

한양대학교 컴퓨터공학부  
컴파일러  
2014년 2학기

# Scanner Construction



# Regular expression → NFA



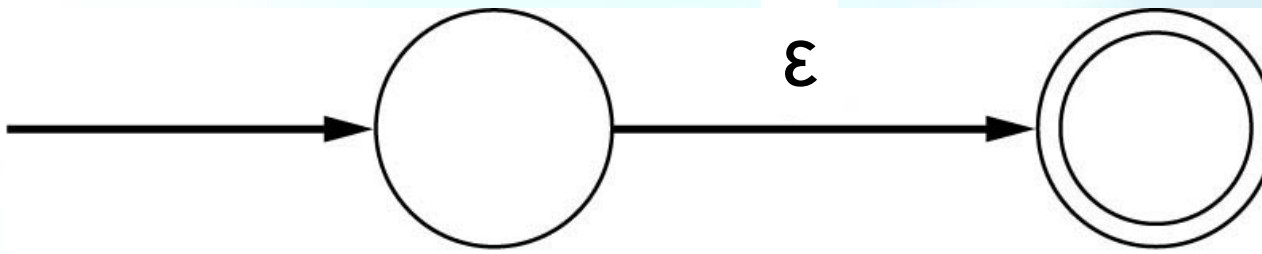
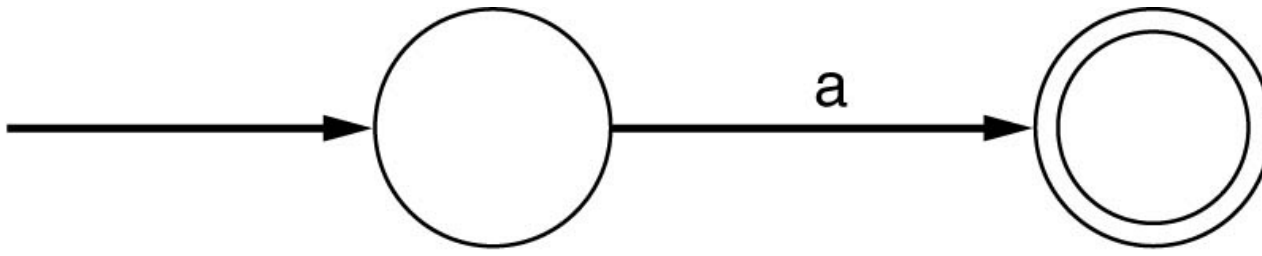
<http://usecurity.hanyang.ac.kr>

- **Thompson's construction**
- Regular expressions
  - Basic regular expressions
  - Concatenation
  - Choice
  - Repetition



# Regular expression $\rightarrow$ NFA

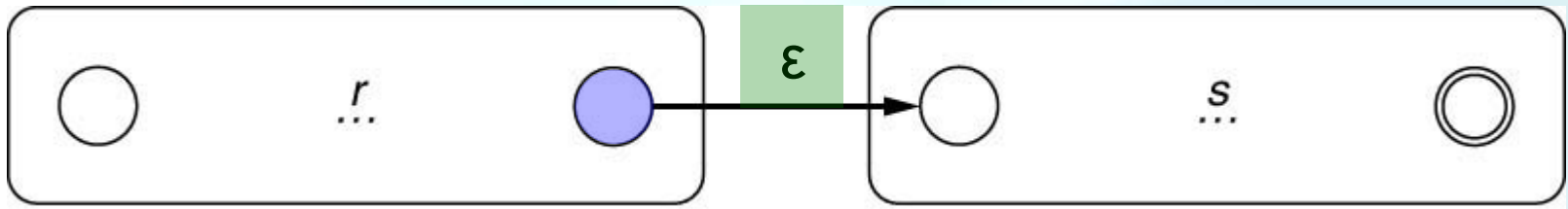
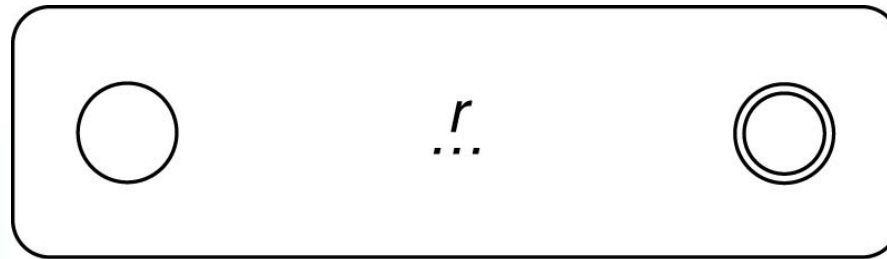
- Basic regular expressions
  - $a$ ,  $\epsilon$ ,  $\Phi$



