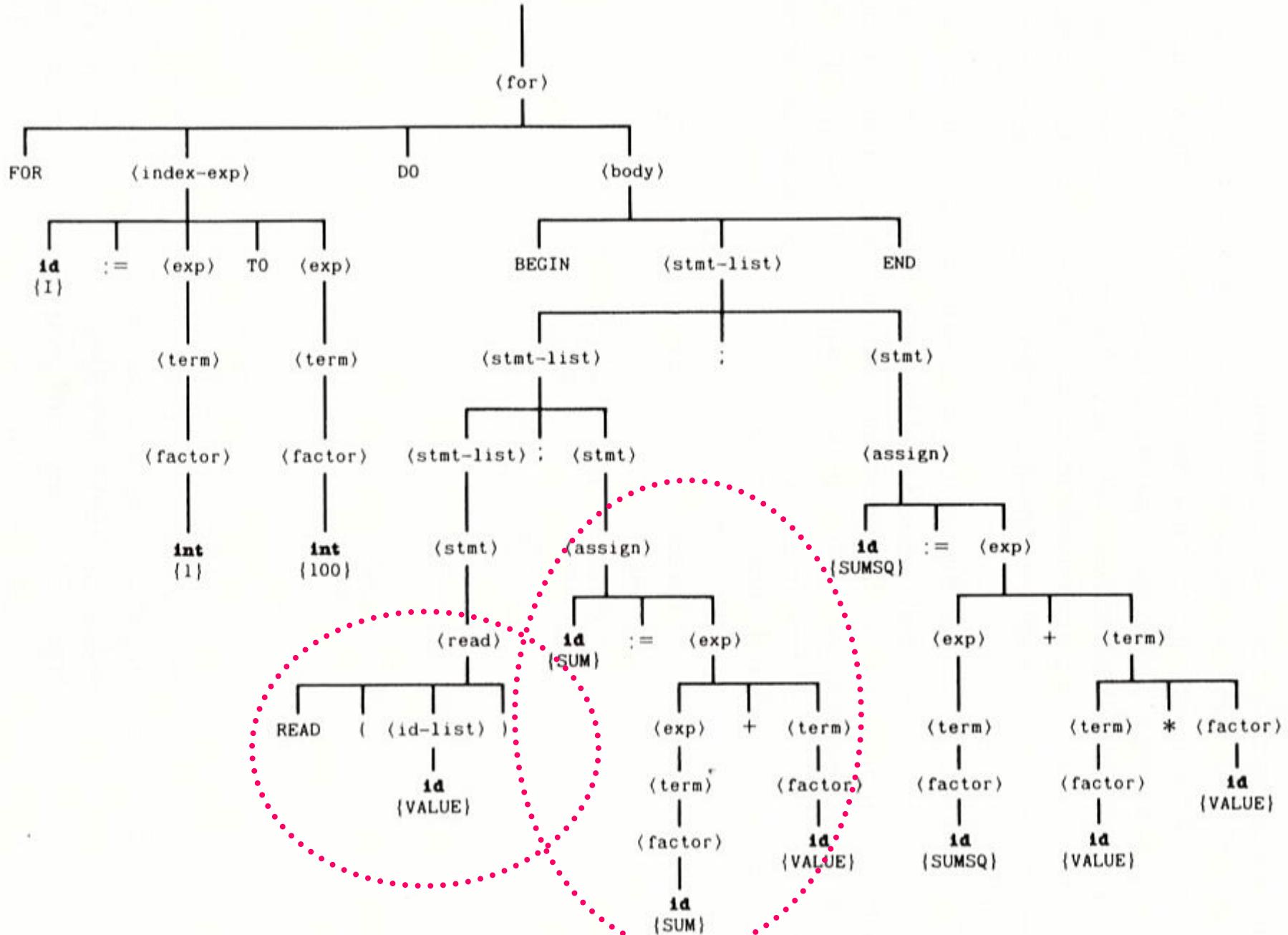


Figure 5.4 Parse tree for the program from Fig. 5.1.

Lexical Analysis

» Function

- Scanning the program to be compiled and recognising the tokens that make up the source statements.

```
<ident>   ::=  <letter> | <ident> <letter> | <ident> <digit>
<letter>   ::=  A | B | C | D | ... | Z
<digit>    ::=  0 | 1 | 2 | 3 | ... | 9
```

» Tokens

- Tokens can be **keywords**, **operators**, **identifiers**, **integers**, **floating-point numbers**, **character strings**, etc.
- Each token is usually represented by some fixed-length code, such as an **integer**, rather than as a variable-length character string (see Figure 5.5)
- Token type, Token specifier (value) (see Figure 5.6)

Scanner Output

- » Token specifier
 - Identifier name, integer value, (type)
- » Token coding scheme
 - Figure 5.5

Token	Code
PROGRAM	1
VAR	2
BEGIN	3
END	4
END.	5
INTEGER	6
FOR	7
READ	8
WRITE	9
TO	10
DO	11
;	12
:	13
,	14
:=	15
+	16
-	17
*	18
DIV	19
(20
)	21
id	22
int	23

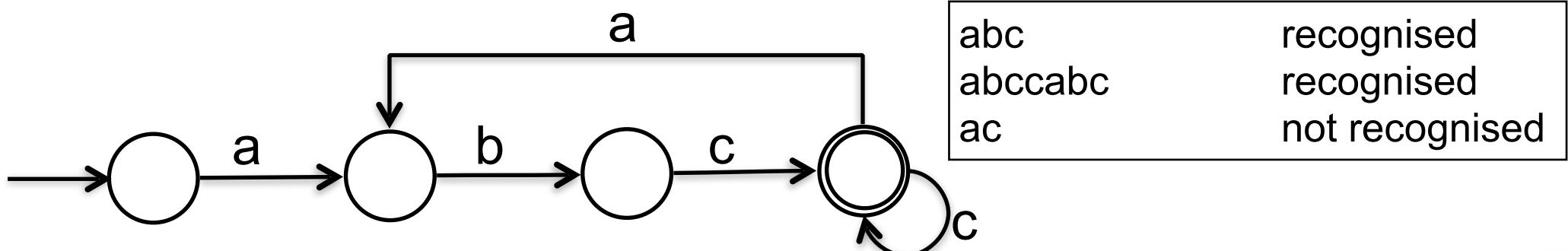
Line	Token type	Token specifier	Line	Token type	Token specifier
1	1		10	22	~SUM
	22	~STATS		15	
2	2			22	~SUM
3	22	~SUM		16	
	14			22	~VALUE
	22	~SUMSQ		12	
	14		11	22	~SUMSQ
	22	~I		15	
	14			22	~SUMSQ
	22	~VALUE		16	
	14			22	~VALUE
	22	~MEAN		18	
	14			22	~VALUE
	22	~VARIANCE	12	4	
	13			12	
	6		13	22	~MEAN
4	3			15	
5	22	~SUM		22	~SUM
	15			19	
	23	#0		23	#100
	12			12	
6	22	~SUMSQ	14	22	~VARIANCE
	15			15	
	23	#0		22	~SUMSQ
	12			19	
7	7			23	#100
	22	~I		17	
	15			22	~MEAN
	23	#1		18	
	10			22	~MEAN
	23	#100		12	
	11		15	9	
8	3			20	
9	8			22	~MEAN
	20			14	
	22	~VALUE		22	~VARIANCE
	21			21	
	12		16	5	

Can be entered directly into a symbol table and a pointer used in the output of the lexicon

Figure 5.6 Lexical scan of the program from Fig. 5.1.

