



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

DCS290, 3/7/2024

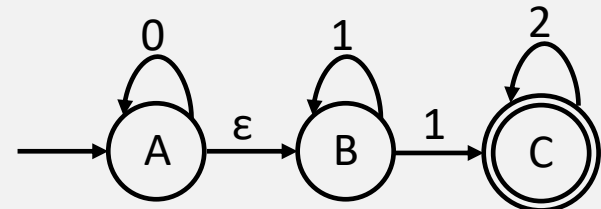


中山大學
SUN YAT-SEN UNIVERSITY



Review Questions

- Q1: RE of binary numbers that are multipliers of 2?
 $(0|1)^*0$
- Q2: meaning of $(a|b)^*bb(a|b)^*$?
Strings of a's and b's with consecutive b's
- Q3: usage of RE and FA in lexical analysis?
RE: specify the token class; FA: implement the token recognizer
- Q4: general workflow from RE to implementation?
RE \rightarrow NFA \rightarrow DFA \rightarrow Table
- Q5: the graph describes NFA or DFA? Why?
NFA. A: ϵ -transition, B: 1-transition

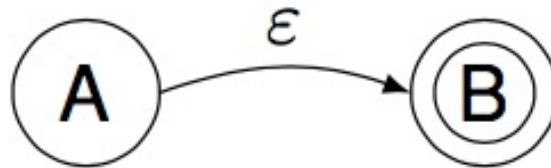


More on Table

- Implementation is efficient[表格是一种高效实现]
 - Table can be automatically generated
 - Need finite memory $O(S \times \Sigma)$
 - Size of transition table
 - Need finite time $O(\text{input length})$
 - Number of state transitions
- Pros and cons of table[表格实现的优劣]
 - **Pro**: can easily find the transitions on a given state and input
 - **Con**: takes a lot of space, when the input alphabet is large, yet most states do not have any moves on most of the input symbols

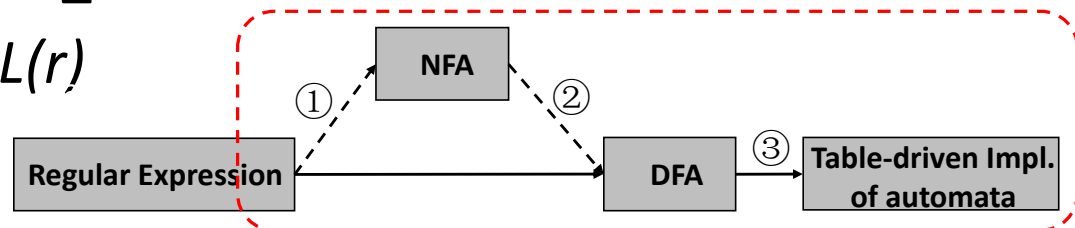
RE \rightarrow NFA

- NFA can have ϵ -moves
 - Edges labelled with ϵ
 - Move from state A to state B without reading any input



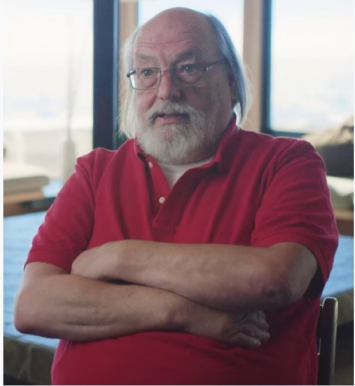
- **M-Y-T algorithm** (Thompson's construction) to convert any RE to an NFA that defines the same language [正则表达式转换到自动机]
 - Input: RE r over alphabet Σ
 - Output: NFA accepting $L(r)$

McNaughton-Yamada-Thompson



Thompson

Ken Thompson



Kenneth Lane Thompson (born February 4, 1943) is an American pioneer of [computer science](#) & Computer Chess Development. Thompson worked at [Bell Labs](#) for most of his career where he designed and implemented the original [Unix operating system](#). He also invented the [B programming language](#), the direct predecessor to the [C programming language](#), and was one of the creators and early developers of the [Plan 9 operating system](#). Since 2006, Thompson has worked at [Google](#), where he co-developed the [Go programming language](#).

Other notable contributions included his work on [regular expressions](#) and early computer text editors [QED](#) and [ed](#), the definition of the [UTF-8 encoding](#), and his work on computer chess that included the creation of [endgame tablebases](#) and the chess machine [Belle](#). He won the [Turing Award](#) in 1983 with his long-term colleague [Dennis Ritchie](#).

In the 1960s, Thompson also began work on [regular expressions](#). Thompson had developed the [CTSS](#) version of the editor [QED](#), which included regular expressions for searching text. QED and Thompson's later editor [ed](#) (the standard text editor on Unix) contributed greatly to the eventual popularity of regular expressions, and regular expressions became pervasive in Unix text processing programs. Almost all programs that work with regular expressions today use some variant of Thompson's notation. He also invented [Thompson's construction algorithm](#) used for converting regular expressions into [nondeterministic finite automata](#) in order to make expression matching faster.^[12]

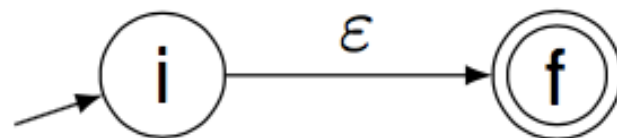
RE \rightarrow NFA (cont.)