# Regular expression $\rightarrow$ NFA

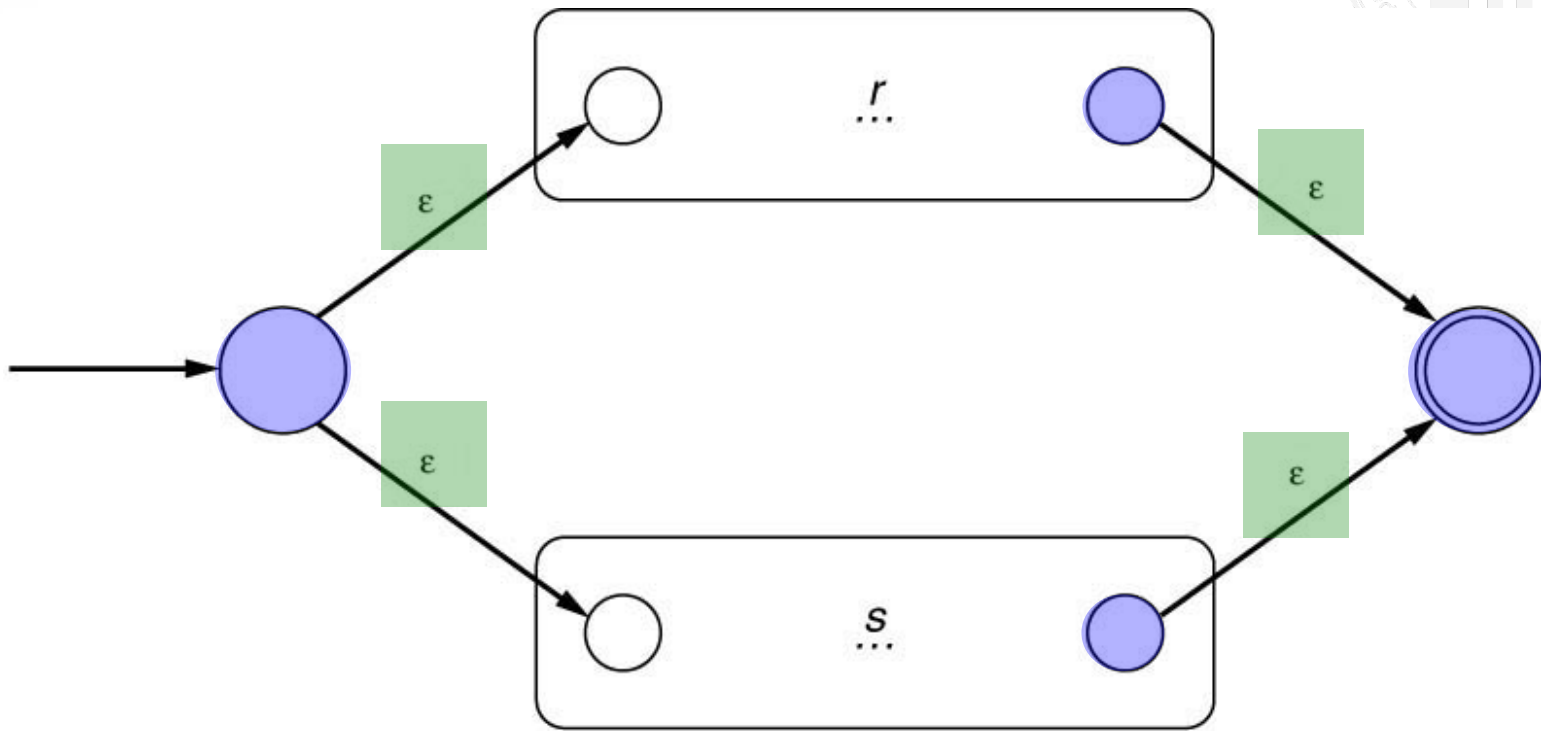
- Concatenation

- $rs$



# Regular expression $\rightarrow$ NFA

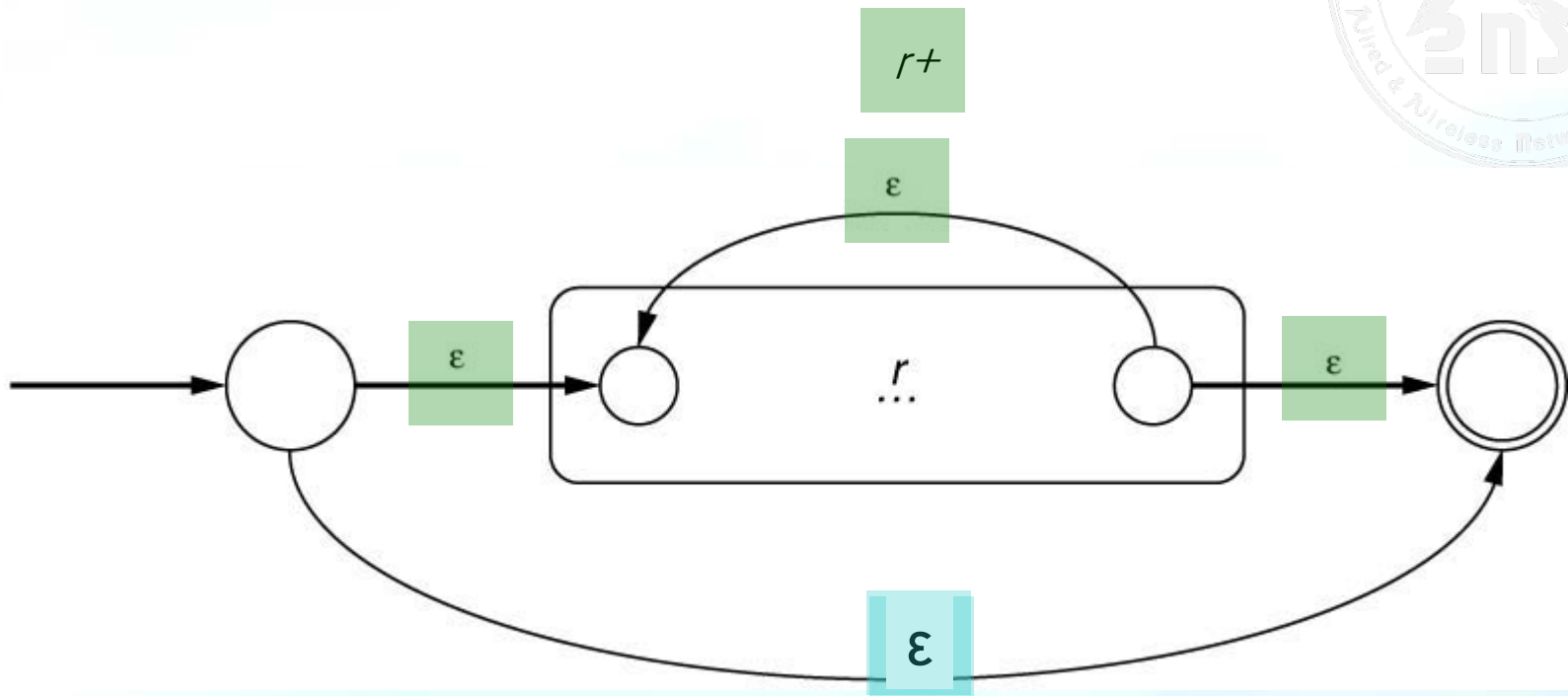
- Choice
  - $r|s$



# Regular expression $\rightarrow$ NFA

- Repetition

- $r^*$

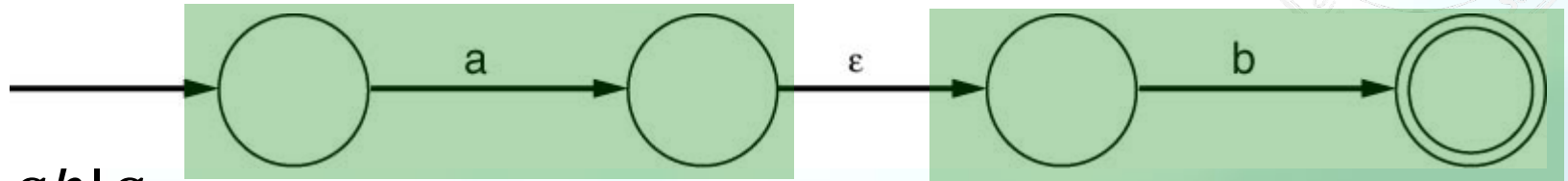
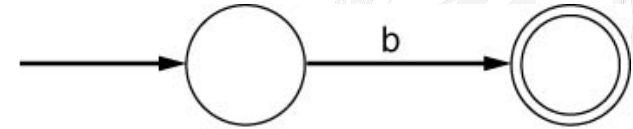
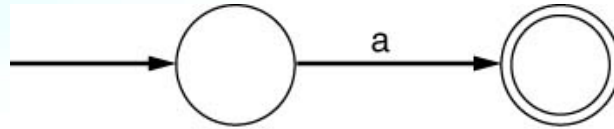


