

# Chapter 3

## Regular Expressions, Regular Languages and Regular Grammars

# Regular Expression

- We have known that a regular language can be described by **some dfa**.
- Now we look at other ways of representing regular languages.

## Definition

Let  $\Sigma$  be a given alphabet. Then

1.  $\phi, \lambda$  and  $a \in \Sigma$  are all regular expressions. These are called primitive regular expressions.
2. If  $r_1$  and  $r_2$  are regular expressions, so are  $r_1 + r_2$ ,  $r_1 \cdot r_2$ ,  $r_1^*$  and  $(r_1)$ .
3. A string is a regular expression if and only if it can be derived from the primitive regular expressions by a finite number of applications of the rules in (2).

# Regular Expressions...

## Regular Expression (RE):

- E is a regular expression over  $\Sigma$  if E is one of:
  - $\varepsilon$
  - $a$ , where  $a \in \Sigma$
  - If  $r$  and  $s$  are **regular expressions (REs)**, then the following expressions also regular:
    - $r \mid s \rightarrow (r \text{ or } s)$
    - $rs \rightarrow (r \text{ followed by } s)$
    - $r^* \rightarrow (r \text{ repeated zero or more times})$
- **Each RE has an equivalent regular language (RL).**

# Regular Expressions...

- **Language Associated with Regular Expressions**

## Definition

The language  $L(r)$  denoted by any regular expression  $r$  is defined by the following rules.

1.  $\phi$  is a regular expression denoting the empty set,
2.  $\lambda$  is a regular expression denoting  $\{\lambda\}$ ,
3. for every  $a \in \Sigma$ ,  $a$  is a regular expression denoting  $\{a\}$

Termination  
condition

If  $r_1$  and  $r_2$  are regular expressions, then

4.  $L(r_1 + r_2) = L(r_1) \cup L(r_2)$ ,

5.  $L(r_1 \cdot r_2) = L(r_1)L(r_2)$ ,

6.  $L((r_1)) = L(r_1)$ ,

7.  $L(r_1^*) = (L(r_1))^*$ .

Reduce  $L(r)$  to simpler components.

# Regular Expressions...

## Regular Language (RL):

- L is a **regular language** over  $\Sigma$  if L is one of:
  - $\emptyset$  empty set
  - $\{\epsilon\}$  a set that contains empty string
  - $\{a\}$  where  $a \in \Sigma$
- If **R** and **S** are **regular languages (RL)**, then the following languages also regular:
  - $R \cup S = \{w \mid w \in R \text{ or } w \in S\}$
  - $RS = \{rs \mid r \in R \text{ and } s \in S\}$
  - $R^* = R^0 \cup R^1 \cup R^2 \cup R^3 \cup \dots$

# Regular Expressions...

## Rules for Specifying Regular Expressions:

1.  $\varepsilon$  is a regular expression  $\rightarrow L = \{\varepsilon\}$
2. If  $a$  is in  $\Sigma$ ,  $a$  is a regular expression  $\rightarrow L = \{a\}$ , the set containing the string  $a$ .
3. Let  $r$  and  $s$  be regular expressions with languages  $L(r)$  and  $L(s)$ . Then

- |  |   |    |  |                              |
|--|---|----|--|------------------------------|
| p<br>r<br>e<br>c<br>e<br>d<br>e<br>n<br>c<br>e | ↑ | a. | $r \mid s$ is a RE                                   | $\Rightarrow L(r) \cup L(s)$ |
|  |   | b. | $rs$ is a RE   | $\Rightarrow L(r) L(s)$      |
|  |   | c. | $r^*$ is a RE  | $\Rightarrow (L(r))^*$       |
|  |   | d. | $(r)$ is a RE $\Rightarrow L(r)$ , extra parenthesis |                              |

# Regular Expressions...

**abc**

- concatenation (“followed by”)

**a | b | c**

- alternation (“or”)

**\***

- zero or more occurrences

**+**

- one or more occurrences

# Regular Expressions...

- **Examples:**

- $\{0, 1\}\{00, 11\} = \{000, 011, 100, 111\}$
- $\{0\}^* = \{\varepsilon, 0, 00, 000, 0000, \dots\}$
- $\{\varepsilon\}^* = \{\varepsilon\}$
- $\{10, 01\}^* = \{\varepsilon, 10, 1010, 101010, \dots,$   
01, 0101, 010101, ...,  
1001, 100101, 10010101, ...,  
0110, 011010, 01101010, \dots\}

- **Notational shorthand:**

- $L^0 = \{\varepsilon\}$
- $L^i = LL^{i-1}$
- $L^+ = LL^*$



# Regular Expressions...

- Let  $L$  be a language over  $\{a, b\}$ , each string in  $L$  contains the substring **bb**
  - $L = \{a, b\}^* \{bb\} \{a, b\}^*$
- $L$  is regular language (RL). **Why?**
  - $\{a\}$  and  $\{b\}$  are RLs
  - $\{a, b\}$  is RL
  - $\{a, b\}^*$  is RL
  - $\{b\}\{b\} = \{bb\}$  is also RL
  - Then  $L = \{a, b\}^* \{bb\} \{a, b\}^*$  is RL

# Regular Expressions...

- Let  $L$  be a language over  $\{a, b\}$ , each string in  $L$ 
  - **begins** and **ends** with an **a**
  - contains at least one **b**
  - $L = \{a\}\{a, b\}^*\{b\}\{a, b\}^*\{a\}$
- $L$  is regular language (RL). **Why?**
  - $\{a\}$  and  $\{b\}$  are RLs
  - $\{a, b\}$  is RL
  - $\{a, b\}^*$  is RL
  - Then  $L = \{a\}\{a, b\}^*\{b\}\{a, b\}^*\{a\}$  is RL

# Regular Expressions...

- $L = \{a, b\}^* \{bb\} \{a, b\}^*$ 
  - $RE = (a|b)^* bb (a|b)^*$
- $L = \{a\} \{a, b\}^* \{b\} \{a, b\}^* \{a\}$ 
  - $RE = a(a|b)^* b(a|b)^* a$
- This  $RE = (a) | ((b)^*(c))$  is equivalent to  $a|b^*c$
- We say REs  $r$  and  $s$  are equivalent ( $r=s$ ), iff  $r$  and  $s$  represent **the same language**.
  - Example:  $r = a|b, s = b|a \rightarrow r = s$  Why?
  - Since  $L(r) = L(s) = \{a, b\}$

# Regular Expressions...

- Let  $\Sigma = \{a, b\}$ 
  - RE  $a|b$   $\rightarrow L = \{a, b\}$
  - RE  $(a|b)(a|b)$   $\rightarrow L = \{aa, ab, ba, bb\}$
  - RE  $aa|ab|ba|bb$  same as above
  - RE  $a^*$   $\rightarrow L = \{\epsilon, a, aa, aaa, \dots\}$
  - RE  $(a|b)^*$   $\rightarrow L = \text{set of all strings of a's and b's including } \epsilon$
  - RE  $(a^*b^*)^*$  same as above
  - RE  $a|a^*b$   $\rightarrow L = \{a, b, ab, aab, aaab, \dots\}$

# Regular Expressions...

---

## Algebraic Properties of regular Expressions

---

$$\mathbf{r \mid s = s \mid r}$$

---

$$\mathbf{r \mid (s \mid t) = (r \mid s) \mid t}$$

---

$$\mathbf{(r \ s) t = r (s \ t)}$$

---

$$\mathbf{r (s \mid t) = r s \mid r t}$$

$$\mathbf{(s \mid t) r = s r \mid t r}$$

---

$$\mathbf{\varepsilon r = r \ \varepsilon = r}$$

---

$$\mathbf{r^* = (r^*)^* = (r \mid \varepsilon)^+ = r^+ \mid \varepsilon}$$

---

$$\mathbf{r^{**} = r^*}$$

---

$$\mathbf{r^+ = r \ r^*}$$

---

# Regular Expression Examples

- All strings of 1s and 0s.

