

# Unit3

# Regular Expression

By Prashant Gautam

# Introduction

- Regular Expression are those algebraic expressions used for representing regular languages, the languages accepted by finite automaton.
- Regular expressions offer a declarative way to express the strings we want to accept. This is what the regular expression offer that the automata do not.
- A regular expression is built up out of simpler regular expression using a set of defining rules. Each regular expression 'r' denotes a language  $L(r)$ .
- The defining rules specify how  $L(r)$  is formed by combining in various ways the languages denoted by the sub-expressions of 'r'.

Let  $\Sigma$  be an alphabet, the regular expression over the alphabet  $\Sigma$  are defined inductively as follows;

- $\Phi$  is a regular expression representing empty language.
- $\epsilon$  is a regular expression representing the language of empty strings.
- if 'a' is a symbol in  $\Sigma$ , then 'a' is a regular expression representing the language {a}.

if 'r' and 's' are the regular expressions representing the language  $L(r)$  and  $L(s)$  then

- $r \cup s$  is a regular expression denoting the language  $L(r) \cup L(s)$ .
- $rs$  is a regular expression denoting the language  $L(r)L(s)$ .
- $r^*$  is a regular expression denoting the language  $(L(r))^*$ .
- $(r)$  is a regular expression denoting the language  $(L(r))$  [defines parenthesis may be placed around regular expressions if we desire].

# Regular operator:

## Union ( $\cup$ / $|$ ):

- If  $L1$  and  $L2$  are any two regular languages then
- $L1 \cup L2 = \{s \mid s \in L1, \text{ or } s \in L2\}$

e.g.

- $L1 = \{00, 11\}$ ,  $L2 = \{\epsilon, 10\}$
- $L1 \cup L2 = \{\epsilon, 00, 11, 10\}$

## Concatenation (.):

- If  $L_1$  and  $L_2$  are any two regular languages then,  $L_1.L_2 = \{l_1.l_2 \mid l_1 \in L_1 \text{ and } l_2 \in L_2\}$
- e.g.  $L_1 = \{00, 11\}$  and  $L_2 = \{\epsilon, 10\}$

## Kleen Closure (\*):

- If L is any regular Language then,

$$L^* = L^i = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \cup \dots$$

## Positive closure (+):

- $L^+ = \bigcup_{i=0}^{\infty} L^i = L^* - \epsilon$



# Precedence of regular operators

- Closure (\*) has highest precedence
- Concatenation (.) has next highest precedence.
- Union (U / | / +) has lowest precedence.

# Regular language

- Let  $\Sigma$  be an alphabet, the class of regular language over  $\Sigma$  is defined inductively as;
  - $\Phi$  is a regular language representing empty language
  - $\{\epsilon\}$  is a regular language representing language of empty strings.
  - For each  $a \in \Sigma$ ,  $\{a\}$  is a regular language.
  - If  $L_1, L_2, \dots, L_n$  is regular languages, then so is  $L_1 \cup L_2 \cup \dots \cup L_n$ .
  - If  $L_1, L_2, L_3, \dots, L_n$  are regular languages, then so is  $L_1 \cdot L_2 \cdot L_3 \cdot \dots \cdot L_n$
  - If  $L$  is a regular language, then so is  $L^*$

# Application of regular languages

- **Validation:** Determining that a string complies with a set of formatting
- **constraints.** Like email address validation, password validation etc.
- **Search and Selection:** Identifying a subset of items from a larger set on the basis of a pattern match.
- **Tokenization:** Converting a sequence of characters into words, tokens (like keywords, identifiers) for later interpretation.

# Algebraic Rules/laws for regular expression

- **Commutativity:** The union of regular expression is commutative but concatenation of regular expression is not commutative. i.e. if  $r$  and  $l$  are regular expressions representing like languages  $L(r)$  and  $L(l)$  then,

$$r + l = l + r \text{ i.e. } r \cup l = l \cup r$$

$$\text{but } r.l \neq l.r$$

- **Associativity:** The unions as well as concatenation of regular expressions are associative.
- i.e. if  $l$ ,  $r$ ,  $s$  are regular expressions representing regular languages  $L(l)$ ,  $L(r)$  and  $L(s)$  then,

$$l+(r+s) = (l+r)+s$$

$$\text{And } l.(r.s) = (l.r).s$$

- **Distributive law:** if  $l, r, s$  are regular expressions representing regular languages  $L(l), L(r)$  and  $L(s)$  then,

$l(m+n) = lm+ln$  ----- left distribution

$(m+n)l = ml+nl$  ----- right distribution.

## Identity law:

- $\Phi$  is identity for union. i.e. for any regular expression  $r$  representing regular expression  $L(r)$ .

$$r + \Phi = \Phi + r = r \text{ i.e. } \Phi \cup r = r$$

$\epsilon$  is identity for concatenation.

$$\text{i.e. } \epsilon.r = r = r.\epsilon$$

**Annihilator:** An annihilator for an operator is a value such that when the operator is applied to the annihilator and some other value, the result is annihilator.

- $\Phi$  is annihilator for concatenation.

$$\text{i.e. } \Phi.r = r.\Phi = \Phi$$



## **Idempotent law of union:**

For any regular expression  $r$  representing the regular language  $L(r)$ ,  $r + r = r$ .

This is the idempotent law of union.

## Law of closure:

for any regular expression  $r$ , representing the regular language  $L(r)$ ,

➤  $(r^*)^* = r^*$

➤ Closure of  $\Phi = \Phi^* = \epsilon$

➤ Closure of  $\epsilon = \epsilon^* = \epsilon$

➤ Positive closure of  $r$ ,  $r^+ = rr^*$ .

