



中山大學  
SUN YAT-SEN UNIVERSITY

计算机学院（软件学院）  
SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

# Compilation Principle 编译原理

---

## 第3讲：词法分析(3)

张献伟

[xianweiz.github.io](https://xianweiz.github.io)

DCS290, 3/2/2023

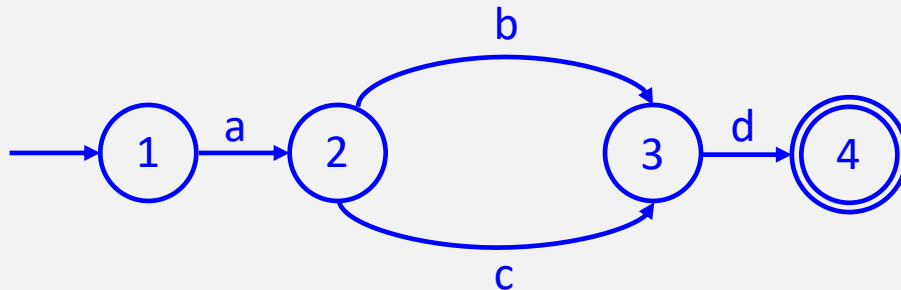
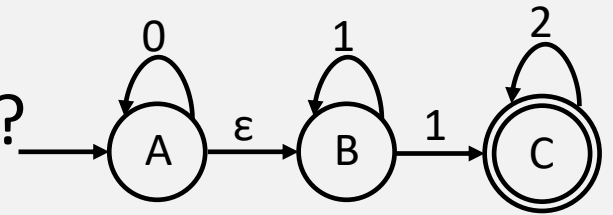


中山大學  
SUN YAT-SEN UNIVERSITY



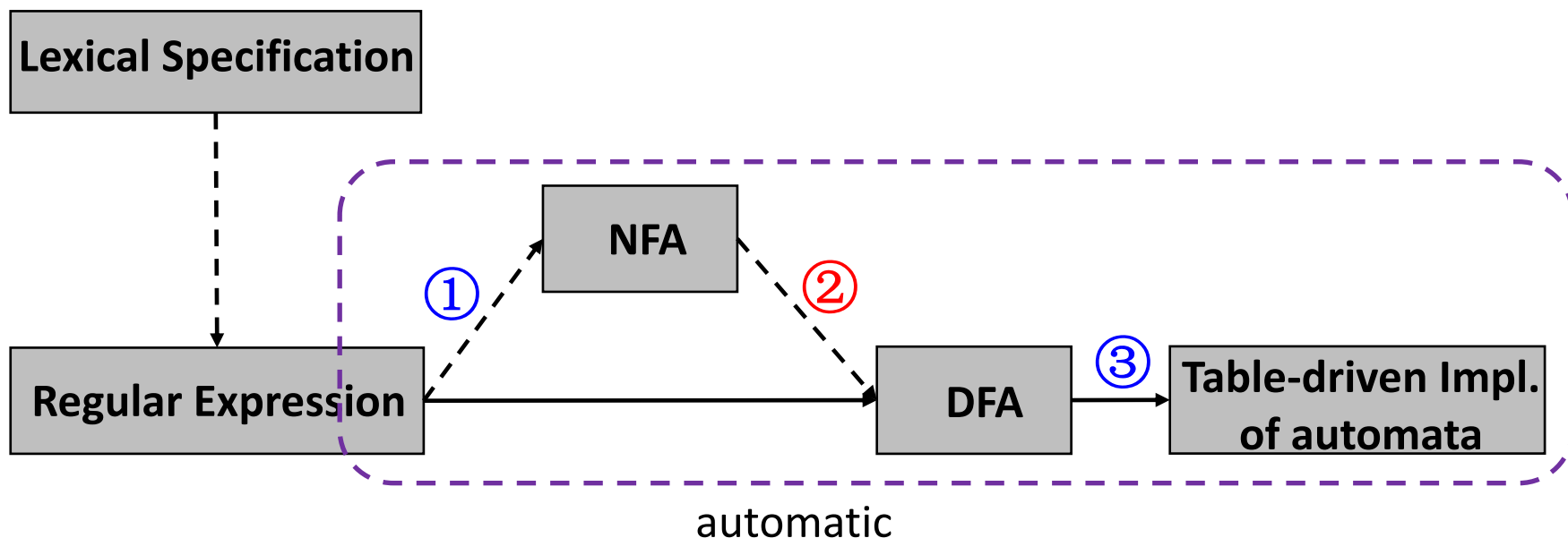
# Review Questions

- Q1: usage of RE and FA in lexical analysis?  
RE: specify the token pattern; FA: implement the token recognizer
- Q2: the general workflow of RE to implementation?  
RE  $\rightarrow$  NFA  $\rightarrow$  DFA  $\rightarrow$  Table
- Q3: the graph describes NFA or DFA? Why?  
NFA. A:  $\epsilon$ -transition, B: 1-transition
- Q4: what's the use of M-Y-T algorithm?  
To convert RE to NFA.
- Q5: FA for the RE  $a(b|c)d$



# The Conversion Flow

- Outline: RE  $\rightarrow$  NFA  $\rightarrow$  DFA  $\rightarrow$  Table-driven Implementation
  - ③ Converting DFAs to table-driven implementations
  - ① Converting REs to NFAs
  - ② Converting NFAs to DFAs



# NFA $\rightarrow$ DFA: Steps

---

- The **initial state** of the DFA is the set of all states the NFA can be in without reading any input
- For any state  $\{q_i, q_j, \dots, q_k\}$  of the DFA and any input  $a$ , the **next state** of the DFA is the set of all states of the NFA that can result as next states if the NFA is in any of the states  $q_i, q_j, \dots, q_k$  when it reads  $a$ 
  - This includes states that can be reached by reading  $a$  followed by any number of  $\epsilon$ -transitions
  - Use this rule to keep adding new states and transitions until it is no longer possible to do so
- The **accepting states** of the DFA are those states that contain an accepting state of the NFA

# NFA $\rightarrow$ DFA: Algorithm

Initially,  $\epsilon\text{-closure}(s_0)$  is the only state in  $Dstates$  and it is unmarked

```
while there is an unmarked state  $T$  in  $Dstates$  do  
    mark  $T$   
    for each input symbol  $a \in \Sigma$  do  
         $U := \epsilon\text{-closure}(\text{move}(T, a))$   
        if  $U$  is not in  $Dstates$  then  
            add  $U$  as an unmarked state to  $Dstates$   
        end if  
         $Dtran[T, a] := U$   
    end do  
end do
```

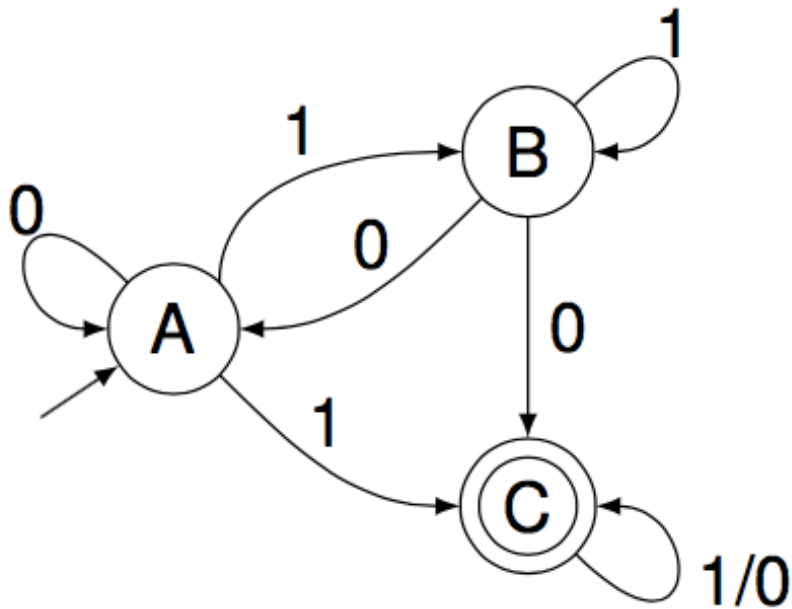
- Operations on NFA states:

- $\epsilon\text{-closure}(s)$ : 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**;  $= \bigcup_{s \in T} \epsilon\text{-closure}(s)$
- $\text{move}(T, a)$ : set of NFA states to which there is a transition on input symbol  $a$  from some state  $s$  in  $T$

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon$ -closure(A) = A
- Keep getting  $\varepsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