# Regular expression $\rightarrow$ NFA

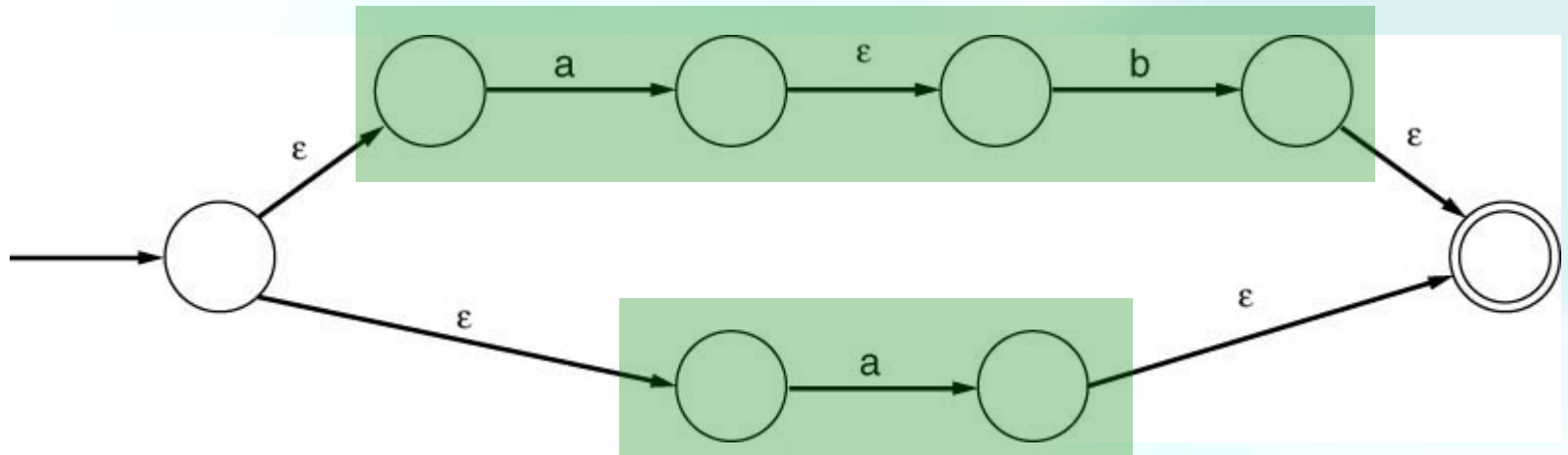
- $ab|a$

- $a, b$

- $ab$



- $ab|a$

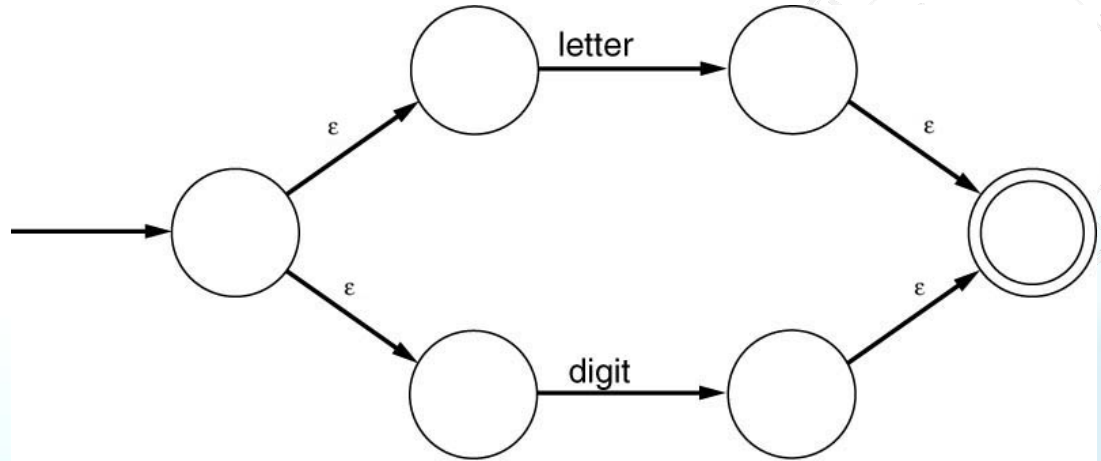




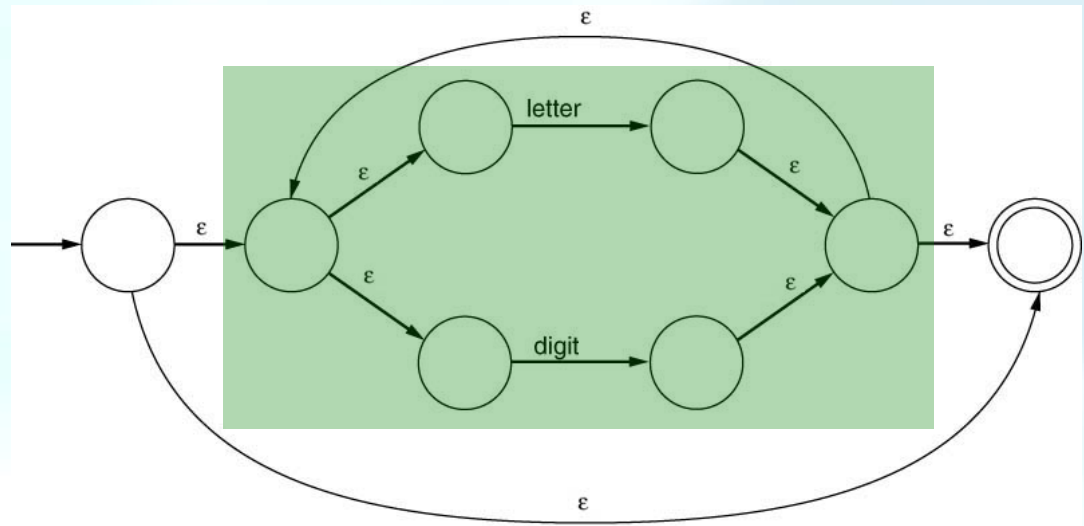
# Regular expression $\rightarrow$ NFA

- $letter(letter/digit)^*$

•  $letter/digit$

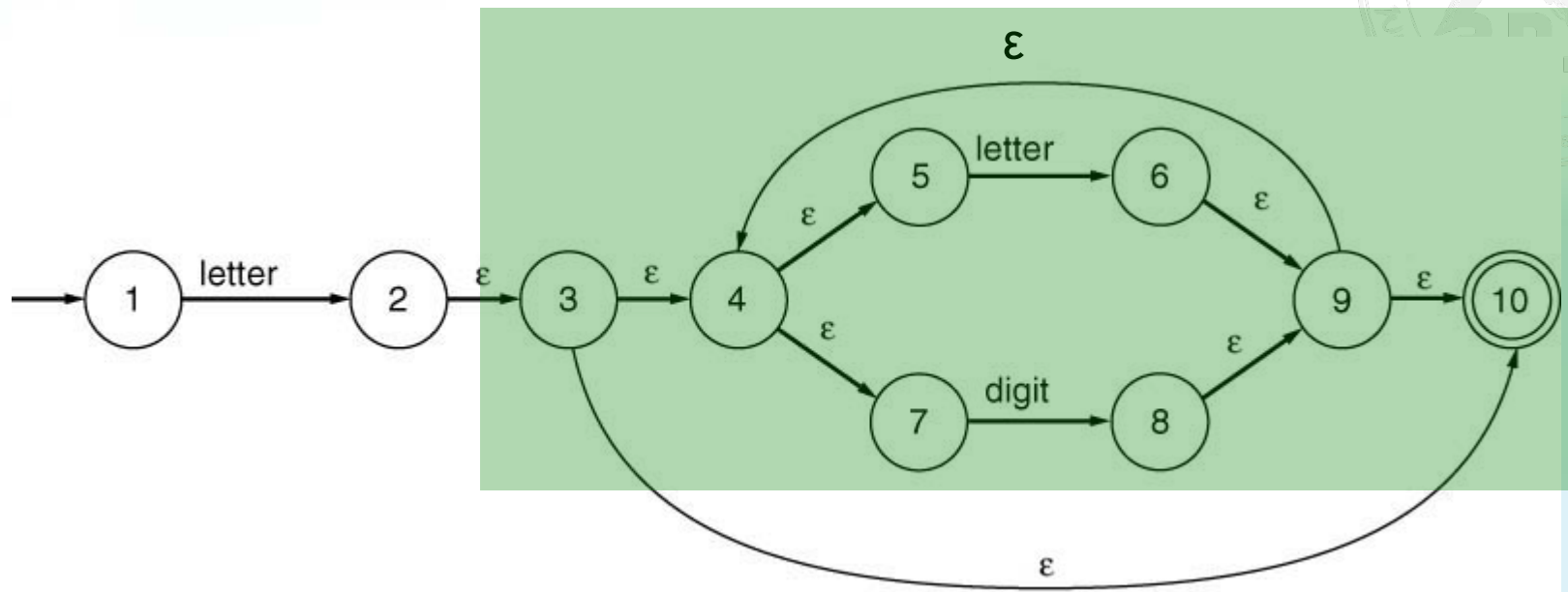


•  $(letter/digit)^*$



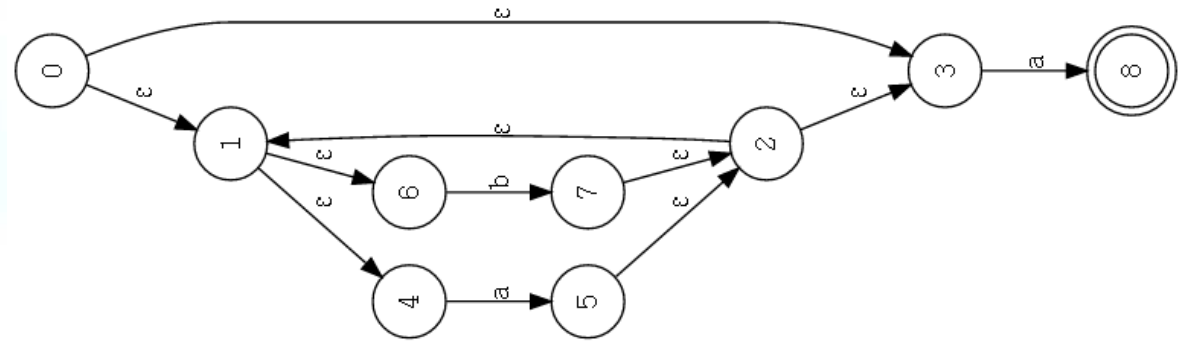
# Regular expression $\rightarrow$ NFA

- $letter(letter/digit)^*$

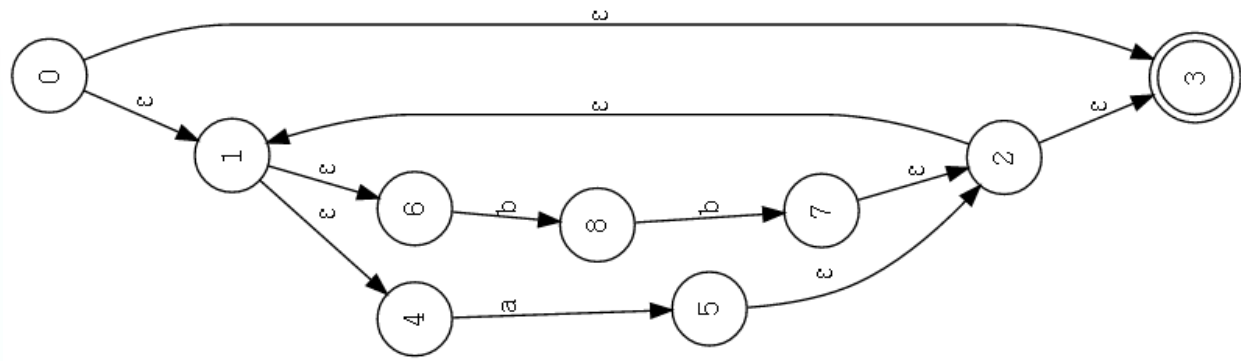


# Examples

- $(a \mid b)^* a$

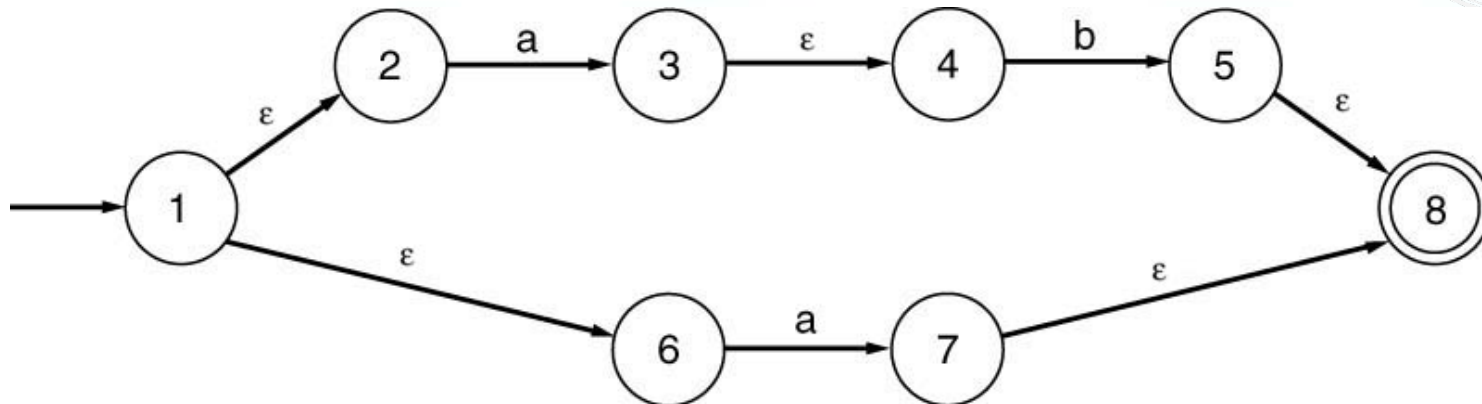


- $(a \mid bb)^*$



## ● The **subset construction** with $\epsilon$ – transitions.

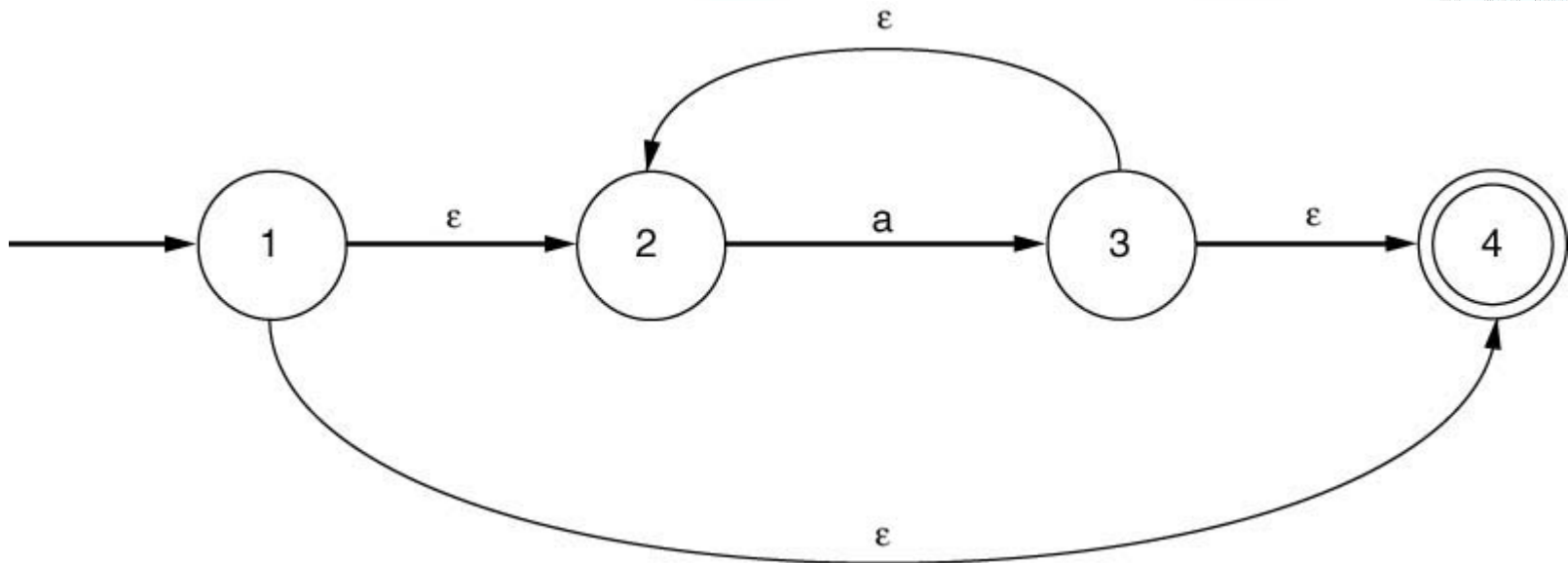
