

---

# Lecture 3: Lexical Analysis Cont.

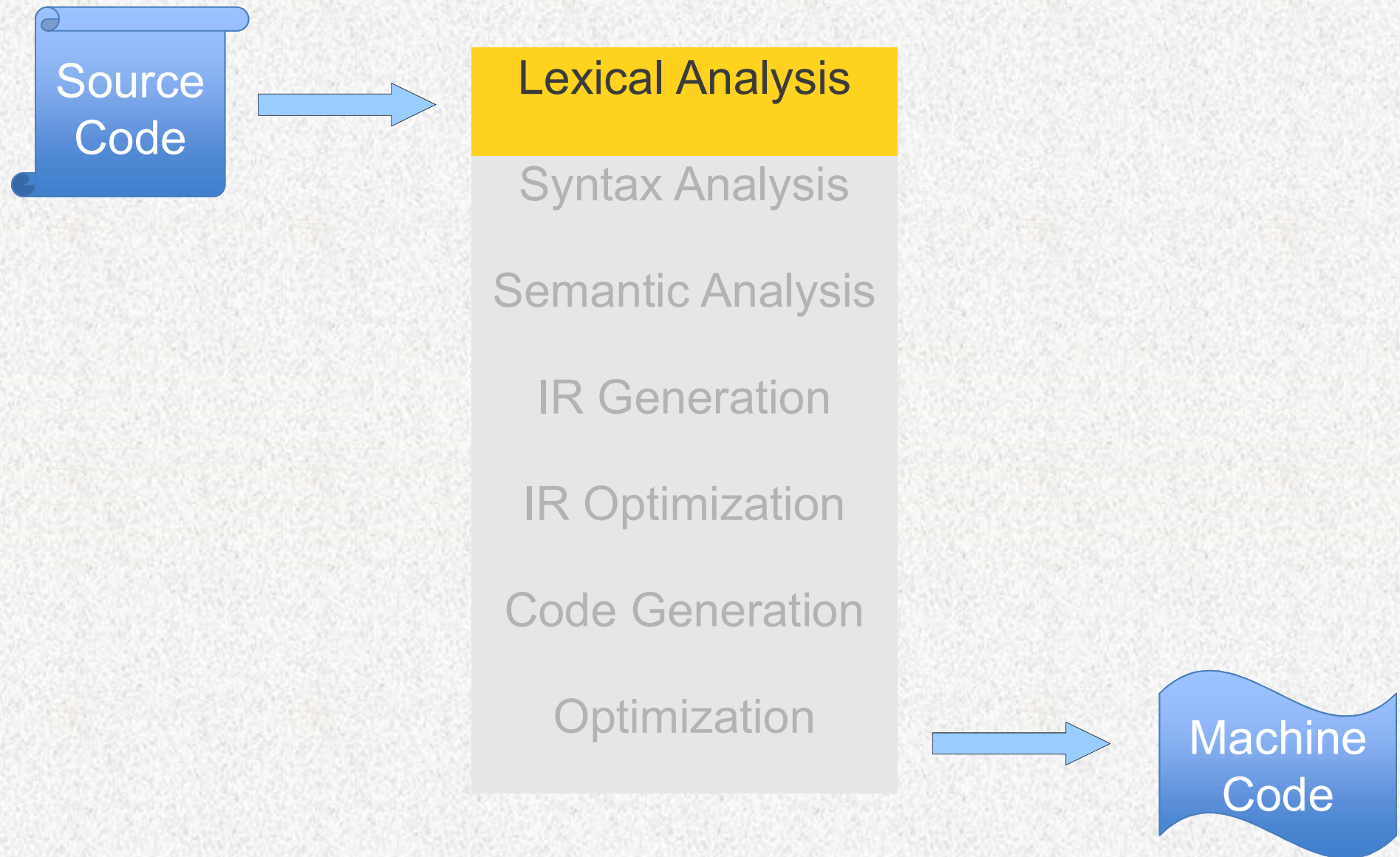
---

Xiaoyuan Xie 谢晓园

[xxie@whu.edu.cn](mailto:xxie@whu.edu.cn)

计算机学院E301

# Where We Are







---

# Formalisms of tokens

---

# Regular Expression

# Finite Automaton



# Implementing Regular Expressions

- Regular expressions can be implemented using **finite automata**.
  - Regular expressions = **specification**
  - Finite automata = **implementation**
- There are two main kinds of finite automata:
  - **NFAs** (**nondeterministic** finite automata), which we'll see in a second, and
  - **DFAs** (**deterministic** finite automata), which we'll see later.

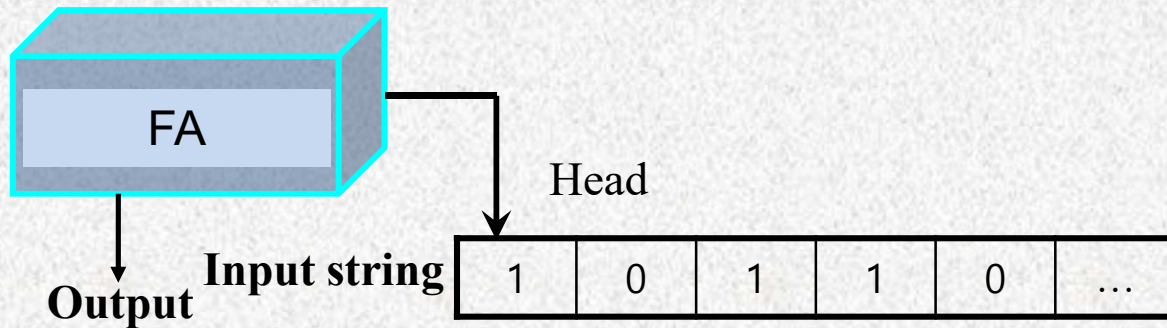


# Finite Automatons

- A finite automaton is a 5-tuple  $(S, \Sigma, \delta, s_0, F)$ 
  - A set of states  $S$  --- nodes
  - An input alphabet  $\Sigma$
  - A transition function  $\delta(S_i, a) = S_j$
  - A start state  $S_0$
  - A set of accepting states  $F \subseteq S$



# Finite Automata



- Input: a string
- Output: accept if the scanning of input string reaches its EOF and the FA reaches an accepting state; reject otherwise

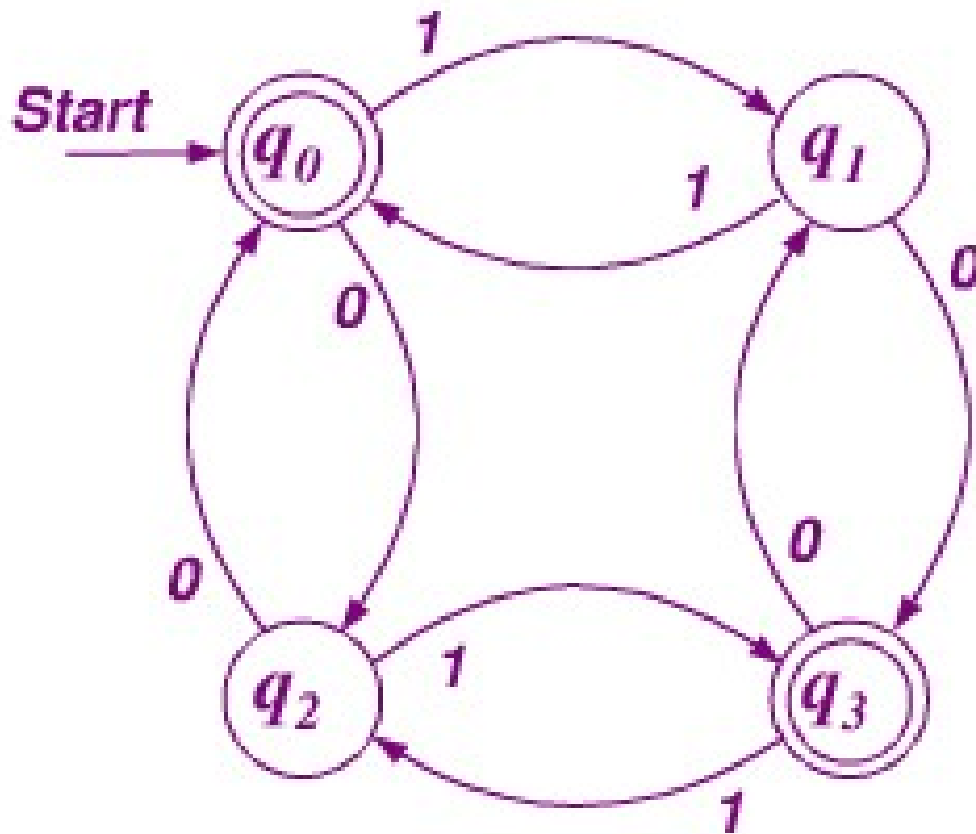


# Strings accepted by an FA

- An FA *accepts an input string  $x$*  iff there is some path with edges labeled with symbols from  $x$  in sequence from the start state to some accepting state in the transition graph
- A state transition from one state to another on the path is called a *move*



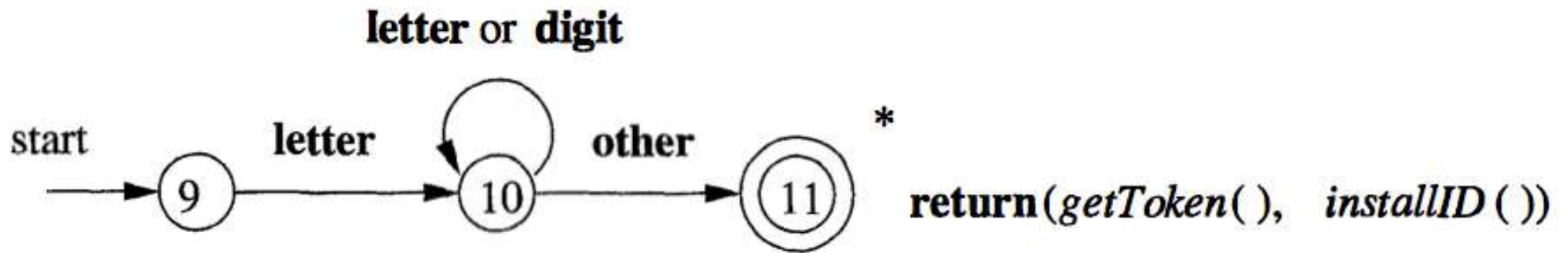
# A More Complex Automaton



“1010”: accept

“101”: reject

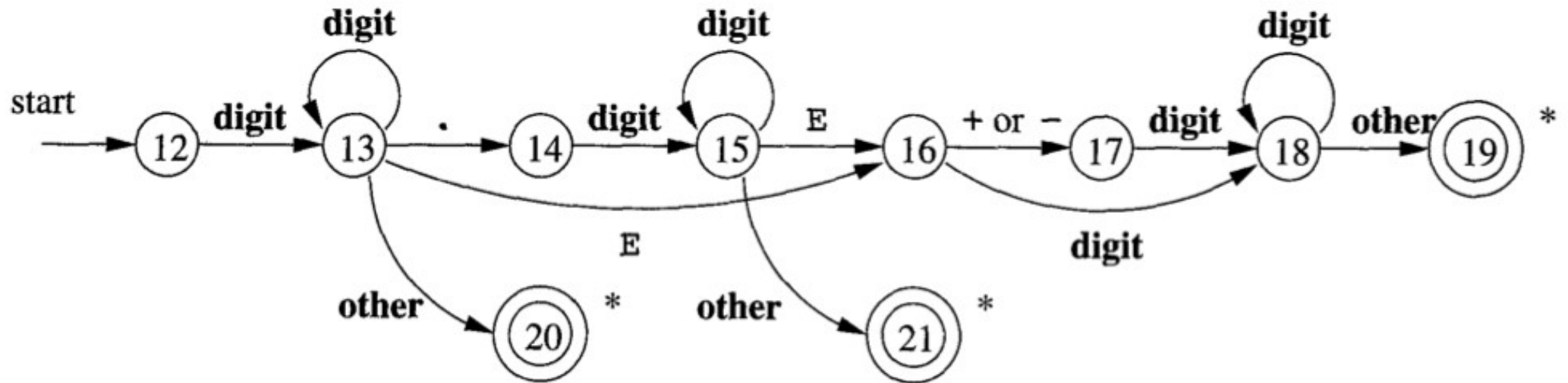
# A More Complex Automaton



h	i	1	2	3
---	---	---	---	---



# A More Complex Automaton



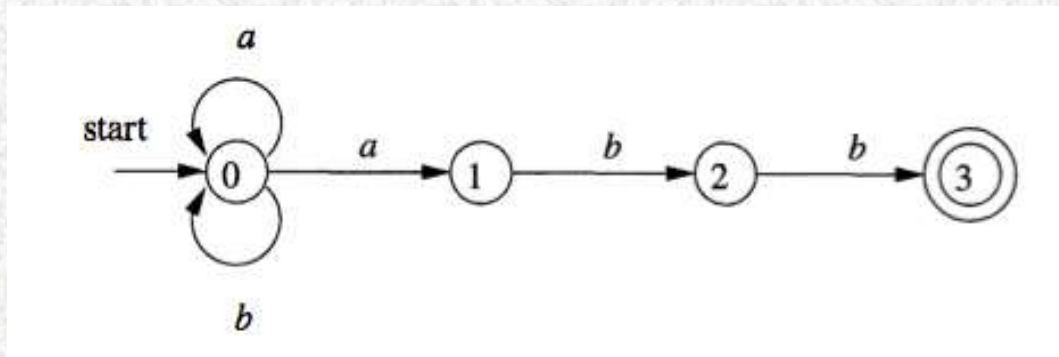
1	2	.	3	7	5
---	---	---	---	---	---

# Strings accepted by an FA

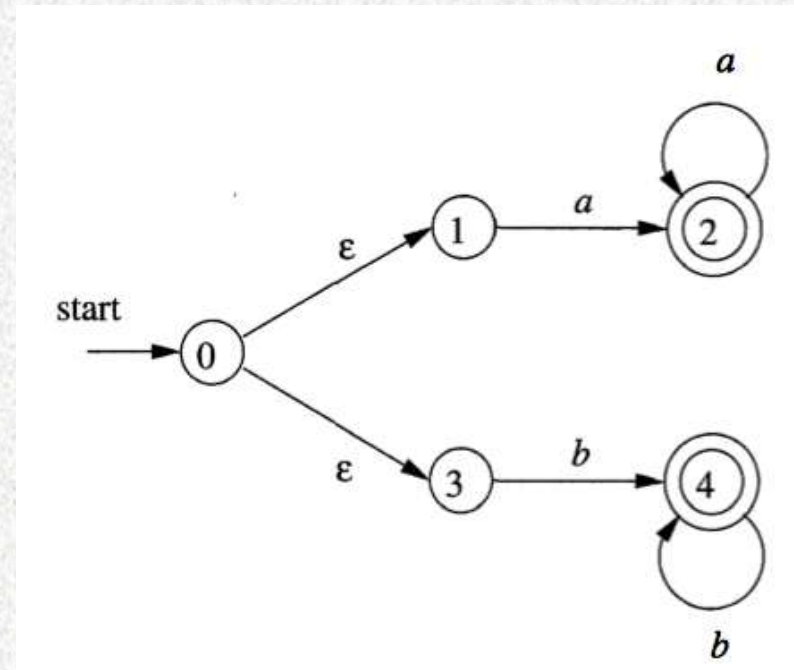
- The *language defined by* an FA is **the set of input strings it accepts**, such as  $(a|b)^*abb$  for the example NFA



# Strings accepted by an FA



$(a|b)^*abb$



$aa^*|bb^*$

# Finite Automata

- Finite automata is a recognizer
- Given an input string, they simply say "yes" or "no" about each possible input string



# Nondeterministic Finite Automata (NFA)

- Definition: an NFA is a 5-tuple  $(S, \Sigma, \delta, s_0, F)$  where
  - $S$  is a finite set of *states*
  - $\Sigma$  is a finite set of *input symbol alphabet*
  - $\delta$  is a *mapping* from  $S \times \Sigma \cup \{\epsilon\}$  to a set of states
  - $S_0 \subseteq S$  is the set of *start states*
  - $F \subseteq S$  is the set of *accepting* (or *final*) states



# Nondeterministic Finite Automata (NFA)

- **Transition Graph**

**Node:** State

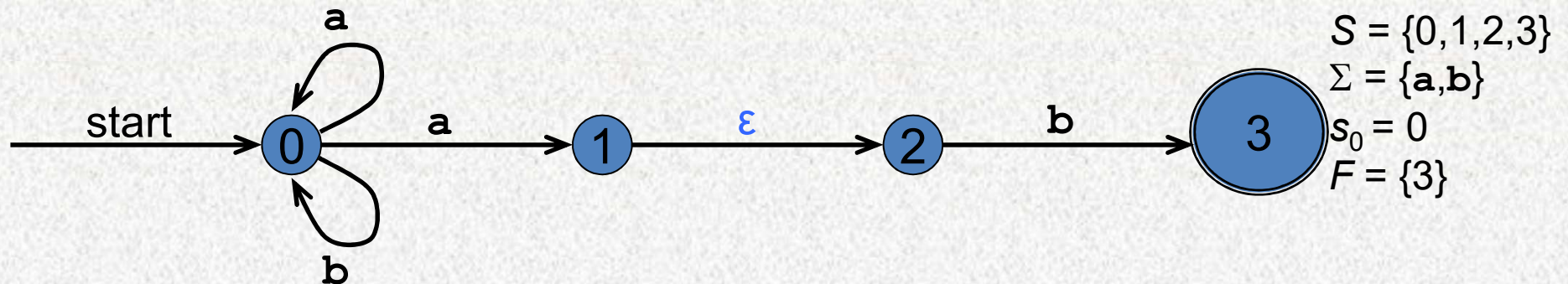
- Non-terminal state: 
- Terminal state: 
- Starting state: 

