

2.3 SPECIFICATION OF TOKENS (REGULAR EXPRESSIONS)

- To specify tokens, Regular Expressions are used.
- It provides convenient and useful notation for representing tokens.
- Regular Expressions define the language accepted by finite Automata (Transition Diagram).
- Regular Expressions are defined over an alphabet Σ .
- If R is a Regular Expression, then $L(R)$ represents language denoted by regular expression.

Language : It is a collection of strings over some fixed alphabet. Empty string can be denoted by ϵ .

E.g. if L (Language) = set of strings of 0's & 1's of length two
then $L = \{00, 01, 10, 11\}$

E.g. → If $L = \{1\}$
 then $L^* = L^0 \cup L^1 \cup L^2 \cup \dots$ Here $*$ can be 0, 1, 2, 3,.....
 $\therefore L^* = \{\epsilon\} \cup \{1\} \cup \{11\} \cup \dots$
 $\therefore L^* = \{\epsilon, 1, 11, 111, \dots\}$

Operations on Languages :

If $L_1 = \{00, 10\}$ & $L_2 = \{01, 11\}$

Operation	Description	Example
Union	$L_1 \cup L_2 = \{\text{set of string in } L_1 \text{ \& strings in } L_2\}$	$L_1 \cup L_2 = \{00, 10, 01, 11\}$
Concatenation	$L_1 L_2 = \{\text{Set of string in } L_1 \text{ followed by strings in } L_2\}$	$L_1 L_2 = \{0001, 0011, 1001, 1011\}$
Kleen closure of L_1 L_1^*	$L_1^* = L_1^0 \cup L_1^1 \cup L_1^2 \cup \dots$ $L_1^* = \bigcup_{i=0}^{\infty} L_1^i$	$L_1^* = \{\epsilon, 00, 10, 1010, 0010, 1000, 0000, 000000, 001000, \dots\}$
Positive Closure L_1^+	$L_1^+ = L_1^1 \cup L_1^2 \cup \dots$ $L_1^+ = \bigcup_{i=1}^{\infty} L_1^i$	$L_1^+ = \{00, 10, 1010, 0010, 1000, 0000, 000000, 001000, \dots\}$

Rules of Regular Expressions

- ϵ is a Regular expression.
- Union of two Regular Expressions R_1 and R_2 .
i.e. $R_1 + R_2$ or $R_1 | R_2$ is also a Regular Expression.
- Concatenation of two Regular Expressions R_1 and R_2
i.e. $R_1 R_2$ is also a Regular Expression.
- Closure of Regular Expression R i.e. , R^* is also a Regular Expression.
- If R is Regular Expression then (R) is also a Regular Expression.

Algebraic Laws

- $R_1 | R_2 = R_2 | R_1$ or $R_1 + R_2 = R_2 + R_1$ (Commutative)
 - $R_1 | (R_2 | R_3) = (R_1 | R_2) | R_3$ (Associative)
- or $R_1 + (R_2 + R_3) = (R_1 + R_2) + R_3$

$$3. \quad R_1 (R_2 R_3) = (R_1 R_2) R_3$$

$$4. \quad R (R_2 \mid R_3) = R_1 R_2 \mid R_1 R_3$$

$$\text{or } R_1 (R_2 + R_3) = R_1 R_2 + R_1 R_3$$

$$5. \quad \varepsilon R = R \varepsilon = R$$

(Associative)
(Distributive)

(Concatenation)

Example 2: Find Regular Expression for following Language.

$$(a) \quad L = \{\varepsilon, 1, 11, 111, \dots\}$$

$$\{\text{Hint : } * = 0, 1, 2, 3, \dots\} \{ \because 1^0 = \varepsilon, 1^1 = 1, 1^2 = 11, 1^3 = 111, \dots \}$$

$$\text{Ans. } 1^*$$

$$(b) \quad L = \{\varepsilon, 11, 1111, 111111, \dots\}$$

$$\text{Ans. } (11)^*$$

$$(c) \quad L = \text{Set of all strings of 0's and 1's} = \{\varepsilon, 0, 1, 01, 11, 00, 000, 101, \dots\}$$

$$\text{Ans. } (0 + 1)^* \text{ or } (0 \mid 1)^*$$

$$(d) \quad L = \text{Set of all strings of 0's and 1's ending with 11.}$$

$$\text{Ans. } (0 + 1)^* 11$$

$$(e) \quad L = \text{Set of all strings of 0's and 1's beginning with 0 and Ending with 1.}$$

$$\text{Ans. } 0 (0 + 1)^* 1$$

Example 3 : Write Regular Expressions for following language over $\Sigma = \{a, b\}$