$(0 \mid 1)^*$

- All strings of 1s and 0s beginning with a 1.

$1(0 \mid 1)^*$

# Regular Expression Examples...

- All strings containing **two** or more **0**s.

$$(1|0)^* \mathbf{0} (1|0)^* \mathbf{0} (1|0)^*$$

- All strings containing an **even** number of **0**s.

$$(1^* \mathbf{0} 1^* \mathbf{0} 1^*)^* | 1^*$$

- All strings of alternating **0**s and **1**s.

$$(\epsilon | 1)(01)^*(\epsilon | 0)$$

- Strings over the alphabet  $\{0, 1\}$  with **no consecutive 0's**

- $(1 | 01)^*(0 | \epsilon)$

- $1^*(01^+)^*(0 | \epsilon)$

- $1^*(011^*)^*(0 | \epsilon)$

# Regular Expression Examples...

- Describe the following in **English**:
  - $(0|1)^*$ 
    - all strings over  $\{0, 1\}$
  - $b^*ab^*ab^*ab^*$  - describe the language?



# Regular Expression Examples...

- **More Examples of RE:**

- $01^*$ 
  - $\{0, 01, 011, 0111, \dots\}$
- $(01^*)(01)$ 
  - $\{001, 0101, 01101, 011101, \dots\}$
- $(0 \mid 1)^* = \{0, 1, 00, 01, 10, 11, \dots\}$ 
  - i.e., all strings of 0 and 1
- $(0 \mid 1)^* 00 (0 \mid 1)^* = \{00, 1001, \dots\}$ 
  - i.e., all 0 and 1 strings containing a “00”

# Regular Expression Examples...

- **More Examples of RE:**

- $(1 \mid 10)^*$ 
  - all strings starting with “1” and containing no “00”
- $(0 \mid 1)^*011$ 
  - all strings ending with “011”
- $0^*1^*$ 
  - all strings with no “0” after “1”
- $00^*11^*$ 
  - all strings with at least one “0” and one “1”, and no “0” after “1”

# Exercise

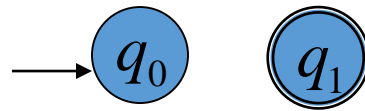
- What languages do the following RE represent?
  - $((0 \mid 1)(0 \mid 1))^* \mid ((0 \mid 1)(0 \mid 1)(0 \mid 1))^*$

# Connection Between RE & RL

- A language  $L$  is called **regular** if and only if there exists some **DFA**  $M$  such that  $L = L(M)$ .
- Since a DFA *has an equivalent* NFA, then
  - A language  $L$  is called **regular** if and only if there exists some NFA  $N$  such that  $L = L(N)$ .
- If we have a RE  $r$ , we can construct an NFA that accept  $L(r)$ .

# Connection Between RE & RL...

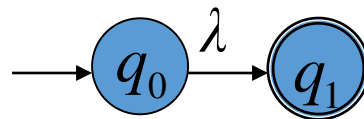
- For  $\emptyset$  in the regular expression, construct NFA



$$L = \{ \} = \emptyset$$

(1) nfa accepts  $\emptyset$ .

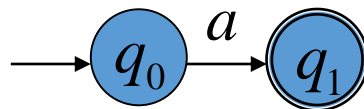
- For  $\epsilon$  in the regular expression, construct NFA



$$L = \{ \epsilon \}$$

(2) nfa accepts  $\{\lambda\}$ .

- For  $a \in \Sigma$  in the regular expression, construct NFA

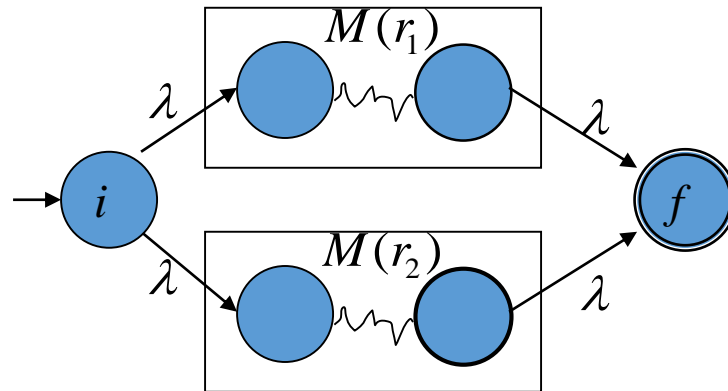


$$L = \{ a \}$$

(3) nfa accepts  $a$ .

# Connection Between RE & RL...

- If  $r_1$  and  $r_2$  are regular expressions,  $M_{r_1}$  and  $M_{r_2}$  are their NFAs.  $r_1|r_2$  has NFA:



$$L = \{L(M_{r_1}) \cup L(M_{r_2})\}$$

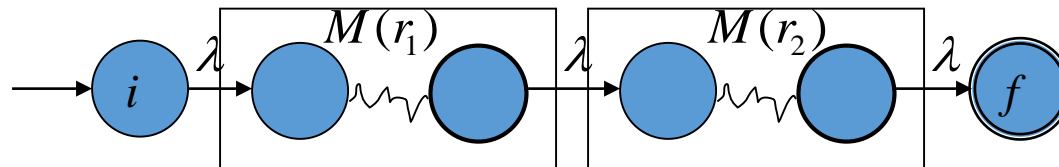
(4) nfa accepts  $L(r_1 + r_2)$ .

where  $i$  and  $f$  are new start / final states, and  $\epsilon$ -moves are introduced from  $i$  to the old start states of  $M_{r_1}$  and  $M_{r_2}$  as well as from all of their final states to  $f$ .

# Connection Between RE & RL...

- If **r1** and **r2** are regular expressions,  $M_{r1}$ ,  $M_{r2}$  their NFAs.
- **r1r2** (concatenation) has NFA:

$$L = \{L(M_{r1})L(M_{r2})\}$$



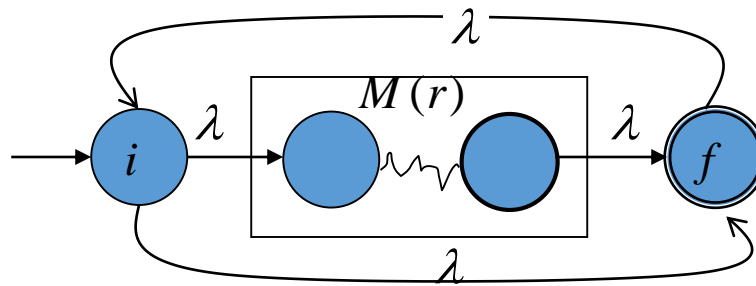
(5) nfa accepts  $L(r_1r_2)$ .

where  **$i$**  is the start state of  $M_{r1}$  (or new under the alternative) and  **$f$**  is the final state of  $M_{r2}$  (or new). Overlap maps final states of  $M_{r1}$  to start state of  $M_{r2}$

# Connection Between RE & RL...

- If  $r$  is a regular expressions and  $M_r$  its NFA,
- $r^*$  (Kleene star) has NFA:

$$L = \{L(M_r)^*\}$$



(6) nfa accepts  $L(r^*)$ .