**Edge:** state transition  $f(S_i, a) = S_j$  



# Transition Graph

- An NFA can be diagrammatically represented by a labeled directed graph called a *transition graph*





# Nondeterministic Finite Automata (NFA)

- **Transit table**

- Line: State
  - Starting state: in general, the first line, or label “+”;
  - Terminal state: “\*” or “-” ;
- Column: All symbols in  $\Sigma$
- Cell: state transition mapping



# Transition Table

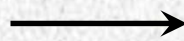
- The mapping  $\delta$  of an NFA can be represented in a *transition table*

$$\delta(0, \mathbf{a}) = \{0, 1\}$$

$$\delta(0, \mathbf{b}) = \{0\}$$

$$\delta(1, \mathbf{b}) = \{2\}$$

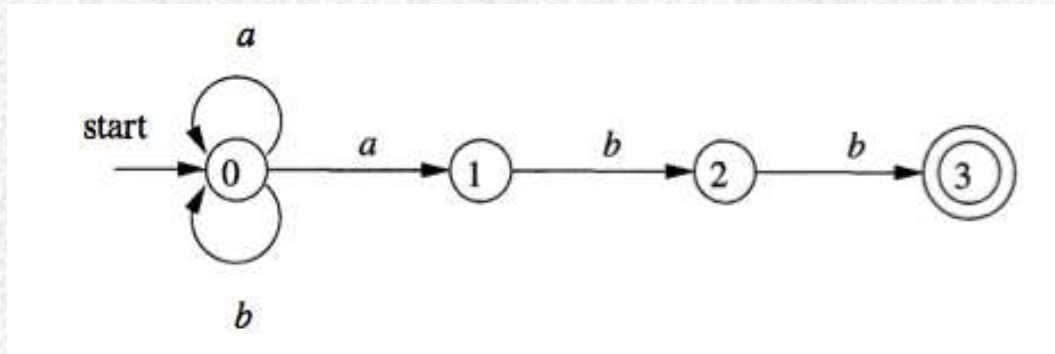
$$\delta(2, \mathbf{b}) = \{3\}$$



<i>State</i>	<i>Input</i> <b>a</b>	<i>Input</i> <b>b</b>
0	{0,1}	{0}
1		{2}
2		{3}



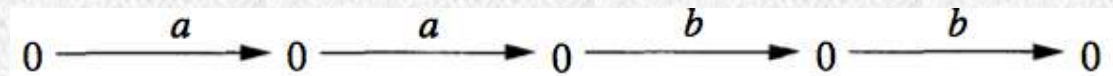
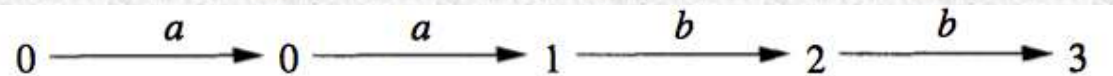
# NFA Example 2



Transition Table

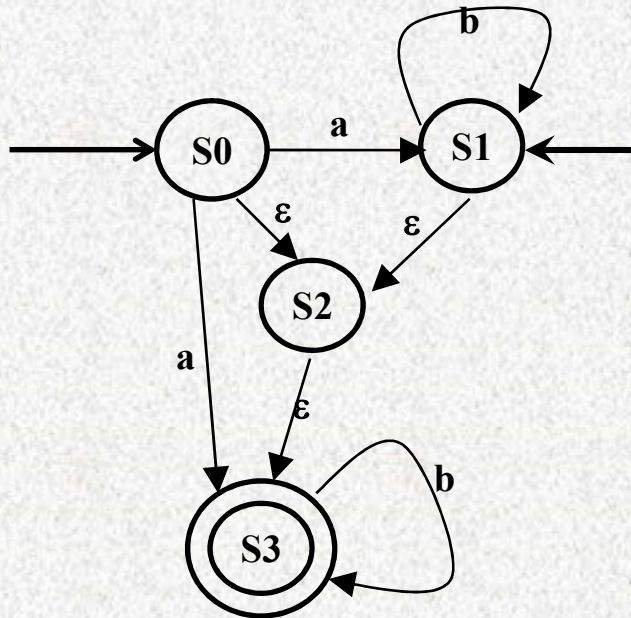
STATE	$a$	$b$	$\epsilon$
0	$\{0, 1\}$	$\{0\}$	$\emptyset$
1	$\emptyset$	$\{2\}$	$\emptyset$
2	$\emptyset$	$\{3\}$	$\emptyset$
3	$\emptyset$	$\emptyset$	$\emptyset$

Acceptance of input strings





# NFA Example 3



	a	b	$\epsilon$
$S0^+$	$\{S1, S3\}$		$\{S2\}$
$S1^+$		$\{S1\}$	$\{S2\}$
$S2$			$\{S3\}$
$S3^-$		$\{S3\}$	



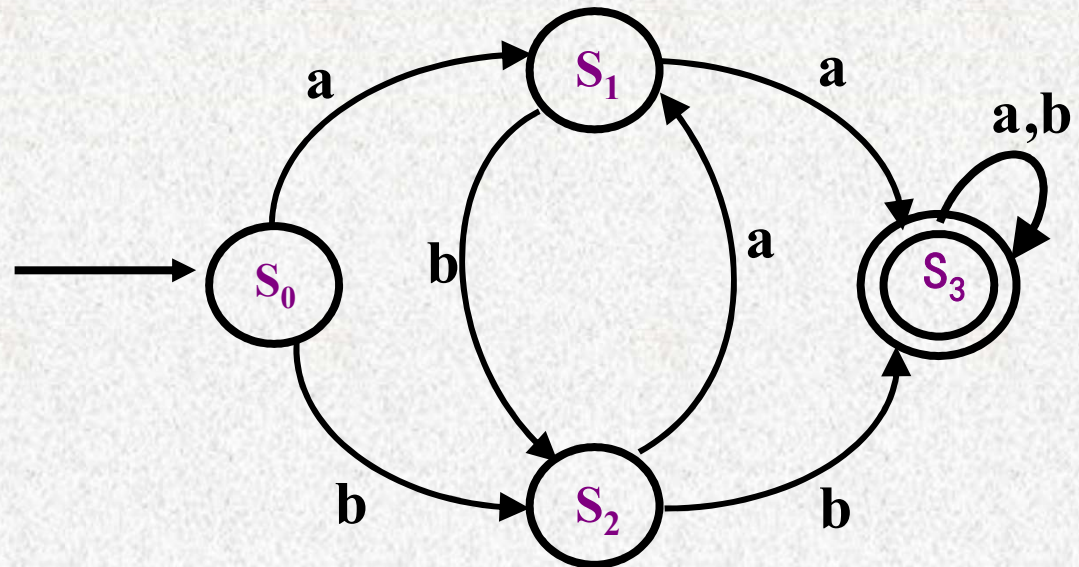
# Deterministic Finite Automata (DFA)

- Definition: an DFA is a 5-tuple  $(S, \Sigma, \delta, s_0, F)$ , is a special case of NFA
  - There are no moves on input  $\varepsilon$ , and
  - For each state  $s$  and input symbol  $a$ , there is exactly one edge out of  $s$  labeled  $a$ .



# Deterministic Finite Automata (DFA)

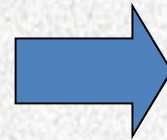
- DFA  $M = (\{S_0, S_1, S_2, S_3\}, \{a, b\}, f, S_0, \{S_3\})$ , :  
     $f(S_0, a) = S_1$                        $f(S_2, a) = S_1$   
     $f(S_0, b) = S_2$                        $f(S_2, b) = S_3$   
     $f(S_1, a) = S_3$                        $f(S_3, a) = S_3$   
     $f(S_1, b) = S_2$                        $f(S_3, b) = S_3$





# Deterministic Finite Automata (DFA)

- For example, DFA  $M = (\{0,1,2,3,4\}, \{a,b\}, \delta, \{0\}, \{3\})$
- $\delta(0, a) = 1$      $\delta(0, b) = 4$   
 $\delta(1, a) = 4$      $\delta(1, b) = 2$   
 $\delta(2, a) = 3$      $\delta(2, b) = 4$   
 $\delta(3, a) = 3$      $\delta(3, b) = 3$   
 $\delta(4, a) = 4$      $\delta(4, b) = 4$

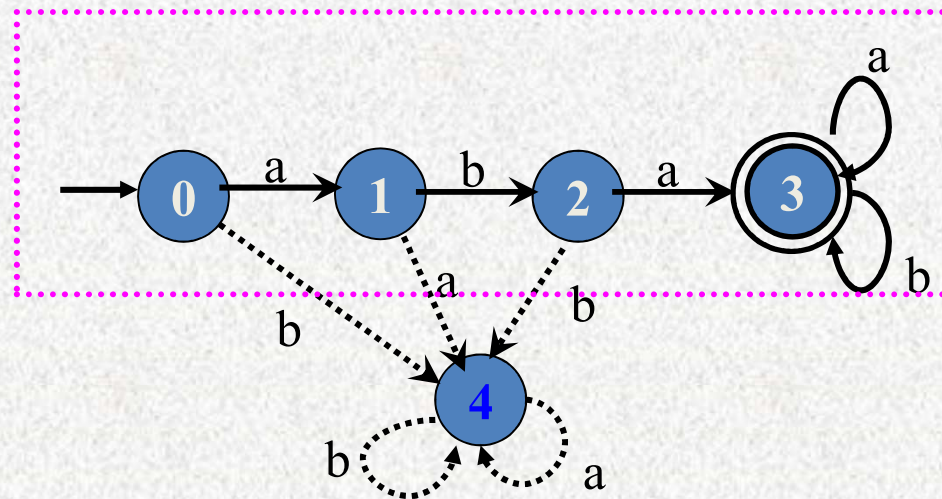


	a	b
0 <sup>+</sup>	1	4
1	4	2
2	3	4
3 <sup>-</sup>	3	3
4	4	4



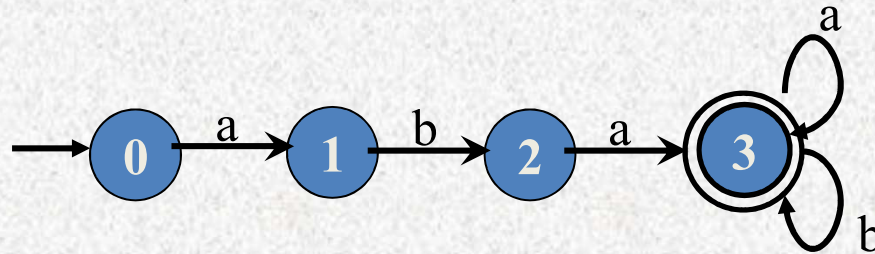
# Deterministic Finite Automata (DFA)

	a	b
0+	1	4
1	4	2
2	3	4
3 <sup>-</sup>	3	3
4	4	4





# Deterministic Finite Automata (DFA)



	a	b
0+	1	⊥
1	⊥	2
2	3	⊥
3-	3	3



	a	b
0+	1	
1		2
2	3	
3-	3	3



# Deterministic Finite Automata (DFA)

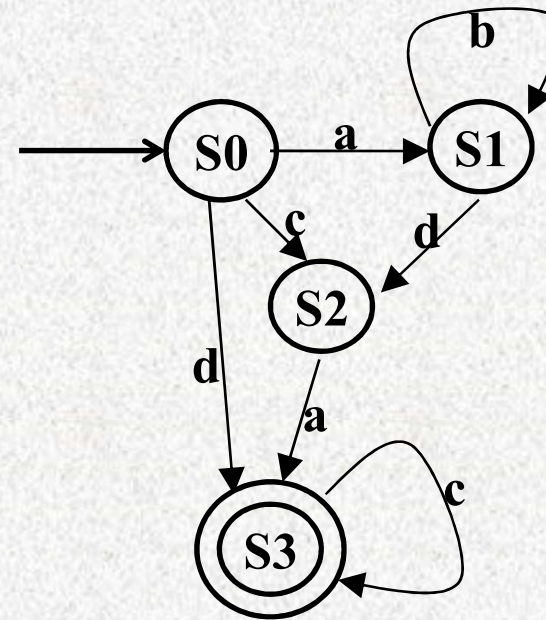
$\Sigma$ : {a, b, c, d}

S: {S0, S1, S2, S3}

Start: S0

Terminal: {S3}

f: {(S0,a)→ S1, (S0,c)→S2,  
(S0,d)→S3, (S1,b)→S1,  
(S1,d)→S2, (S2,a)→S3,  
(S3, c)→S3}





# NFA v.s. DFA




# NFA v.s. DFA

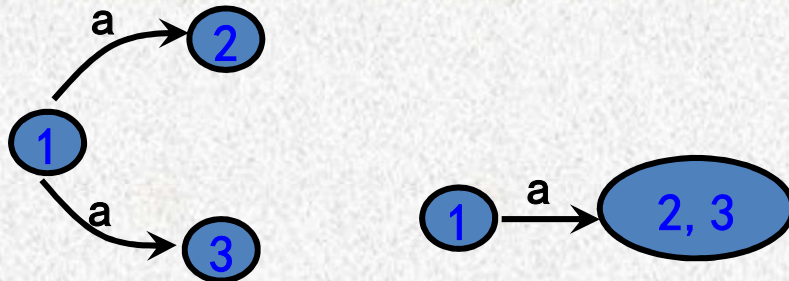
	DFA	NFA
Initial	Single starting state	A set of starting states
$\epsilon$ dege	Not allowed	Allowed
$\delta(S, a)$	$S'$ or $\perp$	$\{S_1, \dots, S_n\}$ or $\perp$
Implementation	Deterministic	Nondeterministic

- DFA accepts an input string with only one path
- NFA accepts an input string with possibly multiple paths



# Construct DFA from NFA

- Construct DFA from NFA
  - For any NFA, **there exists an equivalent DFA**
  - Idea of construction: eliminate the uncertainty
  - Merge N states in NFA into **one single state**
    - Eliminate  $\varepsilon$    
A diagram showing two states, 4 and 5, with a horizontal arrow labeled  $\varepsilon$  pointing from 4 to 5. To the right of this is a blue oval containing the text "4, 5".
    - Eliminate multiple mapping





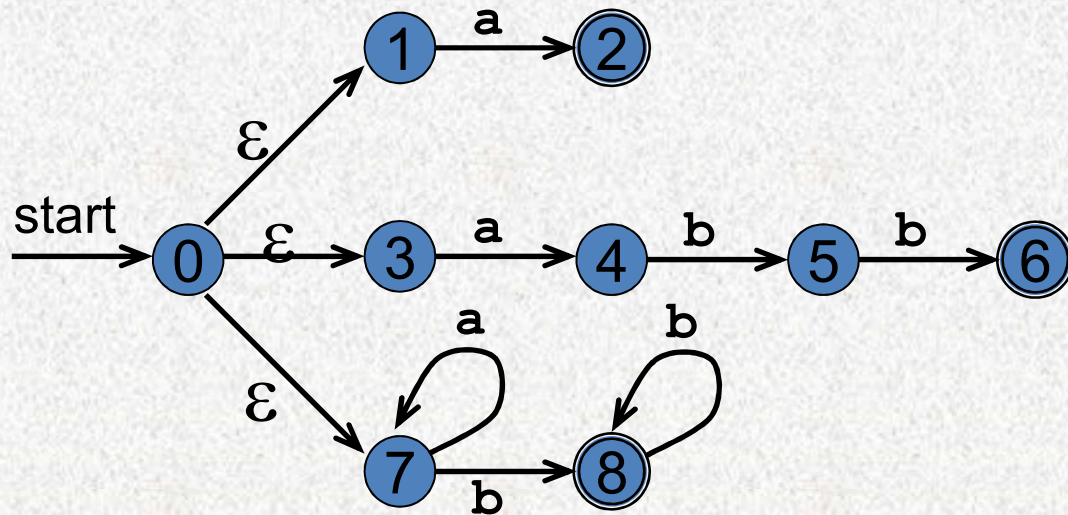
# Construct DFA from NFA

- **INPUT:** An NFA  $N$ .
- **OUTPUT:** A DFA  $D$  accepting the same language as  $N$ .
- **METHOD:** The algorithm constructs a transition table  $D_{\text{tran}}$  for  $D$ . Each state of  $D$  is a set of NFA states, and we construct  $D_{\text{tran}}$  so  $D$  will simulate “in parallel” all possible moves  $N$  can make on a given input string.