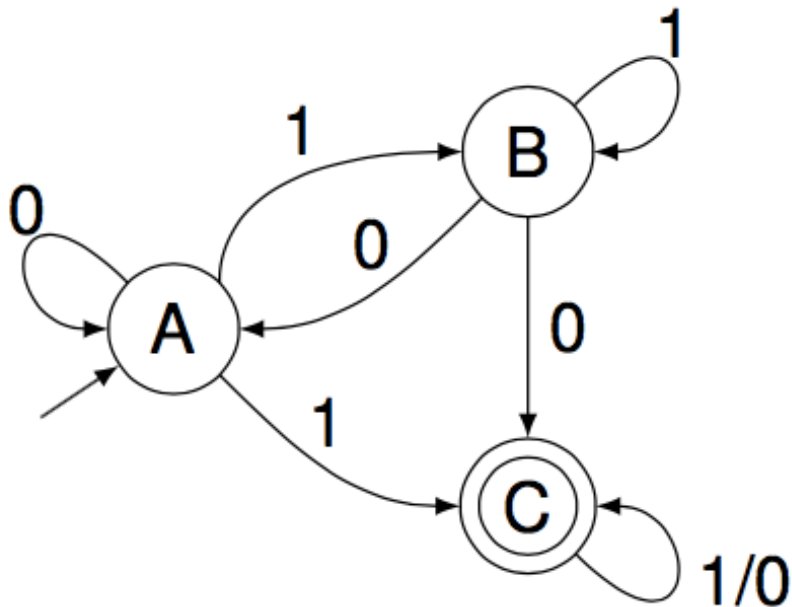
$T: A, a: 0/1$



# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon$ -closure(A) = A
- Keep getting  $\varepsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

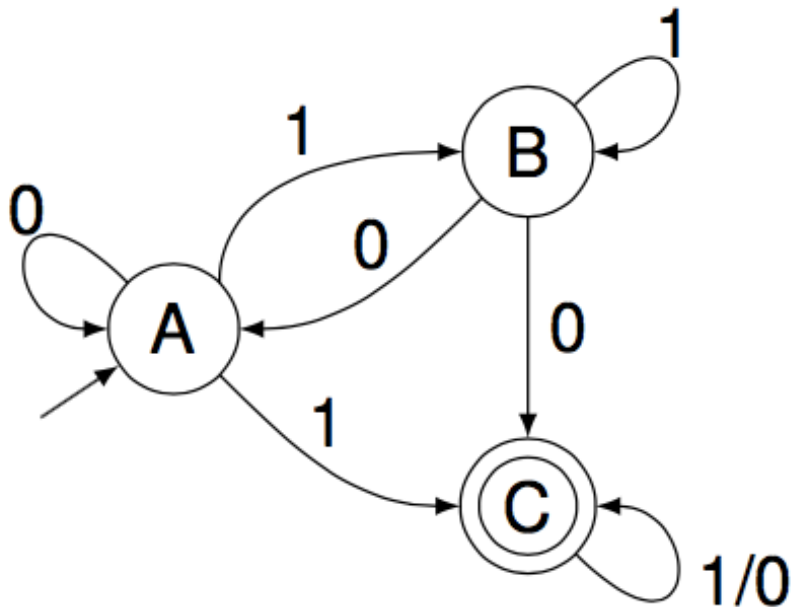


	0	1
A		

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\epsilon$ -closure of the start state
  - $\epsilon$ -closure(A) = A
- Keep getting  $\epsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$



alphabet  $\rightarrow$

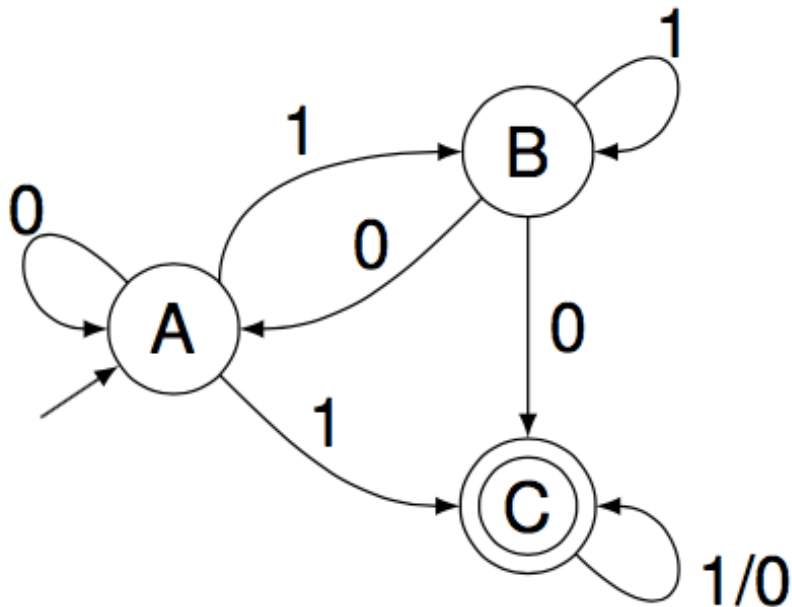
	0	1
A		



# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon$ -closure(A) = A
- Keep getting  $\varepsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

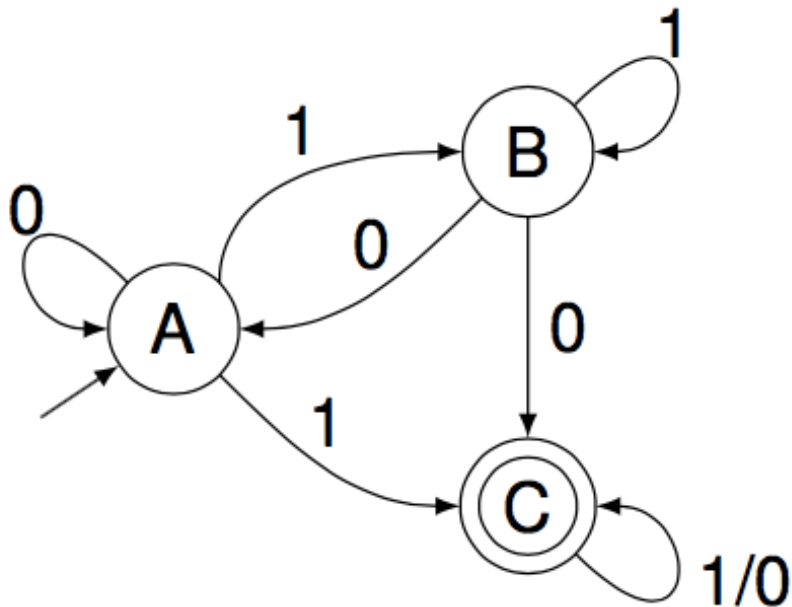


alphabet $\rightarrow$		
state $\downarrow$	0	1
	A	

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon$ -closure(A) = A
- Keep getting  $\varepsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

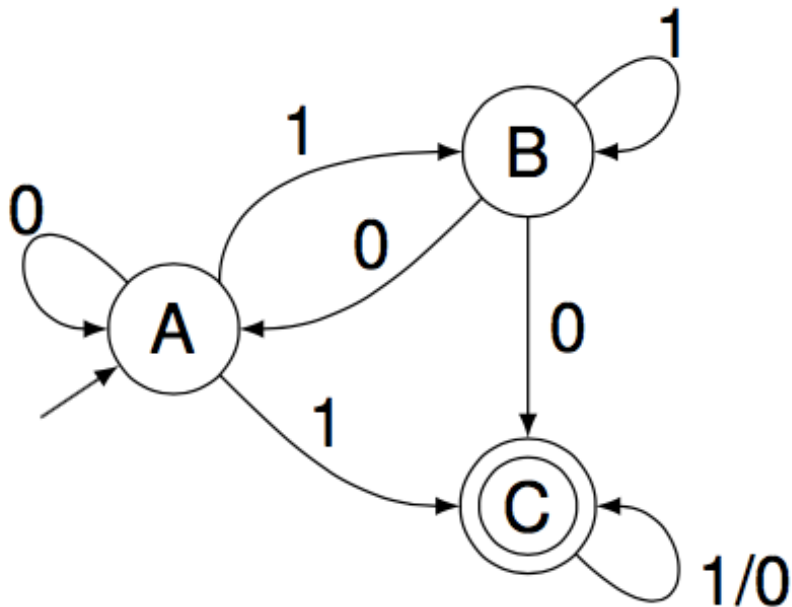


alphabet $\rightarrow$		
state $\downarrow$	0	1
A	A	

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon$ -closure(A) = A
- Keep getting  $\varepsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

