

# **Programming Language Syntax - Scanning -**

Chapter 2, Sections 2.1-2.2

# Regular Expressions

- *Token*: a shortest string of characters with meaning
- Tokens – specified by regular expressions
- An *alphabet*  $\Sigma$  is any finite nonempty set
  - Examples:
    - English:  $\{a, b, \dots, z\}$ ,
    - binary:  $\{0, 1\}$
    - $\{a, b, \dots, z, 0, 1, \dots, 9, \cdot, |, *, \epsilon\}$
- The set of all finite strings over  $\Sigma$  is denoted  $\Sigma^*$
- The *empty* string:  $\epsilon \in \Sigma^*$  (has zero characters)

$\uparrow$   
a blank.  
 $\Sigma = \{a, b, c, \dots, z, " " \}$   
 $\epsilon \neq \epsilon$

$|a| = 1$   
 $|ab| = 2$   
 $|\epsilon| = 0$

# Regular Expressions

- Regular expressions ✓ 331?
- Regular expressions over an alphabet  $\Sigma$  are all strings obtained as follows:

- $\epsilon$  is a regular expression
- any character  $a \in \Sigma$  is a regular expression
- For reg. exp.  $\alpha, \beta$ , the following are reg. exp.:

- $\alpha \cdot \beta$  - concatenation ('.' omitted:  $\alpha\beta$ ) e.g.  $\alpha = a, \beta = ab, \alpha\beta = aab$   $ab \cdot \epsilon = ab$
- $\alpha \mid \beta$  - union ('|' = or) (sometimes denoted  $\alpha + \beta$ )
- $\alpha^*$  - Kleene star (0 or more repetitions)  $\alpha\beta = aab$  it could be infinite.
- $\alpha^+ = \alpha \alpha^*$  (1 or more repetitions)

reg expr  
 $\epsilon$   
 $01110$   
 $0^*$

1 2 3 4  
 $\alpha$   
 $\alpha^2$   
 $\alpha^3$

str expr  
 $\{\epsilon\}$   
 $\{011, 0\}$   
 $\{\epsilon, 0, 00, 000, \dots\}$

# Regular Expressions

- Example: Signed integers:

*sign\_int*  $\rightarrow$  (+ | - |  $\epsilon$ ) (0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9) (0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9)\*

*no sign*  
*sign*  
*one digit.*  
*followed by 0 or more digits.*

- Example: Numerical constants:

*number*  $\rightarrow$  *integer* | *real* *floating-point num.*

*integer*  $\rightarrow$  *digit digit\**

*real*  $\rightarrow$  *integer exponent* | *decimal* (*exponent* |  $\epsilon$ ) *optional exponent*

*decimal*  $\rightarrow$  *digit\** ( *.* *digit* | *digit* *.* ) *digit\**  $\Rightarrow \begin{cases} 1.2 \\ 1. \\ .2 \end{cases}$

*exponent*  $\rightarrow$  ( *e* | *E* ) ( *+* | *-* |  $\epsilon$  ) *integer*

*digit*  $\rightarrow$  0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

*decimal*  
 $\rightarrow$  *digit\** . *digit\** |  
*digit\** . *digit\**

- ' $\rightarrow$ ' means "can be"
- Precedence order: '\*' > '.' > '|'



# Regular Expressions

- Other applications:
  - grep family of tools in Unix
  - many editors
  - scripting languages:
    - Perl
    - Python
    - Ruby
    - awk
    - sed

# Formatting issues

- Upper vs. lower case
  - distinct in some languages: C, Python, Perl
  - same in others: Fortran, Lisp, Ada
- Identifiers: letters, digits, underscore (most languages)
  - camel case: `someIdentifierName`
  - underscore: `some_identifier_name`
- Unicode
  - non-Latin characters have become important
- White spaces
  - usually ignored
  - separate statements: Python, Haskell, Go, Swift
  - indentation important: Python, Haskell

# Context-Free Grammars

## Context Free Grammar (CFG)

- CFG consists of:
  - A set of *terminals*,  $T$
  - A set of *non-terminals*,  $N$
  - A start symbol,  $S \in N$
  - A set of *productions*; subset of  $N \times (N \cup T)^*$
- Example: Balanced parentheses:

$$\begin{array}{l} S \rightarrow \varepsilon \\ S \rightarrow SS \\ S \rightarrow (S) \end{array}$$

*Handwritten notes:*  
 $x \rightarrow \alpha$   
 $LHS \rightarrow RHS$

# Context-Free Grammars

- Example: CFG for arithmetic expressions:

$$expr \rightarrow id \mid number \mid - expr \mid ( expr ) \mid expr op expr$$
$$op \rightarrow + \mid - \mid * \mid / \mid$$

- *Derivation*: start with  $S$ , continue with productions

- replace LHS nonterminal by the RHS

- Example: generate the string: `slope * x + intercept`

give a CFG for a programming language;

$$\underline{expr} \Rightarrow expr op \underline{expr} \quad (S = expr)$$
$$\Rightarrow expr \underline{op} id$$

( $\Rightarrow$  = “derives”)

$$\Rightarrow \underline{expr} + id$$

( $\Rightarrow^*$  = 0 or more steps)

$$\Rightarrow expr op \underline{expr} + id$$
$$\Rightarrow expr \underline{op} id + id$$
$$\Rightarrow \underline{expr} * id + id$$
$$\Rightarrow id * id + id$$

(slope) (x) (intercept)

- Sentential form: any string along the way



# Context-Free Grammars

- *Right-most derivation*: the rightmost nonterminal is replaced

$$\begin{aligned}\underline{expr} &\Rightarrow expr\ op\ \underline{expr} \\ &\Rightarrow expr\ \underline{op}\ id \\ &\Rightarrow \underline{expr} + id \\ &\Rightarrow expr\ op\ \underline{expr} + id \\ &\Rightarrow expr\ \underline{op}\ id + id \\ &\Rightarrow \underline{expr} * id + id \\ &\Rightarrow id * id + id\end{aligned}$$