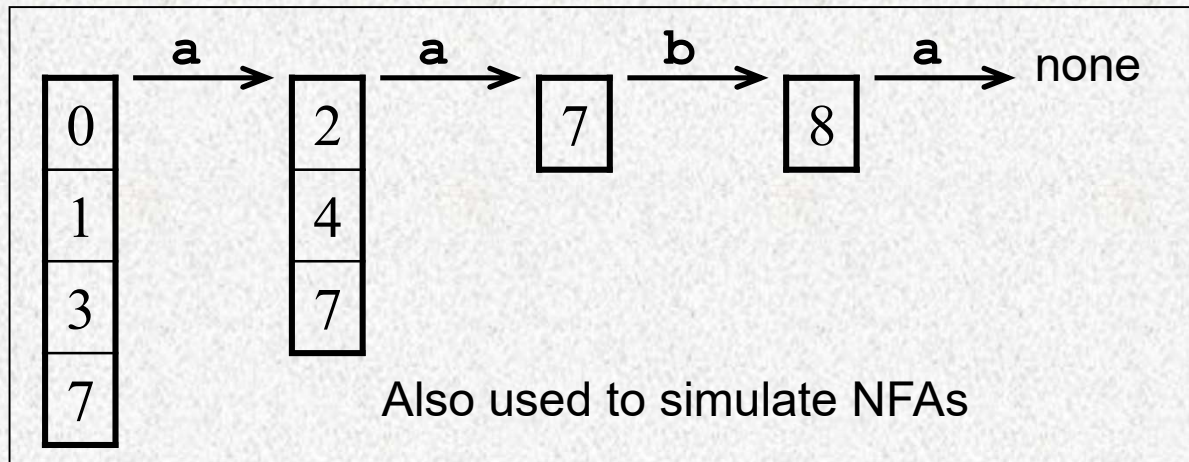
OPERATION	DESCRIPTION
$\epsilon\text{-closure}(s)$	Set of NFA states reachable from NFA state $s$ on $\epsilon$ -transitions alone.
$\epsilon\text{-closure}(T)$	Set of NFA states reachable from some NFA state $s$ in set $T$ on $\epsilon$ -transitions alone; $= \bigcup_{s \text{ in } T} \epsilon\text{-closure}(s)$ .
$\text{move}(T, a)$	Set of NFA states to which there is a transition on input symbol $a$ from some state $s$ in $T$ .



# $\epsilon$ -closure and *move* Examples



$\epsilon$ -closure( $\{0\}$ ) =  $\{0, 1, 3, 7\}$   
 $move(\{0, 1, 3, 7\}, a) = \{2, 4, 7\}$   
 $\epsilon$ -closure( $\{2, 4, 7\}$ ) =  $\{2, 4, 7\}$   
 $move(\{2, 4, 7\}, a) = \{7\}$   
 $\epsilon$ -closure( $\{7\}$ ) =  $\{7\}$   
 $move(\{7\}, b) = \{8\}$   
 $\epsilon$ -closure( $\{8\}$ ) =  $\{8\}$   
 $move(\{8\}, a) = \emptyset$





# The Subset Construction Algorithm

- NFAs can be in many states at once, while DFAs can only be in a single state at a time.
- Key idea: **Make the DFA simulate the NFA.**
- Have the states of the DFA correspond to the *sets of states* of the NFA.
- Transitions between states of DFA correspond to transitions between *sets of states* in the NFA.

# The Subset Construction Algorithm

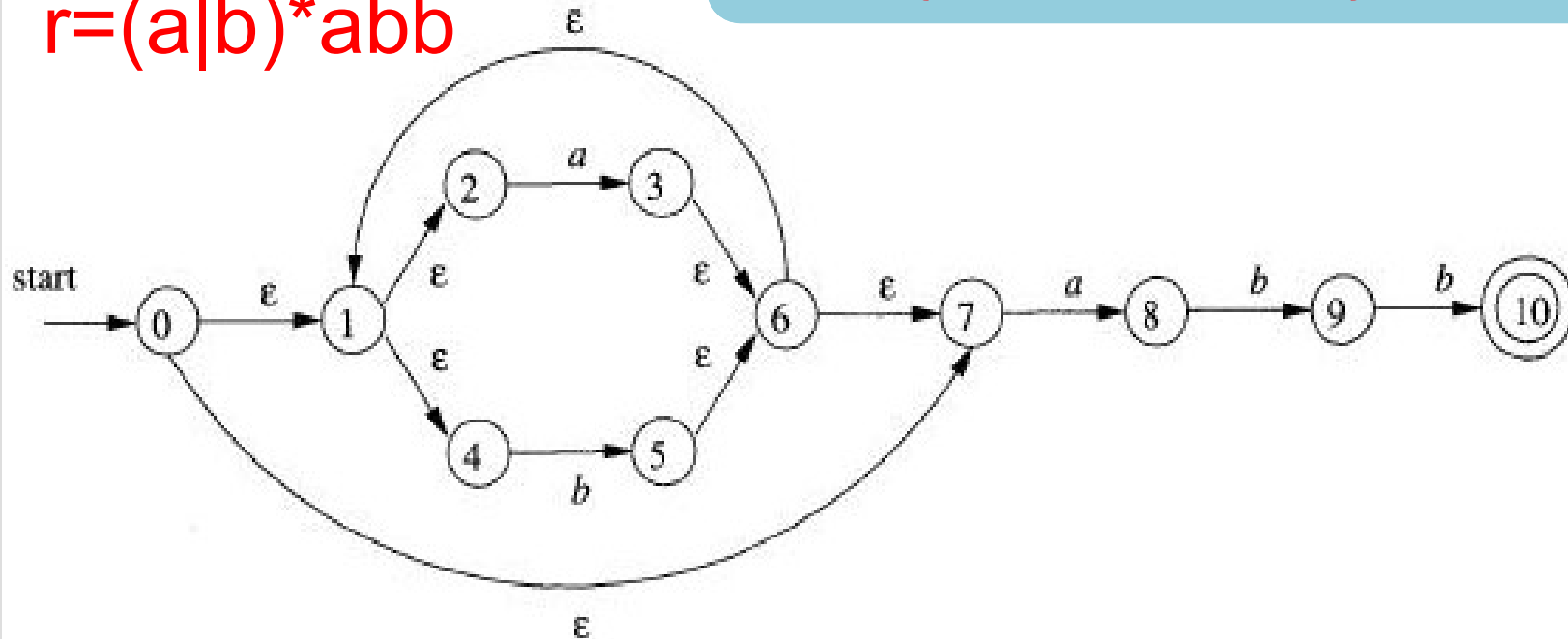
```
initially,  $\epsilon$ -closure( $s_0$ ) is the only state in  $Dstates$ , and it is unmarked;  
while ( there is an unmarked state  $T$  in  $Dstates$  ) {  
    mark  $T$ ;  
    for ( each input symbol  $a$  ) {  
         $U = \epsilon$ -closure( $move(T, a)$ );  
        if (  $U$  is not in  $Dstates$  )  
            add  $U$  as an unmarked state to  $Dstates$ ;  
         $Dtran[T, a] = U$ ;  
    }  
}
```



# Subset Construction Example 1

First, Initial state of NFA is  $\epsilon$ -closure(0),  
i.e.  $A = \{0, 1, 2, 4, 7\}$ ,  $\Sigma = \{a, b\}$

$r = (a|b)^*abb$



$Dtran[A, a] = \epsilon\text{-closure}(\text{move}(A, a)) = \epsilon\text{-closure}(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\}$ ,

Let  $B = Dtran[A, a]$

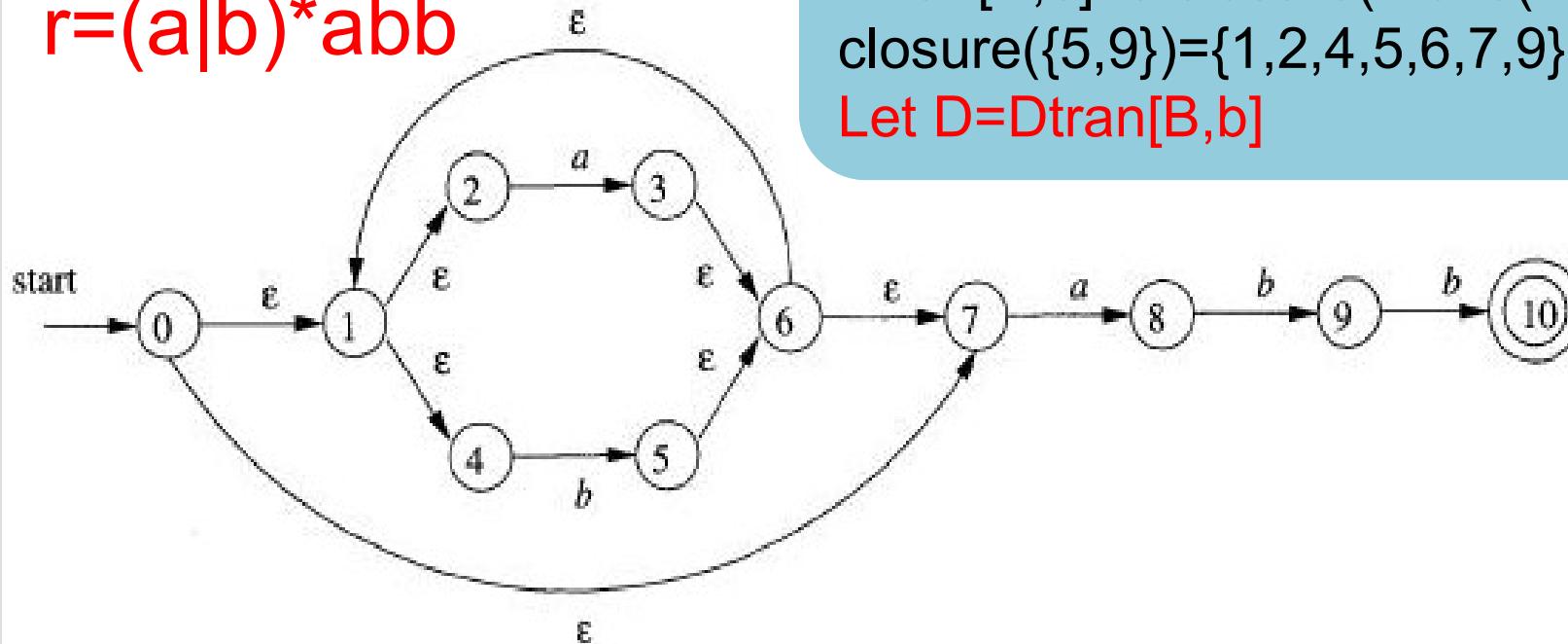
$Dtran[A, b] = \epsilon\text{-closure}(\text{move}(A, b)) = \epsilon\text{-closure}(\{5\}) = \{1, 2, 4, 6, 7\}$ ,

Let  $C = Dtran[A, b]$



# Subset Construction Example 1

$r = (a|b)^*abb$



$Dtran[B, a] = \epsilon\text{-closure}(\text{move}(B, a)) = \epsilon\text{-closure}(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\} = B$

$Dtran[B, b] = \epsilon\text{-closure}(\text{move}(B, b)) = \epsilon\text{-closure}(\{5, 9\}) = \{1, 2, 4, 5, 6, 7, 9\},$

Let  $D = Dtran[B, b]$

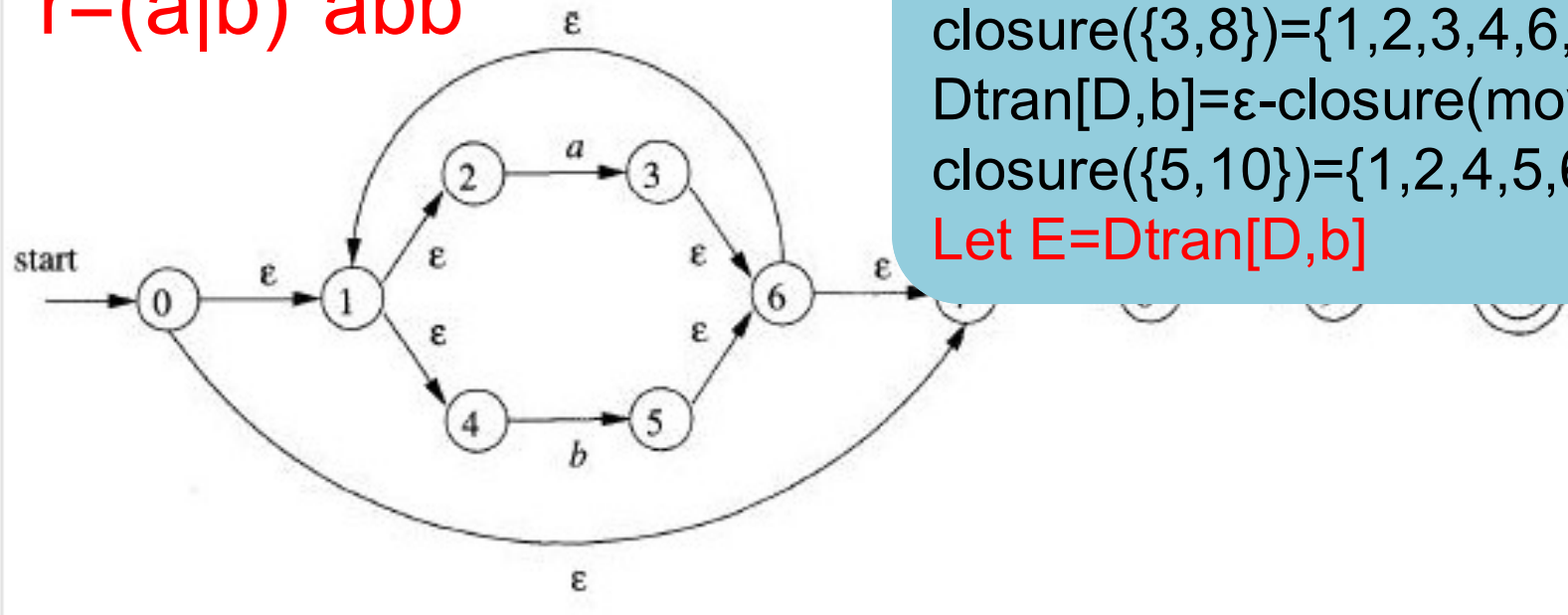
$Dtran[C, a] = \epsilon\text{-closure}(\text{move}(C, a)) = \epsilon\text{-closure}(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\} = B$

$Dtran[C, b] = \epsilon\text{-closure}(\text{move}(C, b)) = \epsilon\text{-closure}(\{5\}) = \{1, 2, 4, 6, 7\} = C$



# Subset Construction Example 1

$r = (a|b)^*abb$



$Dtran[D, a] = \epsilon\text{-closure}(\text{move}(D, a)) = \epsilon\text{-closure}(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\} = B$

$Dtran[D, b] = \epsilon\text{-closure}(\text{move}(D, b)) = \epsilon\text{-closure}(\{5, 10\}) = \{1, 2, 4, 5, 6, 7, 10\},$

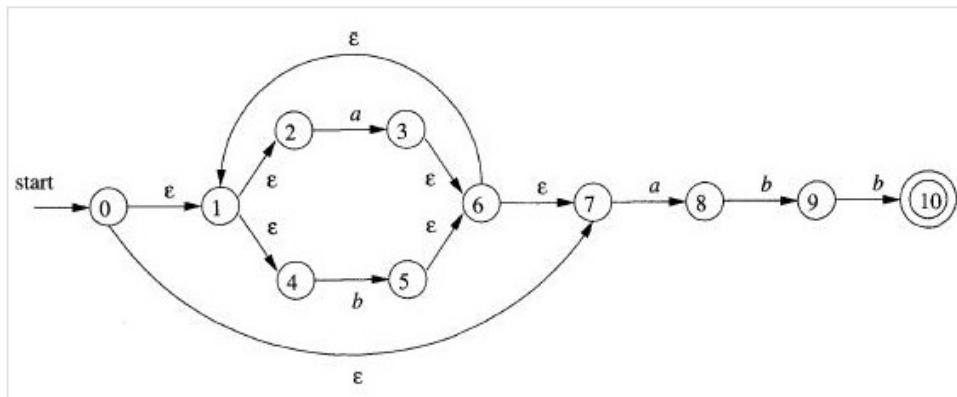
Let  $E = Dtran[D, b]$

$Dtran[E, a] = \epsilon\text{-closure}(\text{move}(E, a)) = \epsilon\text{-closure}(\{3, 8\}) = \{1, 2, 3, 4, 6, 7, 8\} = B$

