



中山大學  
SUN YAT-SEN UNIVERSITY



国家超级计算广州中心  
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

# 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, Spring 2021



中山大學  
SUN YAT-SEN UNIVERSITY



# Review Questions

Q1: Can we have multiple start/accepting states in FA?

start: only one, accepting: multiple

Q2: what are NFA and DFA? How to differentiate?

NFA: non-deterministic FA, DFA: deterministic FA  
 $\epsilon$ -move or multiple transitions per input per state

Q3: how do RE, NFA, DFA relate to each other?

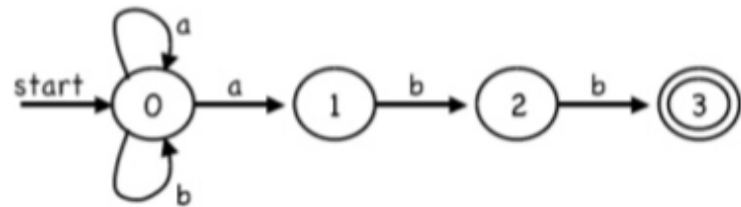
$L(RE) \equiv L(NFA) \equiv L(DFA)$

Q4: the state graph is a NFA or DFA?

NFA, multiple transitions for state '0' on input 'a'

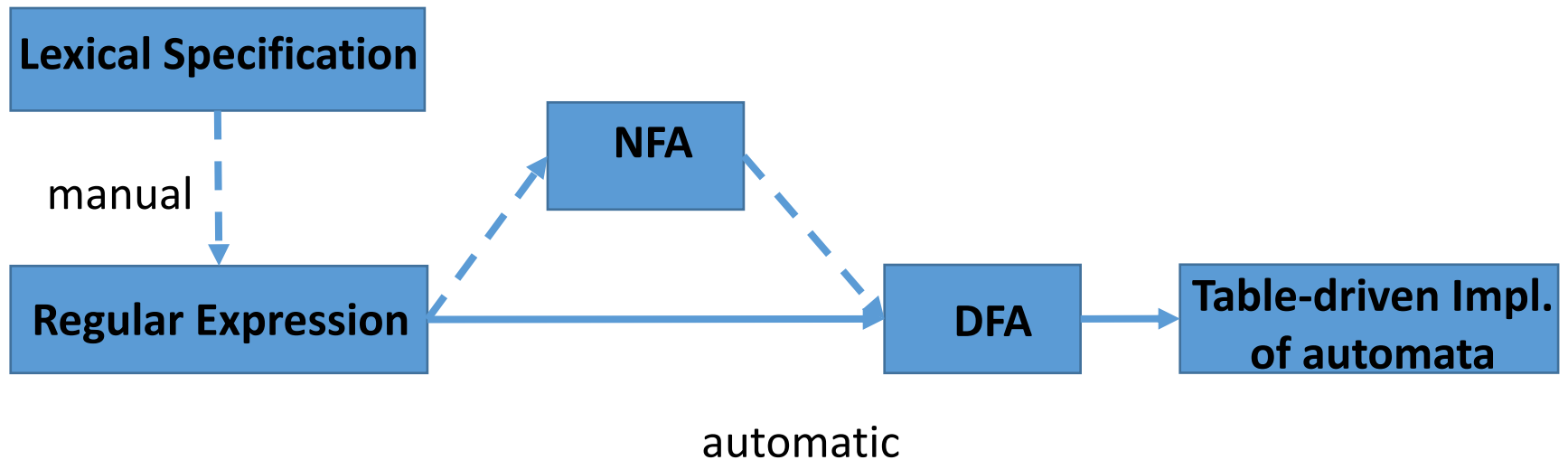
Q5: what's the language it recognizes ?

$(a|b)^*abb$



# Specification to Implementation

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



# NFA $\rightarrow$ DFA: Theory

---

- Question: is  $L(\text{NFA}) \subseteq L(\text{DFA})$ 
  - Otherwise, conversion would be futile
- Theorem:  $L(\text{NFA}) \equiv L(\text{DFA})$ 
  - Both recognize regular languages  $L(\text{RE})$
  - Will show  $L(\text{NFA}) \subseteq L(\text{DFA})$  by construction ( $\text{NFA} \rightarrow \text{DFA}$ )
  - Since  $L(\text{DFA}) \subseteq L(\text{NFA})$ ,  $L(\text{NFA}) \equiv L(\text{DFA})$
- Resulting DFA consumes more memory than NFA
  - Potentially larger transition table
- But DFAs are faster to execute
  - For DFAs, number of transitions == length of input
  - For NFAs, number of potential transitions can be larger
- $\text{NFA} \rightarrow \text{DFA}$  conversion is done because the speed of DFA far outweighs its extra memory consumption

# NFA $\rightarrow$ DFA: Idea

---

- **Subset construction**

- Each state of the constructed DFA corresponds to a set of NFA states
- After reading input  $a_1a_2\dots a_n$ , the DFA is in that state which corresponds to the set of states that the NFA can reach, from its start state, following paths labeled  $a_1a_2\dots a_n$