1. eliminating  $\epsilon$ -transition
2. eliminating multiple transitions on a single input character



# NFA $\rightarrow$ DFA

- $\epsilon$  – closure of a state

- The  $\epsilon$  – closure of a single state  $s$ , denoted by  $\overline{s}$ , is the set of states reachable by only zero or more  $\epsilon$ -transitions.



$$\overline{1} = \{1, 2, 4\}$$

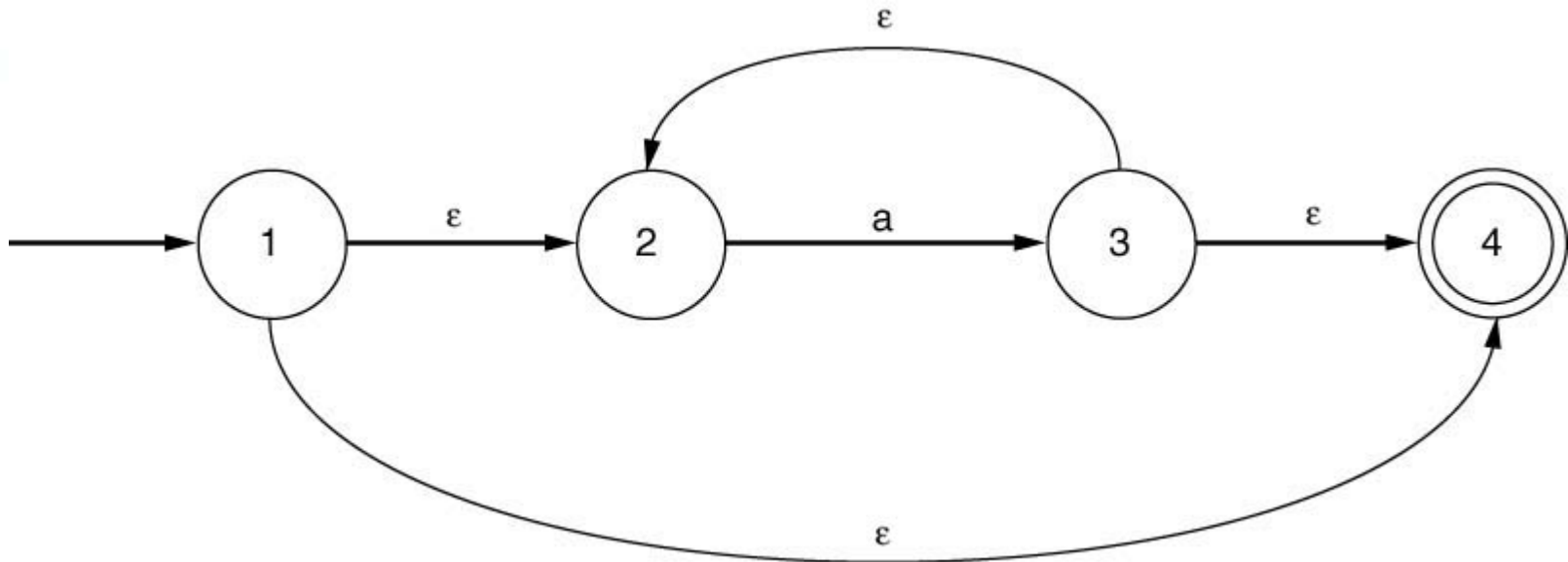
$$\overline{2} = \{2\}$$

$$\overline{3} = \{2, 3, 4\}$$

$$\overline{4} = \{4\}$$

# NFA $\rightarrow$ DFA

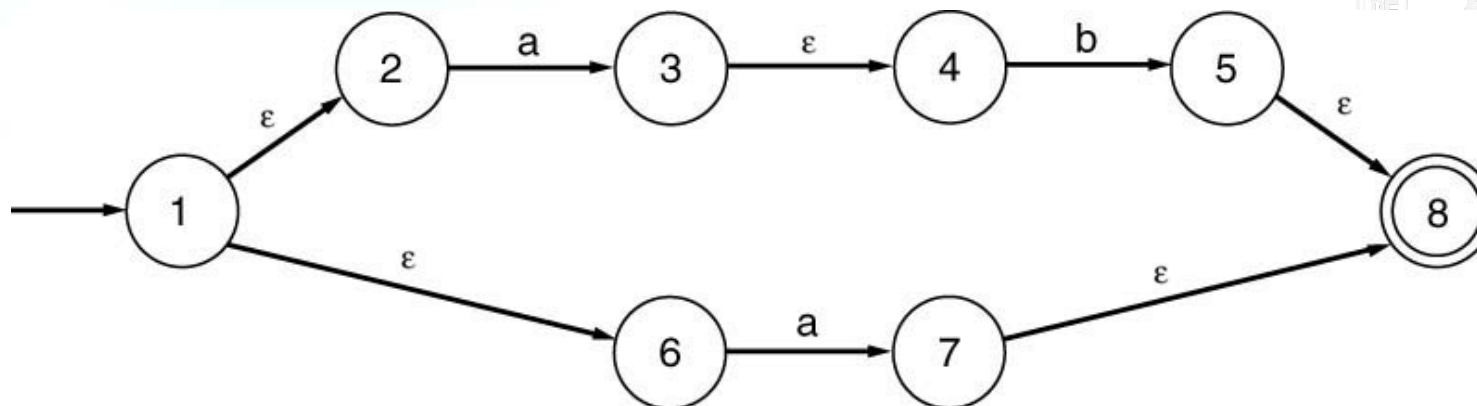
- $\epsilon$  - closure of some states
  - The union of the  $\epsilon$  - closures of each state.