- Step 1: processing atomic REs

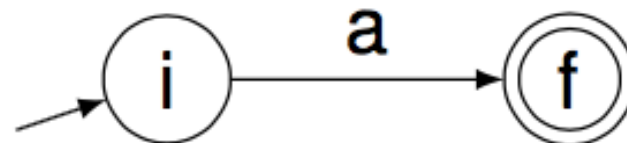
- ϵ expression[空]

- i is a new state, the start state of NFA

- f is another new state, the accepting state of NFA



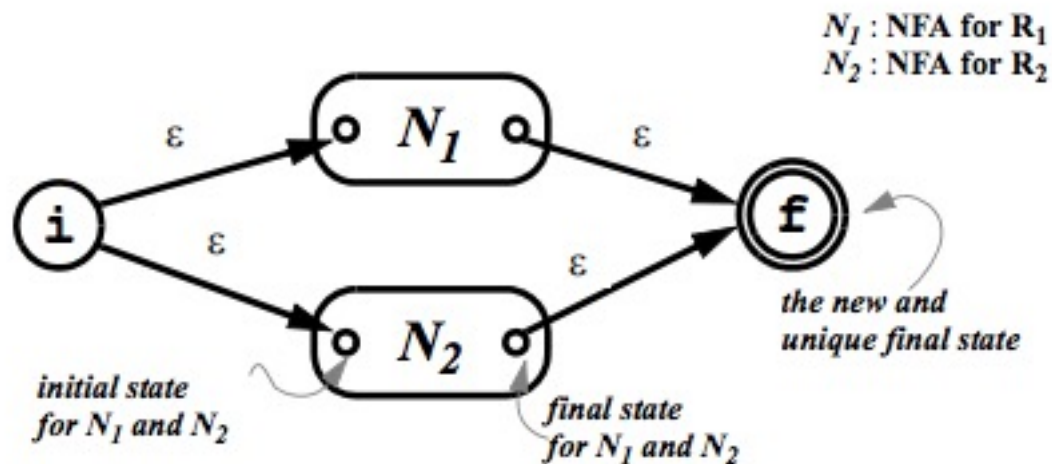
- Single character RE a [单字符]



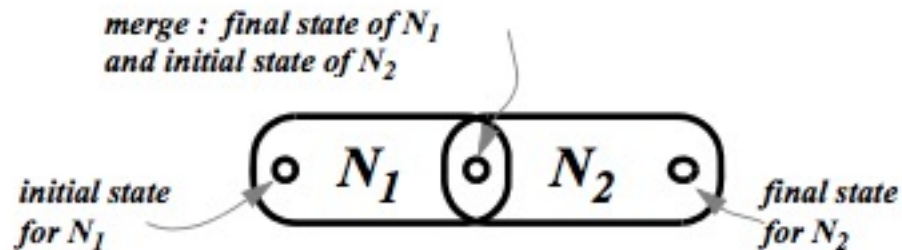
RE \rightarrow NFA (cont.)

- Step 2: processing compound REs[組合]

– $R = R_1 \mid R_2$

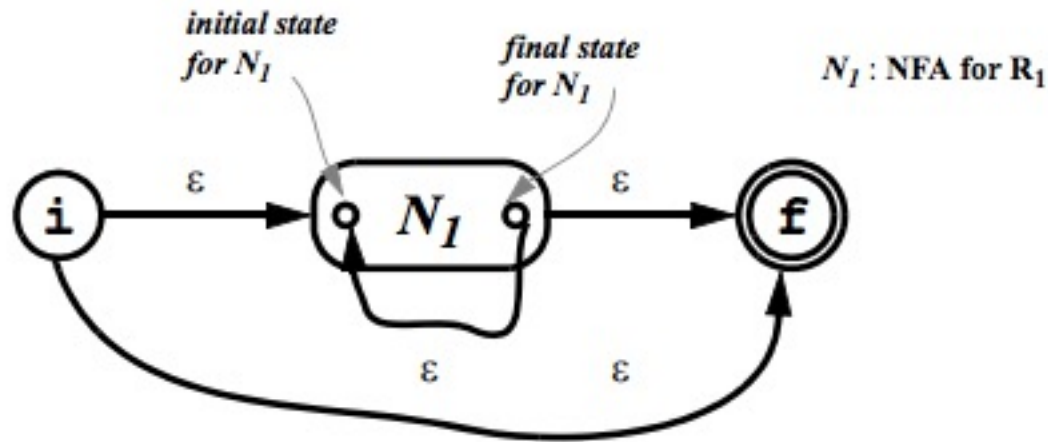


– $R = R_1 R_2$



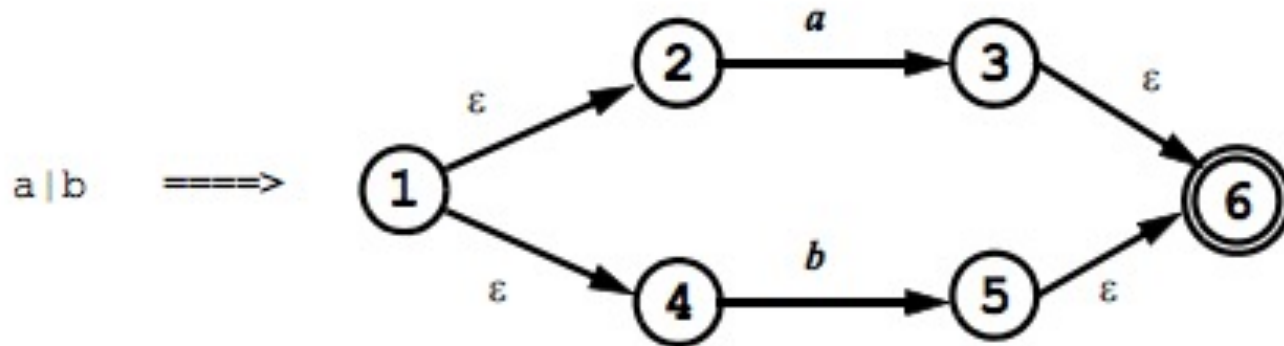
RE \rightarrow NFA (cont.)