- **Algorithm to convert**

- Input: an NFA  $N$
- Output: a DFA  $D$  accepting the same language as  $N$

# NFA $\rightarrow$ DFA: Algorithm

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

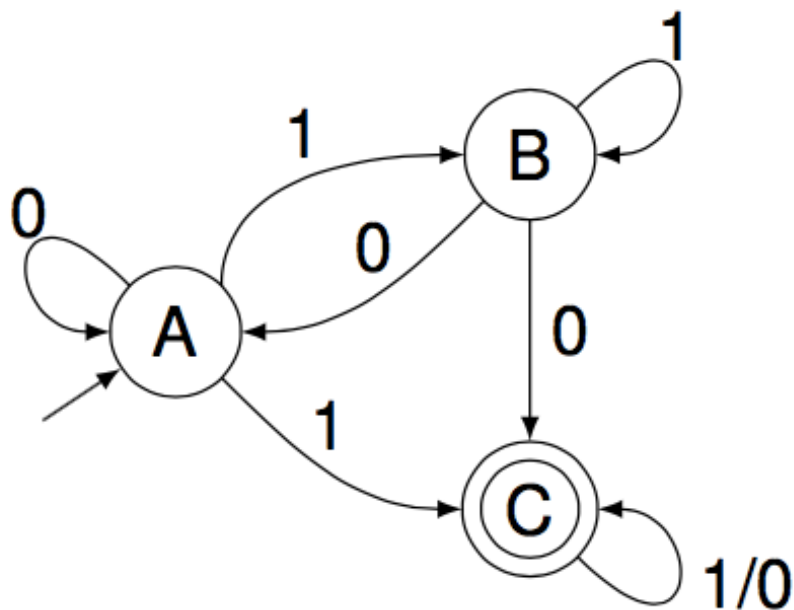
```
while there is an unmarked state  $T$  in  $Dstates$  do  
    mark  $T$   
    for each input symbol  $a \in \Sigma$  do  
         $U := \epsilon\text{-closure}(\text{move}(T, a))$   
        if  $U$  is not in  $Dstates$  then  
            add  $U$  as an unmarked state to  $Dstates$   
    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

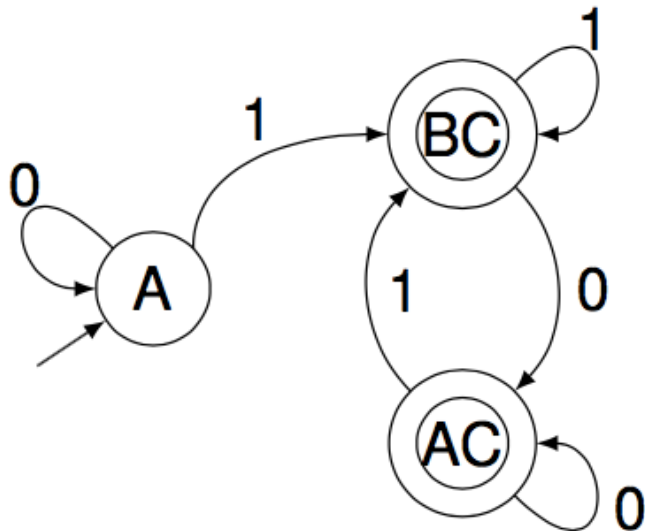
- Start by constructing  $\varepsilon$ -closure of the start state
  - $\varepsilon$ -closure(A) = A
- Keep getting  $\varepsilon$ -closure(move(T, a))
- Stop, when there are no more new states



alphabet			
		0	1
state	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  $N$



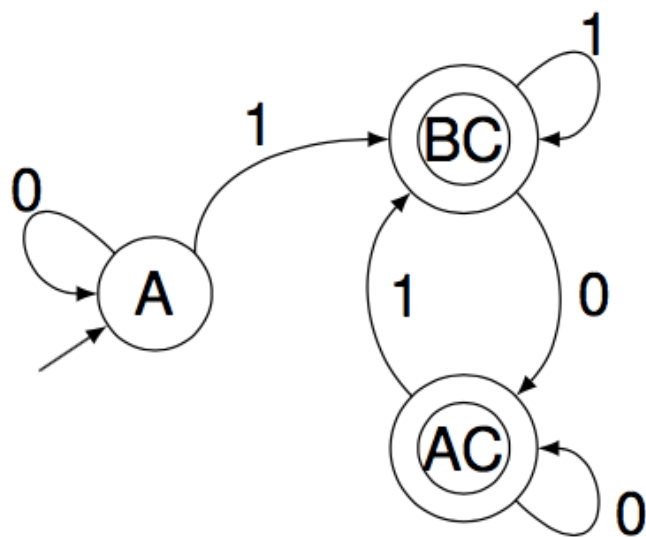
- Is the DFA minimal?