## Practice Session

- Given the alphabet  $A = \{0, 1\}$

1. RE denotes the language of all string that begins with a '1' and ends with a '0'. ⑦

$1(1+0)^*0$

2. RE denotes the language of all strings that ends with 00 (binary number multiple of 4) ⑦

$(1+0)^*00$

3. RE denotes the set of all stings that describe alternating 1s and 0s ⑦

$(01)^* + (10)^*$

4. RE denotes the string having exactly three 1's. ⑦

$(0^* 1 0^* 1 0^* 1 0^*)$

5. RE denotes the string having at most two 0's ⑦

$1^*(0 + \epsilon)1^*(0 + \epsilon)1^*$

6. RE denotes the regular expression to specify the identifier like in C ⑦

$(A | B | C | ..... | Z | a | b | c | ..... | z | \_ ) . ((A | B | C | ..... | Z | a | b | c | ..... | z | \_ ) (0 | 1 | 2 | ..... | 9))^*$

7. RE denotes string having substring 0 0 1 ⑦

$(1+0)^* 001 (1+0)^*$

# Finite Automata and Regular Expression:

- The regular expression approach for describing language is fundamentally different from the finite automaton approach. However, these two notations turn out to represent exactly the same set of languages, which we call regular languages.

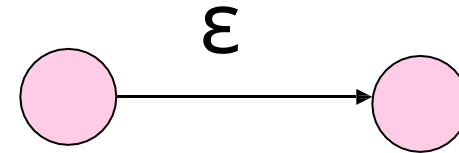
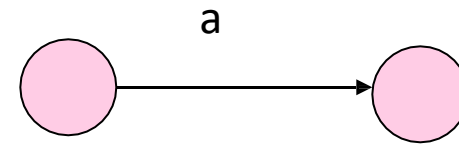
# Regular expression to finite automata

## Theorem: 2

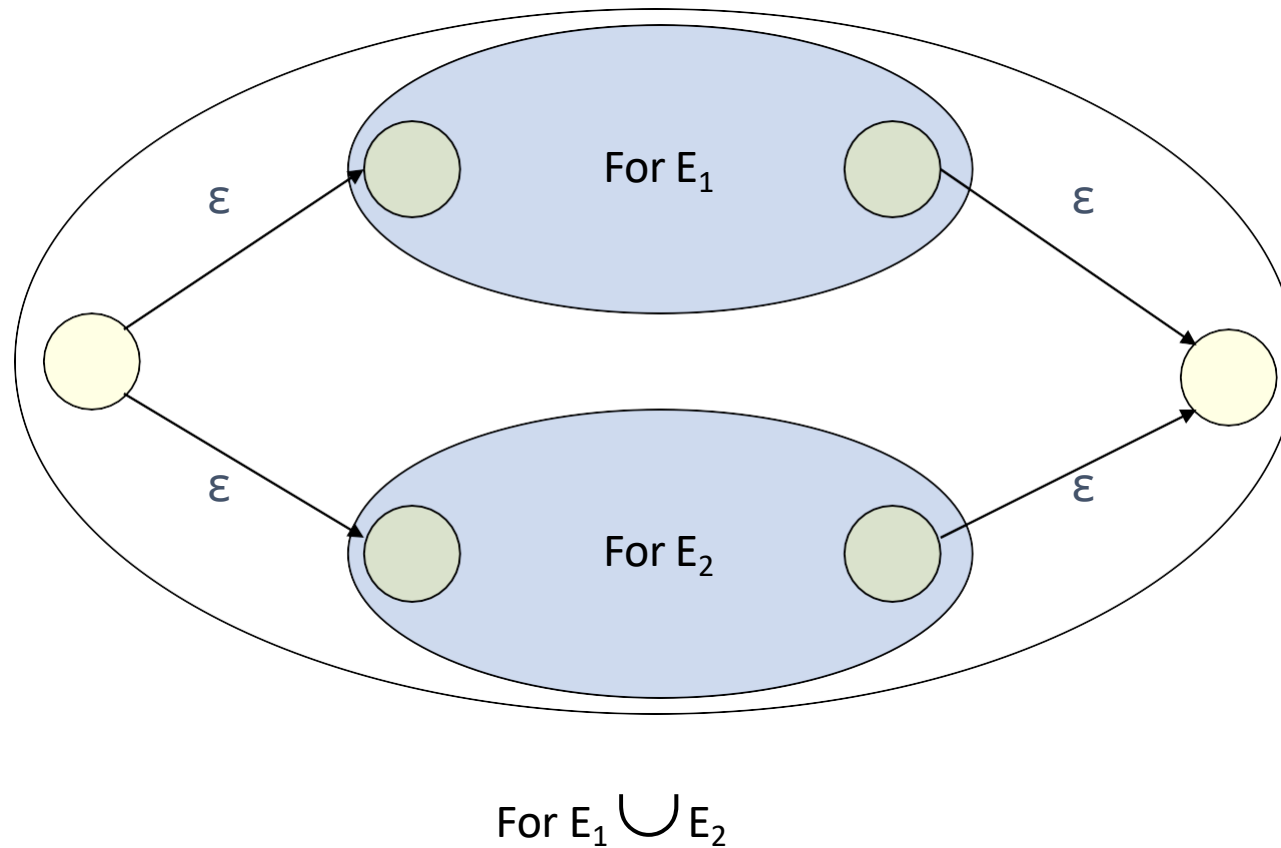
- Every language defined by a regular expression is also defined by a finite automaton.
- [For any regular expression  $r$ , there is an  $\epsilon$ -NFA that accepts the same language represented by  $r$ ].

