

Unit – 2 & 3

# **Regular Languages and Finite Automata**

**&**

# **Lexical Analysis**

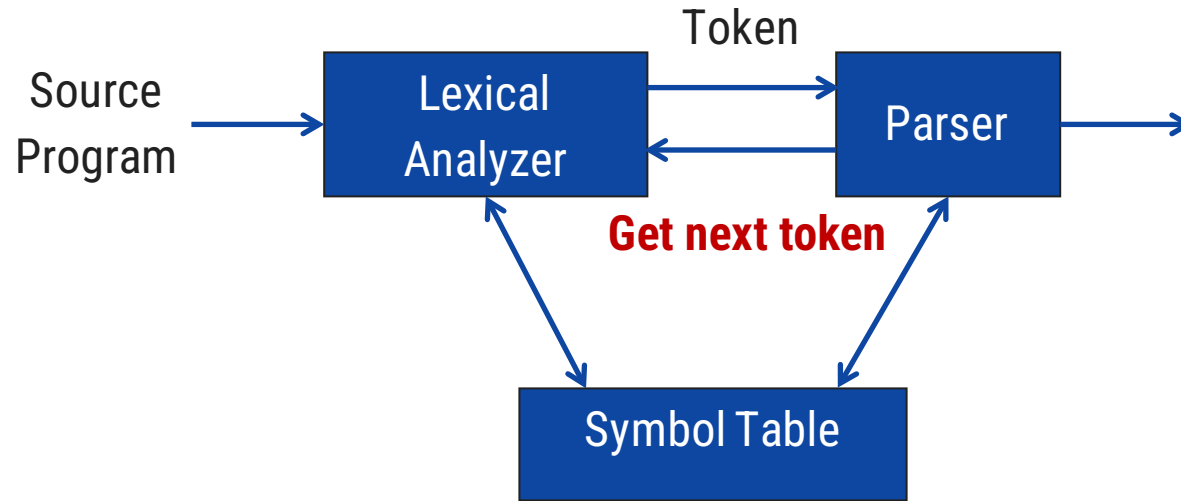
# Topics to be covered



- Interaction of scanner & parser
- Token, Pattern & Lexemes
- Input buffering
- Specification of tokens
- Regular expression & Regular definition
- Transition diagram
- Hard coding & automatic generation lexical analyzers
- Finite automata
- Regular expression to NFA using Thompson's rule
- Conversion from NFA to DFA using subset construction method
- DFA optimization
- Conversion from regular expression to DFA
- An Elementary Scanner Design & It's Implementation

# **Interaction with Scanner & Parser**

# Interaction of scanner & parser



- ▶ Upon receiving a **"Get next token"** command from parser, the lexical analyzer reads the input character until it can identify the next token.
- ▶ Lexical analyzer also stripping out comments and white space in the form of blanks, tabs, and newline characters from the source program.

# Why to separate lexical analysis & parsing?

1. Simplicity in **design**.
2. Improves compiler **efficiency**.
3. Enhance compiler **portability**.

# Token, Pattern & Lexemes

# Token, Pattern & Lexemes

## Token

A **token** is a pair consisting of a token name and an optional attribute value.

The token name is an abstract symbol representing a kind of lexical unit

Categories of Tokens:

1. Identifier
2. Keyword
3. Operator
4. Special symbol
5. Constant


## Pattern

A **pattern** is a description of the form that the lexemes of a token may take.

Example: “*non-empty sequence of digits*”, “*letter followed by letters and digits*”

## Lexemes

The **sequence of character** in a source program **matched with a pattern** for a **token** is called lexeme.

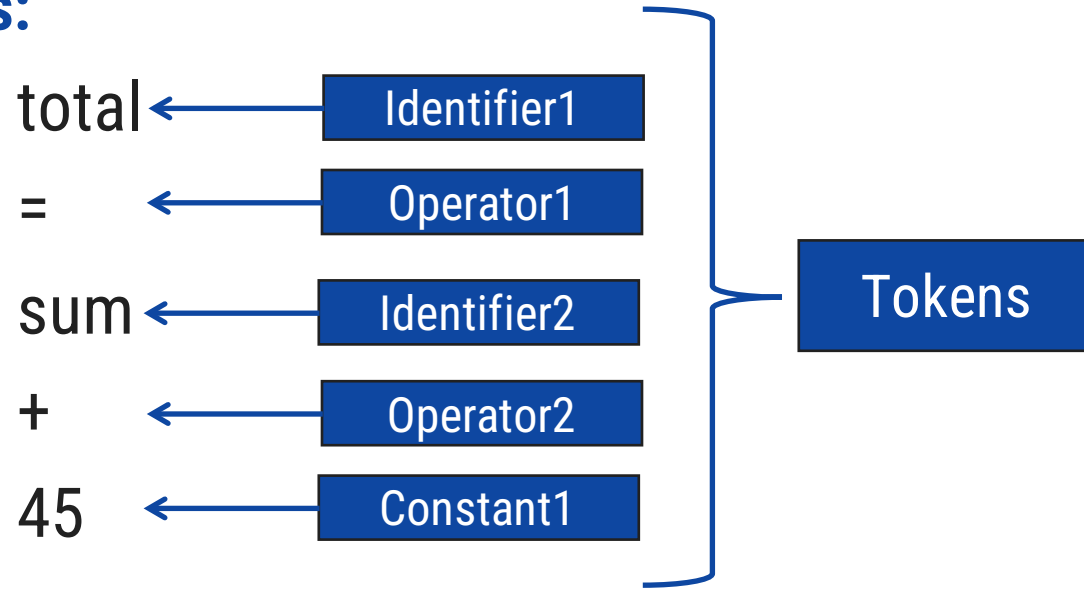


Example: Rate, DIET, count, Flag

# Example: Token, Pattern & Lexemes

Example: total = sum + 45

## Tokens:



## Lexemes

Lexemes of identifier: total, sum

Lexemes of operator: =, +

Lexemes of constant: 45



# Example: Token, Pattern & Lexemes

TOKEN	INFORMAL DESCRIPTION	SAMPLE LEXEMES
<b>if</b>	characters i, f	if
<b>else</b>	characters e, l, s, e	else
<b>comparison</b>	< or > or <= or >= or == or !=	<=, !=
<b>id</b>	letter followed by letters and digits	pi, score, D2
<b>number</b>	any numeric constant	3.14159, 0, 6.02e23
<b>literal</b>	anything but ", surrounded by "'s	"core dumped"

# Input buffering

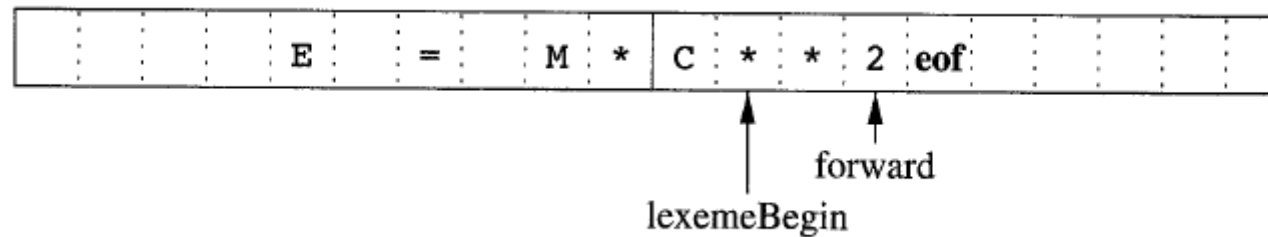
# Input buffering

- ▶ There are mainly two techniques for input buffering:

1. Buffer pairs
2. Sentinels

## Buffer Pair

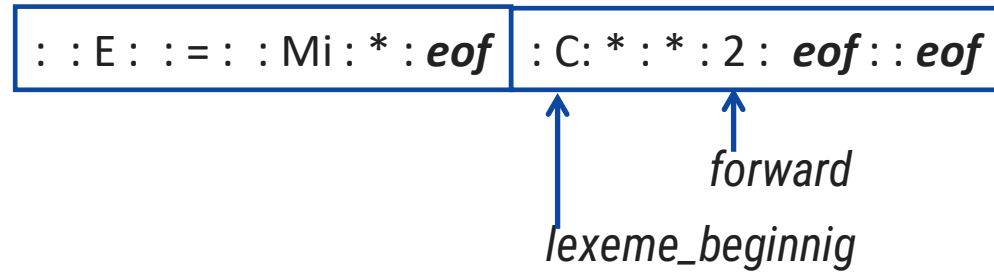
- ▶ The lexical analysis scans the input string from left to right one character at a time.
- ▶ Each buffer is of the same size N, and N is usually the size of a disk block, e.g., 4096 bytes.
- ▶ Using one system read command we can read N characters into a buffer, rather than using one system call per character. If fewer than N characters remain in the input file, then a special character, represented by eof, marks the end of the source file.



# Buffer pairs

- ▶ Two pointers to the input are maintained:
  - ↳ **Pointer lexemeBegin**, marks the beginning of the current lexeme, whose extent we are attempting to determine.
  - ↳ **Pointer forward** scans ahead until a pattern match is found

# Sentinels



- ▶ In buffer pairs we must check, each time we move the forward pointer that we have not moved off one of the buffers.
- ▶ Thus, for each character read, we make two tests.
- ▶ We can combine the buffer-end test with the test for the current character.
- ▶ We can reduce the two tests to one if we extend each buffer to hold a sentinel character at the end.
- ▶ The sentinel is a special character that cannot be part of the source program, and a natural choice is the character **EOF**.

# Specification of tokens

# Strings and Languages

- ▶ Regular expressions are an important notation for specifying lexeme patterns.
- ▶ An **alphabet** is any finite set of symbols.
- ▶ Typical examples of symbols are letters, digits, and punctuation. The set  $\{0,1\}$  is the binary alphabet
- ▶ A **string** over an alphabet is a finite sequence of symbols drawn from that alphabet.
- ▶ The length of a string  $s$ , usually written  $|s|$ , is the number of occurrences of symbols in  $s$ .
- ▶ For example, banana is a string of length six.
- ▶ The empty string, denoted  $\epsilon$ , is the string of length zero.
- ▶ A **language** is any countable set of strings over some fixed alphabet.
- ▶ Abstract languages like  $\emptyset$ , the empty set, or  $\{\epsilon\}$ , the set containing only the empty string, are languages under this definition.

# Strings and languages

Term	Definition
<i>Prefix of s</i>	A string obtained by removing <b>zero or more trailing symbol</b> of string S. e.g., <b>ban</b> is prefix of <b>banana</b> .
<i>Suffix of S</i>	A string obtained by removing <b>zero or more leading symbol</b> of string S. e.g., <b>nana</b> is suffix of <b>banana</b> .
<i>Sub string of S</i>	A string obtained by <b>removing prefix and suffix</b> from S. e.g., <b>nan</b> is substring of <b>banana</b>
<i>Proper prefix, suffix and substring of S</i>	Any nonempty string x that is respectively proper prefix, suffix or substring of S, such that <b><math>s \neq x</math></b> .
<i>Subsequence of S</i>	A string obtained by removing <b>zero or more not necessarily contiguous symbol</b> from S. e.g., <b>baaa</b> is subsequence of <b>banana</b> .



# Operations on languages

Operation	Definition
Union of L and M Written L U M	$L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$
Concatenation of L and M Written LM	$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$
Kleene closure of L Written L*	$L^*$ denotes “zero or more concatenation of” L.
Positive closure of L Written L <sup>+</sup>	$L^+$ denotes “one or more concatenation of” L.

# Regular Expression & Regular Definition

# Regular expression

- ▶ A regular expression is a sequence of characters that define a pattern.

## Notational shorthand's

1. One or more instances: +
2. Zero or more instances: \*
3. Zero or one instances: ?
4. Alphabets:  $\Sigma$

# Rules to define regular expression

1.  $\epsilon$  is a regular expression that denotes  $\{\epsilon\}$ , the set containing empty string.
2. If  $a$  is a symbol in  $\Sigma$  then  $a$  is a regular expression,  $L(a) = \{a\}$
3. Suppose  $r$  and  $s$  are regular expression denoting the languages  $L(r)$  and  $L(s)$ . Then,
  - a.*  $(r)|(s)$  is a regular expression denoting  $L(r) \cup L(s)$
  - b.*  $(r)(s)$  is a regular expression denoting  $L(r)L(s)$
  - c.*  $(r)^*$  is a regular expression denoting  $(L(r))^*$
  - d.*  $(r)$  is a regular expression denoting  $L((r))$

The language denoted by regular expression is said to be a **regular set**.

# Regular expression

► **L = Zero or More Occurrences of a =  $a^*$**



$\epsilon$   
a  
aa  
aaa  
aaaa  
aaaaa.....

Infinite .....

# Regular expression

► **L = One or More Occurrences of a =  $a^+$**



a  
aa  
aaa  
aaaa  
aaaaa.....

**Infinite .....**

# Precedence and associativity of operators

Operator	Precedence	Associative
Kleene *	1	left
Concatenation	2	left
Union	3	left

# Regular expression examples

1. 0 or 1

Strings: 0, 1

R. E. =  $0 \mid 1$

2. 0 or 11 or 111

Strings: 0, 11, 111

R. E. =  $0 \mid 11 \mid 111$

3. String having zero or more  $a$ .

Strings:  $\epsilon$ ,  $a$ ,  $aa$ ,  $aaa$ ,  $aaaa$  ....

R. E. =  $a^*$

4. String having one or more  $a$ .

Strings:  $a$ ,  $aa$ ,  $aaa$ ,  $aaaa$  ....

R. E. =  $a^+$

5. Regular expression over  $\Sigma = \{a, b, c\}$  that represent all string of length 3.

Strings:  $abc$ ,  $bca$ ,  $bbb$ ,  $cab$ ,  $aba$  ....

R. E. =  $(a|b|c)(a|b|c)(a|b|c)$

6. All binary string

Strings: 0, 11, 101, 10101, 1111 ...

R. E. =  $(0 \mid 1)^+$



# Regular expression examples

7. 0 or more occurrence of either a or b or both

*Strings:*  $\epsilon$ , *a*, *aa*, *abab*, *bab* ...

*R.E.* =  $(a \mid b)^*$

8. 1 or more occurrence of either a or b or both

*Strings:* *a*, *aa*, *abab*, *bab*, *bbbaaa* ...

*R.E.* =  $(a \mid b)^+$

9. Binary no. ends with 0

*Strings:* *0*, *10*, *100*, *1010*, *11110* ...

*R.E.* =  $(0 \mid 1)^* 0$

10. Binary no. ends with 1

*Strings:* *1*, *101*, *1001*, *10101*, ...

*R.E.* =  $(0 \mid 1)^* 1$

11. Binary no. starts and ends with 1

*Strings:* *11*, *101*, *1001*, *10101*, ...

*R.E.* =  $1(0 \mid 1)^* 1$

12. String starts and ends with same character

*Strings:* *00*, *101*, *aba*, *baab* ...

*R.E.* =  $1(0 \mid 1)^* 1 \mid 0(0 \mid 1)^* 0$   
 $a(a \mid b)^* a \mid b(a \mid b)^* b$

# Regular expression examples

13. All string of a and b starting with a

*Strings: a, ab, aab, abb...*

*R.E. =  $a(a \mid b)^*$*

14. String of 0 and 1 ends with 00

*Strings: 00, 100, 000, 1000, 1100...*

*R.E. =  $(0 \mid 1)^* 00$*

15. String ends with abb

*Strings: abb, babb, ababb...*

*R.E. =  $(a \mid b)^* abb$*

16. String starts with 1 and ends with 0

*Strings: 10, 100, 110, 1000, 1100...*

*R.E. =  $1(0 \mid 1)^* 0$*

17. All binary string with at least 3 characters and 3<sup>rd</sup> character should be zero

*Strings: 000, 100, 1100, 1001...*

*R.E. =  $(0 \mid 1)(0 \mid 1)0(0 \mid 1)^*$*

18. Language which consist of exactly two b's over the set  $\Sigma = \{a, b\}$

*Strings: bb, bab, aabb, abba...*

*R.E. =  $a^* b a^* b a^*$*

# Regular expression examples

19. The language with  $\Sigma = \{a, b\}$  such that 3<sup>rd</sup> character from right end of the string is always a.

*Strings: aaa, aba, aaba, abb...*

*R.E. =  $(a | b)^* a(a|b)(a|b)$*

20. Any no. of  $a$  followed by any no. of  $b$  followed by any no. of  $c$

*Strings:  $\epsilon$ , abc, aabbcc, aabc, abb...*

*R.E. =  $a^* b^* c^*$*

21. String should contain at least three 1

*Strings: 111, 01101, 0101110...*

*R.E. =  $(0|1)^* 1 (0|1)^* 1 (0|1)^* 1 (0|1)^*$*

22. String should contain exactly two 1

*Strings: 11, 0101, 1100, 010010, 100100...*

*R.E. =  $0^* 1 0^* 1 0^*$*

23. Length of string should be at least 1 and at most 3

*Strings: 0, 1, 11, 01, 111, 010, 100...*

*R.E. =  $(0|1) | (0|1)(0|1) | (0|1)(0|1)(0|1)$*

24. No. of zero should be multiple of 3

*Strings: 000, 010101, 110100, 000000, 100010010...*

*R.E. =  $(1^* 01^* 01^* 01^*)^*$*

# Regular expression examples

25. The language with  $\Sigma = \{a, b, c\}$  where  $a$  should be multiple of 3

*Strings: aaa, baaa, bacaba, aaaaaa. .* **R.E. =  $((b|c)^*a(b|c)^*a(b|c)^*a(b|c)^*)^*$**

26. Even no. of 0

*Strings: 00, 0101, 0000, 100100....* **R.E. =  $(1^*01^*01^*)^*$**

27. String should have odd length

*Strings: 0, 010, 110, 000, 10010....* **R.E. =  $(0|1)((0|1)(0|1))^*$**

28. String should have even length

*Strings: 00, 0101, 0000, 100100....* **R.E. =  $((0|1)(0|1))^*$**

29. String start with 0 and has odd length

*Strings: 0, 010, 010, 000, 00010....* **R.E. =  $(0)((0|1)(0|1))^*$**

30. String start with 1 and has even length

*Strings: 10, 1100, 1000, 100100....* **R.E. =  $1(0|1)((0|1)(0|1))^*$**

31. All string begins or ends with 00 or 11

*Strings: 00101, 10100, 110, 01011 ...* **R.E. =  $(00|11)(0|1)^*|(0|1)^*(00|11)$**

# Regular expression examples

32. Language of all string containing both 11 and 00 as substring

*Strings:* 0011, 1100, 100110, 010011 ... *R.E.* =  $((0|1)^*00(0|1)^*11(0|1)^*) \mid ((0|1)^*11(0|1)^*00(0|1)^*)$

33. String ending with 1 and not contain 00

*Strings:* 011, 1101, 1011 .... *R.E.* =  $(1|01)^+$

34. Language of identifier

*Strings:* area, i, redious, grade1 .... *R.E.* =  $(\_ + L)(\_ + L + D)^*$

*where L is Letter & D is digit*

# Regular definition

- ▶ A regular definition gives names to certain regular expressions and uses those names in other regular expressions.

- ▶ Regular definition is a sequence of definitions of the form:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

.....

$$d_n \rightarrow r_n$$

Where  $d_i$  is a **distinct name** &  $r_i$  is a **regular expression**.

- Example: Regular definition for identifier

**letter**  $\rightarrow$  A|B|C|.....|Z|a|b|.....|z

**digit**  $\rightarrow$  0|1|.....|9|

**id**  $\rightarrow$  **letter** (**letter** | **digit**)\*

# Regular definition example

## ► Example: Unsigned Pascal numbers

3

5280

39.37

6.336E4

1.894E-4

2.56E+7

## Regular Definition

digit  $\rightarrow 0|1|...|9$

digits  $\rightarrow \text{digit digit}^*$

optional\_fraction  $\rightarrow \text{.digits} | \epsilon$

optional\_exponent  $\rightarrow (\text{E}(+|-|\epsilon)\text{digits})|\epsilon$

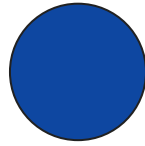
num  $\rightarrow \text{digits optional\_fraction optional\_exponent}$

# Transition Diagram



# Transition Diagram

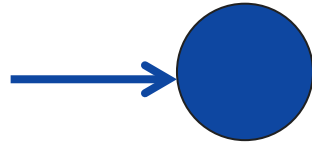
- ▶ A stylized flowchart is called transition diagram.



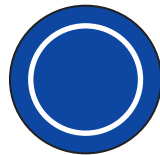
is a state



is a transition

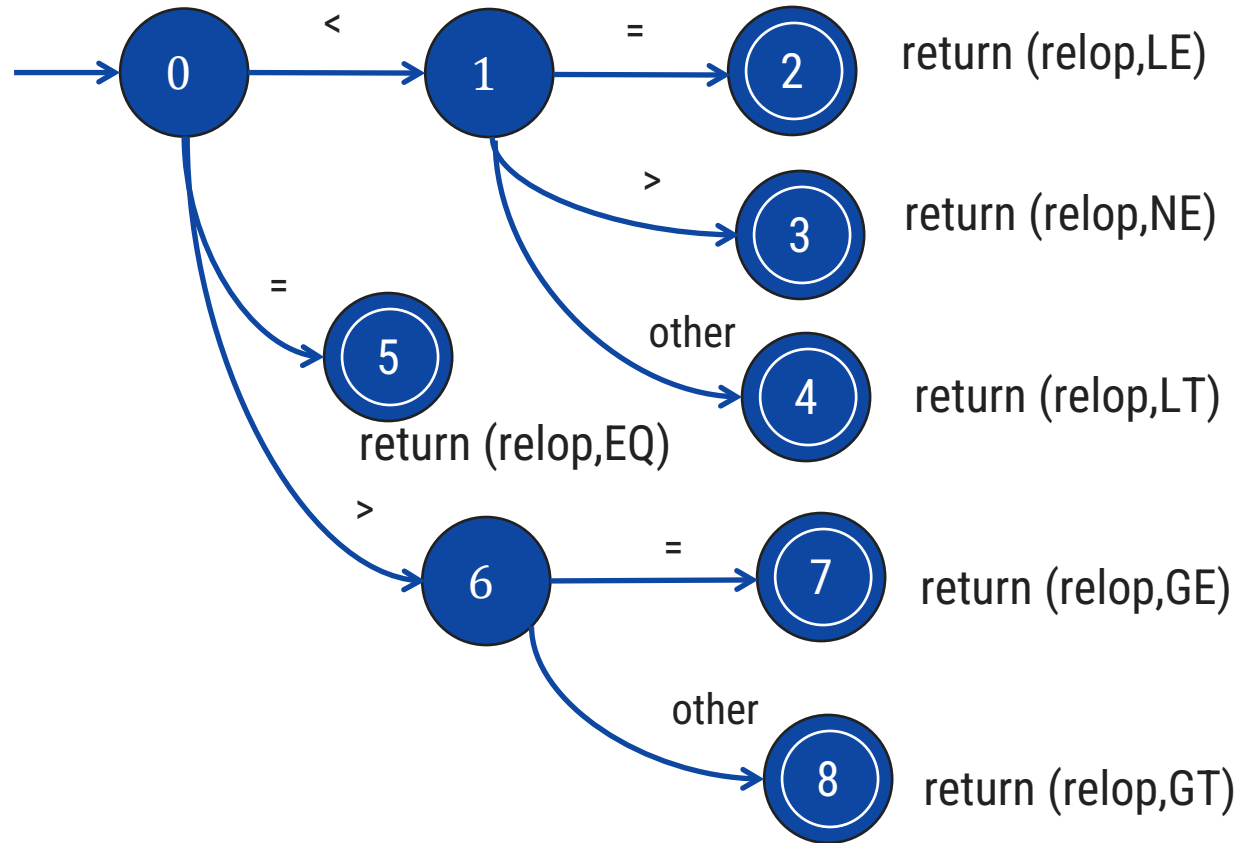


is a start state

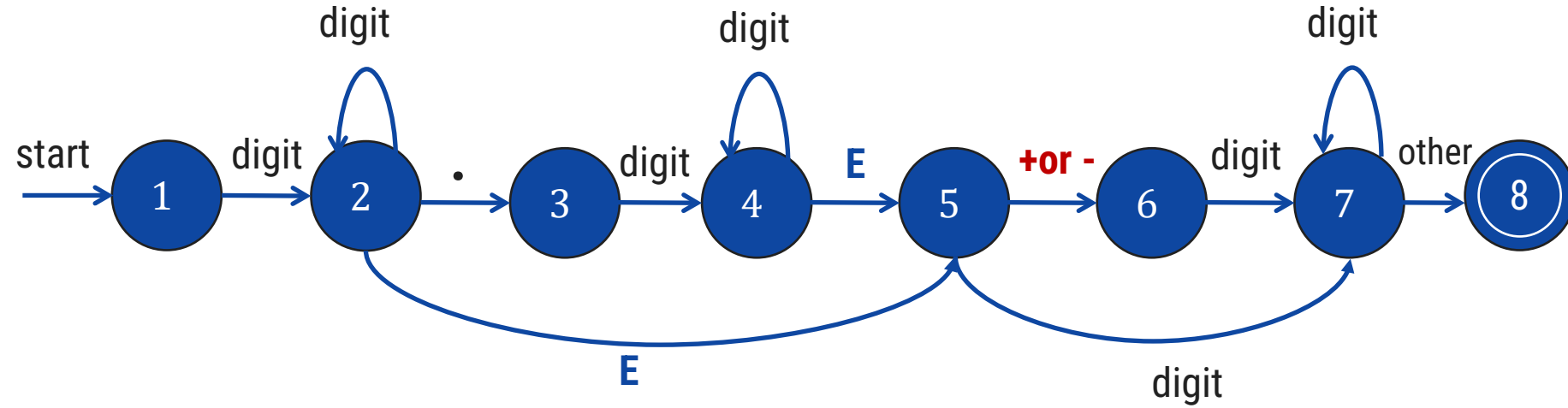


is a final state

# Transition Diagram : Relational operator



# Transition diagram : Unsigned number



3

5280

39.37

1.894 E - 4

2.56 E + 7

45 E + 6

96 E 2

# Finite Automata

# Finite Automata

- ▶ Finite Automata are recognizers.
  - ↳ FA simply say “Yes” or “No” about each possible input string.
- ▶ Finite Automata is a mathematical model consist of:
  1. Set of states  $S$
  2. Set of input symbol  $\Sigma$
  3. A transition function *move*
  4. Initial state  $S_0$
  5. Final states or accepting states  $F$

# Finite Automata

- ▶ A **finite automaton**, or **finite state machine** is a 5-tuple  $(Q, \Sigma, q_0, A, \delta)$  where
  - ↳  $Q$  is finite set of states;
  - ↳  $\Sigma$  is finite alphabet of *input symbols*;
  - ↳  $q_0 \in Q$  (*initial state*);
  - ↳  $A \subseteq Q$  (the set of *accepting states*);
  - ↳  $\delta$  is a function from  $Q \times \Sigma$  to  $Q$  (the *transition function*).
- ▶ For any element  $q$  of  $Q$  and any symbol  $a \in \Sigma$ , we interpret  $\delta(q, a)$  as the state to which the FA moves, if it is in state  $q$  and receives the input  $a$ .