- *Left-most derivation*: the leftmost nonterminal is replaced

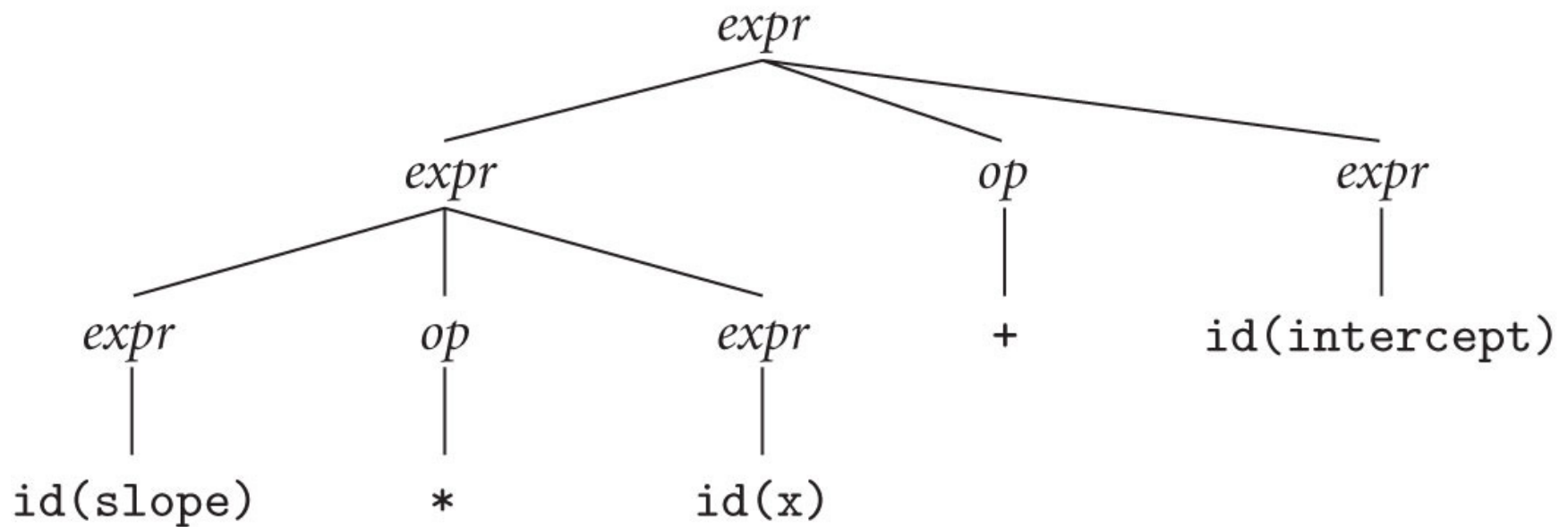
$$\begin{aligned}\underline{expr} &\Rightarrow \underline{expr}\ op\ expr \\ &\Rightarrow \underline{expr}\ op\ expr\ op\ expr \\ &\Rightarrow id\ \underline{op}\ expr\ op\ expr \\ &\Rightarrow id * \underline{expr}\ op\ expr \\ &\Rightarrow id * id\ \underline{op}\ expr \\ &\Rightarrow id * id + \underline{expr} \\ &\Rightarrow id * id + id\end{aligned}$$

# Context-Free Grammars

## Parse Tree

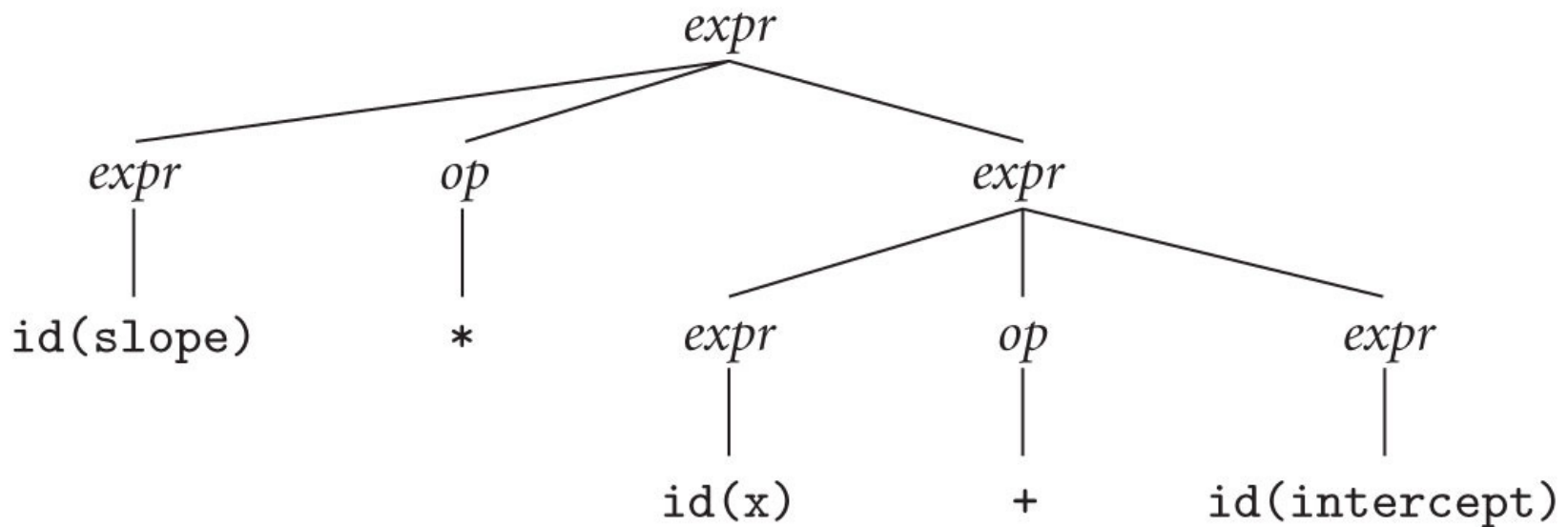
- Represents a derivation graphically
- Example: Parse tree for the string:

slope \* x + intercept



# Context-Free Grammars

- Different parse tree for: `slope * x + intercept`
- Tree allowed by the grammar but incorrect for the expression