  - As few states as possible

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



# NFA $\rightarrow$ DFA: Minimization

- Any DFA can be converted to its minimum-state equivalent DFA
  - Partitioning the states of a DFA into groups of states that **cannot be distinguished**
  - Each groups of states is then merged into a single state of the min-state DFA



Initial: {A}, {BC, AC}

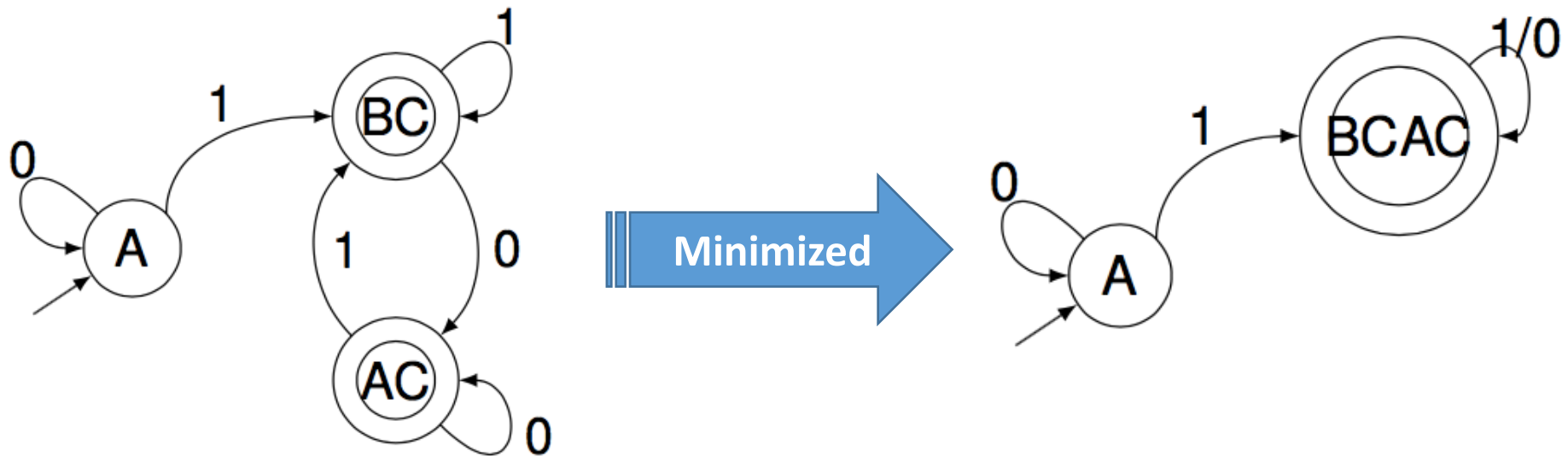
For {BC, AC}

- 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'

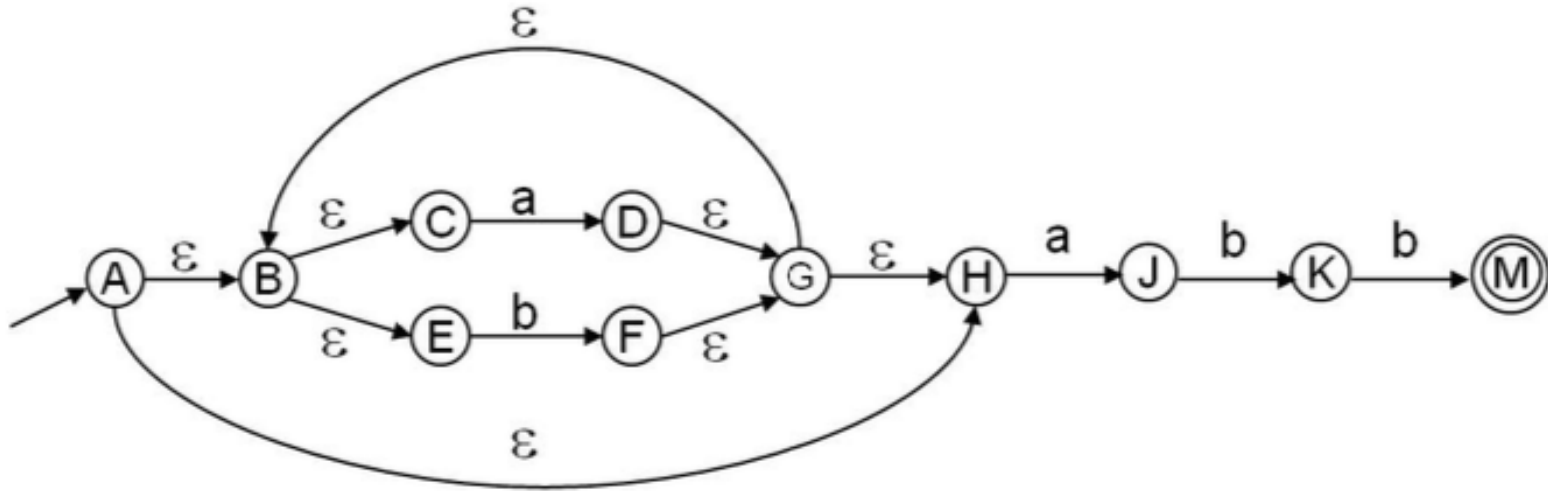
Final: {A}, {BCAC}

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

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

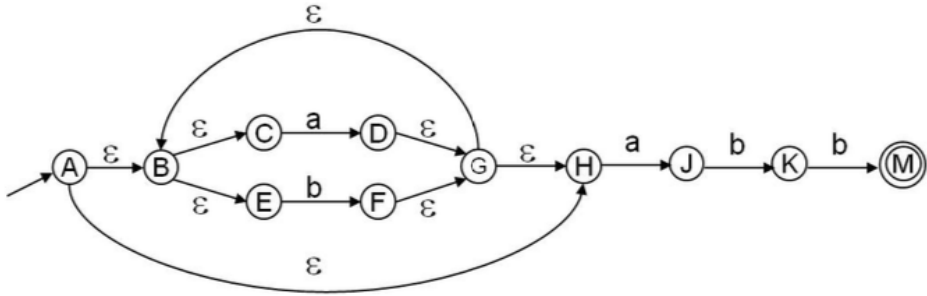


# NFA $\rightarrow$ DFA: More Example



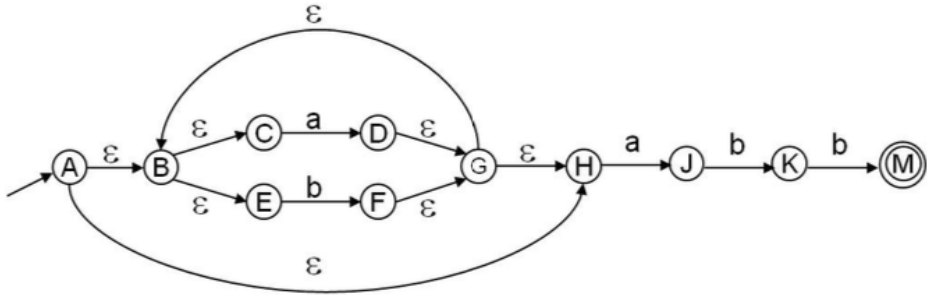
- Start state of the equivalent DFA
  - $\epsilon$ -closure(A) = {A, B, C, E, H} = A'
- $\epsilon$ -closure(move(A', a)) =  $\epsilon$ -closure({D, J}) = {B, C, D, E, H, G, J} = B'