# Types of finite automata

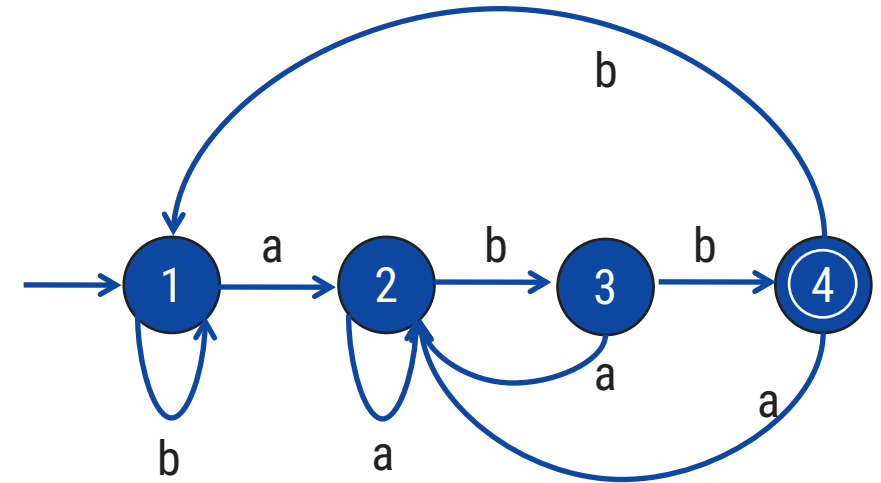
- Types of finite automata are:

## DFA

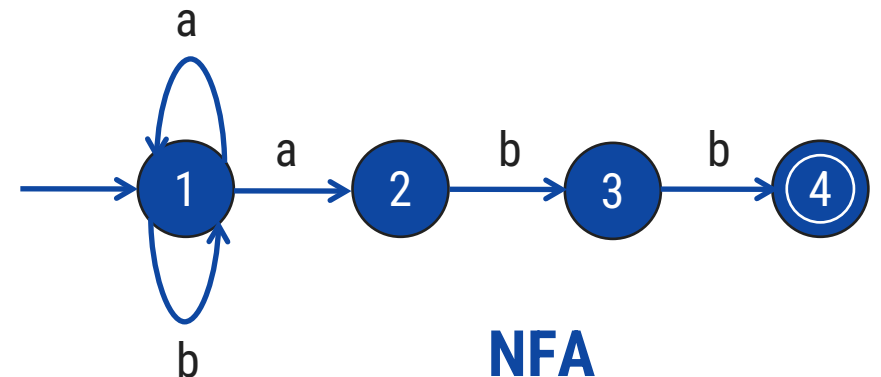
- Deterministic finite automata (DFA): have for each state **exactly one edge** leaving out for **each symbol**.

## NFA

- Nondeterministic finite automata (NFA): There are **no restrictions** on the edges leaving a state. There can be **several with the same symbol as label** and some edges can be labeled with  $\epsilon$ .



DFA



NFA

# Example: Finite Automata

►  $M = (Q, \Sigma, q_0, A, \delta)$

↳  $Q = \{q_0, q_1, q_2\}$

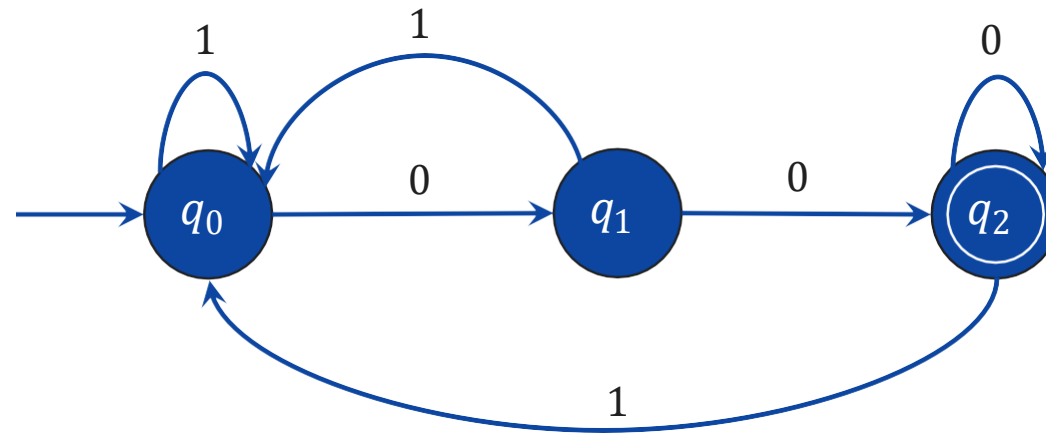
↳  $\Sigma = \{0,1\}$

↳  $q_0 = q_0$

↳  $A = \{q_2\}$

↳  $\delta$  is defined as

$\delta$	Input	
State	0	1
$q_0$	$q_1$	$q_0$
$q_1$	$q_2$	$q_0$
$q_2$	$q_2$	$q_0$



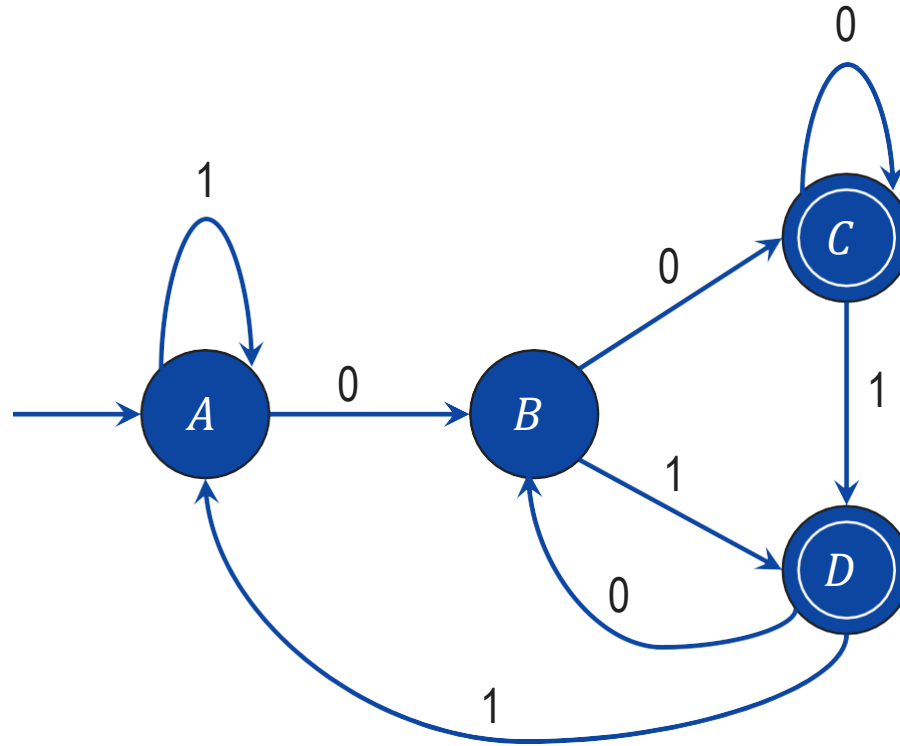


# Applications of FA

- ▶ Lexical analysis phase of a compiler.
- ▶ Design of digital circuit.
- ▶ String matching.
- ▶ Communication Protocol for information exchange.

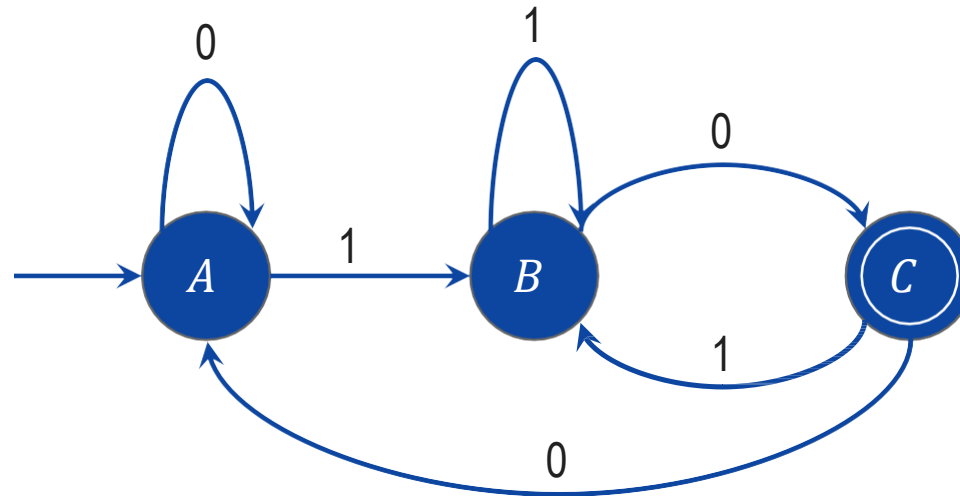
# FA Examples

- The string with next to last symbol as 0.



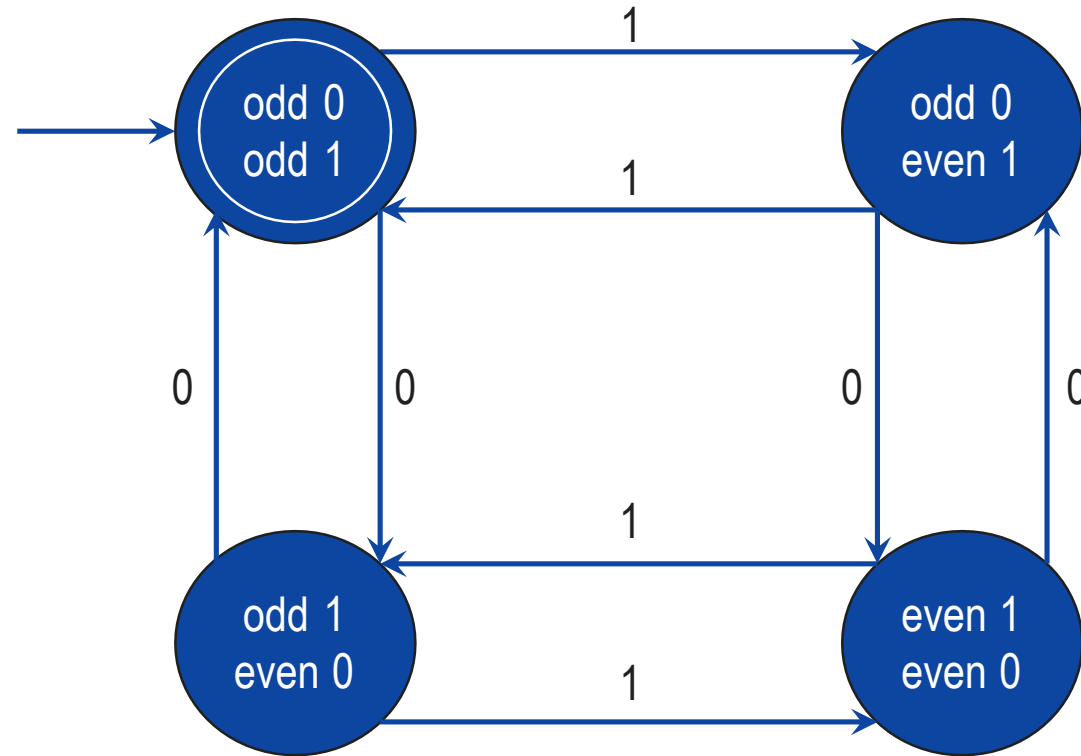
# FA Examples

- The strings ending with 10.



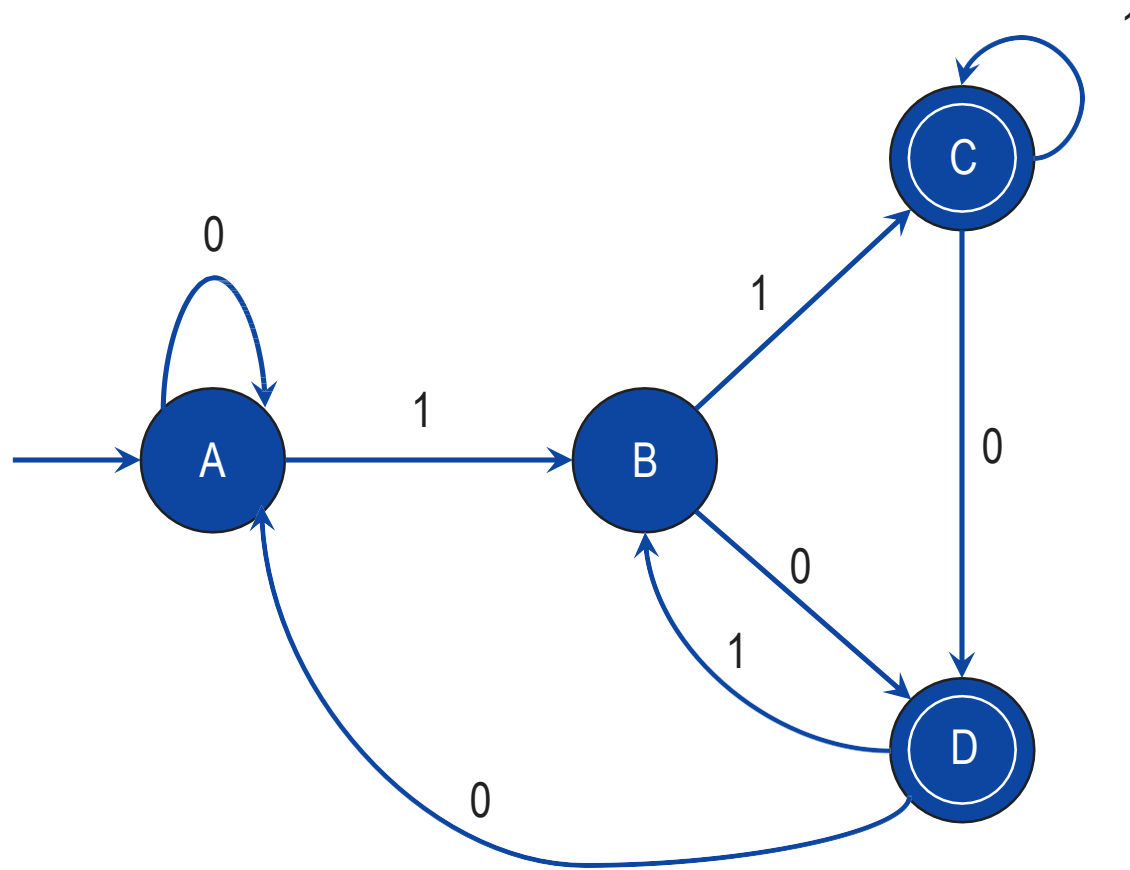
# FA Examples

- ▶ The string with number of 0's and number of 1's are odd



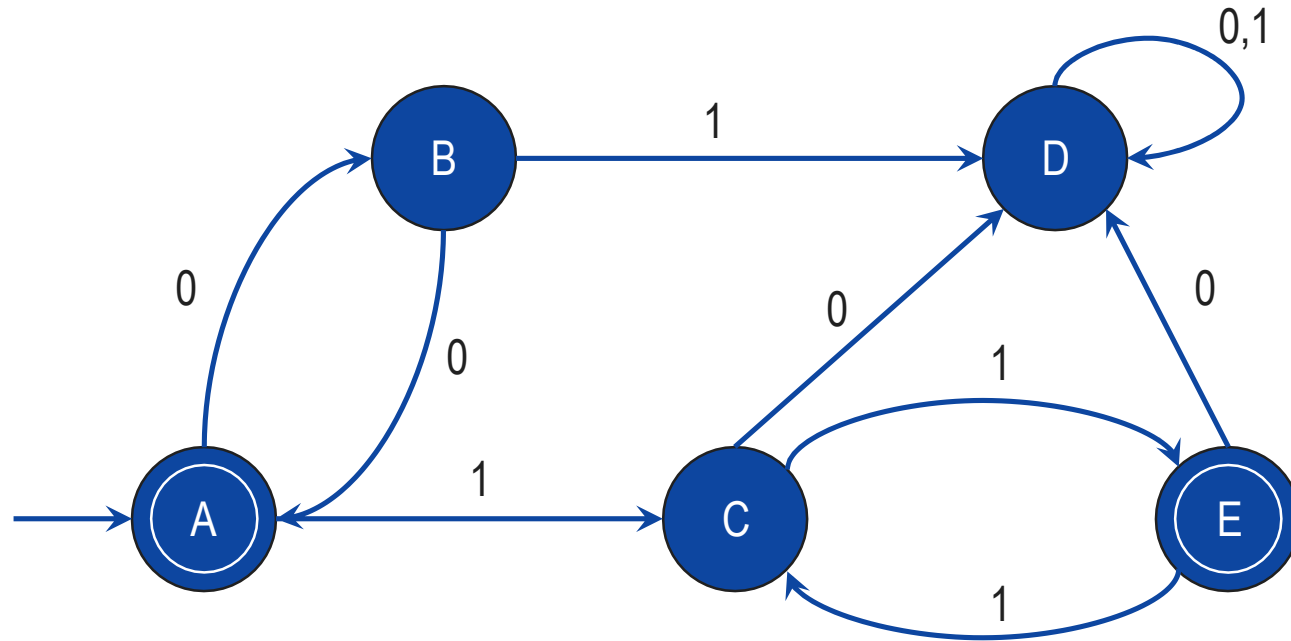
# FA Examples

- The string ending in 10 or 11



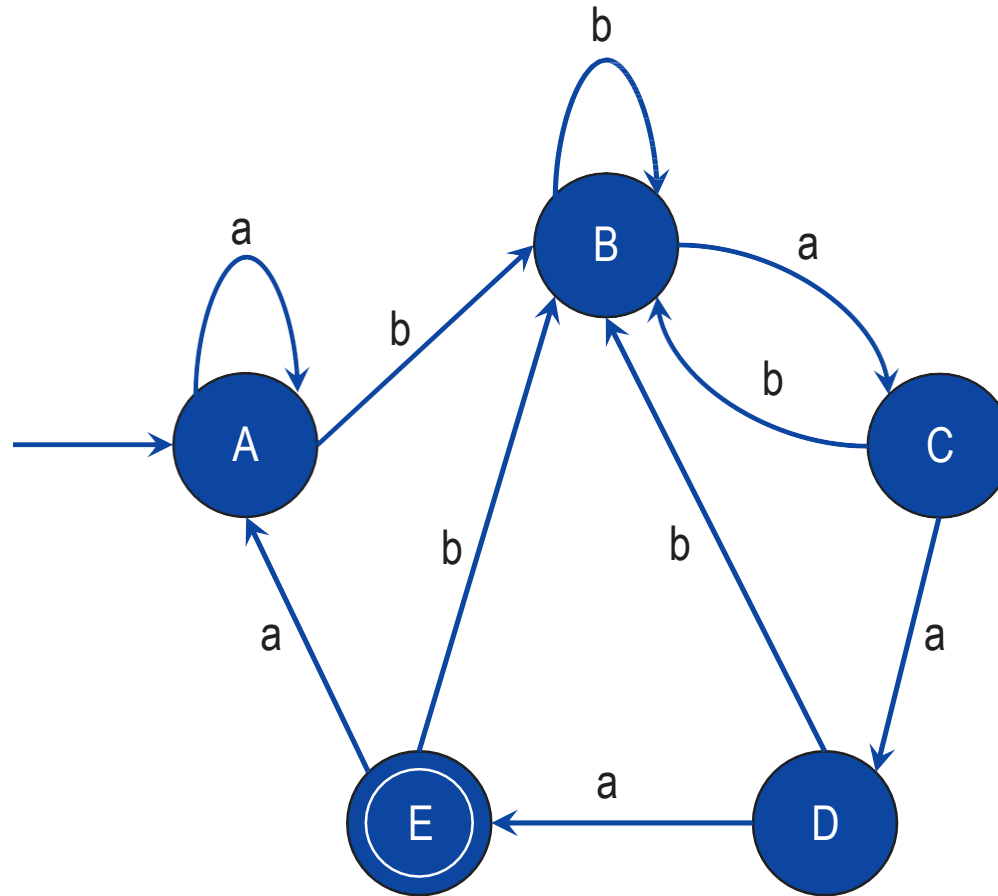
# FA Examples

- ▶ The string corresponding to Regular expression  $\{00\}^*\{11\}^*$



# FA Examples

►  $(a+b)^*baaa$



# Extended Transition Function $\delta^*$ for FA

► Let  $M = (Q, \Sigma, q_0, A, \delta)$  be an Finite Automata. We define the function

$$\delta^*: Q \times \Sigma^* \rightarrow Q$$

as follows:

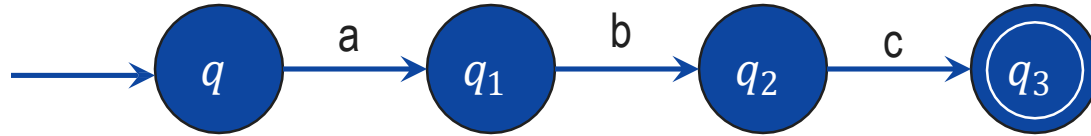
1. For any  $q \in Q$ ,  $\delta^*(q, \Lambda) = q$
2. For any  $q \in Q$ ,  $y \in \Sigma^*$ , and  $a \in \Sigma$ ,

$$\delta^*(q, ya) = \delta(\delta^*(q, y), a)$$



# $\delta^*$ Example

## ► Consider FA



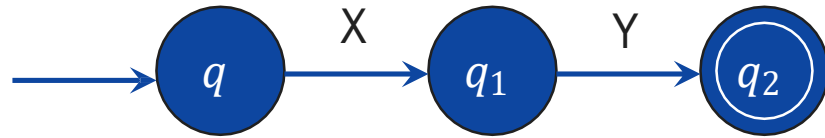
$\delta^*(q, \Lambda)$	$= q$
$\delta^*(q, ya)$	$= \delta(\delta^*(q, y), a)$

## ► Calculate $\delta^*(q, abc)$

$$\begin{aligned}\delta^*(q, abc) &= \delta(\delta^*(q, ab), c) \\ &= \delta(\delta(\delta^*(q, a), b), c) \\ &= \delta(\delta(\delta^*(q, \Lambda a), b), c) \\ &= \delta(\delta(\delta(\underline{\delta^*(q, \Lambda)}, a), b), c) \\ &= \delta(\delta(\delta(\underline{q}, a), b), c) \\ &= \delta(\delta(\underline{q_1}, b), c) \\ &= \delta(\underline{q_2}, c) \\ &= \underline{q_3}\end{aligned}$$

# Exercise

Consider following FA



1. calculate  $\delta^*(q, XY)$

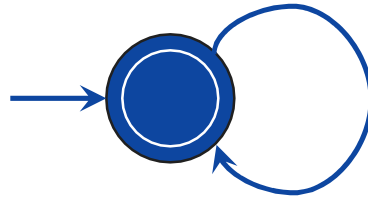
$\delta^*(q, \Lambda)$	$= q$
$\delta^*(q, ya)$	$= \delta(\delta^*(q, y), a)$

# Acceptance by FA

- ▶ Let  $M = (Q, \Sigma, q_0, A, \delta)$  be an FA. A string  $x \in \Sigma^*$  is **accepted** by  $M$  if  $\delta^*(q_0, x) \in A$ . If string is not accepted, we can say it is **rejected** by  $M$ . The language accepted by  $M$ , or the language recognized by  $M$ , is the set

$$L(M) = \{x \in \Sigma^* \mid x \text{ is accepted by } M\}$$

- ▶ If  $L$  is any language over  $\Sigma$ ,  $L$  is accepted or recognized by  $M$  if and only if  $L = L(M)$ .