*bad*

- Ambiguous grammar: two different parse trees for one string
  - Ambiguity is a problem for parsers
  - We want unambiguous grammars

# Context-Free Grammars

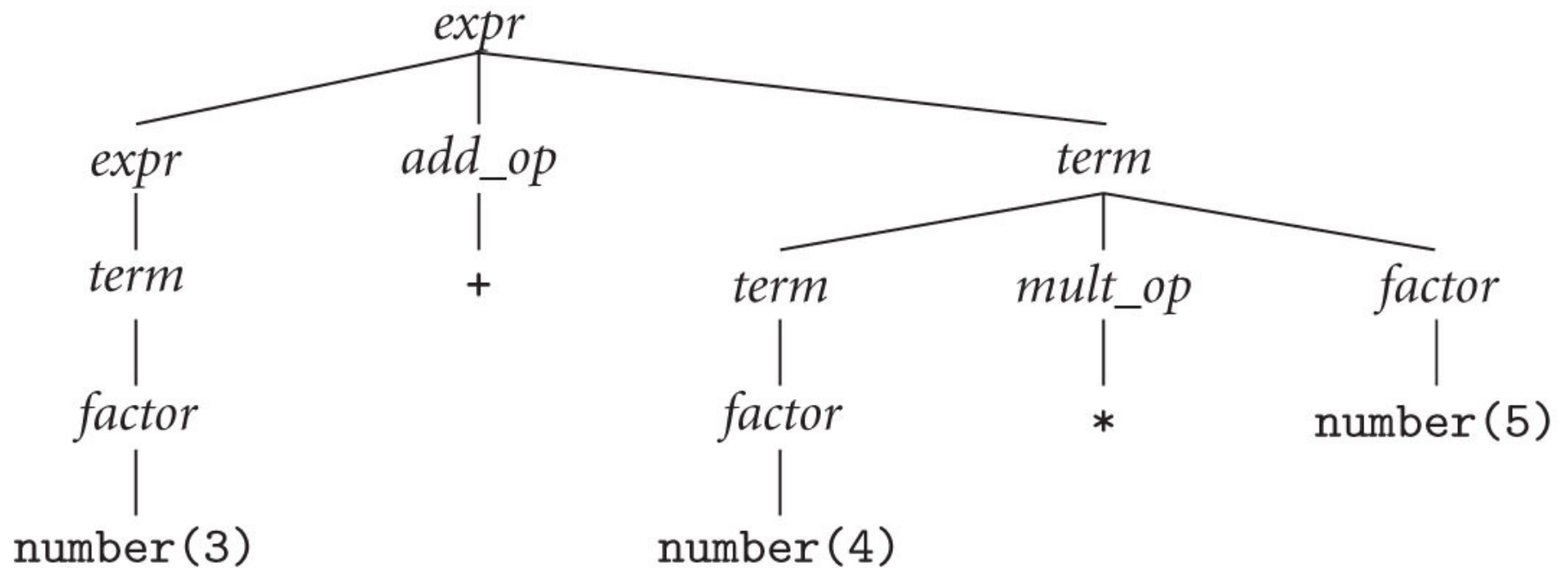
- Better version – unambiguous
- Captures *associativity* and *precedence*

$$\text{expr} \rightarrow \text{term} \mid \text{expr add\_op term}$$
$$\text{term} \rightarrow \text{factor} \mid \text{term mult\_op factor}$$
$$\text{factor} \rightarrow \text{id} \mid \text{number} \mid - \text{factor} \mid ( \text{expr} )$$
$$\text{add\_op} \rightarrow + \mid -$$
$$\text{mult\_op} \rightarrow * \mid /$$