# RE to $\epsilon$ -NFA: Basis

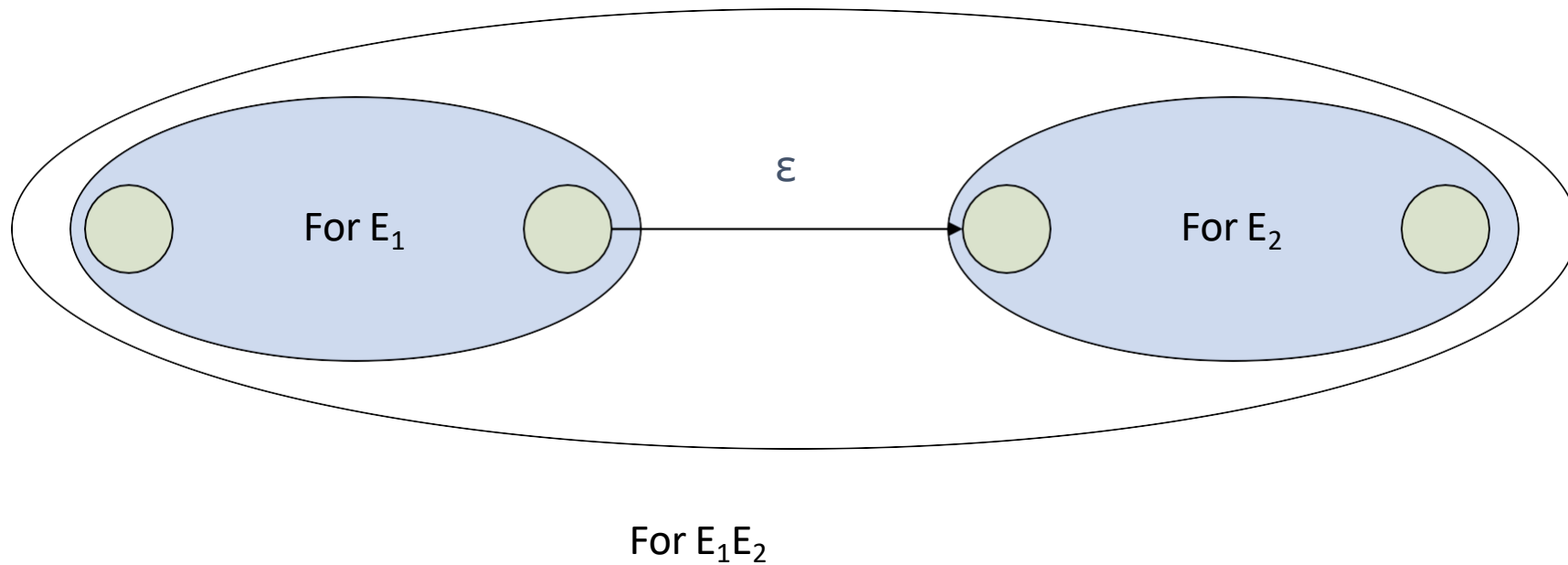
- Symbol  $a$ :
- $\epsilon$ :
- $\emptyset$ :



# RE to $\epsilon$ -NFA: Induction 1 – Union

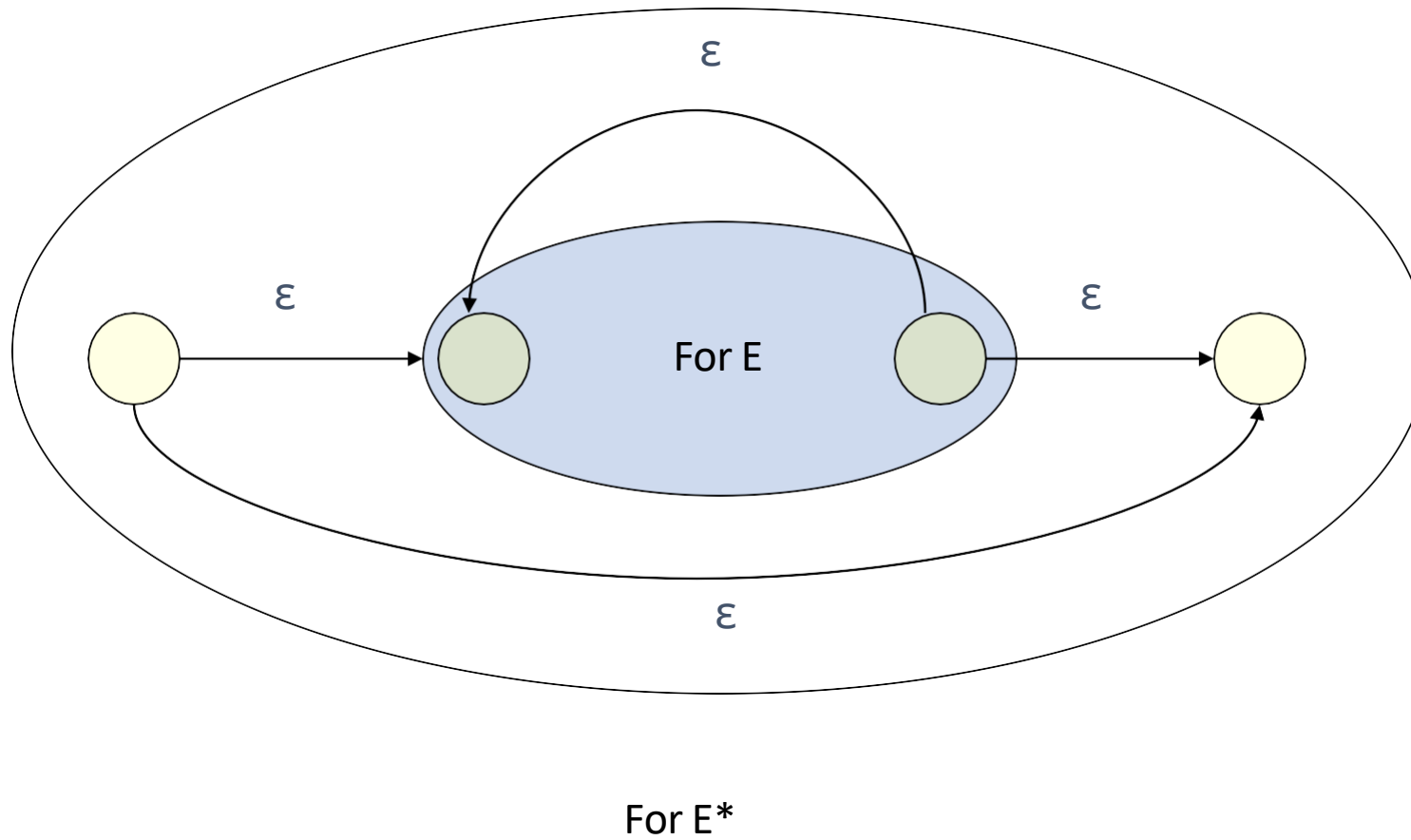


# REto $\epsilon$ -NFA: Induction 2 – Concatenation





# REto $\epsilon$ -NFA: Induction 3 – Closure



# Conversion: Regular Expression to e-NFA

- Thomson's Construction

Thomson's Construction is simple and systematic method. It guarantees that the resulting NFA will have exactly one final state, and one start state.