# Union, Intersection & Complement of Languages

- Suppose  $M_1 = (Q_1, \Sigma, q_1, A_1, \delta_1)$  and  $M_2 = (Q_2, \Sigma, q_2, A_2, \delta_2)$  accepts languages  $L_1$  and  $L_2$ , respectively. Let  $M$  be an FA defined by  $M = (Q, \Sigma, q_0, A, \delta)$ , where

$$Q = Q_1 \times Q_2$$

$$q_0 = (q_1, q_2)$$

and the transition function  $\delta$  is defined by the formula

$$\delta((p, q), a) = (\delta_1(p, a), \delta_2(q, a))$$

for any  $p \in Q_1$  and  $q \in Q_2$  and  $a \in \Sigma$  then

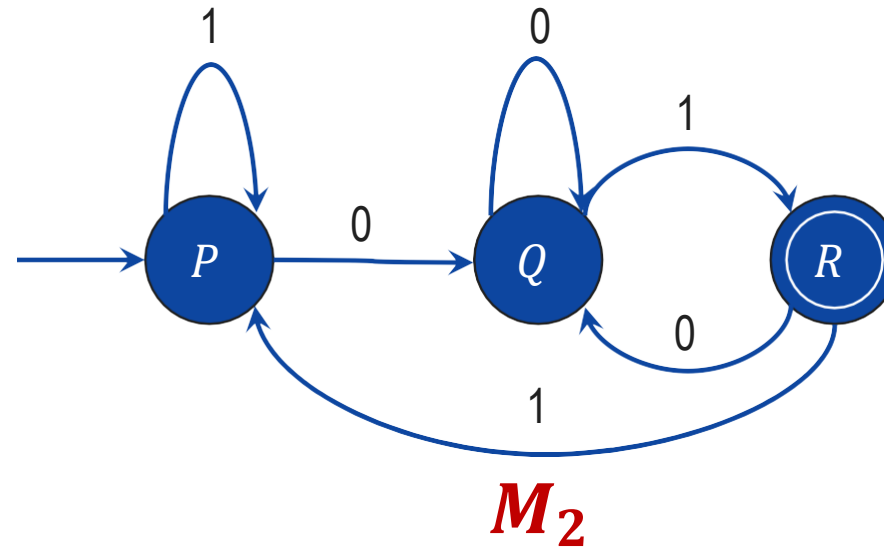
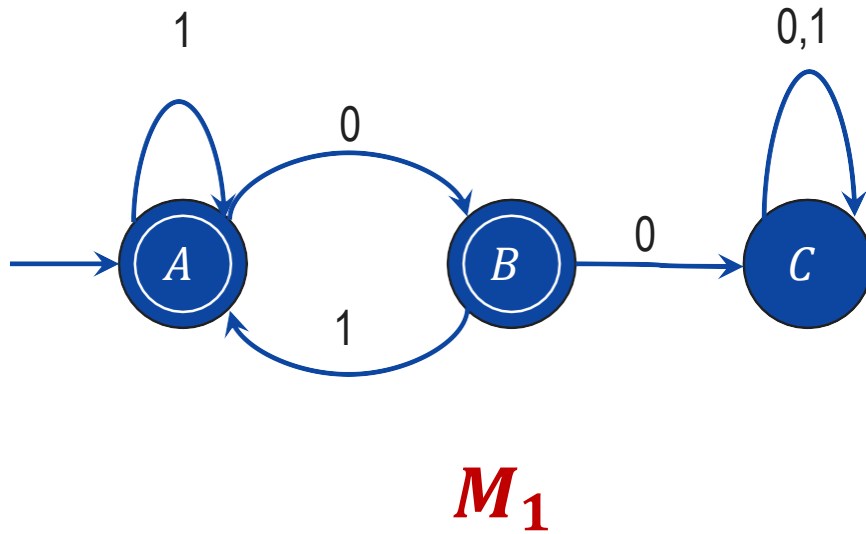
1. if  $A = \{(p, q) \mid \mathbf{p} \in \mathbf{A_1} \text{ or } \mathbf{q} \in \mathbf{A_2}\}$ ,  $M$  accepts the language  $\mathbf{L_1 \cup L_2}$ ;
2. if  $A = \{(p, q) \mid \mathbf{p} \in \mathbf{A_1} \text{ and } \mathbf{q} \in \mathbf{A_2}\}$ ,  $M$  accepts the language  $\mathbf{L_1 \cap L_2}$ ;
3. if  $A = \{(p, q) \mid \mathbf{p} \in \mathbf{A_1} \text{ and } \mathbf{q} \notin \mathbf{A_2}\}$ ,  $M$  accepts the language  $\mathbf{L_1 - L_2}$ ;

# Example

Draw Finite Automata for following languages:

1.  $L_1 = \{x/x \text{ 00 is not substring of } x, x \in \{0,1\}^*\}$
2.  $L_2 = \{x/x \text{ ends with 01, } x \in \{0,1\}^*\}$

Draw FA for  $L_1 \cup L_2$ ,  $L_1 \cap L_2$  and  $L_1 - L_2$

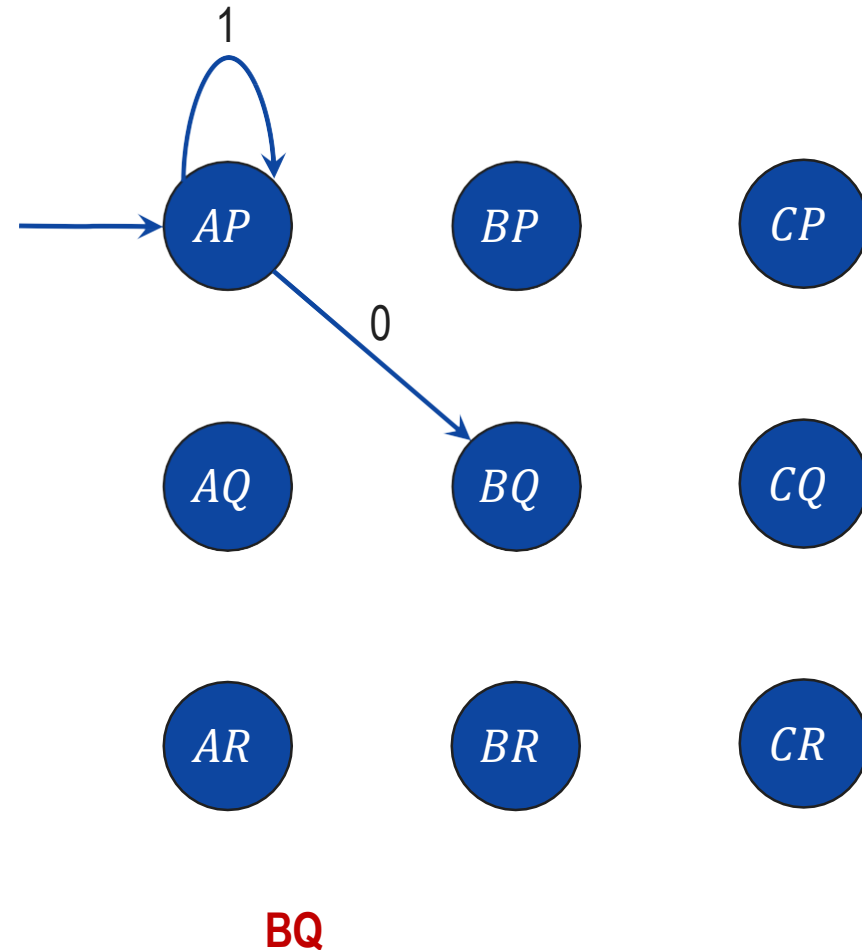
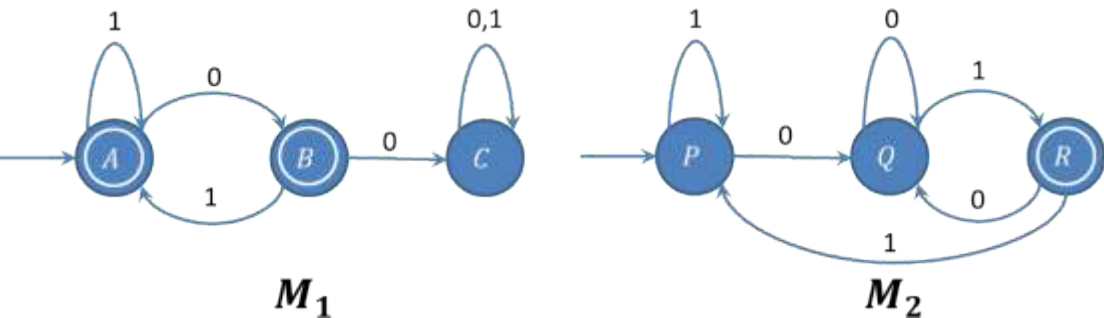


# Computation for $L_1 \cup L_2$ , $L_1 \cap L_2$ and $L_1 - L_2$

- ▶ Here  $Q_1 = \{A, B, C\}$  and  $Q_2 = \{P, Q, R\}$
- ▶ So,  $Q = Q_1 \times Q_2 = \{AP, AQ, AR, BP, BQ, BR, CP, CQ, CR\}$
- ▶  $q_0 = (q_1, q_2) = (A, P)$
- ▶ Computing Transition Function  $\delta$

$$\begin{aligned}\delta((A, P), 0) &= (\delta(A, 0), \delta(P, 0)) \\ &= BQ\end{aligned}$$

$$\begin{aligned}\delta((A, P), 1) &= (\delta(A, 1), \delta(P, 1)) \\ &= AP\end{aligned}$$



# Computation for $L_1 \cup L_2$ , $L_1 \cap L_2$ and $L_1 - L_2$

$$\begin{aligned}\delta((B, Q), 0) &= (\delta(B, 0), \delta(Q, 0)) \\ &= CQ\end{aligned}$$

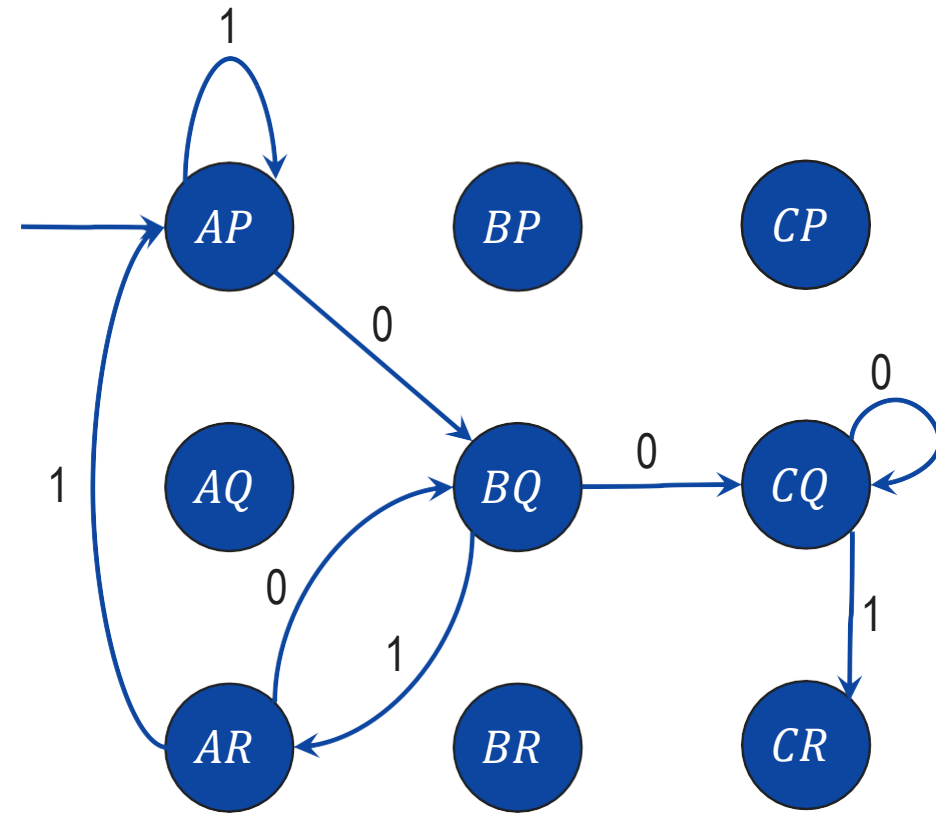
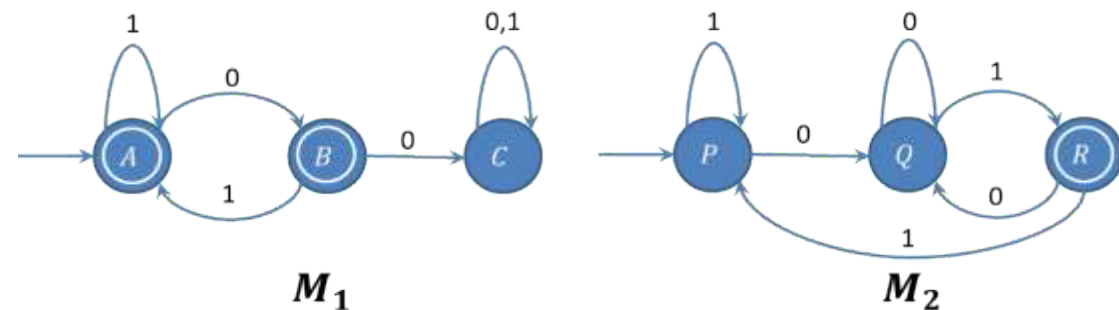
$$\begin{aligned}\delta((B, Q), 1) &= (\delta(B, 1), \delta(Q, 1)) \\ &= AR\end{aligned}$$

$$\begin{aligned}\delta((C, Q), 0) &= (\delta(C, 0), \delta(Q, 0)) \\ &= CQ\end{aligned}$$

$$\begin{aligned}\delta((C, Q), 1) &= (\delta(C, 1), \delta(Q, 1)) \\ &= CR\end{aligned}$$

$$\begin{aligned}\delta((A, R), 0) &= (\delta(A, 0), \delta(R, 0)) \\ &= BQ\end{aligned}$$

$$\begin{aligned}\delta((A, R), 1) &= (\delta(A, 1), \delta(R, 1)) \\ &= AP\end{aligned}$$



~~$BQ$~~   ~~$CQ$~~   ~~$AR$~~   $CR$

# Computation for $L_1 \cup L_2$ , $L_1 \cap L_2$ and $L_1 - L_2$

$$\delta((C, R), 0) = (\delta(C, 0), \delta(R, 0))$$

$$= CQ$$

$$\delta((C, R), 1) = (\delta(C, 1), \delta(R, 1))$$

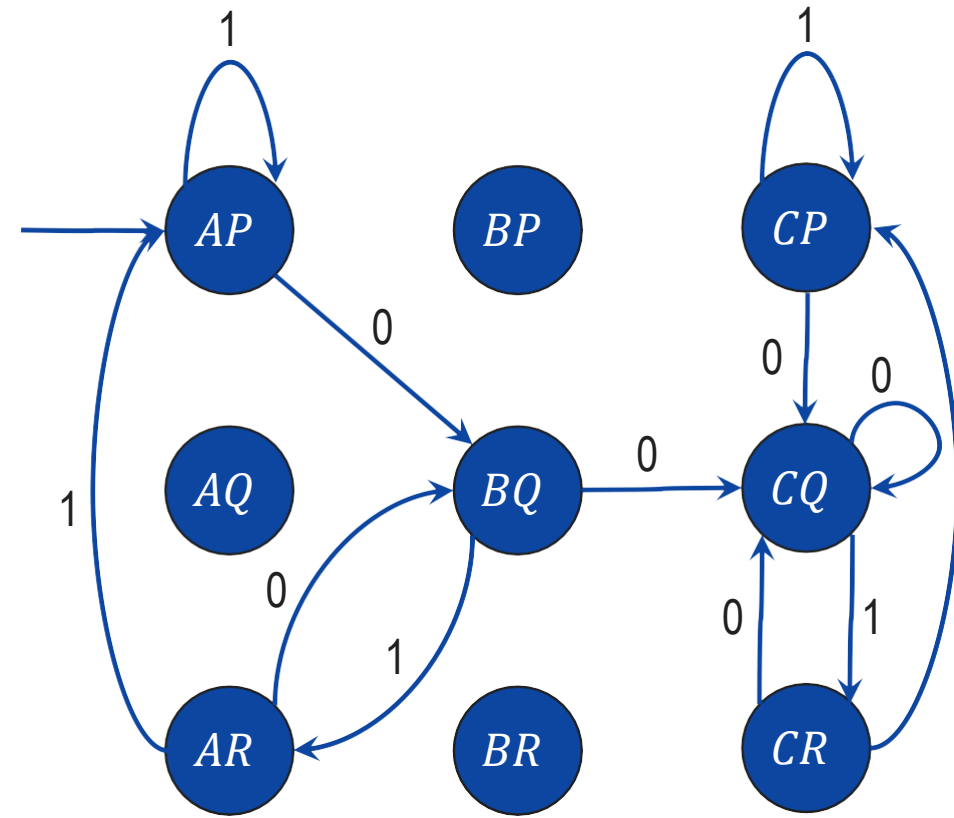
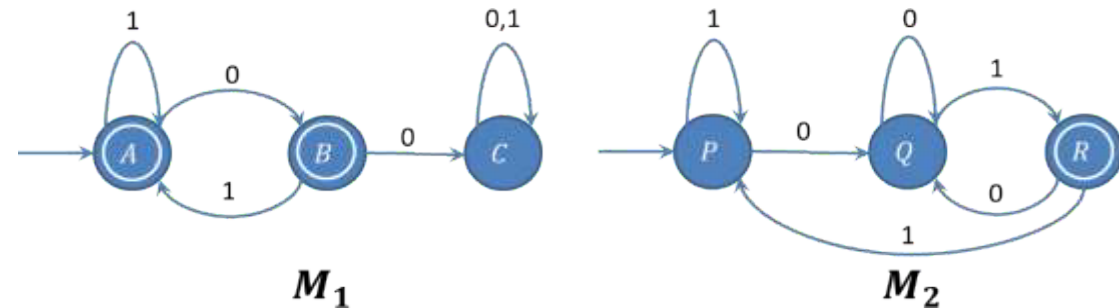
$$= CP$$

$$\delta((C, P), 0) = (\delta(C, 0), \delta(P, 0))$$

$$= CQ$$

$$\delta((C, P), 1) = (\delta(C, 1), \delta(P, 1))$$

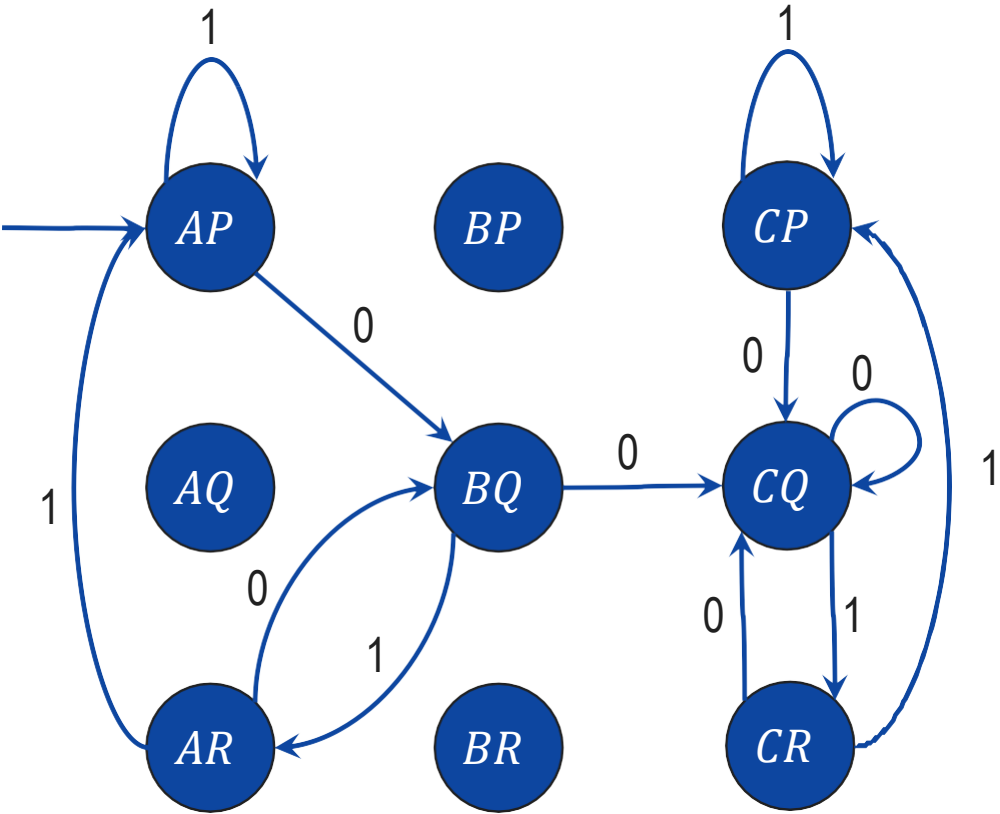
$$= CP$$



~~$BQ$~~   ~~$CQ$~~   ~~$AR$~~   ~~$CR$~~   ~~$CP$~~



# Removing Unconnected States

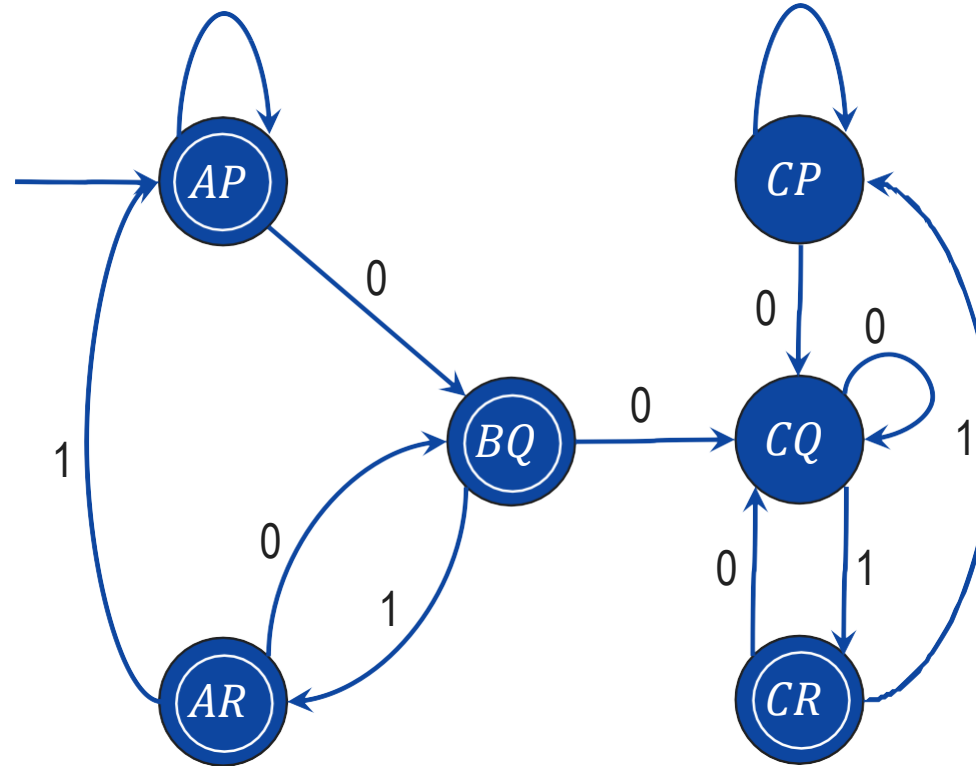


$$L_1 \cup L_2$$

$$A_1 = \{A, B\}$$

$$A_2 = \{R\}$$

As per theorem stated earlier,  
the states which consists A  
or B or R will be Accepting  
states in the resultant FA.



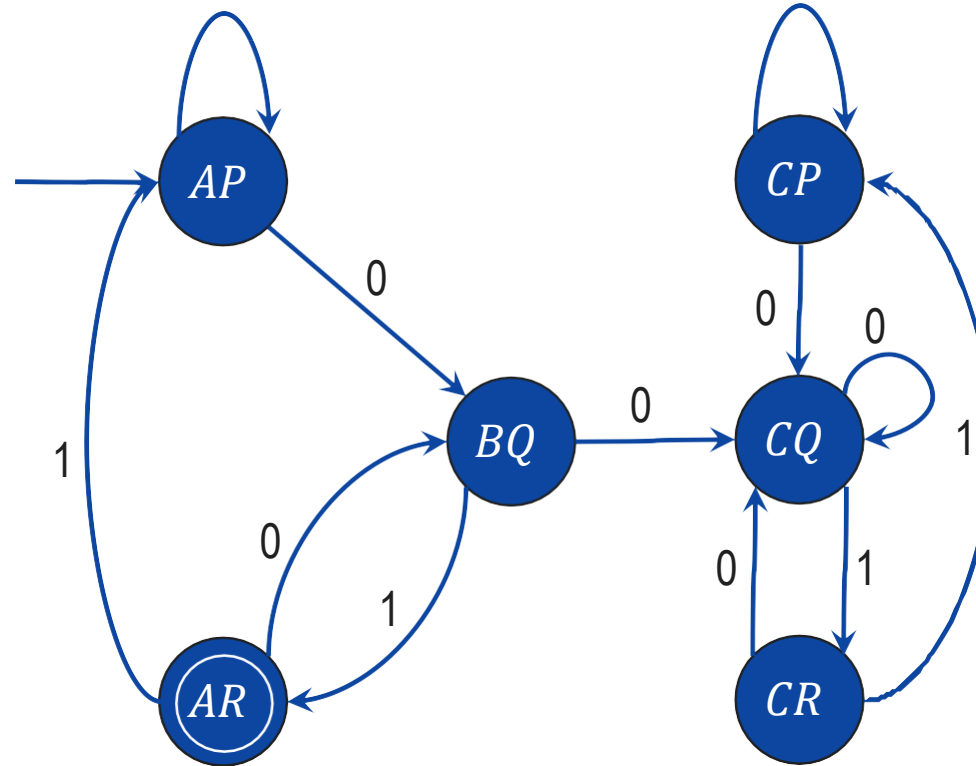
Accepting States  $A = \{AP, AR, BQ, CR\}$

$$L_1 \cap L_2$$

$$A_1 = \{A, B\}$$

$$A_2 = \{R\}$$

Therefore, as per theorem stated earlier, the states which consists (A and R) and (B and R) will be Accepting states in the resultant FA.



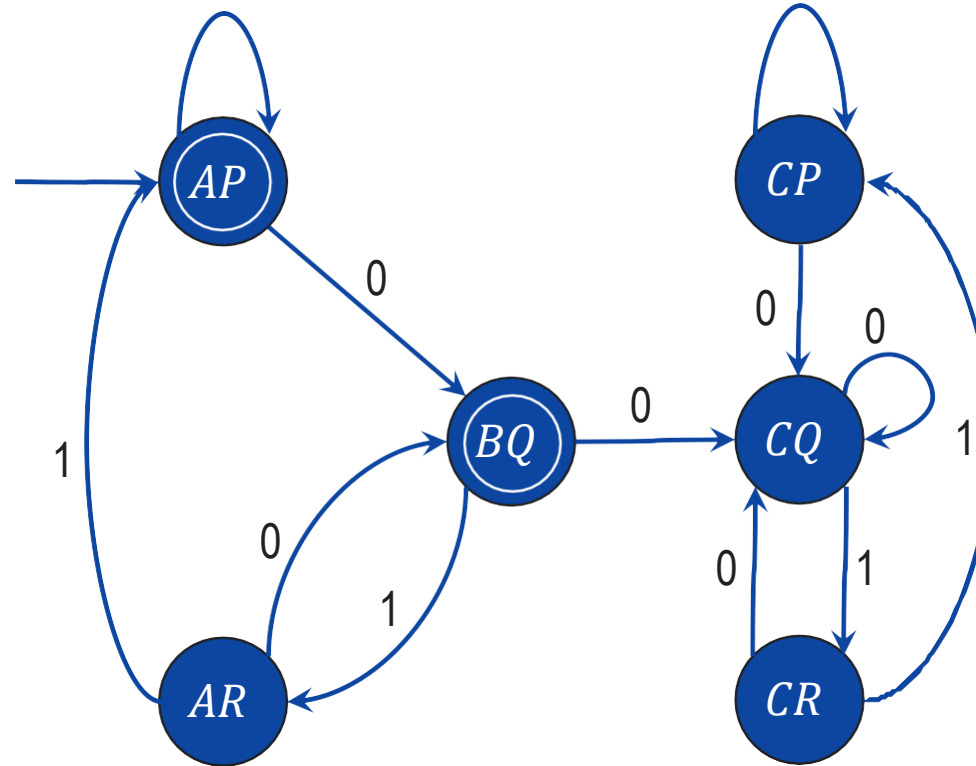
Accepting States  $A = \{AR\}$

$$L_1 - L_2$$

$$A_1 = \{A, B\}$$

$$A_2 = \{R\}$$

Therefore, as per theorem stated earlier, the states which consists A or B but should not contain R with them will be Accepting states in the resultant FA.