$$\overline{\{1,3\}} = \overline{1} \cup \overline{3} = \{1,2,4\} \cup \{2,3,4\} = \{1,2,3,4\}$$

# NFA → DFA

## ● The Subset construction



string: **ab**

$$\overline{1} = \{1, 2, 6\}, \overline{2} = \{2\}, \overline{3} = \{3, 4\}, \overline{4} = \{4\}, \overline{5} = \{5, 8\}, \overline{6} = \{6\}, \overline{7} = \{7, 8\}, \overline{8} = \{8\}$$

$$\overline{\overline{1}} = \{1, 2, 6\} \xrightarrow[3. \text{"a"}]{a} \overline{\overline{3, 7}} = \{3, 4, 7, 8\} \xrightarrow[6. \text{"b"}]{b} \overline{\overline{5}} = \{5, 8\}$$

1. start state  
ε-closure

2. {1,2,6}  
symbol transition

4. "a"  
state

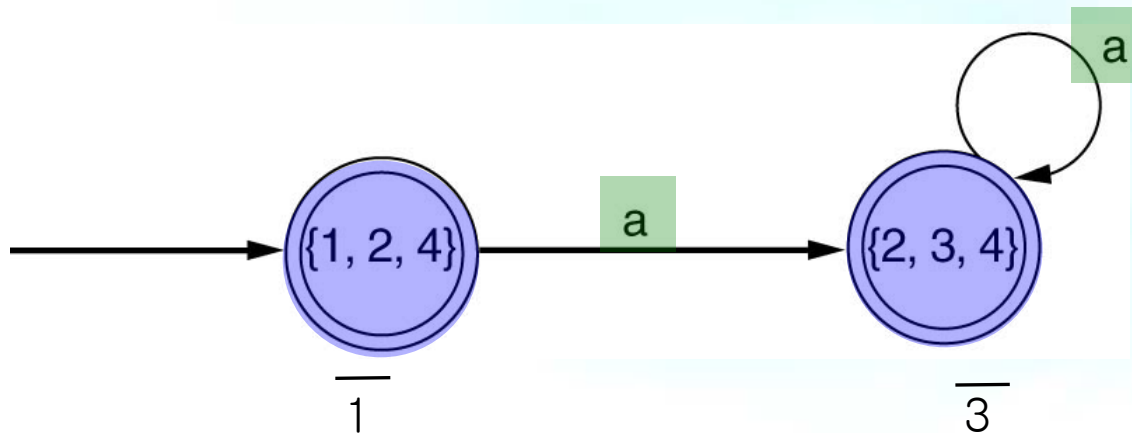
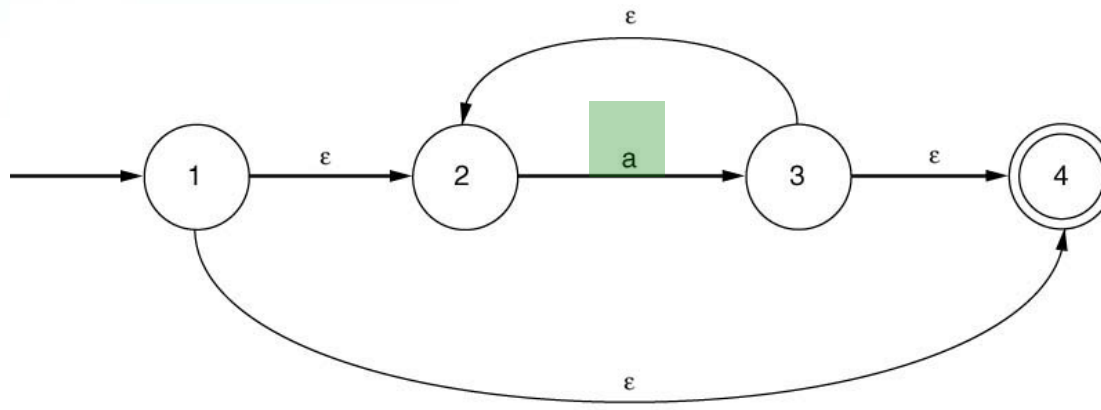
5. {3,4,7,8}  
symbol transition

7. "b"  
state

8. accept state

# NFA $\rightarrow$ DFA

- NFA  $\rightarrow$  DFA ( $a^*$ )



$$\overline{1} = \{1, 2, 4\}$$

$$\overline{2} = \{2\}$$

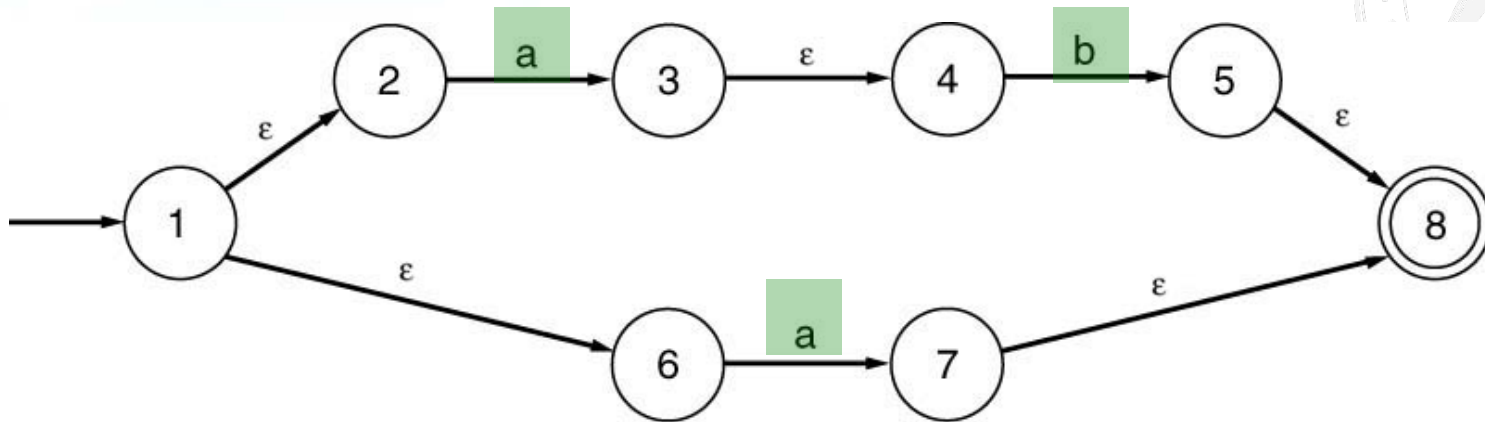
$$\overline{3} = \{2, 3, 4\}$$

$$\overline{4} = \{4\}$$



# NFA $\rightarrow$ DFA

- $ab|a$

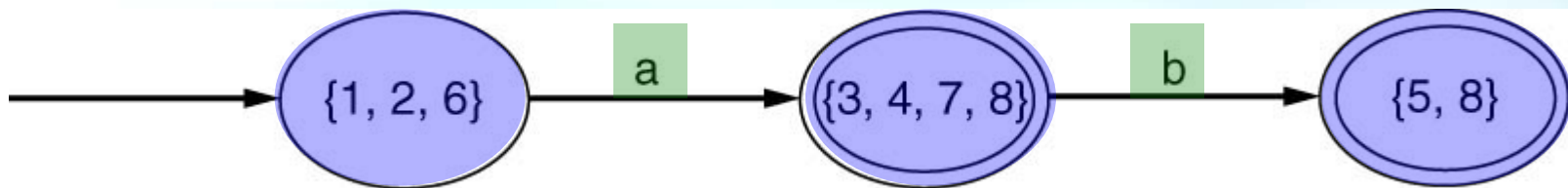


$$\overline{1} = \{1, 2, 6\}, \overline{2} = \{2\}, \overline{3} = \{3, 4\}, \overline{4} = \{4\}, \overline{5} = \{5, 8\}, \overline{6} = \{6\}, \overline{7} = \{7, 8\}, \overline{8} = \{8\}$$

