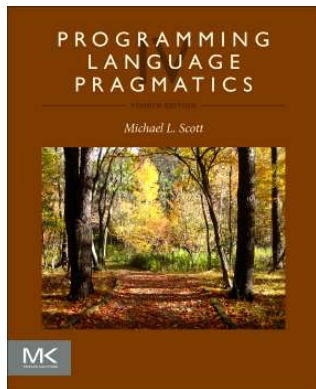


# Syntax and Lexical Analysis

*17-363/17-663: Programming Language Pragmatics*

---



Reading: PLP chapter 2 through section 2.2



Prof. Jonathan Aldrich



# Specifying Syntax

- Let's start by specifying the idea of a *digit*:

$$\textit{digit} \longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

- From this we can build *natural numbers*:

$$\textit{non\_zero\_digit} \longrightarrow 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$
$$\textit{natural\_number} \longrightarrow \textit{non\_zero\_digit} \textit{digit}^*$$

- Simple concepts like these can be expressed with *regular expressions*

# Regular Expressions

- A regular expression is one of the following:
  - A character
  - The empty string, denoted by  $\varepsilon$
  - Two regular expressions concatenated
  - Two regular expressions separated by  $|$  (i.e., or)
  - A regular expression followed by the Kleene star  $*$  (concatenation of zero or more strings)

# Regular Expressions

- Numerical constants accepted by a simple hand-held calculator:

*number*  $\longrightarrow$  *integer* | *real*

*integer*  $\longrightarrow$  *digit* *digit*\*

*real*  $\longrightarrow$  *integer* *exponent* | *decimal* (*exponent* |  $\epsilon$ )

*decimal*  $\longrightarrow$  *digit*\* ( *.* *digit* | *digit* *.* ) *digit*\*

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

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



# Practice with Regular Expressions

- Define a regular expression for C-style comments
  - You may use abbreviations like *non-\** or *newline*
  - You may use Kleene + (1 or more) in addition to Kleene \*

# Practice with Regular Expressions

- Define a regular expression for C-style comments
  - You may use abbreviations like *non-\** or *newline*
  - You may use Kleene + (1 or more) in addition to Kleene \*
- One solution (from the textbook)

*comment* → */\* (non-\* | \* non-/) \*+ /*  
*| // (non-newline)\* newline*

# From Tokens to Grammar

- Regular expressions are great for describing *tokens*
  - The smallest meaningful units of syntax – numbers, symbols, keywords, and identifiers
  - These constructs have no interesting recursive structure
- But real programs have recursive structure, even in expressions like  $2 * (x + (y / 3))$
- To capture higher-level syntax we need *context-free grammars*

# Context-Free Grammars

- A calculator expression grammar is recursive:

$$\text{expr} \longrightarrow \text{id} \mid \text{number} \mid - \text{expr} \mid ( \text{expr} )$$
$$\mid \text{expr op expr}$$
$$\text{op} \longrightarrow + \mid - \mid * \mid /$$

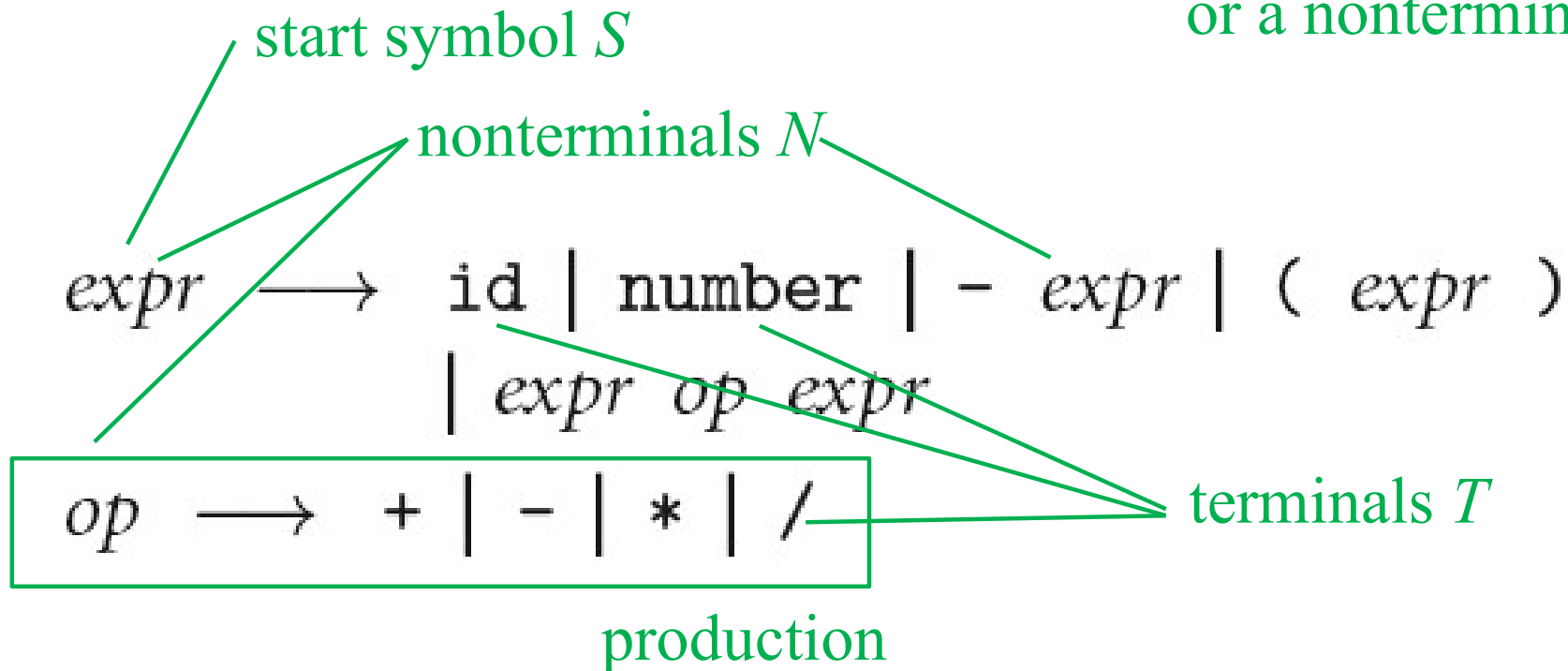
*expr* is defined in terms of itself!