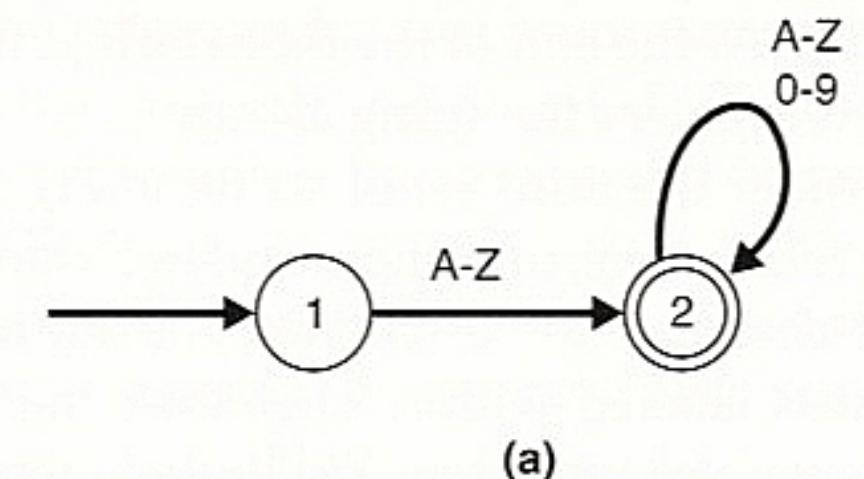
Finite Automata (FA)

- » Starting State
- » Ending State(s)
- » Transitions upon scanning a character
- » The sequence of characters are recognised (accepted) by the FA if it moves it from a start state to an end state, otherwise, it is not recognised (rejected)



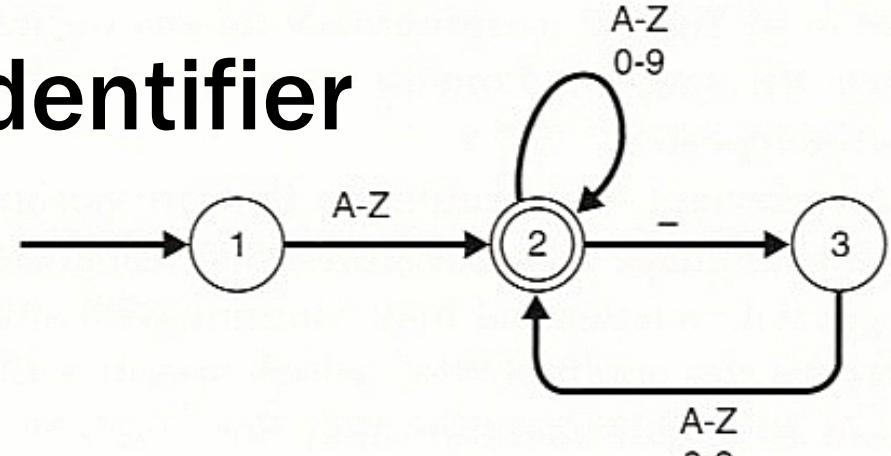
Token Recogniser

- » By grammar
 - $\text{<ident>} ::= \text{<letter>} \mid \text{<ident>} \text{<letter>} \mid \text{<ident>} \text{<digit>}$
 - $\text{<letter>} ::= \text{A} \mid \text{B} \mid \text{C} \mid \text{D} \mid \dots \mid \text{Z}$
 - $\text{<digit>} ::= \text{0} \mid \text{1} \mid \text{2} \mid \text{3} \mid \dots \mid \text{9}$
- » By scanner - modelling as finite automata
 - Figure 5.8 (a)



Recognising Identifier

```
get first Input_Character
if Input_Character in ['A'...'Z'] then
begin
    while Input_Character in ['A'...'Z', '0'...'9']
    begin
        get next Input_Character
        if Input_Character = '_' then
            begin
                get next Input_Character
                Last_Char_Is_Underscore := true
            end {if '_'}
        else
            Last_Char_Is_Underscore := false
    end {while}
    if Last_Char_Is_Underscore then
        return (Token_Error)
    else
        return (Valid_Token)
end {if first in ['A'...'Z']}
else
    return (Token_Error)
```



(b)

Identifiers allowing underscore (_)

State	A-Z	0-9	-	
1	2			{starting state}
2	2	2	3	{final state}
3	2	2		

(b)

Figure 5.10 Token recognition using (a) algorithmic code and (b) tabular representation of finite automaton.

Recognising Integer

- » Allowing leading zeroes
 - Figure 5.8 (c)
- » Disallowing leading zeroes
 - Figure 5.8 (d)