Method:

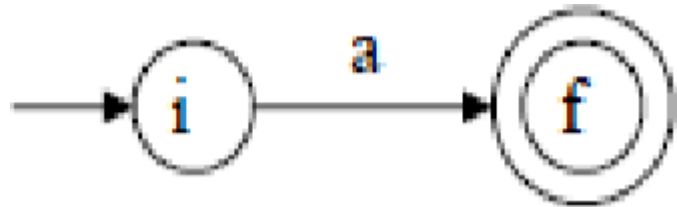
- First parse the regular expression into sub-expressions
- Construct NFA's for each of the basic symbols in regular expression
- Finally combine all NFA's of sub-expressions and we get required NFA of given regular expression.

# Rules

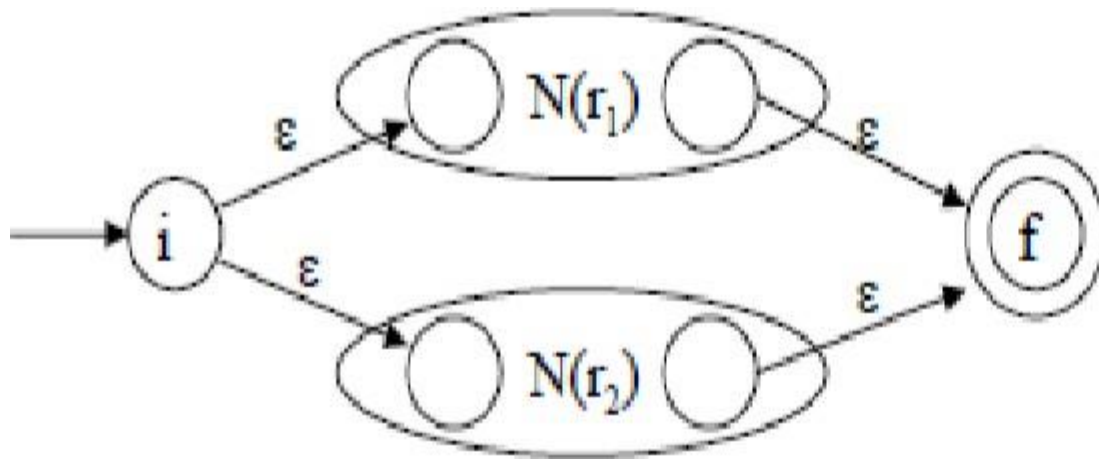
1. To recognize an empty string  $\epsilon$



2. To recognize a symbol  $a$  in the alphabet  $\Sigma$



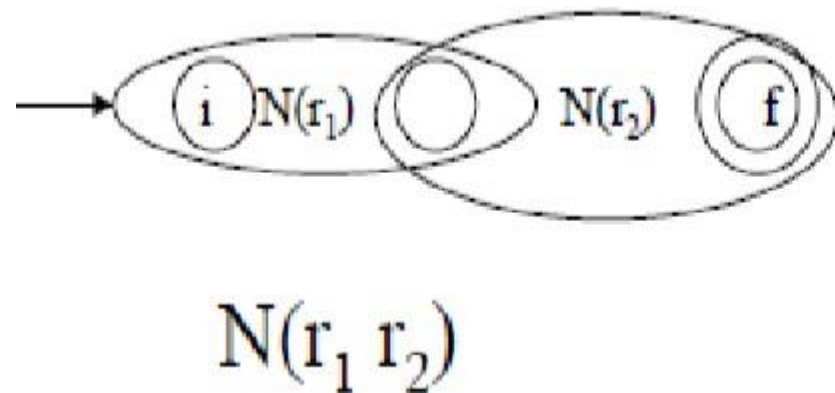
3. If  $N(r_1)$  and  $N(r_2)$  are NFAs for regular expressions  $r_1$  and  $r_2$
- a) For regular expression  $r_1 + r_2$



$N(r_1 + r_2)$

3.

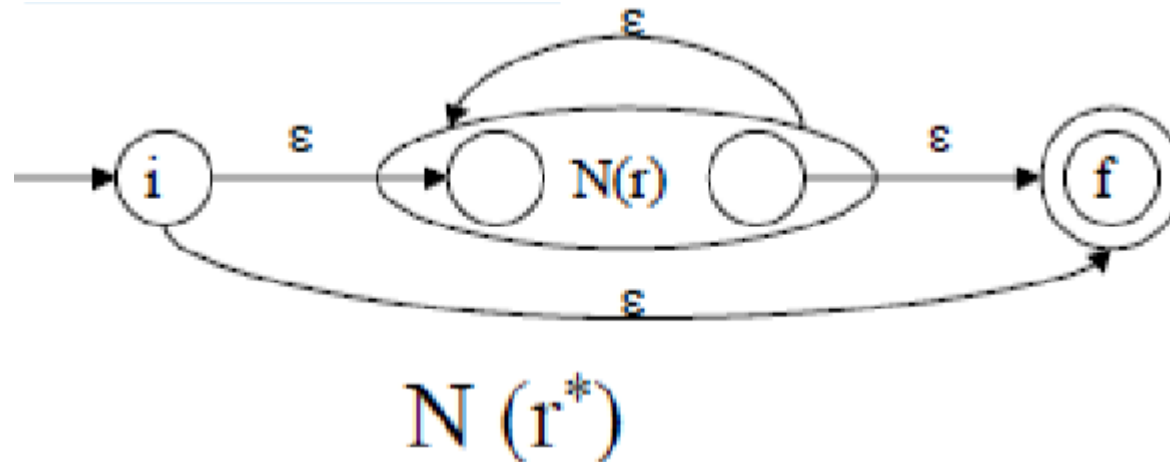
b) For regular expression  $r_1 r_2$



The start state of  $N(r_1)$  becomes the start state of  $N(r_1 r_2)$  and final state of  $N(r_2)$  become final state of  $N(r_1 r_2)$

3.