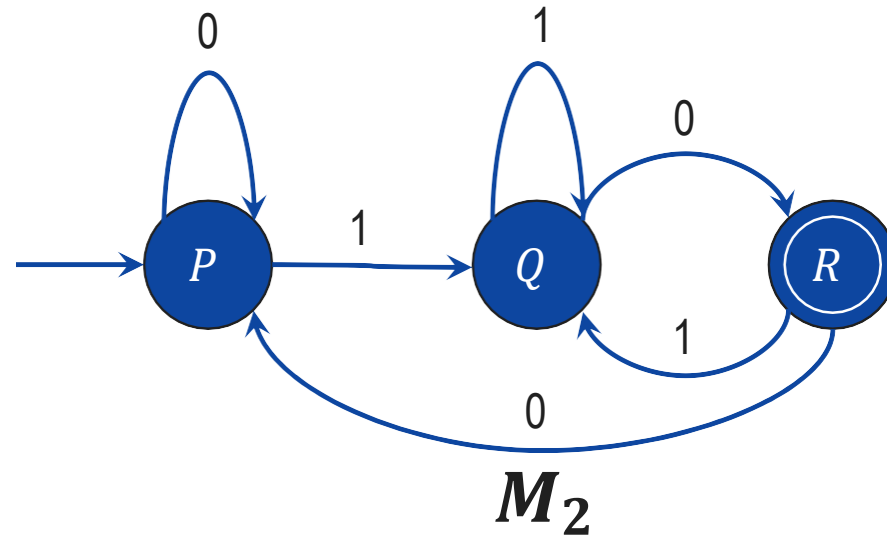
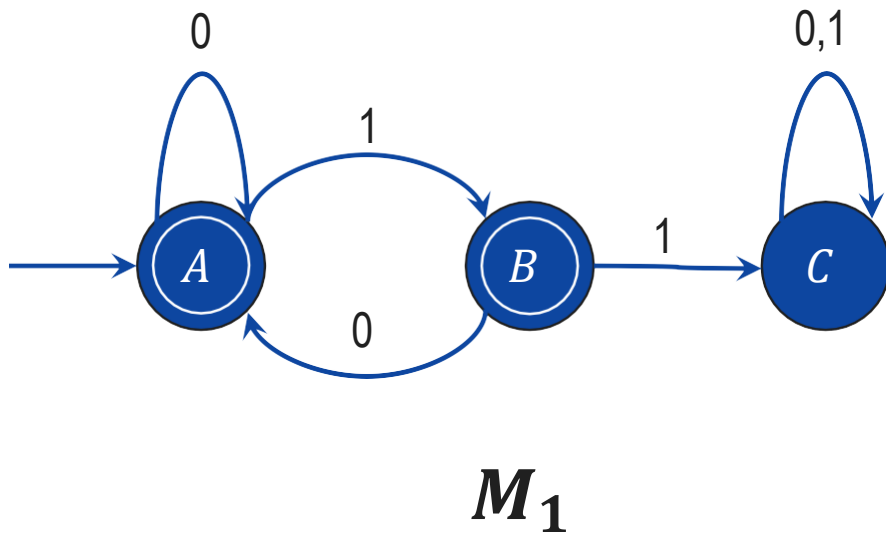
Accepting States  $A = \{AP, BQ\}$

# Exercise

Draw Finite Automata for following languages:

1.  $L_1 = \{x/x \text{ 11 is not substring of } x, x \in \{0,1\}^*\}$
2.  $L_2 = \{x/x \text{ ends with } 10, x \in \{0,1\}^*\}$

Draw FA for  $L_1 \cup L_2$ ,  $L_1 \cap L_2$  and  $L_1 - L_2$ .



# Nondeterministic Finite Automata (NFA)

- ▶ A nondeterministic finite automaton is a 5-tuple  $M = (Q, \Sigma, q_0, A, \delta)$  where  $Q$  and  $\Sigma$  are nonempty finite sets,  $q_0 \in Q$ ,  $A \subseteq Q$  and
$$\delta : Q \times \Sigma \rightarrow 2^Q$$

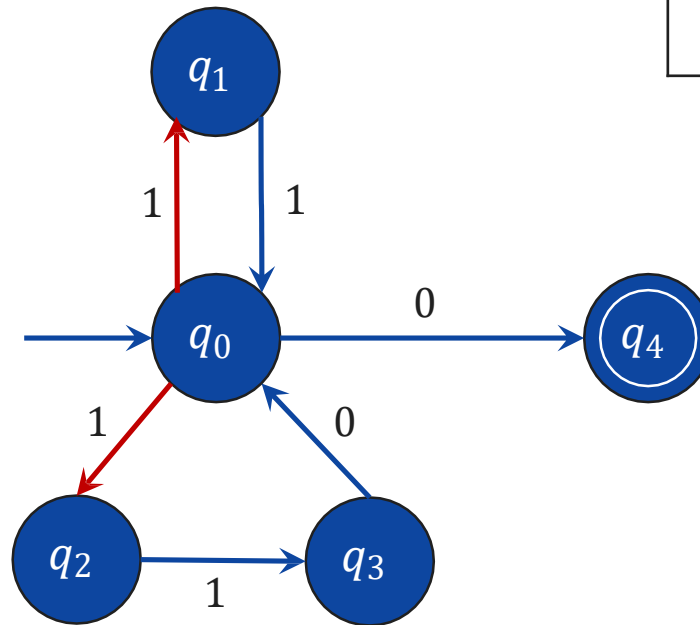
$Q$  is the set of states,  $\Sigma$  is the alphabet,  $q_0$  is the initial state and  $A$  is the set of accepting states.

# Example of NFA for $\{11,110\}^*\{0\}$

$$M = (Q, \Sigma, q_0, A, \delta)$$

- $Q = \{q_0, q_1, q_2, q_3, q_4\}$
- $\Sigma = \{0,1\}$
- $q_0 = q_0$
- $A = \{q_4\}$
- $\delta$  is defined as

$\delta$	Input	
State	0	1
$q_0$	$\{q_4\}$	$\{q_1, q_2\}$
$q_1$	$\emptyset$	$\{q_0\}$
$q_2$	$\emptyset$	$\{q_3\}$
$q_3$	$\{q_0\}$	$\emptyset$
$q_4$	$\emptyset$	$\emptyset$



# Nonrecursive Definition of $\delta^*$ for NFA

- ▶ For an NFA  $M = (Q, \Sigma, q_0, A, \delta)$ , and any  $p \in Q$ ,  $\delta^*(p, \Lambda) = \{p\}$ . For any  $p \in Q$  and any  $x = a_1a_2 \dots a_n \in \Sigma^*$  (with  $n \geq 1$ ),  $\delta^*(p, x)$  is the set of all states  $q$  for which there is a sequence of states  $p = p_0, p_1, \dots, p_{n-1}, p_n = q$  satisfying
$$p_i \in \delta(p_{i-1}, a_i) \text{ for each } i \text{ with } 1 \leq i \leq n$$



# Recursive Definition of $\delta^*$ for NFA

► Let  $M = (Q, \Sigma, q_0, A, \delta)$  be an NFA. The function  $\delta^*: Q \times \Sigma^* \rightarrow 2^Q$  is defined as follows.

1. For any  $q \in Q$ ,  $\delta^*(q, \Lambda) = \{q\}$ .
2. For any  $q \in Q$ ,  $y \in \Sigma^*$ , and  $a \in \Sigma$ ,

$$\delta^*(q, ya) = \bigcup_{r \in \delta^*(q, y)} \delta(r, a)$$

# Acceptance by NFA

- ▶ Let  $M = (Q, \Sigma, q_0, A, \delta)$  be an NFA. The string  $x \in \Sigma^*$  is accepted by  $M$  if  $\delta^*(q_0, x) \cap A \neq \emptyset$ . The language recognized, or accepted, by  $M$  is the set  $L(M)$  of all string accepted by  $M$ . For any language  $L \subseteq \Sigma^*$ ,  $L$  is recognized by  $M$  if  $L = L(M)$ .

# Example: Recursive Definition of $\delta^*$ in NFA

►  $M = (Q, \Sigma, q_0, A, \delta)$

→  $Q = \{q_0, q_1, q_2, q_3\}$

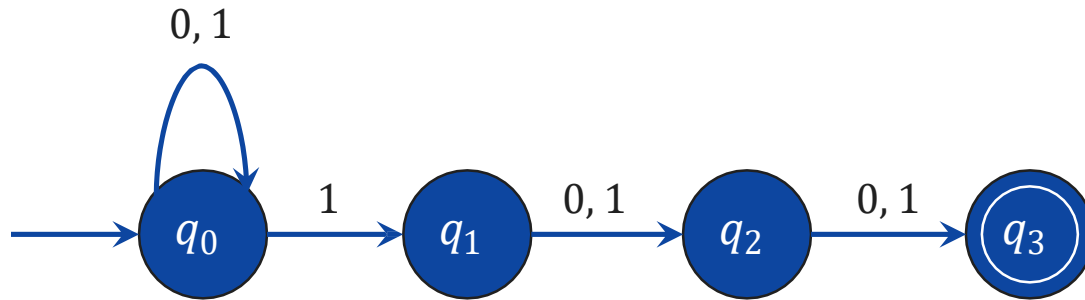
→  $\Sigma = \{0,1\}$

→  $q_0 = q_0$

→  $A = \{q_3\}$

→  $\delta$  is defined as

$\delta$	Input	
State	<b>0</b>	<b>1</b>
$q_0$	$\{q_0\}$	$\{q_0, q_1\}$
$q_1$	$\{q_2\}$	$\{q_2\}$
$q_2$	$\{q_3\}$	$\{q_3\}$
$q_3$	$\emptyset$	$\emptyset$



# Example: Recursive Definition of $\delta^*$ in NFA

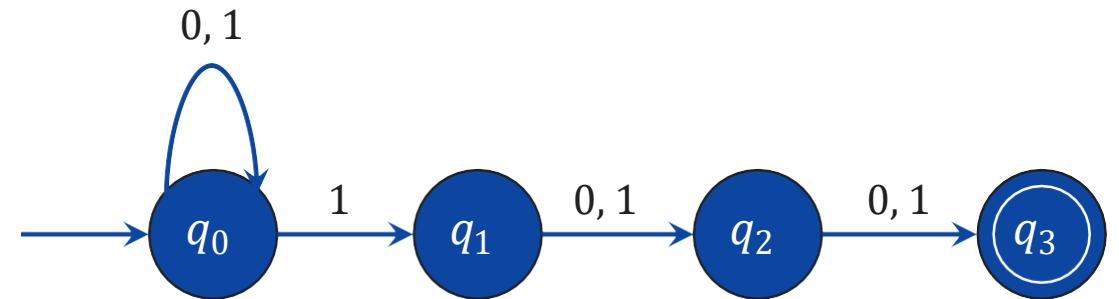
$$\delta^*(q_0, 0) = \bigcup_{r \in \delta^*(q_0, \Lambda)} \delta(r, 0)$$

$$= \bigcup_{r \in \{q_0\}} \delta(r, 0)$$

$$= \delta(q_0, 0)$$

$$= \{q_0\}$$

$\delta^*(q, \Lambda)$	$= q$
$\delta^*(q, ya)$	$= \bigcup_{r \in \delta^*(q, y)} \delta(r, a)$



# Example: Recursive Definition of $\delta^*$ in NFA

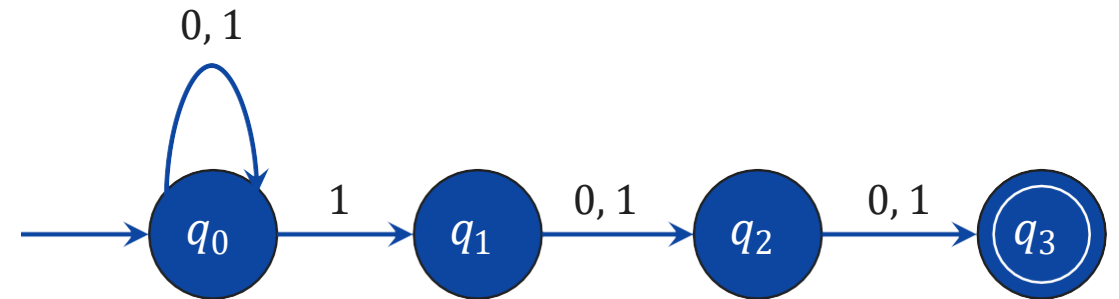
$$\delta^*(q_0, 1) = \bigcup_{r \in \delta^*(q_0, \Lambda)} \delta(r, 1)$$

$$= \bigcup_{r \in \{q_0\}} \delta(r, 1)$$

$$= \delta(q_0, 1)$$

$$= \{q_0, q_1\}$$

$\delta^*(q, \Lambda)$	$= q$
$\delta^*(q, ya)$	$= \bigcup_{r \in \delta^*(q, y)} \delta(r, a)$



# Example: Recursive Definition of $\delta^*$ in NFA

$$\delta^*(q_0, 11) = \bigcup_{r \in \delta^*(q_0, 1)} \delta(r, 1)$$

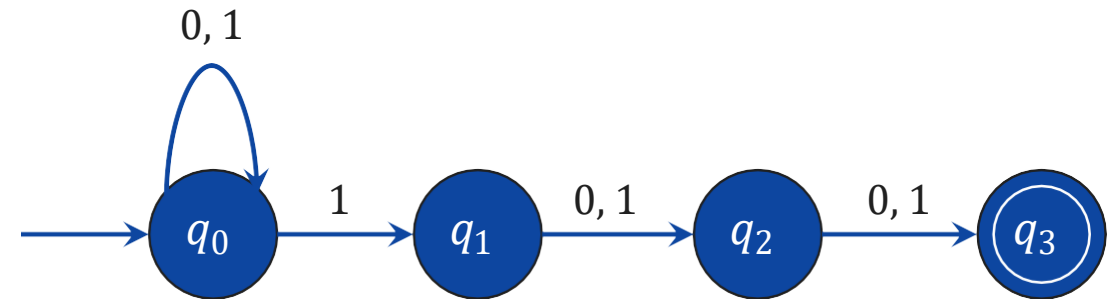
$$= \bigcup_{r \in \{q_0, q_1\}} \delta(r, 1)$$

$$= \delta(q_0, 1) \cup \delta(q_1, 1)$$

$$= \{q_0, q_1\} \cup \{q_2\}$$

$$= \{q_0, q_1, q_2\}$$

$\delta^*(q, \Lambda)$	$= q$
$\delta^*(q, ya)$	$= \bigcup_{r \in \delta^*(q, y)} \delta(r, a)$



# Example: Recursive Definition of $\delta^*$ in NFA

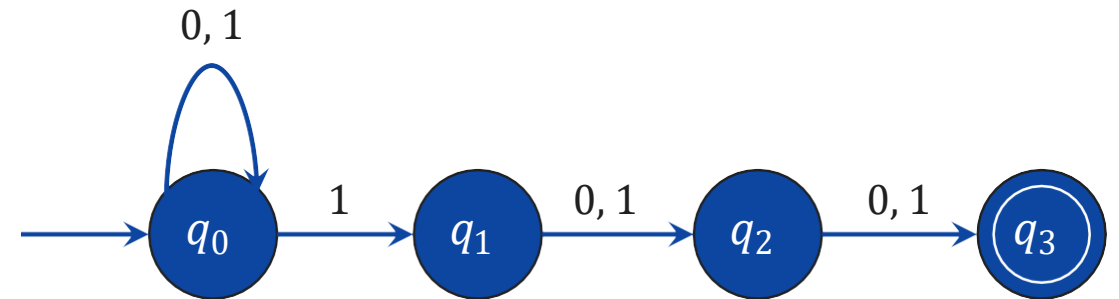
$$\delta^*(q_0, 01) = \bigcup_{r \in \delta^*(q_0, 0)} \delta(r, 1)$$

$$= \bigcup_{r \in \{q_0\}} \delta(r, 1)$$

$$= \delta(q_0, 1)$$

$$= \{q_0, q_1\}$$

$\delta^*(q, \Lambda)$	$= q$
$\delta^*(q, ya)$	$= \bigcup_{r \in \delta^*(q, y)} \delta(r, a)$



# Example: Recursive Definition of $\delta^*$ in NFA

$$\delta^*(q_0, 111) = \bigcup_{r \in \delta^*(q_0, 11)} \delta(r, 1)$$

$$= \bigcup_{r \in \{q_0, q_1, q_2\}} \delta(r, 1)$$

$$= \delta(q_0, 1) \cup \delta(q_1, 1) \cup \delta(q_2, 1)$$

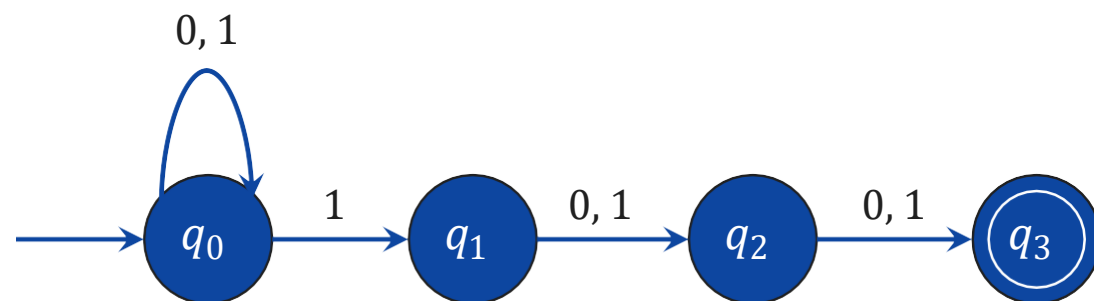
$$= \{q_0, q_1, q_2, q_3\}$$

$$\delta^*(q_0, 011) = \bigcup_{r \in \delta^*(q_0, 01)} \delta(r, 1)$$

$$= \delta(q_0, 1) \cup \delta(q_1, 1)$$

$$= \{q_0, q_1, q_2\}$$

$\delta^*(q, \Lambda)$	$= q$
$\delta^*(q, ya)$	$= \bigcup_{r \in \delta^*(q, y)} \delta(r, a)$

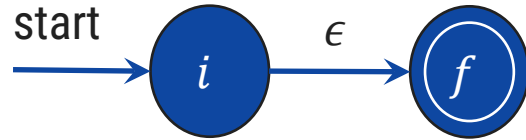




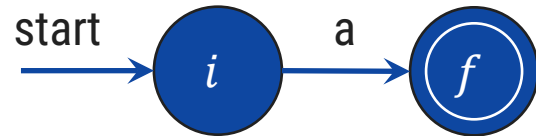
# Regular expression to NFA using Thompson's rule

# Regular expression to NFA using Thompson's rule

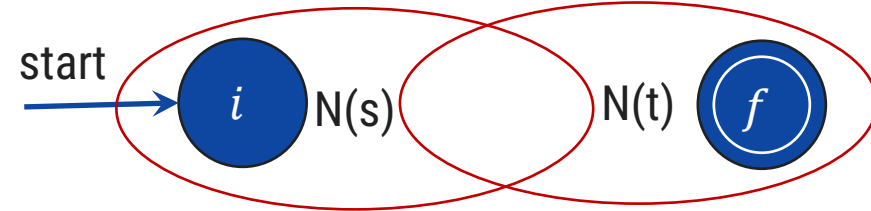
1. For  $\epsilon$ , construct the NFA



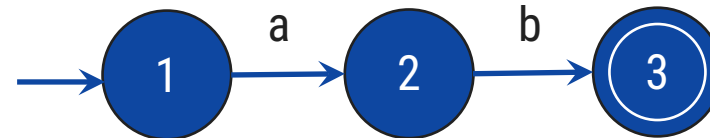
2. For  $a$  in  $\Sigma$ , construct the NFA