- $\epsilon$ -closure(move(A', b)) =  $\epsilon$ -closure({F}) = {B, C, E, F, G, H} = C'
- ... ..

# Step 1: Construct the NFA Table



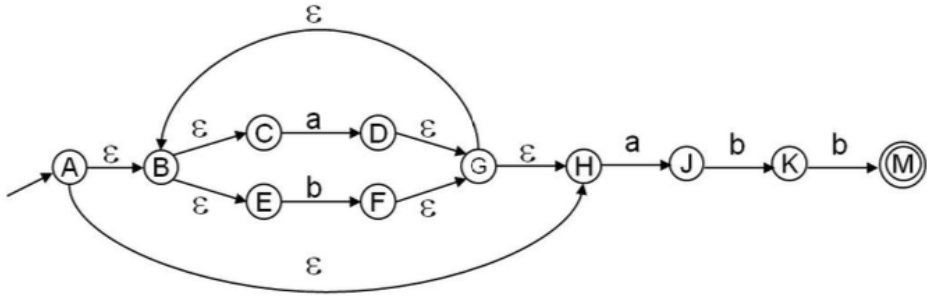
	$\epsilon$	a	b
A	BH		
B	CE		
C		D	
D	G		
E			F
F	G		
G	BH		
H		J	
I			
J			K
K			M
M			

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



	$\epsilon$	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 $\epsilon$ -closure



	$\epsilon$	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

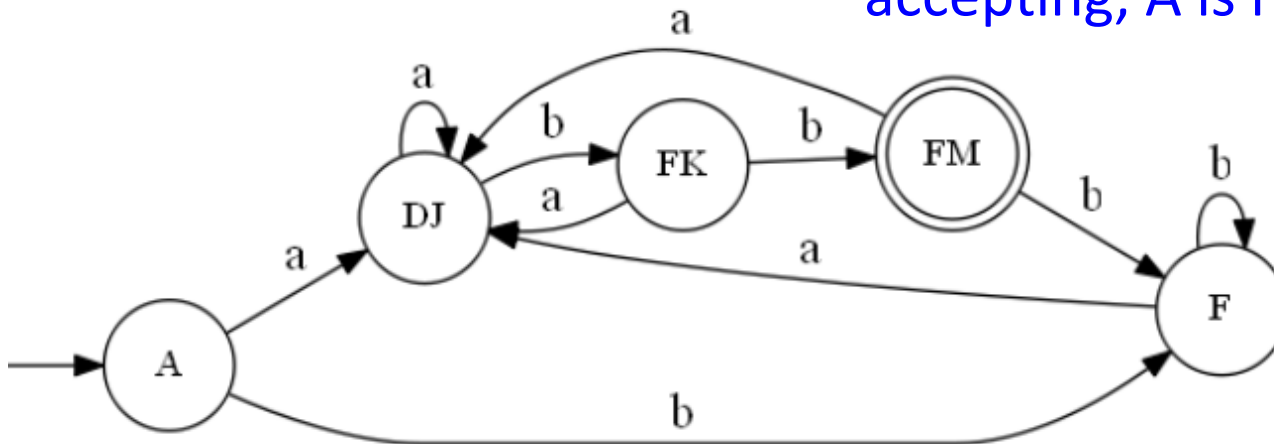
	a	b
A	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

	$\epsilon$	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(cont.)

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

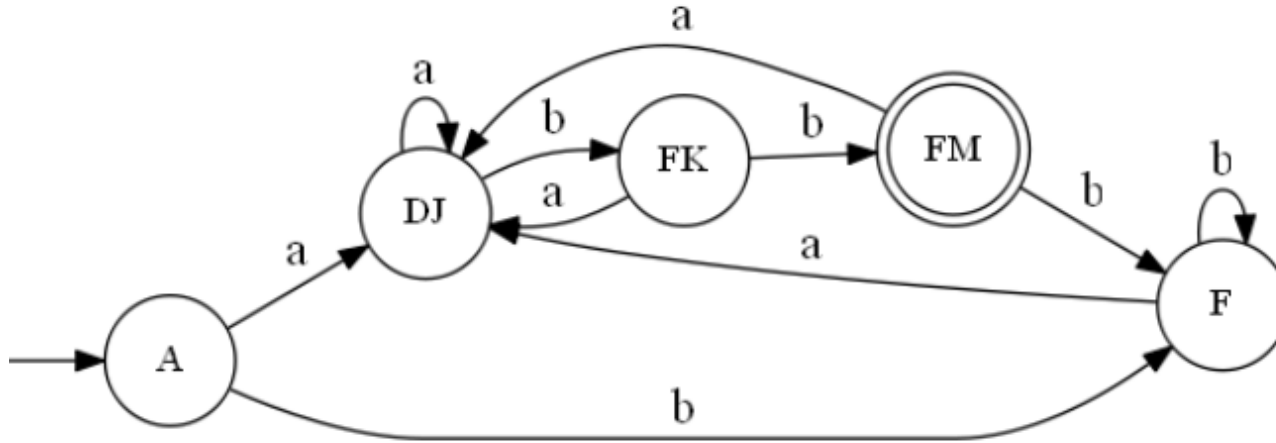
- Is the DFA minimal?
  - States A and F should be merged
- Should we merge states A and FM?
  - NO. A and FM are in different sets from the very beginning (FM is accepting, A is not).



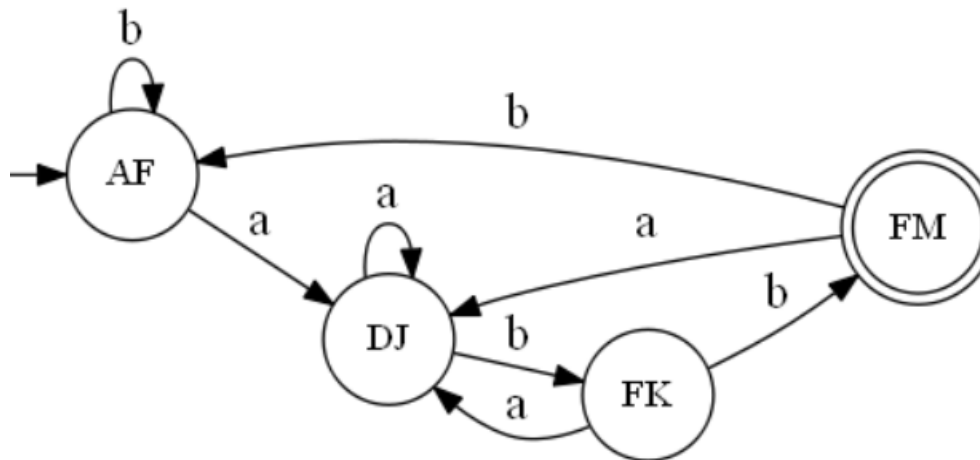


# Step 5: (Optional) Minimize DFA

- Original DFA: before merging A and F