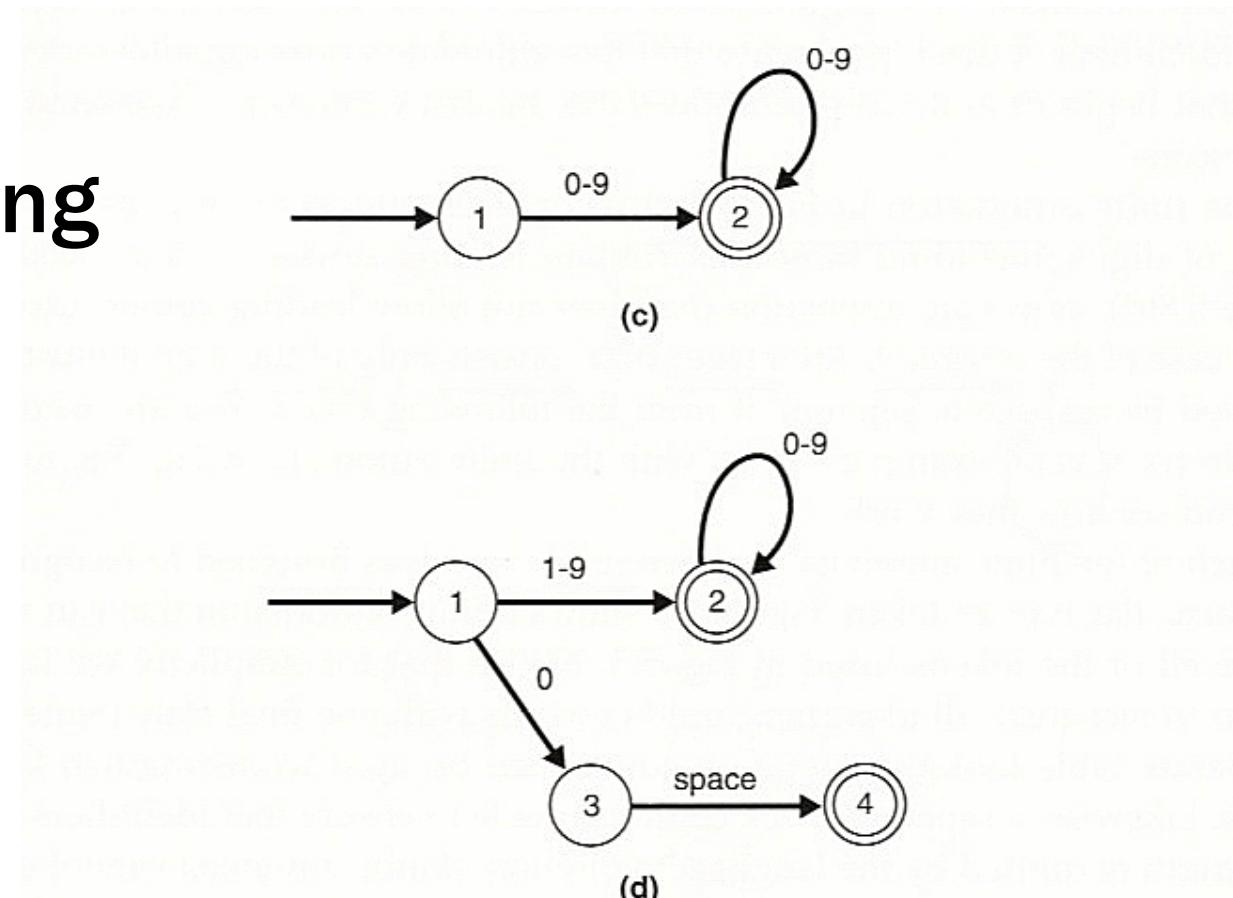


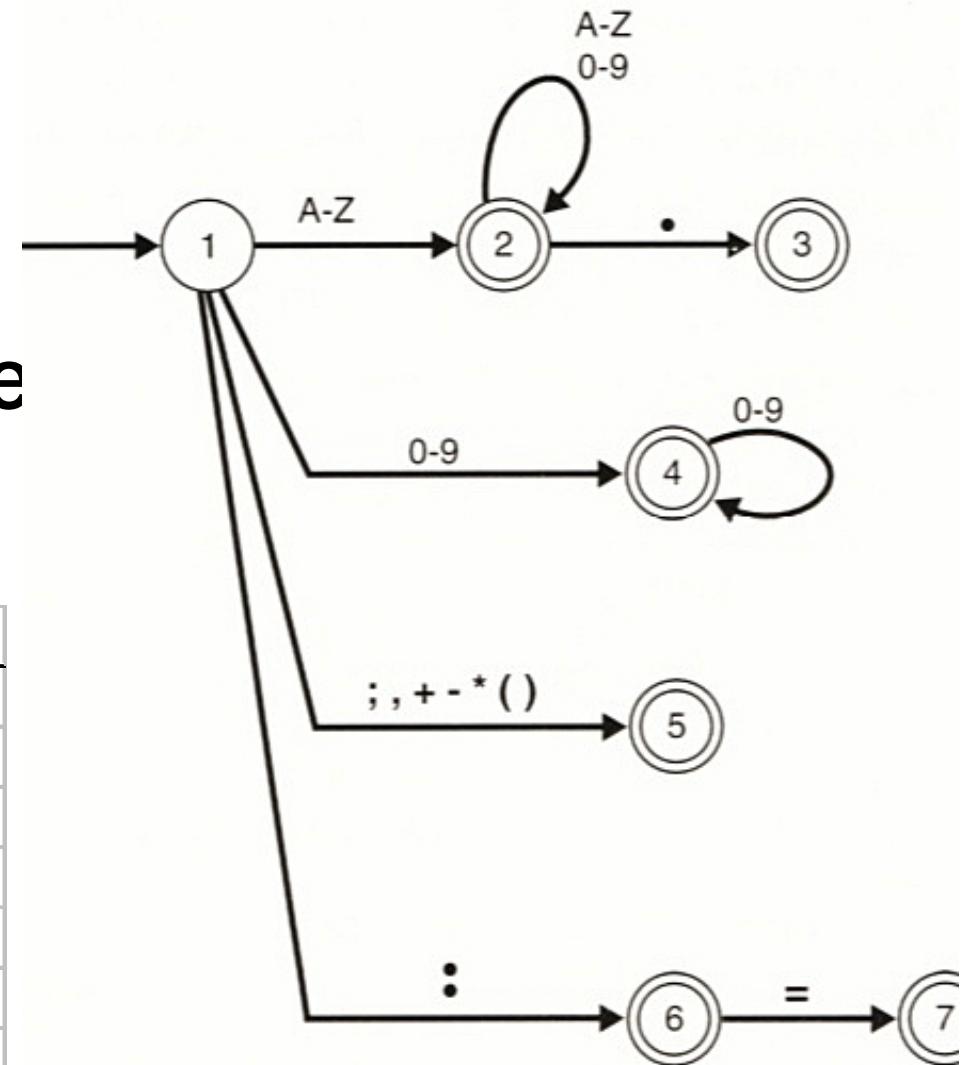
Figure 5.8 Finite automata for typical programming language tokens.

Scanner - Implementation

- » **Figure 5.10 (a)**
 - Algorithmic code for identifier recognition

- » Tabular representation of finite automaton for Figure 5.9.

State	A-Z	0-9	;,+-*()	:	=	.
1	2	4	5	6		
2	2	2				3
3						
4		4				
5						
6				7		
7						



» 5.9 Finite automaton to recognize tokens from Fi

Syntactic Analysis

- » Recognise source statements as language constructs or build the parse tree for the statements.
 - Bottom-up
 - Operator-precedence parsing
 - Shift-reduce parsing
 - LR(0) parsing
 - LR(1) parsing
 - SLR(1) parsing
 - LALR(1) parsing
 - Top-down
 - Recursive-descent parsing
 - LL(1) parsing

Operator-Precedence Parsing

- » Operator

- Any terminal symbol (or any token)

- » Precedence

- * > +

- + < *

- » Operator-precedence

- Precedence relations between operators

A + B * C - D
< >

PROGRAM ÷ VAR

and

BEGIN < FOR

; > END

but

END > ;

Precedence Matrix for the Fig. 5.2

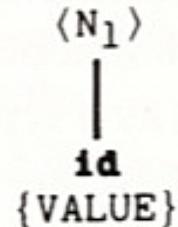
Figure 5.11 Precedence matrix for the grammar from Fig. 5.2.

Operator-Precedence Parse Example

BEGIN READ (VALUE);

(i) ... BEGIN READ (**id**)
 \Leftarrow \doteq $\Leftarrow \dots \Rightarrow$

(ii) ... BEGIN READ ($\langle N_1 \rangle$) ;
 \Leftarrow \doteq \doteq \Rightarrow



(iii) ... BEGIN $\langle N_2 \rangle$;
 $\dots \dots$

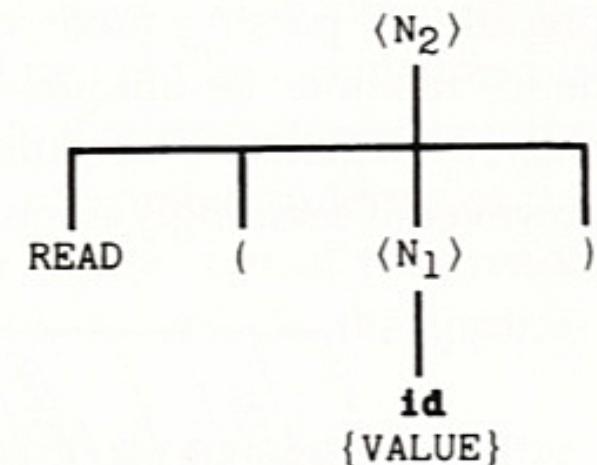
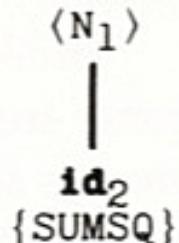


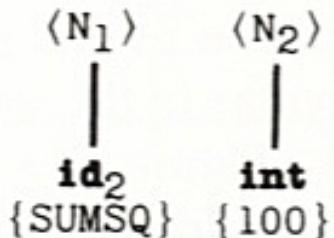
Figure 5.12 Operator-precedence parse of a READ statement.

Operator-Precedence Parse Example

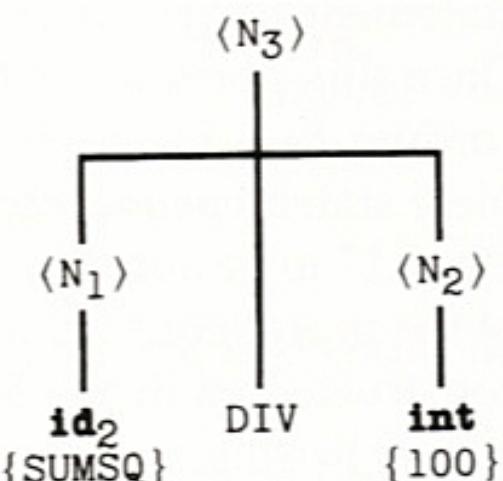
... **id₁** := **id₂** DIV
 < ÷ < >



... **id₁** := ⟨N₁⟩ DIV **int**
 < ÷ < >



... **id₁** := ⟨N₁⟩ DIV ⟨N₂⟩
 < ÷ < >

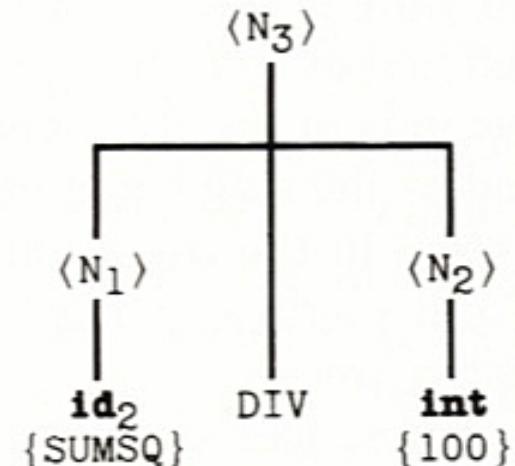


... **id₁** := ⟨N₃⟩ - < **id₃** >

Operator-Precedence Parse Example

... **id₁** := ⟨N₃⟩ - ⟨N₄⟩ * **id₄** ;
 < ÷ ⟨ < ⟨ > >

.....



... **id₁** := ⟨N₃⟩ - ⟨N₄⟩ * ⟨N₅⟩ ;
 < ÷ ⟨ < >

.....