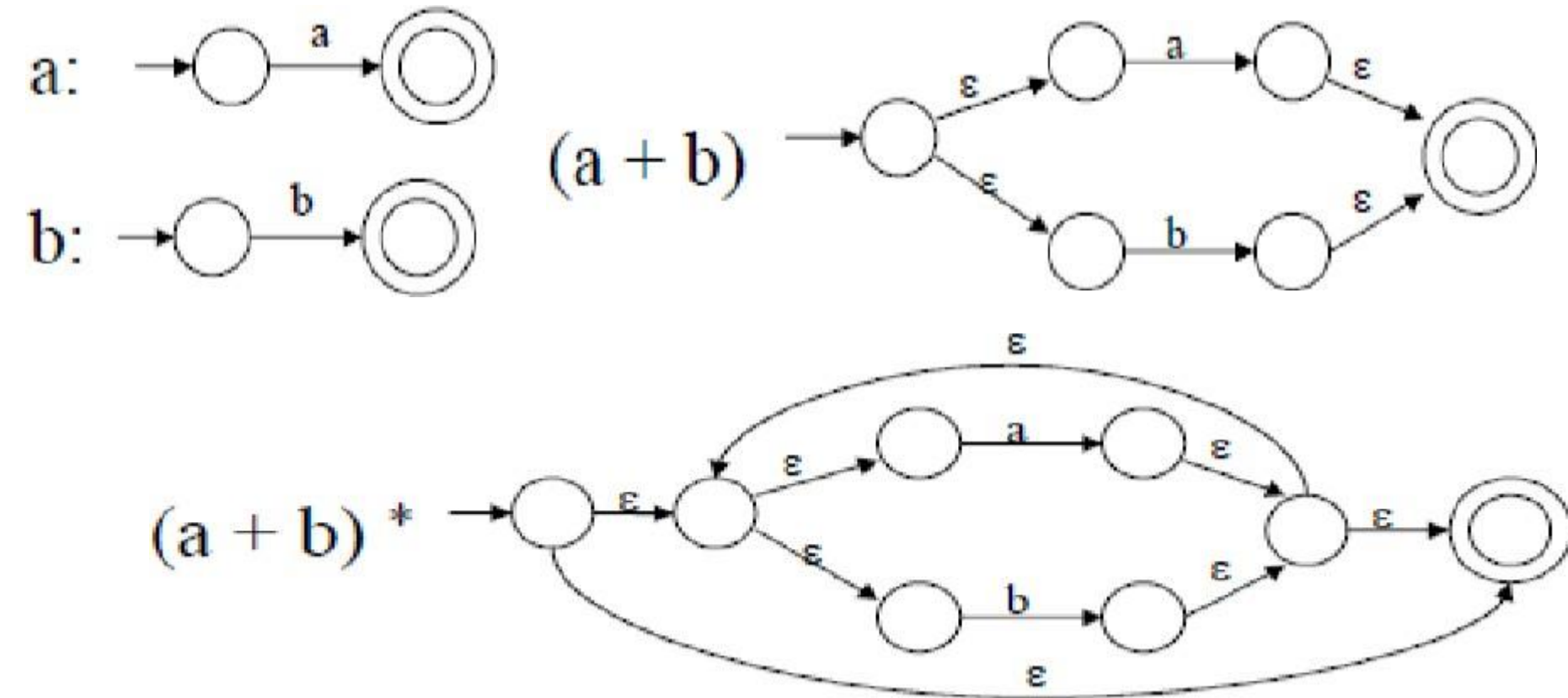
c) For regular expression  $r^*$

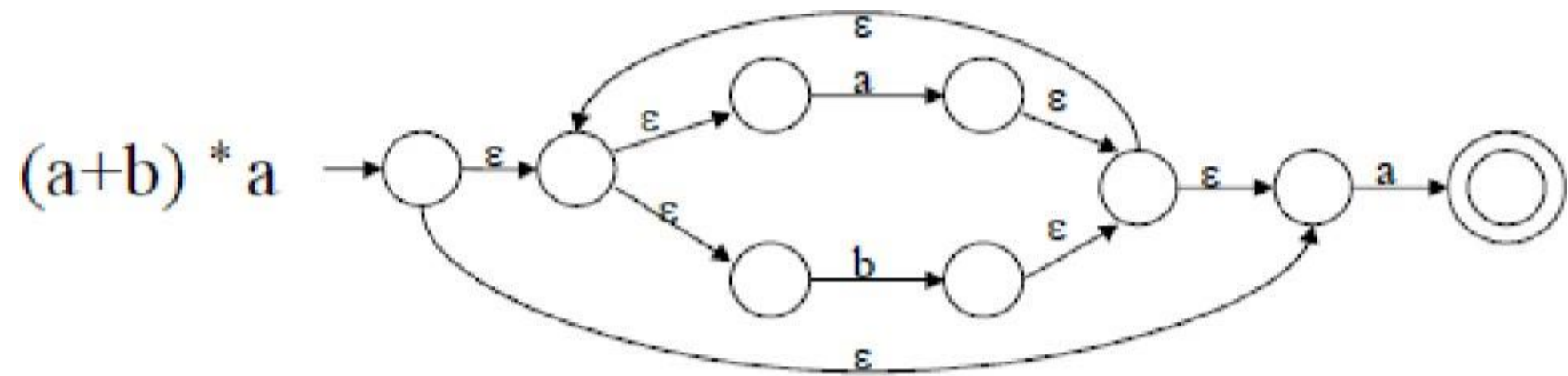


Using rule 1 and 2 we construct NFA's for each basic symbol in the expression, we combine these basic NFA using rule 3 to obtain the NFA for entire expression.

# Example

- NFA construction of RE  $(a + b)^* a$







# TRY Yourself

$(ab + a)^*$

# Assignment: RE to e-NFA

- $(0+1)^*$
- $(0+1)^*1(0+1)$
- $(00+1)^*10$
- $(00+1)^*10$
- $(00+1)^*10$
- $(00+1)^*10$

# DFA to Regular Expression: Arden's Theorem

- Let's assume that **P** and **Q** be two regular expressions.
- In case if **P** does not include a  $\epsilon$  (null string), then

**R** = **Q** + **RP** has a unique solution that is

$$\mathbf{R} = \mathbf{QP}^*$$

# Proof

- $R = Q + (Q + RP)P$

[After putting the value  $R = Q + RP$ ]