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



# NFA $\rightarrow$ DFA: Space Complexity

---

- 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 #states

# NFA $\rightarrow$ DFA: Time Complexity

---

- DFA execution

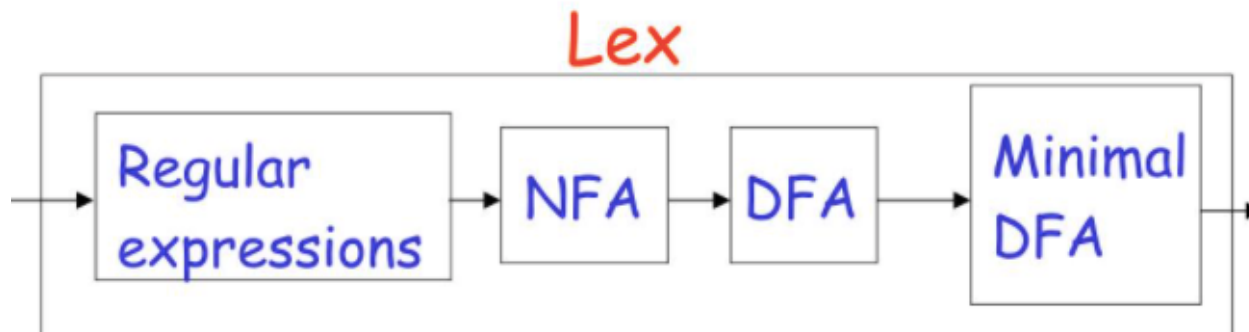
- Requires  $O(|X|)$  steps, where  $|X|$  is the input length
- Each step takes constant time
  - If current state is  $S$  and input is  $c$ , then read  $T[S, c]$
  - Update current state to state  $T[S, c]$
- Time complexity =  $O(X)$

- NFA execution

- Requires  $O(|X|)$  steps, where  $|X|$  is the input length
- 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 =  $O(|X| * N^2)$

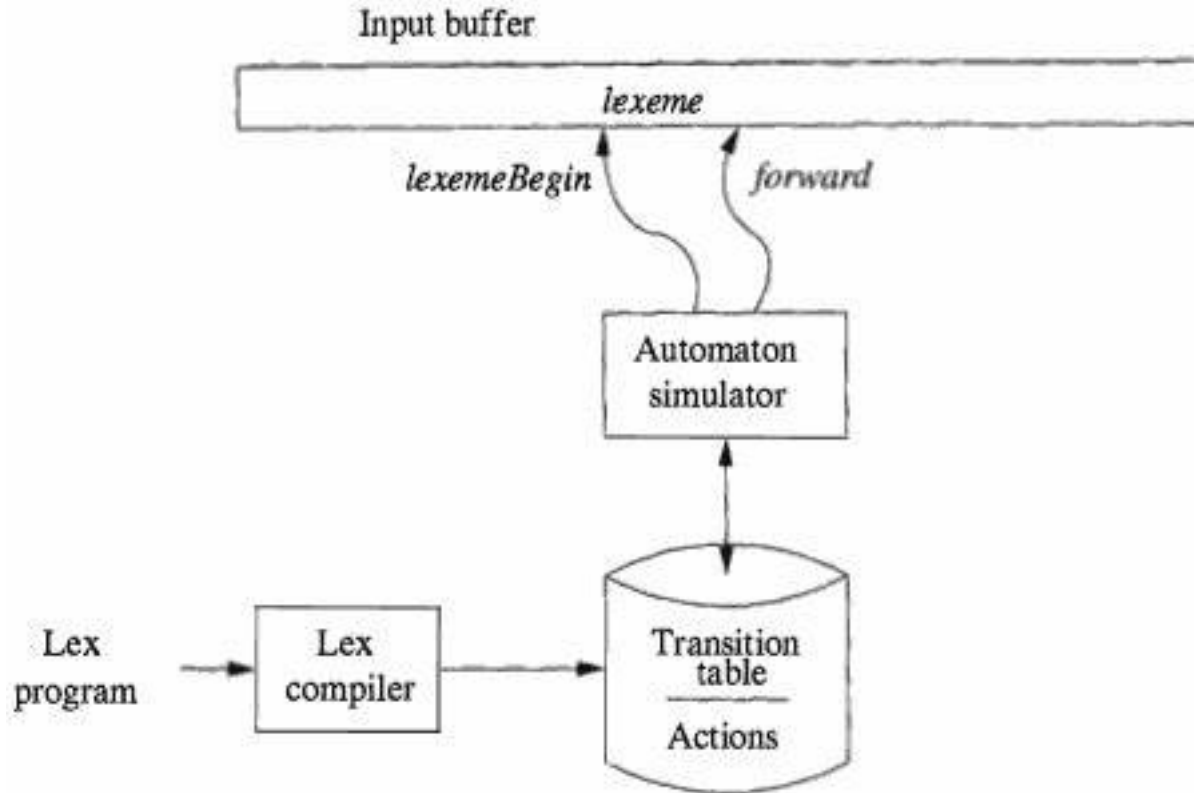
# Implementation in Practice

- 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



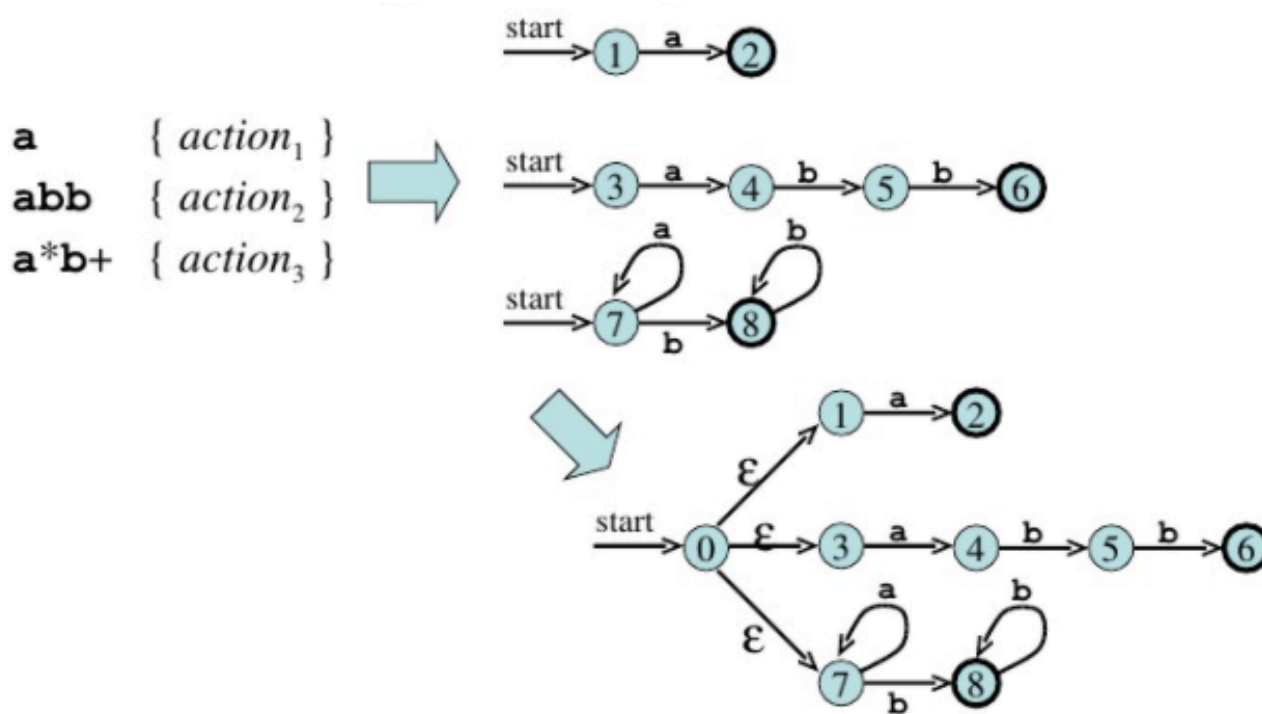
# Lexical Analyzer Generated by Lex

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



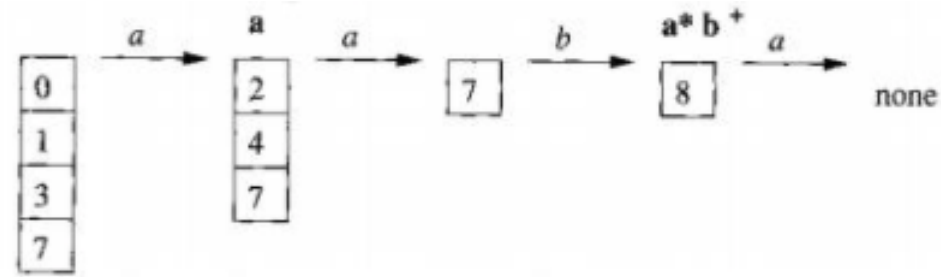