- Step 2: processing compound REs
 - $R = R_1^*$



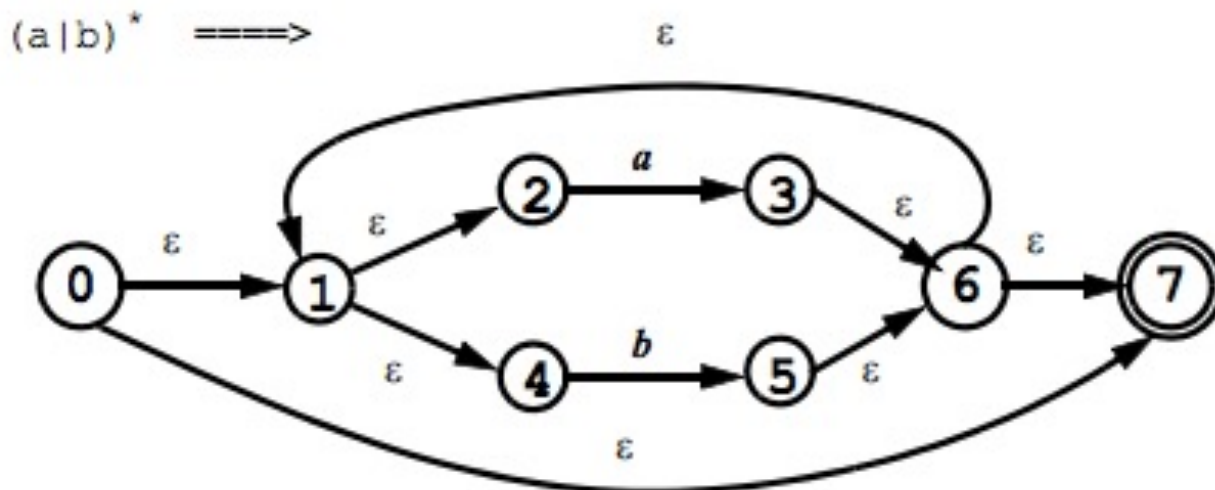
Example

- Convert “ $(a|b)^*abb$ ” to NFA

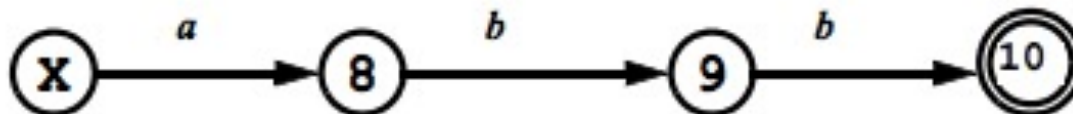


Example (cont.)

- Convert “ $(a|b)^*abb$ ” to NFA

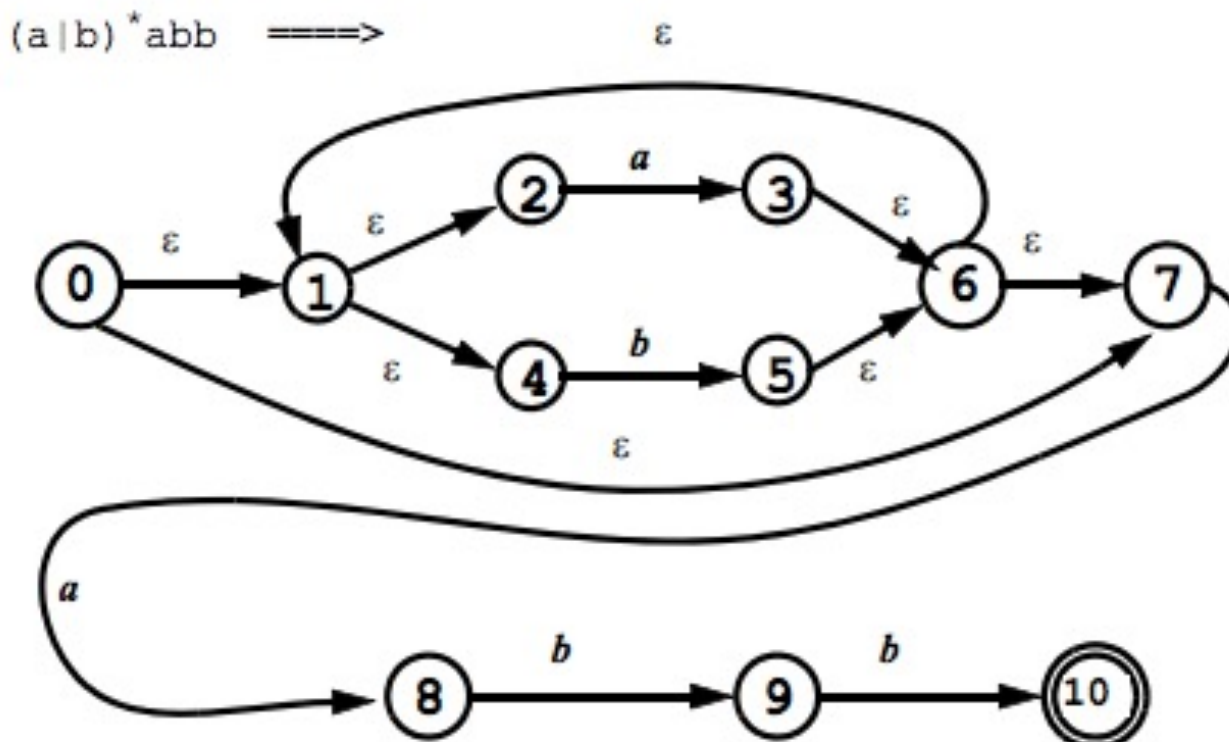


$abb \implies$ (several steps are omitted)



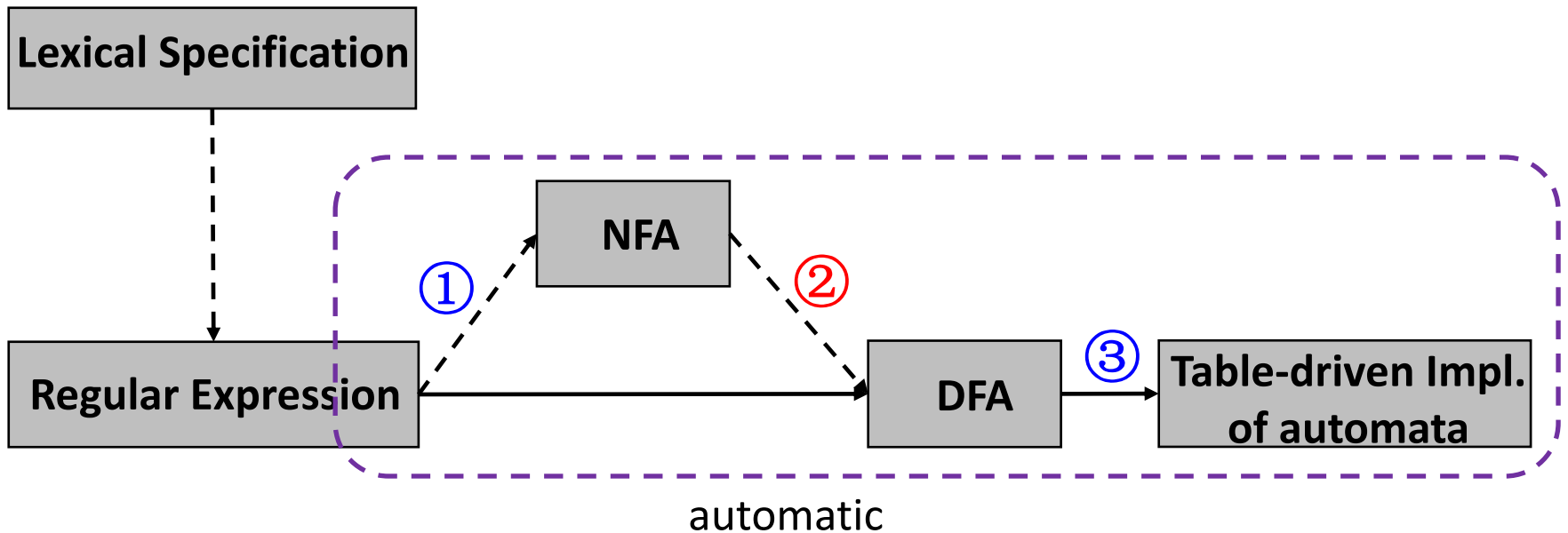
Example (cont.)