where :  $i$  is new start state and  $f$  is new final state

$\epsilon$ -move  $i$  to  $f$  (to accept null string)

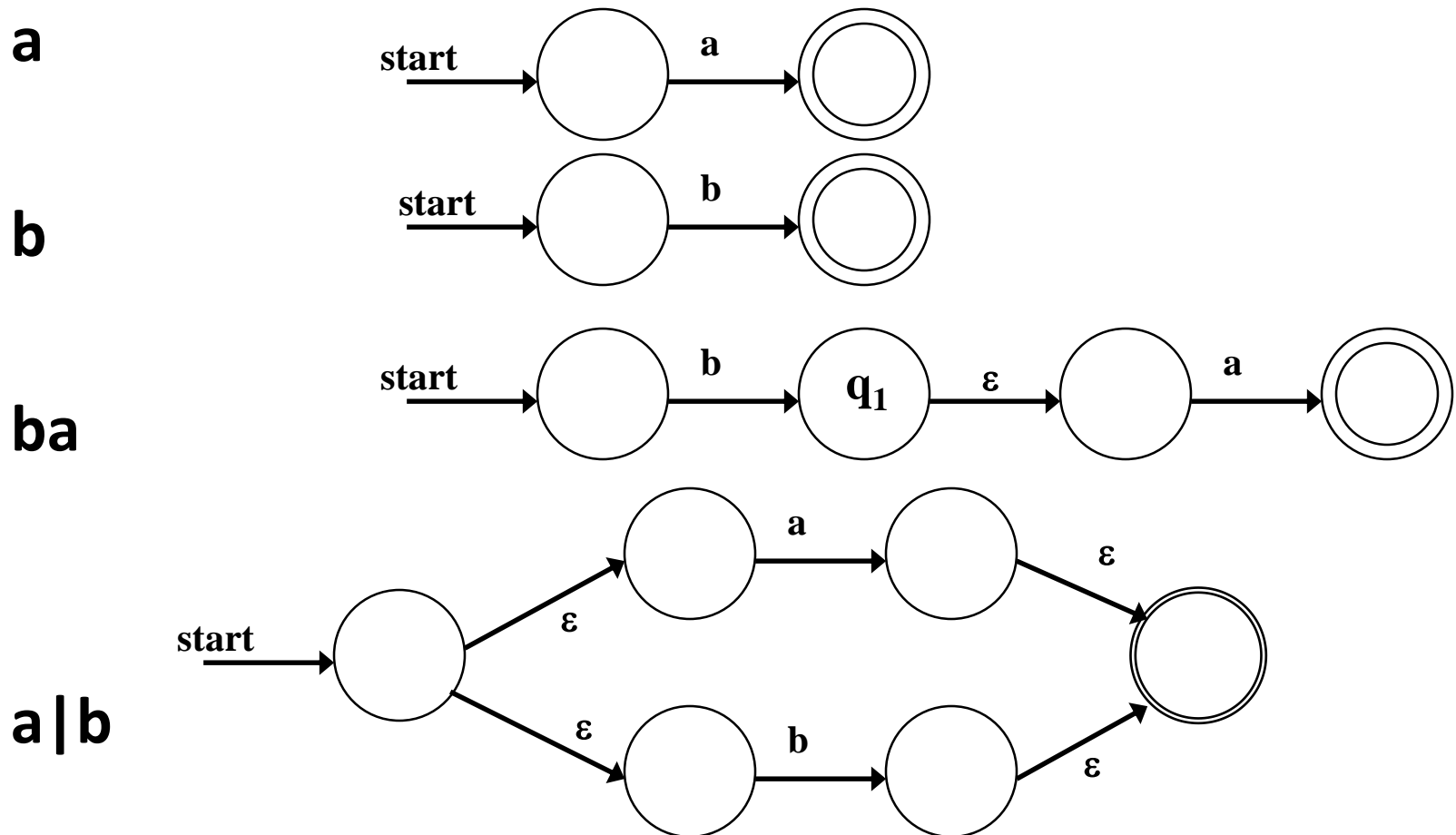
$\epsilon$ -moves  $i$  to old start, old final(s) to  $f$

$\epsilon$ -move  $f$  to  $i$  (**WHY?**)



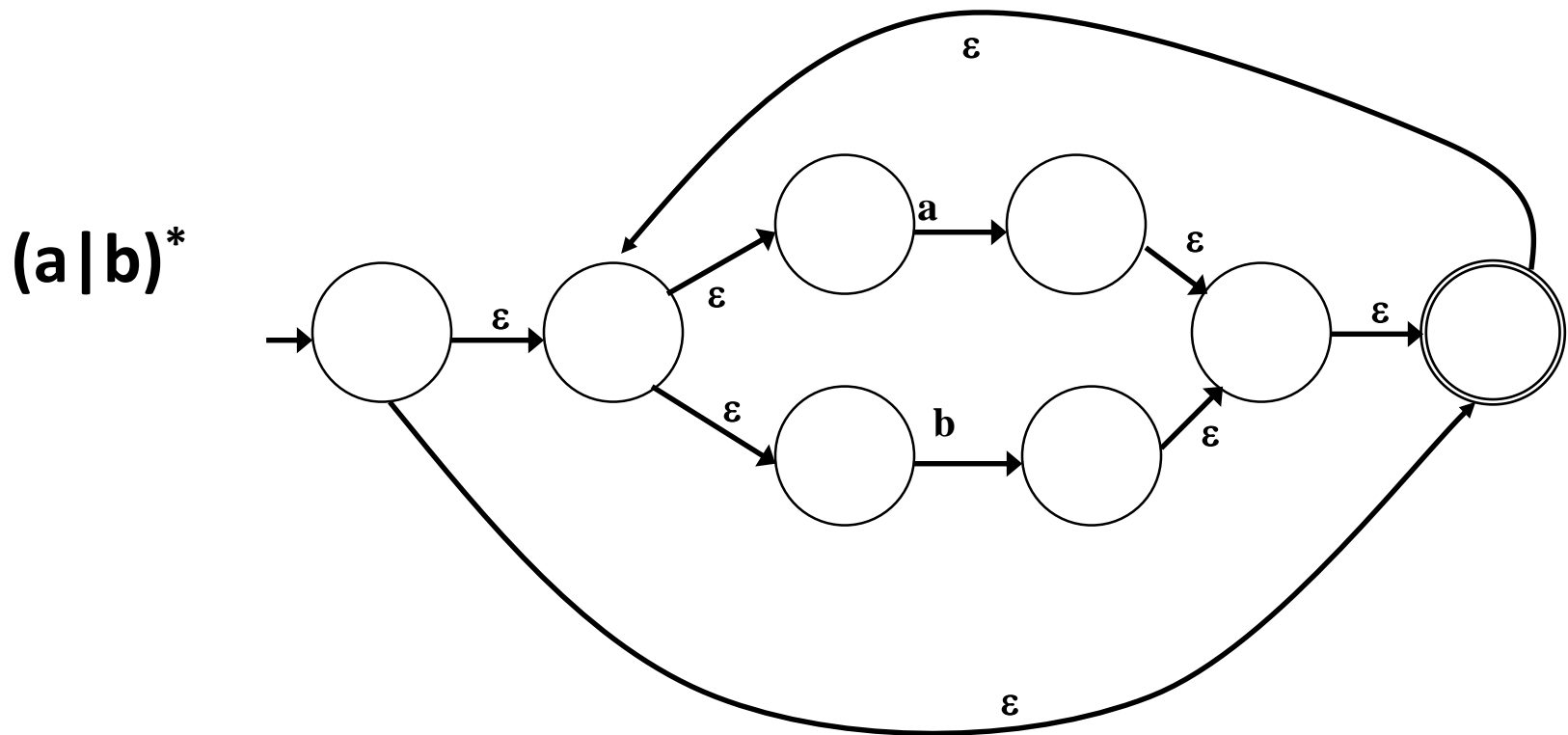
# Connection Between RE & RL: **Example**