$$(a) \quad \text{Strings of length zero or one.}$$

$$\text{Ans. } \varepsilon \mid a \mid b \text{ or } (\varepsilon + a + b)$$

$$(b) \quad \text{Strings of length two.}$$

$$\text{Ans. } aa \mid ab \mid ba \mid bb \text{ or } (aa + ab + ba + bb)$$

$$(c) \quad \text{Strings of Even length.}$$

$$\text{Ans. } (aa \mid ab \mid ba \mid bb)^* \text{ or } (aa + ab + ba + bb)^*$$

$$(d) \quad \text{Set of all strings of a's and b's having atleast two occurrences of aa.}$$

$$\text{Ans. } (a + b)^* aa (a + b)^* aa (a + b)^*$$

Example 4 : Write Regular Expressions for all types of Tokens.

$$\text{Ans. (a) Keywords} = \text{begin} \mid \text{end} \mid \text{if} \mid \text{then} \mid \text{else}$$

$$(b) \quad \text{Identifier} = \text{letter} (\text{letter} + \text{digit})^*$$

$$(c) \quad \text{Constant} = \text{digit} \text{ digit}^* = \text{digit}^+$$

$$(d) \quad \text{Relation-op} = < \mid > \mid <= \mid >= \mid <>$$

Example 5. Write Regular Expression in which second letter from right end of string is 1 where $\Sigma = \{0, 1\}$.

$$\text{Sol. } (0 + 1)^* 1 (0 + 1)$$

Example 6. Write Regular Expression for $\Sigma = \{a, b\}$

(a) $L =$ set of strings having atleast one occurrence of double letter

Ans. $(a + b)^* (aa + bb) (a + b)^*$

(b) $L =$ set of strings having double letter at Beginning and Ending of string.

Ans. $(aa + bb) (a + b)^* (aa + bb)$

(c) $L =$ set of string having double letter at beginning or on ending of string.

Ans. $(aa + bb) (a + b)^* + (a + b)^* (aa + bb) + (aa + bb) (a + b)^* (aa + bb)$.

Example 7. Write regular expression for set of words having a, e, i, o, u appearing in the order although not necessarily consecutively.

Ans. If $\alpha = ([b - d] + [f - h] + [j - n] + [p - t] + [v - z] + \text{Blankspace})^*$

\therefore Required Regular Expression will be

$a a \alpha e a i \alpha o \alpha u \alpha$

2.4. RECOGNITION OF TOKEN (FINITE AUTOMATA)

Regular Expression are used to represent various type of Tokens. From these Regular Expression, Finite Automata is generated. **Finite Automata is basically other name of Lexical Analyzer** which checks or recognize whether a string is Token or not.

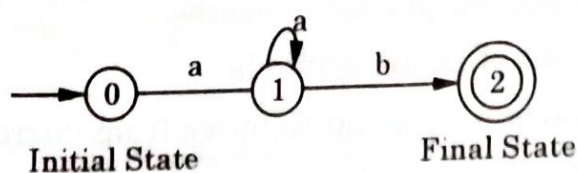
Finite Automata :

- It is a Machine or a Recognizer for a language that is used to check whether string is accepted by a language or not.
- It gives answer "Yes" if string is accepted else "no".
- In **Finite Automata**, **Finite** means finite number of states & **Automata** means, the Automatic Machine which works without any interference of human being.

Representation of Finite Automata : There are 2 ways to represent finite Automata.

1. **Transition Diagram :** It is a directed graph or flow chart having states & edges.

E.g.



0, 1, 2 \rightarrow States

0 \rightarrow Initial state

2 \rightarrow Final State

$a, b \rightarrow$ Input symbols.

2. **Transition Table :**

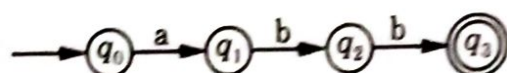
States \ Input		
	a	b
0	1	-
1	1	2
2	-	-

• **Finite Automata can be Represented by 5 tuple $(Q, \Sigma, \delta, q_0, F)$**

- (a) Q is finite non-empty set of states.
- (b) Σ is finite set of input symbols.
- (c) δ is transition function.
- (d) $q_0 \in Q$ is initial state.
- (e) $F \subseteq Q$ is the set of final states.

Example 8: Design Finite Automata which accepts string "abb"

Ans.