- Convert “ $(a|b)^*abb$ ” to NFA



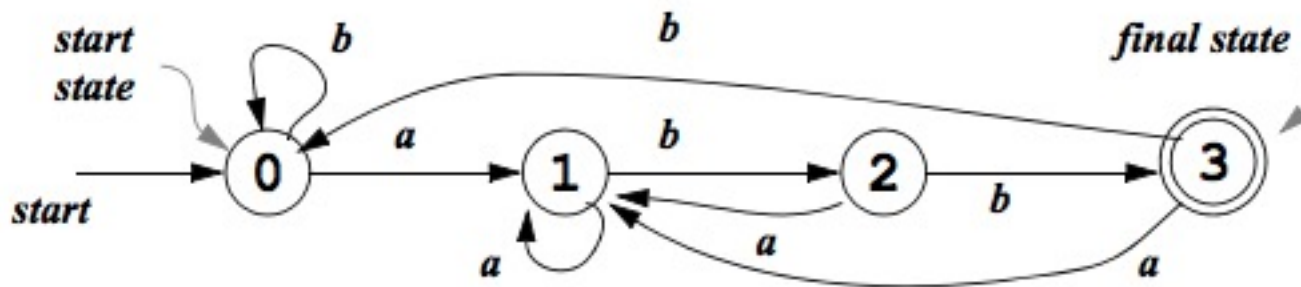
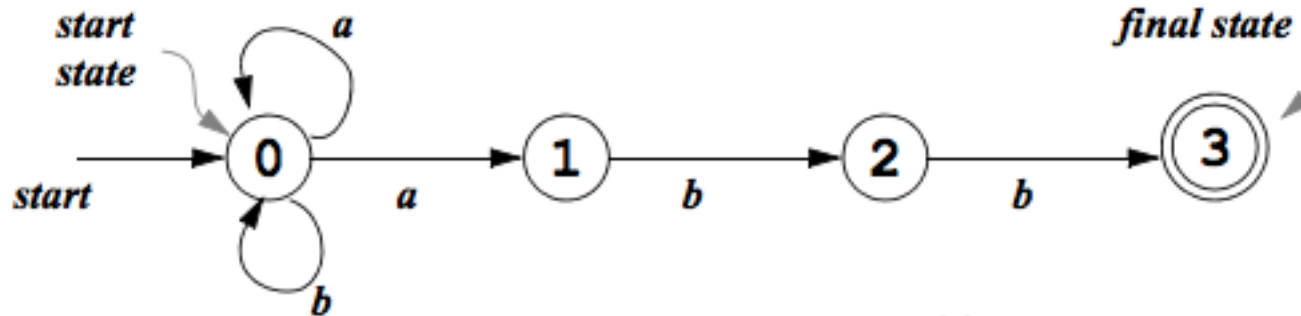
The Conversion Flow[转换流程]

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



NFA \rightarrow DFA: Same[等价]

- NFA and DFA are equivalent



To show this, we must prove every DFA can be converted into an NFA which accepts the same language, and vice-versa

NFA \rightarrow DFA: Theory[相关理论]

- 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 (NFA \rightarrow DFA)
 - Since $L(\text{DFA}) \subseteq L(\text{NFA})$, $L(\text{NFA}) \equiv L(\text{DFA})$
- Resulting DFA consumes more memory than NFA (Any DFA can be easily changed into NFA (e.g., add ϵ moves))
 - Potentially larger transition table as shown later
- But DFAs are faster to execute
 - For DFAs, number of transitions == length of input
 - For NFAs, number of potential transitions can be larger
- NFA \rightarrow DFA conversion is done because the speed of DFA far outweighs its extra memory consumption

NFA \rightarrow DFA: Idea

- Algorithm to convert[转换算法]
 - Input: an NFA N
 - Output: a DFA D accepting the same language as N
- **Subset construction**[子集构建]
 - Each state of the constructed DFA corresponds to a set of NFA states[一个DFA状态对应多个NFA状态]
 - Hence, the name ‘subset construction’
 - 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 $a_1a_2\dots a_n$

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[初始状态]
- 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 [下一状态]
 - This includes states that can be reached by reading a followed by any number of ϵ -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[接收状态]

NFA \rightarrow DFA: Algorithm

Initially, ε -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**

$U := \varepsilon$ -closure(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:

- ε -closure(s): set of NFA states reachable from NFA state s on ε -transitions **alone**