- Build an **NFA- $\epsilon$**  that accepts  **$(a|b)^*ba$**



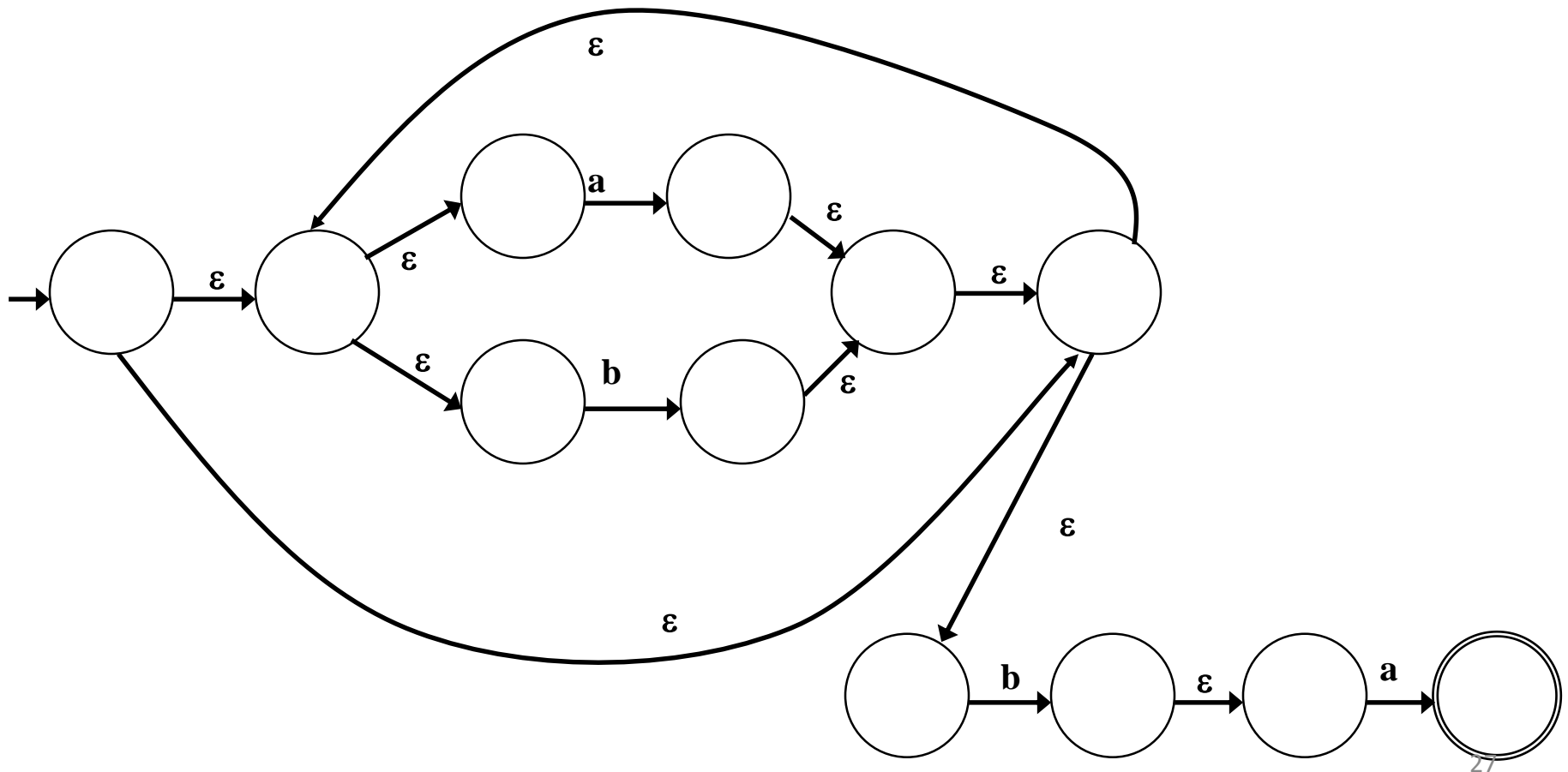
# Example...

- Build an **NFA- $\epsilon$**  that accepts  **$(a|b)^*ba$**



# Example...

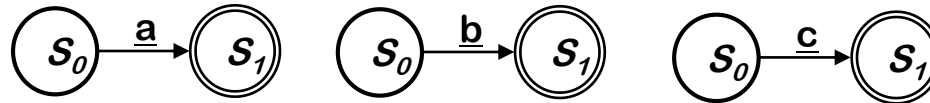
- Build an **NFA- $\epsilon$**  that accepts  **$(a|b)^*ba$**



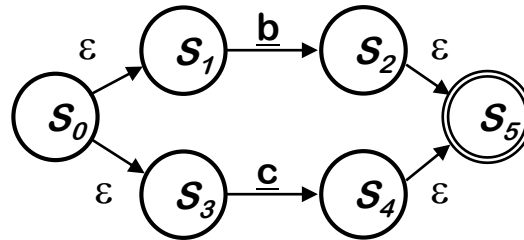
# Example...

- Let's try  $\underline{a}(\underline{b} \mid \underline{c})^*$

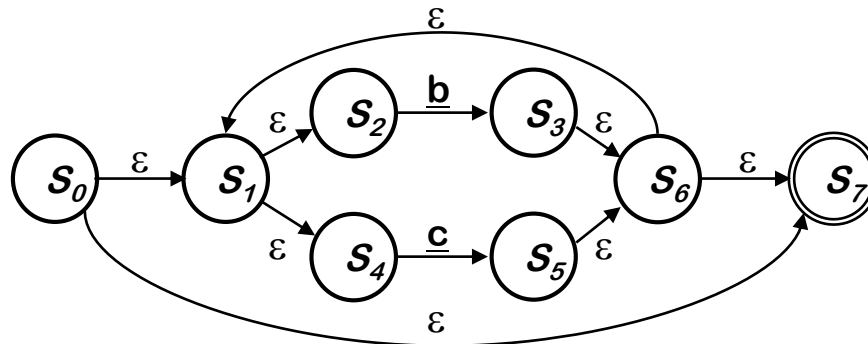
1.  $\underline{a}$ ,  $\underline{b}$ , &  $\underline{c}$



2.  $\underline{b} \mid \underline{c}$

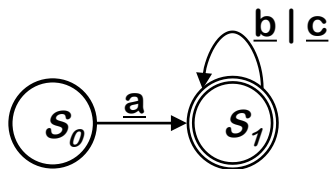
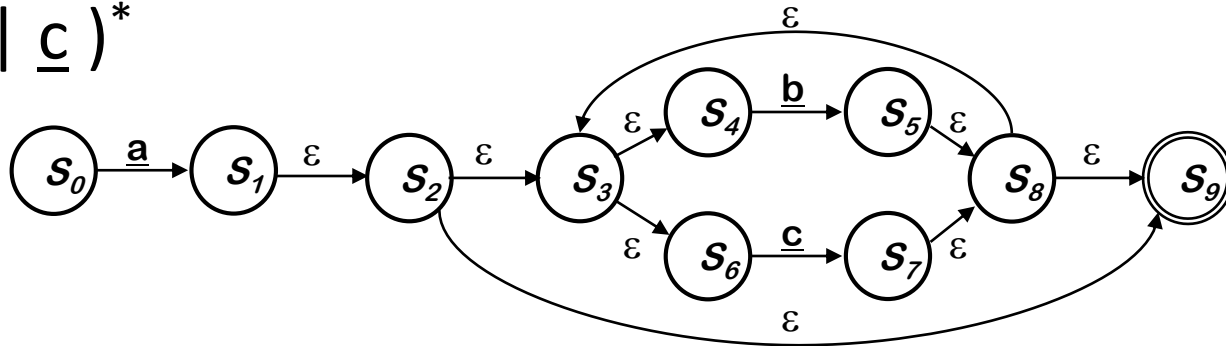