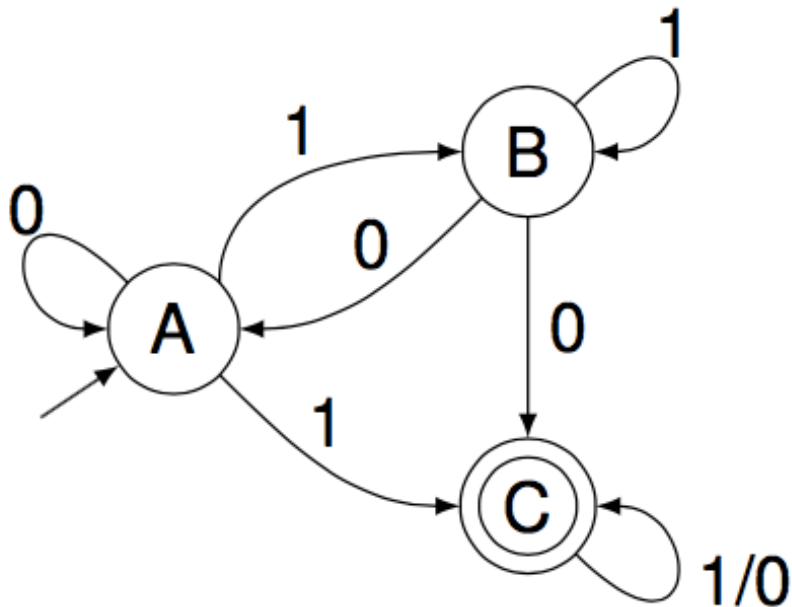


alphabet $\rightarrow$		
state $\downarrow$	0	1
	A	BC

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon$ -closure(A) = A
- Keep getting  $\varepsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

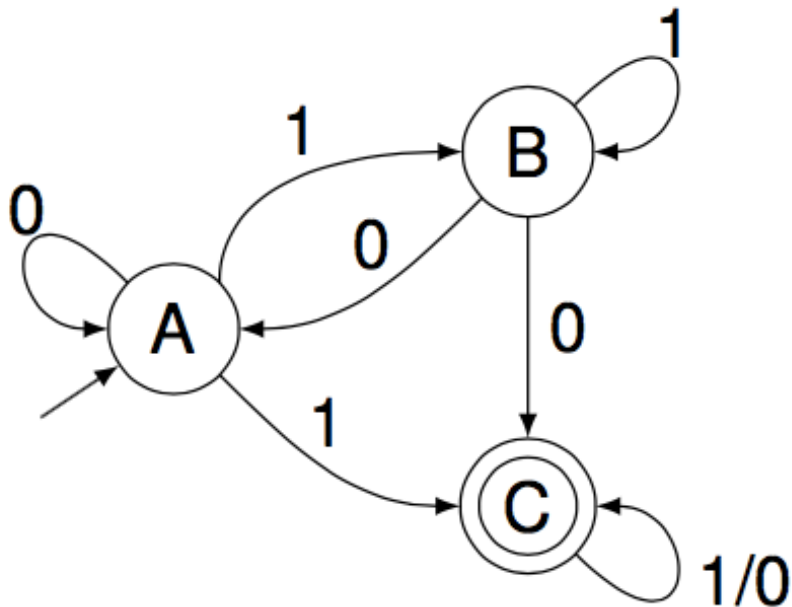


alphabet $\rightarrow$		
state $\downarrow$	0	1
A	A	BC

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\epsilon$ -closure of the start state
  - $\epsilon$ -closure(A) = A
- Keep getting  $\epsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

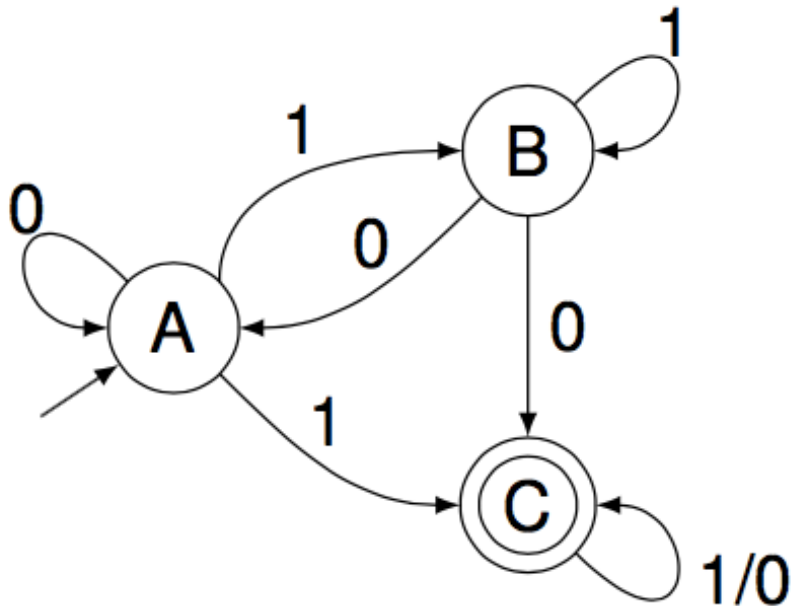


alphabet			
		0	1
state	A	A	BC
	BC		

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\epsilon$ -closure of the start state
  - $\epsilon$ -closure(A) = A
- Keep getting  $\epsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

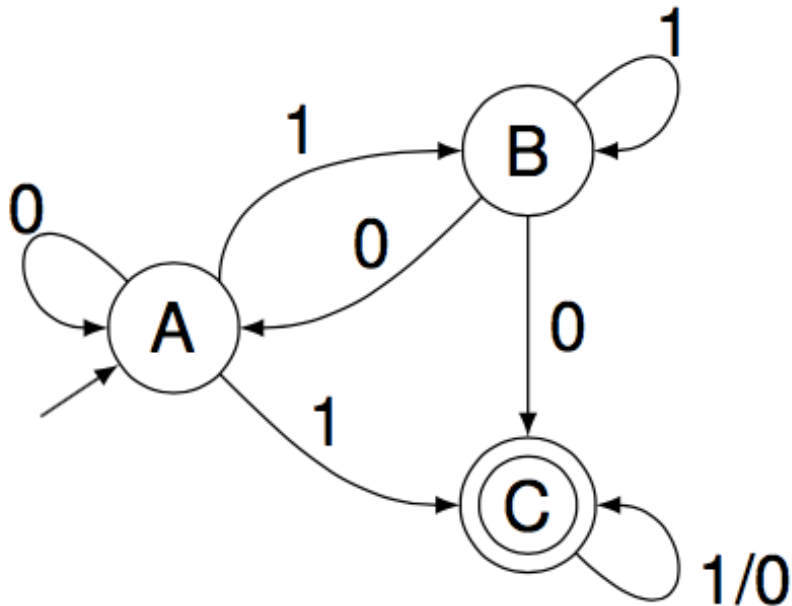


alphabet			
		0	1
state	A	A	BC
	BC	AC	

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\epsilon$ -closure of the start state
  - $\epsilon$ -closure(A) = A
- Keep getting  $\epsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

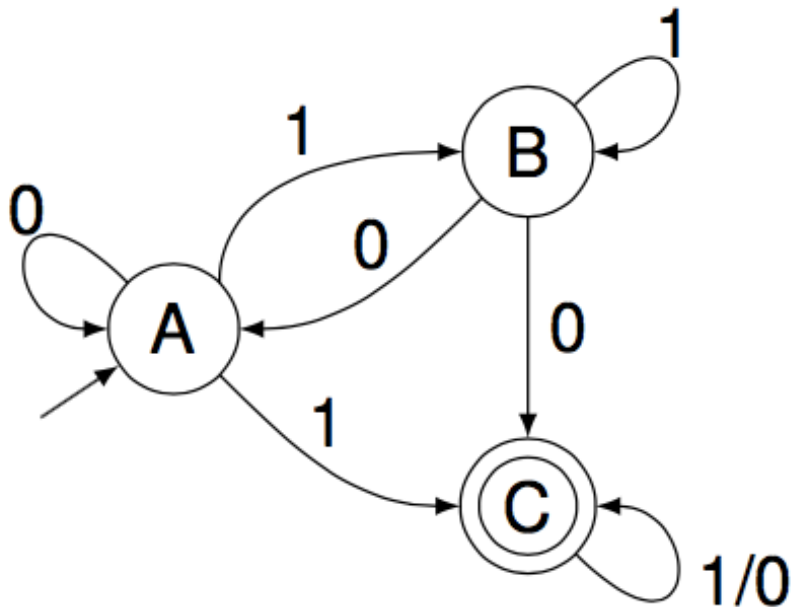


alphabet			
		0	1
state	A	A	BC
	BC	AC	BC

# NFA $\rightarrow$ DFA: Example

- Start by constructing  $\epsilon$ -closure of the start state
  - $\epsilon$ -closure(A) = A
- Keep getting  $\epsilon$ -closure( $move(T, a)$ )
- Stop, when there are no more new states