# Context-Free Grammars (CFGs)

- Anatomy of a CFG
  - In Backus-Naur Form (BNF)

*A symbol is a terminal or a nonterminal*



# Context-Free Grammars

- In this grammar,  
generate the string

$$\begin{aligned} \text{expr} &\longrightarrow \text{id} \mid \text{number} \mid - \text{expr} \mid ( \text{expr} ) \\ &\quad \mid \text{expr op expr} \\ \text{op} &\longrightarrow + \mid - \mid * \mid / \end{aligned}$$

**"slope \* x + intercept"**

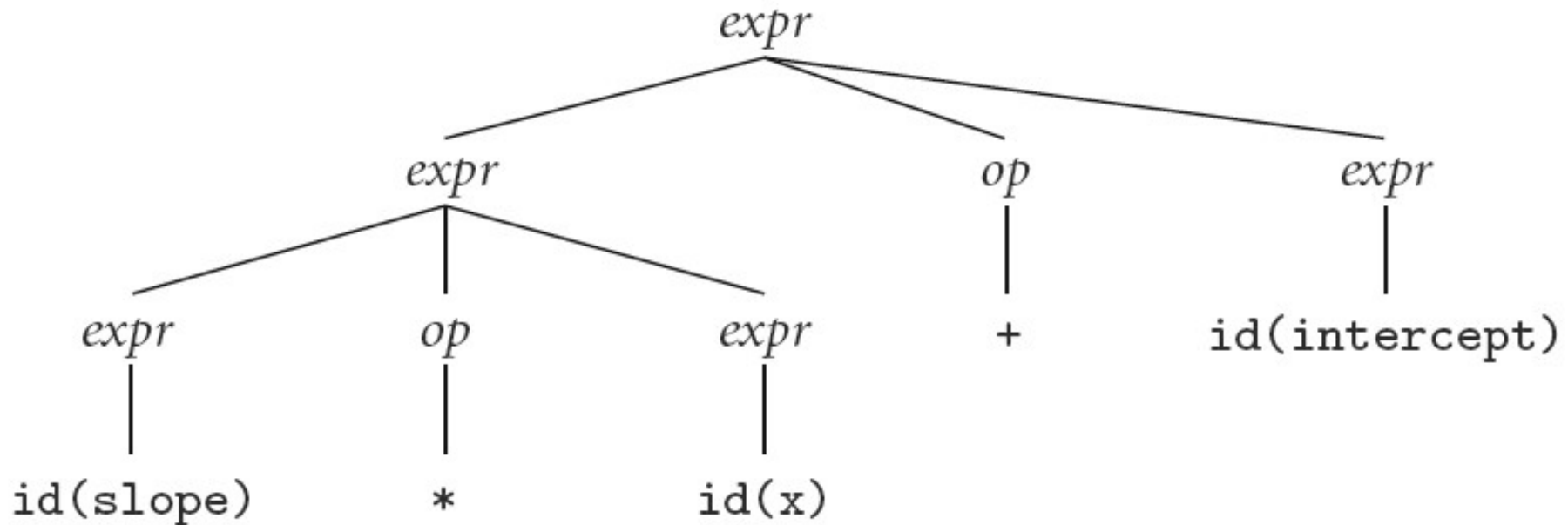
$$\begin{aligned} \text{expr} &\implies \text{expr op } \underline{\text{expr}} \\ &\implies \text{expr } \underline{\text{op}} \text{ id} \\ &\implies \underline{\text{expr}} + \text{id} \\ &\implies \text{expr op } \underline{\text{expr}} + \text{id} \\ &\implies \text{expr } \underline{\text{op}} \text{ id} + \text{id} \\ &\implies \underline{\text{expr}} * \text{id} + \text{id} \\ &\implies \text{id} * \text{id} + \text{id} \\ &\quad (\text{slope}) \quad (\text{x}) \quad (\text{intercept}) \end{aligned}$$