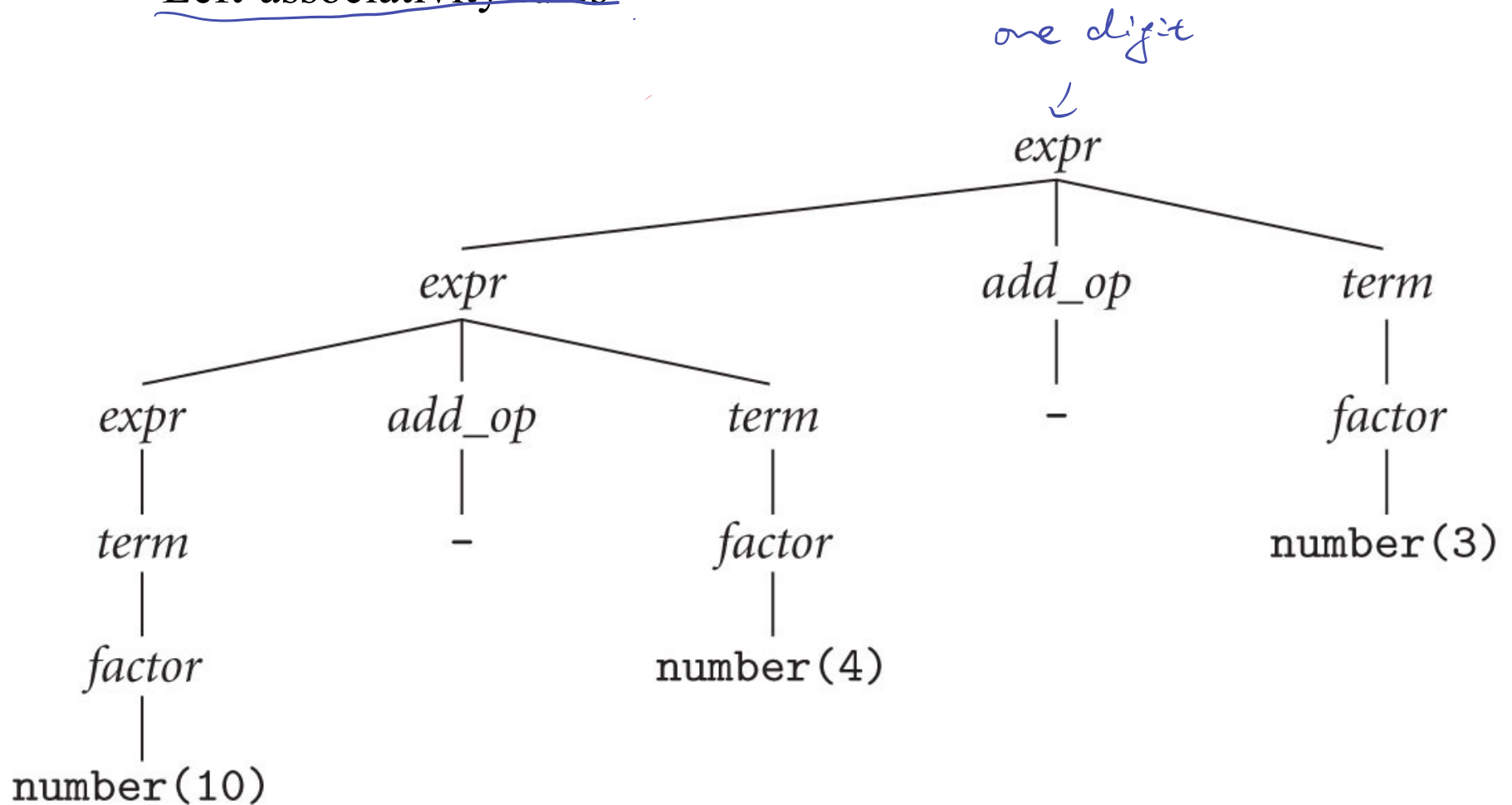
# Context-Free Grammars

- Parse tree for:  $3 + 4 * 5$ 
  - Precedence rules



# Context-Free Grammars

- Parse tree for:  $10 - 4 - 3$ 
  - Left-associativity rules



# Scanning

*Scanning = Lexical Analysis*

- tokenizing source
- removing comments
- saving text of identifiers, numbers, strings
- saving source locations (file, line, column) for error messages

# Scanning

- Example: simple calculator language

*assign*  $\rightarrow$  `:=`

(Algol style; C has ‘=’)

*plus*  $\rightarrow$  `+`

*minus*  $\rightarrow$  `-`

*times*  $\rightarrow$  `*`

*div*  $\rightarrow$  `/`

*lparen*  $\rightarrow$  `(`

*rparen*  $\rightarrow$  `)`

*id*  $\rightarrow$  *letter* ( *letter* | *digit* )<sup>\*</sup> (except for `read` and `write`)

*number*  $\rightarrow$  *digit*<sup>int</sup> *digit*<sup>\*</sup> | *digit*<sup>float</sup> ( *digit* | *digit* . ) *digit*<sup>\*</sup>

*comment*  $\rightarrow$  `/*` ( *non-\** | *\* non-/* )<sup>\*</sup> *\*+ /*  
| `///` ( *non-newline* )<sup>\*</sup> *newline*



# Scanning

## Ad-hoc scanner

- Longest possible token extracted
- White spaces are delimiters

finite automaton:

- vertices: states

- edges: oriented/labelled by a single char

0: start

@: final state -

skip any initial white space (spaces, tabs, and newlines)

if  $\text{cur\_char} \in \{ '(', ')', '+', '-', '*' \}$  *tokens*

return the corresponding single-character token

if  $\text{cur\_char} = ':'$

read the next character

if it is '=' then return *assign* else announce an error

if  $\text{cur\_char} = '/'$

peek at the next character

if it is '\*' or '/' *comment*.