$T: A, a: 0/1$

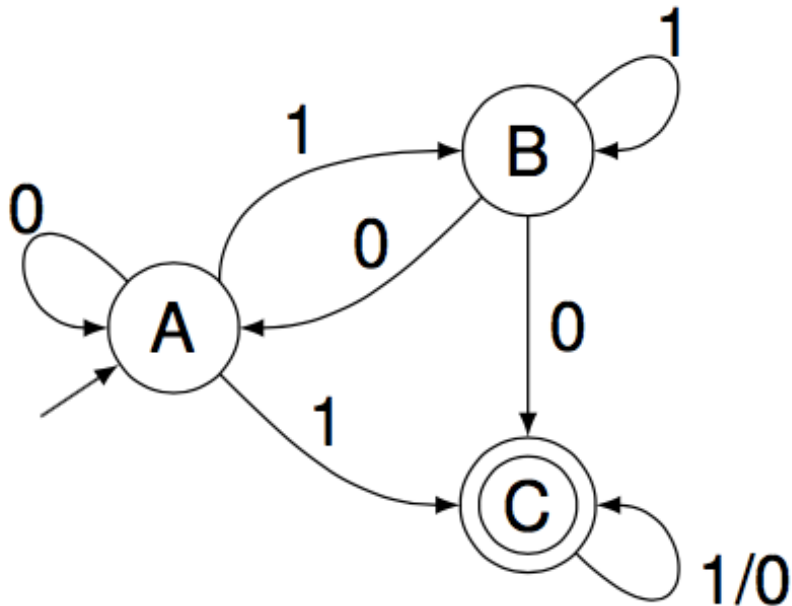


alphabet			
		0	1
state	A	A	BC
	BC	AC	BC



# NFA $\rightarrow$ DFA: Example

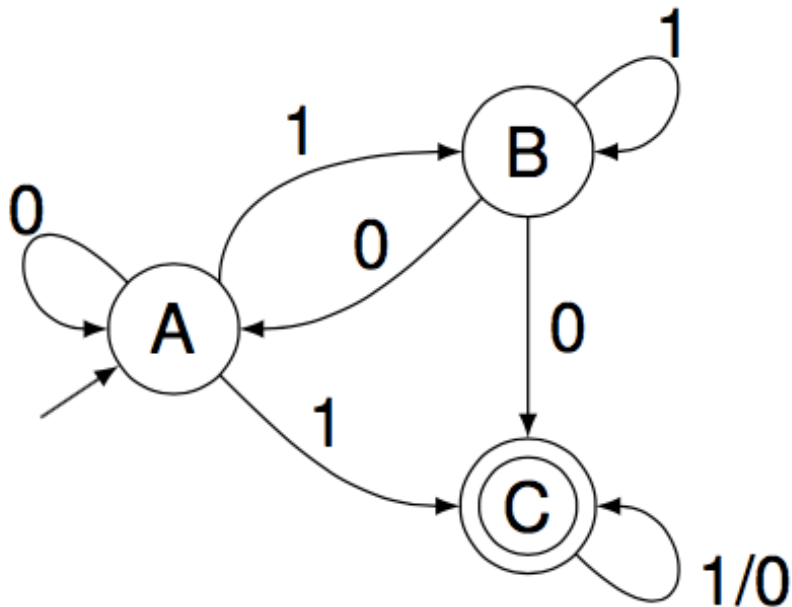
- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon\text{-closure}(A) = A$
- Keep getting  $\varepsilon\text{-closure}(\text{move}(T, a))$   $T: A, a: 0/1$
- Stop, when there are no more new states



alphabet $\rightarrow$		0	1
state $\downarrow$	A	A	BC
	BC	AC	BC
	AC		

# NFA $\rightarrow$ DFA: Example

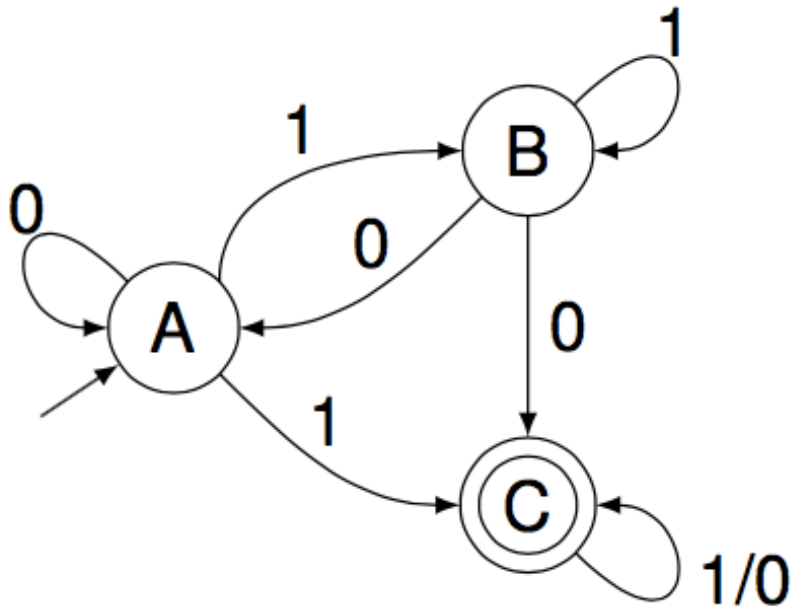
- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon\text{-closure}(A) = A$
- Keep getting  $\varepsilon\text{-closure}(\text{move}(T, a))$   $T: A, a: 0/1$
- Stop, when there are no more new states



alphabet			
		0	1
state	A	A	BC
	BC	AC	BC
	AC	AC	

# NFA $\rightarrow$ DFA: Example

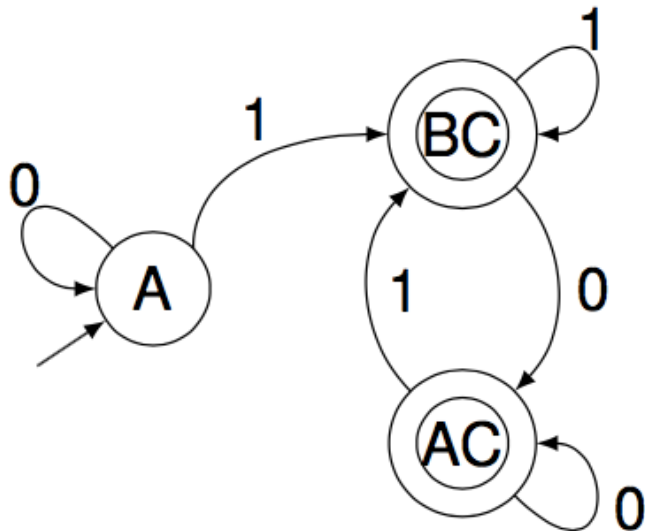
- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon\text{-closure}(A) = A$
- Keep getting  $\varepsilon\text{-closure}(\text{move}(T, a))$   $T: A, a: 0/1$
- Stop, when there are no more new states



alphabet $\rightarrow$		0	1
state $\downarrow$	A	A	BC
	BC	AC	BC
	AC	AC	BC

# NFA $\rightarrow$ DFA: Example (cont.)

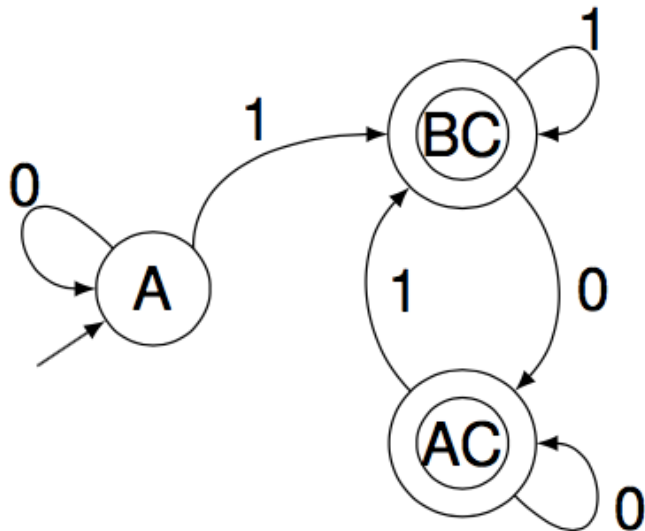
- Mark the final states of the DFA
  - The accepting states of  $D$  are all those sets of  $N$ 's states that include at least one accepting state of  $N$



alphabet $\rightarrow$		0	1
state $\downarrow$	A	A	BC
	BC	AC	BC
	AC	AC	BC

# NFA $\rightarrow$ DFA: Example (cont.)

- Mark the final states of the DFA
  - The accepting states of  $D$  are all those sets of  $N$ 's states that include at least one accepting state of  $N$

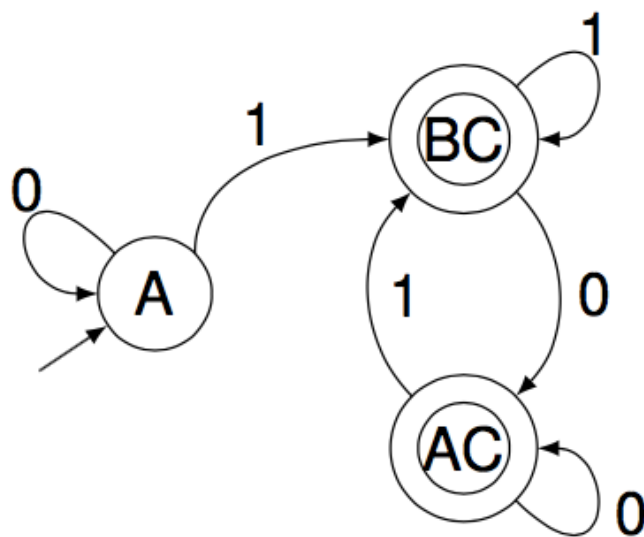


- Is the DFA minimal?
  - As few states as possible

alphabet $\rightarrow$		0	1
state $\downarrow$	A	A	BC
	BC	AC	BC
	AC	AC	BC

# NFA $\rightarrow$ DFA: Minimization[最小化]

- Any DFA can be converted to its minimum-state equivalent DFA
  - Discover sets of equivalent states
  - Represent each such set with just one state
- Two states are equivalent if and only if:
  - $\forall \alpha \in \Sigma$ , transitions on  $\alpha$  lead to equivalent states
  - $\alpha$ -transitions to distinct sets  $\Rightarrow$  states must be in distinct sets