This is called a  
*derivation*



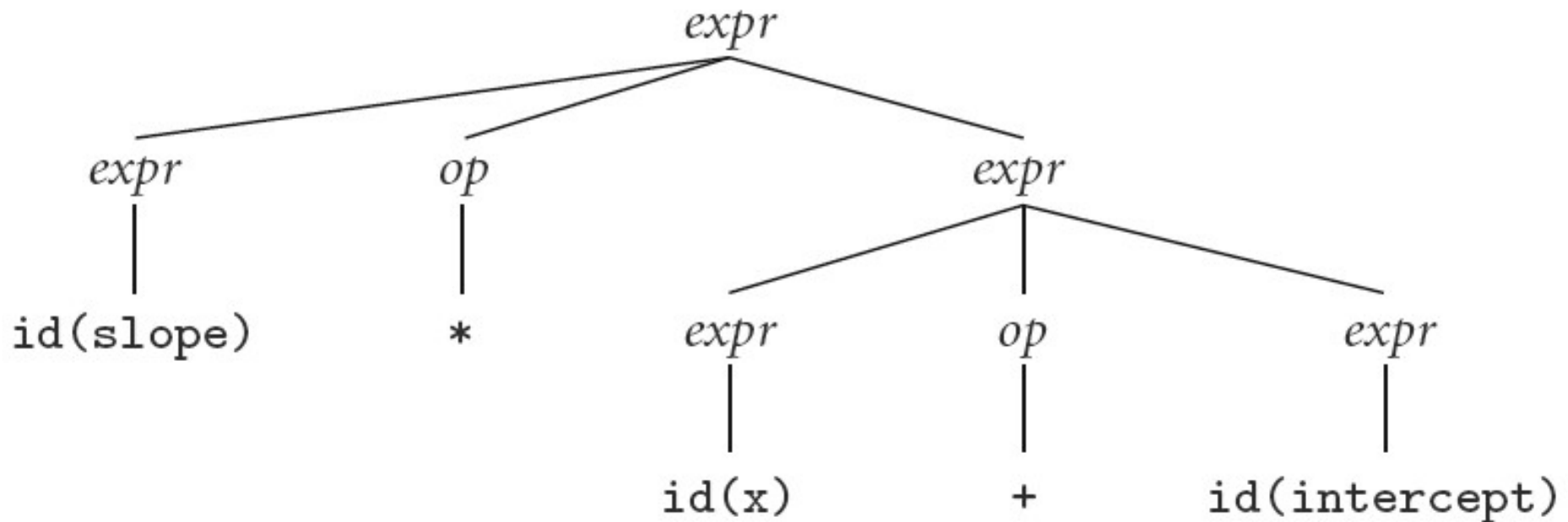
# Context-Free Grammars

- Parse tree for expression grammar for **"slope \* x + intercept"**



# Context-Free Grammars

- Alternate (Incorrect) Parse tree for "slope \* x + intercept"
- Our grammar is *ambiguous*

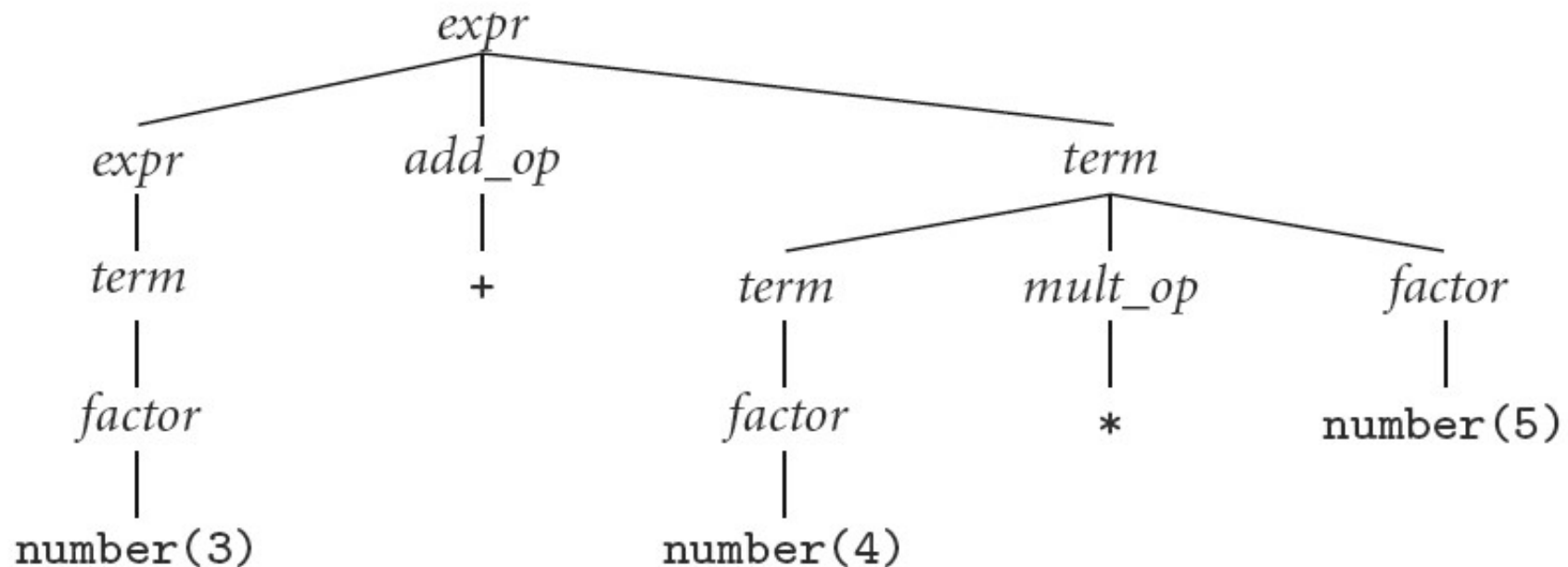


# Context-Free Grammars

- A better version because it is unambiguous and captures precedence
  1.  $expr \longrightarrow term \mid expr \textit{ add\_op } term$
  2.  $term \longrightarrow factor \mid term \textit{ mult\_op } factor$
  3.  $factor \longrightarrow id \mid number \mid - factor \mid ( expr )$
  4.  $add\_op \longrightarrow + \mid -$
  5.  $mult\_op \longrightarrow * \mid /$