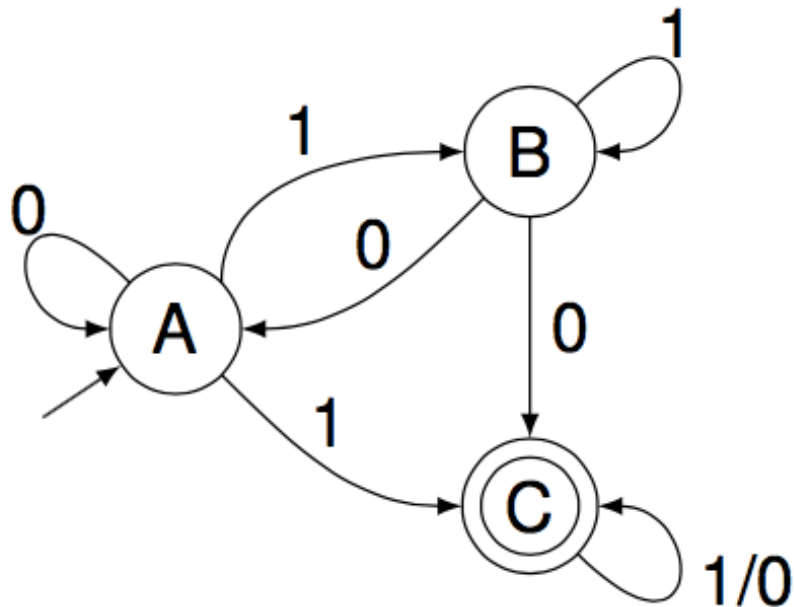
- ε -closure(T): set of NFA states reachable from some NFA state s in set T on ε -transitions **alone**; $= \bigcup_{s \in T} \varepsilon$ -closure(s)

- 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 ϵ -closure of the start state [初始状态]
 - ϵ -closure(A) = A
- Keep getting ϵ -closure($move(T, a)$) [更多状态]
- Stop, when there are no more new states

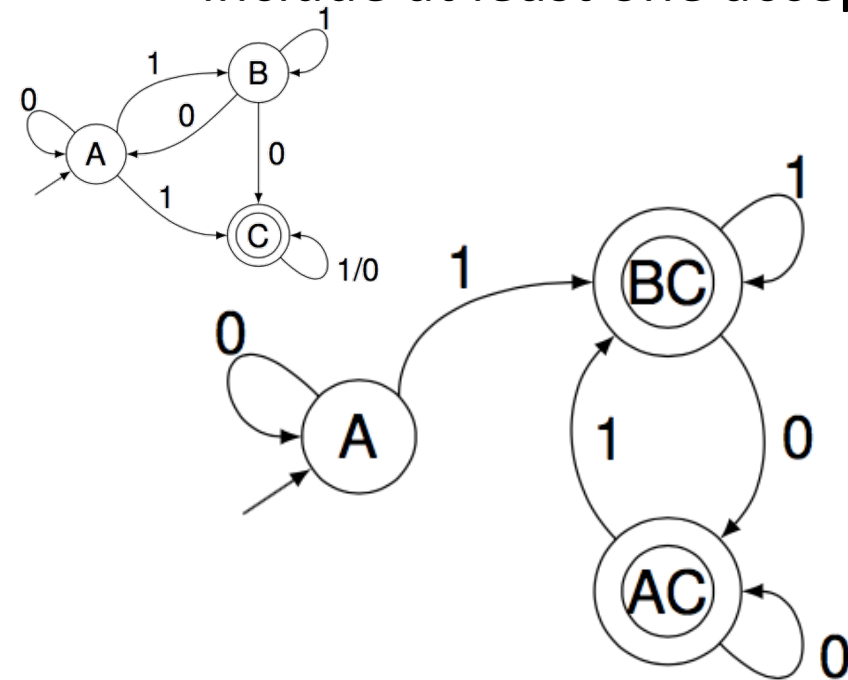
$T: A, a: 0/1$



		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 [终止状态]
 - The accepting states of D are all those sets of N 's states that include at least one accepting state of N



alphabet \rightarrow

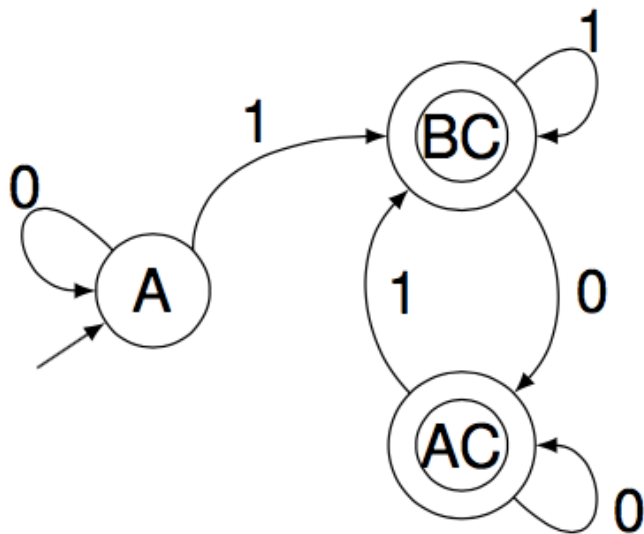
state \downarrow

	0	1
A	A	BC
BC	AC	BC
AC	AC	BC

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

NFA \rightarrow DFA: Minimization[最小化]