Initial: {A}, {BC, AC}

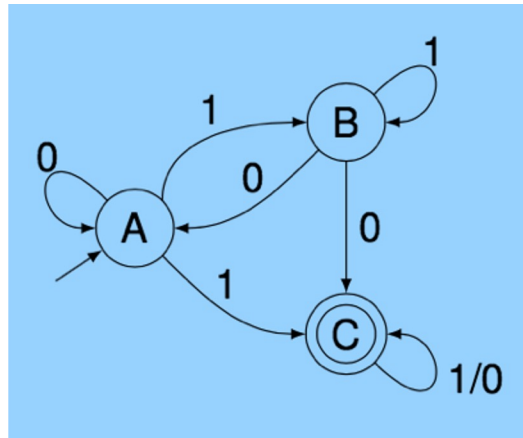
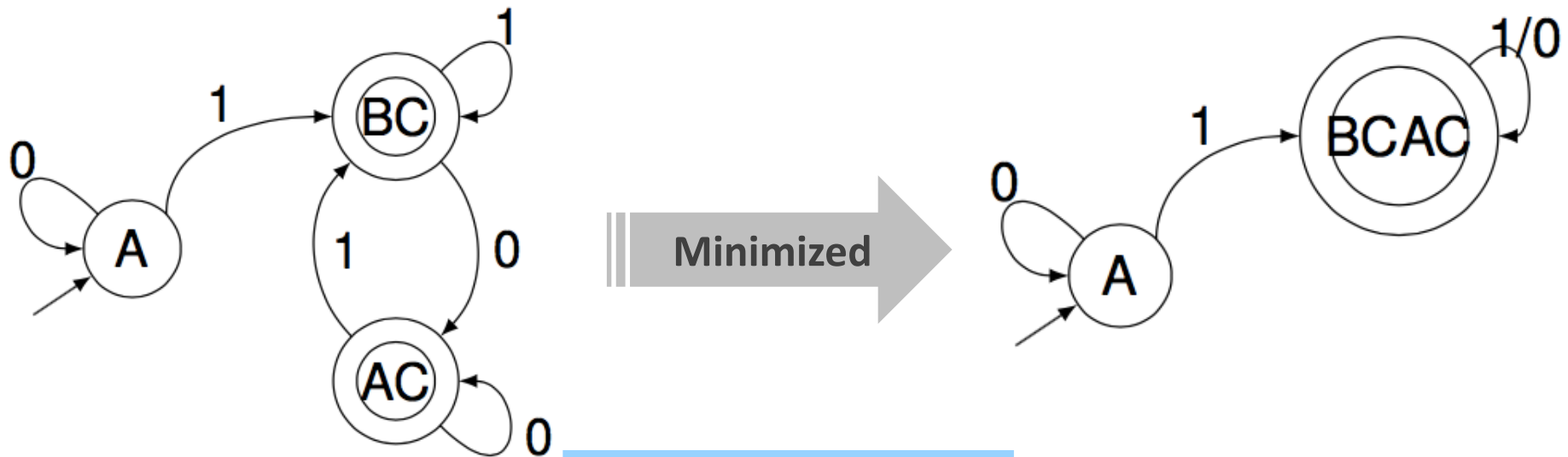
For {BC, AC} Initial sets:  
{non-accepting states}, {accepting states}

- BC on '0'  $\rightarrow$  AC, AC on '0'  $\rightarrow$  AC
- BC on '1'  $\rightarrow$  BC, AC on '1'  $\rightarrow$  BC
- No way to distinguish BC from AC on any string starting with '0' or '1'

Final: {A}, {BCAC}

# NFA $\rightarrow$ DFA: Minimization (cont.)

- States *BC* and *AC* do not need differentiation
  - Should be merged into one



# Minimization Algorithm

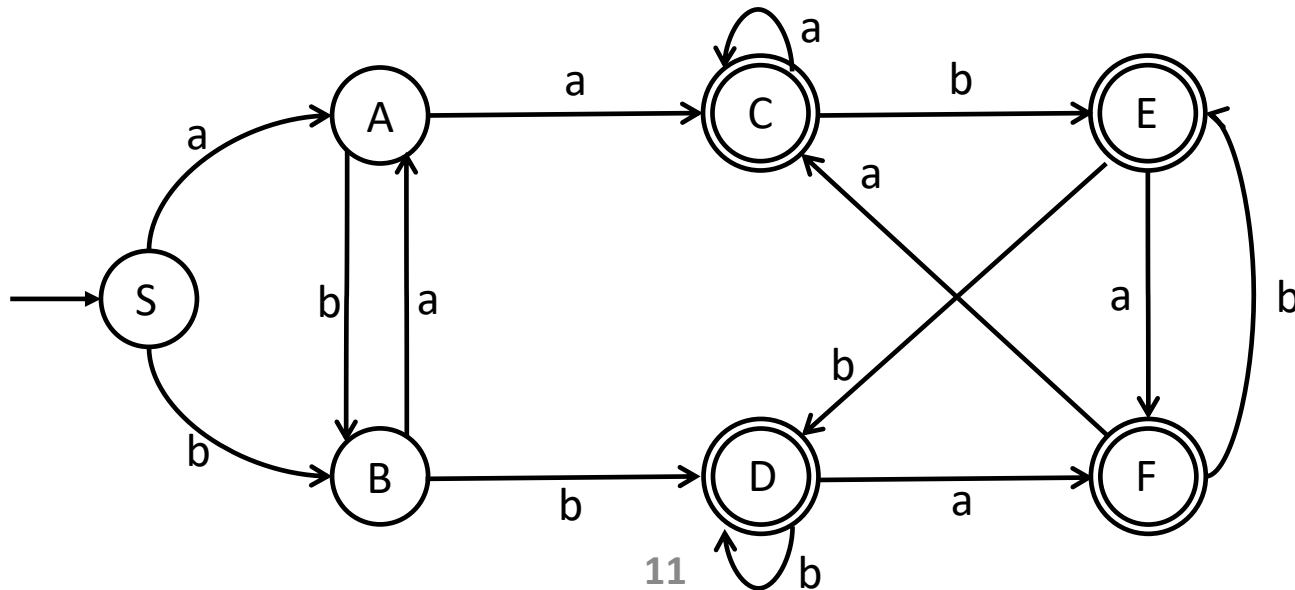
- The algorithm
  - Partitioning the states of a DFA into groups of states that **cannot be distinguished (i.e., equivalent)**
  - Each groups of states is then merged into a single state of the min-state DFA
- For a DFA  $(\Sigma, S, n, F, \delta)$ 
  - The initial partition  $P_0$ , has two sets  $\{F\}$  and  $\{S - F\}$
  - Splitting a set (i.e., partitioning a set  $s$  by input symbol  $\alpha$ )
    - Assume  $q_a$  and  $q_b \in s$ , and  $\delta(q_a, \alpha) = q_x$  and  $\delta(q_b, \alpha) = q_y$
    - If  $q_x$  and  $q_y$  are not in the same set, then  $s$  must be split (i.e.,  $\alpha$  splits  $s$ )
    - One state in the final DFA cannot have two transitions on  $\alpha$

```
P <- {F}, {S - F}
while (P is still changing)
  T <- {}
  for each state s ∈ P
    for each α ∈ Σ
      partition s by α into s1 and s2
      T <- T ∪ s1 ∪ s2
  if T ≠ P then
    P <- T
```