3.  $(\underline{b} \mid \underline{c})^*$



# Example...

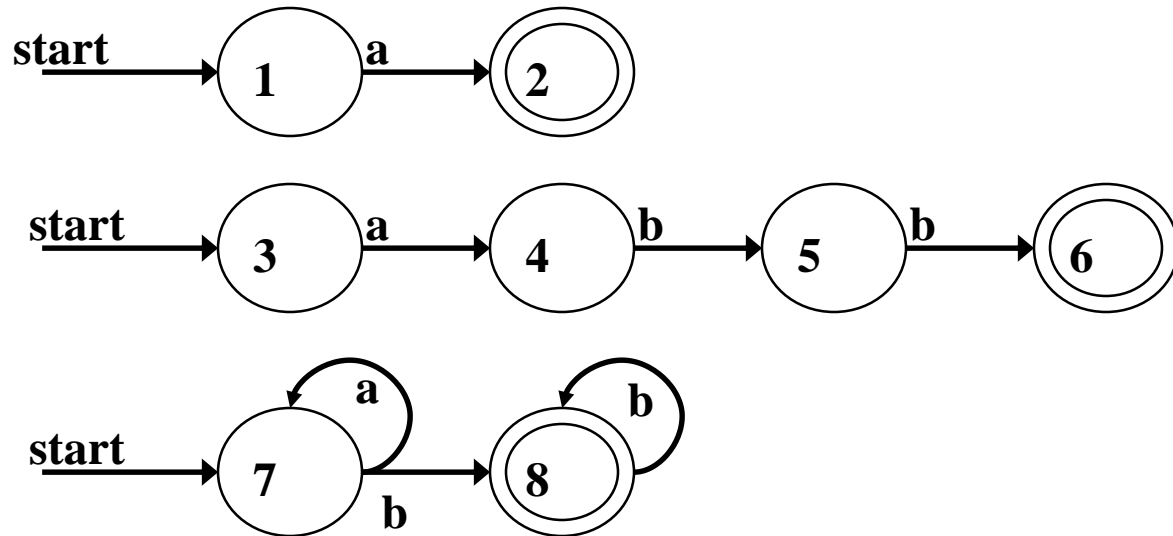
4.  $\underline{a}(\underline{b} \mid \underline{c})^*$



# Example...

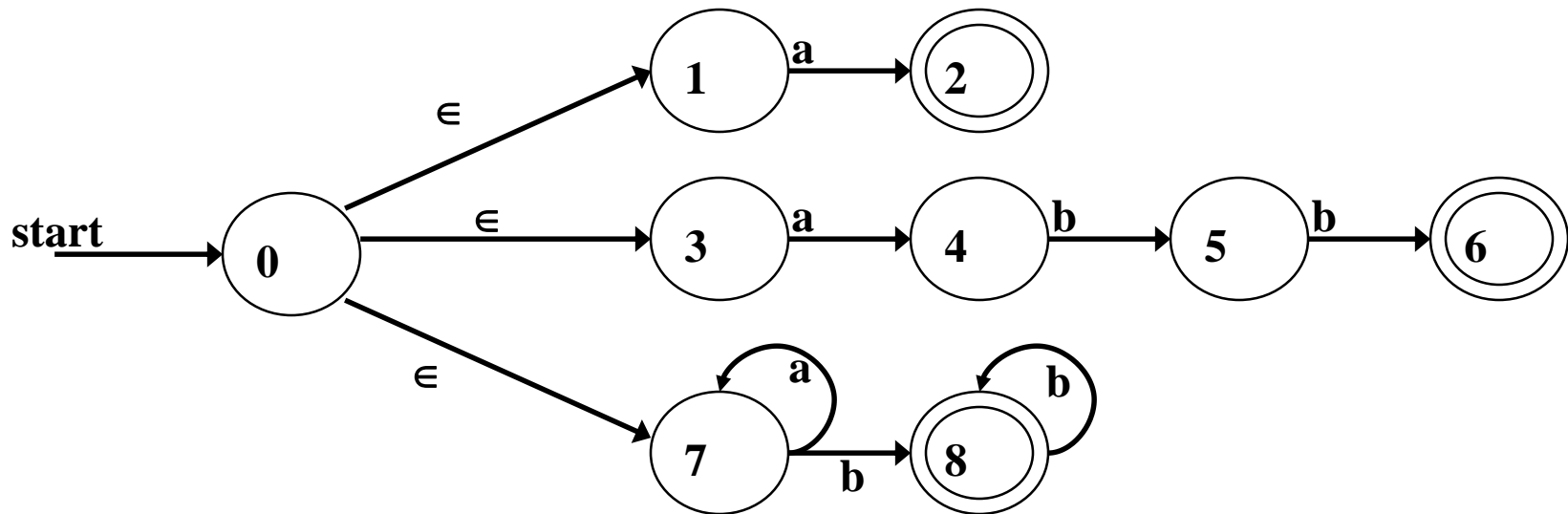
Let :  $a$   
       $abb$   
       $a^*b^+$  } 3 patterns

NFA's :