# Context-Free Grammars

- Parse tree for expression grammar (with left associativity) for  $3 + 4 * 5$



# Practice with CFGs

- Add  $\&\&$  and  $\|$  to this grammar

- Left-associative

- Precedence:  $+$  over  $\&\&$  over  $\|$

1.  $expr \longrightarrow term \mid expr \textit{ add\_op } term$

2.  $term \longrightarrow factor \mid term \textit{ mult\_op } factor$

3.  $factor \longrightarrow id \mid number \mid - factor \mid ( expr )$

4.  $add\_op \longrightarrow + \mid -$

5.  $mult\_op \longrightarrow * \mid /$



# Practice with CFGs

- One solution

$orexpr \rightarrow andexpr \mid oexpr \parallel andexpr$

$andexpr \rightarrow expr \mid andexpr \&\& expr$

1.  $expr \rightarrow term \mid expr \text{ add\_op } term$

2.  $term \rightarrow factor \mid term \text{ mult\_op } factor$

3.  $factor \rightarrow id \mid number \mid - factor \mid ( expr )$

4.  $add\_op \rightarrow + \mid -$

5.  $mult\_op \rightarrow * \mid /$

Also replace  
with *orexpr*

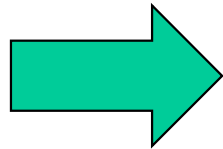




# Lexical Analysis (or “Scanning”)

- Divides source code into tokens
- Removes comments
- Saves text of identifiers, strings, numbers
- Tags tokens with line numbers, for error messages

```
y := x;  
z := 1;  
while y > 1 do  
    z := z * y;  
    y := y - 1  
od
```



```
y := x ; z := 1 ; while y  
> 1 do z := z * y ; y :=  
y - 1 od
```

# Scanning

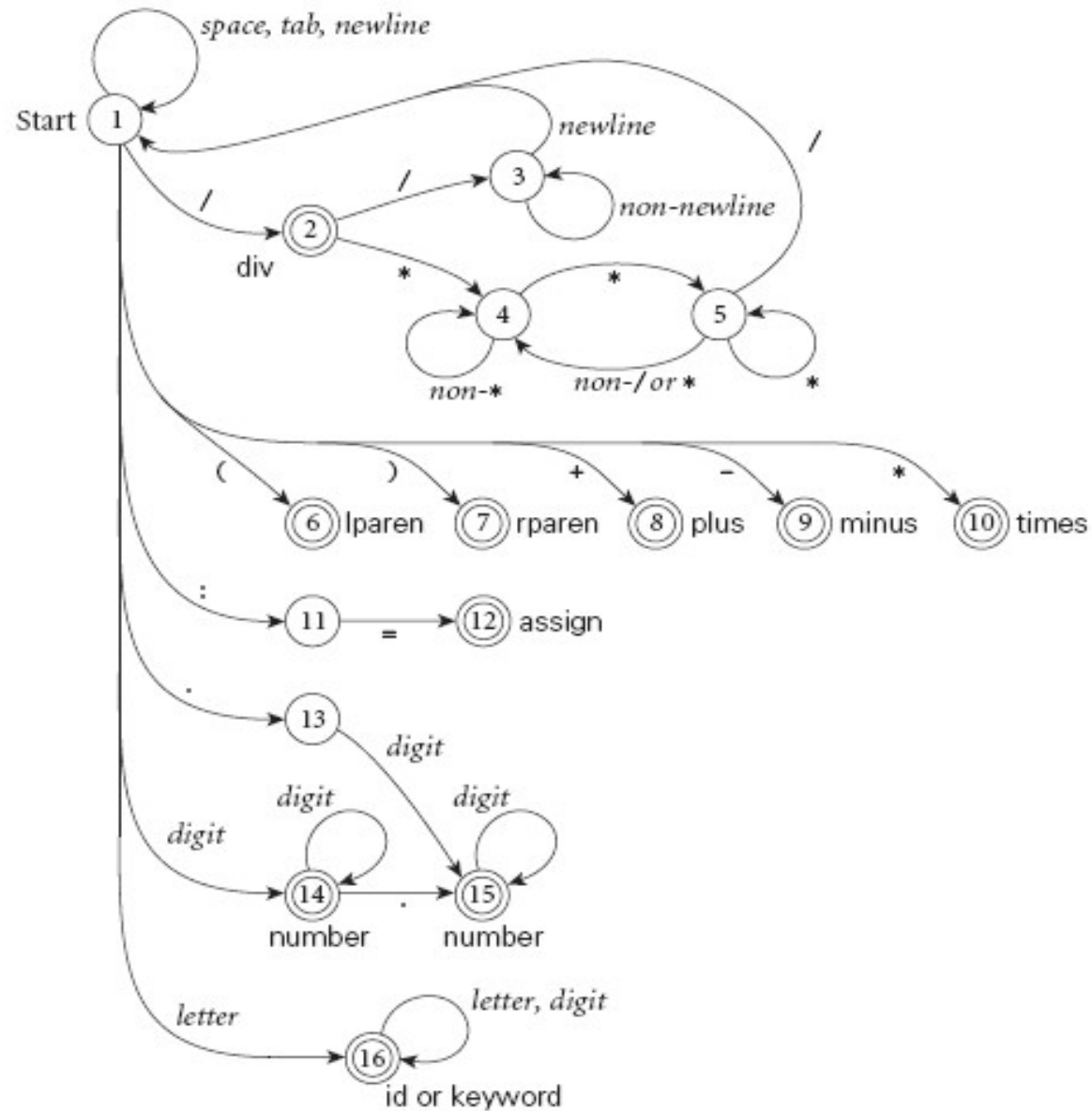
- Suppose we are building an ad-hoc (hand-written) scanner for a calculator language:
  - We read the characters one at a time with look-ahead
- If it is one of the one-character tokens  
( ) + - \* /  
we announce that token
- If it is a digit, we keep reading digits until we can't anymore, then announce a number
- If it is a letter, we keep reading letters and digits and maybe underscores until we can't anymore, then announce an identifier