$\overline{1}$

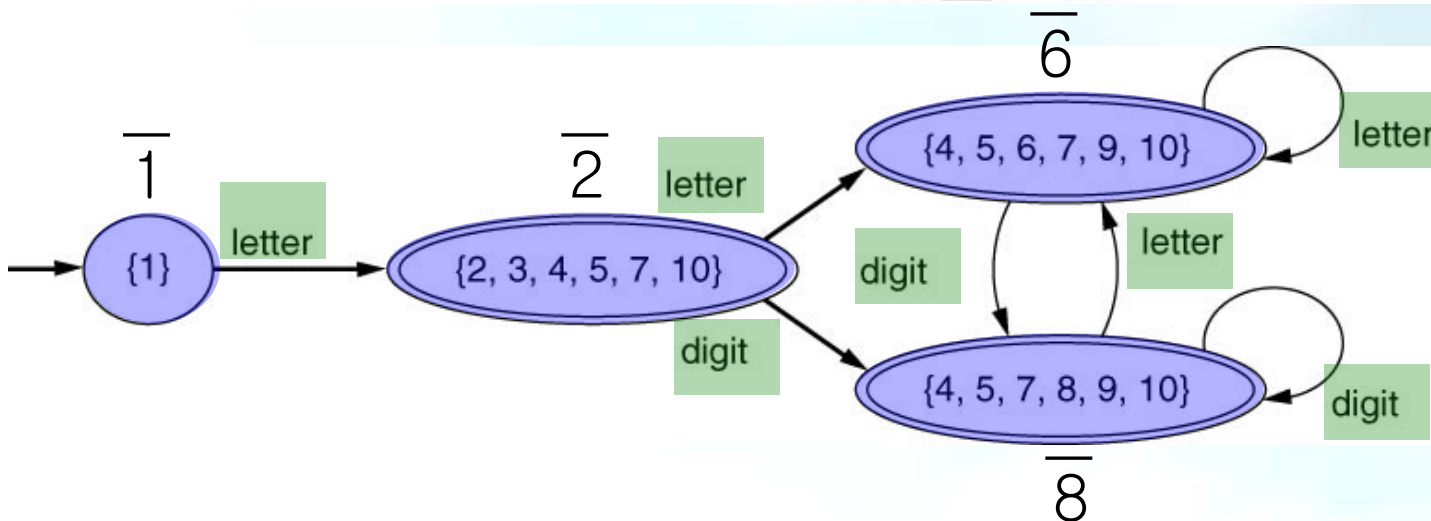
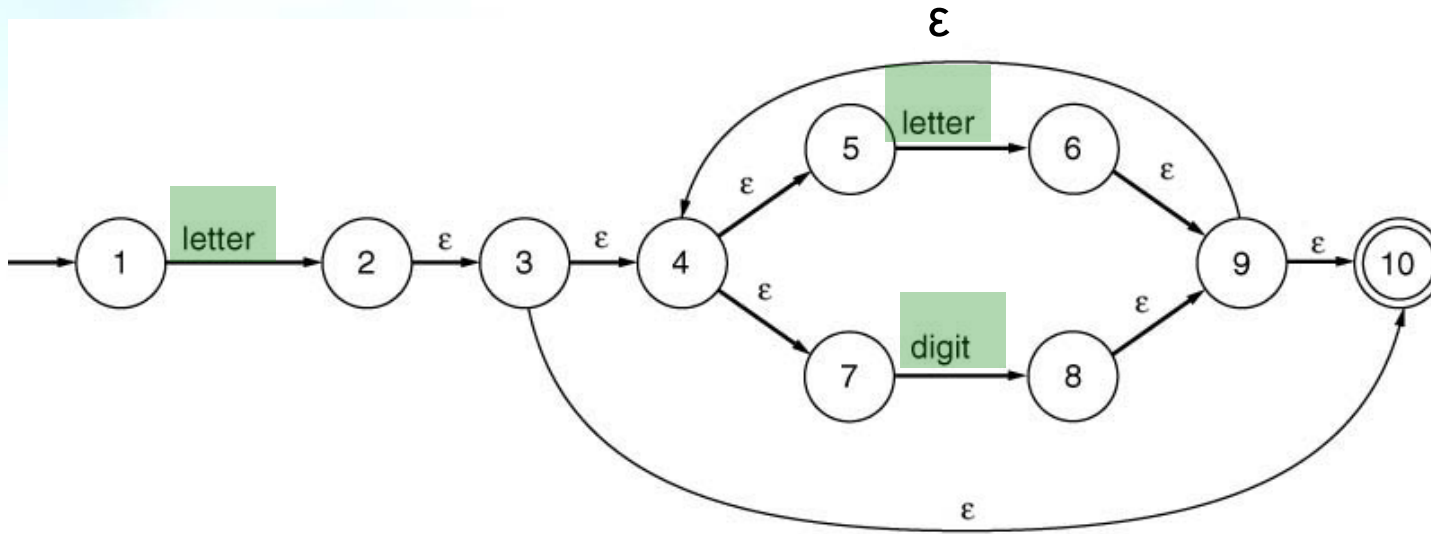
$\overline{3} \cup \overline{7}$