# Example...

NFA for :  $a \mid abb \mid a^*b^+$

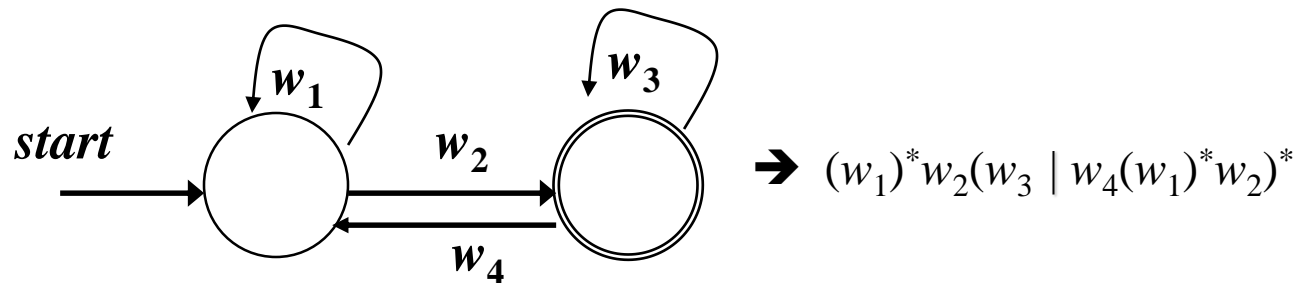
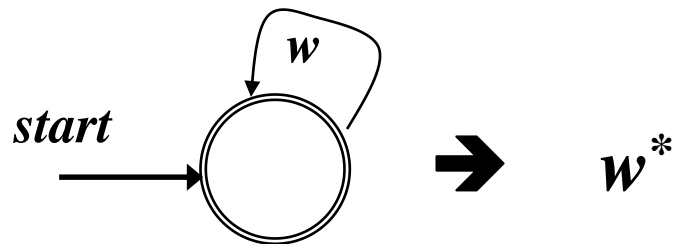
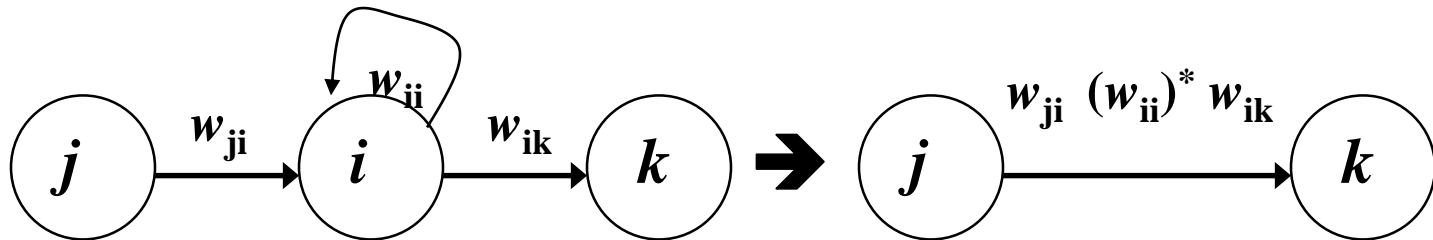
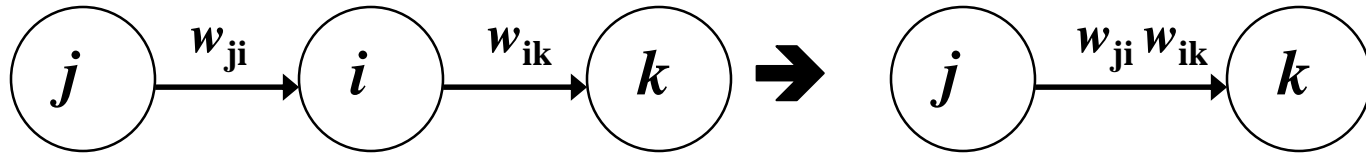


# NFA to RE

- If  $L$  is accepted by some **NFA- $\epsilon$** , then  $L$  is represented by some ***regular expression***.
- An ***expression graph*** is like a state diagram but it can have ***regular expressions*** as labels on arcs.
- An NFA- $\epsilon$  is an ***expression graph***.
- An expression graph can be ***reduced to one*** with just ***two states***.
- If we reduce an NFA- $\epsilon$  in this way, the arc label then corresponds to the regular expression representing it.

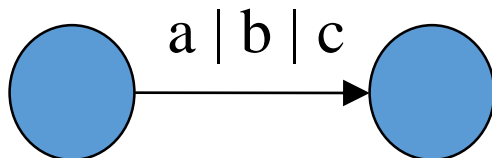
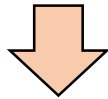
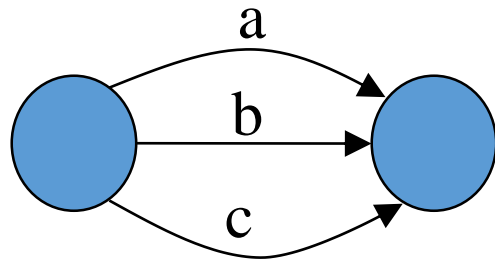


# NFA to RE...

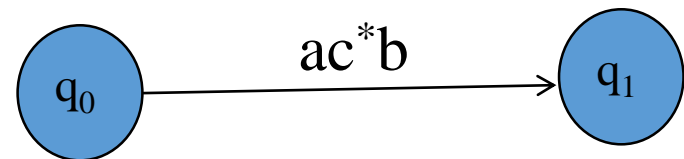
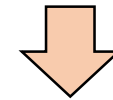
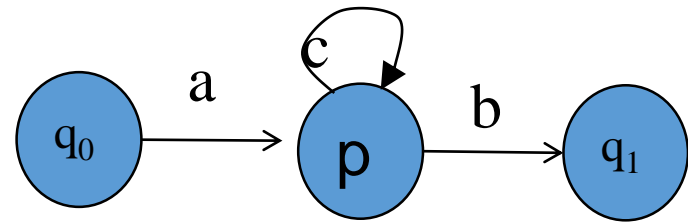


# NFA to RE...

- Merge Edges :



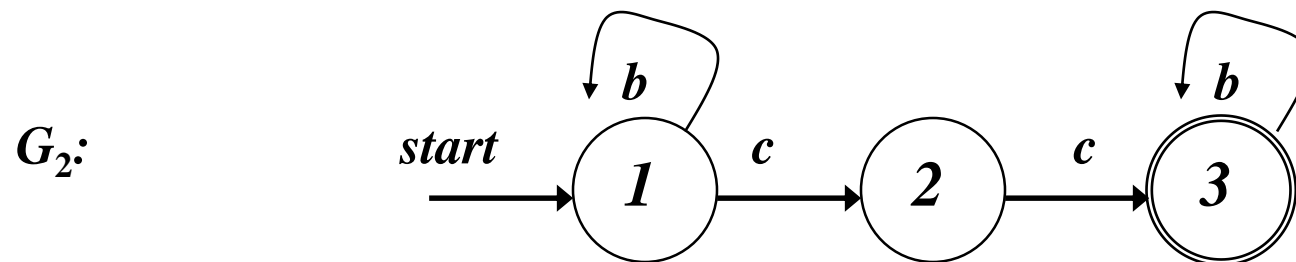
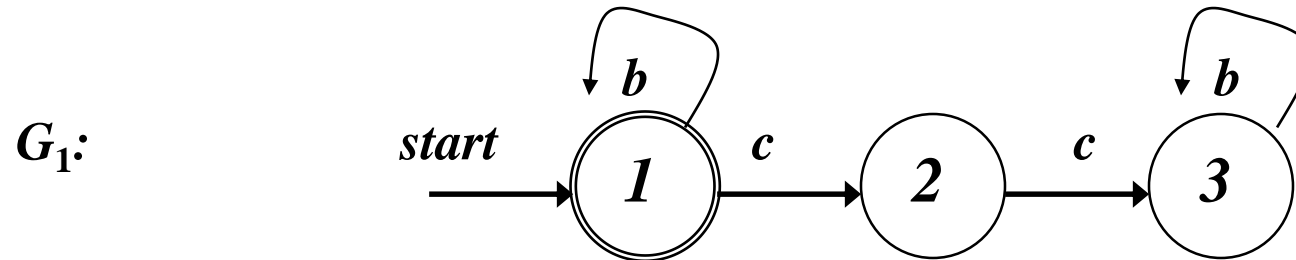
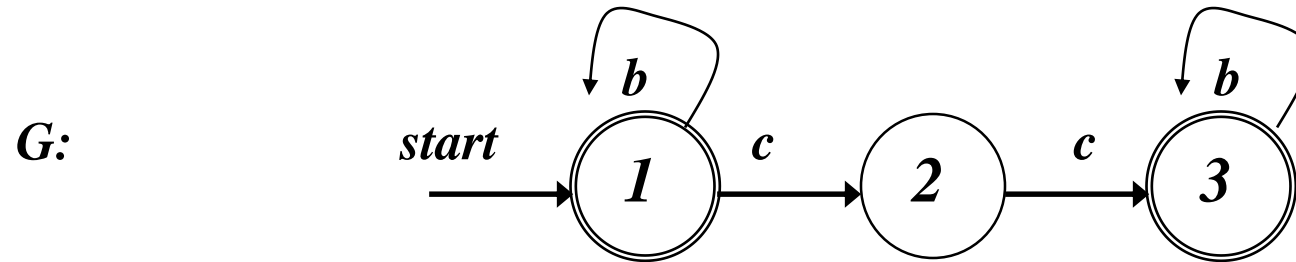
- Replace state by Edges