read additional characters until "\*" or "/" or *newline* is seen, respectively

jump back to top of code

else return *div*

if  $\text{cur\_char} = .$

read the next character

if it is a digit

read any additional digits

return *number*

else announce an error

if  $\text{cur\_char}$  is a digit

read any additional digits and at most one decimal point

return *number*

if  $\text{cur\_char}$  is a letter

read any additional letters and digits

check to see whether the resulting string is **read** or **write**

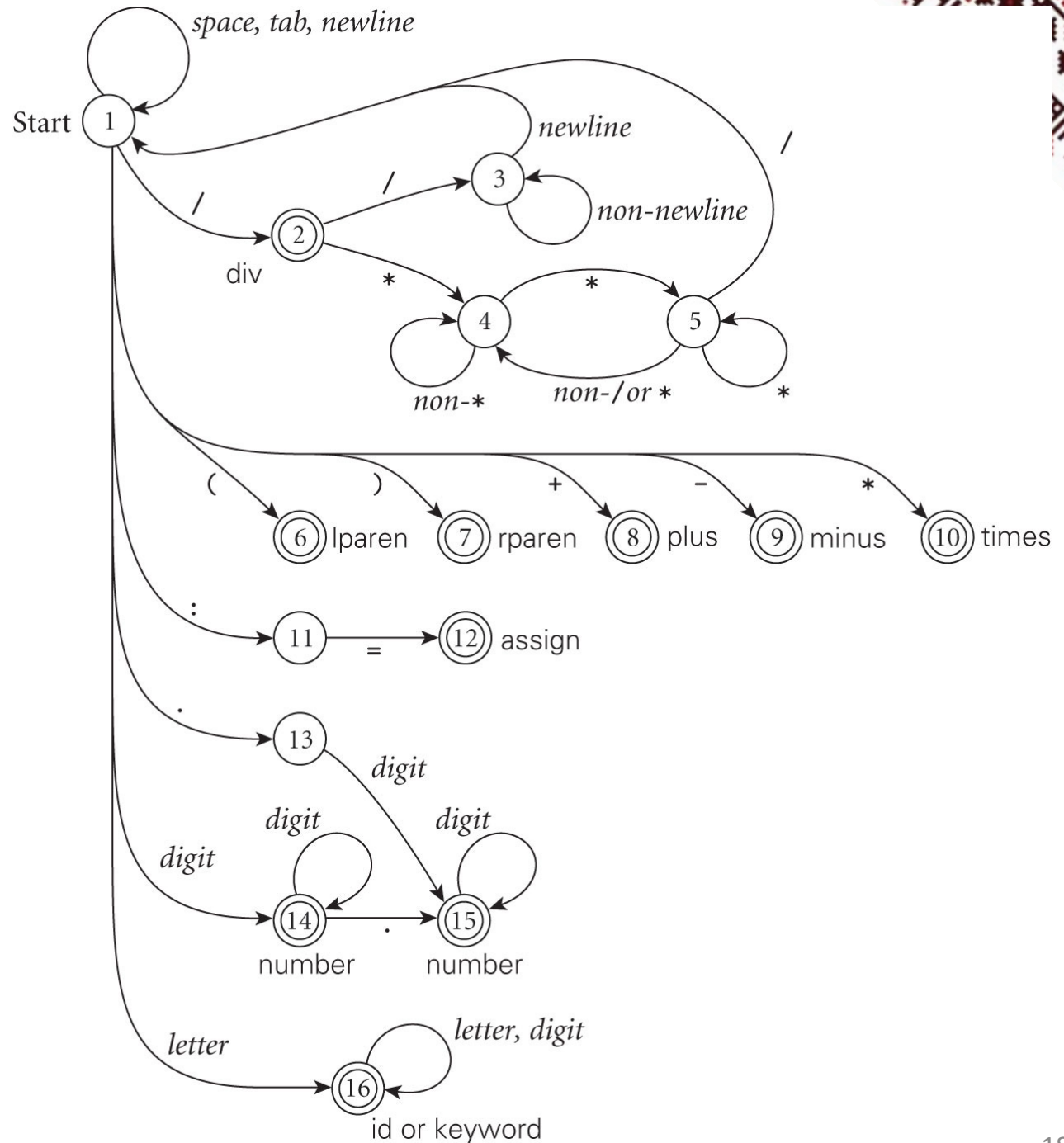
if so then return the corresponding token

else return *id*

else announce an error

# Scanning

- Structured scanner
- DFA – Deterministic Finite Automaton
- Separate final state for each token category



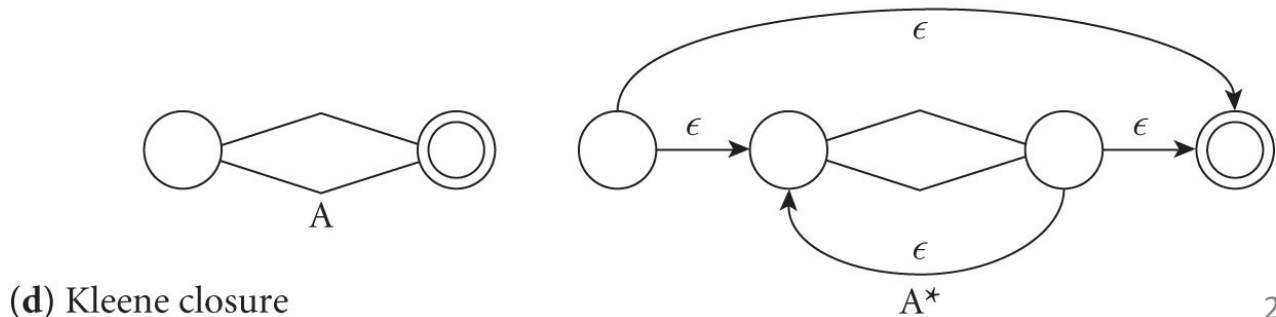
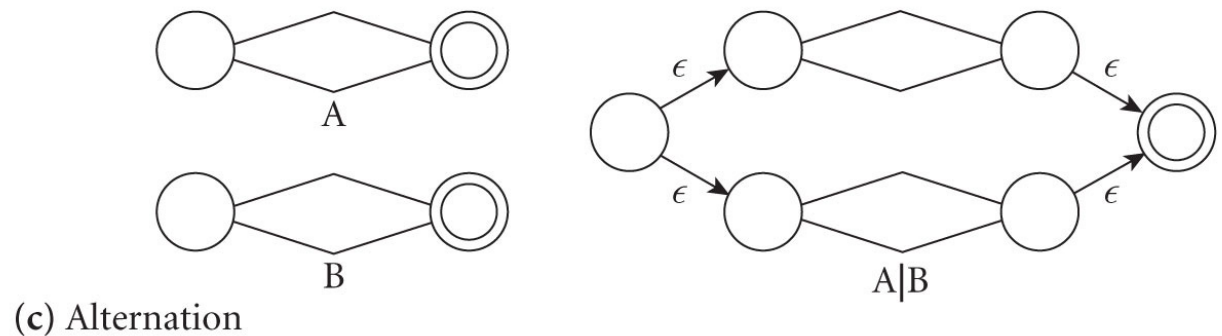
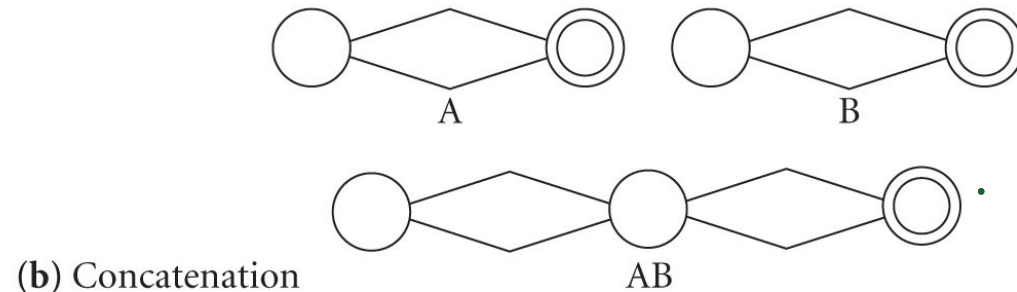
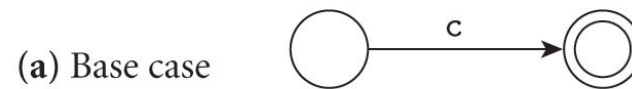
# Scanning