- $R = Q + QP + RPP$

- If you include the value of R recursively repetitively then you will get the following equation –

$$R = Q + QP + QP^2 + QP^3 \dots$$

$$R = Q (\epsilon + P + P^2 + P^3 + \dots)$$

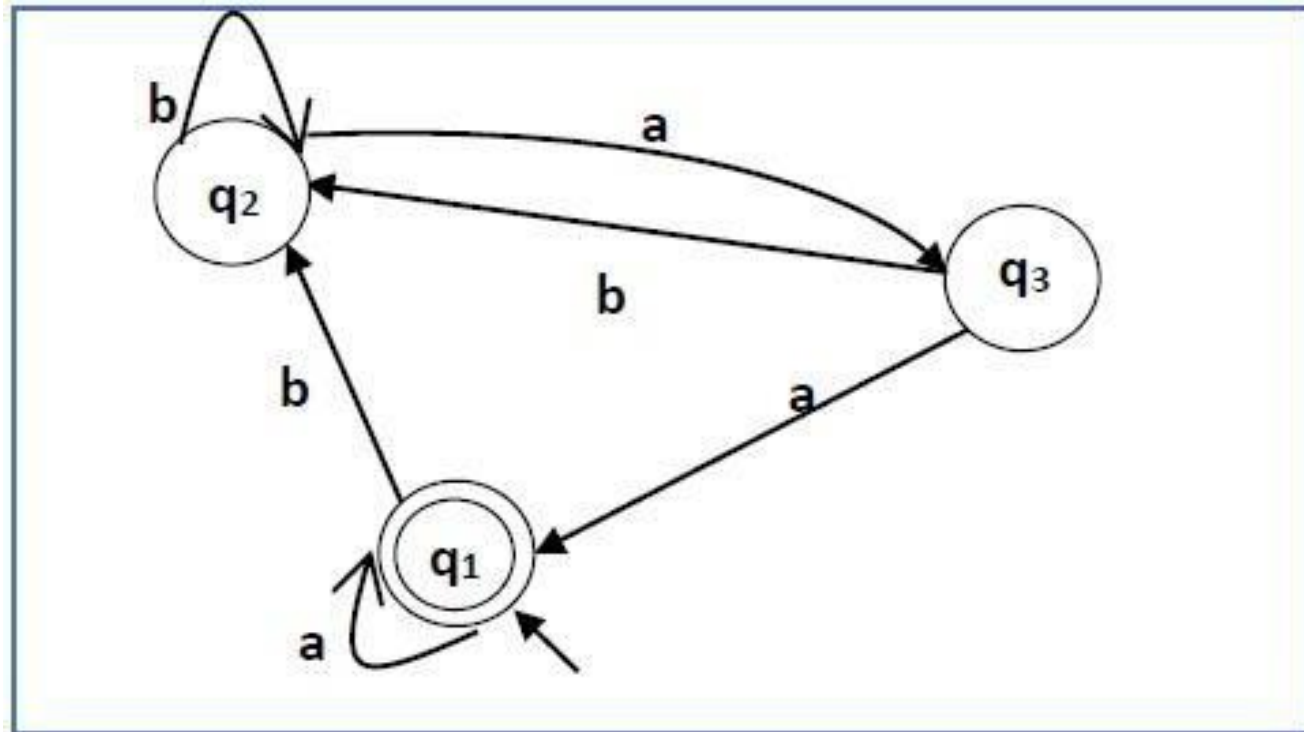
$$R = QP^* \quad [\text{As } P^* \text{ represents } (\epsilon + P + P^2 + P^3 + \dots)]$$

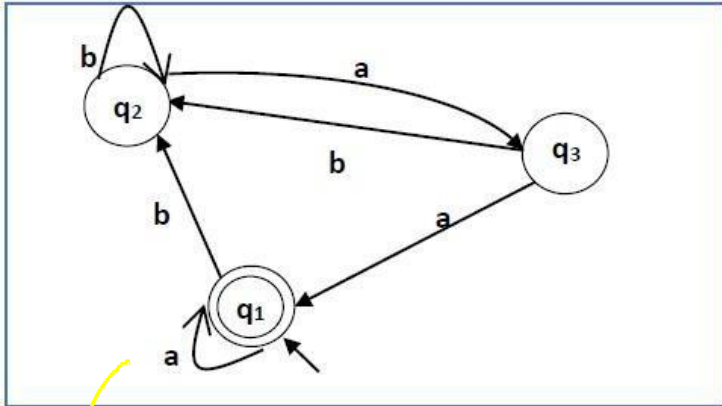
# Assumptions for Applying Arden's Theorem

- The transition diagram must not have  $\epsilon$  transitions

## Problem

Construct a regular expression related to the automata given below.





$$R = Q + RP$$

$$R = QP^*$$

solving.

$$r_2 = r_1 b + r_2 b + r_3 b$$

$$r_2 = r_1 b + r_1 b + r_2 a b \quad [\text{from (3)}]$$

$$\frac{r_2}{R} = \frac{r_1 b}{Q} + \frac{r_2}{R} \frac{(b + ab)}{P}$$

$$r_2 = r_1 b (b + ab)^* \quad \text{--- (4)}$$

Now substituting value of  $r_2$  &  $r_3$  in eq (1)

$$r_1 = r_1 a + r_2 a a + \epsilon$$

$$r_1 = r_1 a + r_1 b (b + ab)^* \cdot \underline{aa} + (\epsilon)$$

$$\xrightarrow{R} r_1 = \epsilon + \underline{r_1} \left[ \underline{a + b (b + ab)^* \cdot aa} \right]$$

$$r_1 = \epsilon \cdot (a + b (b + ab)^* \cdot aa)^*$$

$$\Rightarrow (a + b (b + ab)^* \cdot aa)^*$$

# Solution

- Here the initial state is **q2** and the final state is **q2**.
- The equations for the three states q1, q2, and q3 are as follows –

$$q_1 = q_1a + q_3a + \varepsilon \quad (\varepsilon \text{ move is because } q_1 \text{ is the initial state.})$$

$$q_2 = q_1b + q_2b + q_3b$$

$$q_3 = q_2a$$



Now, we will solve these three equations –

$$q_2 = q_1b + q_2b + q_3b$$

$$= q_1b + q_2b + (q_2a)b \quad (\text{Substituting value of } q_3)$$

$$= q_1b + q_2(b + ab)$$

$$= q_1b (b + ab)^* \quad (\text{Applying Arden's Theorem})$$

$$q_1 = q_1a + q_3a + \varepsilon$$

$$= q_1a + q_2aa + \varepsilon \quad (\text{Substituting value of } q_3)$$

$$= q_1a + q_1b(b + ab^*)aa + \varepsilon \quad (\text{Substituting value of } q_2)$$