$\overline{5}$

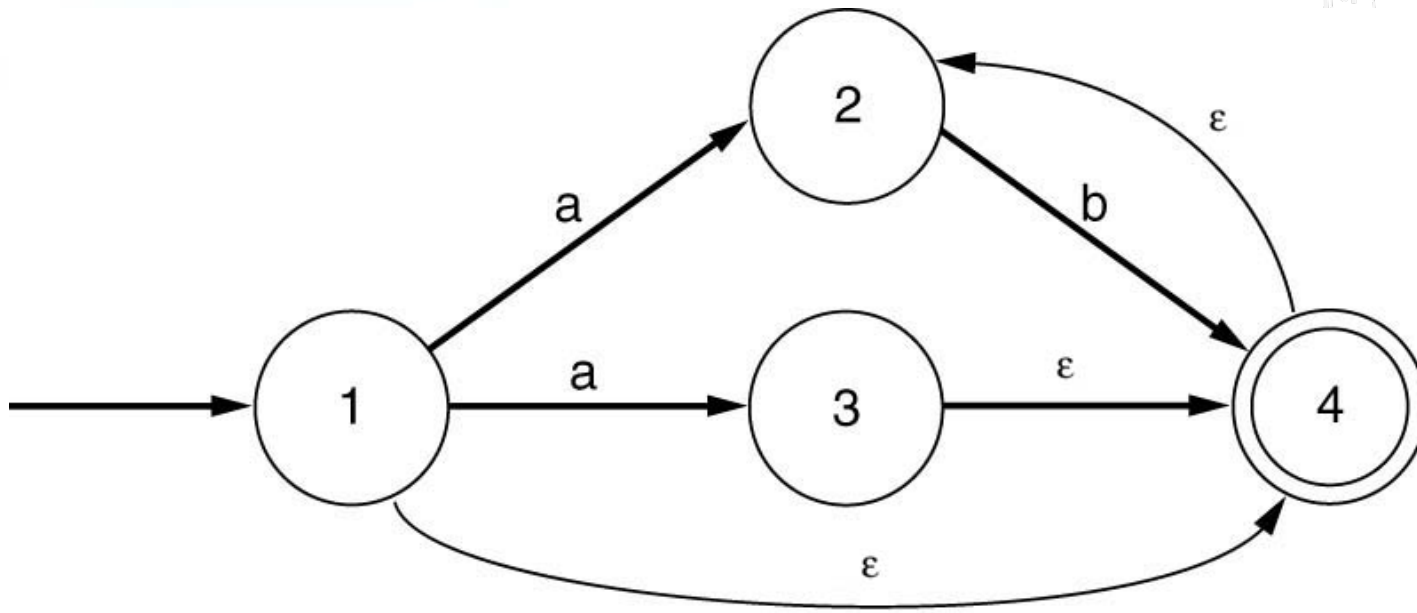


# NFA $\rightarrow$ DFA

•  $letter(letter|digit)^*$

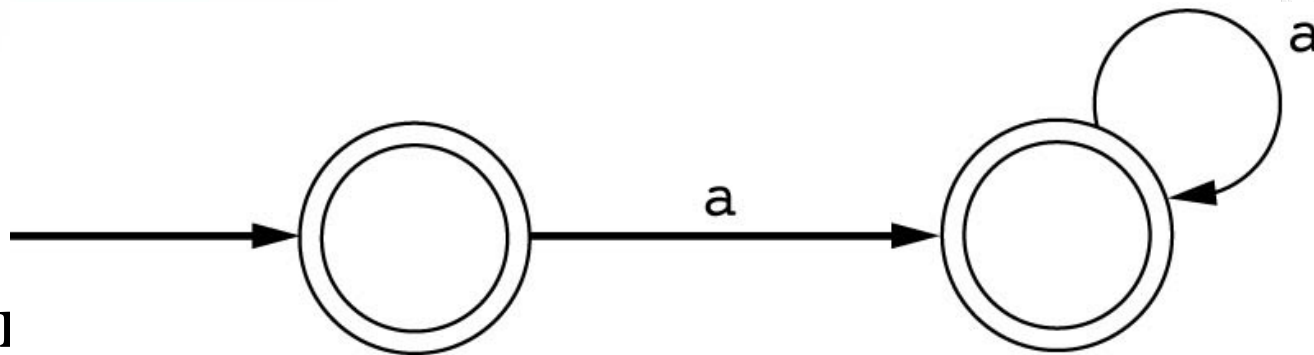


## Example 2.10

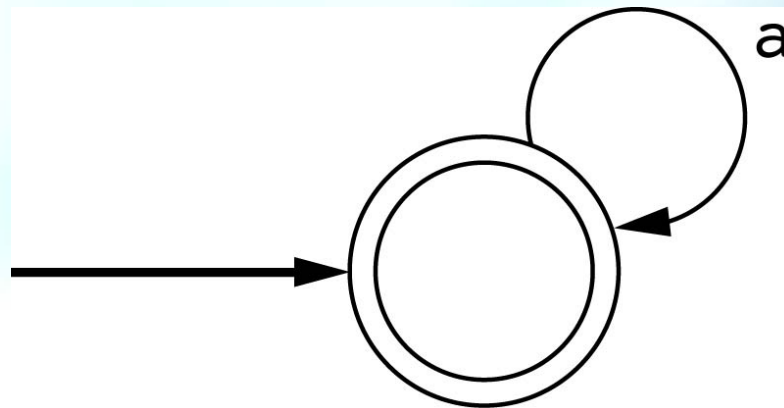


# Minimizing DFA

- DFA for  $a^*$  that is constructed by subset construction

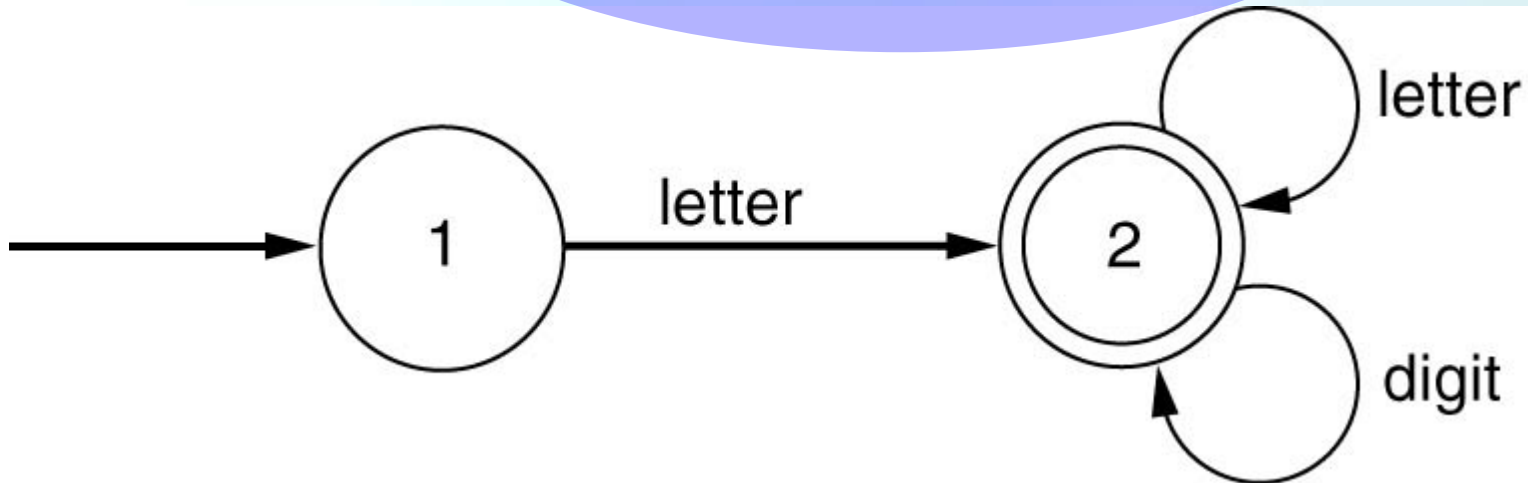
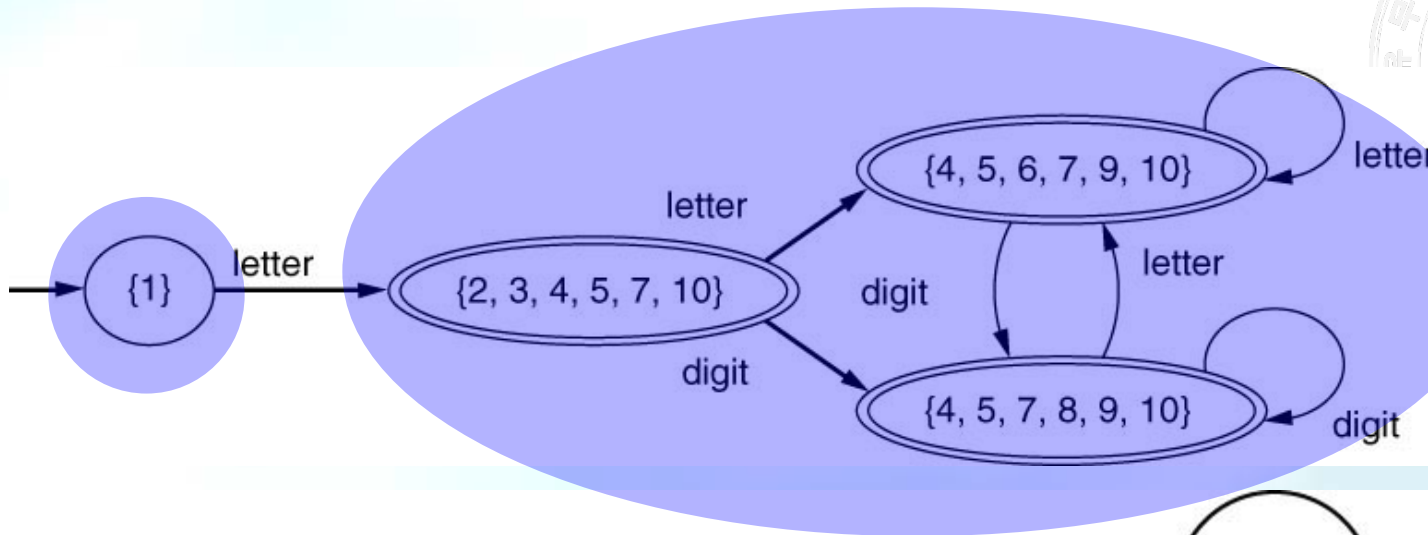


- ca1



# Minimizing DFA

- $letter(letter/digit)^*$



**Algorithm** given any DFA, there is an equivalent DFA containing a minimum number of states, and this minimum state is unique