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



Initial: {A}, {BC, AC}

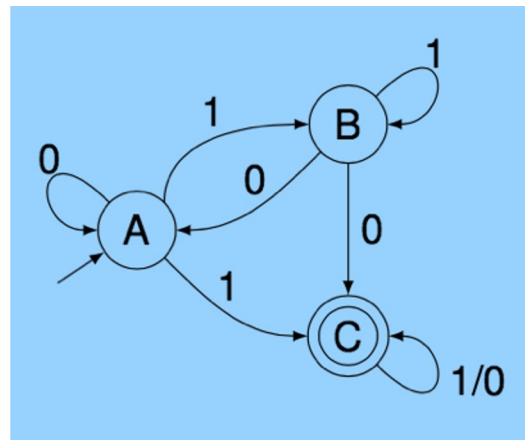
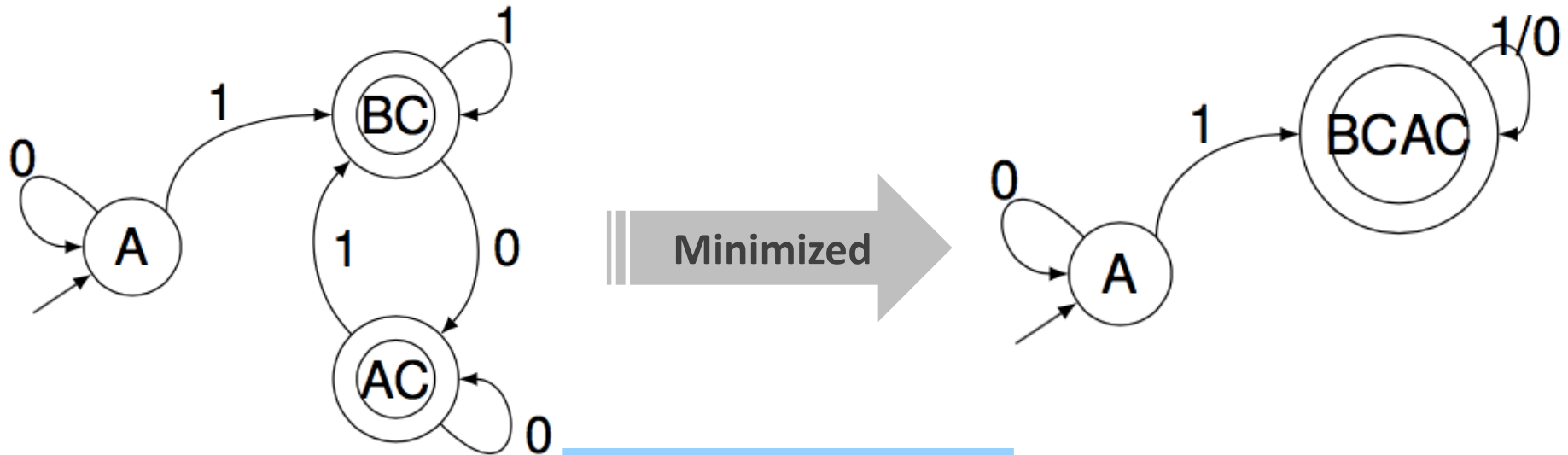
For {BC, AC} Initial sets:
{non-accepting states}, {accepting states}

- 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



	0	1
A	A	BC
BC	AC	BC
AC	AC	BC

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 $(\Sigma, S, n, F, \delta)$

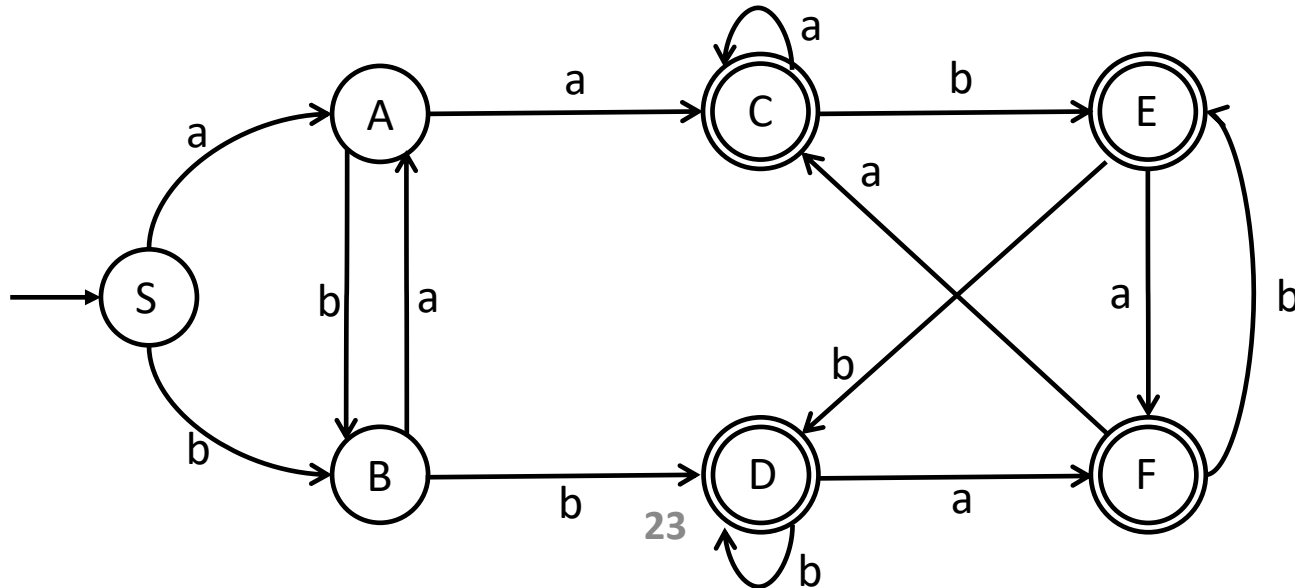
- The initial partition P_0 , has two sets and $\{S - F\}$
- Splitting a set (i.e., partitioning a set by input symbol α)

```
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 & s2
      T ← T ∪ s1 ∪ s2
  if T ≠ P then
    P ← T
```

- 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., α splits S)
- One state in the final DFA cannot have two transitions on α