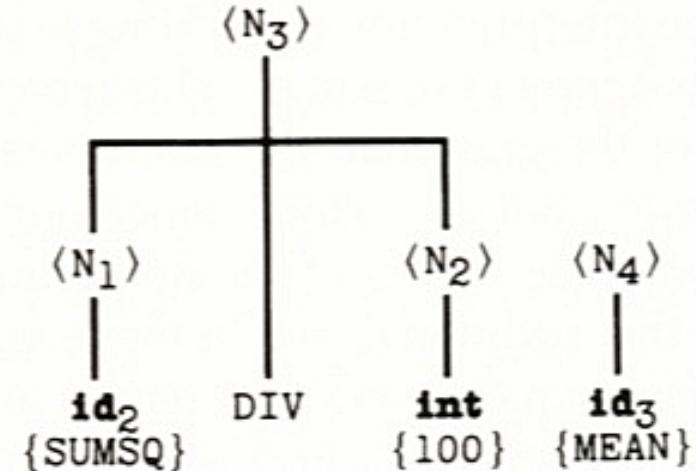
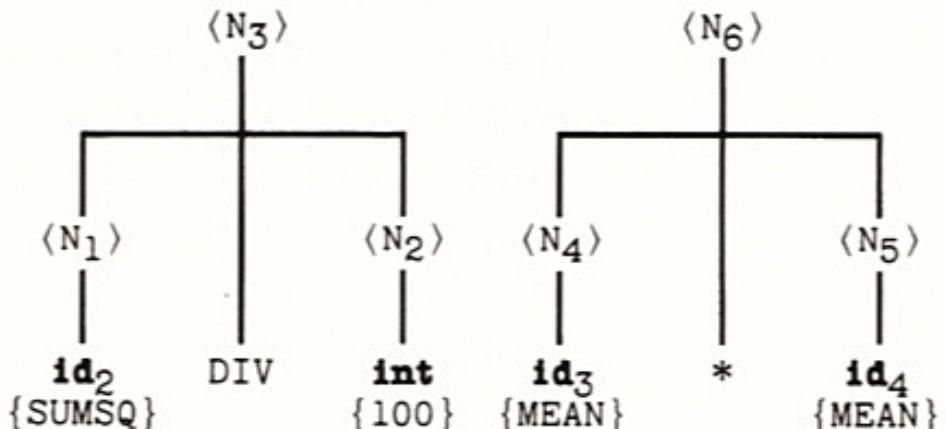


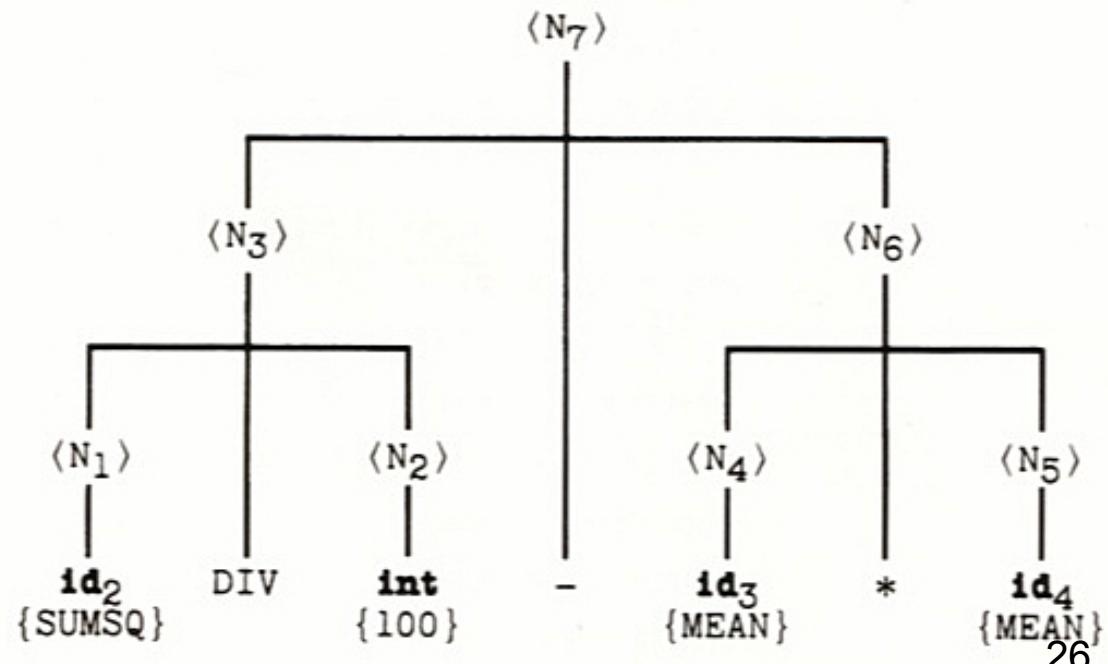
Figure 5.13 Operator-precedence parse of an assignment statement.

Operator-Precedence Parse Example

(vii) ... **id₁** := ⟨N₃⟩ - ⟨N₆⟩ ;
 \ll \doteq \ll \gg



(viii) ... **id₁** := ⟨N₇⟩ ;
 \ll \doteq \gg



Operator-Precedence Parsing

- » Bottom-up parsing
 - » Generating precedence matrix
 - Aho et al. (1988)

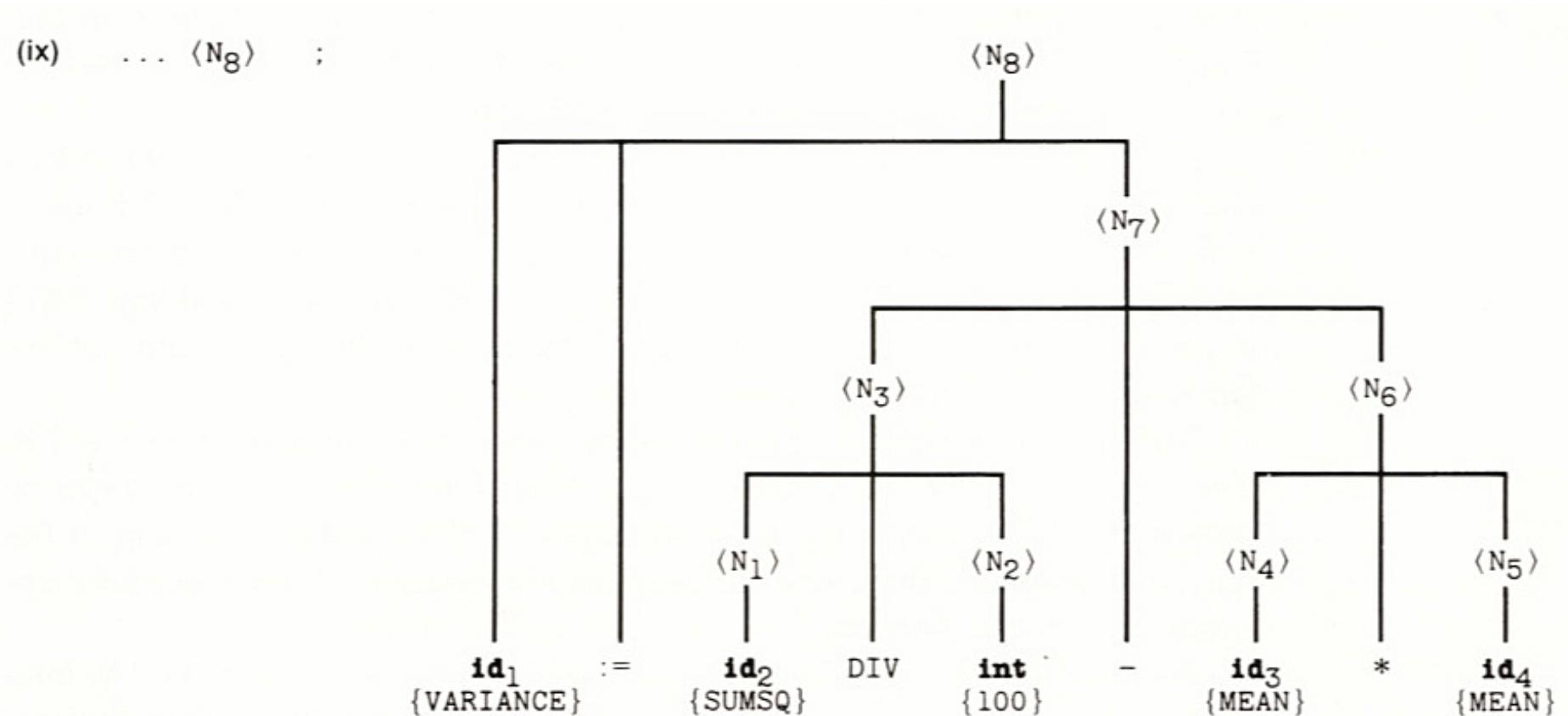


Figure 5.13 (cont'd)