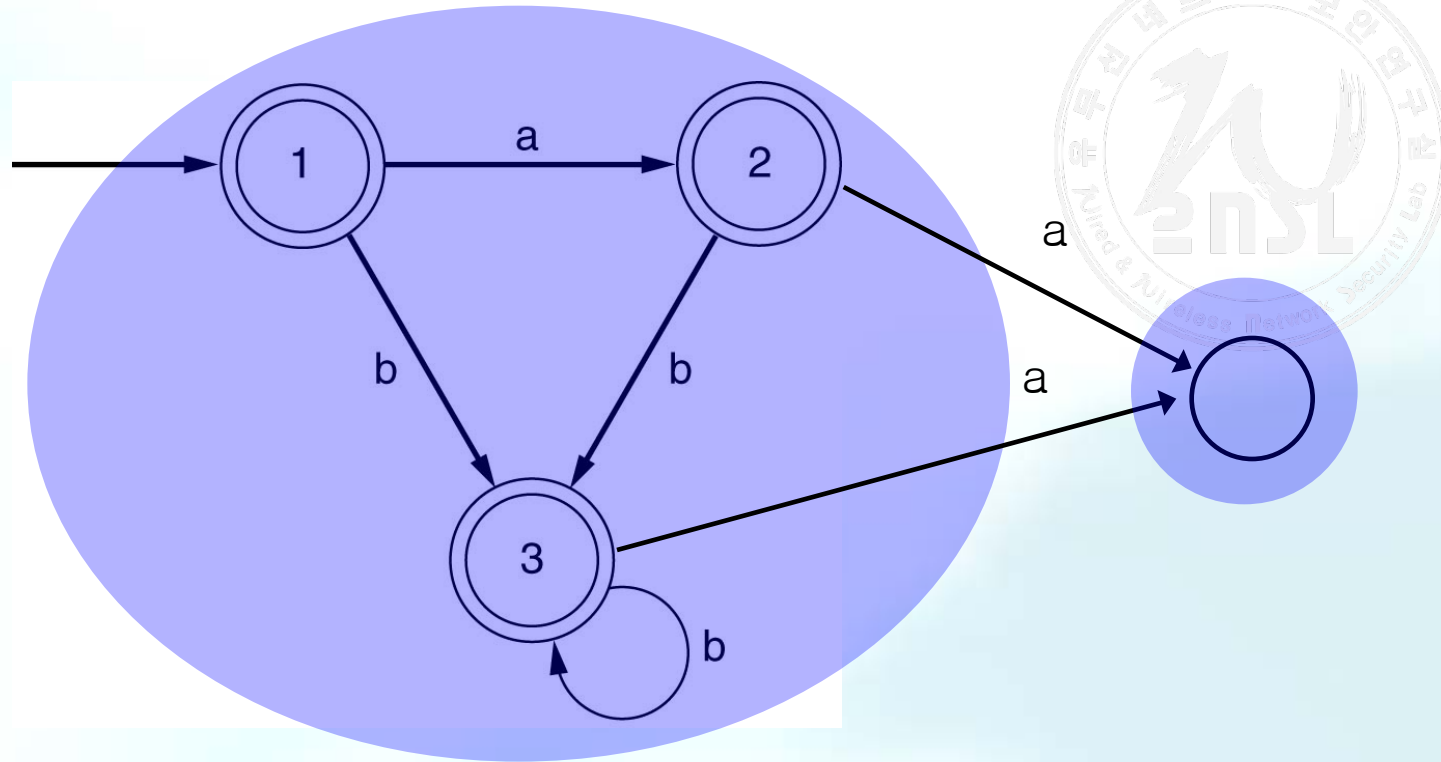
<http://usecurity.hanyang.ac.kr>

1. unified into single states
  1. one with all the accepting states
  2. the other with all the nonaccepting states
2. Consider transitions on each character  $a$  of the alphabet
3. if there are two accepting/nonaccepting states that have transitions on  $a$  that land in different states,  $a$  **distinguishes** the states and the set of states must be split.



# Minimizing DFA

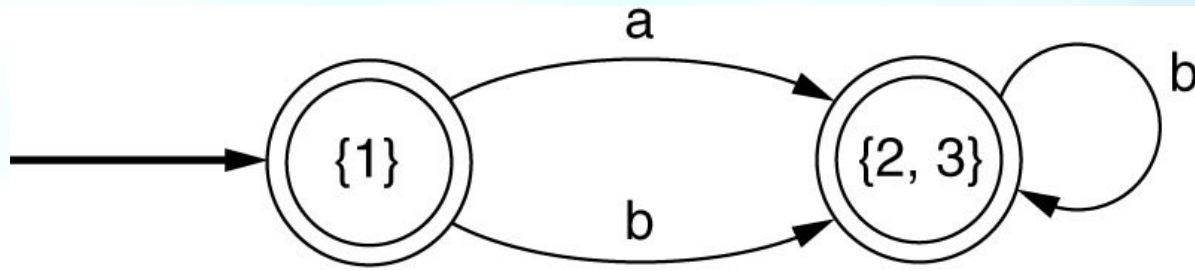
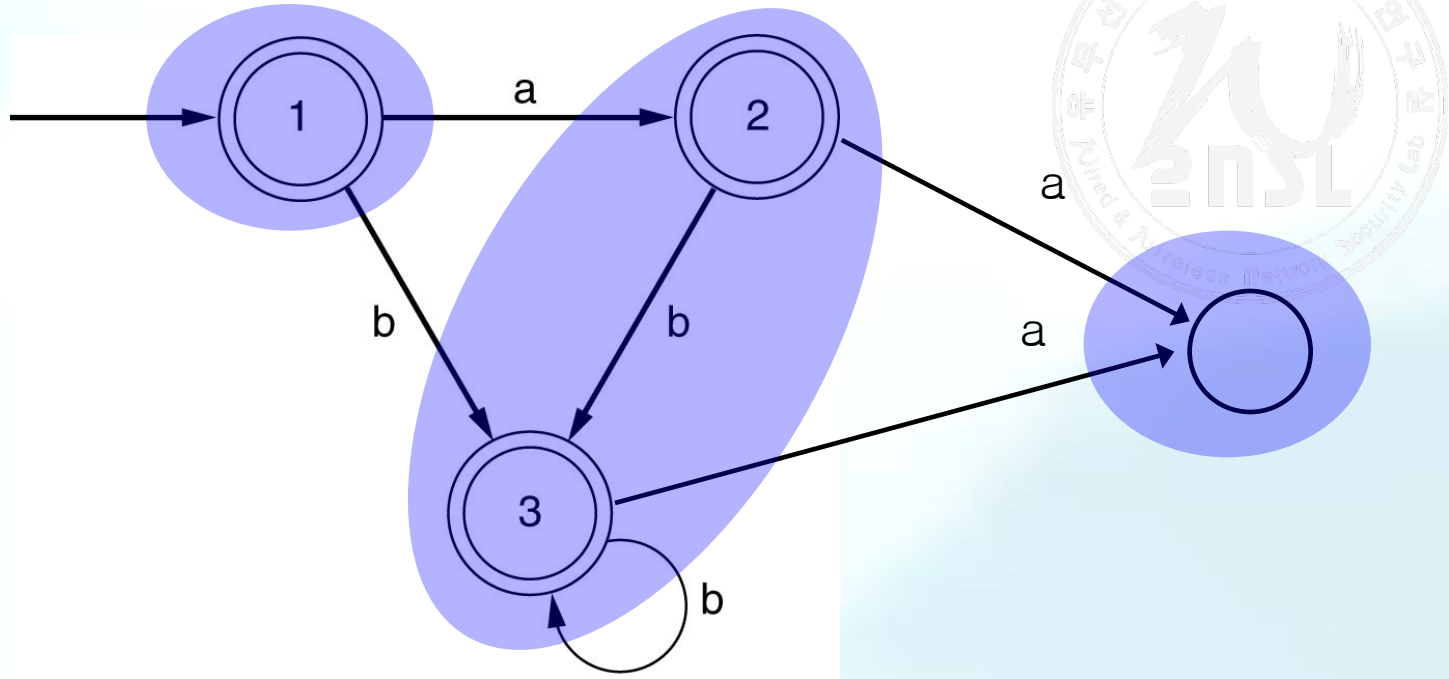
- $(a|\epsilon)b^*$



- $T(1,a) \neq T(2,a)$

# Minimizing DFA

- $(a|\epsilon)b^*$





# Use of Lex to generate a scanner automatically



<http://usecurity.hanyang.ac.kr>

- Lex program
  - Input: a text file containing
    - Regular expressions
    - Actions to be taken when each expression is matched
  - Output: C source code (lex.yy.c or lexyy.c)
    - Defining a procedure yylex
      - that is a table-driven implementation of a DFA
      - that operates like a getToken procedure



# Metacharacter conventions in Lex

- Table 2.2

Pattern	Meaning
<b>a</b>	the character <i>a</i>
<b>"a"</b>	the character <i>a</i> , even if <i>a</i> is a metacharacter
<b>\a</b>	the character <i>a</i> when <i>a</i> is a metacharacter
<b>a*</b>	zero or more repetitions of <i>a</i>
<b>a+</b>	one or more repetitions of <i>a</i>
<b>a?</b>	an optional <i>a</i>
<b>a b</b>	<i>a</i> or <i>b</i>
<b>(a)</b>	<i>a</i> itself
<b>[abc]</b>	any of the characters <i>a</i> , <i>b</i> , or <i>c</i>
<b>[a-d]</b>	any of the characters <i>a</i> , <i>b</i> , <i>c</i> , or <i>d</i>
<b>[^ab]</b>	any character except <i>a</i> or <i>b</i>
<b>.</b>	any character except a newline
<b>{xxx}</b>	the regular expression that the name <i>xxx</i> represents



# Format of a Lex input file



<http://usecurity.hanyang.ac.kr>

{definitions}

%%

{rules}

%%

{auxiliary routines}



# Table 2.3 Some Lex internal names



Lex Internal Name	Meaning/Use
<b>lex.yy.c</b> or <b>lexyy.c</b>	Lex output file name
<b>yylex</b>	Lex scanning routine
<b>yytext</b>	string matched on current action
<b>yyin</b>	Lex input file (default: <b>stdin</b> )
<b>yyout</b>	Lex output file (default: <b>stdout</b> )
<b>input</b>	Lex buffered input routine
<b>ECHO</b>	Lex default action (print <b>yytext</b> to <b>yyout</b> )

# Homework #1



<http://usecurity.hanyang.ac.kr>

- Exercises
  - 2.1, 2.2, 2.12, 2.13, 2.16
- 주의 사항
  - Handwriting으로 제출
  - Cover page는 생략. 첫 번째 페이지에 homework 번호, 학번, 이름 기입