$$= q_1(a + b(b + ab)^*aa) + \varepsilon$$

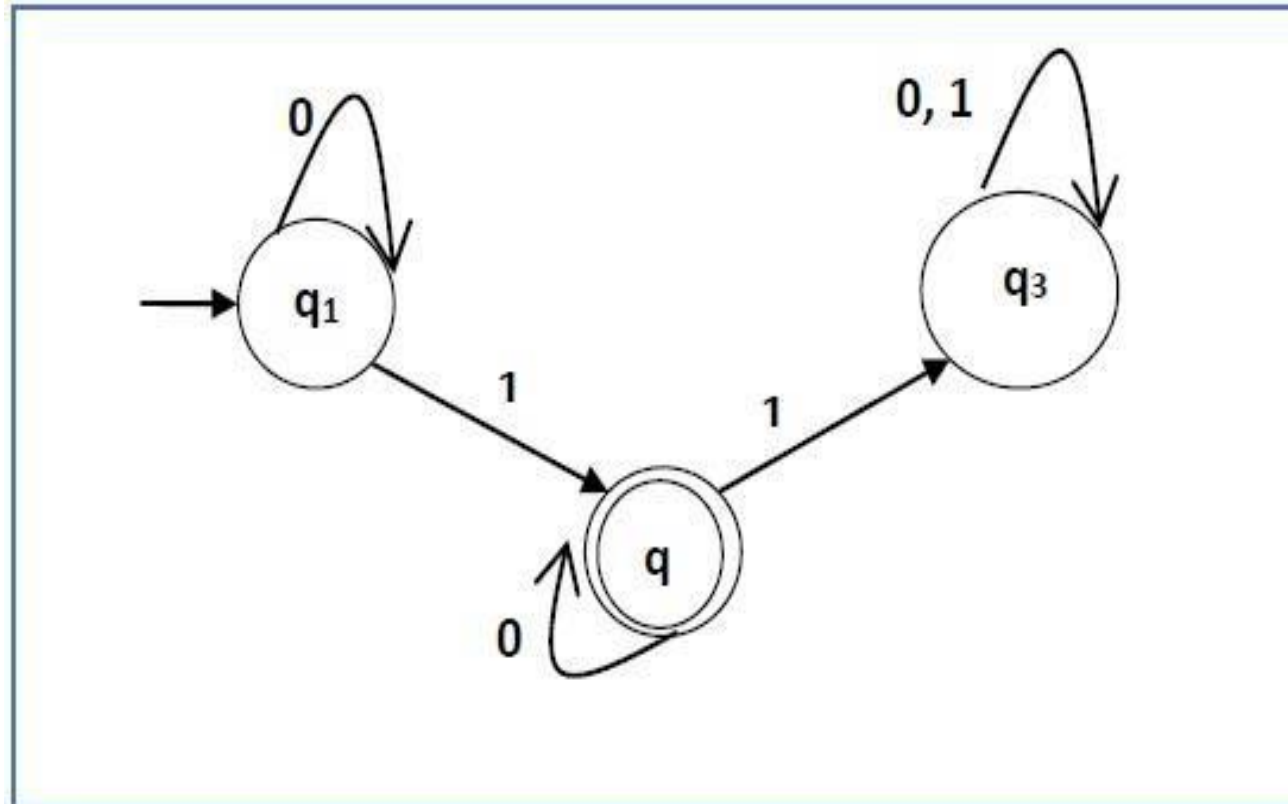
$$= \varepsilon (a + b(b + ab)^*aa)^*$$

$$= (a + b(b + ab)^*aa)^*$$

Hence, the regular expression is  $(a + b(b + ab)^*aa)^*$ .

## Problem 2 : Assignment

construct a regular expression related to the automata given below –



# Solution

- Here the initial state is  $q_1$  and the final state is  $q_2$
- Let's write the below equations –