# Scanning with floating point

- If it is a digit, we keep reading until we find a non-digit
  - if that is not a . we announce an integer
  - otherwise, we keep looking for a real number
  - if the character after the . is not a digit we announce an integer and reuse the . and the look-ahead

# Scanning

- Pictorial representation of a scanner for calculator tokens, in the form of a finite automaton



# Scanning

- This is a deterministic finite automaton (DFA)
  - Lex, scangen, etc. build these things automatically from a set of regular expressions
  - Specifically, they construct a machine that accepts the language

```
identifier | int const
| real const | comment | symbol
| ...
```

# Scanning

- We run the machine over and over to get one token after another
  - Nearly universal rule:
    - always take the longest possible token from the input thus foobar is foobar and never f or foo or foob
    - more to the point, 3.14159 is a real const and never 3, ., and 14159
- Regular expressions "generate" a regular language; DFAs "recognize" it

# Scanning

- Scanners tend to be built three ways
  - ad-hoc
  - semi-mechanical pure DFA  
(usually realized as nested case statements)
  - table-driven DFA
- Ad-hoc generally yields the fastest, most compact code by doing lots of special-purpose things, though good automatically-generated scanners come very close

# Scanning

- Writing a pure DFA as a set of nested case statements is a surprisingly useful programming technique
  - though it's often easier to use perl, awk, sed
  - for details see Example 2.16
- Table-driven DFA is what lex and scangen produce
  - lex/lalrpop in the form of C/Rust code
  - scangen in the form of numeric tables and a separate driver (for details see Figure 2.11-2.12)



# Scanning

- Note that the rule about longest-possible tokens means you return only when the next character can't be used to continue the current token
  - the next character will generally need to be saved for the next token
- In some cases, you may need to peek at more than one character of look-ahead in order to know whether to proceed
  - In Pascal, for example, when you have a 3 and you see a dot
    - do you proceed (in hopes of getting 3.14)?  
or
    - do you stop (in fear of getting 3..5)?

# Scanning

- In messier cases, you may not be able to get by with any fixed amount of look-ahead. In Fortran, for example, we have

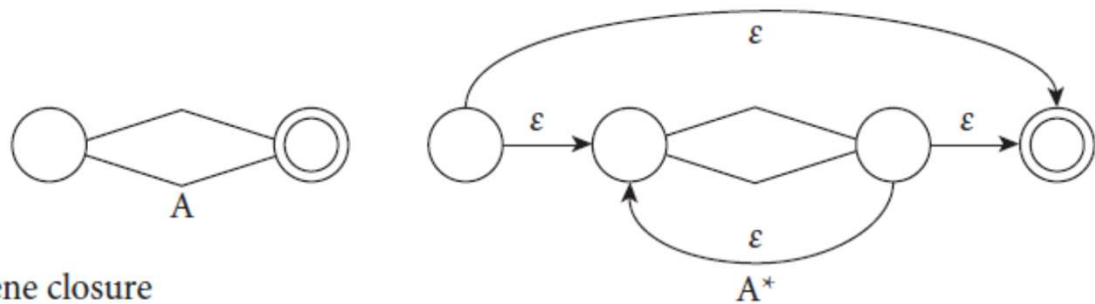
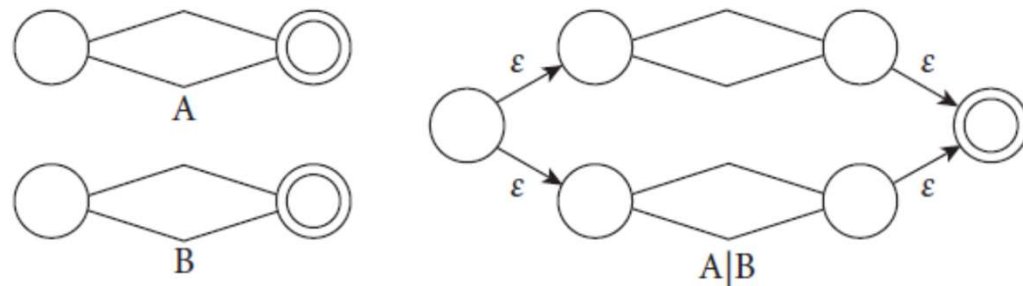
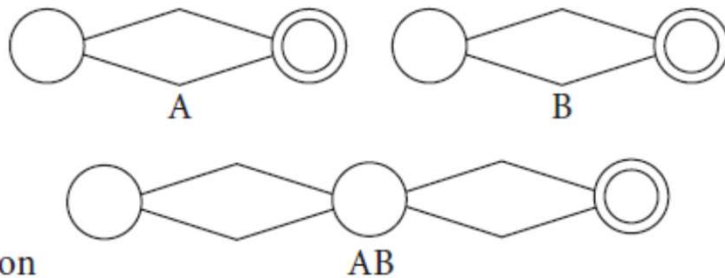
```
DO 5 I = 1, 25    loop
DO 5 I = 1.25    assignment
                  (to DO5I)
```