3. For regular expression  $st$

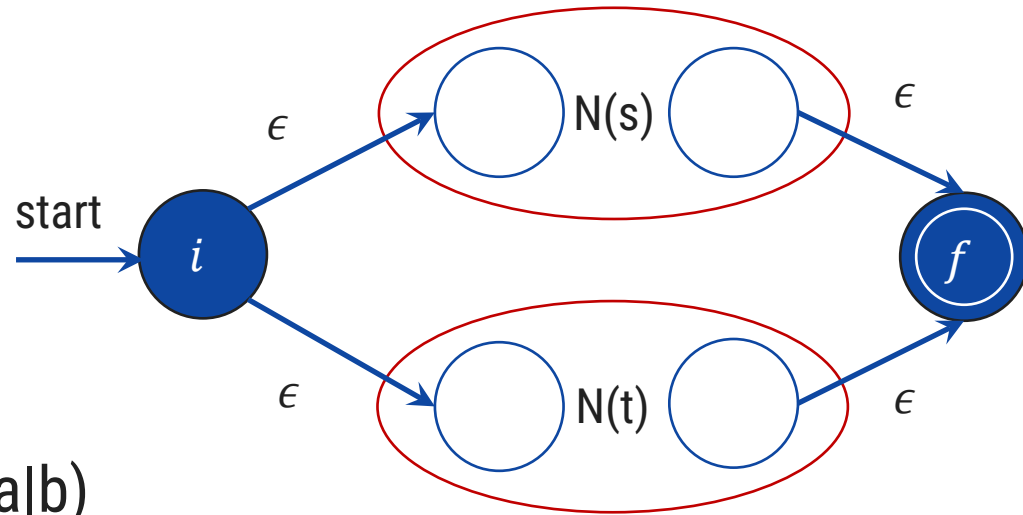


Ex:  $ab$

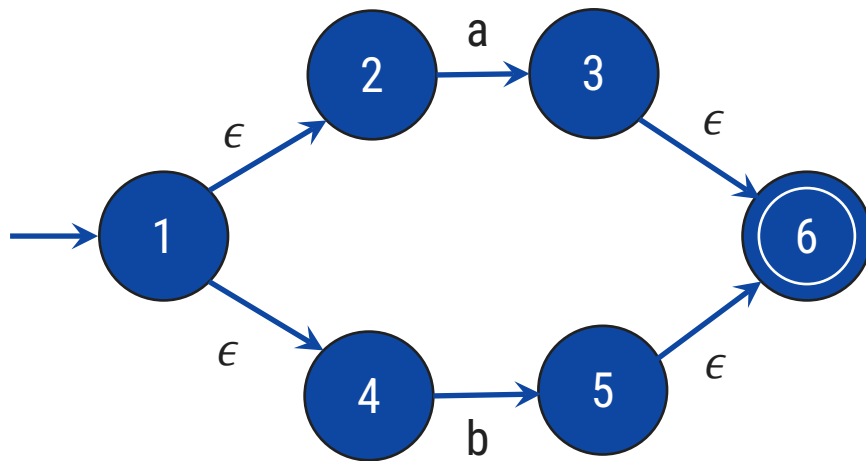


# Regular expression to NFA using Thompson's rule

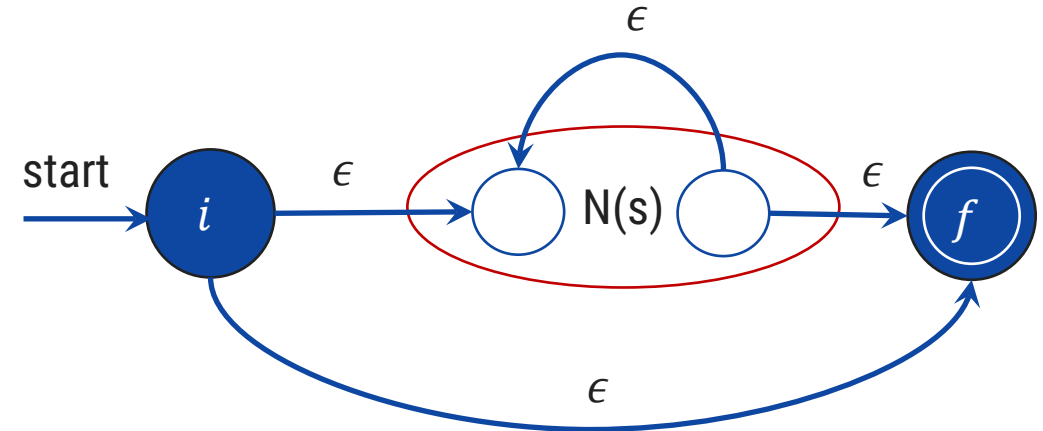
4. For regular expression  $s|t$



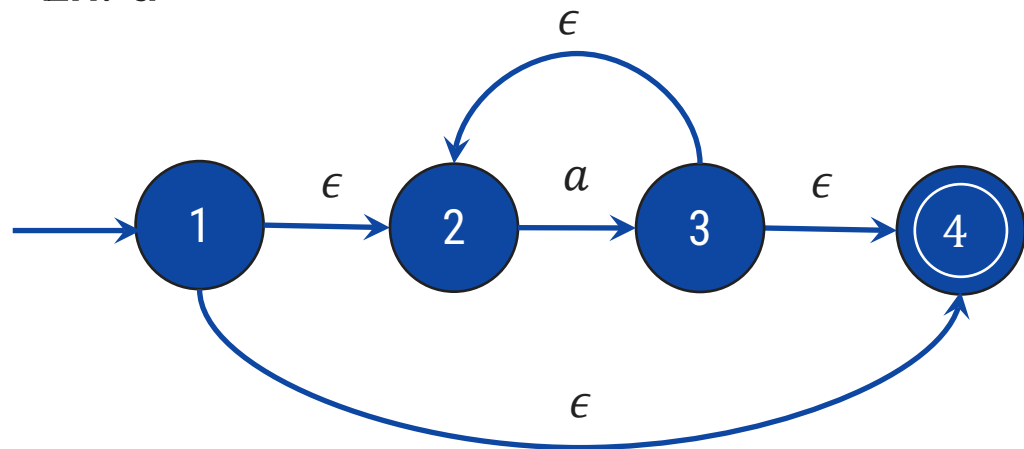
Ex:  $(a|b)$



5. For regular expression  $s^*$

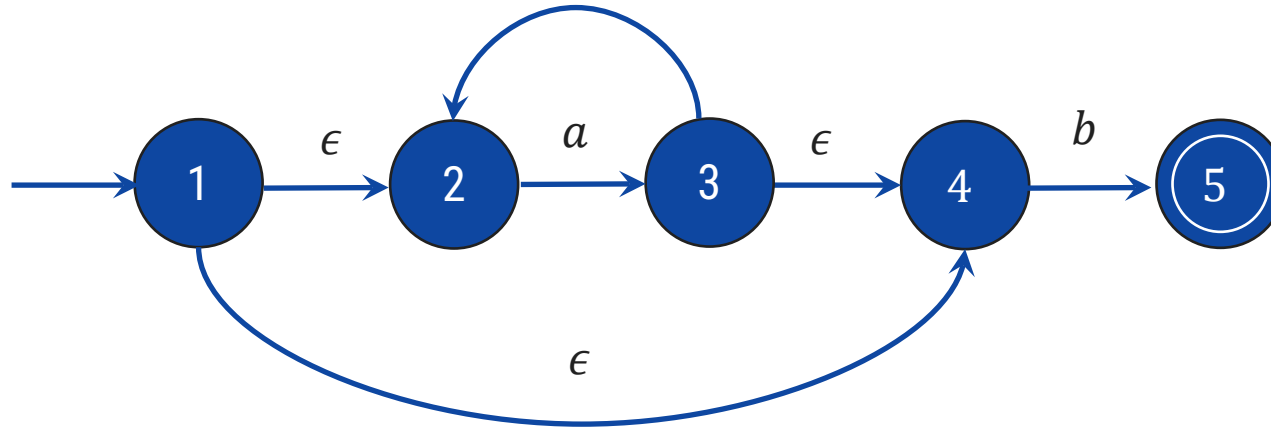


Ex:  $a^*$

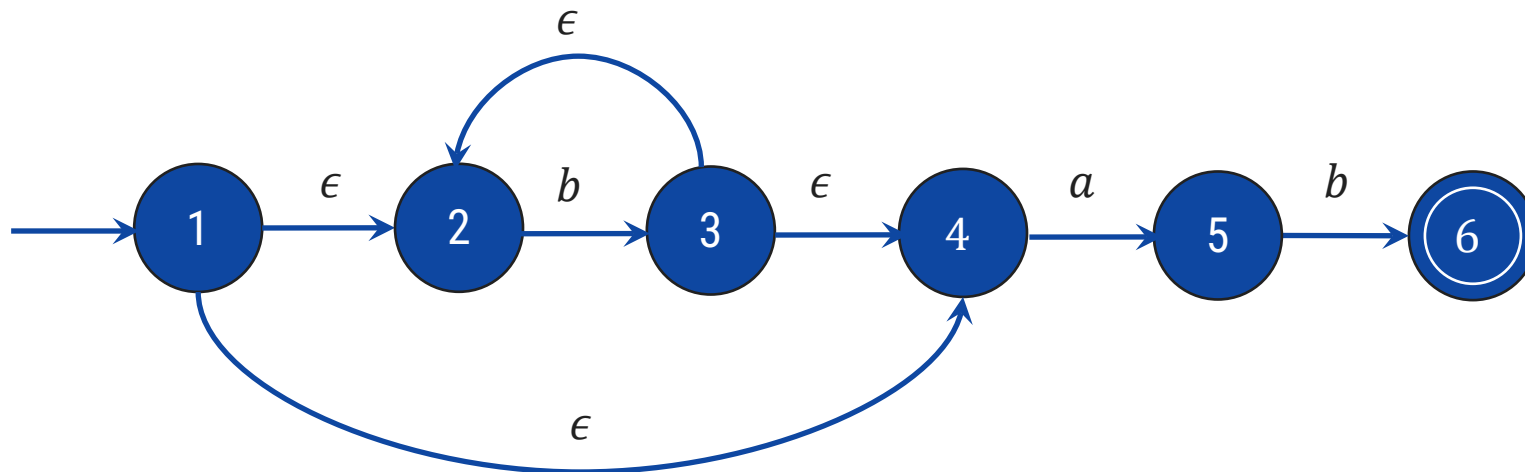


# Regular expression to NFA using Thompson's rule

►  $a^*b$



►  $b^*ab$



# Exercise

Convert following regular expression to NFA:

1. abba
2.  $bb(a)^*$
3.  $(a|b)^*$
4.  $a^* | b^*$
5.  $a(a)^*ab$
6.  $aa^*+ bb^*$
7.  $(a+b)^*abb$
8.  $10(0+1)^*1$
9.  $(a+b)^*a(a+b)$
10.  $(0+1)^*010(0+1)^*$
11.  $(010+00)^*(10)^*$
12.  $100(1)^*00(0+1)^*$

# **Conversion from NFA to DFA using subset construction method**

# Subset construction algorithm

**Input:** An NFA  $N$ .

**Output:** A DFA  $D$  accepting the same language.

**Method:** Algorithm construct a transition table  $Dtran$  for  $D$ . We use the following operation:

OPERATION	DESCRIPTION
$\epsilon - closure(s)$	Set of NFA states reachable from NFA state $s$ on $\epsilon -$ transition alone.
$\epsilon - closure(T)$	Set of NFA states reachable from some NFA state $s$ in $T$ on $\epsilon -$ transition alone.
$Move(T, a)$	Set of NFA states to which there is a transition on input symbol $a$ from some NFA state $s$ in $T$ .

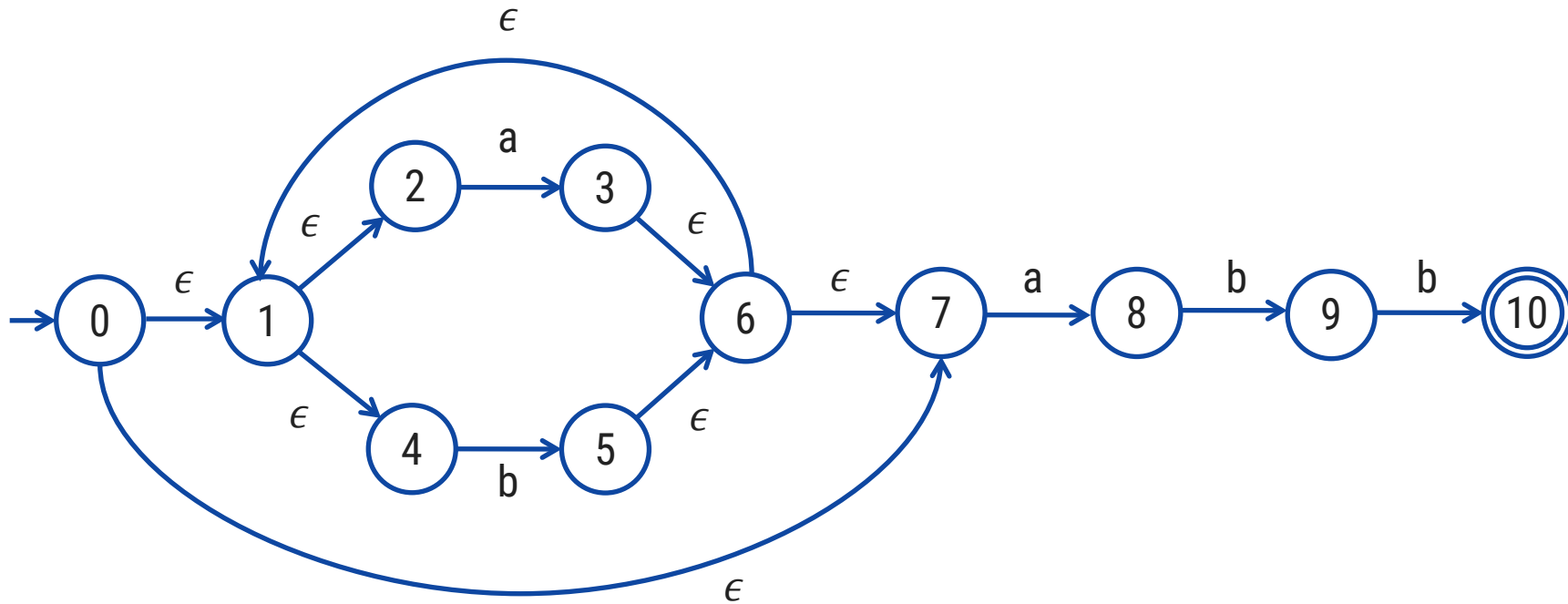
# Subset construction algorithm

```
initially  $\epsilon - closure(s_0)$  be the only state in  $Dstates$  and it is unmarked;  
while there is unmarked states  $T$  in  $Dstates$  do begin  
    mark  $T$ ;  
    for each input symbol  $a$  do begin  
         $U = \epsilon - closure(move(T, a))$ ;  
        if  $U$  is not in  $Dstates$  then  
            add  $U$  as unmarked state to  $Dstates$ ;  
         $Dtran[T, a] = U$   
    end  
end
```

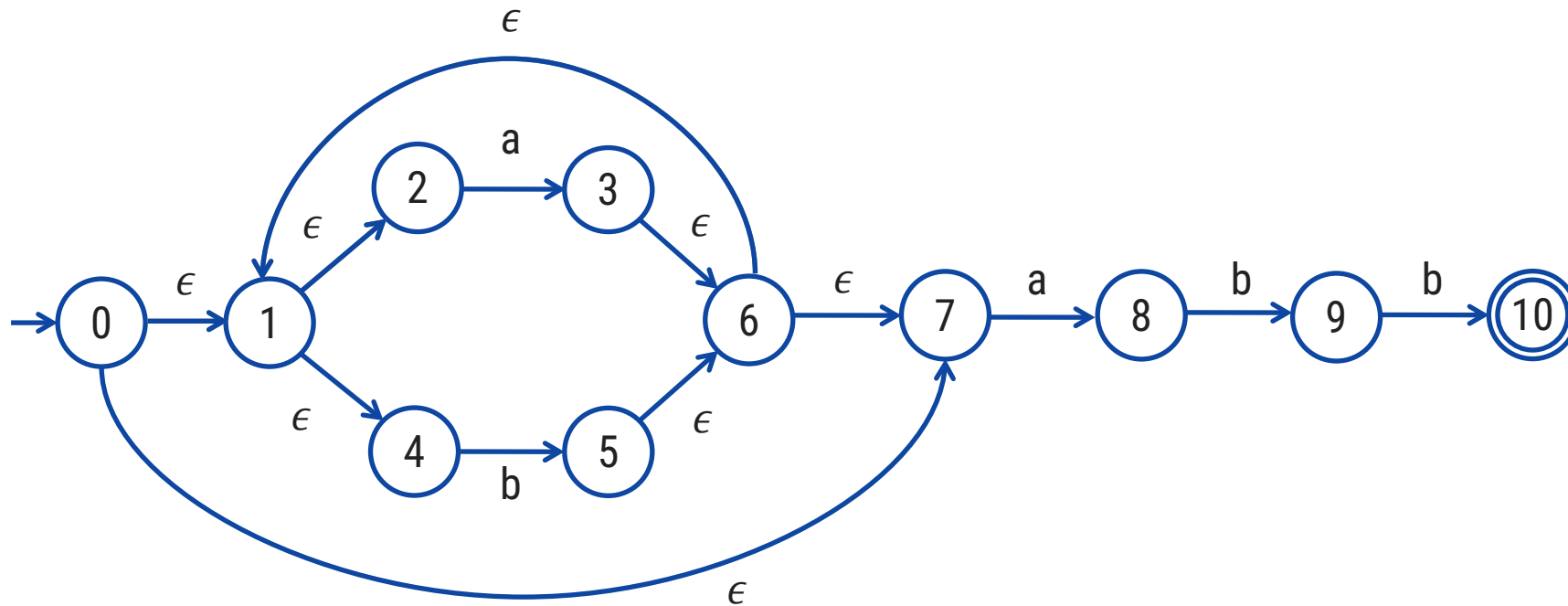


# Conversion from NFA to DFA

$(a|b)^*abb$



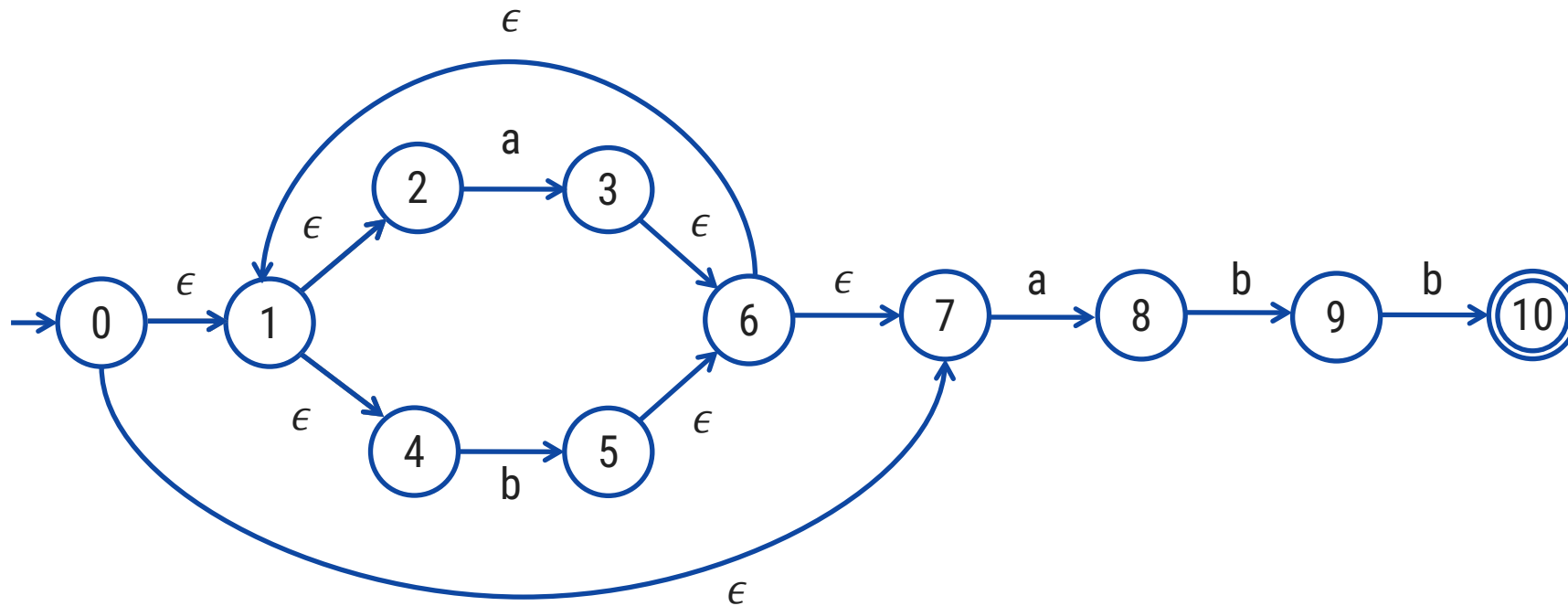
# Conversion from NFA to DFA



$\epsilon$ - Closure(0)=

= {0,1,2,4,7} ---- **A**

# Conversion from NFA to DFA



**$A = \{0, 1, \underline{2}, 4, \underline{7}\}$**

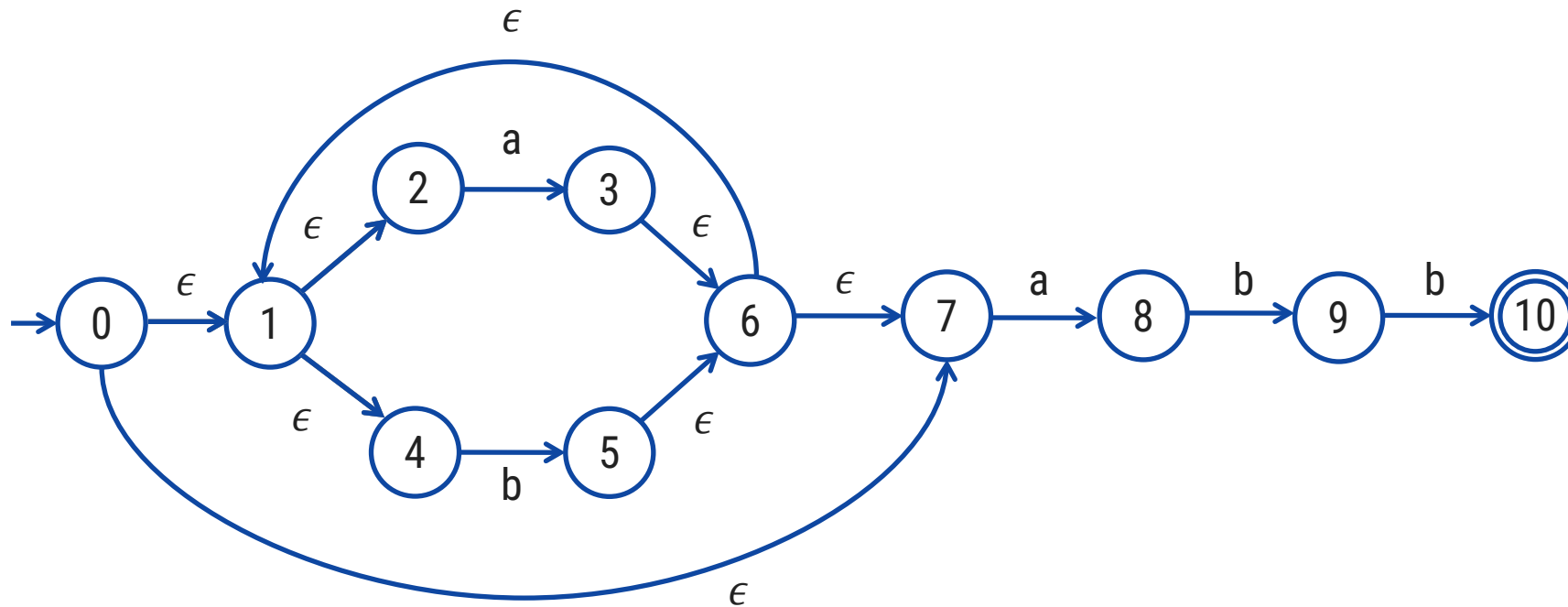
$\text{Move}(A, a) = \{3, 8\}$

$\epsilon$ - Closure( $\text{Move}(A, a)$ )

$= \{1, 2, 3, 4, 6, 7, 8\}$  ----- **B**

States	a	b
<b>A</b> = {0,1,2,4,7}		
<b>B</b> = {1,2,3,4,6,7,8}		

# Conversion from NFA to DFA



**A = {0, 1, 2, 4, 7}**

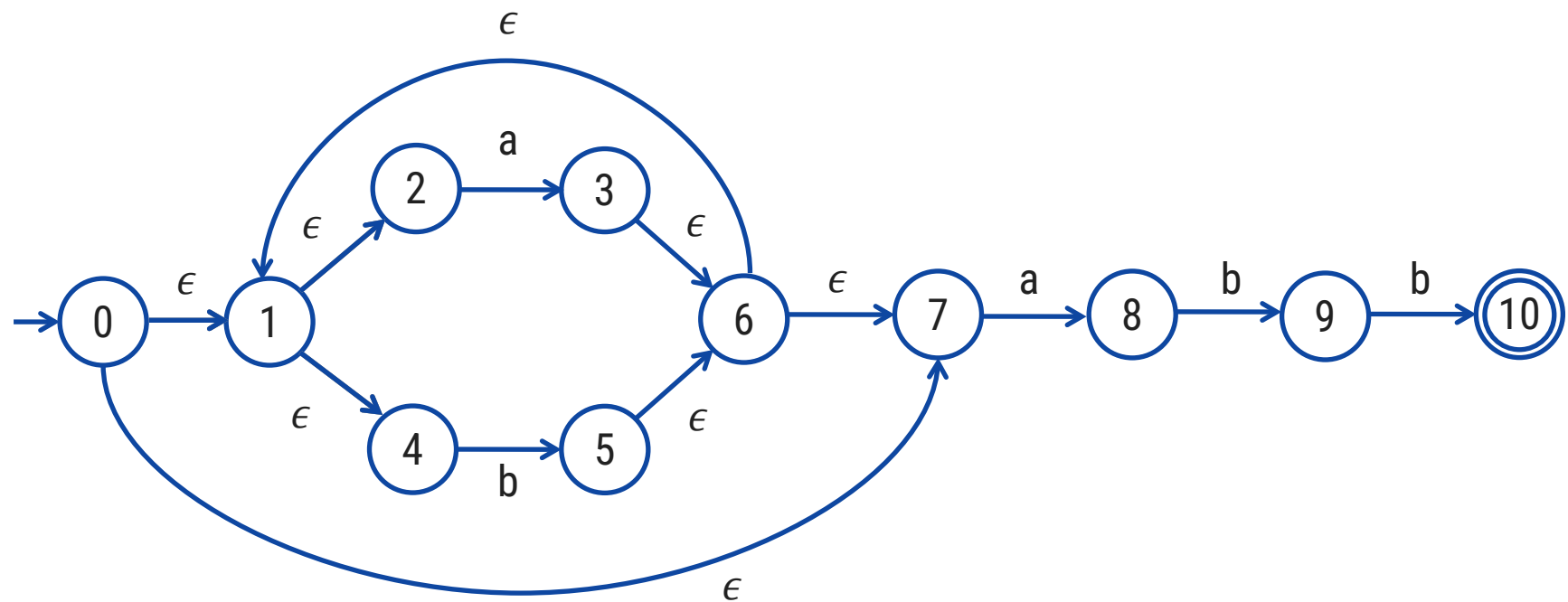
Move(A,b) =

$\epsilon$ - Closure(Move(A,b)) =

= {1,2,4,5,6,7} ---- **C**

States	a	b
<b>A = {0,1,2,4,7}</b>	<b>B</b>	
<b>B = {1,2,3,4,6,7,8}</b>		
<b>C = {1,2,4,5,6,7}</b>		

# Conversion from NFA to DFA



$B = \{1, \underline{2}, 3, 4, 6, \underline{7}, 8\}$

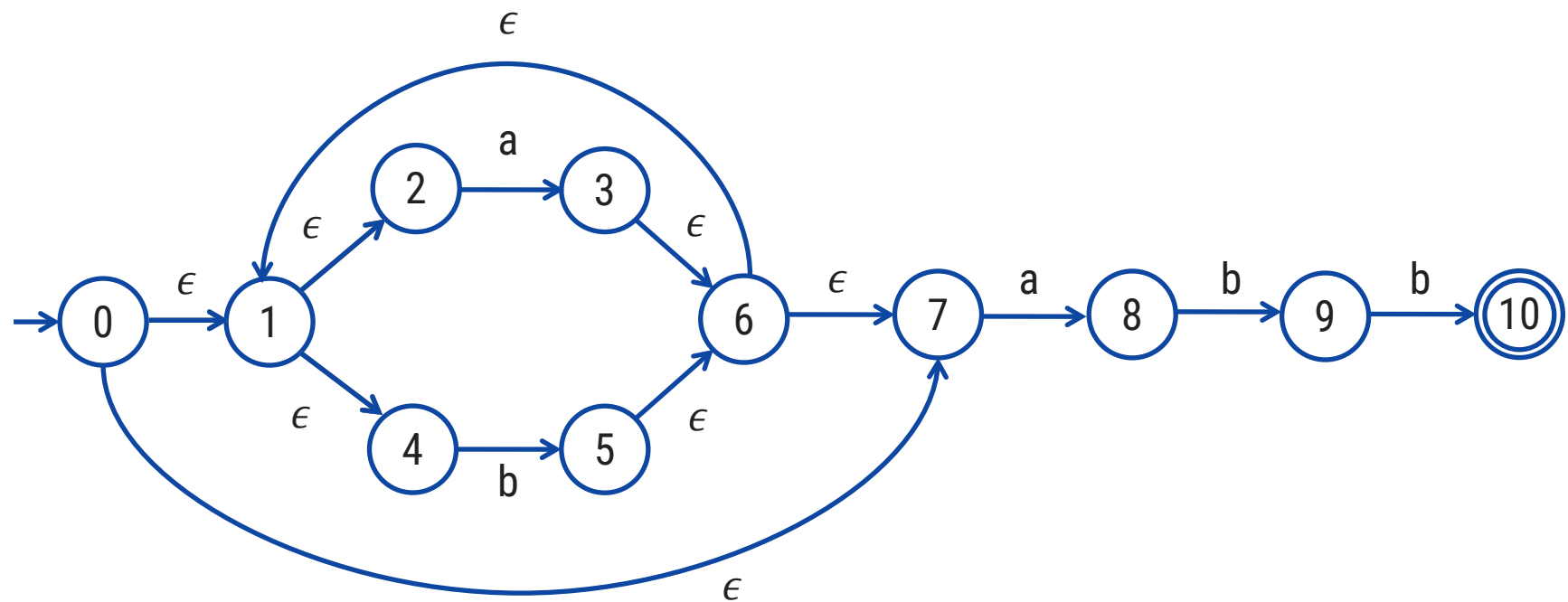
$\text{Move}(B,a) = \{3,8\}$

$\epsilon\text{-Closure}(\text{Move}(B,a))$

$= \{1,2,3,4,6,7,8\} \text{ ---- } B$

States	a	b
A = {0,1,2,4,7}	B	C
B = {1,2,3,4,6,7,8}		
C = {1,2,4,5,6,7}		

# Conversion from NFA to DFA



**B = {1, 2, 3, 4, 6, 7, 8}**

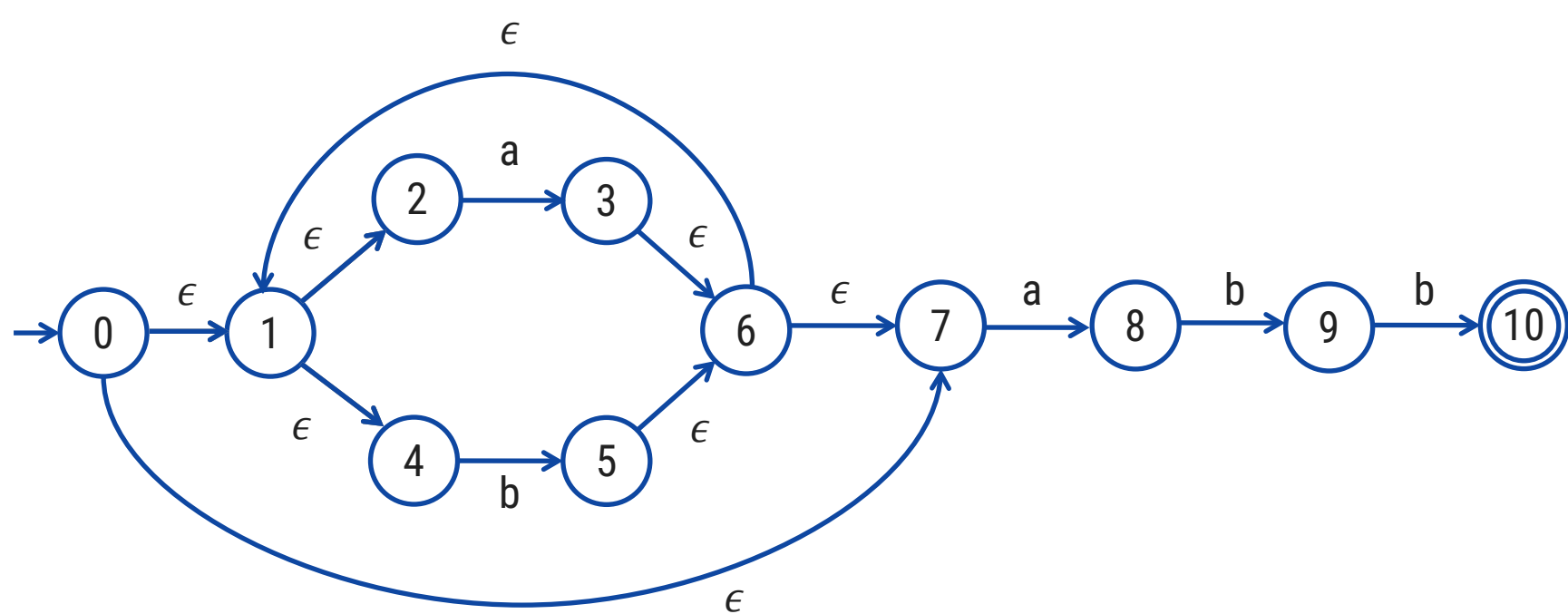
Move(B,b) = {5,9}

$\epsilon$ - Closure(Move(B,b))