Shift-reduce Parsing with Stack

» Figure 5.14

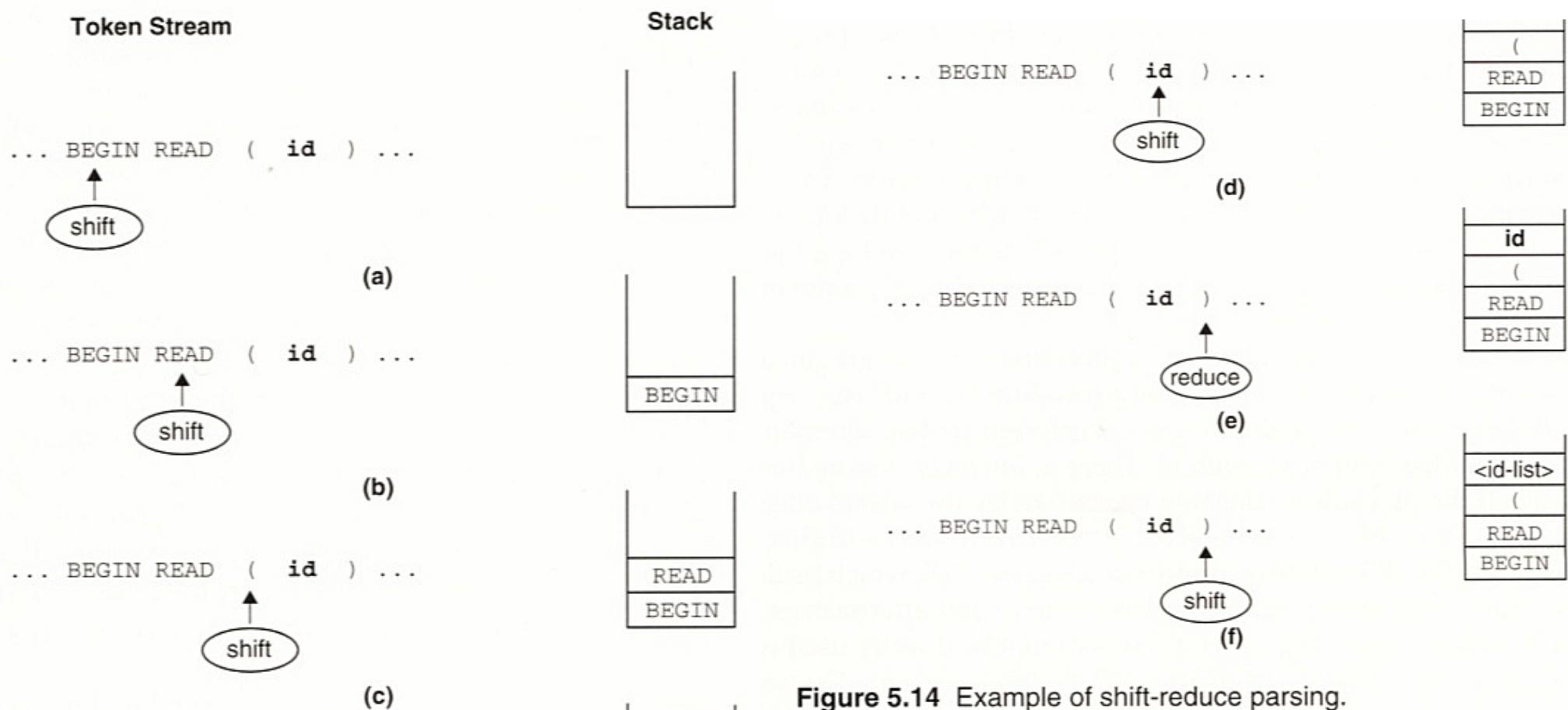


Figure 5.14 Example of shift-reduce parsing.

Recursive-Descent Parsing

- » Each **nonterminal symbol** in the grammar is associated with a **procedure**.
 - » $\langle \text{read} \rangle ::= \text{READ } (\langle \text{id-list} \rangle)$
 - » $\langle \text{id-list} \rangle ::= \text{id} \mid \langle \text{id-list} \rangle, \text{id}$
- » Testing for alternatives:
 - » $\langle \text{stmt} \rangle ::= \langle \text{assign} \rangle \mid \langle \text{read} \rangle \mid \langle \text{write} \rangle \mid \langle \text{for} \rangle$

Recursive-Descent Parsing of READ

```
procedure READ
begin
    FOUND := FALSE
    if TOKEN = 8 {READ} then
        begin
            advance to next token
            if TOKEN = 20 { ( ) } then
                begin
                    advance to next token
                    if IDLIST returns success then
                        if TOKEN = 21 { ) } then
                            begin
                                FOUND := TRUE
                                advance to next token
                            end {if ) }
                        end {if ( )}
                    end {if READ}
                if FOUND = TRUE then
                    return success
                else
                    return failure
            end {READ}
```

Recursive-Descent Parsing of IDLIST

```
procedure IDLIST
begin
    FOUND := FALSE
    if TOKEN = 22 {id} then
        begin
            FOUND := TRUE
            advance to next token
            while (TOKEN = 14 {,}) and (FOUND = TRUE) do
                begin
                    advance to next token
                    if TOKEN = 22 {id} then
                        advance to next token
                    else
                        FOUND := FALSE
                end {while}
            end {if id}
            if FOUND = TRUE then
                return success
            else
                return failure
        end {IDLIST}
```

(a)

Figure 5.16 Recursive-descent parse of a READ statement.

Recursive-Descent Parsing (cont'd.)

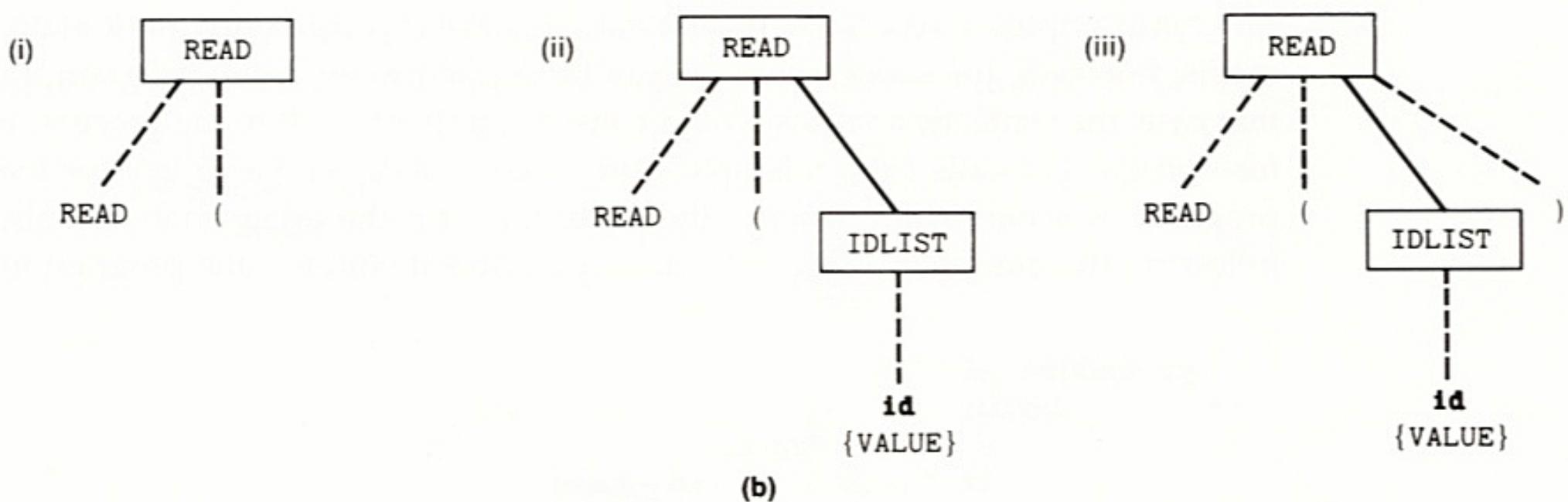


Figure 5.16 (cont'd)