- DFA
  - Built automatically from regular expressions
  - Tools: `lex`, `flex`, `scagen`
  - Difficult to build directly
    - build first an NFA – Nondeterministic FA
    - convert to DFA
    - minimize DFA (smallest number of states)



# Scanning

- Reg.exp. to NFA
- Follows the structural definition of regular expressions



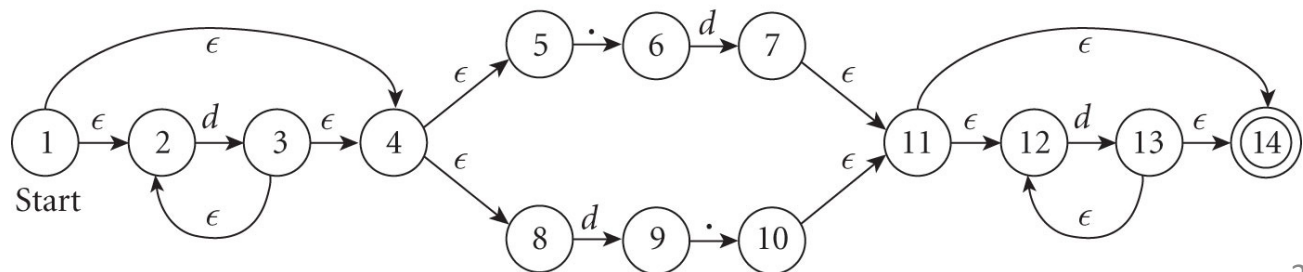
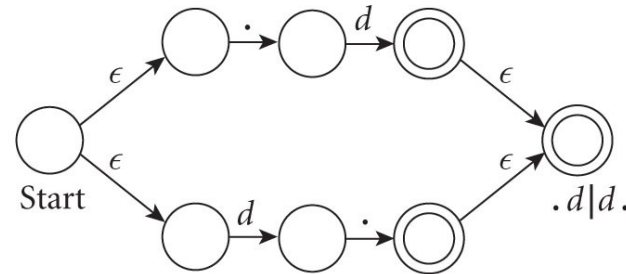
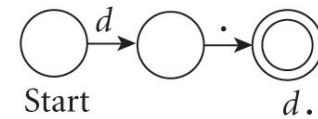
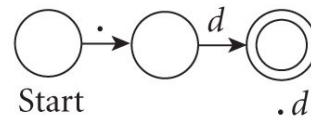
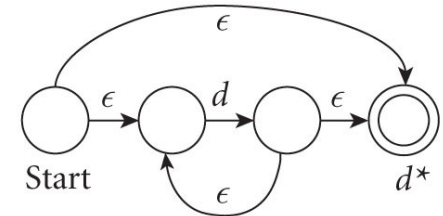
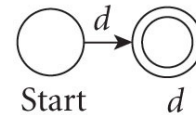
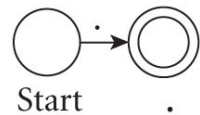


# Scanning

- Reg.exp. to NFA

- Example:

