# Contents

# 1.  Introduction to Automata Theory & Compiler Design

The theory of computation is a field of study that deals with the nature and limits of computation. It encompasses various mathematical and theoretical aspects related to the capabilities and limitations of computers and computational systems.

Automata theory is a branch of theoretical computer science that deals with the study of abstract machines and the problems they can solve. It has applications in various areas of computer science, including formal languages, compiler design, artificial intelligence, and hardware design

It mainly focuses on what kind of things you can compute mechanically, how fast, and how much space it takes.

Examples:

1.  Design a machine that accepts Binary Strings end with 0

    110010110 – LAST DIGIT O ACCEPT

    00000111 – REJECT DIGIT BECAUSE LAST DIGIT 1

2.  Design a machine that accepts Java Code and converts to the binary that the string written is Valid / Invalid

    e.g: -JAVA CODE → COMPILER ->BINARY –NO SYNTAX ERROR –VALID
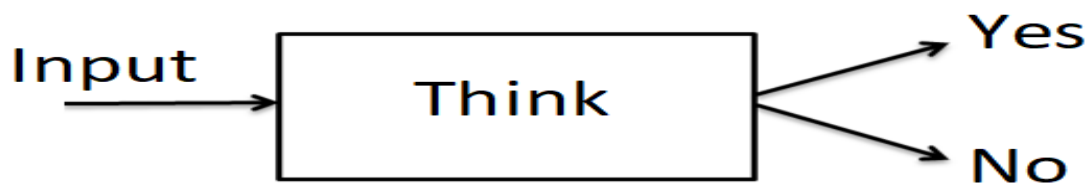
    SYNATAX ERROR -INVALID



Figure 1.1 General Concept of theory of computation design

Theory of Computations is classified as Finite State Machine(FSM), Context Free Language(CFL), Turing Machine, and Undecidable.    This can be depicted as below diagram
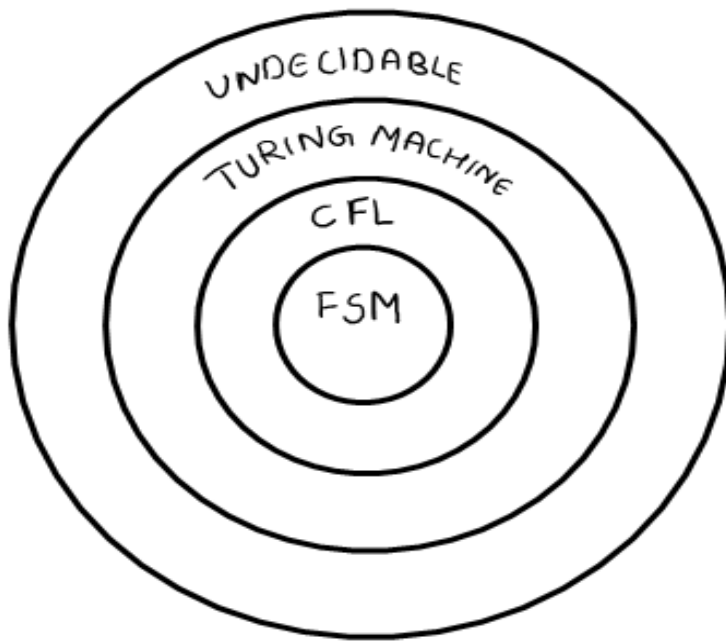
Figure 1.2: Levels of Theory of computations applied

These computations used as below areas such as

## 1.1 Applications

### I.  FINITE STATE MACHINE APPLICATIONS

A Finite State Machine (FSM) is a computational model consisting of states, transitions, and an initial state. It processes inputs and transitions between states based on predefined rules, making it suitable for modeling sequential logic in various applications, such as control systems and parsing algorithms. FSMs are defined by a finite set of states, inputs, transitions, and an optional set of final or accepting states. This low computing model used in below areas

- **Lexical Analysis in Compilers:**   FSAs are used in the lexical analysis phase of compiler design to recognize and tokenize the input source code. Each token in the source code can be identified using a state transition in the FSA.