Eliminating Left-Recursion

- » **Left Recursion in lines 3, 6, 7, 10 and 11**

```
<id-list> ::= id { , id }
```

- » **Modification:**

```
6 <id-list>      ::= id | <id-list> , id
```

- **<dec-list> ::= <dec> | <dec-list>;<dec>**

- » **Modification:**

- **<dec-list> ::= <dec> {;<dec>}**

Recursive-Descent Parsing (cont'd.)

```
1  <prog>      ::= PROGRAM <prog-name> VAR <dec-list> BEGIN <stmt-list> END.
2  <prog-name>  ::= id
3a <dec-list>   ::= <dec> { ; <dec> }
4   <dec>       ::= <id-list> : <type>
5   <type>      ::= INTEGER
6a <id-list>   ::= id { , id }
7a <stmt-list> ::= <stmt> { ; <stmt> }
8   <stmt>      ::= <assign> | <read> | <write> | <for>
9   <assign>    ::= id := <exp>
10a <exp>       ::= <term> { + <term> | - <term> }
11a <term>      ::= <factor> { * <factor> | DIV < factor> }
12   <factor>   ::= id | int | ( <exp> )
13   <read>      ::= READ ( <id-list> )
14   <write>     ::= WRITE ( <id-list> )
15   <for>        ::= FOR <index-exp> DO <body>
16   <index-exp> ::= id := <exp> TO <exp>
17   <body>       ::= <stmt> | BEGIN <stmt-list> END
```

Figure 5.15 Simplified Pascal grammar modified for recursive-descent parse.

Recursive-Descent Parsing of ASSIGN

```
procedure ASSIGN
begin
    FOUND := FALSE
    if TOKEN = 22 {id} then
        begin
            advance to next token
            if TOKEN = 15 { := } then
                begin
                    advance to next token
                    if EXP returns success then
                        FOUND := TRUE
                end {if := }
            end {if id}
            if FOUND = TRUE then
                return success
            else
                return failure
        end {ASSIGN}
```

Recursive-Descent Parsing of EXP

```
procedure EXP
begin
    FOUND := FALSE
    if TERM returns success then
        begin
            FOUND := TRUE
            while ((TOKEN = 16 {+}) or (TOKEN = 17 {-}))
                and ( FOUND = TRUE ) do
                begin
                    advance to next token
                    if TERM returns failure then
                        FOUND := FALSE
                end {while}
            end {if TERM}
            if FOUND = TRUE then
                return success
            else
                return failure
        end {EXP}
```

Figure 5.17 Recursive-descent parse of an assignment statement.

Recursive-Descent Parsing of TERM

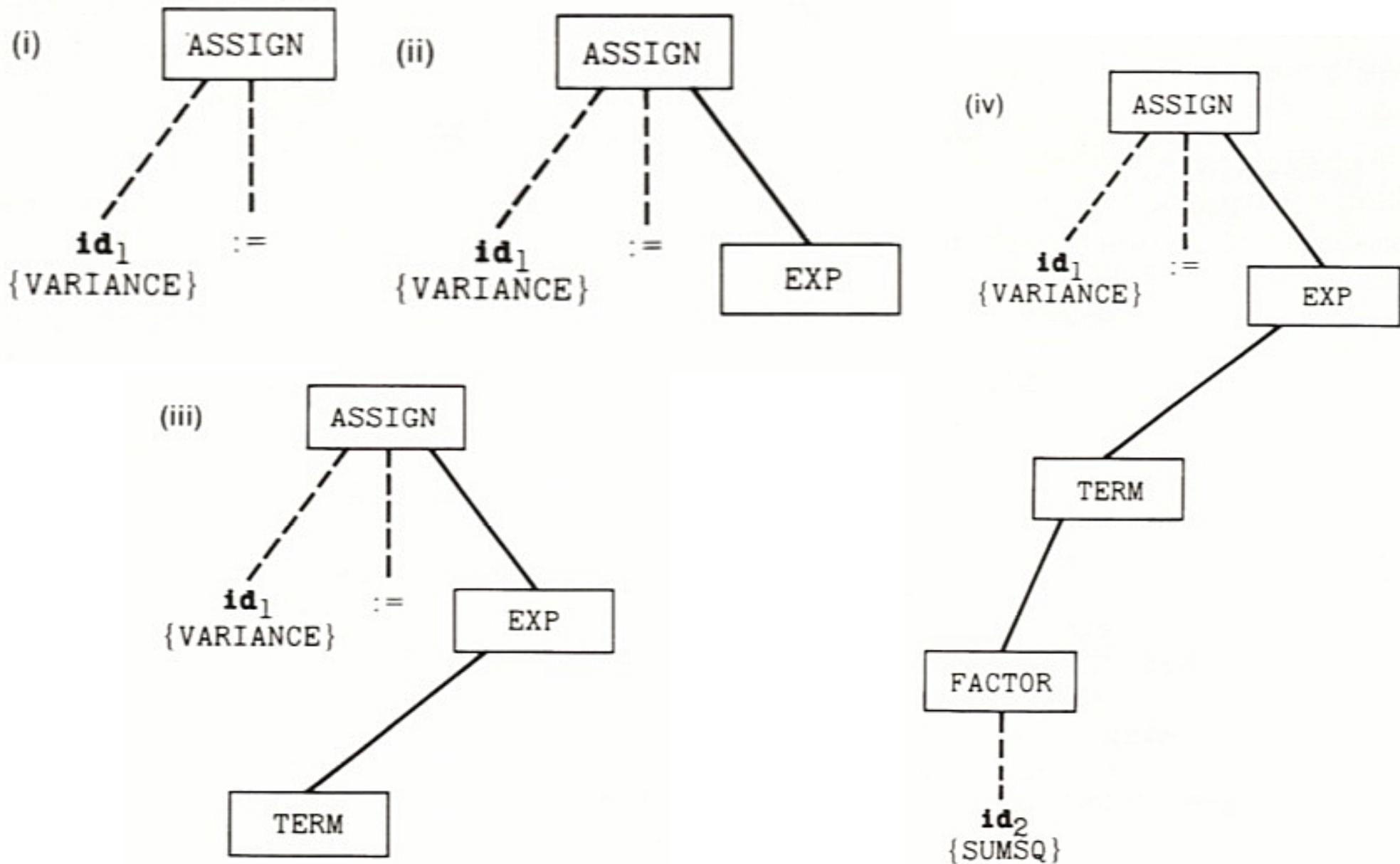
```
procedure TERM
begin
    FOUND := FALSE
    if FACTOR returns success then
        begin
            FOUND := TRUE
            while ({TOKEN = 18 {*}}) or (TOKEN = 19 {DIV
                and (FOUND = TRUE) do
                    begin
                        advance to next token
                        if FACTOR returns failure then
                            FOUND := FALSE
                    end {while}
            end {if FACTOR}
            if FOUND = TRUE then
                return success
            else
                return failure
        end {TERMINATOR}
    end {if FOUND = TRUE}
```