$d^* ( \cdot d \mid d \cdot ) d^*$

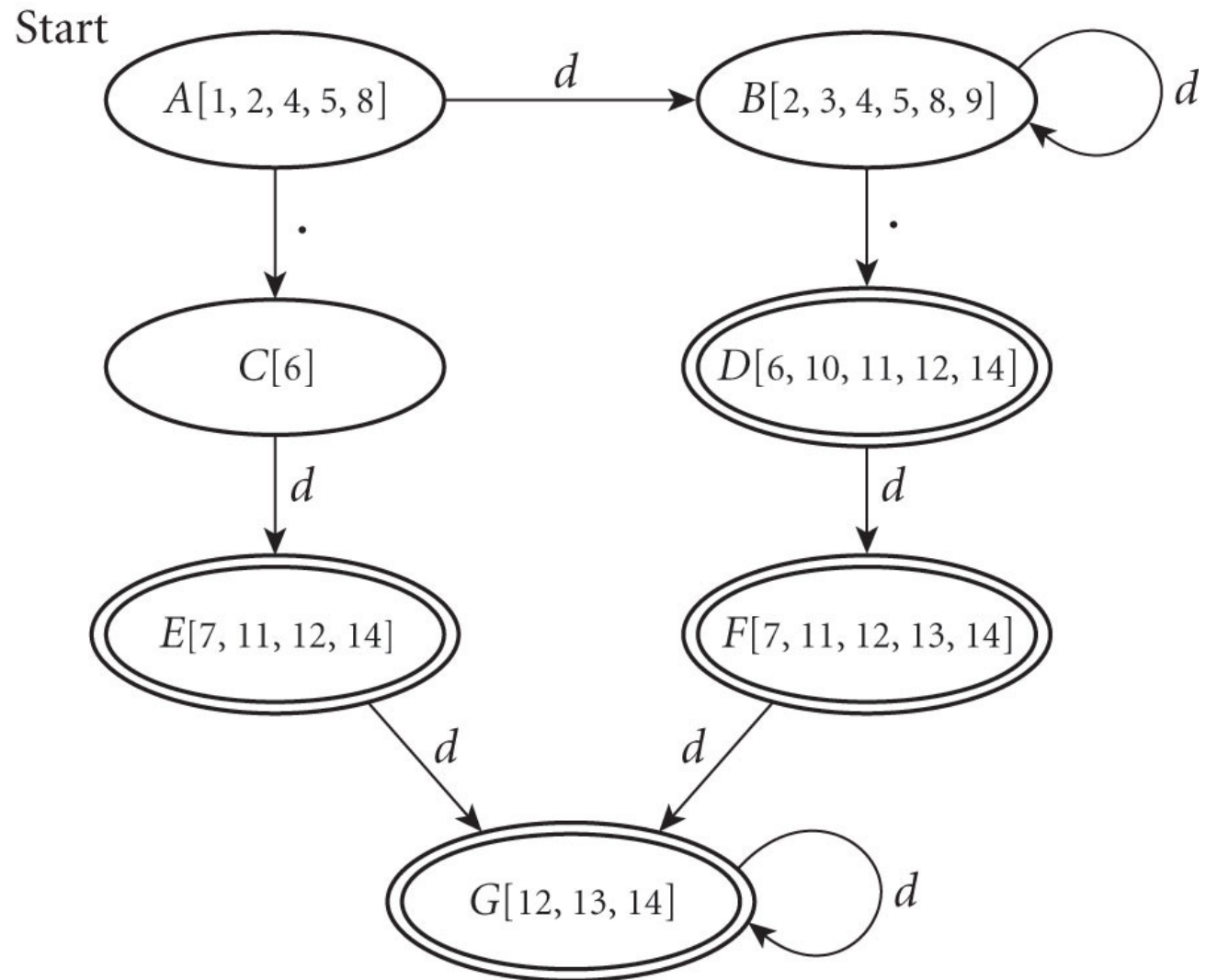


# Scanning

- NFA to DFA

- Example:

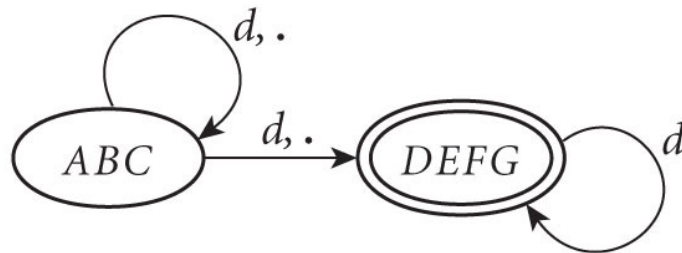
$d^* ( \cdot d \mid d \cdot ) d^*$



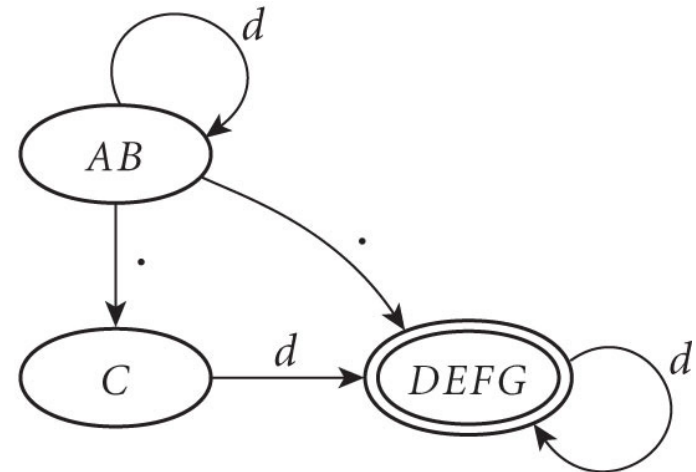
# Scanning

- DFA minimization
- Example:  $d^* ( \cdot d \mid d \cdot ) d^*$

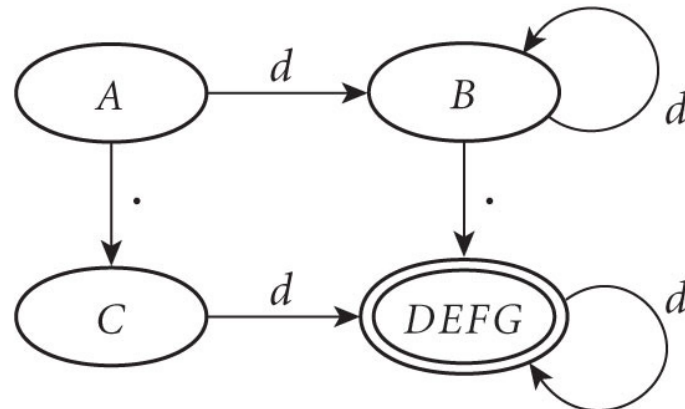
(a)  
Start



(b)  
Start



(c)  
Start



# Scanning

- Scanners are built three ways:
  - ad-hoc:
    - fastest, most compact code
  - semi-mechanical pure DFA
    - nested **case** statements
  - table-driven DFA
    - automatically-generated scanners
- “Longest-possible token” rule
  - return only when next character cannot continue current token
  - the next character needs to be saved for the next token
- Keywords
  - DFA would need many states to identify
  - Better treat keywords as exceptions to the identifier rule



# Scanning

## ■ Nested case statement DFA

```
state := 1
loop
  read cur_char
  case state of
    1: case cur_char of
      ‘ ‘, ‘\t’, ‘\n’:    ...
      ‘a’ ... ‘z’:        ...
      ‘0’ ... ‘9’:        ...
      ‘>’:                ...
      ...
    2: case cur_char of
      ...
    :
    n: case cur_char of
      ...
```