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

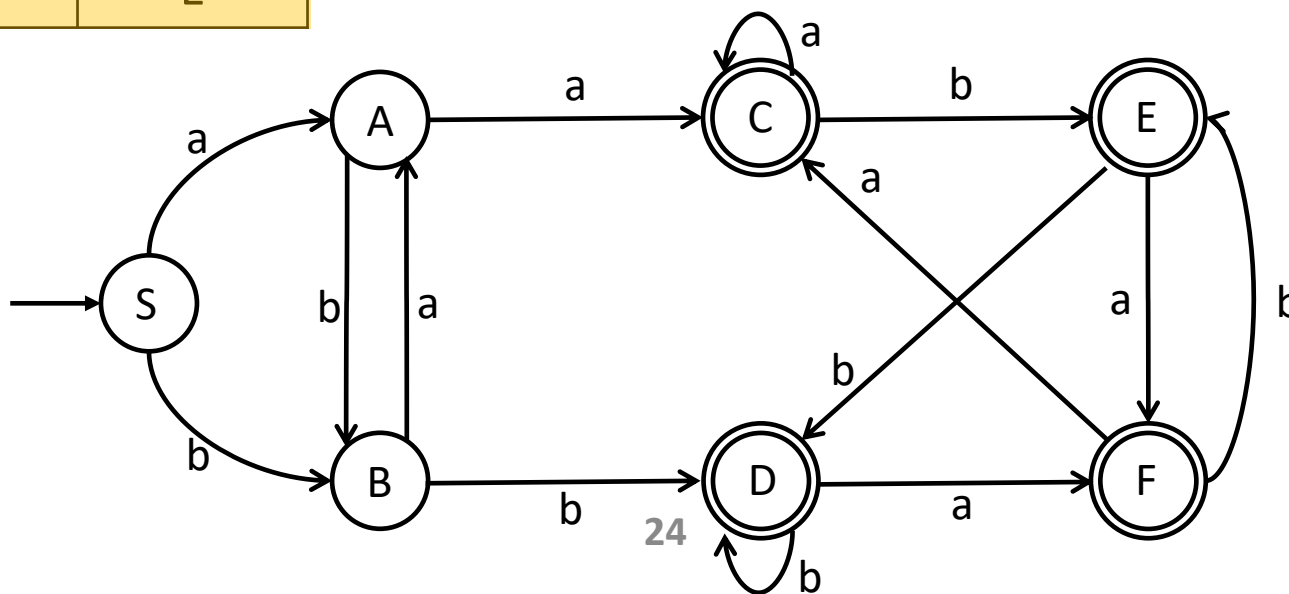


Example (cont.)

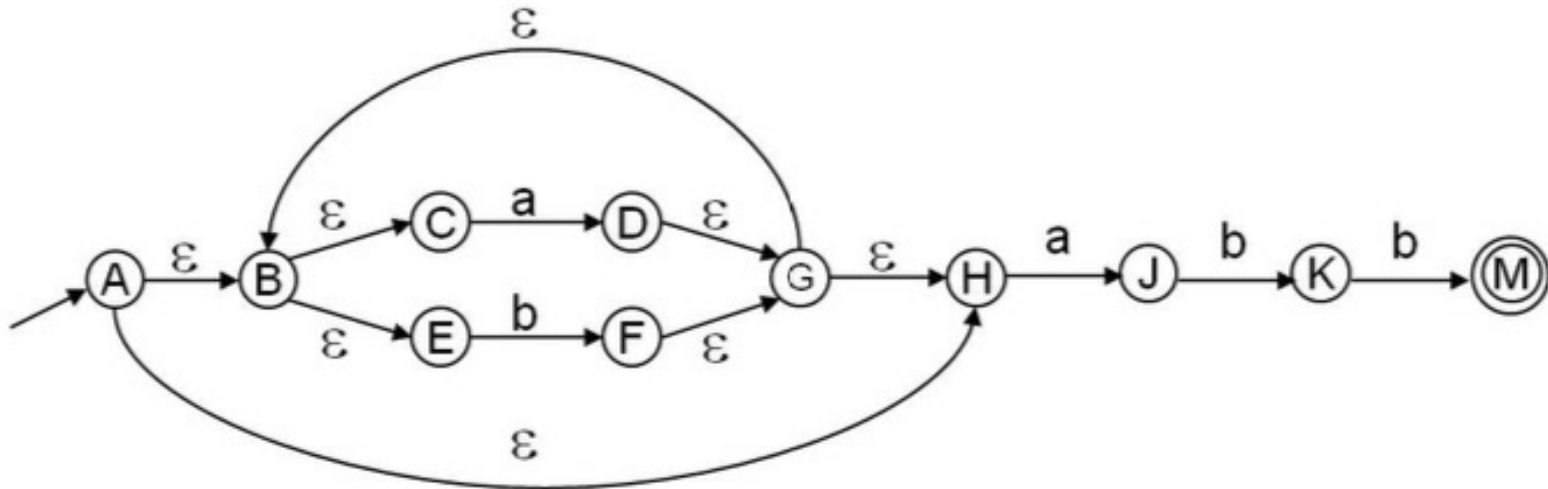
	a	b
S	A	B
A	C	B
B	A	D
C	C	E
D	F	D
E	F	D
F	C	E

	a	b
S	A	B
A	C	B
B	A	D
CF	C	E
DE	F	D

	a	b
S	A	B
A	C	B
B	A	D
CFDE	CF	DE



NFA \rightarrow DFA: More Example

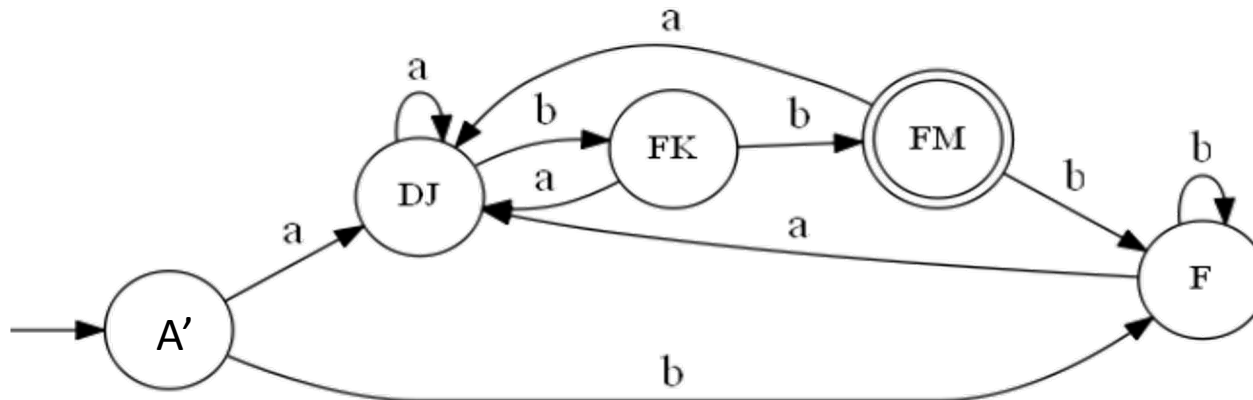


- Start state of the equivalent DFA
 - ϵ -closure(A) = {A, B, C, E, H} = A'
- ϵ -closure(move(A', a)) = ϵ -closure({D, J}) = {B, C, D, E, H, G, J} = B'
- ϵ -closure(move(A', b)) = ϵ -closure({F}) = {B, C, E, F, G, H} = C'
-

NFA \rightarrow DFA: More Example (cont.)