States	:	$Q = \{q_0, q_1, q_2, q_3\}$
Input symbols	:	$\Sigma = \{a, b\}$
Transition Function δ	:	$\{\delta(q_0, a) = q_1, \delta(q_1, b) = q_2, \delta(q_2, b) = q_3\}$
Initial State	:	q_0
Final State (F)	:	$\{q_3\}$

2.4.1 Types of Finite Automata

- (a) Deterministic Finite Automata (DFA)
- (b) Non-Deterministic Finite Automata (NFA)

(a) Deterministic Finite Automata (DFA) : "Deterministic" means on each input there is one and only one state to which the automata can transition from its current state.

• **DFA is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$**

where

- (a) Q is finite non-empty set of states
- (b) Σ is finite set of input-symbols
- (c) δ is a transition function to move from current state to next state
 $\therefore \delta : Q \times \Sigma \rightarrow Q$
- (d) $q_0 \in Q$ is initial state
- (e) $F \subseteq Q$ is set of final states

(b) **Non-Deterministic Finite Automata (NFA)** : "Non-Deterministic" means there can be several possible transitions. So, output is non-deterministic for a given input.

• **NFA is 5-tuple $(\theta, \Sigma, \delta, q_0, F)$ where**

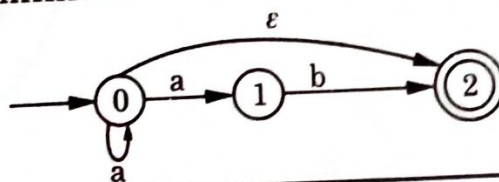
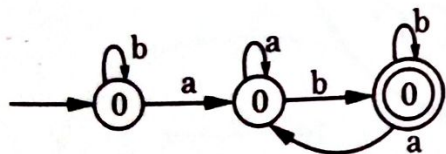
- (a) Q is finite non-empty set of states.
- (b) Σ is finite set of input symbols.
- (c) δ is transition function to move from current state to next state.

$$\delta : Q \times \Sigma \rightarrow 2^Q$$

- (d) $q_0 \in Q$ is initial state
- (e) $F \subseteq Q$ is set of final states.

• **Difference Between DFA and NFA**

DFA	NFA
<ol style="list-style-type: none"> Every transition from one state to other is unique & deterministic in nature. Null transitions (ϵ) are not allowed. Transition function $\delta : Q \times \Sigma \rightarrow Q$. Requires less memory as transitions & states are less E.g DFA <p>In fig. Each State has unique input symbol $\Sigma = \{a, b\}$ outgoing from it.</p>	<p>There can be multiple transitions for an input i.e. non-deterministic.</p> <p>Null Transitions (ϵ) are allowed means transition from current state to next state without any input.</p> <p>Transition function $\delta : Q \times \Sigma \rightarrow 2^Q$.</p> <p>Requires more memory.</p> <p>E.g NFA</p> <p>In fig. As, there are multiple transitions of input symbol 'a' on state 0. & no transition from state 2, i.e. non-deterministic.</p>



2.5 CONVERSION OF REGULAR EXPRESSION TO FINITE AUTOMATA (NFA)

A Regular Expression is basically a representation of Tokens. But to Recognize a token, we need token Recognizer which is nothing but a Finite Automata (NFA). So, we can convert Regular Expression into NFA.

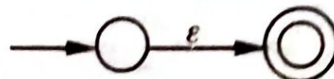
Algorithm for conversion of RE to NFA:

Input : A Regular Expression R

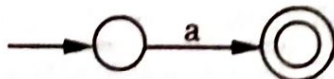
Output : NFA accepting language denoted by R

Method :