# Scanning

- Look-ahead
- May need to peek at more than one character
- *look-ahead* – characters necessary to decide
- Example: Pascal
  - have 3 so far and see ‘.’
  - 3.14 or 3..5 may follow
- Example: Fortran
  - arbitrarily long look-ahead
  - DO 5 I = 1,25
    - execute statements up to 5 for I from 1 to 25
  - DO 5 I = 1.25
    - assign 1.25 to the variable DO5I
  - NASA’s Mariner 1 may have been lost due to ‘.’ i.o ‘,’
  - Fortran 77 has better syntax: DO 5,I = 1,25

# Scanning

## ■ Table-driven scanning

```
state = 0 .. number_of_states
token = 0 .. number_of_tokens
scan_tab : array [char, state] of record
    action : (move, recognize, error)
    new_state : state
token_tab : array [state] of token      -- what to recognize
keyword_tab : set of record
    k_image : string
    k_token : token
-- these three tables are created by a scanner generator tool
```

(continued on next slide)

# Scanning

## Table-driven scanning (cont'd)

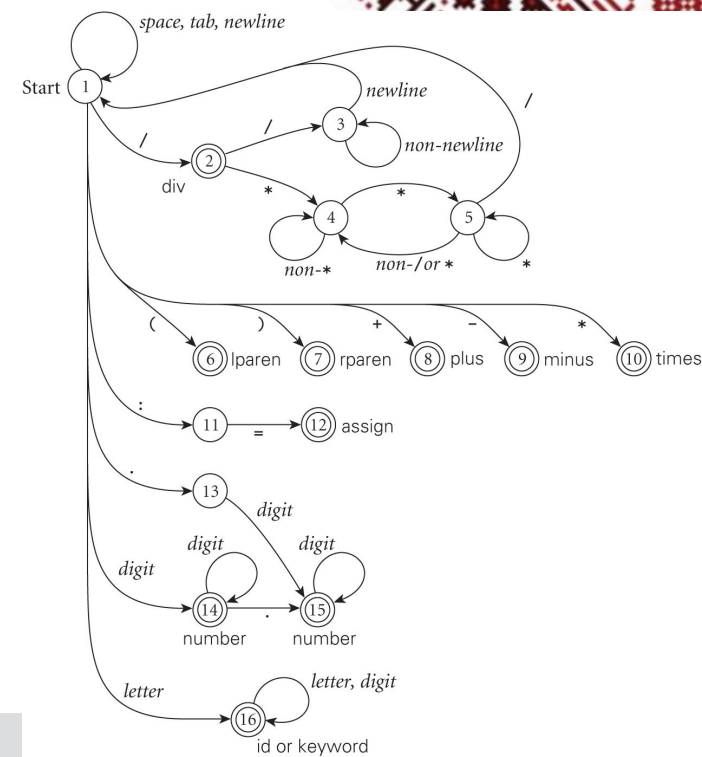
```
state = 0 .. number_of_states
token = 0 .. number_of_tokens
scan_tab : array [char, state] of record
    action : (move, recognize, error)
    new_state : state
token_tab : array [state] of token
keyword_tab : set of record
    k_image : string
    k_token : token
```

```
tok : token
cur_char : char
remembered_chars : list of char
repeat
    cur_state : state := start_state
    image : string := null
    remembered_state : state := 0      -- none
    loop
        read cur_char
        case scan_tab[cur_char, cur_state].action
            move:
                if token_tab[cur_state]  $\neq$  0
                    -- this could be a final state
                    remembered_state := cur_state
                    remembered_chars :=  $\epsilon$ 
                    add cur_char to remembered_chars
                    cur_state := scan_tab[cur_char, cur_state].new_state
                recognize:
                    tok := token_tab[cur_state]
                    unread cur_char      -- push back into input stream
                    exit inner loop
            error:
                if remembered_state  $\neq$  0
                    tok := token_tab[remembered_state]
                    unread remembered_chars
                    remove remembered_chars from image
                    exit inner loop
                -- else print error message and recover; probably start over
        append cur_char to image
    -- end inner loop
until tok  $\notin$  {white_space, comment}
look image up in keyword_tab and replace tok with appropriate keyword if found
return (tok, image)
```



# Scanning

- Scanner table used by previous code
  - state 17: white spaces; state 18: comments
  - scan\_tab: entire table but last column
  - token\_tab: last column
  - keyword\_tab = {read, write}



State	Current input character													
	space, tab	newline	/	*	(	)	+	-	:	=	.	digit	letter	
1	17	17	2	10	6	7	8	9	11	-	13	14	16	-
2	-	-	3	4	-	-	-	-	-	-	-	-	-	div
3	3	18	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	5	4	4	4	4	4	4	4	4	4	4
5	4	4	18	5	4	4	4	4	4	4	4	4	4	4
6	-	-	-	-	-	-	-	-	-	-	-	-	-	lparen
7	-	-	-	-	-	-	-	-	-	-	-	-	-	rparen
8	-	-	-	-	-	-	-	-	-	-	-	-	-	plus
9	-	-	-	-	-	-	-	-	-	-	-	-	-	minus
10	-	-	-	-	-	-	-	-	-	-	-	-	-	times
11	-	-	-	-	-	-	-	-	-	12	-	-	-	-
12	-	-	-	-	-	-	-	-	-	-	-	-	-	assign
13	-	-	-	-	-	-	-	-	-	-	-	15	-	-
14	-	-	-	-	-	-	-	-	-	-	15	14	-	number
15	-	-	-	-	-	-	-	-	-	-	-	15	-	number
16	-	-	-	-	-	-	-	-	-	-	-	16	16	identifier
17	17	17	-	-	-	-	-	-	-	-	-	-	-	white_space
18	-	-	-	-	-	-	-	-	-	-	-	-	-	comment

# Scanning

- Lexical errors
- Very few – most strings correspond to some token
- Should recover to enable the compiler to detect more errors
  - throw away the current, invalid, token
  - skip forward to the next possible beginning of a new token
  - restart the scanning algorithm
  - count on the error-recovery mechanism of the parser to cope with a syntactically invalid sequence of tokens