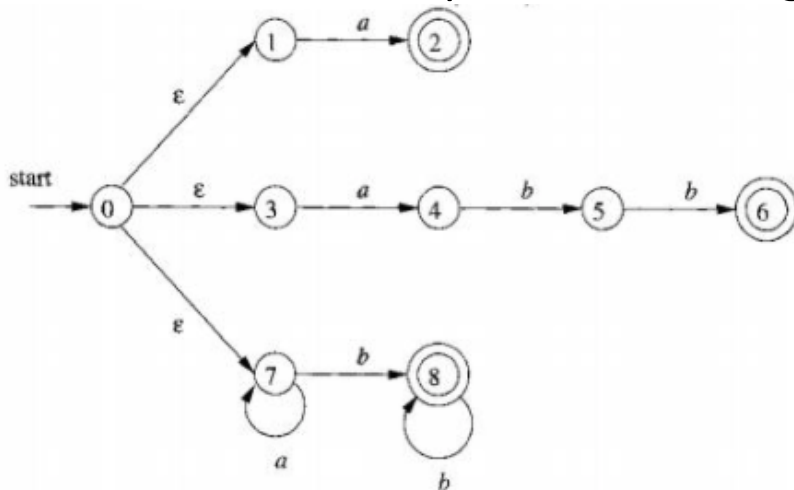
# Lex: Example

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



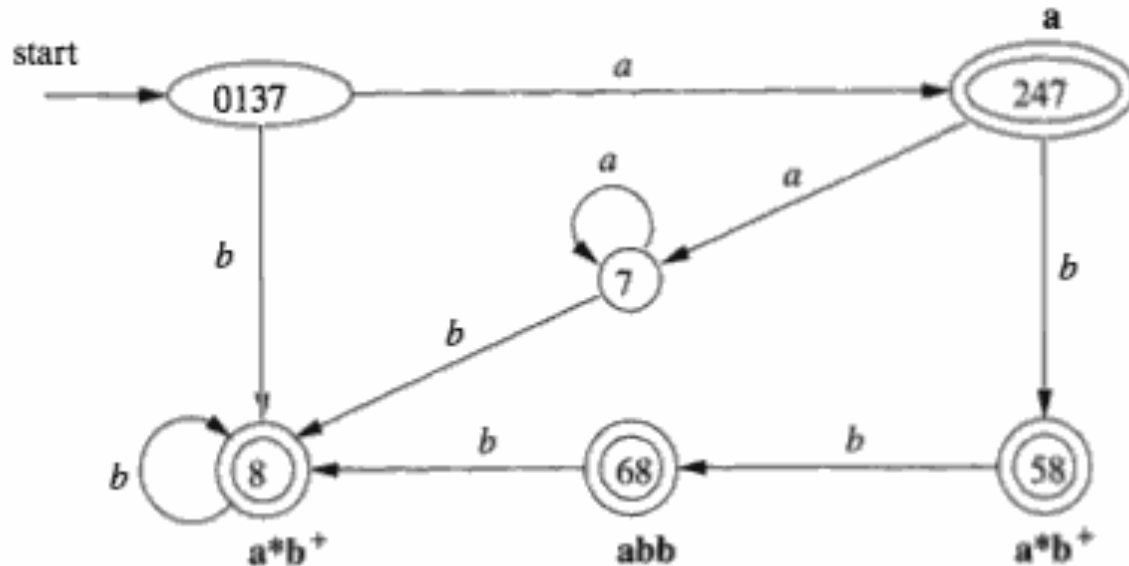
# Lex: Example (cont.)

- NFA's for lexical analyzer
- Input: **aaba**
  - $\epsilon$ -closure(0) = {0, 1, 3, 7}
  - Empty states after reading the fourth input symbol
    - There are no transitions out of state 8
    - Back up, looking for a set of states that include an accepting state
  - State 8:  $a^*b^+$  has been matched
    - Select **aab** as the lexeme, execute action  $A_3$
    - Return to parser indicating that token w/ pattern  $p_3=a^*b^+$  has been found



# Lex: Example (cont.)

- DFA's for lexical analyzer
- Input: **abba**
  - Sequence of states entered:  $0137 \rightarrow 247 \rightarrow 58 \rightarrow 68$
  - At the final a, there is no transition out of state 68
    - 68 itself is an accepting state that reports pattern  $p_2 = \text{abb}$





# How Much Should We Match?

- In general, find the **longest match** possible