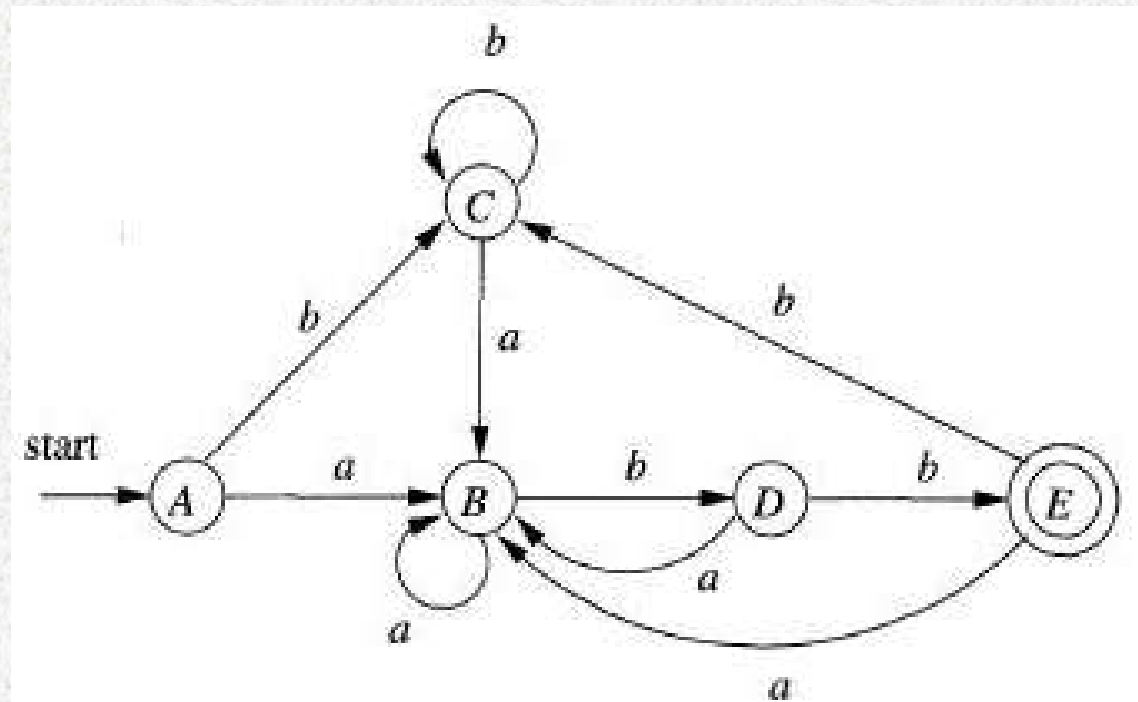
$Dtran[E, b] = \epsilon\text{-closure}(\text{move}(E, b)) = \epsilon\text{-closure}(\{5\}) = \{1, 2, 4, 6, 7\} = C$



# Subset Construction Example 1

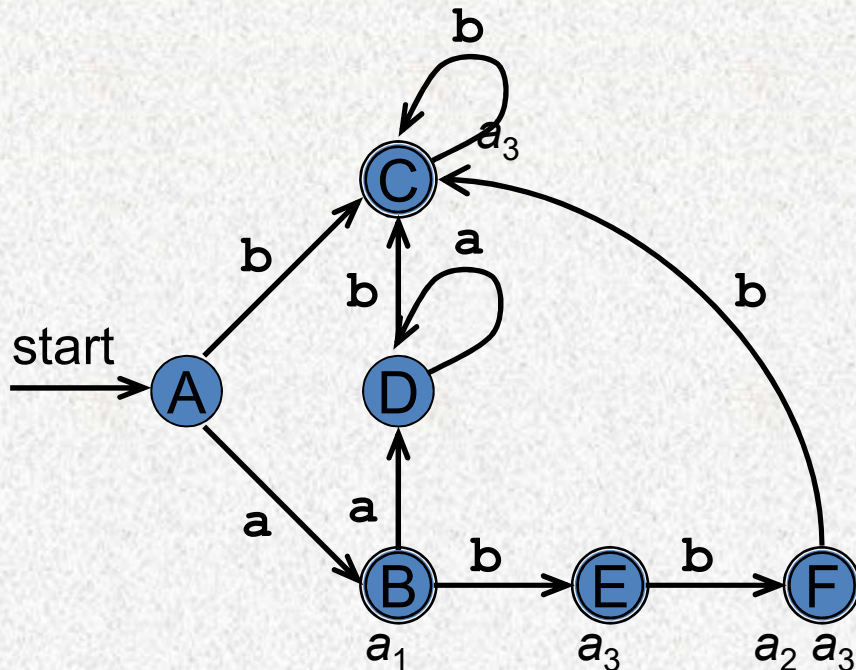
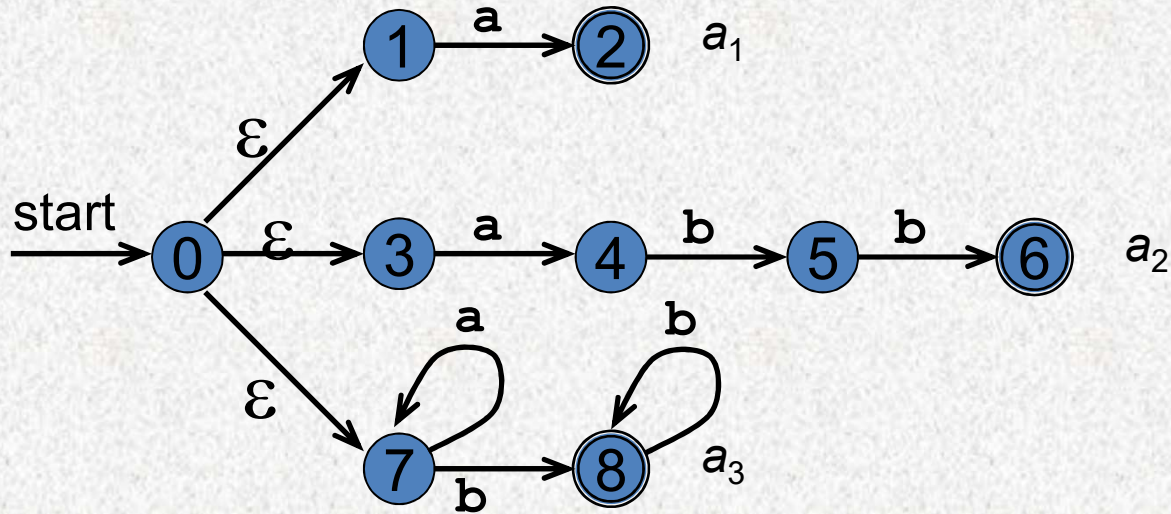


NFA STATE	DFA STATE	<i>a</i>	<i>b</i>
{0, 1, 2, 4, 7}	<i>A</i>	<i>B</i>	<i>C</i>
{1, 2, 3, 4, 6, 7, 8}	<i>B</i>	<i>B</i>	<i>D</i>
{1, 2, 4, 5, 6, 7}	<i>C</i>	<i>B</i>	<i>C</i>
{1, 2, 4, 5, 6, 7, 9}	<i>D</i>	<i>B</i>	<i>E</i>
{1, 2, 3, 5, 6, 7, 10}	<i>E</i>	<i>B</i>	<i>C</i>





# Subset Construction Example 2



*Dstates*

A = {0, 1, 3, 7}

B = {2, 4, 7}

C = {8}

D = {7}

E = {5, 8}

F = {6, 8}

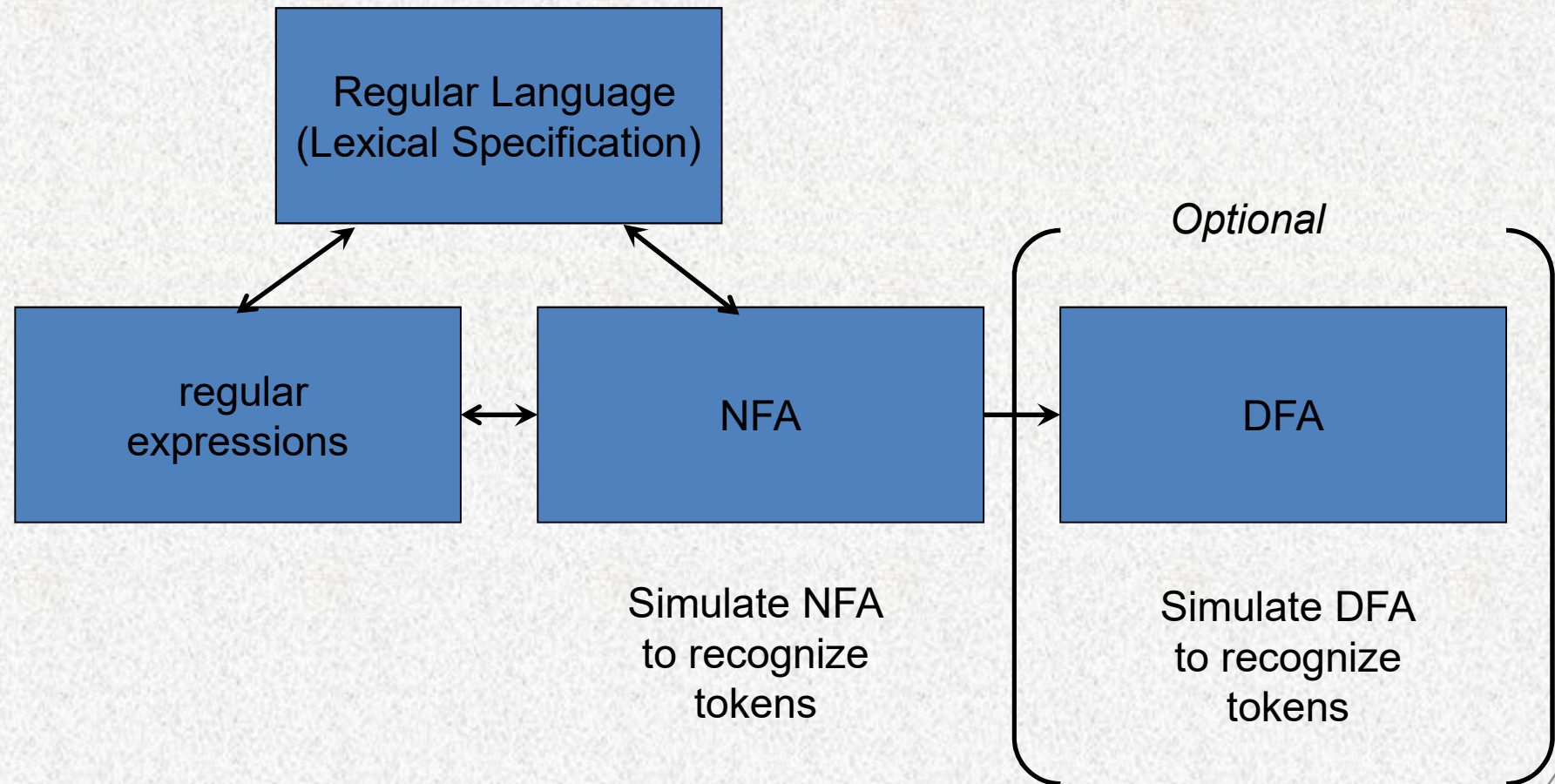


# RE to NFA/DFA



# Design of a Lexical Analyzer Generator

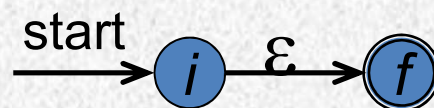
- Translate regular expressions to NFA
- Translate NFA to an efficient DFA



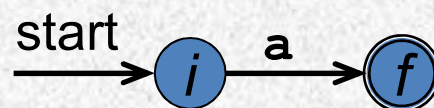


# From Regular Expression to NFA (Thompson's Construction)

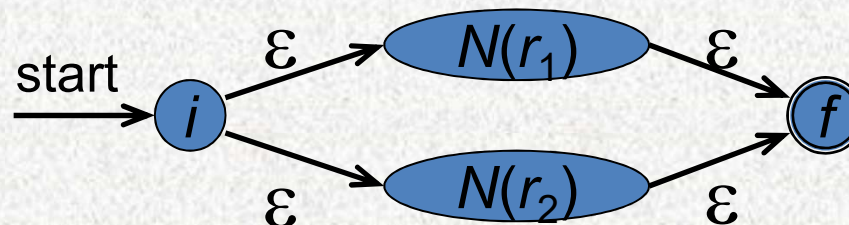
$\epsilon$



$a \in \Sigma$



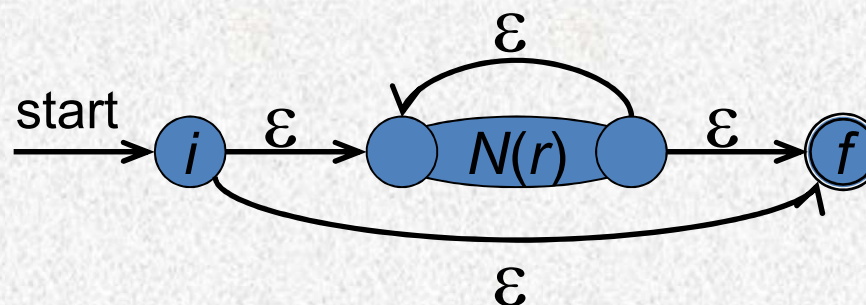
$r_1 \mid r_2$



$r_1 r_2$



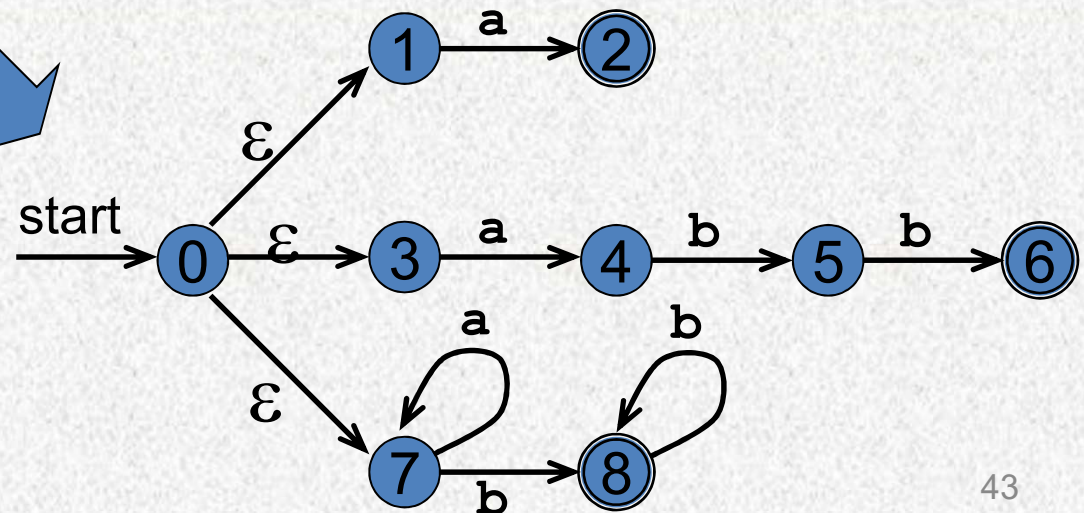
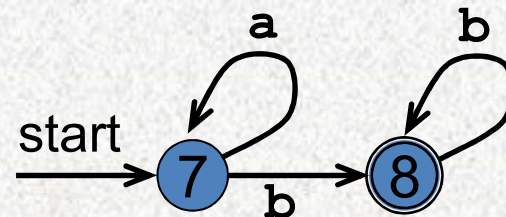
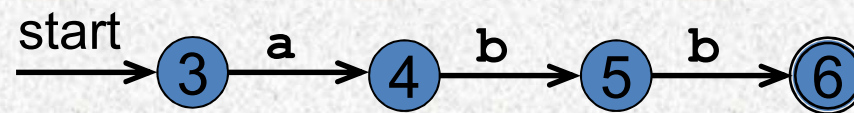
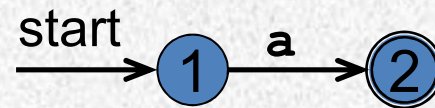
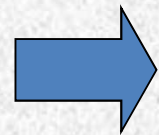
$r^*$





# Combining the NFAs of a Set of Regular Expressions

$a$        $\{ action_1 \}$   
 $abb$      $\{ action_2 \}$   
 $a^*b^+$     $\{ action_3 \}$

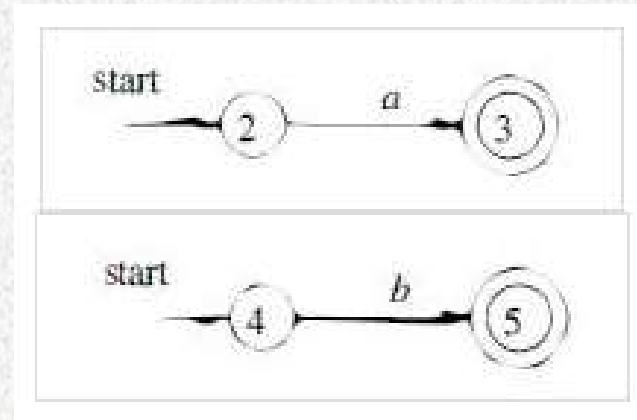




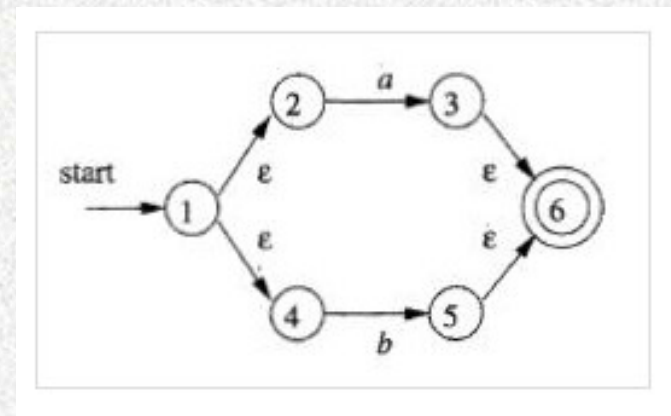
# Combining the NFAs of a Set of Regular Expressions

$r = (a|b)^*abb$

$r_1 = a$ ,  $r_2 = b$ , we have NFA:



$r_3 = r_1|r_2$ , we have NFA:

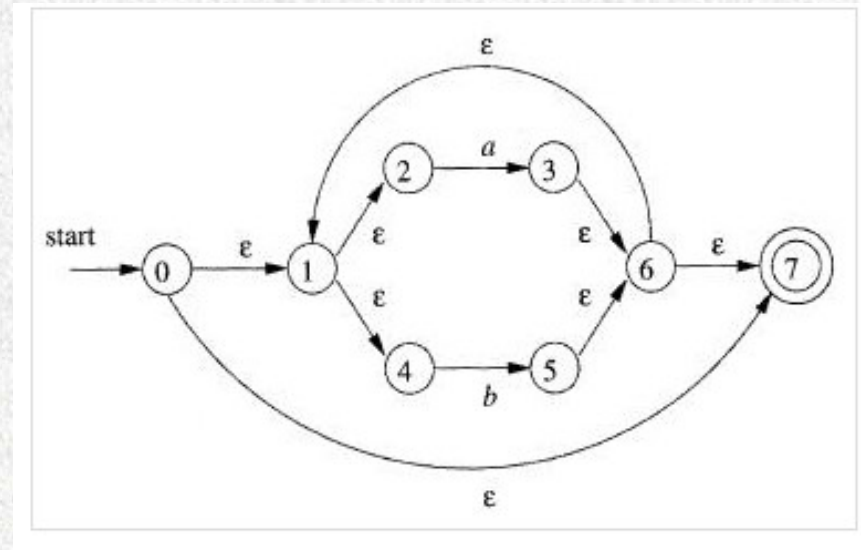




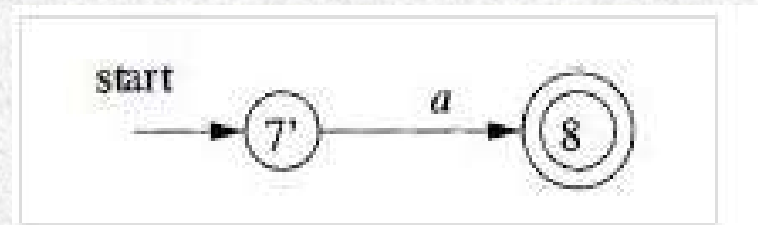
# Combining the NFAs of a Set of Regular Expressions

$r = (a|b)^*abb$

$r_5 = r_3^*$ , we have NFA:



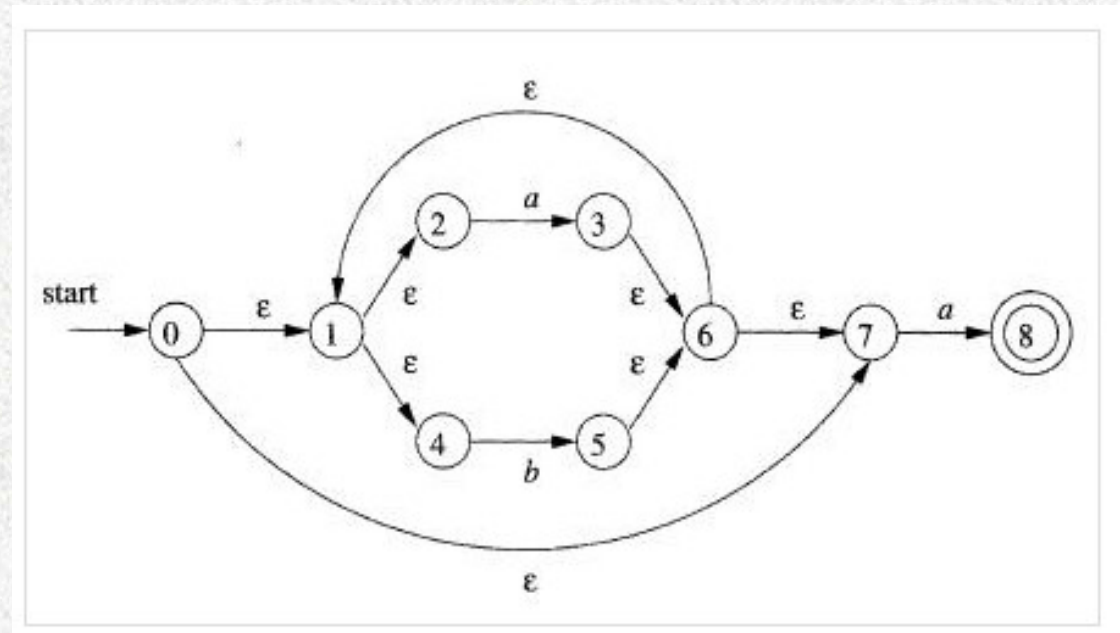
$r_6 = a$ , we have NFA:



# Combining the NFAs of a Set of Regular Expressions

$r = (a|b)^*abb$

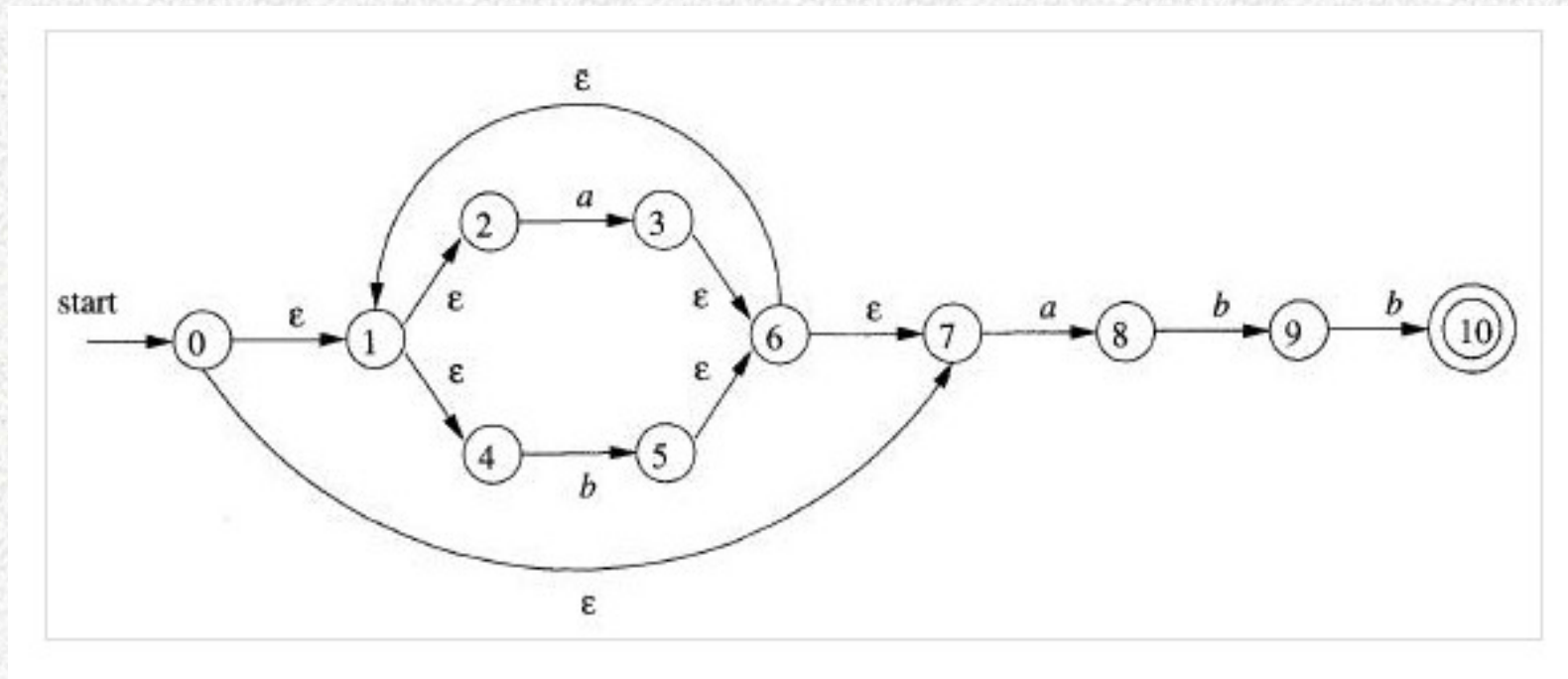
$r_7 = r_5 r_6$ , we have NFA:





# Combining the NFAs of a Set of Regular Expressions

$r = (a|b)^*abb$



# Simulating the NFA

**Algorithm 3.22:** Simulating an NFA.

**INPUT:** An input string  $x$  terminated by an end-of-file character **eof**. An NFA  $N$  with start state  $s_0$ , accepting states  $F$ , and transition function  $move$ .

**OUTPUT:** Answer “yes” if  $M$  accepts  $x$ ; “no” otherwise.

**METHOD:** The algorithm keeps a set of current states  $S$ , those that are reached from  $s_0$  following a path labeled by the inputs read so far. If  $c$  is the next input character, read by the function  $nextChar()$ , then we first compute  $move(S, c)$  and then close that set using  $\epsilon$ -closure(). The algorithm is sketched in Fig. 3.37.

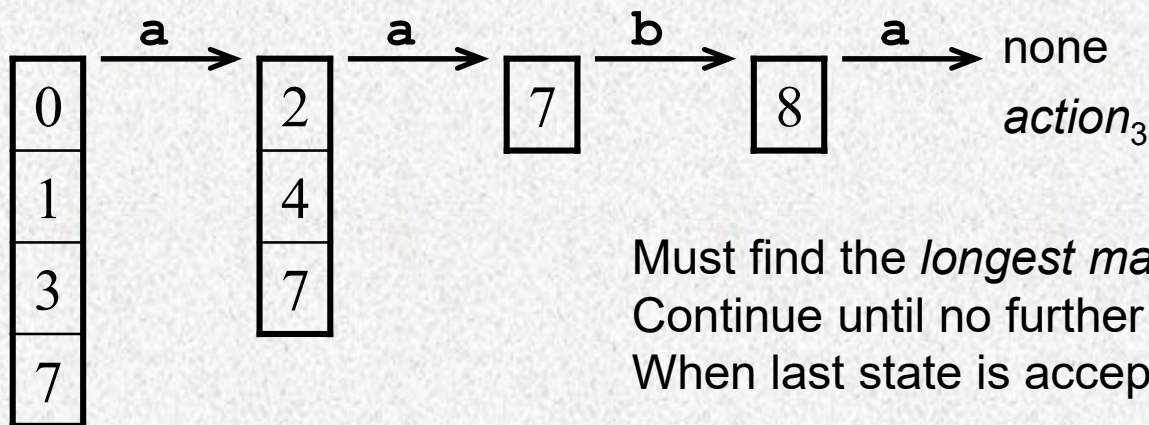
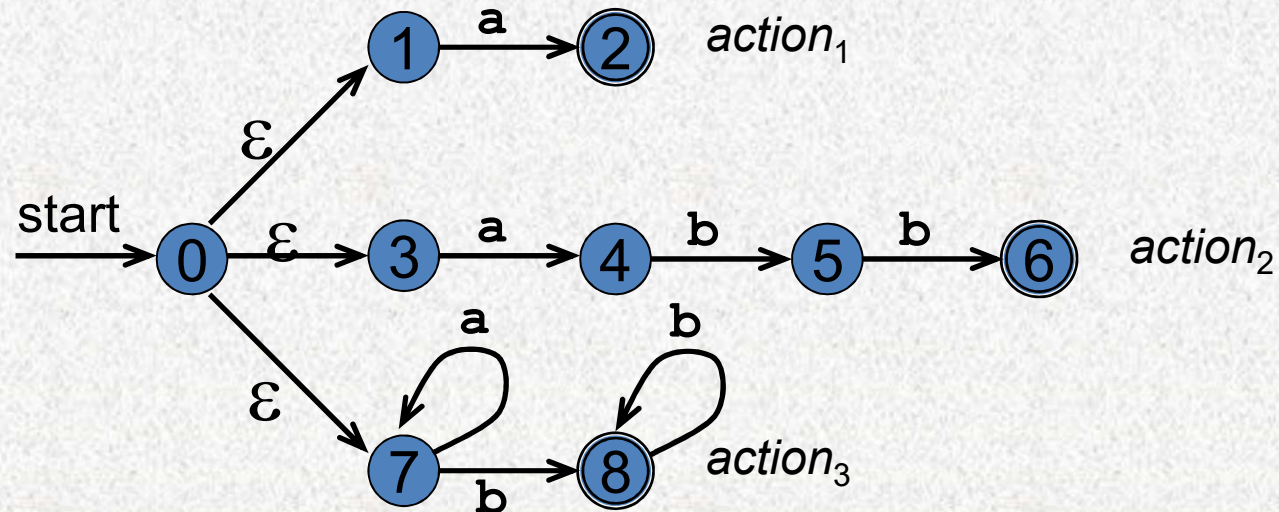
□

```
1)   $S = \epsilon\text{-closure}(s_0);$ 
2)   $c = nextChar();$ 
3)  while (  $c \neq eof$  ) {
4)       $S = \epsilon\text{-closure}(move(S, c));$ 
5)       $c = nextChar();$ 
6)  }
7)  if (  $S \cap F \neq \emptyset$  ) return "yes";
8)  else return "no";
```

Figure 3.37: Simulating an NFA



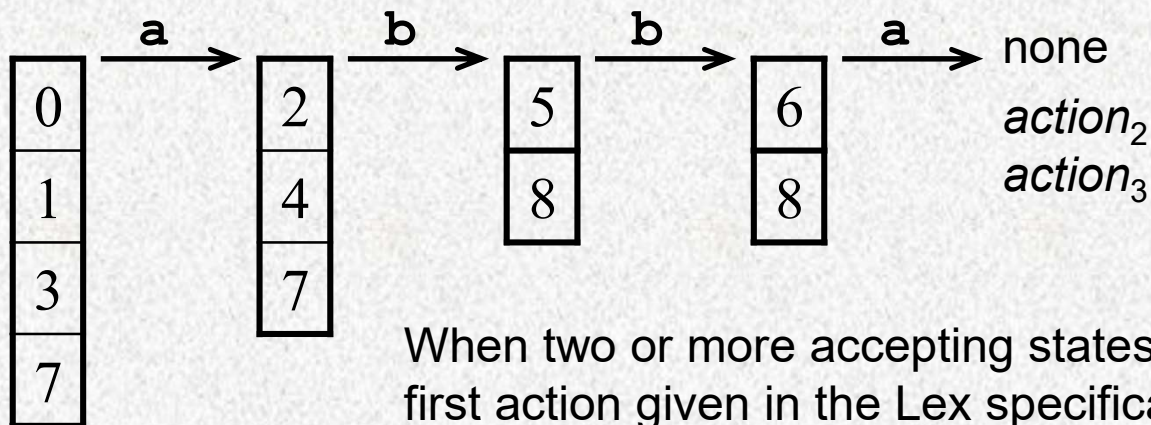
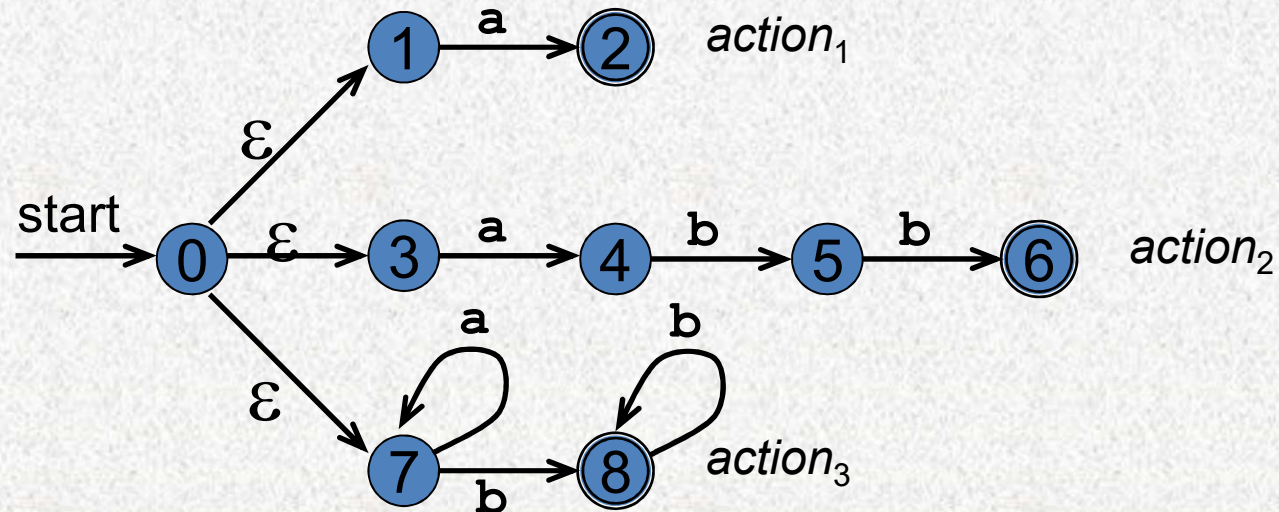
# Simulating a NFA Example 1



Must find the *longest match*:  
 Continue until no further moves are possible  
 When last state is accepting: execute action



# Simulating a NFA Example 2



When two or more accepting states are reached, the first action given in the Lex specification is executed