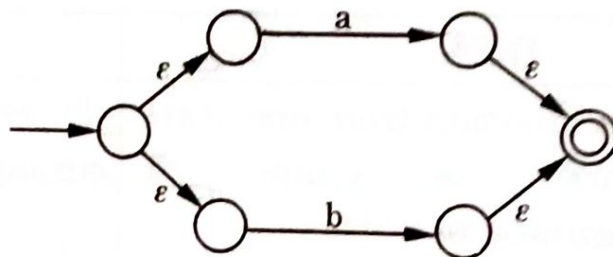
1. For ϵ , NFA is



2. For a, NFA is



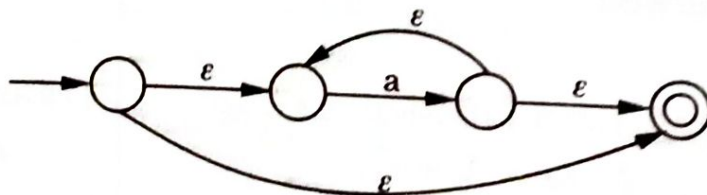
3. For $a + b$, or $a | b$ NFA is



4. For ab , NFA is

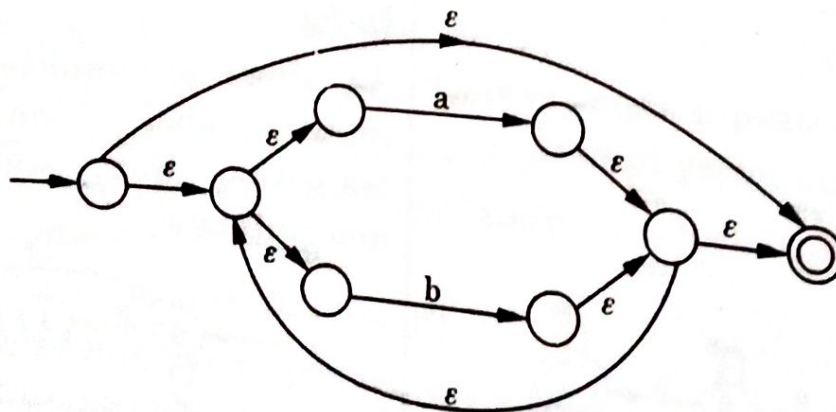


5. For a^* , NFA is



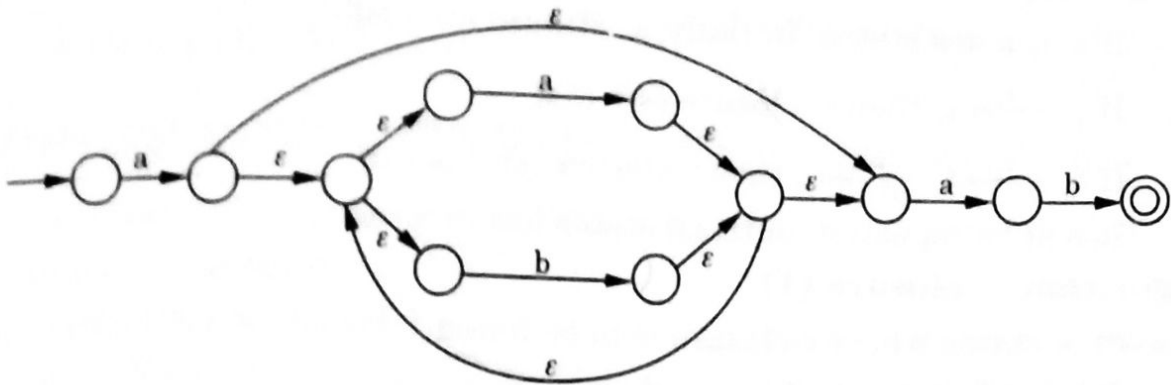
Example 9: Draw NFA for Regular Expression $(a+b)^*$

Ans.



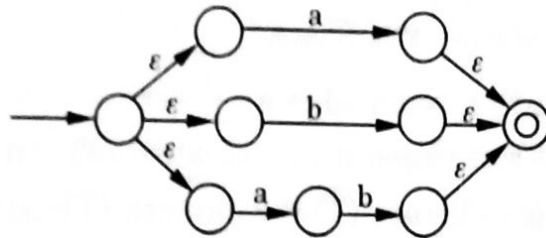
Example 10: Draw NFA for Regular Expression $a(a+b)^*ab$

Ans.



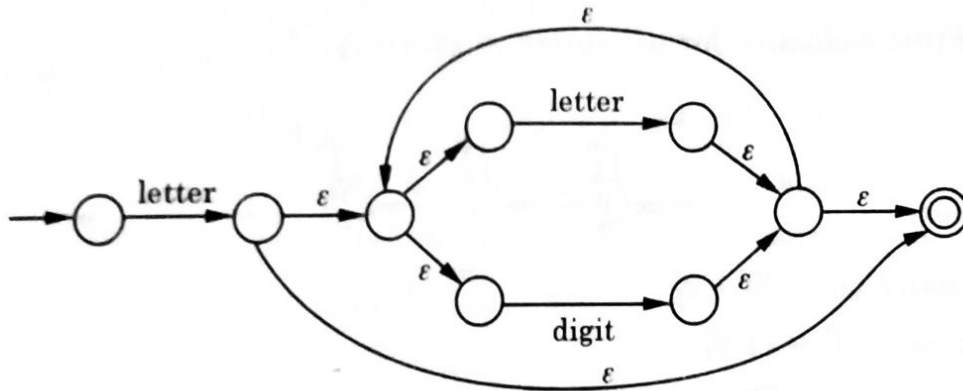
Example 11: Draw NFA for $a+b+ab$

Ans.



Example 12 : Draw NFA for $\text{letter}(\text{letter} + \text{digit})^*$

Ans.



Example 13: Draw NFA corresponding to $(0+1)^*1(0+1)$

Ans.