procedure FACTOR

begin

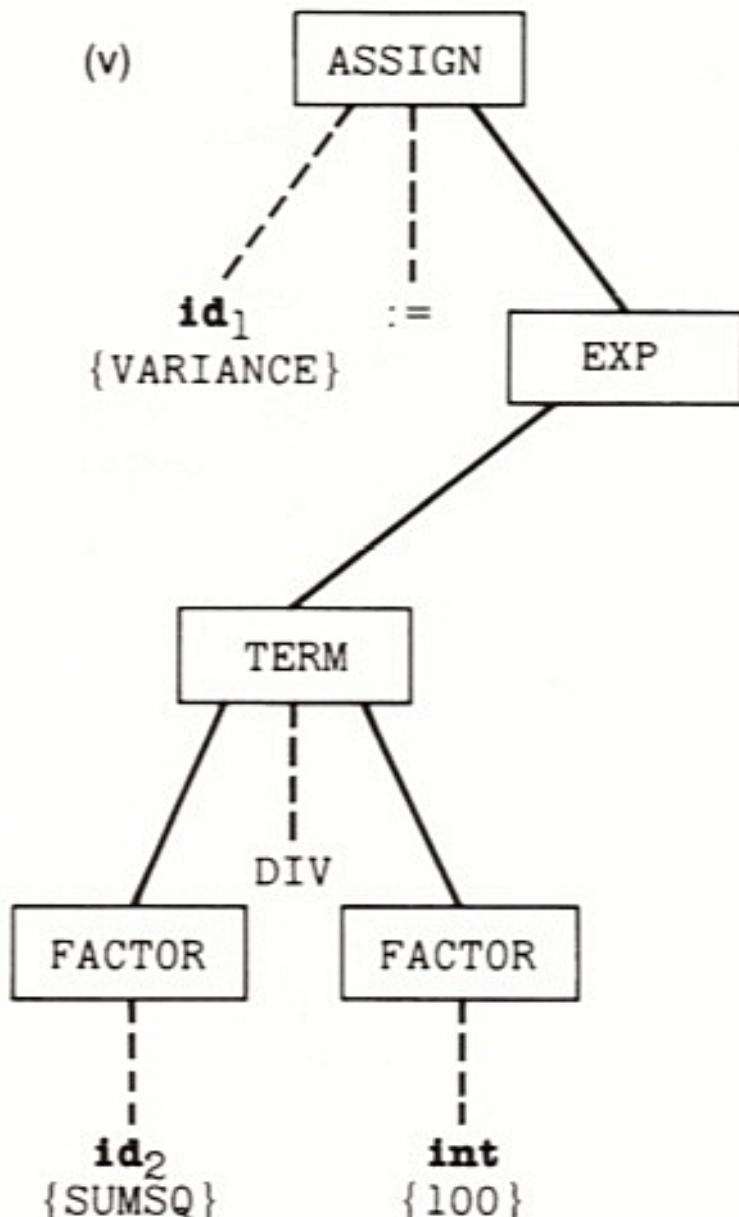
```
    FOUND := FALSE
    if (TOKEN = 22 {id}) or (TOKEN = 23 {int}) then
        begin
            FOUND := TRUE
            advance to next token
        end {if id or int}
    else
        if TOKEN = 20 { ( ) } then
            begin
                advance to next token
                if EXP returns success then
                    if TOKEN = 21 { ) } then
                        begin
                            FOUND := TRUE
                            advance to next token
                        end {if )
                    end {if ( )
                if FOUND = TRUE then
                    return success
            else
                return failure
    end {FACTOR}
```

Recursive-Descent Parsing (cont'd.)

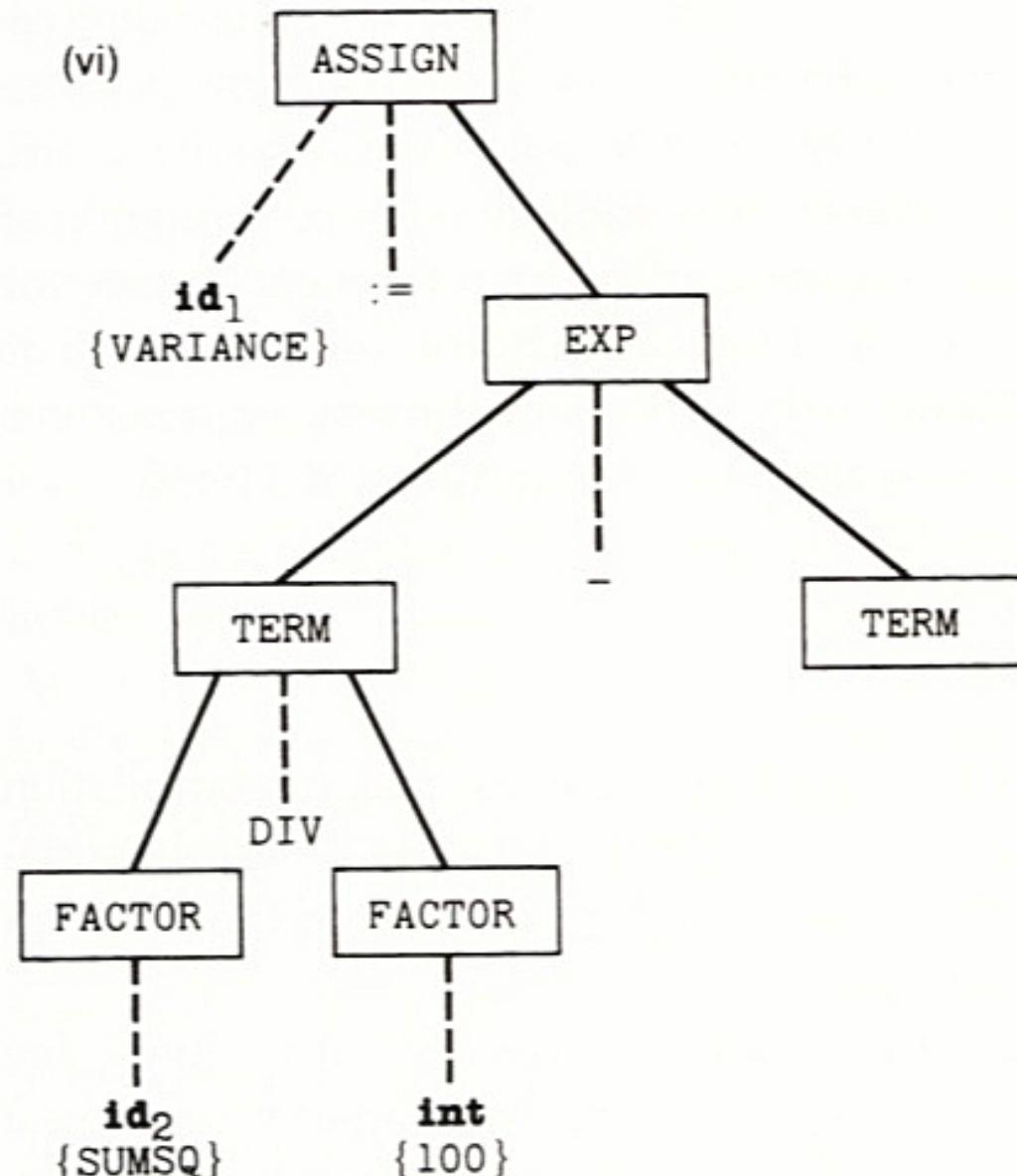


Recursive-Descent Parsing (cont'd.)