$$q_1 = q_1 0 + \varepsilon$$

$$q_2 = q_1 1 + q_2 0$$

$$q_3 = q_2 1 + q_3 0 + q_3 1$$

Now, we will solve these three equations –

$$q_1 = \varepsilon 0^* [As, \varepsilon R = R]$$

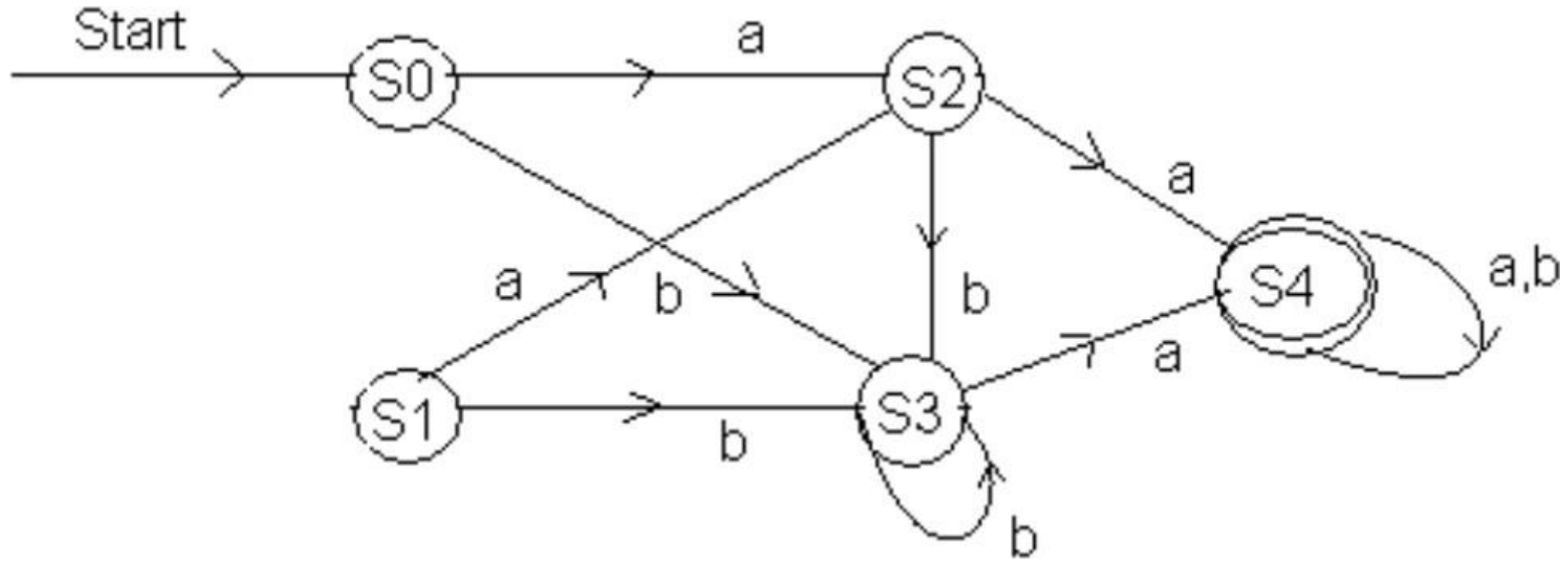
So,  $q_1 = 0^*$

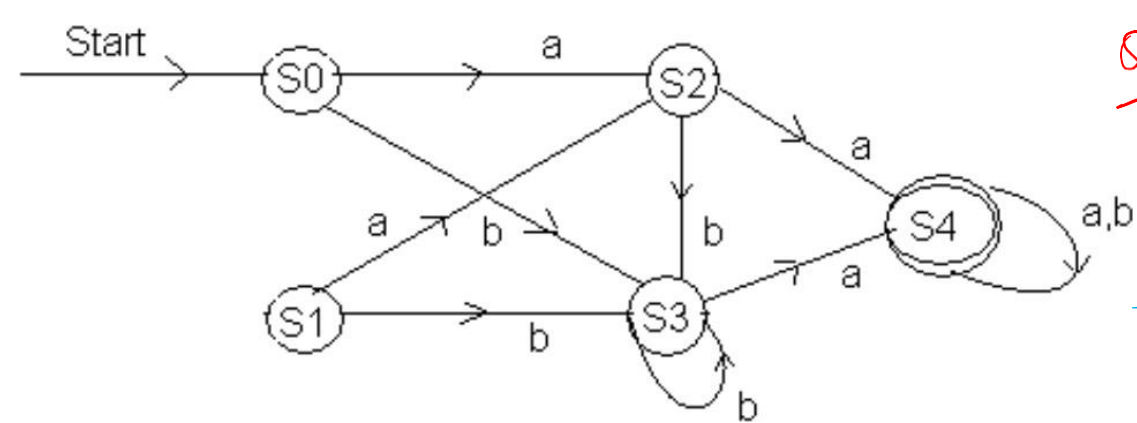
$$q_2 = 0^*1 + q_2 0$$

So,  $q_2 = 0^*1(0)^*$  [By Arden's theorem]

Hence, the regular expression is  **$0^*10^*$**

# 1. DFA State Minimization





Q. Minimize

(reduce the no. of states as far as possible.)

Sol<sup>n</sup>:

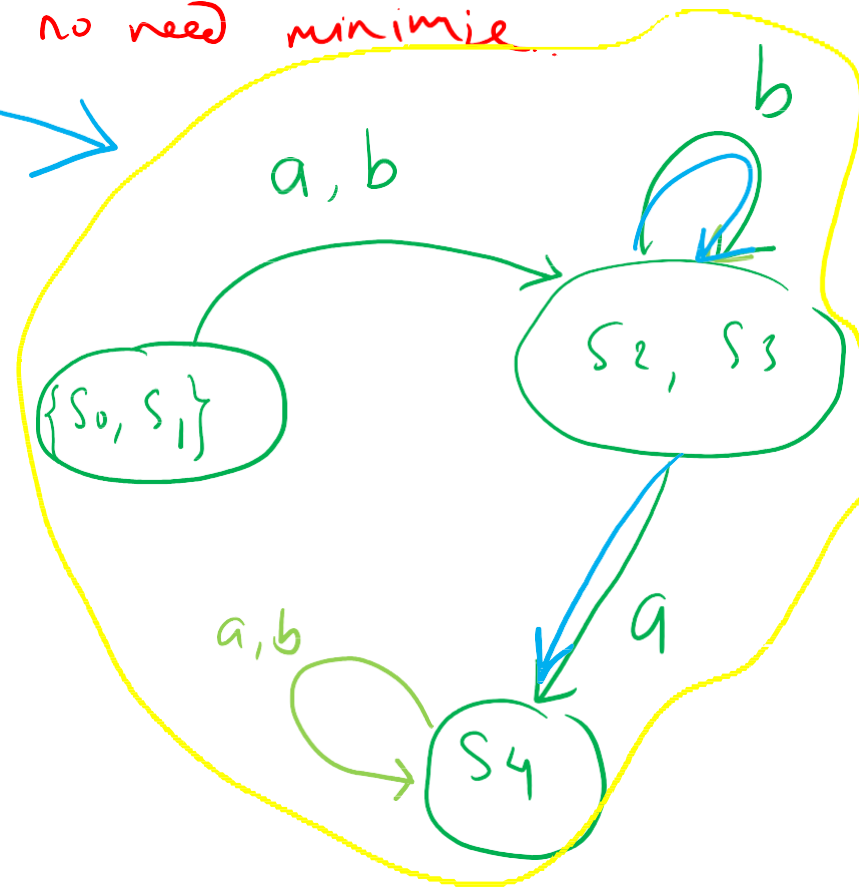
Acc. state.

	a	b
S <sub>0</sub>	S <sub>2</sub>	S <sub>3</sub>
S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>
S <sub>2</sub>	S <sub>4</sub>	S <sub>3</sub>
S <sub>3</sub>	S <sub>4</sub>	S <sub>3</sub>

	a	b
S <sub>0</sub> , S <sub>1</sub>	S <sub>2,3</sub>	S <sub>3,4</sub>
S <sub>2</sub> , S <sub>3</sub>	S <sub>4</sub>	S <sub>3,4</sub>

	a	b
S <sub>4</sub>	S <sub>4</sub>	S <sub>4</sub>

1) only one Acc. state so  
no need minimize



## 2. DFA State Minimization (Assignment)