= {1,2,4,5,6,7,9} ----- **D**

States	a	b
A = {0,1,2,4,7}	B	C
B = {1,2,3,4,6,7,8}	B	
C = {1,2,4,5,6,7}		
D = {1,2,4,5,6,7,9}		

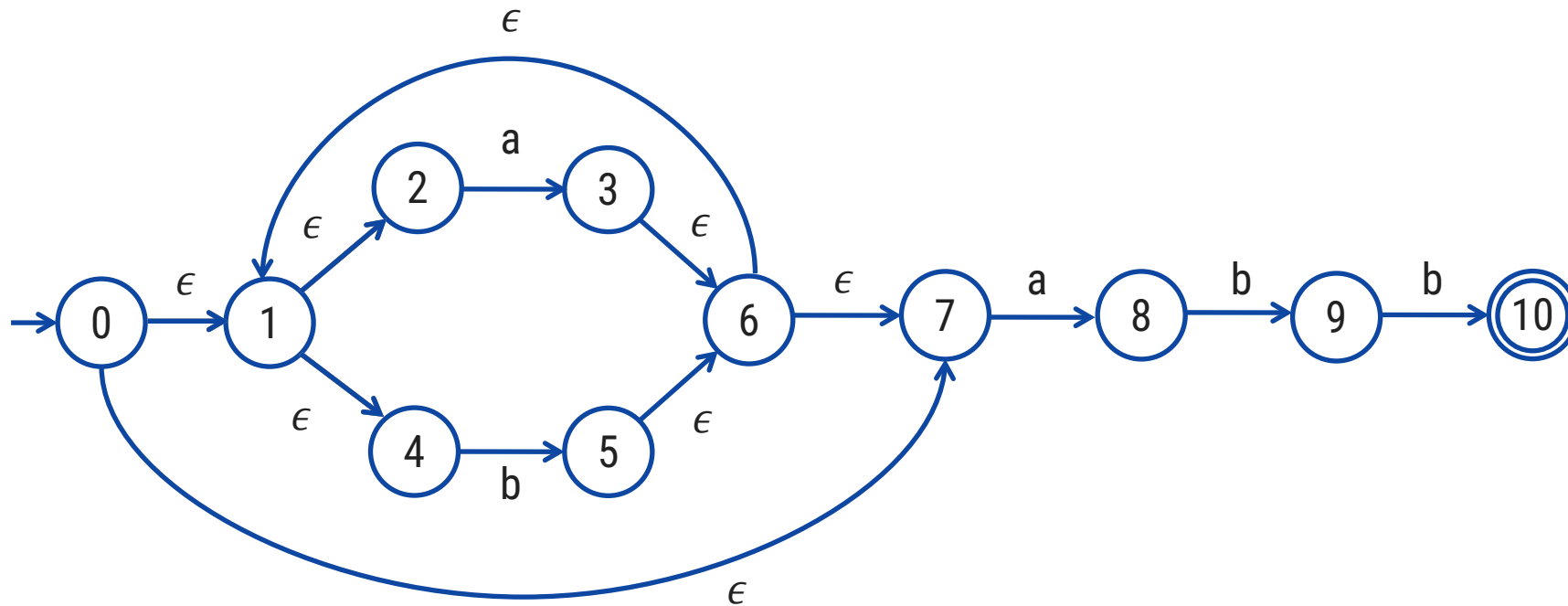
# Conversion from NFA to DFA



**C = {1, 2, 4, 5, 6, 7}**  
Move(C,a) = {3,8}  
 $\epsilon$ - Closure(Move(C,a))  
= {1,2,3,4,6,7,8} ----- **B**

States	a	b
A = {0,1,2,4,7}	B	C
B = {1,2,3,4,6,7,8}	B	D
C = {1,2,4,5,6,7}		
D = {1,2,4,5,6,7,9}		

# Conversion from NFA to DFA



**C = {1, 2, 4, 5, 6, 7}**

Move(C,b) =

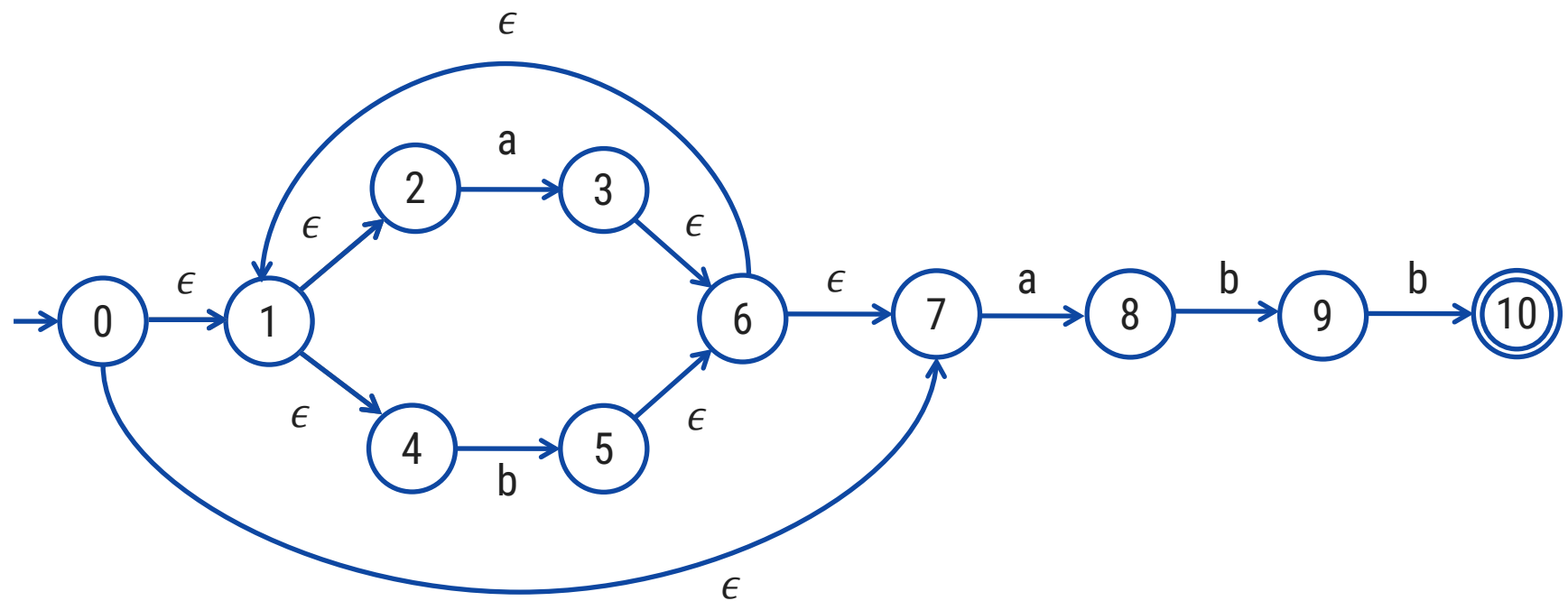
$\epsilon$ - Closure(Move(C,b))=

= {1,2,4,5,6,7} ---- **C**

States	a	b
A = {0,1,2,4,7}	B	C
B = {1,2,3,4,6,7,8}	B	D
C = {1,2,4,5,6,7}	B	
D = {1,2,4,5,6,7,9}		



# Conversion from NFA to DFA



$D = \{1, \underline{2}, 4, 5, 6, \underline{7}, 9\}$

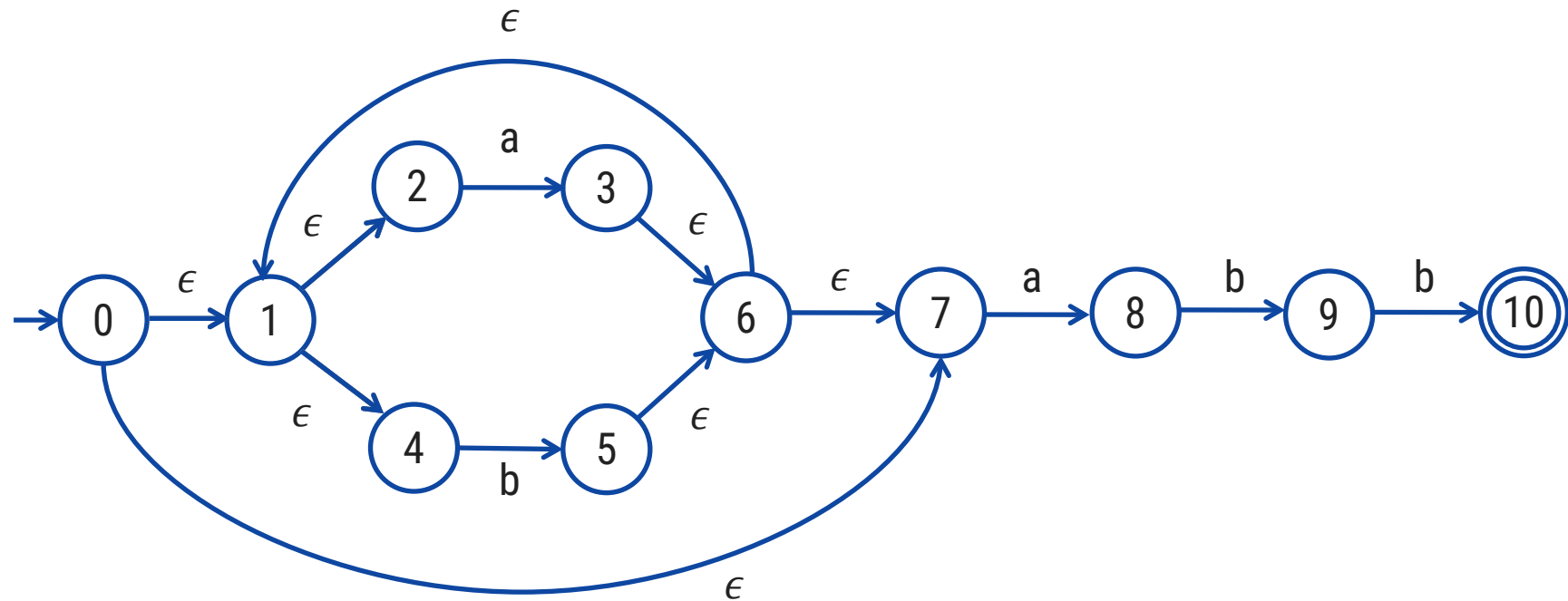
$\text{Move}(D, a) = \{3, 8\}$

$\epsilon\text{-Closure}(\text{Move}(D, a))$

$= \{1, 2, 3, 4, 6, 7, 8\}$  ----- **B**

States	a	b
$A = \{0, 1, 2, 4, 7\}$	<b>B</b>	<b>C</b>
$B = \{1, 2, 3, 4, 6, 7, 8\}$	<b>B</b>	<b>D</b>
$C = \{1, 2, 4, 5, 6, 7\}$	<b>B</b>	<b>C</b>
$D = \{1, 2, 4, 5, 6, 7, 9\}$		

# Conversion from NFA to DFA



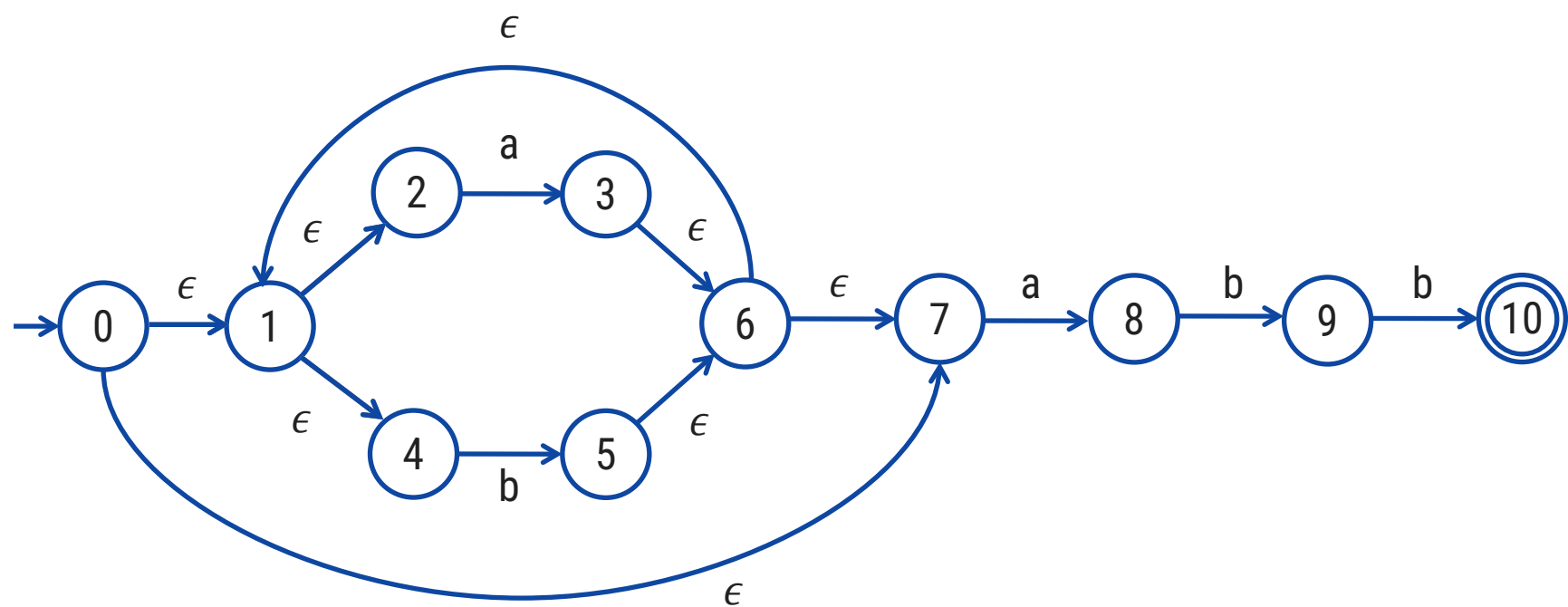
$D = \{1, 2, \underline{4}, 5, 6, 7, \underline{9}\}$

$\text{Move}(D, b) = \{5, 10\}$

$\epsilon\text{-Closure}(\text{Move}(D, b))$   
 $= \{1, 2, 4, 5, 6, 7, 10\}$  ----- **E**

States	a	b
$A = \{0, 1, 2, 4, 7\}$	B	C
$B = \{1, 2, 3, 4, 6, 7, 8\}$	B	D
$C = \{1, 2, 4, 5, 6, 7\}$	B	C
$D = \{1, 2, 4, 5, 6, 7, 9\}$	B	
$E = \{1, 2, 4, 5, 6, 7, 10\}$		

# Conversion from NFA to DFA



$E = \{1, \underline{2}, 4, 5, 6, \underline{7}, 10\}$

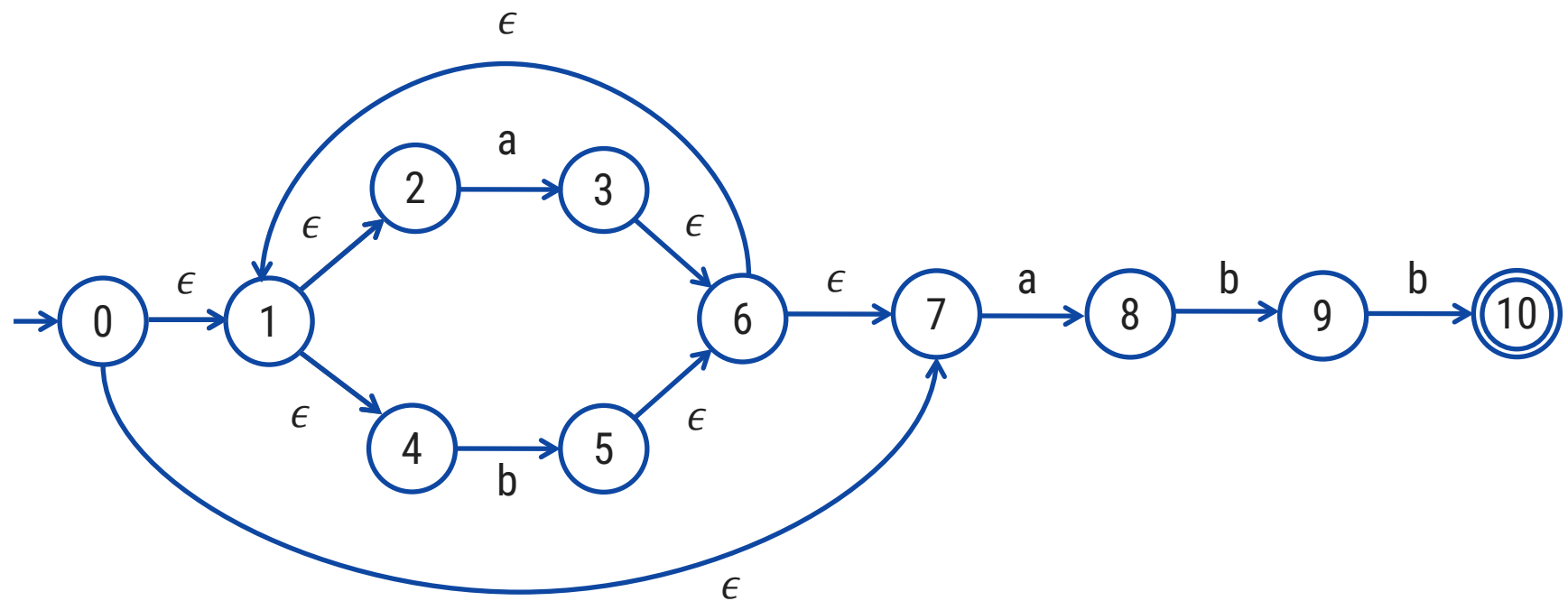
$\text{Move}(E, a) = \{3, 8\}$

$\epsilon\text{-Closure}(\text{Move}(E, a))$

$= \{1, 2, 3, 4, 6, 7, 8\} \text{ ---- B}$

States	a	b
A = {0,1,2,4,7}	B	C
B = {1,2,3,4,6,7,8}	B	D
C = {1,2,4,5,6,7}	B	C
D = {1,2,4,5,6,7,9}	B	E
E = {1,2,4,5,6,7,10}		

# Conversion from NFA to DFA



$E = \{1, 2, \underline{4}, 5, 6, 7, 10\}$

$\text{Move}(E, b) =$

$\epsilon\text{-Closure}(\text{Move}(E, b)) =$

$= \{1, 2, 4, 5, 6, 7\} \text{ ---- C}$

States	a	b
A = {0,1,2,4,7}	B	C
B = {1,2,3,4,6,7,8}	B	D
C = {1,2,4,5,6,7}	B	C
D = {1,2,4,5,6,7,9}	B	E
E = {1,2,4,5,6,7,10}	B	

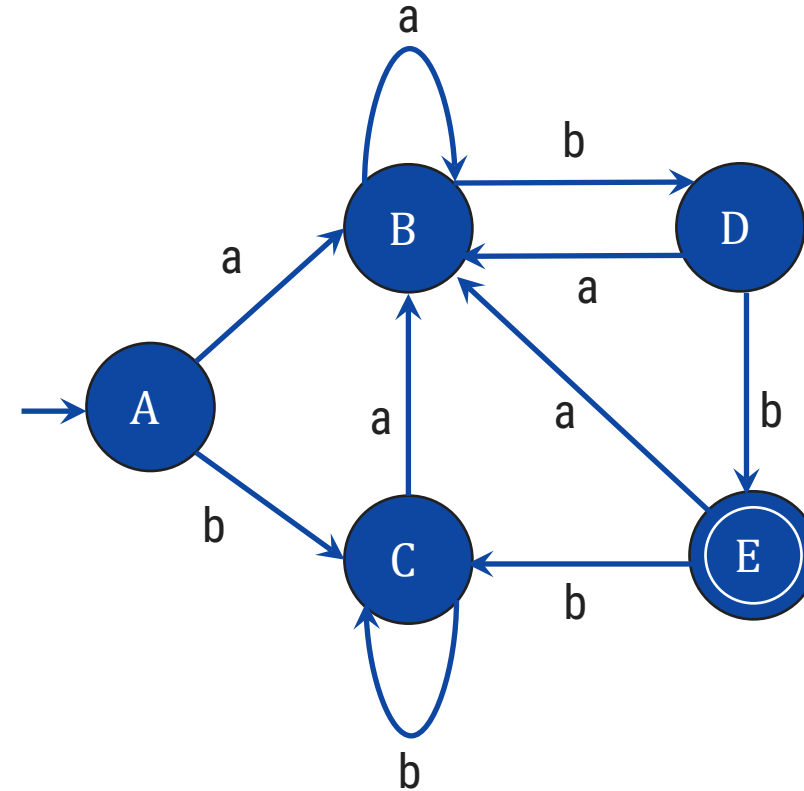
# Conversion from NFA to DFA

States	a	b
A = {0,1,2,4,7}	B	C
B = {1,2,3,4,6,7,8}	B	D
C = {1,2,4,5,6,7}	B	C
D = {1,2,4,5,6,7,9}	B	E
E = {1,2,4,5,6,7,10}	B	C

**Transition Table**

**Note:**

- **Accepting state in NFA is 10**
- **10 is element of E**
- **So, E is acceptance state in DFA**



**DFA**

# Exercise

- ▶ Convert following regular expression to DFA using subset construction method:
  1.  $(a+b)^*a(a+b)$
  2.  $(a+b)^*ab^*a$

# DFA optimization

# DFA optimization

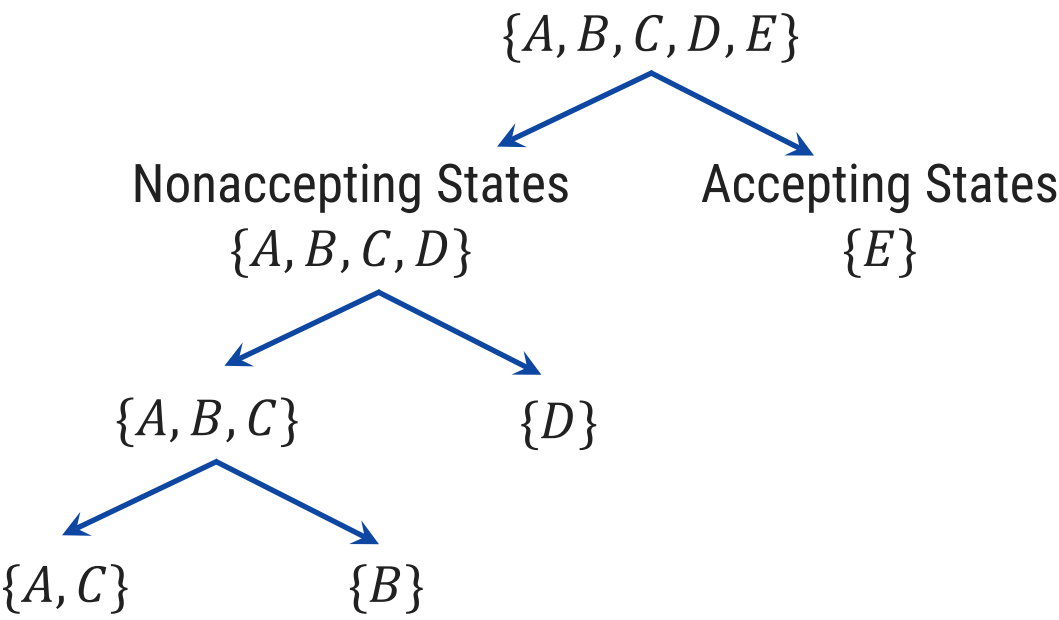
1. Construct an initial partition  $\Pi$  of the set of states with two groups: the accepting states  $F$  and the non-accepting states  $S - F$ .
2. Apply the repartition procedure to  $\Pi$  to construct a new partition  $\Pi_{new}$ .
3. If  $\Pi_{new} = \Pi$ , let  $\Pi_{final} = \Pi$  and continue with step (4). Otherwise, repeat step (2) with  $\Pi = \Pi_{new}$ .  
    **for** each group  $G$  of  $\Pi$  **do begin**  
        partition  $G$  into subgroups such that two states  $s$  and  $t$   
        of  $G$  are in the same subgroup if and only if for all  
        input symbols  $a$ , states  $s$  and  $t$  have transitions on  $a$   
        to states in the same group of  $\Pi$ .  
        replace  $G$  in  $\Pi_{new}$  by the set of all subgroups formed.  
    **end**



# DFA optimization

4. Choose one state in each group of the partition  $\Pi_{final}$  as the representative for that group. The representatives will be the states of  $M'$ . Let  $s$  be a representative state, and suppose on input  $a$  there is a transition of  $M$  from  $s$  to  $t$ . Let  $r$  be the representative of  $t$ 's group. Then  $M'$  has a transition from  $s$  to  $r$  on  $a$ . Let the start state of  $M'$  be the representative of the group containing start state  $s_0$  of  $M$ , and let the accepting states of  $M'$  be the representatives that are in  $F$ .
5. If  $M'$  has a dead state  $d$ , then remove  $d$  from  $M'$ . Also remove any state not reachable from the start state.