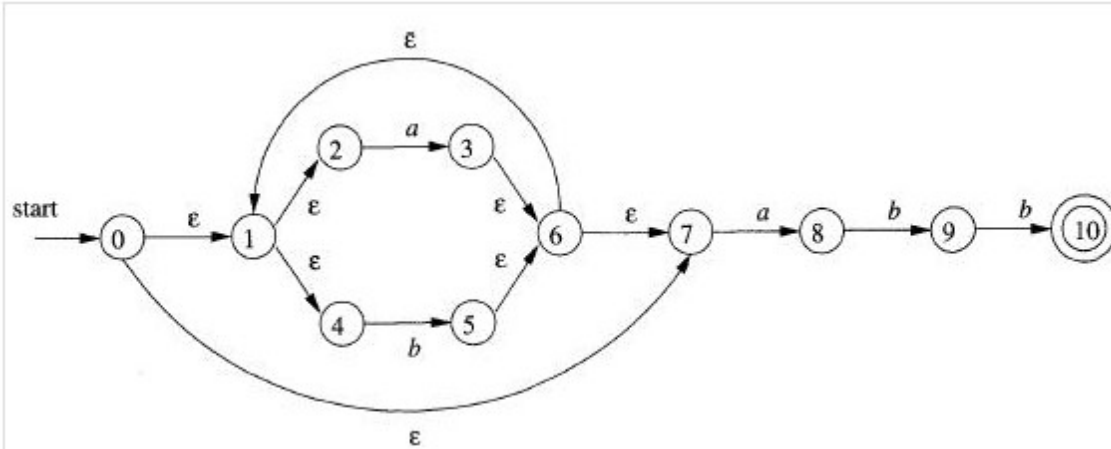


# From NFA to DFA

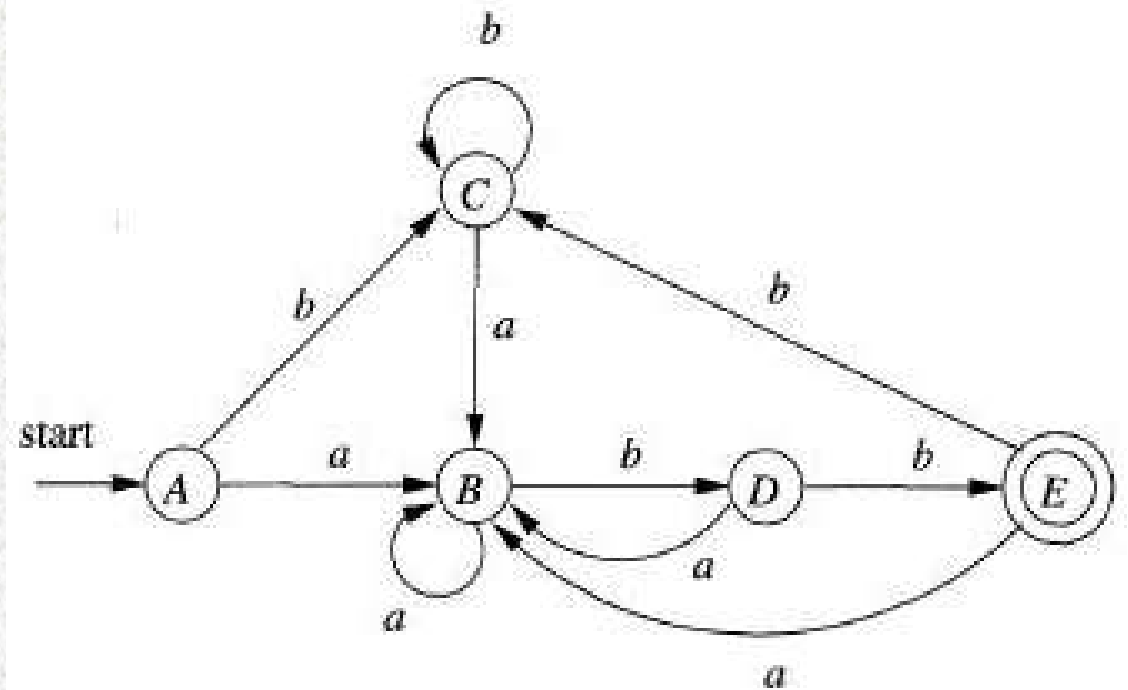
## The Subset Construction Algorithm

```
initially,  $\epsilon$ -closure( $s_0$ ) is the only state in  $Dstates$ , and it is unmarked;  
while ( there is an unmarked state  $T$  in  $Dstates$  ) {  
    mark  $T$ ;  
    for ( each input symbol  $a$  ) {  
         $U = \epsilon$ -closure( $move(T, a)$ );  
        if (  $U$  is not in  $Dstates$  )  
            add  $U$  as an unmarked state to  $Dstates$ ;  
         $Dtran[T, a] = U$ ;  
    }  
}
```

# From NFA to DFA



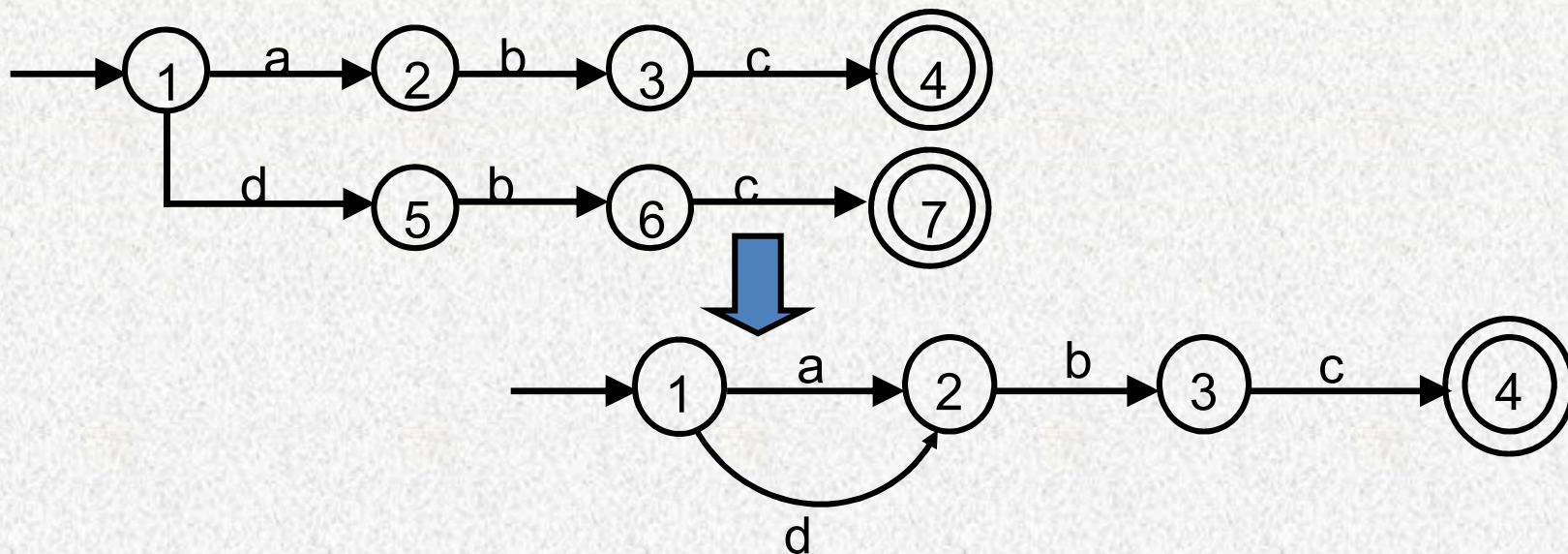
NFA STATE	DFA STATE	a	b
{0, 1, 2, 4, 7}	A	B	C
{1, 2, 3, 4, 6, 7, 8}	B	B	D
{1, 2, 4, 5, 6, 7}	C	B	C
{1, 2, 4, 5, 6, 7, 9}	D	B	E
{1, 2, 3, 5, 6, 7, 10}	E	B	C





# Minimizing DFA

After conversion from NFA, the DFA may contain some equivalent states, which lead to low efficiency in the analysis





# Minimizing DFA

- ***Lots*** of methods
- All involve finding **equivalent states**:
  - States that go to equivalent states under all inputs (sounds recursive)
- We will use the ***Partitioning Method***



# Minimizing DFA

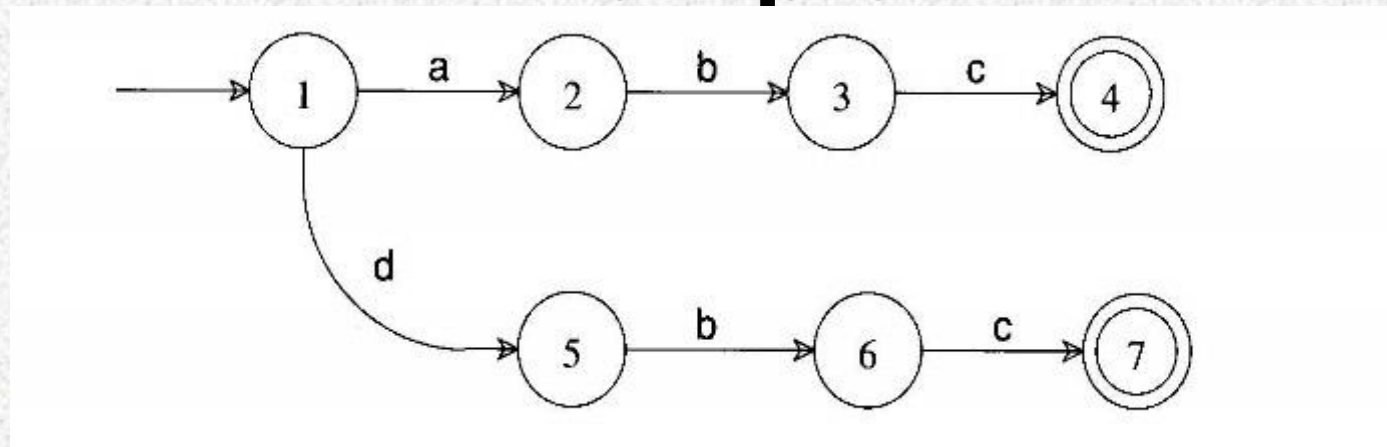
- Step 1
  - Start with an initial partition  $\Pi$  with two group: F and S-F (accepting and nonaccepting)
- Step 2
  - Split Procedure
- Step 3
  - If (  $\Pi_{\text{new}} = \Pi$  )  
     $\Pi_{\text{final}} = \Pi$  and continue step 4  
    else  
         $\Pi = \Pi_{\text{new}}$  and go to step 2
- Step 4
  - Construct the minimum-state DFA by  $\Pi_{\text{final}}$  group.
  - Delete the dead state

# Split Procedure

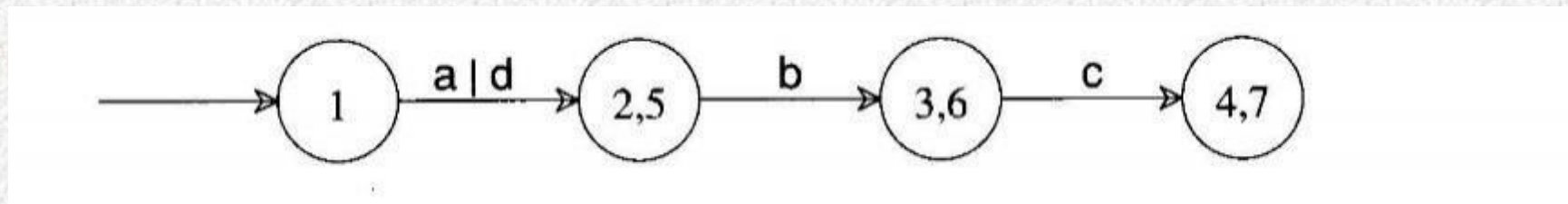
```
initially, let  $\Pi_{\text{new}} = \Pi$ ;  
for ( each group  $G$  of  $\Pi$  ) {  
    partition  $G$  into subgroups such that two states  $s$  and  $t$   
        are in the same subgroup if and only if for all  
        input symbols  $a$ , states  $s$  and  $t$  have transitions on  $a$   
        to states in the same group of  $\Pi$ ;  
    /* at worst, a state will be in a subgroup by itself */  
    replace  $G$  in  $\Pi_{\text{new}}$  by the set of all subgroups formed;  
}
```



# Example



- initially, two sets  $\{1, 2, 3, 5, 6\}, \{4, 7\}$ .
- $\{1, 2, 3, 5, 6\}$  splits  $\{1, 2, 5\}, \{3, 6\}$  on  $c$ .
- $\{1, 2, 5\}$  splits  $\{1\}, \{2, 5\}$  on  $b$ .



# Minimizing the DFA

- Major operation: partition states into equivalent classes according to
  - final / non-final states
  - transition functions

	a	b
A	B	C
B	B	D
C	B	C
D	B	E
Ⓔ	B	C



(A B C D E)  
 (A B C D)(E)  
 (A B C)(D)(E)  
 (A C)(B)(D)(E)



	a	b
A C	B	A C
B	B	D
D	B	E
Ⓔ	B	A C



# Minimizing the DFA

- DFA  $D = (\{0,1,2,3,4,5\}, \{a,b\}, \delta, 0, \{0,1\})$ , 其中  $\delta$  见表

states	a	b
0	1	2
1	1	4
2	1	3
3	3	2
4	0	5
5	5	4

Step 1:  
 $A = \{0,1\}$ ,  $B = \{2,3,4,5\}$ 。

States	partition	a	b
0	A	1(A)	2(B)
1	A	1(A)	4(B)
2	B	1(A)	3(B)
3	B	3(B)	2(B)
4	B	0(A)	5(B)
5	B	5(B)	4(B)

# Minimizing the DFA

- DFA  $D = (\{0,1,2,3,4,5\}, \{a,b\}, \delta, 0, \{0,1\})$ ,

Cannot be divided any more

stat es	parti tion	a	b
0	A	1(A)	2(B)
1	A	1(A)	4(B)
2	B	1(A)	3(B)
3	B	3(B)	2(B)
4	B	0(A)	5(B)
5	B	5(B)	4(B)



sta tes	a	b
0	1	2
1	1	4
2	1	3
3	3	2
4	0	5
5	5	4

stat es	partit ion	a	b
0	A	1(A)	2(B)
1	A	1(A)	4(B)
2	B	1(A)	3(C)
3	C	3(C)	2(B)
4	B	0(A)	5(C)
5	C	5(C)	4(B)



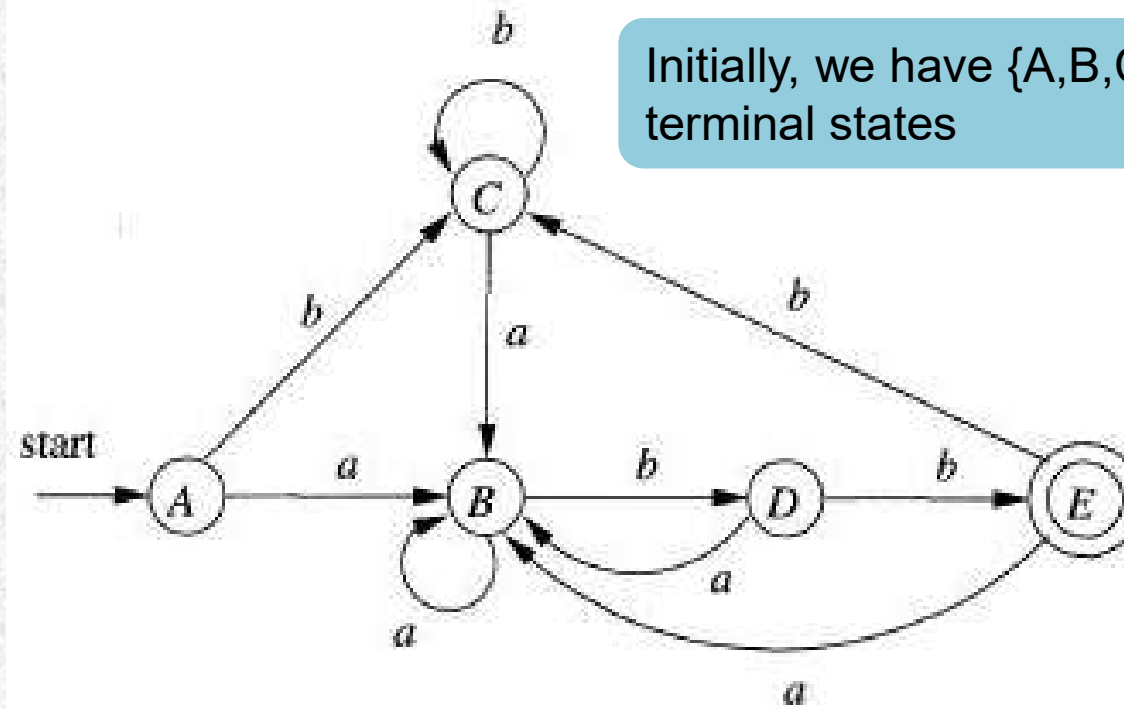
# Minimizing the DFA

- DFA  $D = (\{0,1,2,3,4,5\}, \{a,b\}, \delta, 0, \{0,1\})$  is minimized to:  
DFA  $D' = (\{A,B,C\}, \{a,b\}, \delta, A, \{A\})$ , where  $\delta$  is defined as follows

state	a	b
A	A	B
B	A	C
C	C	B

# Minimizing the DFA-Example

- $r=(a|b)^*abb$



Initially, we have  $\{A,B,C,D\},\{E\}$ , which are for non-terminal and terminal states