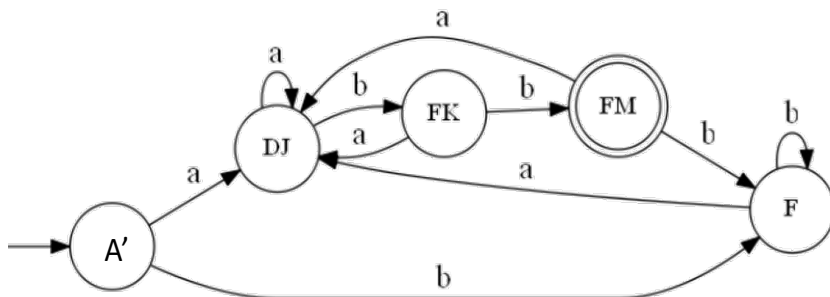
	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

- 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).



NFA \rightarrow DFA: More Example (cont.)

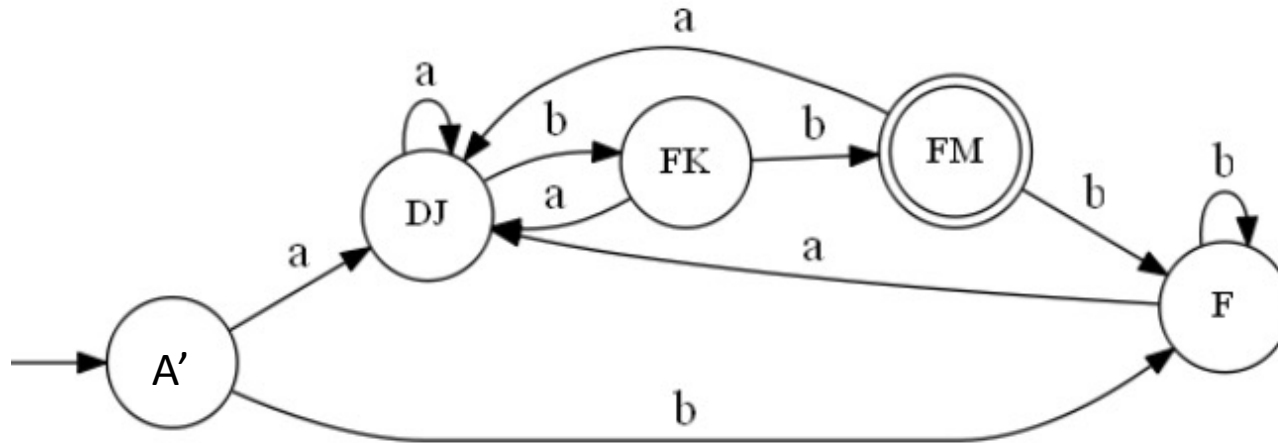
- P0: $s_1 = \{A', DJ, FK, F\}$, $s_2 = \{FM\}$
- For s_1 , further splits into $\{A', DJ, F\}$, $\{FK\}$
 - a: $A' \rightarrow DJ \in s_1$, $DJ \rightarrow DJ \in s_1$, $FK \rightarrow DJ \in s_1$, $F \rightarrow DJ \in s_1 \Rightarrow$ a doesn't distinct
 - b: $A' \rightarrow F \in s_1$, $DJ \rightarrow FK \in s_1$, $FK \rightarrow FM \in s_2$, $F \rightarrow F \in s_1 \Rightarrow$ b distincts $s_1 \Rightarrow s_{11} = \{A', DJ, F\}$, $s_{12} = \{FK\}$
- For s_{11} , further splits into $\{A', DJ, F\}$, $\{FK\}$
 - a: $A' \rightarrow DJ \in s_{11}$, $DJ \rightarrow DJ \in s_{11}$, $F \rightarrow DJ \in s_{11} \Rightarrow$ a doesn't distinct
 - b: $A' \rightarrow F \in s_{11}$, $DJ \rightarrow FK \in s_{12}$, $F \rightarrow DJ \in s_{11} \Rightarrow$ b distincts $s_{11} \Rightarrow s_{111} = \{A', F\}$, $s_{112} = \{DJ\}$
- For s_{111} , impossible to further split
- Final states: $S_{111} = \{A', F\}$, $S_{112} = \{DJ\}$, $S_{12} = \{FK\}$, $S_2 = \{FM\}$



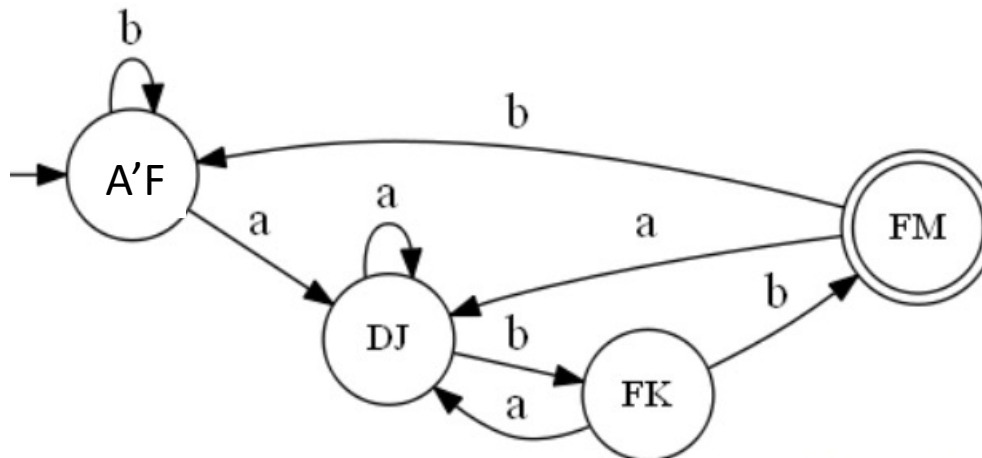
	a	b
A'	DJ	F
DJ	DJ	FK
F	DJ	F
FK	DJ	FM
FM	DJ	F

NFA \rightarrow DFA: More Example (cont.)

- 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