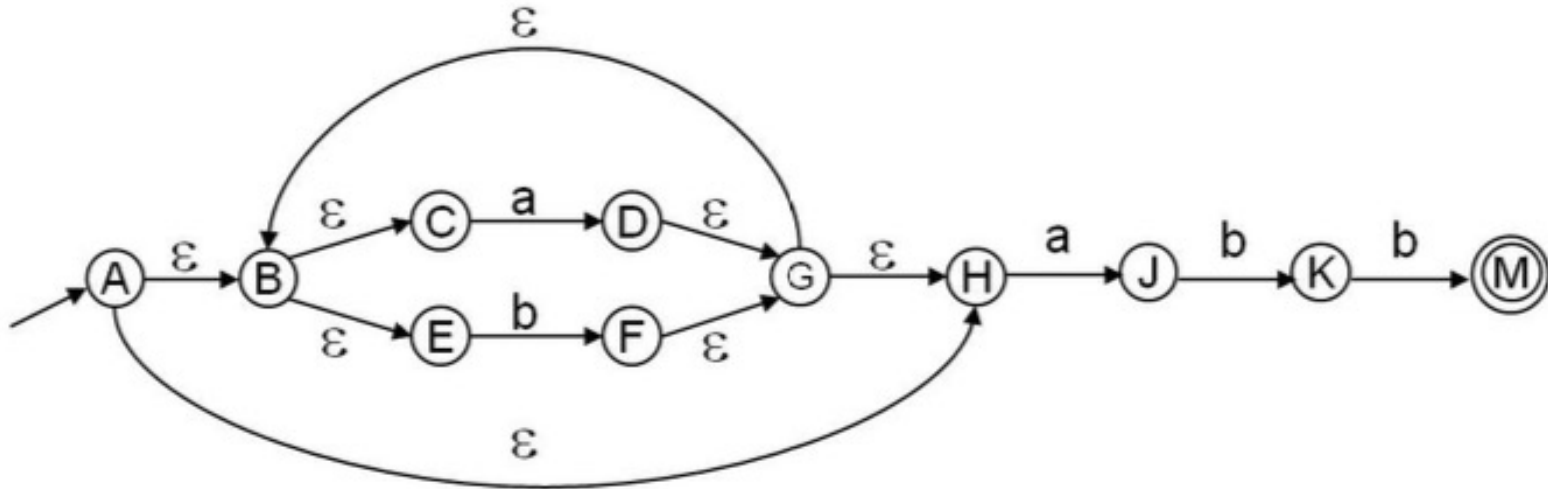


# Example

- P0:  $s_1 = \{S, A, B\}$ ,  $s_2 = \{C, D, E, F\}$
- For  $s_1$ , further splits into  $\{S\}$ ,  $\{A\}$ ,  $\{B\}$ 
  - a:  $S \rightarrow A \in s_1$ ,  $A \rightarrow C \in s_2$ ,  $B \rightarrow A \in s_1 \Rightarrow a$  distincts  $s_1 \Rightarrow \{S, B\}, \{A\}$
  - b:  $S \rightarrow B \in s_1$ ,  $A \rightarrow B \in s_1$ ,  $B \rightarrow D \in s_2 \Rightarrow b$  distincts  $s_1 \Rightarrow \{S\}, \{B\}, \{A\}$
- For  $s_2$ , all states are equivalent
  - a:  $C \rightarrow C \in s_2$ ,  $D \rightarrow F \in s_2$ ,  $E \rightarrow F \in s_2$ ,  $F \rightarrow C \in s_2 \Rightarrow a$  doesn't
  - b:  $C \rightarrow E \in s_2$ ,  $D \rightarrow D \in s_2$ ,  $E \rightarrow D \in s_2$ ,  $F \rightarrow E \in s_2 \Rightarrow b$  doesn't

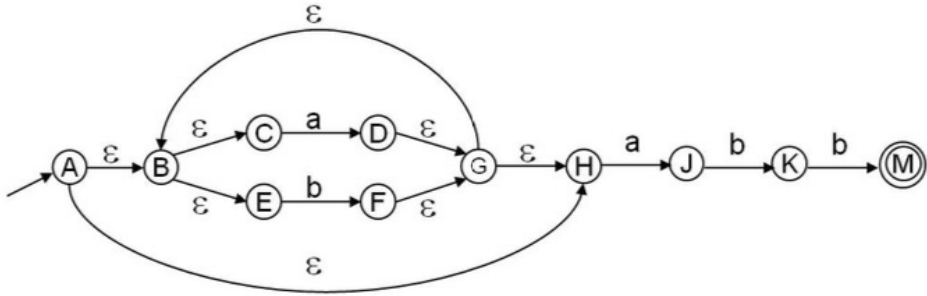


# NFA $\rightarrow$ DFA: More Example



- Start state of the equivalent DFA
  - $\epsilon$ -closure(A) = {A, B, C, E, H} = A'
- $\epsilon$ -closure(move(A', a)) =  $\epsilon$ -closure({D, J}) = {B, C, D, E, H, G, J} = B'
- $\epsilon$ -closure(move(A', b)) =  $\epsilon$ -closure({F}) = {B, C, E, F, G, H} = C'
- ... ..

# Step 1: Construct the NFA Table

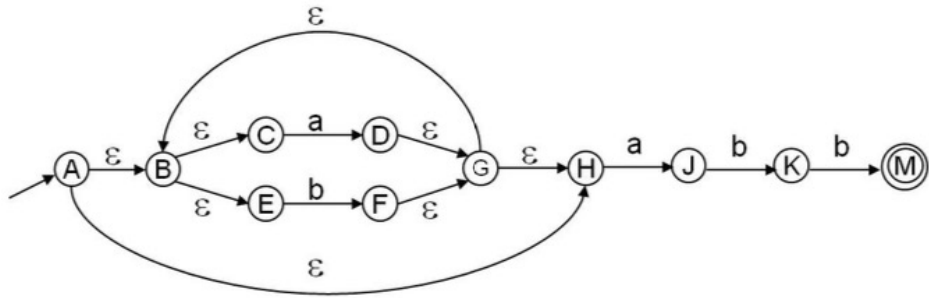


State table of the NFA  $\Rightarrow$

NFA最对应的转换表

	$\epsilon$	a	b
A	BH		
B	CE		
C		D	
D	G		
E			F
F	G		
G	BH		
H		J	
I			
J			K
K			M
M			

# Step 2: Update $\varepsilon$ Column to $\varepsilon$ -closure



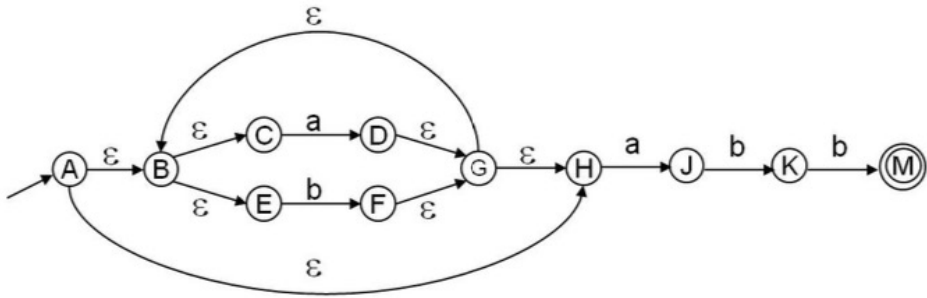
$\varepsilon$ -closure of the NFA state  
e.g.,  $\varepsilon\text{-closure}(D) = \{D, B, H, C, E\}$

求取 $\varepsilon$ 闭包



	$\varepsilon$	a	b
A	ABHCE		
B	BCE		
C		D	
D	DBHCE		
E			F
F	FGBHCE		
G	GBHCE		
H		J	
I			
J			K
K			M
M			

# Step 3: Update other cols based on the $\epsilon$ col



Get the transitions of the  $\epsilon$  col  
e.g.,  $\{D, B, H, C, E\} \xrightarrow{a} \{D, J\}$



	$\epsilon$	a	b
A	ABHCE	DJ	F
B	BCE	D	F
C		D	
D	DBHCE	DJ	F
E			F
F	FGBHCE	DJ	F
G	GBHCE	DJ	F
H		J	
I			
J			K
K			M
M			

# Step 4: Construct the DFA Table

	$\epsilon$	a	b
A	ABHCE <i>A' for short</i>	DJ	F
B	BCE	D	F
C		D	
D	DBHCE	DJ	F
E			F
F	FGBHCE	DJ	F
G	GBHCE	DJ	F
H		J	
I			
J			K
K			M
M			

# Step 4: Construct the DFA Table

	a	b
A'	DJ	F

	$\epsilon$	a	b
A	ABHCE <i>A' for short</i>	DJ	F
B	BCE	D	F
C		D	
D	DBHCE	DJ	F
E			F
F	FGBHCE	DJ	F
G	GBHCE	DJ	F
H		J	
I			
J			K
K			M
M			

# Step 4: Construct the DFA Table

	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F

	$\epsilon$	a	b
A	ABHCE <i>A' for short</i>	DJ	F
B	BCE	D	F
C		D	
D	DBHCE	DJ	F
E			F
F	FGBHCE	DJ	F
G	GBHCE	DJ	F
H		J	
I			
J			K
K			M
M			



# Step 4: Construct the DFA Table

	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM

	$\epsilon$	a	b
A	ABHCE <i>A' for short</i>	DJ	F
B	BCE	D	F
C		D	
D	DBHCE	DJ	F
E			F
F	FGBHCE	DJ	F
G	GBHCE	DJ	F
H		J	
I			
J			K
K			M
M			