# NFA to RE...

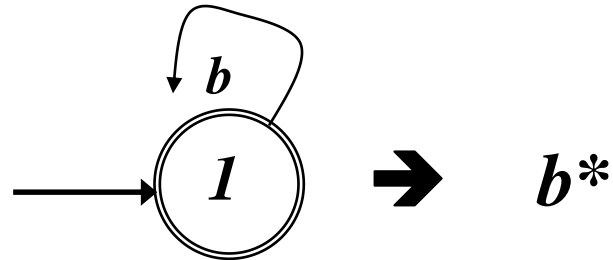
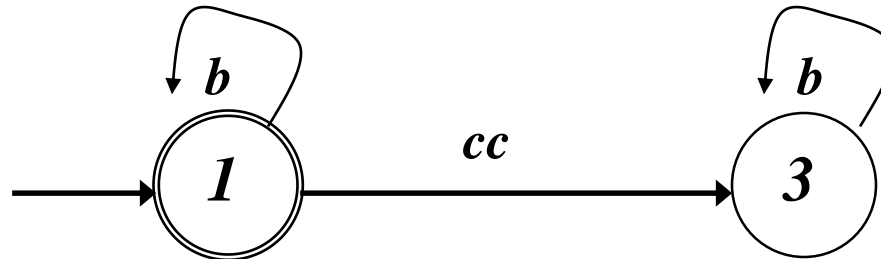
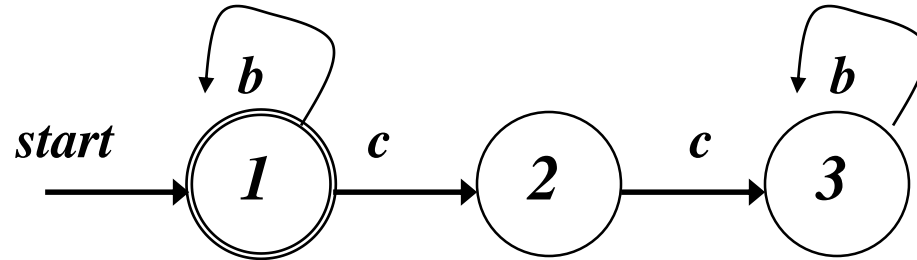
- Let  $\mathbf{G}$  be the state diagram of a finite automata.
- Let  $m$  be the *number of final states* of  $\mathbf{G}$ .
- Make  $m$  copies of  $\mathbf{G}$ , each of which has ***one final state***.
- Call these graphs  $\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_m$
- For each  $\mathbf{G}_t$ 
  - Repeat
    - Do the steps in the previous slide
    - Until the only states in  $\mathbf{G}_t$  are the start state and the single final state.
    - Determine the RE of  $\mathbf{G}_t$
- The RE of  $\mathbf{G}$  is obtained by joining RE's of each  $\mathbf{G}_t$  by  $\cup$  or  $|$ .

# NFA to RE...



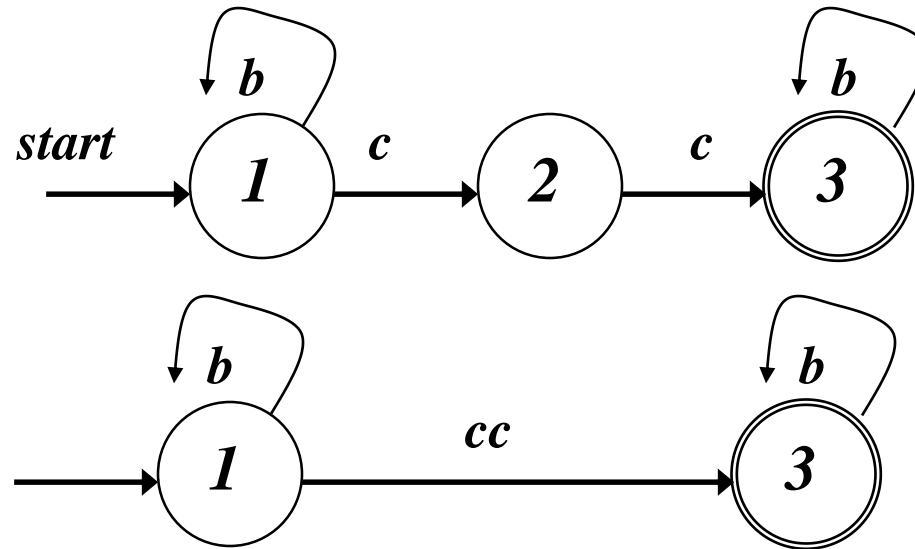
# NFA to RE...

$G_1$ :



# NFA to RE...

$G_2$ :

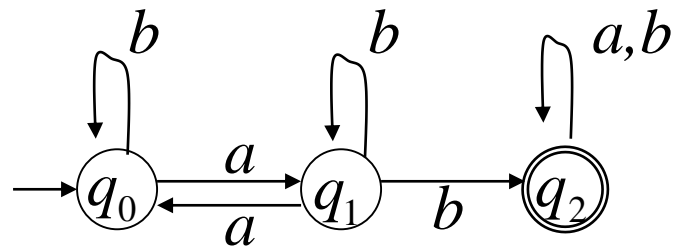


$\Rightarrow b^*ccb^*$

RE for  $G \Rightarrow b^* \mid b^*ccb^*$

# NFA to RE: Exercise

- Given the following NFA, find the regular expression it accepts.



# Applications of Regular Expressions

- In many programming languages the set of integer constants is defined by the regular expression ***sdd\****, where ***s*** stands for the sign, with possible values from  $\{+, -, \lambda\}$ , and ***d*** stands for the digits 0 to 9.
- An application of pattern matching occurs in ***text editing***.
- All text editors allow files to be scanned for the occurrence of a given string; most editors extend this to permit searching for patterns.
  - **For example**, the ***vi*** editor in the ***UNIX operating system*** recognizes the command ***/aba\*c/*** as an instruction to search the file for the first occurrence of the string ***ab***, followed by an arbitrary number of ***a***'s, followed by ***ac***.



# Regular Grammars

- The ***third way*** of describing regular languages is by means of certain simple grammars. (*what are the other two ways?*)

## Right- and Left-Linear Grammars

A grammar  $G = (V, T, S, P)$  is said to be right - linear if all productions are of the form

$$A \rightarrow xB,$$

$$A \rightarrow x,$$

where  $A, B \in V$ , and  $x \in T^*$ . A grammar is said to be left - linear if all productions are of the form

$$A \rightarrow Bx,$$

$$A \rightarrow x.$$

A regular grammar is one that is either right - linear or left - linear.