# DFA optimization



- ▶ Now no more splitting is possible.
- ▶ If we chose A as the representative for group (AC), then we obtain reduced transition table

States	a	b
A	B	C
B	B	D
C	B	C
D	B	E
E	B	C

States	a	b
A	B	A
B	B	D
D	B	E
E	B	A

**Optimized  
Transition Table**

# Conversion from regular expression to DFA

# Rules to compute nullable, firstpos, lastpos

## ▶ nullable( $n$ )

↪ The subtree at node  $n$  generates languages including the empty string.

## ▶ firstpos( $n$ )

↪ The set of positions that can match the first symbol of a string generated by the subtree at node  $n$ .

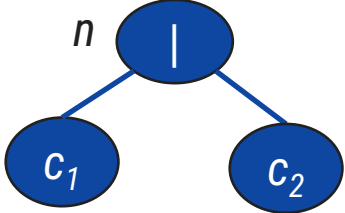
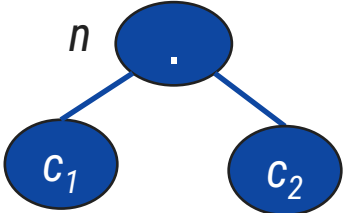
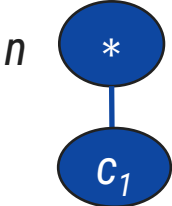
## ▶ lastpos( $n$ )

↪ The set of positions that can match the last symbol of a string generated by the subtree at node  $n$ .

## ▶ followpos( $i$ )

↪ The set of positions that can follow position  $i$  in the tree.

# Rules to compute nullable, firstpos, lastpos

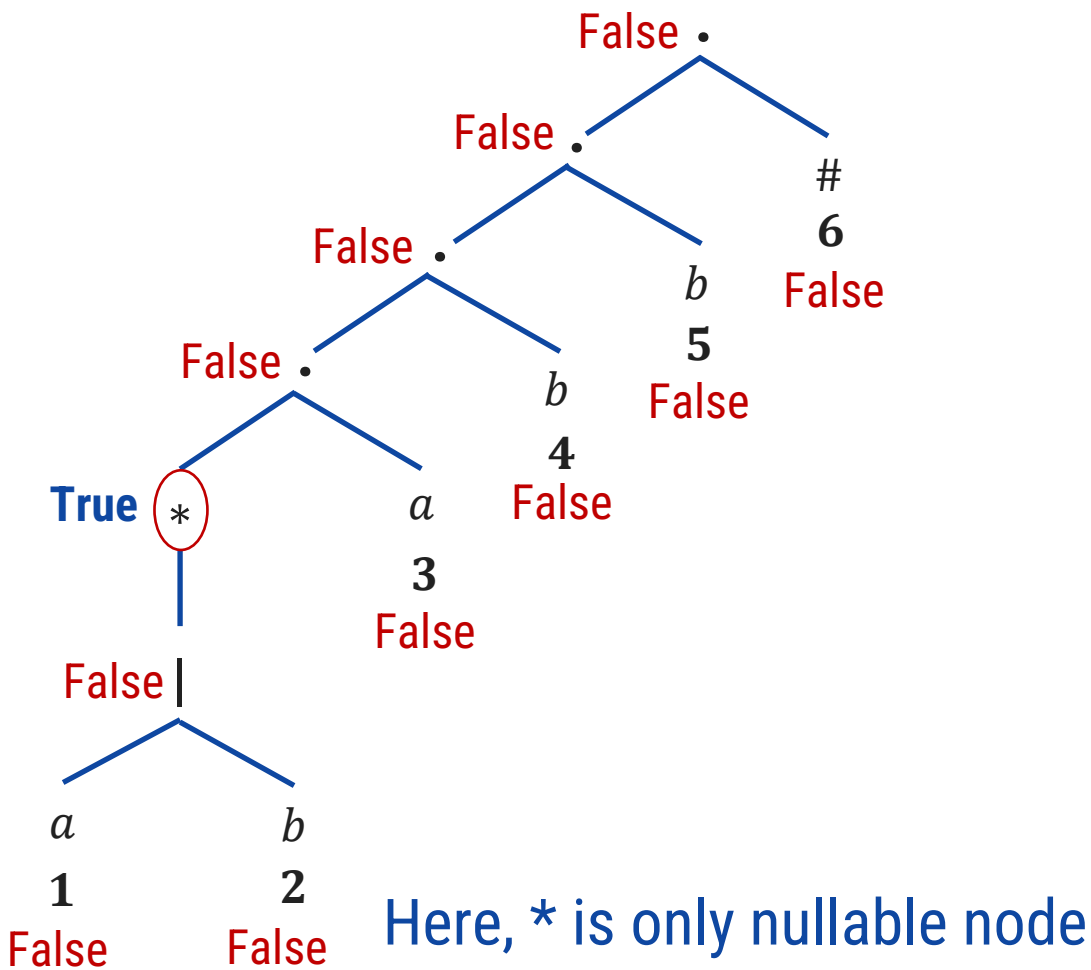
Node n	nullable(n)	firstpos(n)	lastpos(n)
A leaf labeled by $\epsilon$	<b>true</b>	$\emptyset$	$\emptyset$
A leaf with position <b>i</b>	<b>false</b>	$\{i\}$	$\{i\}$
	nullable( $c_1$ ) <b>or</b> nullable( $c_2$ )	firstpos( $c_1$ ) $\cup$ firstpos( $c_2$ )	lastpos( $c_1$ ) $\cup$ lastpos( $c_2$ )
	nullable( $c_1$ ) <b>and</b> nullable( $c_2$ )	<b>if</b> (nullable( $c_1$ )) <b>then</b> firstpos( $c_1$ ) $\cup$ firstpos( $c_2$ ) <b>else</b> firstpos( $c_1$ )	<b>if</b> (nullable( $c_2$ )) <b>then</b> lastpos( $c_1$ ) $\cup$ lastpos( $c_2$ ) <b>else</b> lastpos( $c_2$ )
	<b>true</b>	firstpos( $c_1$ )	lastpos( $c_1$ )

# Rules to compute followpos

1. If  $n$  is **concatenation** node with left child  $c1$  and right child  $c2$  and  $i$  is a position in  $\text{lastpos}(c1)$ , then all position in  $\text{firstpos}(c2)$  are in  $\text{followpos}(i)$
2. If  $n$  is  $*$  node and  $i$  is position in  $\text{lastpos}(n)$ , then all position in  $\text{firstpos}(n)$  are in  $\text{followpos}(i)$

# Conversion from regular expression to DFA

$(a|b)^*abb\#$

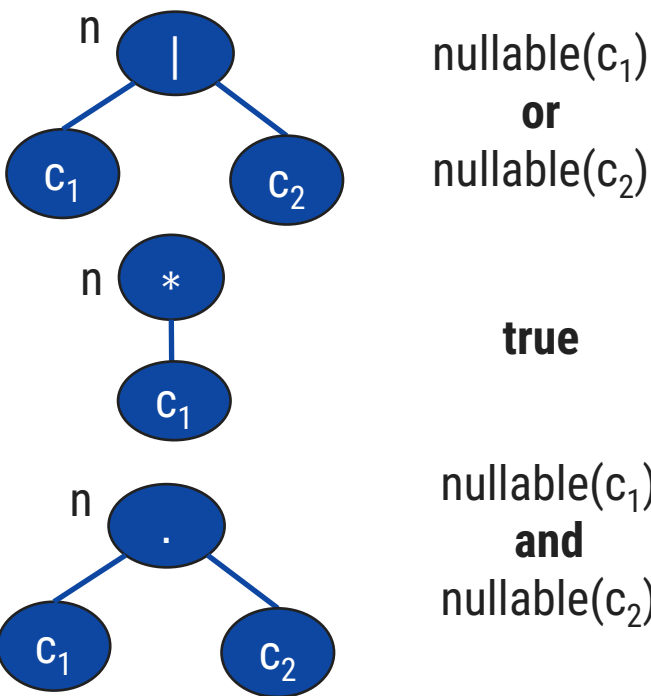


Step 1: Construct Syntax Tree

Step 2: Nullable node

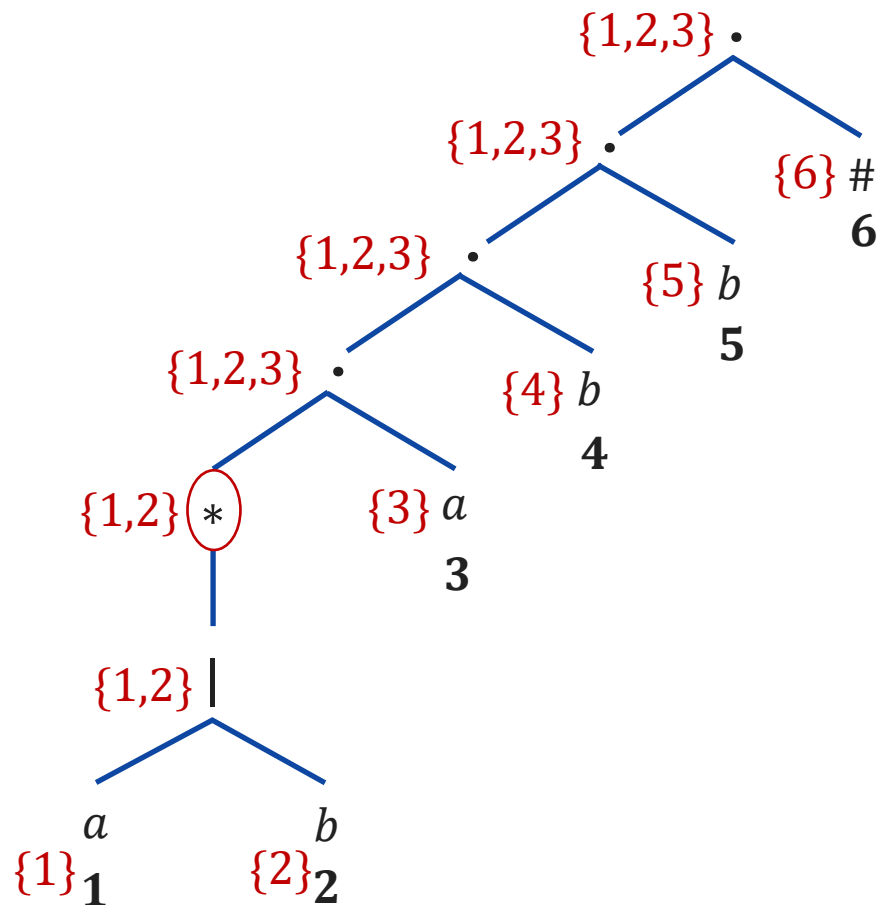
A leaf labeled by  $\epsilon = \text{True}$

A leaf with position  $i = \text{false}$



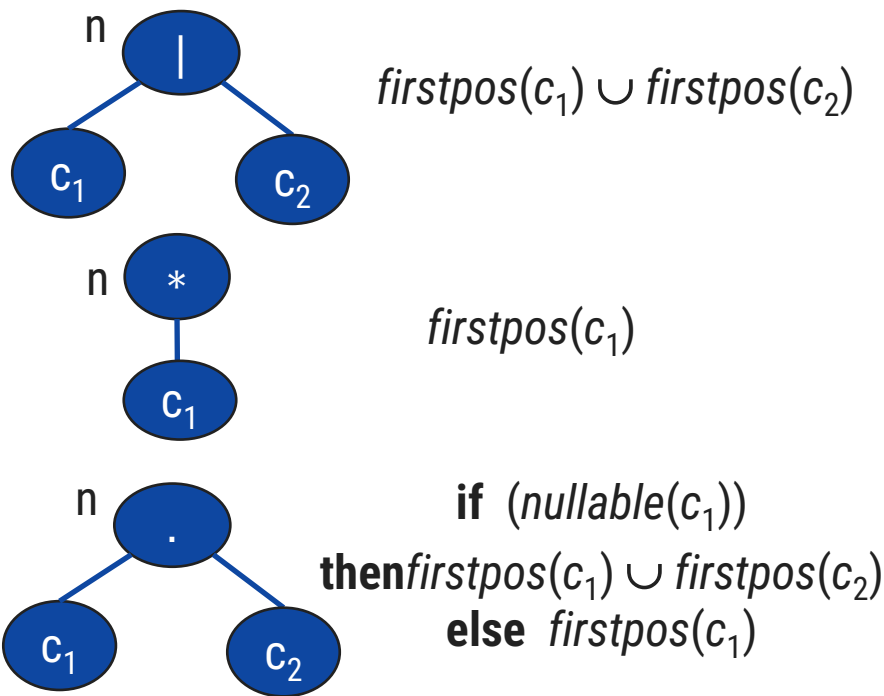
# Conversion from regular expression to DFA

## Step 3: Calculate firstpos



### Firstpos —

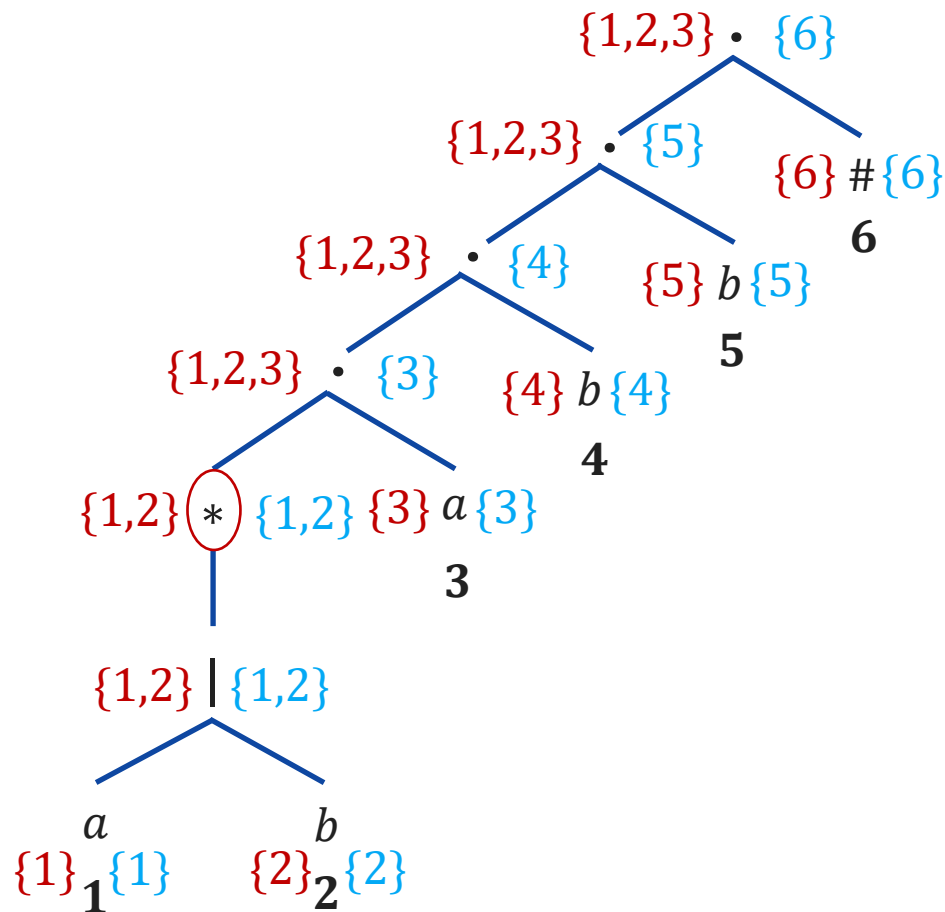
A leaf with position  $i = \{i\}$



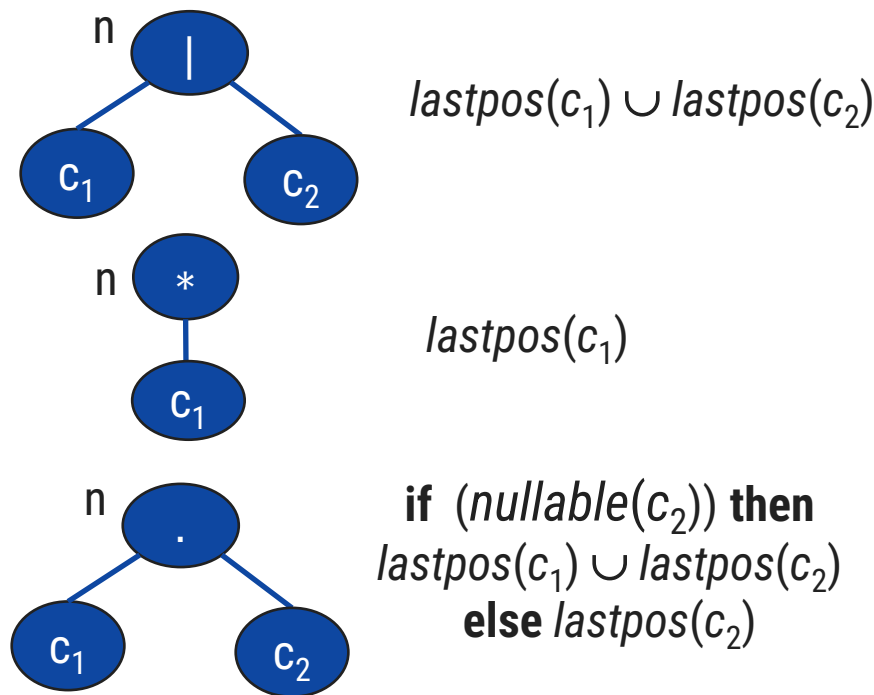


# Conversion from regular expression to DFA

## Step 3: Calculate lastpos



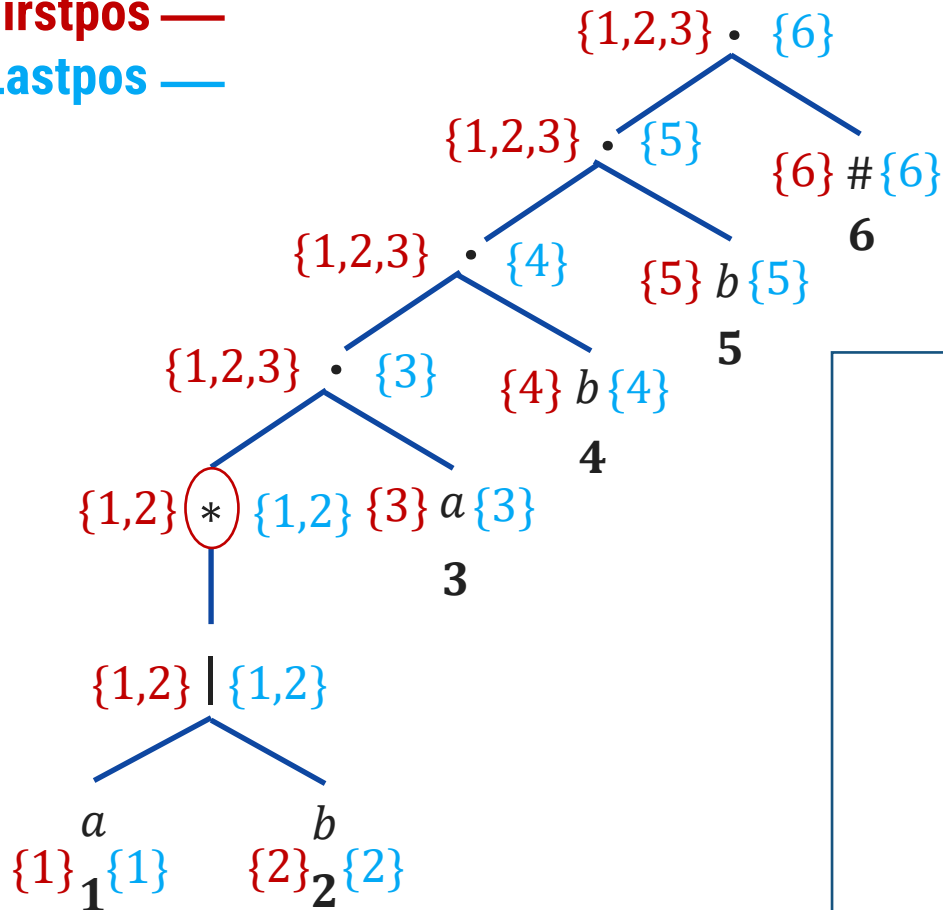
A leaf with position  $i = \{i\}$



# Conversion from regular expression to DFA

## Step 4: Calculate followpos

**Firstpos** —  
**Lastpos** —



Position	followpos
5	
4	
3	
2	1,2,
1	1,2,

### Rule:

If  $n$  is \* node and  $i$  is position in  $\text{lastpos}(n)$ , then all position in  $\text{firstpos}(n)$  are in  $\text{followpos}(i)$

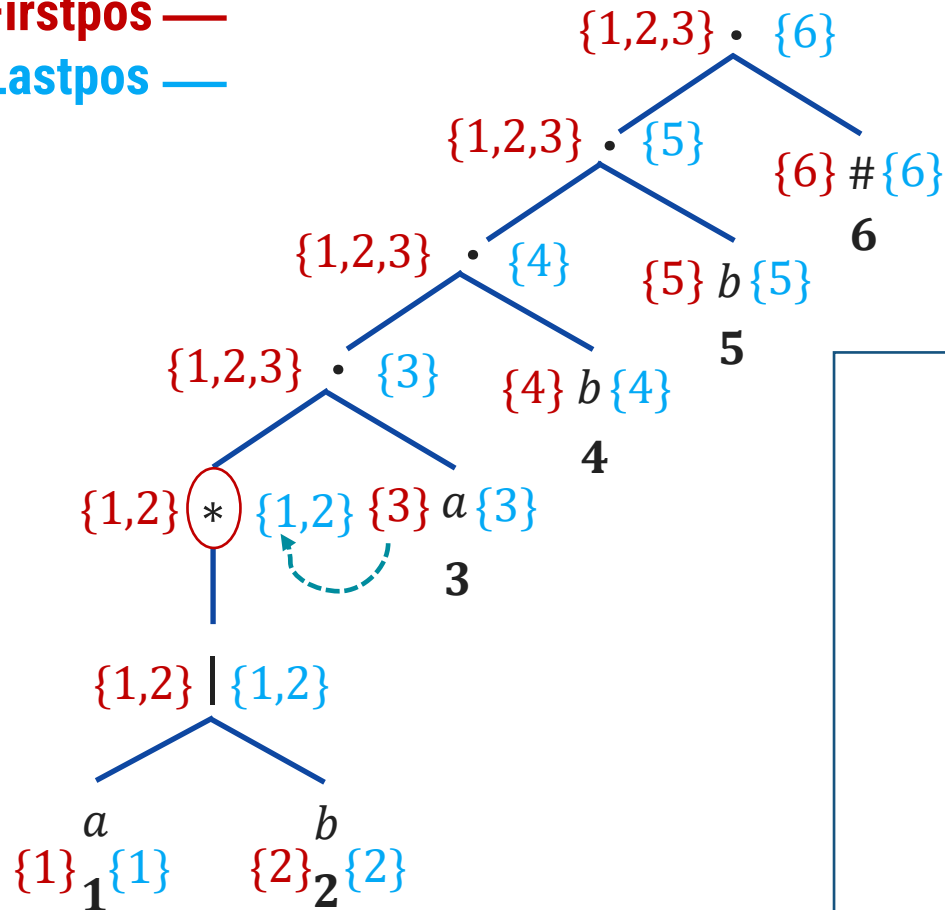
$$\{1,2\} \circledast \{1,2\}$$

$$\begin{aligned} i &= \text{lastpos}(n) = \{1,2\} \\ \text{firstpos}(n) &= \{1,2\} \\ \text{followpos}(1) &= \{1,2\} \\ \text{followpos}(2) &= \{1,2\} \end{aligned}$$

# Conversion from regular expression to DFA

## Step 4: Calculate followpos

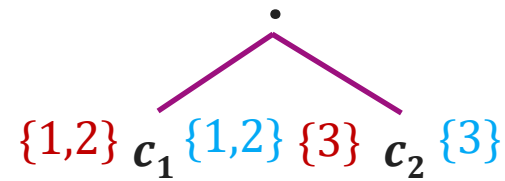
**Firstpos** —  
**Lastpos** —



Position	followpos
5	
4	
3	
2	1,2,3
1	1,2,3

### Rule:

If  $n$  is **concatenation** node  
with left child  $c_1$  and right  
child  $c_2$  and  $i$  is a position in  
 $\text{lastpos}(c_1)$ , then all position  
in  $\text{firstpos}(c_2)$  are in  
 $\text{followpos}(i)$