- Here, we need to remember we were in a potentially final state, and save enough information that we can back up to it, if we get stuck later

# Converting a RE to a DFA

1. Write regular expressions for each construct
  - Except keywords – special case of identifiers
2. Construct NFA from REs
3. Convert NFA to a DFA (set of subsets)
4. Minimize DFA (find equivalence classes)
5. Fix up the result
  - Longest-possible token rule
  - Discard whitespace and comments
  - Distinguish keywords from identifiers
  - Save text, token location
  - Return a special EOF token at end of file

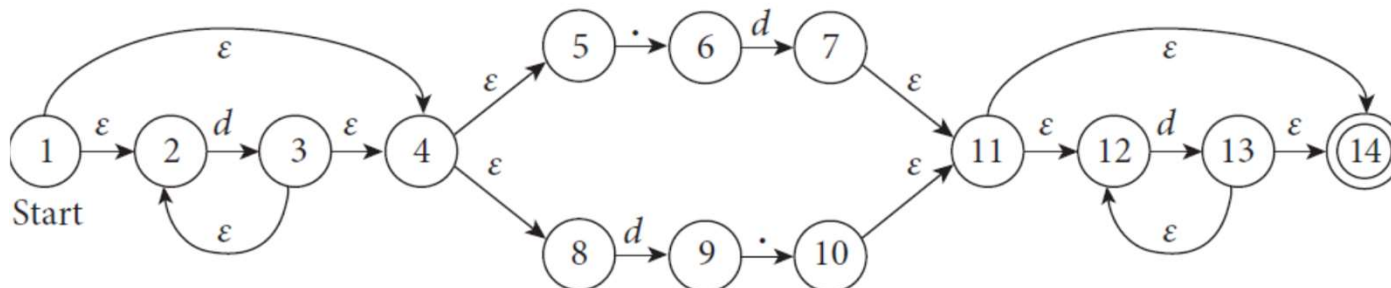
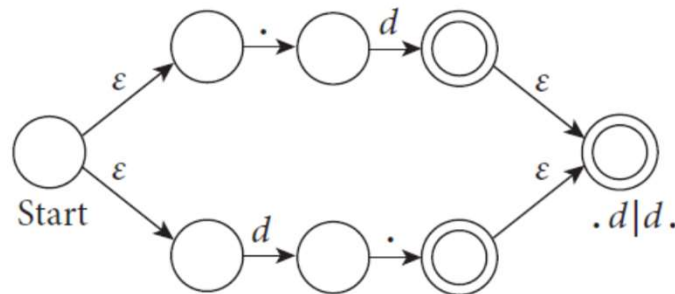
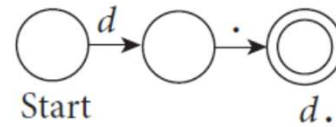
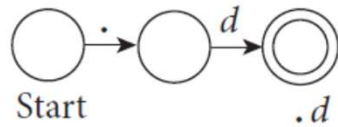
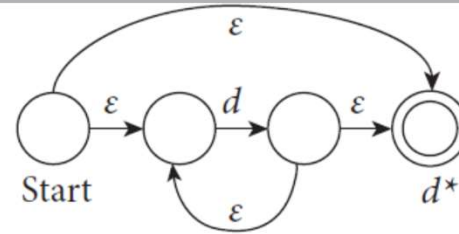
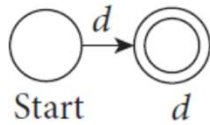
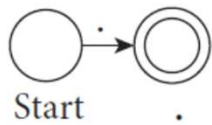
# RE to NFA Construction