# Regular Grammars...

## Example

The grammar  $G_1 = (\{S\}, \{a, b\}, S, P_1)$ , with  $P_1$  given as

$$S \rightarrow abS/a \quad \text{is right linear.}$$

The grammar  $G_2 = (\{S, S_1, S_2\}, \{a, b\}, S, P_2)$ , with  $P_2$  given as

$$S \rightarrow S_1ab, \quad S_1 \rightarrow S_1ab/S_2, \quad S_2 \rightarrow a \quad \text{is left linear.}$$

Both  $G_1$  and  $G_2$  are regular grammars.

$$\begin{array}{l} \longrightarrow L(G_1) = L((ab)^*a) \\ L(G_2) = L(aab(ab)^*) \end{array}$$

## Example

The grammar  $G_1 = (\{S, A, B\}, \{a, b\}, S, P)$ , with  $P$  given as

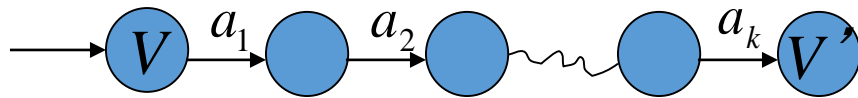
$$S \rightarrow A, \quad A \rightarrow aB/\lambda, \quad B \rightarrow Ab \quad \text{is not regular.}$$

# Right-linear Grammars Generate Regular Languages

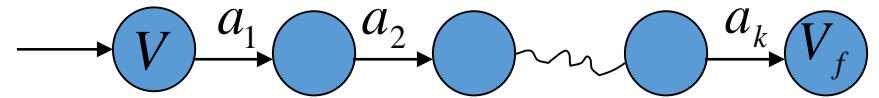
- Right-linear grammar  $\longleftrightarrow$  Regular language
- Left-linear grammar  $\longleftrightarrow$  Regular language

**Theorem 3.3** Let  $G=(V,T,S,P)$  be a right-linear grammar. Then  $L(G)$  is a regular language.

Policy of Proof: Construct an *nfa* which accepts  $L(G)$ .



(1) *nfa* for  $V \rightarrow a_1 a_2 \dots a_k V'$

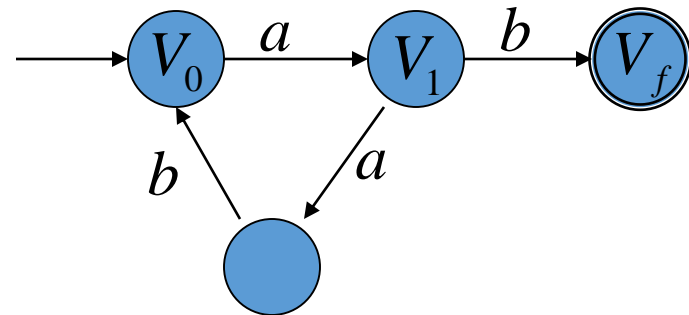


(2) *nfa* for  $V \rightarrow a_1 a_2 \dots a_k$

Example: Construct a finite automaton that accepts the language generated by the grammar

$$V_0 \rightarrow aV_1,$$

$$V_1 \rightarrow abV_0/b.$$



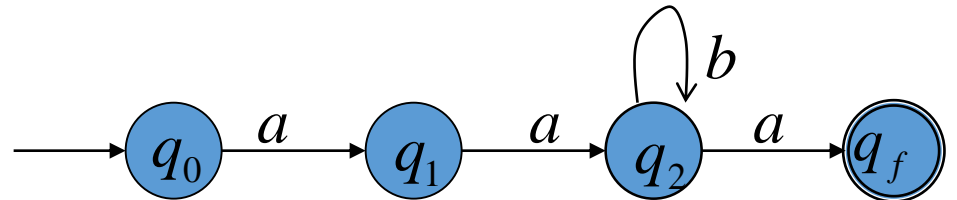
**Theorem 3.4** If  $L$  is a regular language on an alphabet  $\Sigma$  then there exists a right-linear grammar  $G=(V, \Sigma, S, P)$  such that  $L=L(G)$ .

Policy of Proof: Let  $M$  be the dfa that accepts  $L$ . Construct the right-linear grammar  $G$  from  $M$  such that  $G$  generates  $L$ .

Example: Construct a right - linear grammar for  $L(aab^*a)$ .

**Solution**

A dfa which accepts  $L(aab^*a)$ .



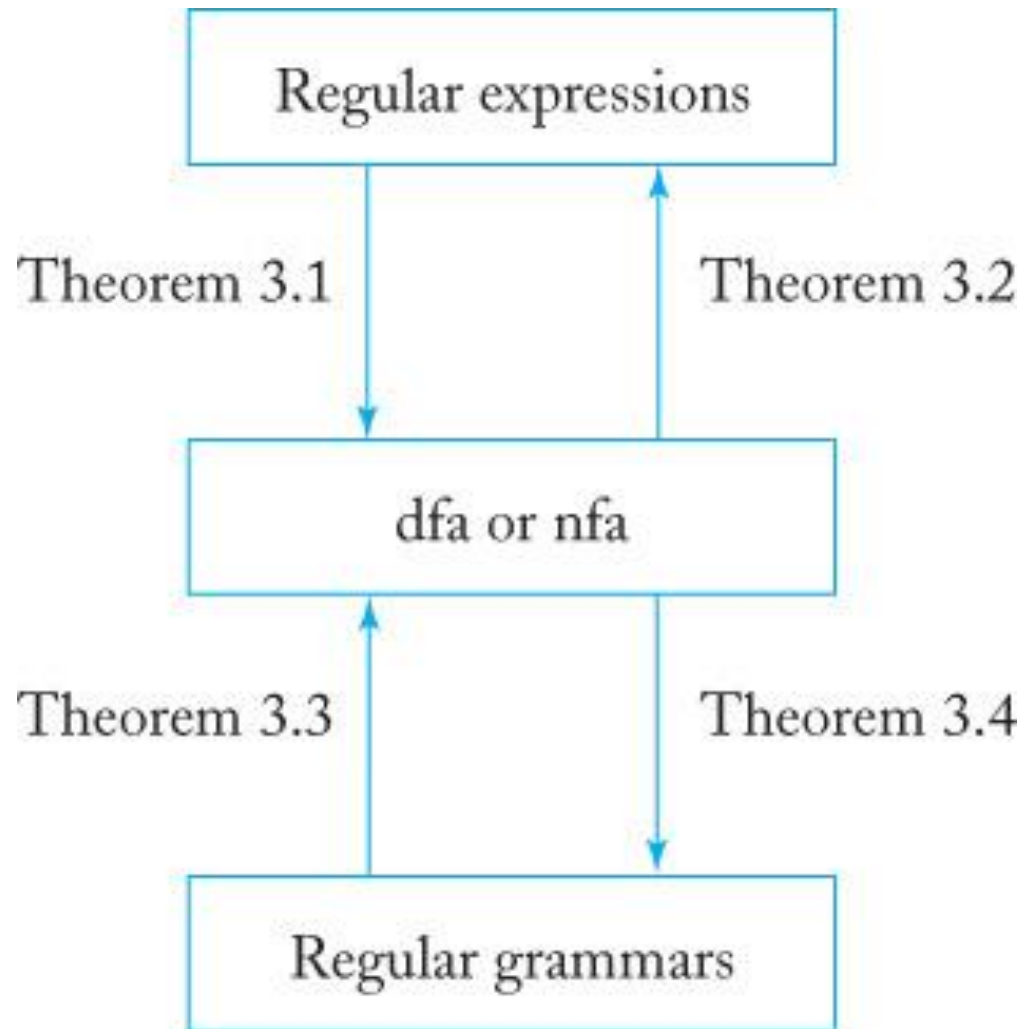
A grammar which generates  $L(aab^*a)$ :

$$G = (V, \Sigma, S, P)$$

where  $V = \{q_0, q_1, q_2, q_f\}$ ,  $\Sigma = \{a, b\}$ ,  $S = q_0$  and  $P$  is as follows:

$$q_0 \rightarrow aq_1, q_1 \rightarrow aq_2, q_2 \rightarrow bq_2, q_2 \rightarrow aq_f, q_f \rightarrow \lambda.$$

# Regular Grammars...



# Reading (Self Study)

- Conversion of Right Linear Grammars into Left Linear Grammars and vice versa.