# Step 4: Construct the DFA Table

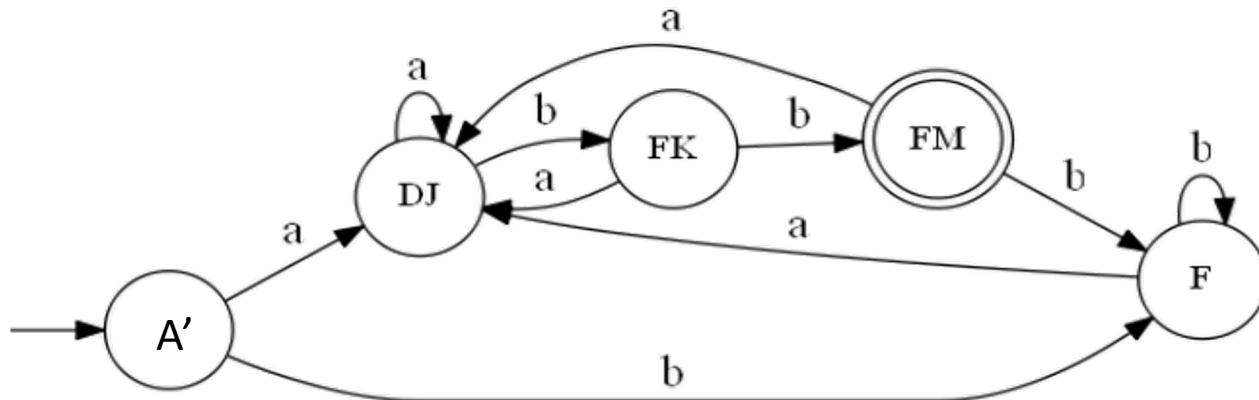
	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

	$\epsilon$	a	b
A	ABHCE <i>A' for short</i>	DJ	F
B	BCE	D	F
C		D	
D	DBHCE	DJ	F
E			F
F	FGBHCE	DJ	F
G	GBHCE	DJ	F
H		J	
I			
J			K
K			M
M			

# Step 4: Construct the DFA Table(cont.)

	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

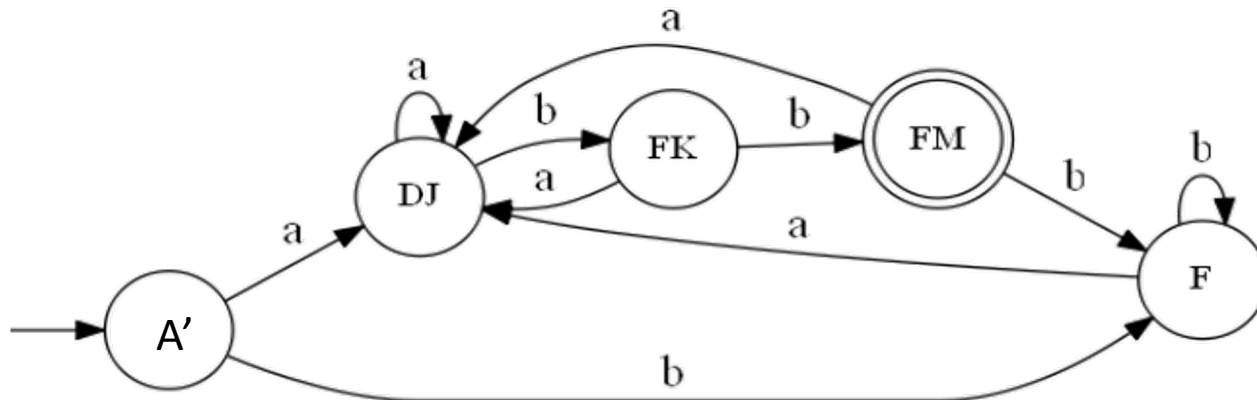
- Is the DFA minimal?



# Step 4: Construct the DFA Table(cont.)

	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

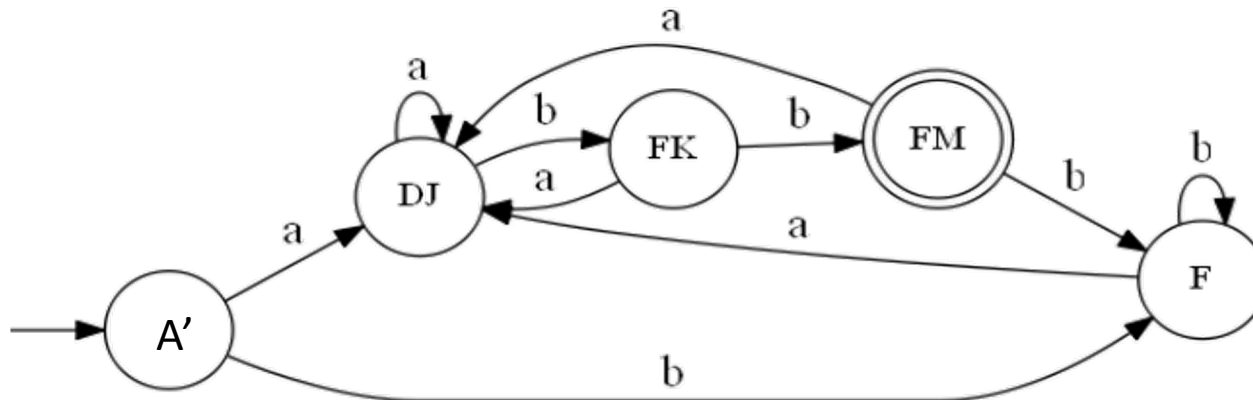
- Is the DFA minimal?
  - States A' and F should be merged



# Step 4: Construct the DFA Table(cont.)

	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

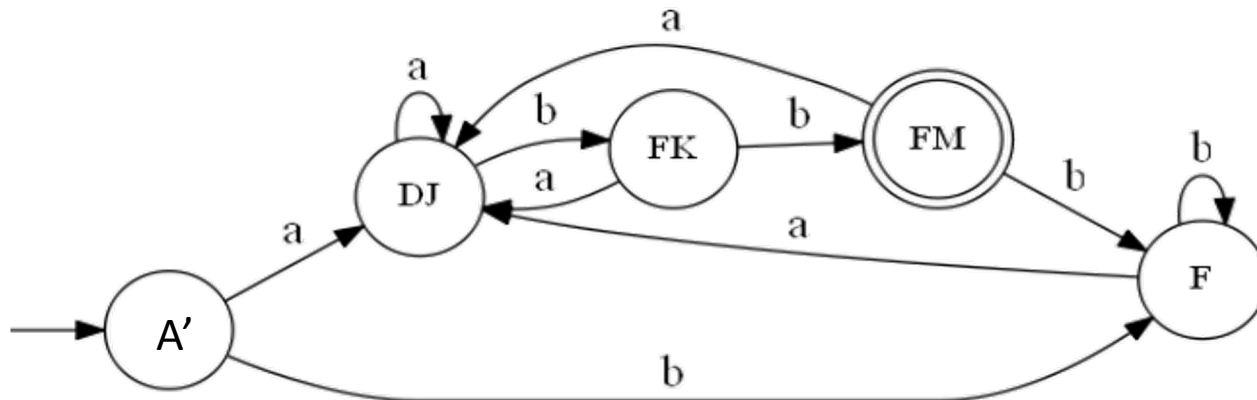
- Is the DFA minimal?
  - States A' and F should be merged
- Should we merge states A' and FM?



# Step 4: Construct the DFA Table(cont.)

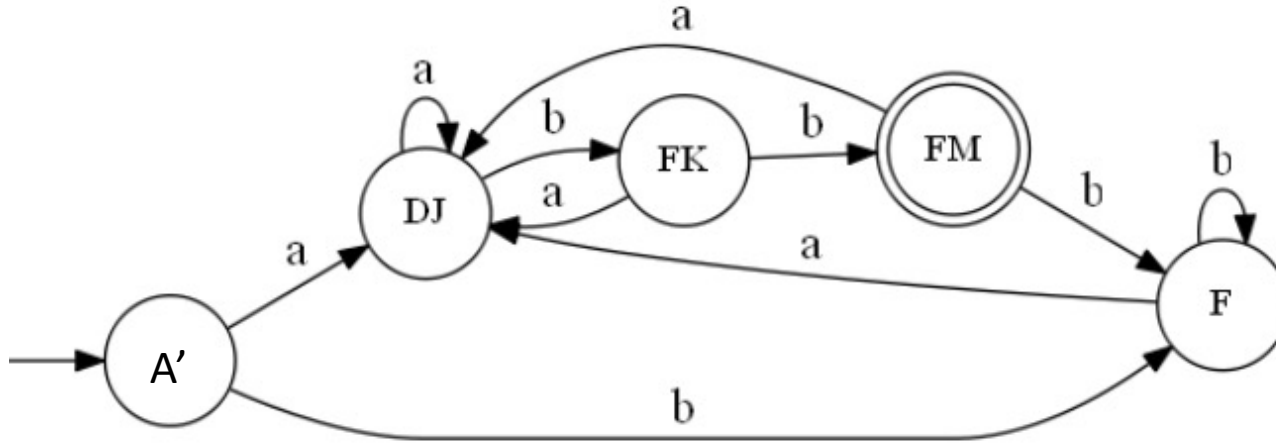
	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

- Is the DFA minimal?
  - States A' and F should be merged
- Should we merge states A' and FM?
  - NO. A' and FM are in different sets from the very beginning (FM is accepting, A' is not).