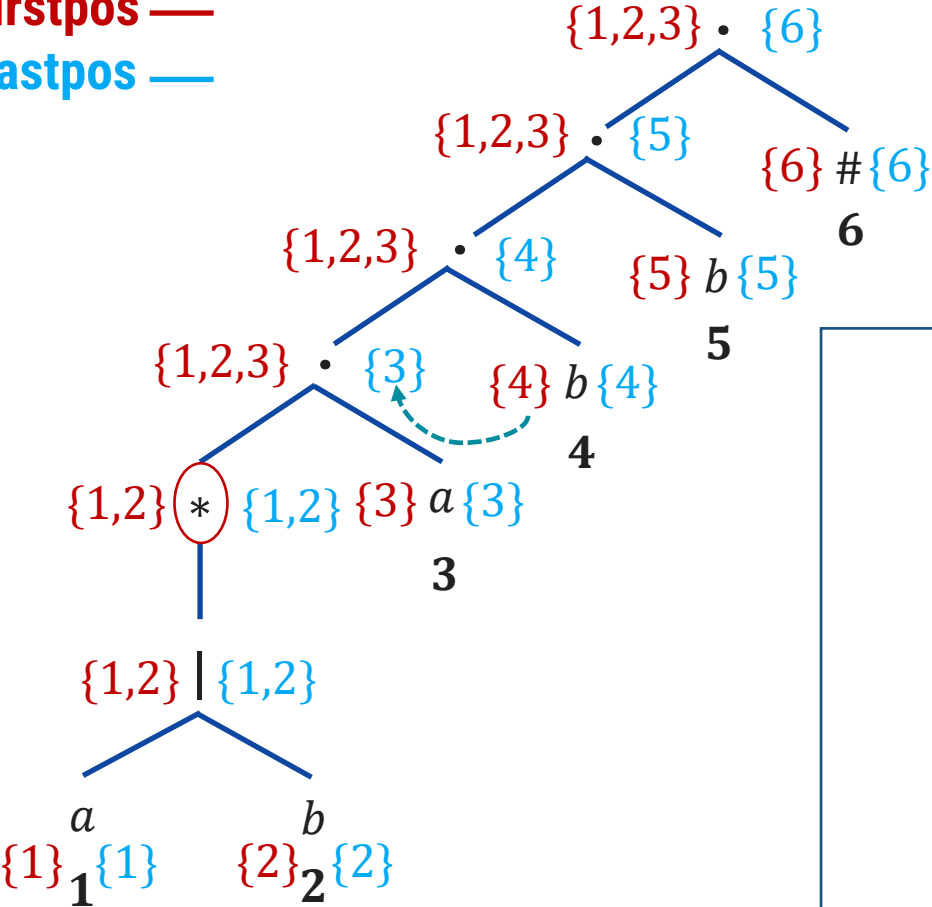


$$\begin{aligned} i &= \text{lastpos}(c_1) = \{1,2\} \\ \text{firstpos}(c_2) &= \{3\} \\ \text{followpos}(1) &= \{3\} \\ \text{followpos}(2) &= \{3\} \end{aligned}$$

# Conversion from regular expression to DFA

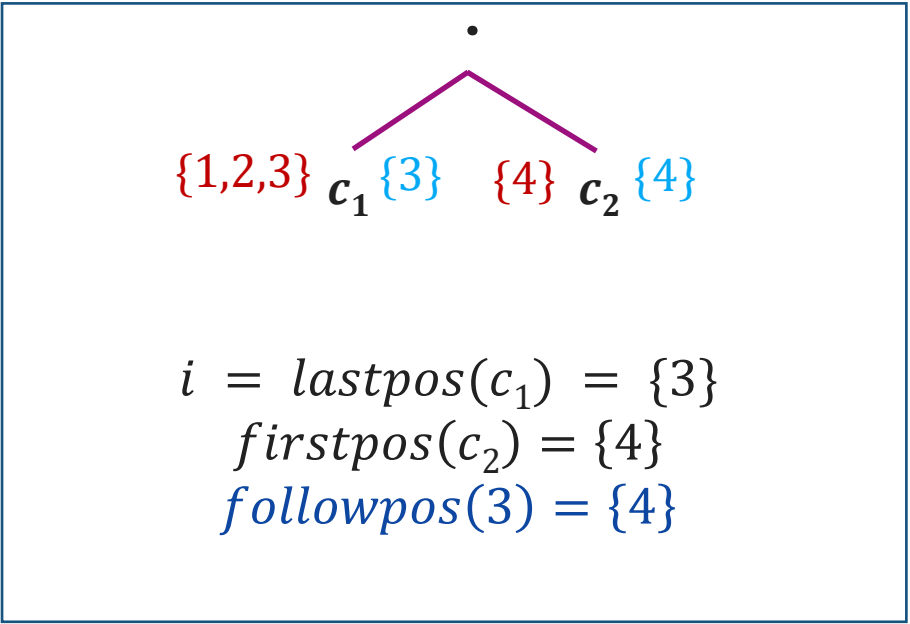
Step 4: Calculate followpos

**Firstpos** —  
**Lastpos** —



Position	followpos
5	
4	
3	4
2	1,2,3
1	1,2,3

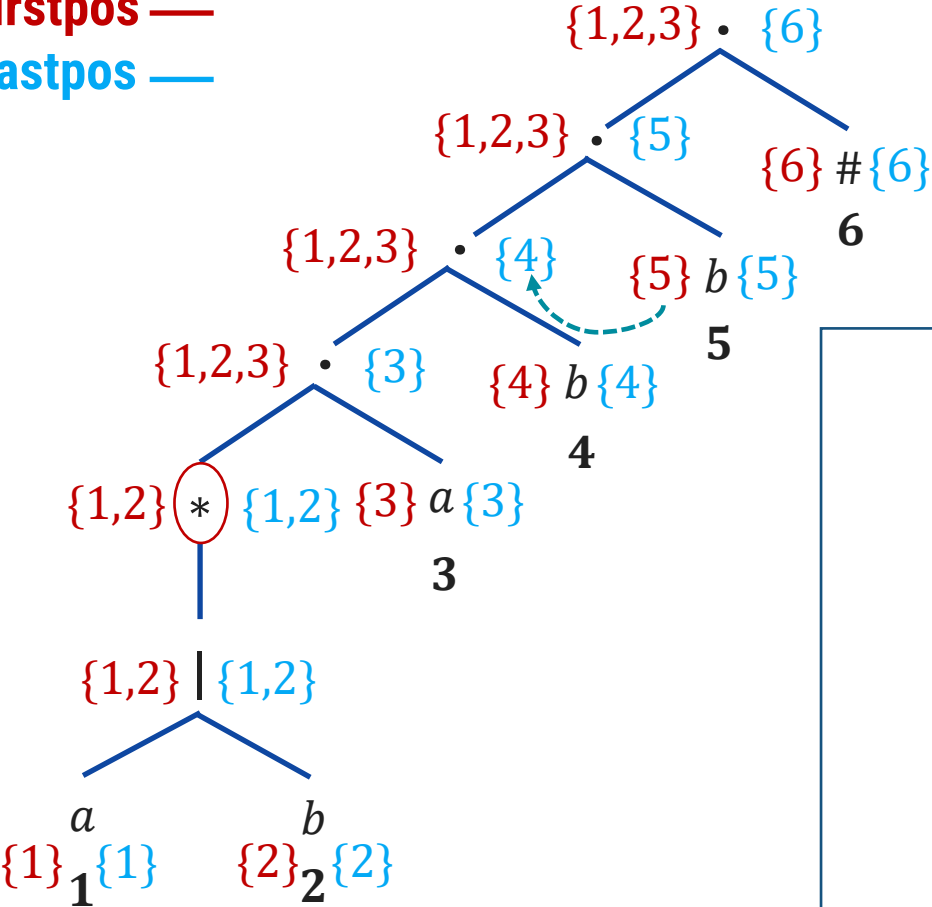
**Rule:**  
If  $n$  is **concatenation** node with left child  $c_1$  and right child  $c_2$  and  $i$  is a position in  $lastpos(c_1)$ , then all position in  $firstpos(c_2)$  are in  $followpos(i)$



# Conversion from regular expression to DFA

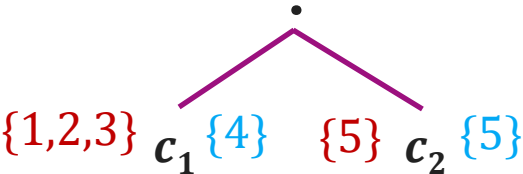
Step 4: Calculate followpos

**Firstpos** —  
**Lastpos** —



Position	followpos
5	
4	5
3	4
2	1,2,3
1	1,2,3

**Rule:**  
If  $n$  is **concatenation** node with left child  $c_1$  and right child  $c_2$  and  $i$  is a position in  $\text{lastpos}(c_1)$ , then all position in  $\text{firstpos}(c_2)$  are in  $\text{followpos}(i)$

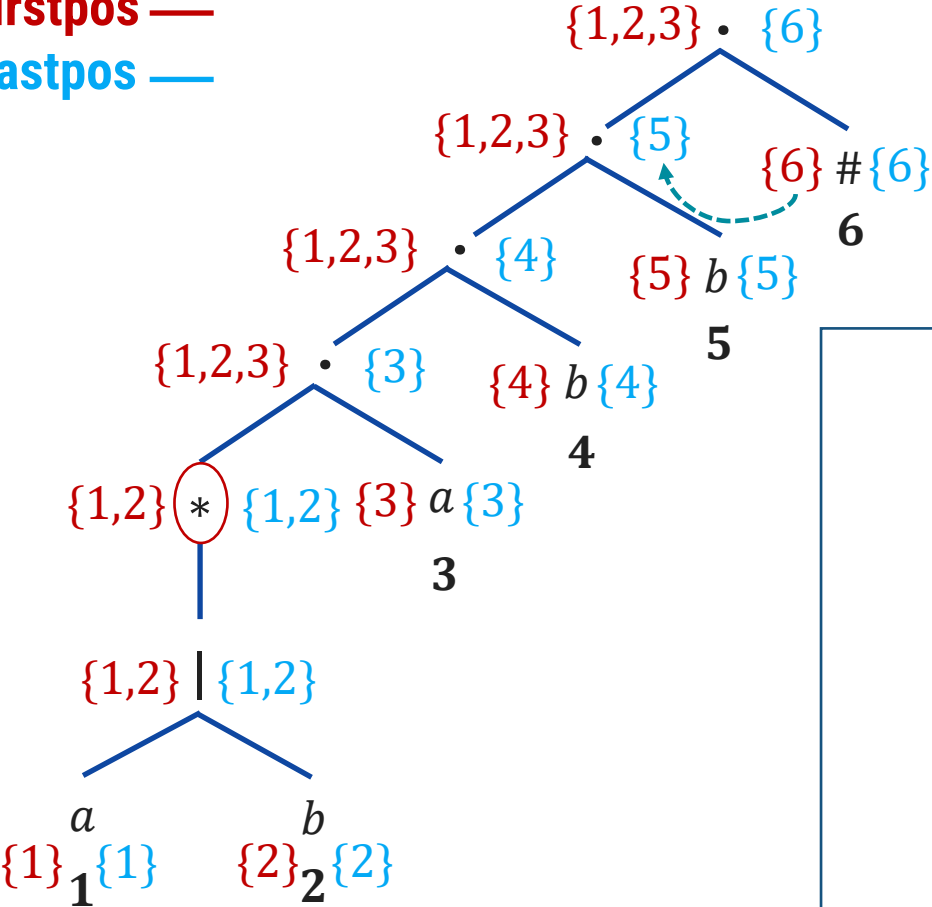


$$\begin{aligned} i &= \text{lastpos}(c_1) = \{4\} \\ \text{firstpos}(c_2) &= \{5\} \\ \text{followpos}(4) &= \{5\} \end{aligned}$$

# Conversion from regular expression to DFA

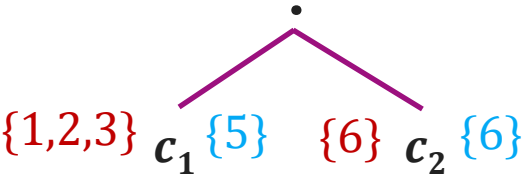
Step 4: Calculate followpos

**Firstpos** —  
**Lastpos** —



Position	followpos
5	6
4	5
3	4
2	1,2,3
1	1,2,3

**Rule:**  
If  $n$  is **concatenation** node with left child  $c_1$  and right child  $c_2$  and  $i$  is a position in  $\text{lastpos}(c_1)$ , then all position in  $\text{firstpos}(c_2)$  are in  $\text{followpos}(i)$



$$\begin{aligned} i &= \text{lastpos}(c_1) = \{5\} \\ \text{firstpos}(c_2) &= \{6\} \\ \text{followpos}(5) &= \{6\} \end{aligned}$$

# Conversion from regular expression to DFA

Initial state = *firstpos* of root = {1,2,3} ----- A

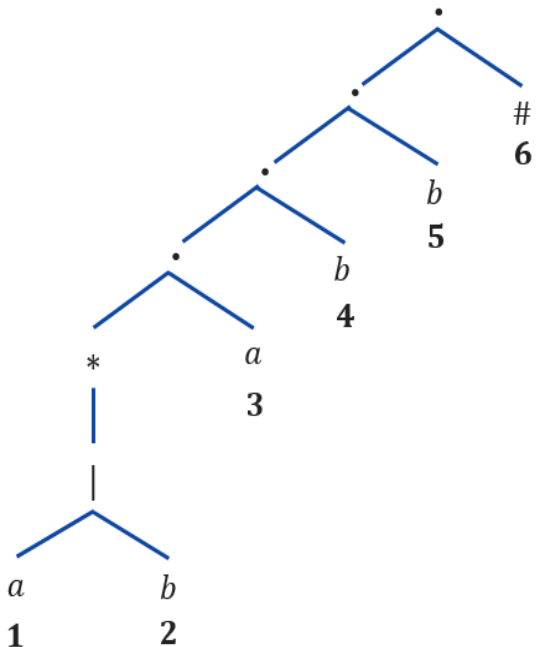
## State A

$$\begin{aligned} \delta((1,2,3),a) &= \text{followpos}(1) \cup \text{followpos}(3) \\ &= (1,2,3) \cup (4) = \{1,2,3,4\} \text{ ----- B} \end{aligned}$$

$$\begin{aligned} \delta((1,2,3),b) &= \text{followpos}(2) \\ &= (1,2,3) \text{ ----- A} \end{aligned}$$

Position	followpos
5	6
4	5
3	4
2	1,2,3
1	1,2,3

States	a	b
A={1,2,3}		
B={1,2,3,4}		



# Conversion from regular expression to DFA

## State B

$\delta((1,2,3,4),a) = \text{followpos}(1) \cup \text{followpos}(3)$   
 $= (1,2,3) \cup (4) = \{1,2,3,4\}$  ----- B

$\delta((1,2,3,4),b) = \text{followpos}(2) \cup \text{followpos}(4)$   
 $= (1,2,3) \cup (5) = \{1,2,3,5\}$  ----- C

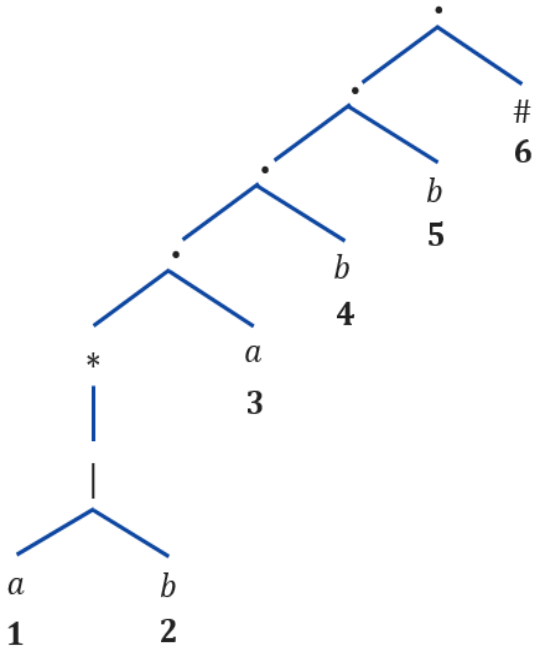
## State C

$\delta((1,2,3,5),a) = \text{followpos}(1) \cup \text{followpos}(3)$   
 $= (1,2,3) \cup (4) = \{1,2,3,4\}$  ----- B

$\delta((1,2,3,5),b) = \text{followpos}(2) \cup \text{followpos}(5)$   
 $= (1,2,3) \cup (6) = \{1,2,3,6\}$  ----- D

Position	followpos
5	6
4	5
3	4
2	1,2,3
1	1,2,3

States	a	b
A={1,2,3}	B	A
B={1,2,3,4}		
C={1,2,3,5}		
D={1,2,3,6}		



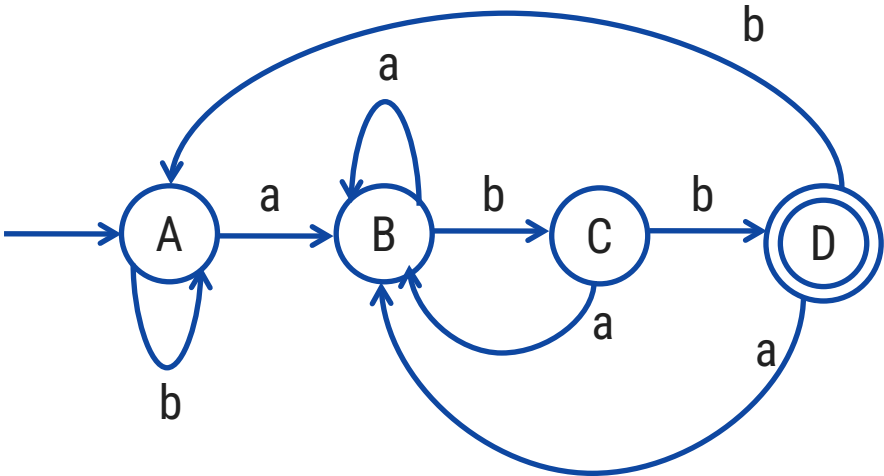


# Conversion from regular expression to DFA

## State D

$\delta((1,2,3,6),a) = \text{followpos}(1) \cup \text{followpos}(3)$   
 $= (1,2,3) \cup (4) = \{1,2,3,4\}$  ----- B

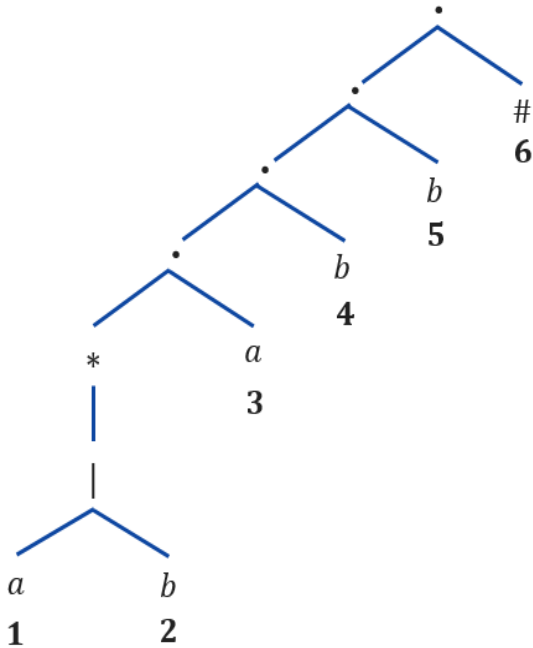
$\delta((1,2,3,6),b) = \text{followpos}(2)$   
 $= (1,2,3)$  ----- A



DFA

Position	followpos
5	6
4	5
3	4
2	1,2,3
1	1,2,3

States	a	b
A={1,2,3}	B	A
B={1,2,3,4}	B	C
C={1,2,3,5}	B	D
D={1,2,3,6}		



# **An Elementary Scanner Design & It's Implementation**

# An Elementary Scanner Design & It's Implementation

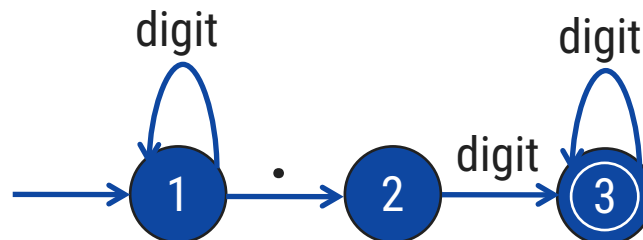
## Tasks of Scanner

1. The main purpose of the scanner is to return the next input token to the parser.
2. The scanner must identify the complete token and sometimes differentiate between keywords and identifiers.
3. The scanner may perform symbol-table maintenance, inserting identifiers, literals, and constants into the tables.
4. The scanner also eliminate the white spaces.

**Regular Expression:** Tokens can be **Specified** using regular expression.

Example: `id`  $\rightarrow$  `letter(letter | digit)*`

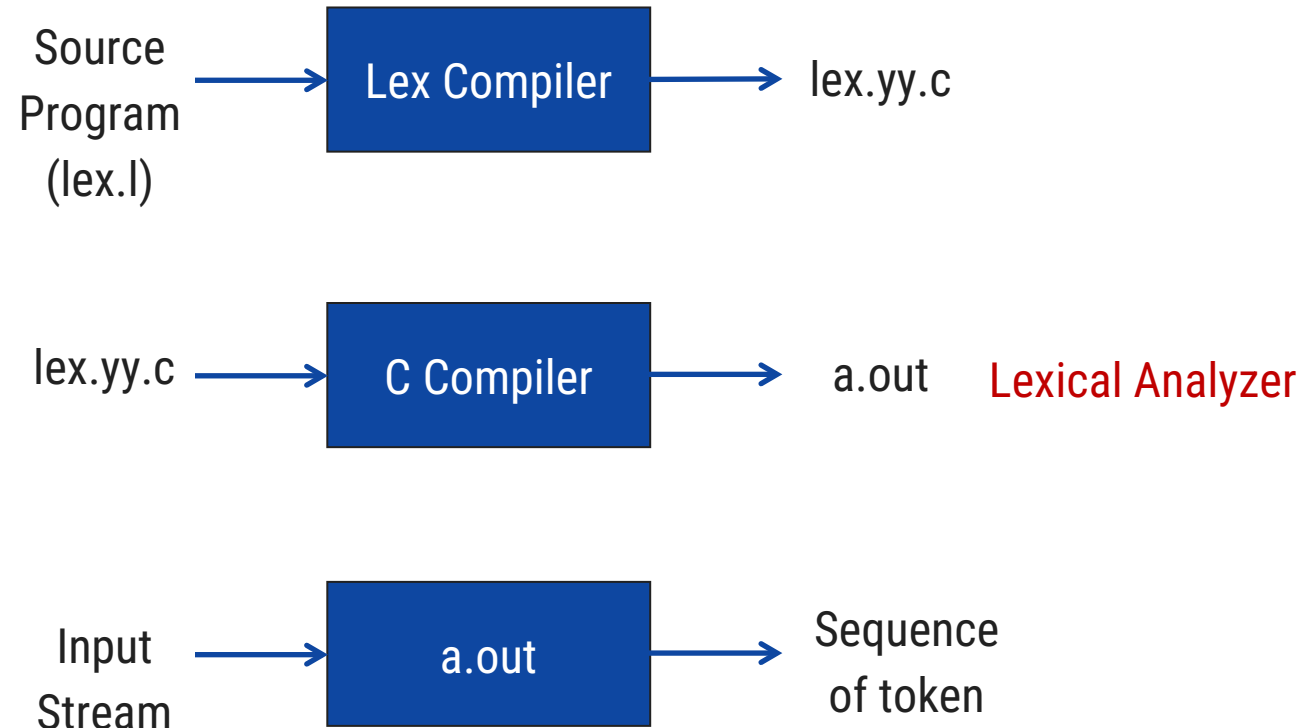
**Transition Diagram:** Finite-state diagrams or transition diagrams are often used to recognize a token



# Implementation of Lexical Analyzer (Lex)

- ▶ Lex is tool or a language which is useful for generating a lexical Analyzer and it specifies the regular expression
- ▶ Regular expression is used to represent the patterns for a token.

## Creating Lexical Analyzer with LEX



# Structure of Lex Program

► Any lex program contains mainly three sections

1. Declaration
2. Translation rules
3. Auxiliary Procedures

## Structure of Program

**Declaration** → It is used to declare variables, constant & regular definition  
Syntax: Pattern {Action}

%%  
%  
Example:  
%%

**Translation rule** → pattern1 {Action1}  
%% pattern2 {Action2}  
% pattern3 {Action3}

**Auxiliary Procedures**  $\xrightarrow{\text{L}} \text{%%}$  All the function needed are specified over here.

# Example: Lex Program

► Program: Write Lex program to recognize identifier, keywords, relational operator and numbers

---

```
/* Declaration */
```

```
%{  
    /* Lex program for recognizing tokens */
```

```
%}
```

```
Letter      [a-z A-z]
```

```
Digit       [0-9]
```

```
Id           {Letter}({Letter}|{Digit})*
```

```
Numbers      {Digit}+ (. {Digit}+)? (E[+ -]? Digit+)?
```

```
/* Auxiliary Procedures */
```

```
install_id()
```

```
{
```

```
    /* procedure to lexeme into the symbol table and return a pointer */
```

```
}
```

```
/* Translation rule */
```

```
%%
```

```
{Id}      {printf("%s is an identifier",yytext);}
```

```
If        {printf("%s is a keyword",yytext);}
```

```
else      {printf("%s is a keyword",yytext);}
```

```
"<"      {printf("%s is a less then operator",yytext);}
```

```
">="     {printf("%s is a greater then equal to operator",yytext);}
```

```
{Numbers} {printf("%s is a number",yytext);}
```

```
%%
```

Input string: If year < 2021

# References

## Books:

### **1. Compilers Principles, Techniques and Tools, PEARSON Education (Second Edition)**

Authors: Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman

### **2. Compiler Design, PEARSON (for Gujarat Technological University)**

Authors: Alfred V. Aho, Ravi Sethi, Jeffrey D. Ullman