  - We have seen examples
  - One more example: input string **aabbb ...**
    - Have many prefixes that match the third pattern
    - Continue reading *b*'s until another *a* is met
    - Report the lexeme to be the initial *a*'s followed by as many *b*'s as there are
- If same length, rule appearing first takes precedence
  - String **abb** matches both the second and third
  - We consider it as a lexeme for  $p_2$ , since that pattern listed first

<b>a</b>	{ <i>action</i> <sub>1</sub> }
<b>abb</b>	{ <i>action</i> <sub>2</sub> }
<b>a*b+</b>	{ <i>action</i> <sub>3</sub> }

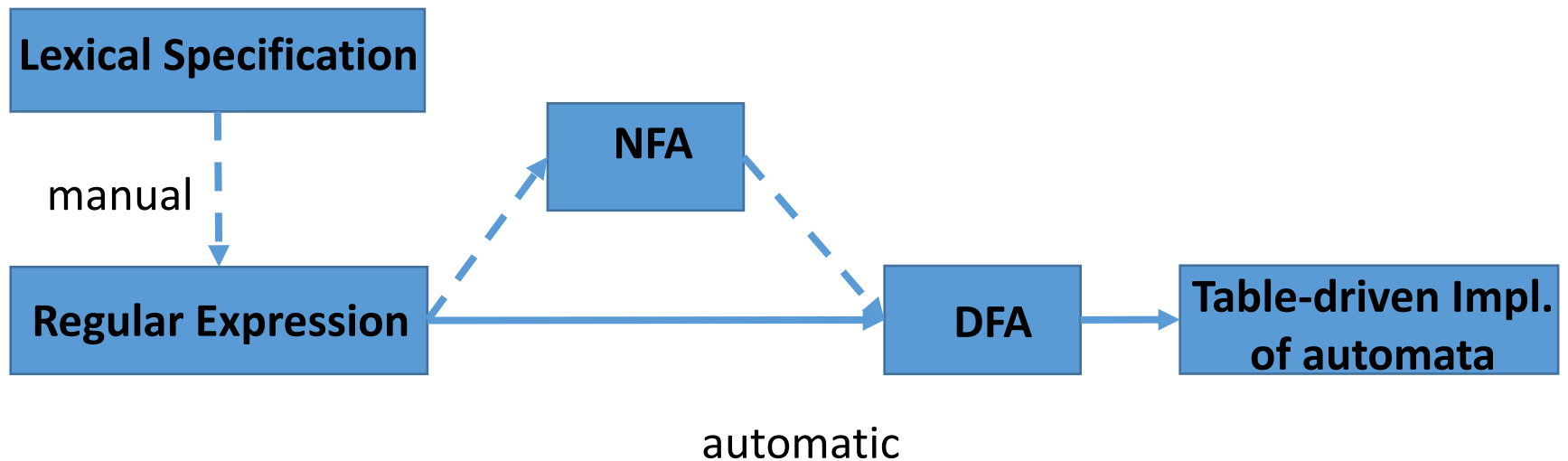
# How to Match Keywords?

---

- Example: to recognize the following tokens
  - Identifiers: `letter(letter|digit)*`
  - Keywords: `if, then, else`
- **Approach 1:** Make REs for keywords and place them before REs for identifiers so that they will take precedence
  - Will result in more bloated finite state machine
- **Approach 2:** Recognize keywords and identifiers using same RE but differentiate using special keyword table
  - Will result in more streamlined finite state machine
  - But extra table lookup is required
- Usually approach 2 is more efficient than 1, but you can implement approach 1 in your projects for simplicity

# Conversion Flow

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



# Beyond Regular Languages

---

- Regular languages are expressive enough for tokens
  - Can express identifiers, strings, comments, etc.
- However, it is the weakest (least expressive) language
  - Many languages are not regular
  - C programming language is not
    - The language matching braces “{{{...}}}” is also not
  - Finite automata cannot count # of times char encountered
    - $L = \{a^n b^n \mid n \geq 1\}$
    - Crucial for analyzing languages with nested structures (e.g. nested for loop in C language)
- We need a more powerful language for parsing
  - Later, we will discuss context-free languages (CFGs)