Let's apply this to

$d^* (.d | d.) d^*$

# RE to NFA Construction

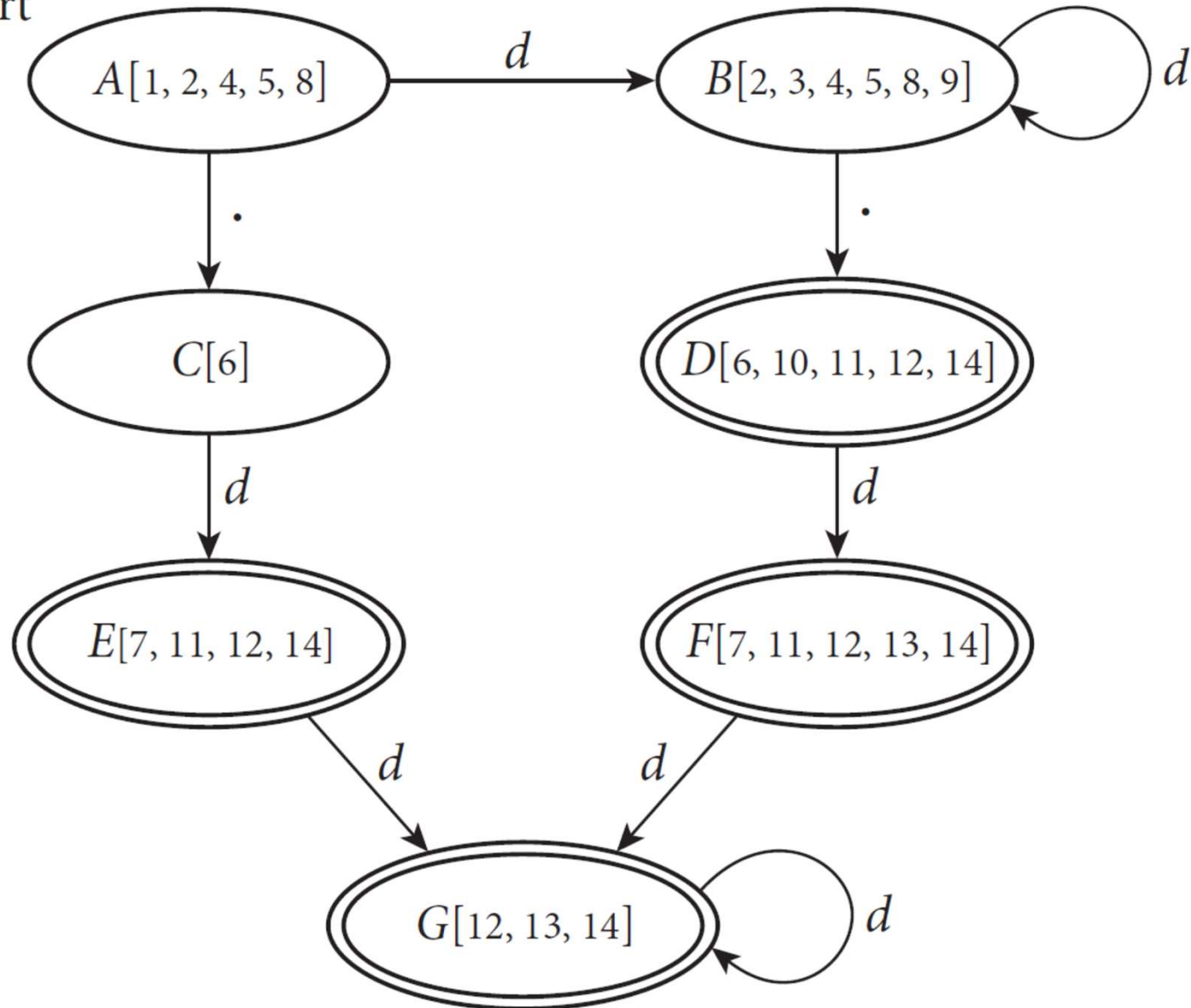


# NFA to DFA Construction

- Each state in the DFA is a set of NFA states
  - “Set of subsets”
- The start DFA state contains the start NFA state, plus all states reachable through  $\epsilon$ -transitions
- For each input that can be consumed from one of those NFA states, we create another DFA state with the set of destination states (plus states from  $\epsilon$ -transitions)