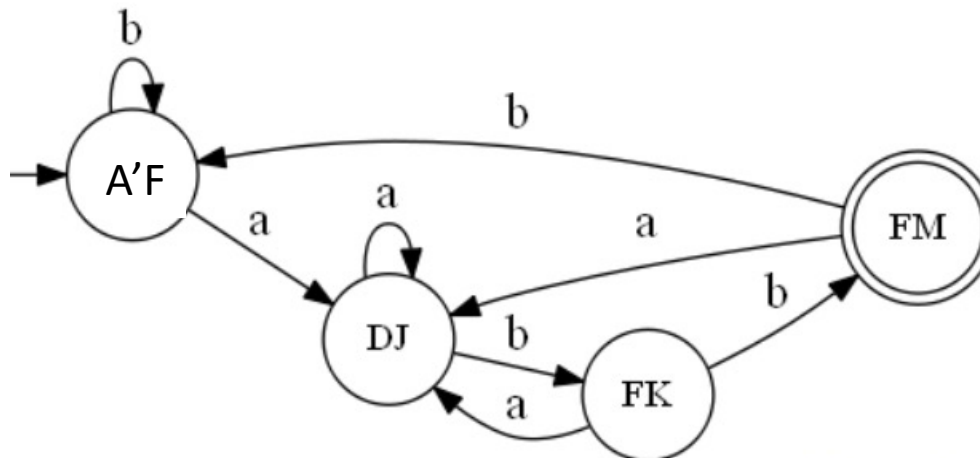


# Step 5: (Optional) Minimize DFA

- Original DFA: before merging A' and F



- Minimized DFA: Do you see the original RE  $(a|b)^*abb$



# NFA $\rightarrow$ DFA: Space Complexity[空间复杂度]

---

- NFA may be in many states at any time
- How many different possible states in DFA?
  - If there are  $N$  states in NFA, the DFA must be in some subset of those  $N$  states
  - How many non-empty subsets are there?
    - $2^N - 1$
- The resulting DFA has  $O(2^N)$  space complexity, where  $N$  is number of original states in NFA
  - For real languages, the NFA and DFA have about same number of states



# NFA $\rightarrow$ DFA: Time Complexity[时间复杂度]

- DFA execution

- Requires  $O(|X|)$  steps, where  $|X|$  is the input length
- Each step takes constant time
  - If current state is  $S$  and input is  $c$ , then read  $T[S, c]$
  - Update current state to state  $T[S, c]$
- Time complexity =  $O(|X|)$

**Deterministic:**  
unique transition

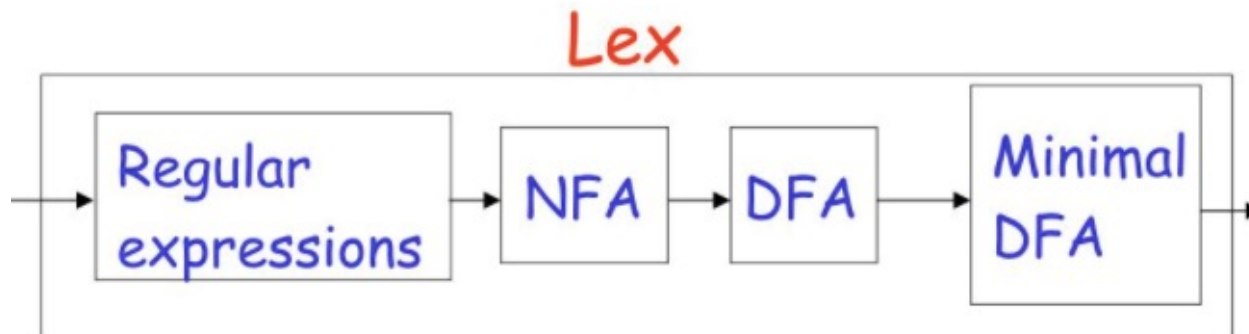
- NFA execution

- Requires  $O(|X|)$  steps, where  $|X|$  is the input length
  - Anyway, the input symbols should be completely processed
- Each step takes  $O(N^2)$  time, where  $N$  is the number of states
  - Current state is a set of potential states, up to  $N$
  - On input  $c$ , must union all  $T[S_{\text{potential}}, c]$ , up to  $N$  times
    - Each union operation takes  $O(N)$  time
- Time complexity =  $O(|X| * N^2)$

**Non-deterministic:**  
from current state,  
you can transit to any  
(including itself)

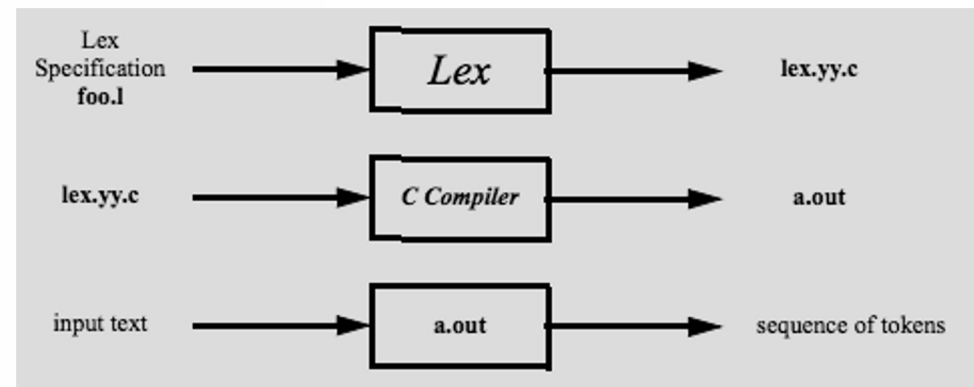
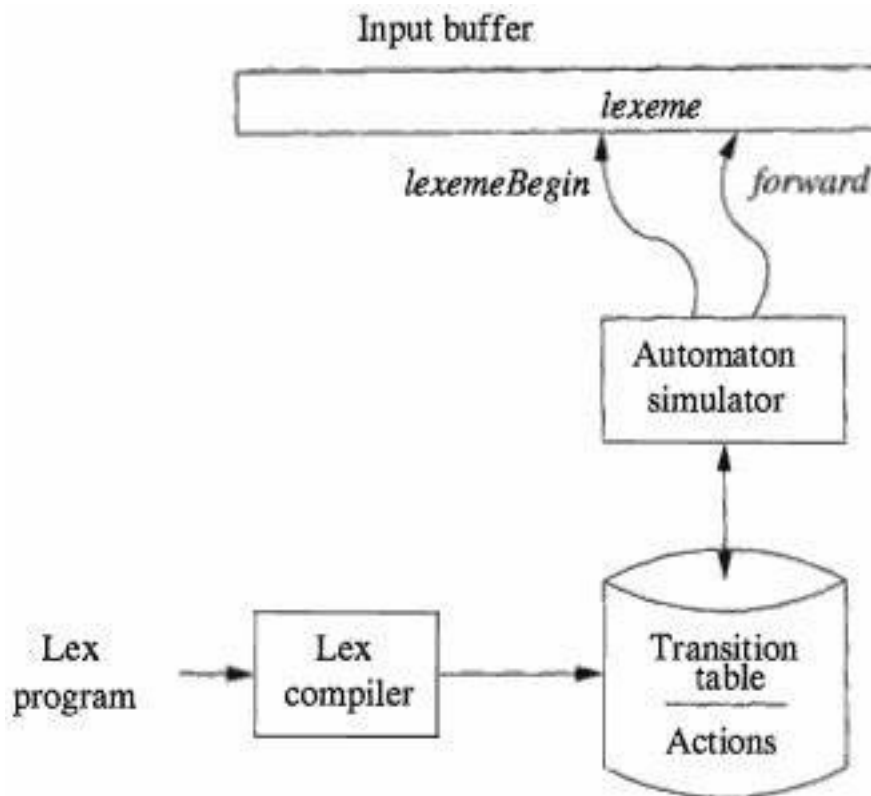
# Implementation in Practice[实际实现]

- Lex: RE  $\rightarrow$  NFA  $\rightarrow$  DFA  $\rightarrow$  Table
  - Converts regular expressions to NFA
  - Converts NFA to DFA
  - Performs DFA state minimization to reduce space
  - Generate the transition table from DFA
  - Performs table compression to further reduce space
- Most other automated lexers also choose DFA over NFA
  - Trade off space for speed



# Lexical Analyzer Generated by Lex

- A Lex program is turned into a transition table and actions, which are used by a FA simulator
- Automaton recognizes matching any of the patterns



# Lex: Example

- Three patterns, three NFAs
- Combine three NFAs into a single NFA
  - Add start state 0 and  $\epsilon$ -transitions