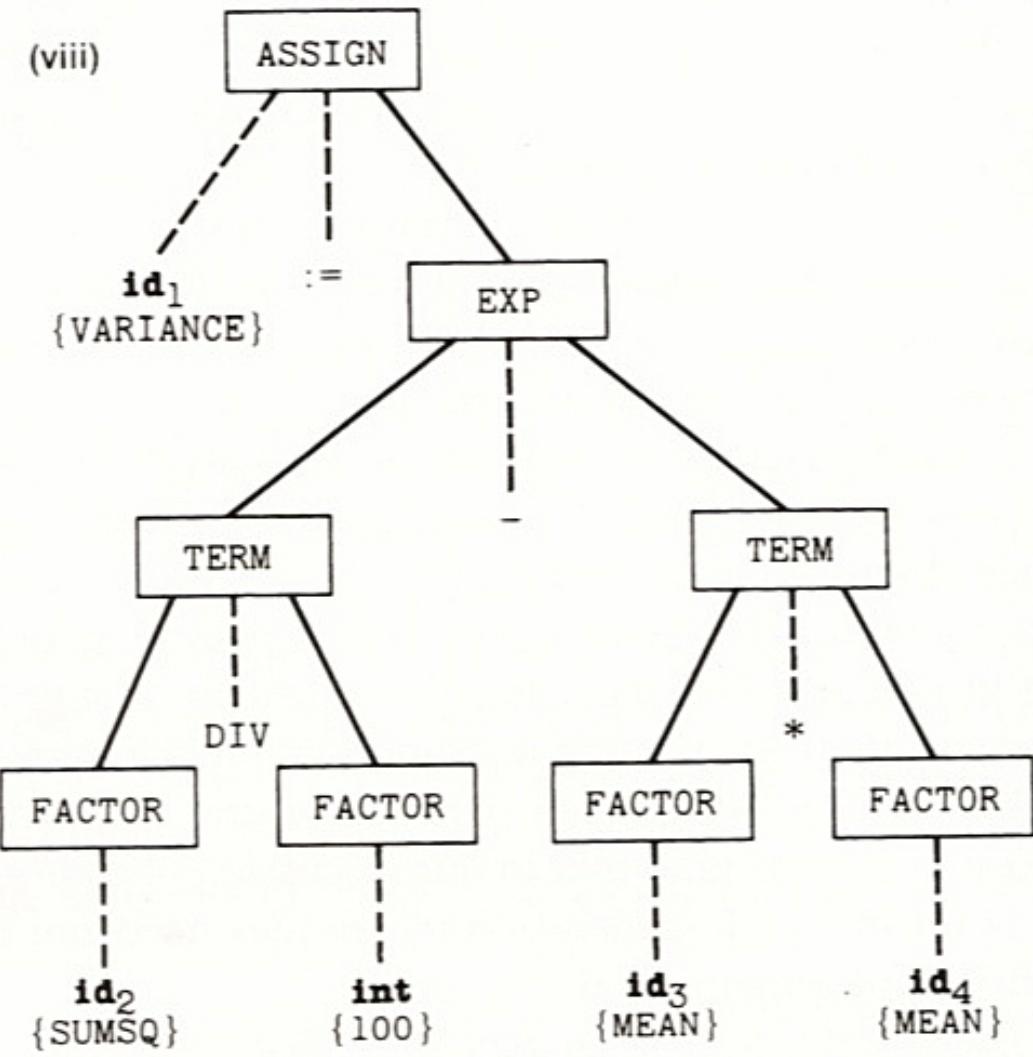
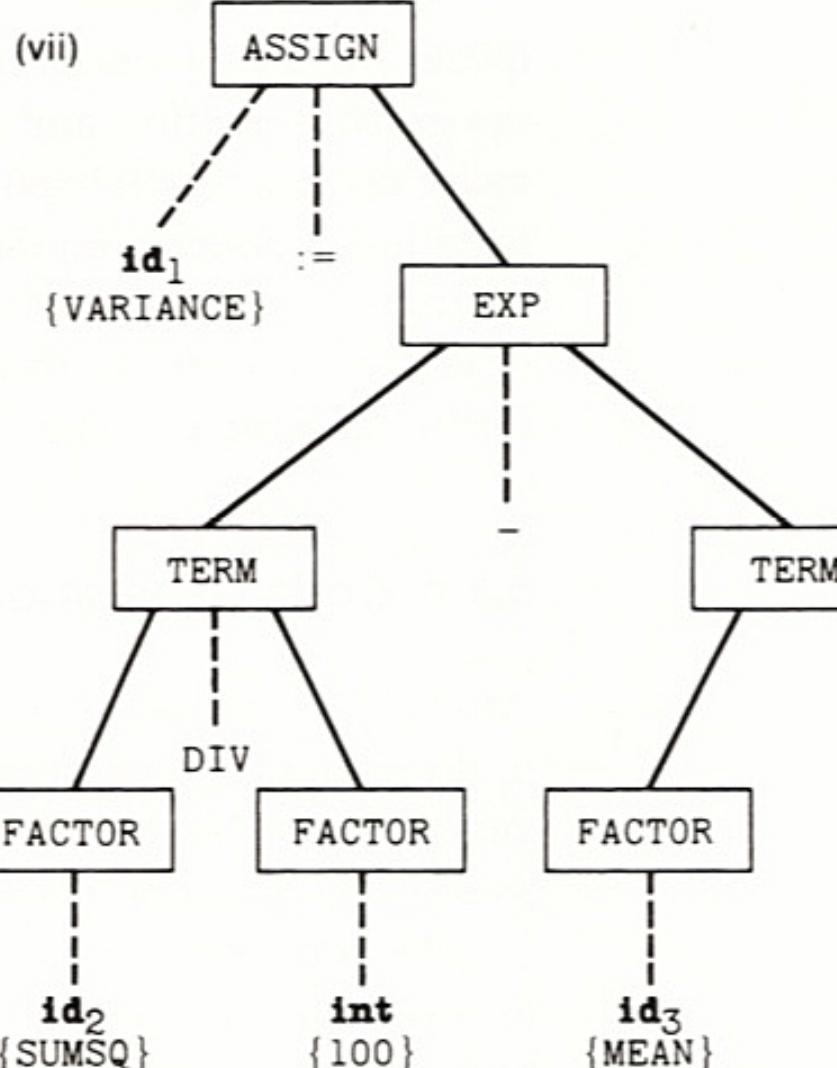
(v)



(vi)



Recursive-Descent Parsing (cont'd.)



(b)

EXAMPLE:

Fahrenheit := 32 + celsius * 1.8

| f | a | h | r | e | n | h | e | i | t | : | = | 3 | 2 | + | c | e | l | s | i | o | u | s | * | 1 | . | 8 | ; |

Getchar()

Lexical analyser (scanner) (converts from character stream into a stream of tokens.)

Symbol Table

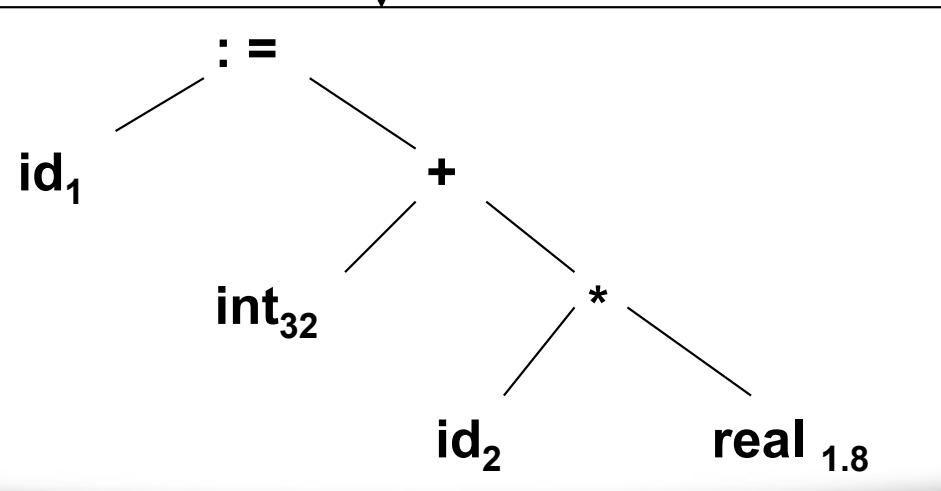
1	fahrenheit	real
2	celsius	real

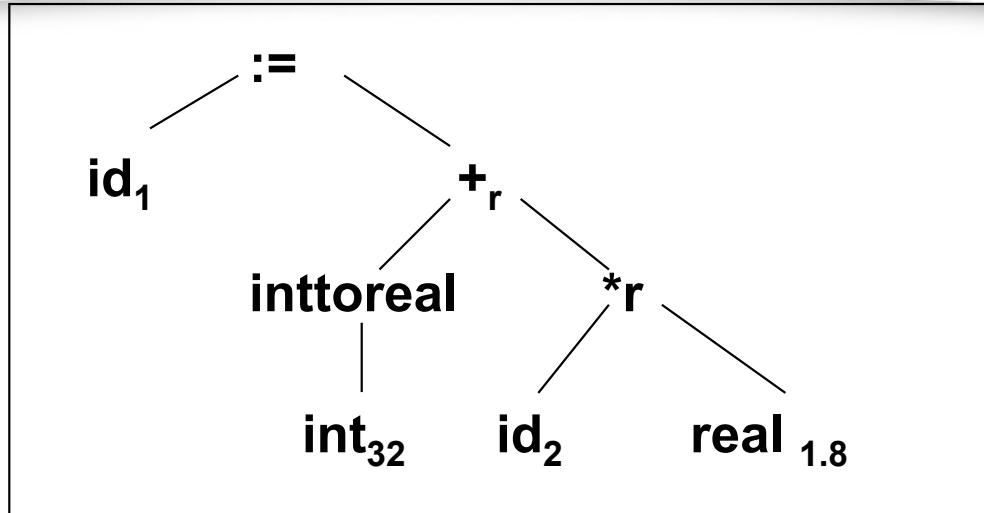
↑ ↑
name attribute

[id, 1][:=][int, 32][+][id, 2][*][int, 1.8][;]

index in symbol table

Syntax analyser (parser)
(Construct syntactic structure of the program)





Symbol Table

1	fahrenheit	real
2	celsius	real

Intermediate code generator

```

Temp1 := inttoreal(32)
Temp2 := id2
Temp2 := Temp2 * 1.8
Temp1 := Temp1 + Temp2
id1 := Temp1
  
```

Intermediate code