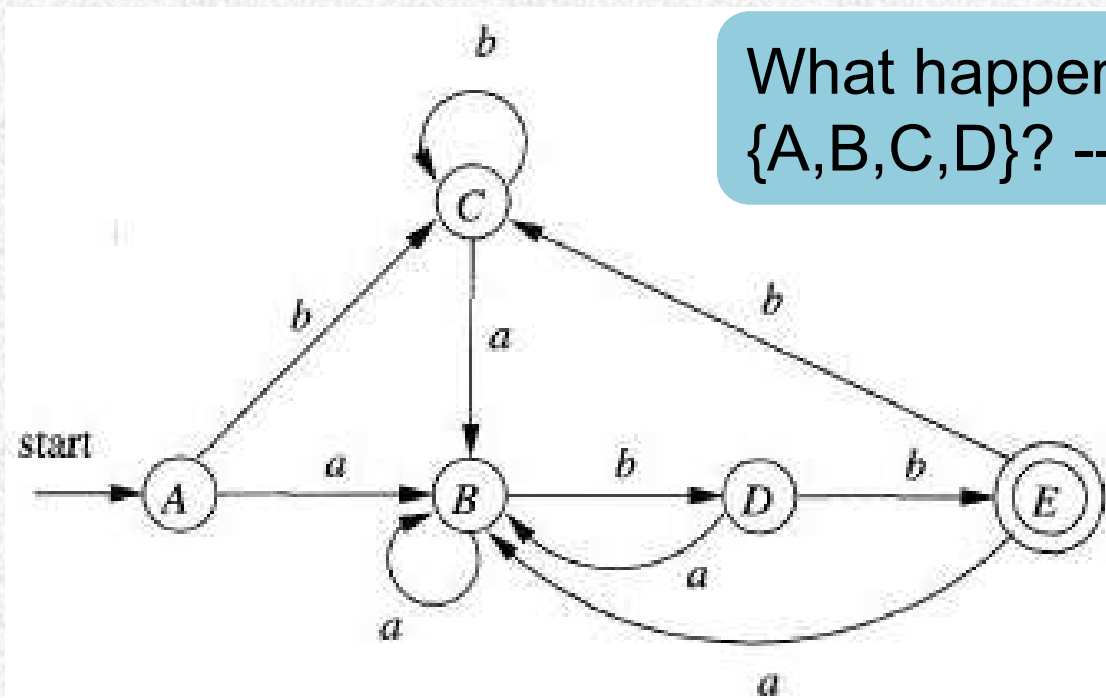
$\{E\}$  is not dividable, so we only consider  $\{A, B, C, D\}$



# Minimizing the DFA-Example

Is  $\{A, B, C, D\}$  dividable?

■  $r = (a|b)^*abb$



What happens when take in a under  $\{A, B, C, D\}$ ? --- still with  $\{A, B, C, D\}$

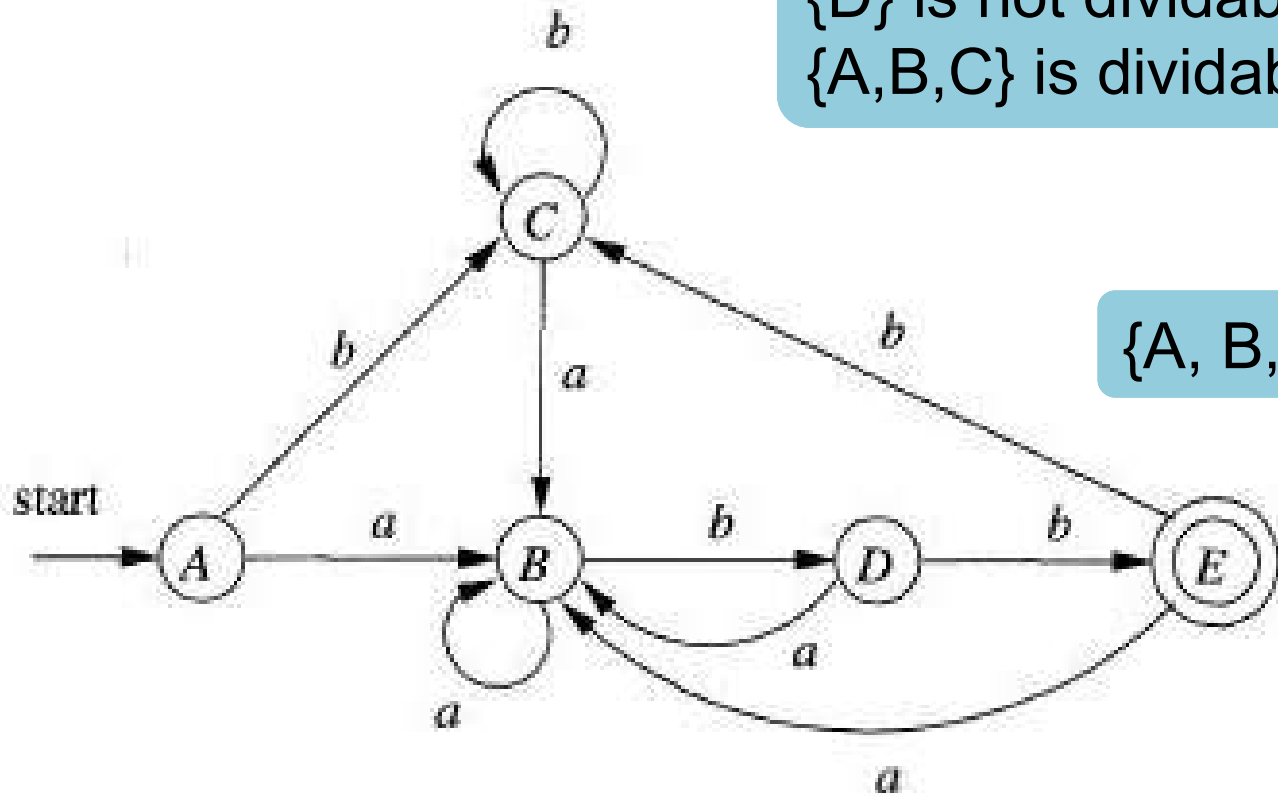
What happens when take in b under  $\{A, B, C, D\}$ ? --- becomes  $\{A, B, C\}, \{D\}$

# Minimizing the DFA-Example

■  $r = (a|b)^*abb$

{D} is not dividable, so let us see whether {A,B,C} is dividable?

{A, B, C} becomes {A,C},{B}

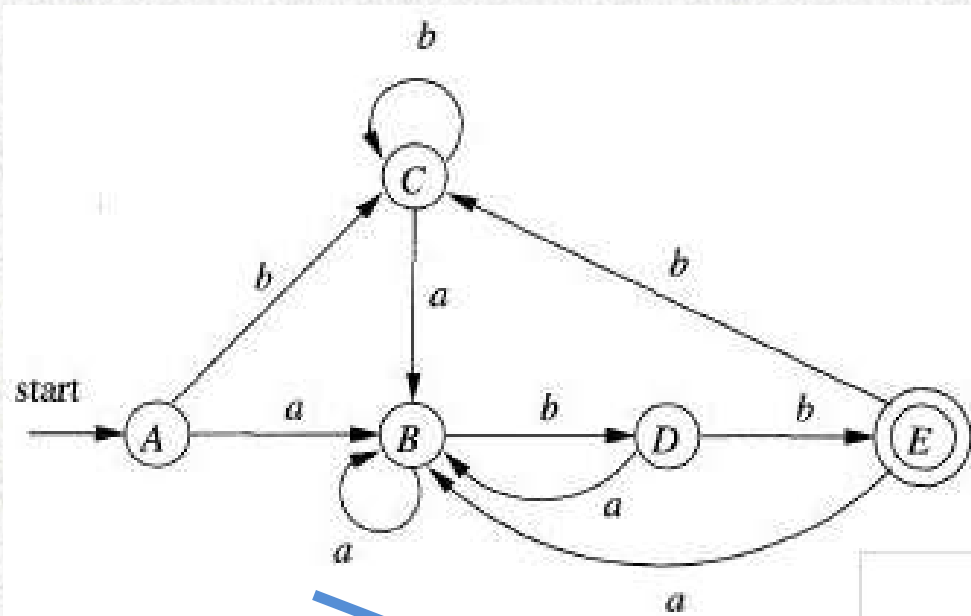




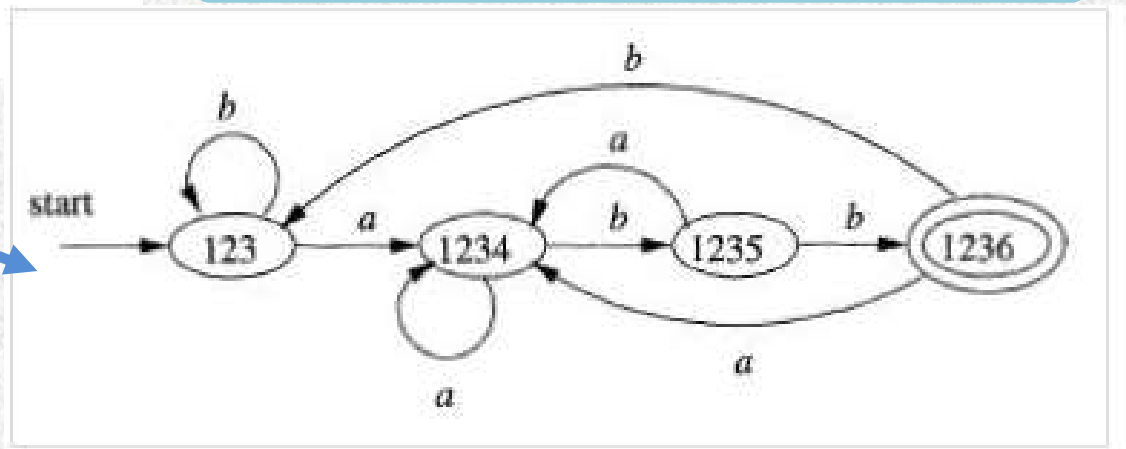
# Minimizing the DFA-Example

■  $r=(a|b)^*abb$

{A,C} is dividable?

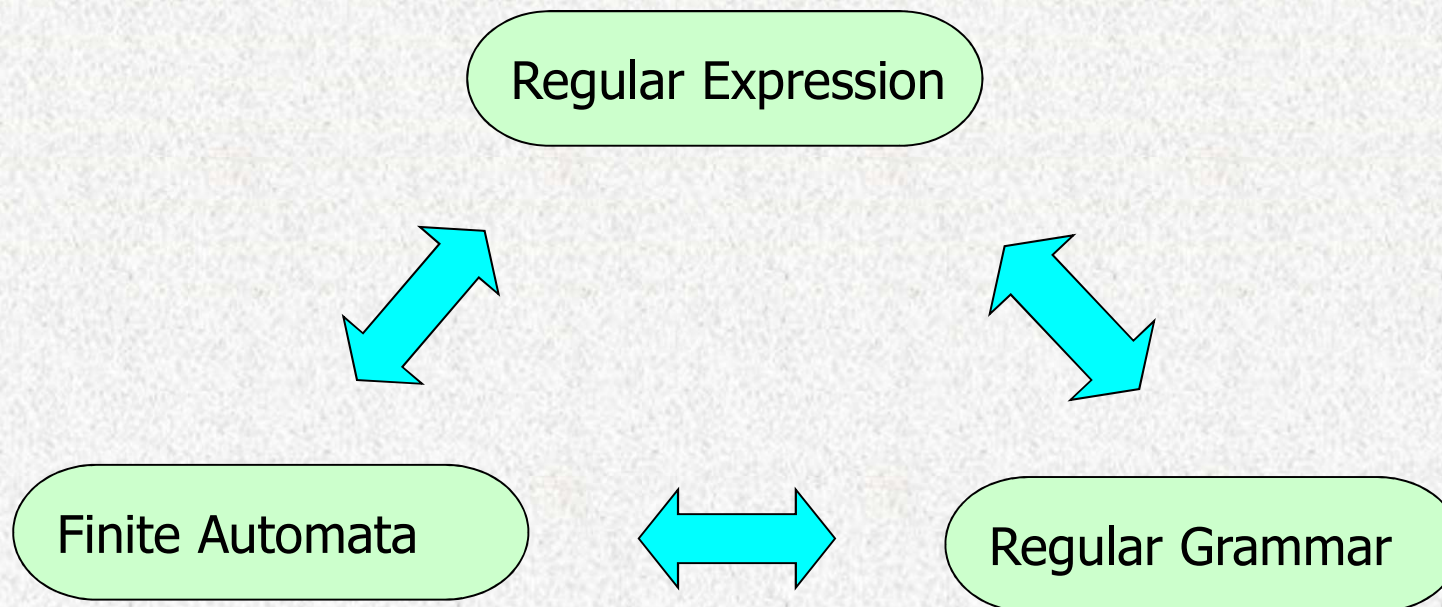


Finally, we have {A,C},{B},{D},{E}



# RE v.s. NFA/DFA

- RE, DFA(NFA),  $L(RE)$  are equivalent to each other





# Exercise

## ■ Given an NFA N

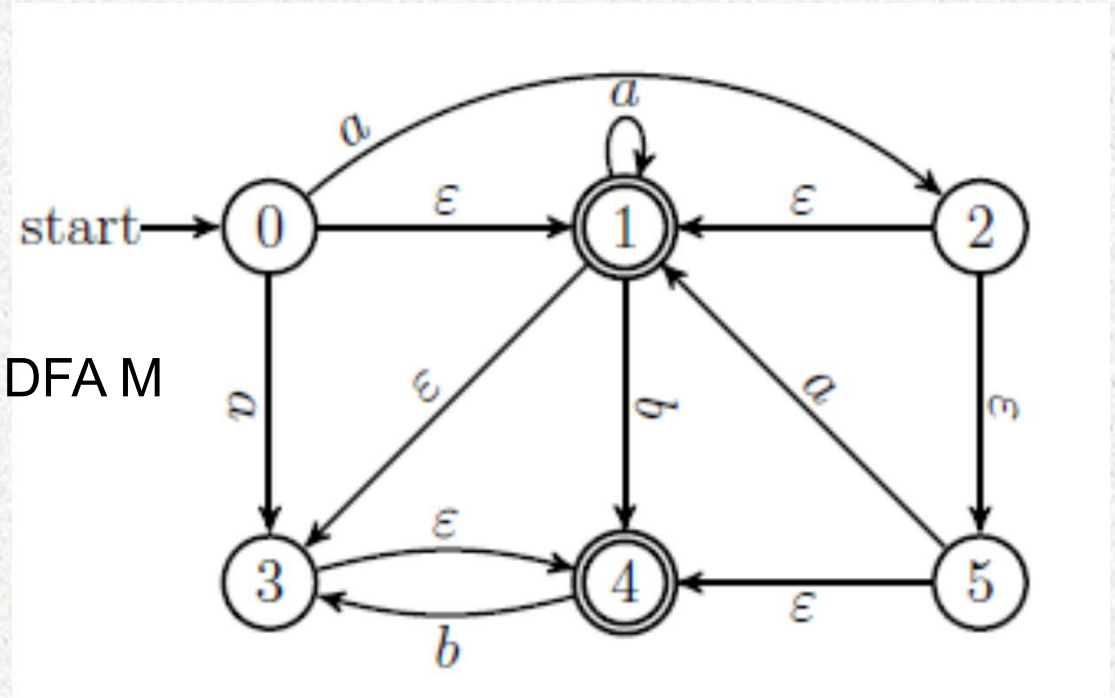
(1) Simulate the NFA on input “aaabb”

(2) Convert the NFA N to its equivalent DFA M

(3) Minimize the DFA M

(4) Describe what can this DFA/NFA accept in natural language

(5) Write down the regular expression  $re$ , such that  $L(re) = L(N)$





---

# Homework-W3

---



# Homework – week 3

- pp. 125, Exercise 3.3.5 (c)(d)(f)(h)
- pp.152, Exercise 3.6.5
- pp. 166, Exercise 3.7.1 (b), Exercise 3.7.2 (b), Exercise 3.7.3 (d)
- pp. 172, Exercise 3.8.1
- pp.187, Exercise 3.9.4