# NFA to DFA Construction (example)

Start

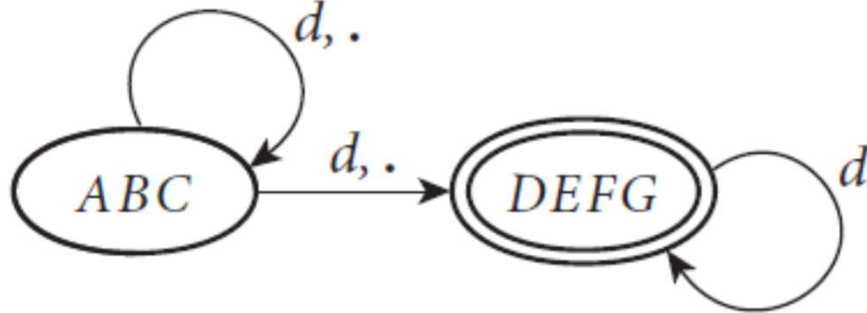


# DFA Minimization

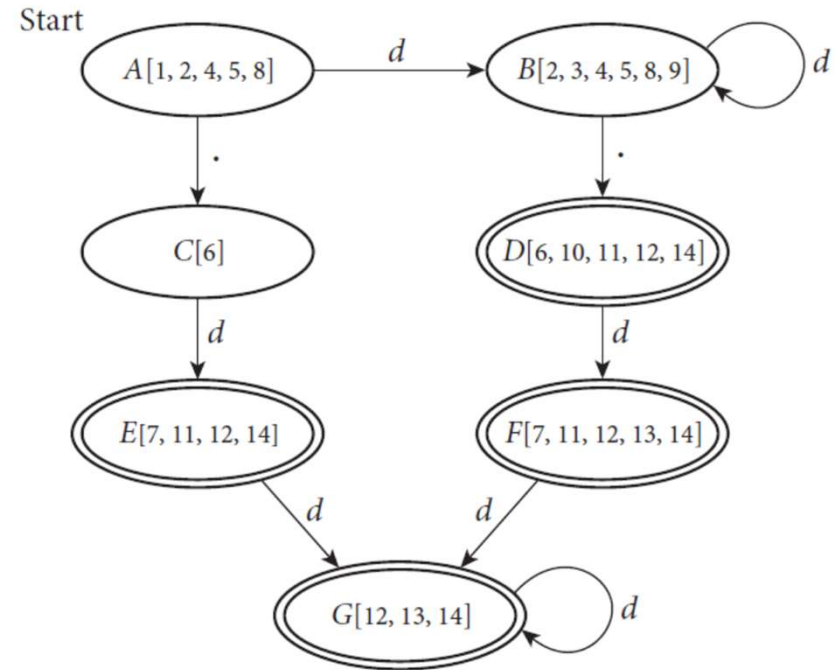
- Start by merging all DFA states into two equivalence classes: final and non-final
- Iteratively split the equivalence classes as follows:
  - Identify a class  $X$  with a nondeterministic transitions on some character  $c$  that lead to  $k$  other equivalence classes
  - Split  $X$  into  $k$  equivalence classes according to the destinations of these transitions



# DFA Minimization

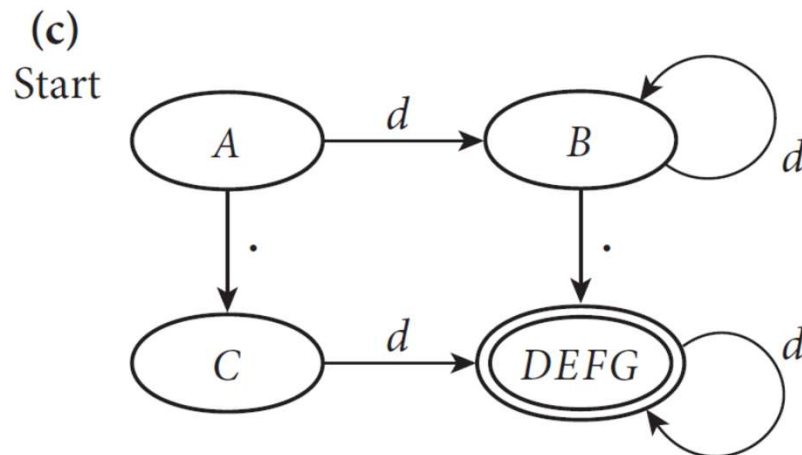
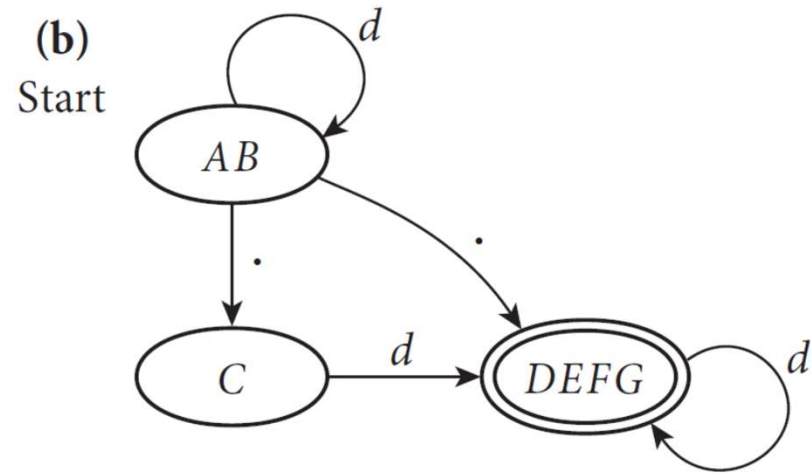
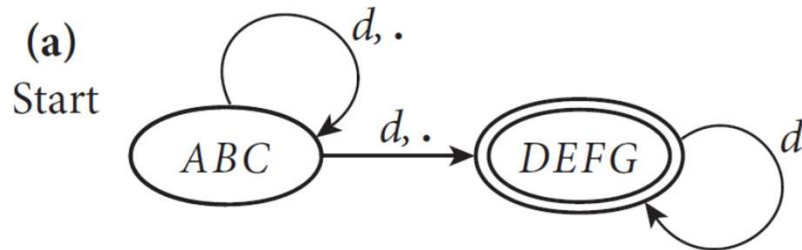


- Example: Consider the diagram on the left, derived by merging states from the one on the upper right.
- The transition from  $ABC$  on  $d$  is nondeterministic
- So we split  $ABC$  into two parts:
  - $AB$  (where  $d$  transitions back to  $ABC$ )
  - $C$  (where  $d$  transitions to  $DEFG$ )
- The result is (b) on the next page



# DFA Minimization (example)

- From state (b) we can now make the  $\cdot$  transition deterministic by splitting  $AB$  into  $A$  and  $B$ . Now we have a DFA.



# Syntax and Lexical Analysis

- We use regular expressions to define tokens
  - Concatenation, alternation, repetition
- A scanner uses a DFA to recognize tokens
  - Often the DFA is machine-generated
- Context-free grammars define higher-level structure
  - Must structure the right way to avoid ambiguity
  - Interesting parsing challenges – future lecture!