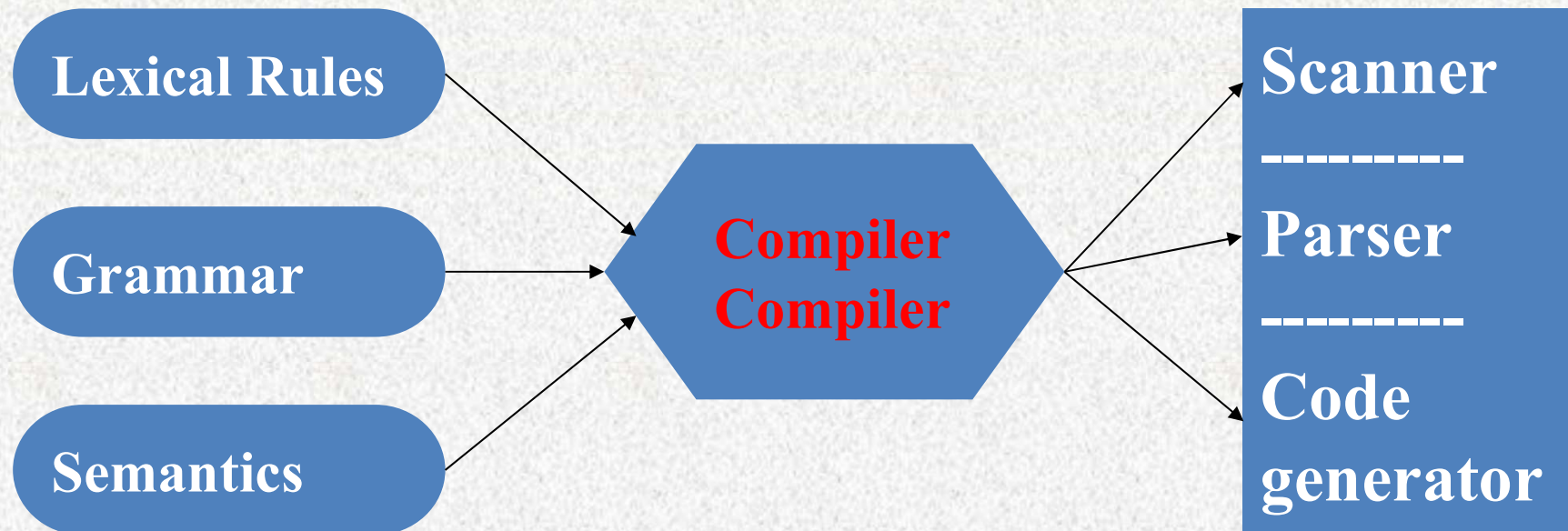


# Lexical Analyzer Implementation

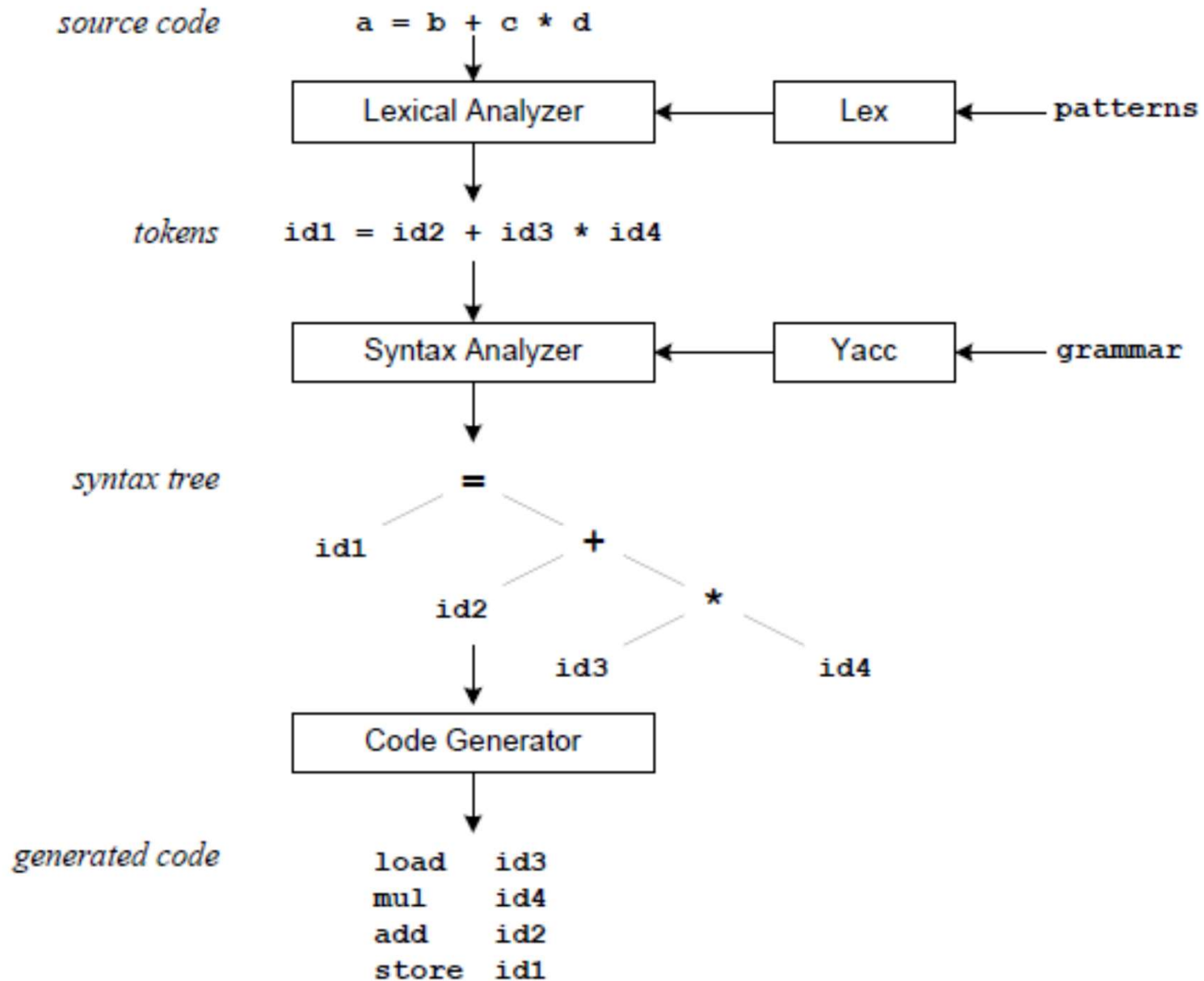


# Overview

- Writing a compiler is difficult requiring lots of time and effort
- Construction of the scanner and parser is routine enough that the process may be automated



# Overview





# LEX

- Lex is a scanner generator
  - Input is description of patterns and actions
  - Output is a C program which contains a function `yylex()` which, when called, matches patterns and performs actions per input
  - Typically, the generated scanner performs lexical analysis and produces tokens for the (YACC-generated) parser

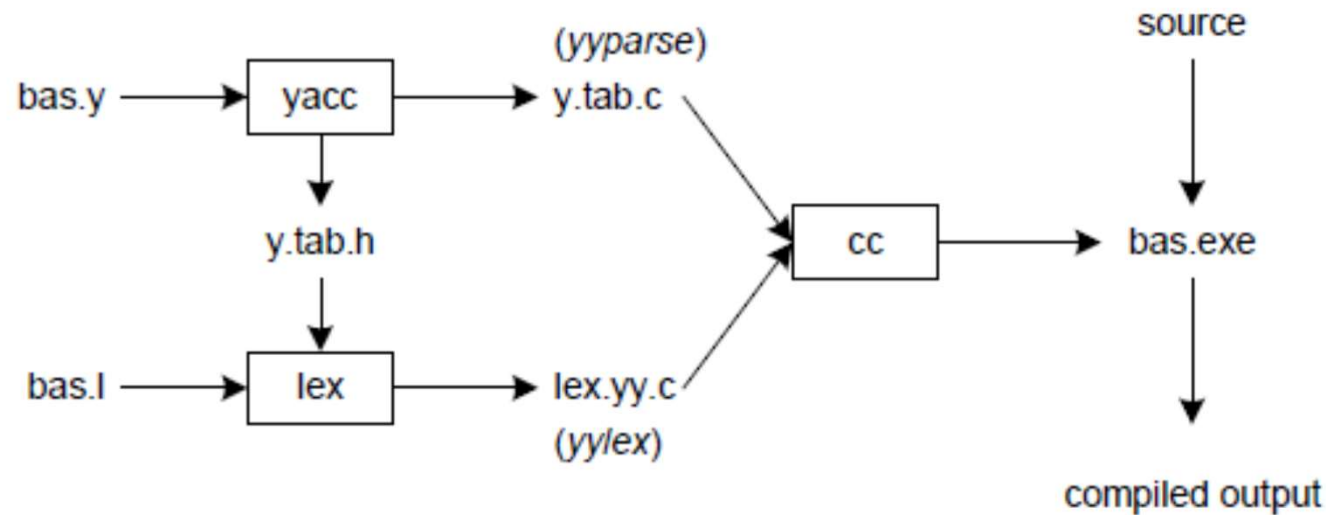


# YACC

- What is **YACC** ?
  - **Tool which will produce a parser for a given grammar.**
  - YACC (Yet Another Compiler Compiler) is a program designed to compile a LALR(1) grammar and to produce the source code of the syntactic analyzer of the language produced by this grammar
  - Input is a grammar (rules) and actions to take upon recognizing a rule
  - Output is a C program and optionally a header file of tokens



# LEX and YACC: a team



**Figure 2:** Building a Compiler with Lex/Yacc

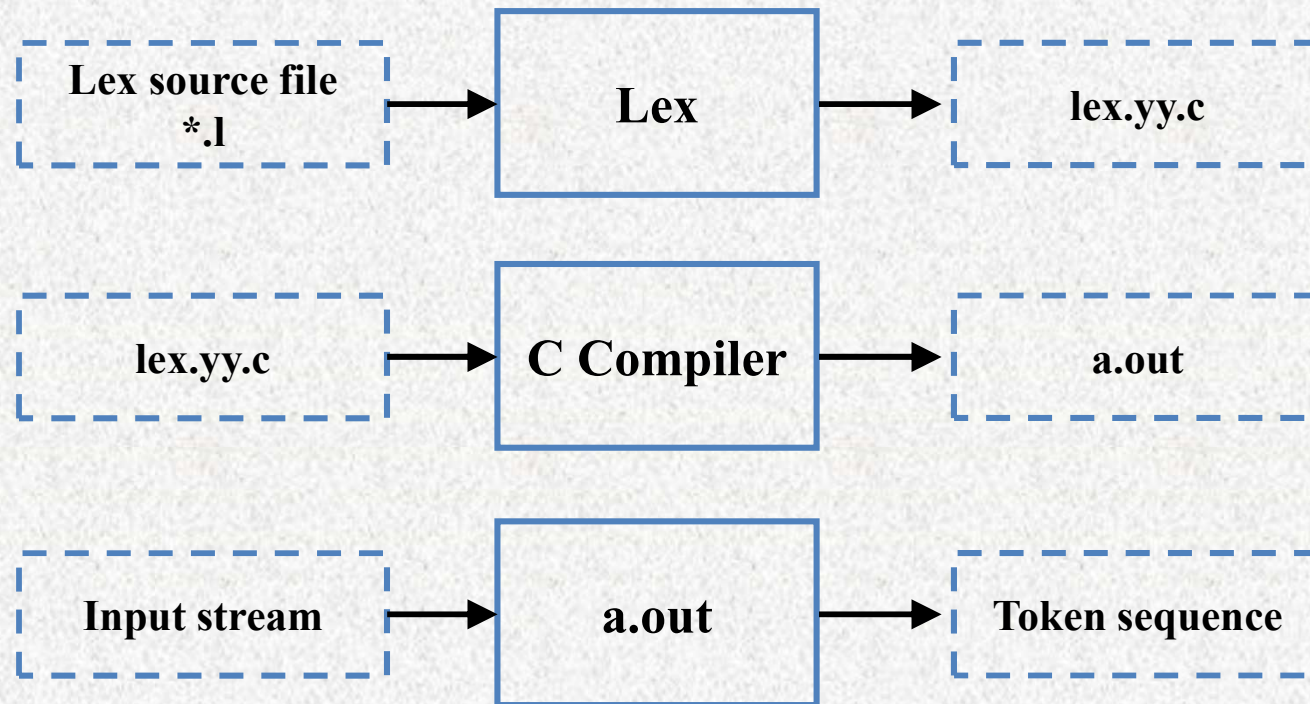
```
yacc -d bas.y          # create y.tab.h, y.tab.c
lex bas.l              # create lex.yy.c
cc lex.yy.c y.tab.c -obas.exe  # compile/link
```

# Availability

- lex, yacc on most UNIX systems
- bison: a yacc replacement from GNU
- flex: *fast lexical* analyzer
- BSD yacc
- Windows/MS-DOS versions exist



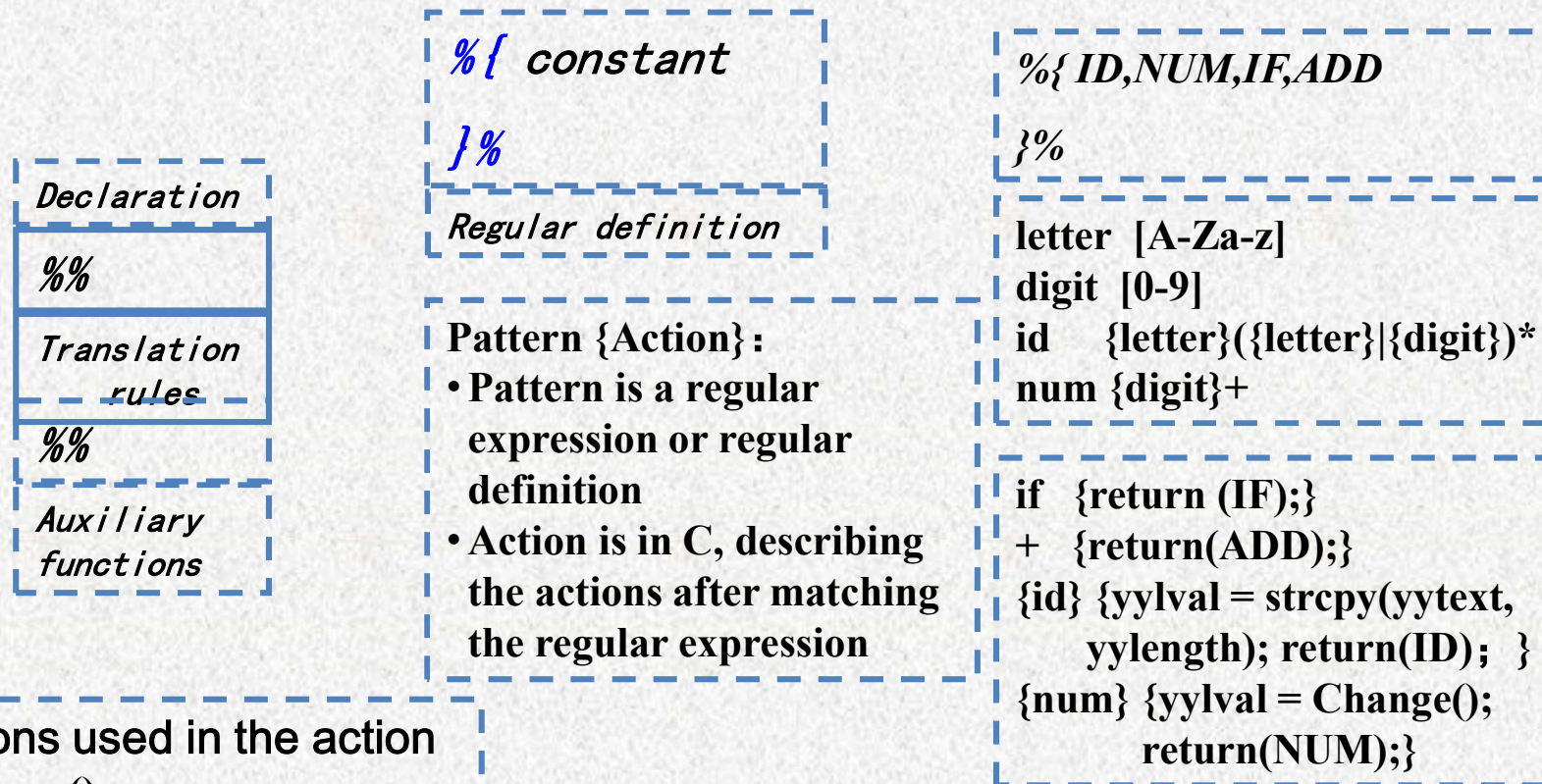
# Lex



Create your lexical analyzer with Lex



# Structure of Lex source file



Functions used in the action

```
int Change()
{ /*Convert string into integer*/
}
```

**yylval:** value of the token  
**yytext:** lexeme of the token  
**yylen:** length of the lexeme



# Example: LEX

```
%{
#include <stdio.h>
#include "y.tab.h"
%}

id      [_a-zA-Z][_a-zA-Z0-9]*
wspc    [ \t\n]+
semi    [;]
comma   [,]

%%

int      { return INT; }
char     { return CHAR; }
float    { return FLOAT; }
{comma}  { return COMMA; }      /* Necessary? */
{semi}   { return SEMI; }
{id}     { return ID; }
{wspc}   { ; }
```



# Example: Definitions

```
%{  
#include <stdio.h>  
#include <stdlib.h>  
%}  
%start  line  
%token  CHAR, COMMA, FLOAT, ID, INT, SEMI  
%%
```



# Example: Rules

```
/* This production is not part of the "official"  
 * grammar. It's primary purpose is to recover from  
 * parser errors, so it's probably best if you leave  
 * it here. */
```

```
line : /* lambda */  
      | line decl  
      | line error {  
          printf("Failure :-(\n");  
          yyerrok;  
          yyclearin;  
      }  
      ;
```



# Example: Rules

```
decl : type ID list { printf("Success!\n"); } ;
```

```
list : COMMA ID list  
      | SEMI  
      ;
```

```
type : INT | CHAR | FLOAT  
      ;
```

```
%%
```



# Example: Supplementary Code

```
extern FILE *yyin;
main()
{
    do {
        yyparse();
    } while(!feof(yyin));
}
yyerror(char *s)
{
    /* Don't have to do anything! */
}
```

# Next Time