- **Regular Expression Matching:** FSAs are closely related to regular expressions, and they are used to implement pattern-matching algorithms. Regular expressions can be converted into equivalent FSAs, making them a fundamental tool for text processing and searching.

- **Natural Language Processing (NLP):** FSAs are applied in various NLP tasks, such as part-of-speech tagging and named entity recognition. They can model simple grammatical structures and help in parsing and analyzing linguistic patterns.

- **Network Protocol Specification:** FSAs are used in the specification and modeling of network protocols. They can represent the different states that a communication system can be in and the transitions between these states based on inputs.

- **Digital Circuit Design:** In digital electronics, FSAs are used to design and describe the behavior of sequential circuits. Finite state machines are employed to model the control units of digital systems.

- **Control Systems:** FSAs find applications in modeling and controlling various systems, such as traffic lights, vending machines, and industrial processes. The discrete nature of state transitions makes them suitable for modeling these systems.

- **Robotics:** In robotics, FSAs can be used to model the behavior of robots. They can represent different states of the robot and transitions between these states based on sensor inputs or commands.

- **Game Design:** FSAs are employed in game design to model the behavior of characters, game states, and events. They can be used to define the rules and logic governing the game.

- **Communication Protocols:** Finite state machines are used in the design of communication protocols to model the different phases of communication and the allowable transitions between these phases.

- **VLSI Design:** In Very Large-Scale Integration (VLSI) design, FSAs are used for modeling and verification of digital circuits. They help in ensuring that the designed circuits operate correctly based on specified conditions.


## II.    CONTEXT FREE LANGUAGE APPLICATIONS

A context-free language is a type of formal language in formal language theory. It is described by a context-free grammar, where production rules define the language's syntax. Context-free languages are recognized by pushdown automata and are widely used in the design and analysis of programming languages. This concept of computation used in following areas

- **Compiler Design:** Context-free grammars are extensively used in the syntax analysis phase of compiler design. They help define the structure of programming languages, allowing compilers to recognize and parse source code.

- **Syntax Highlighting and Code Formatting:** Context-free grammars are employed in text editors and Integrated Development Environments (IDEs) for syntax highlighting and code formatting. They assist in visually distinguishing different elements of the code based on their syntactic roles**.**

- **Natural Language Processing (NLP):** Context-free grammars are used in certain aspects of natural language processing, such as parsing sentences and recognizing syntactic structures in languages. They help in understanding the grammatical relationships within sentences**.**

- **XML and HTML Parsing:** Context-free grammars play a role in parsing and validating XML (eXtensible Markup Language) and HTML (HyperText Markup Language) documents. These languages have hierarchical structures that can be described using context-free grammars**.**

- **Database Query Languages:** Query languages, such as SQL (Structured Query Language), often have a context-free structure. Context-free grammars are used to parse and process queries for relational databases**.**

- **Programming Language Design:** When designing new programming languages, context-free grammars are employed to specify the syntax of the language. This is crucial for creating compilers or interpreters for the language.

- **Expression Evaluation:** Context-free grammars are used to describe the syntax of mathematical expressions. They are applied in parsing and evaluating expressions in computer algebra systems, calculators, and other applications.

- **Markup Languages:** Context-free grammars are used to define the syntax of markup languages like LaTeX, which is widely used for typesetting scientific documents. They help in specifying how documents should be structured.

- **Data Interchange Formats:** Some data interchange formats, such as JSON (JavaScript Object Notation), have context-free grammars that define the structure of data. These grammars are used for parsing and generating data in a standardized format**.**

- **Pattern Matching and String Analysis:** Context-free grammars are employed in pattern matching algorithms and string analysis. They help in identifying patterns within strings, which is useful in various applications, including data mining and information retrieval**.**


## III.     TURING MACHINE APPLICATIONS

Turing Machines, the theoretical model of computation introduced by Alan Turing, serve as a foundation for understanding the limits and capabilities of computation. While they are more abstract and not directly used in practical applications like Finite State Automata or Context-Free Grammars, their concepts have influenced various aspects of computer science

- **Algorithm Analysis and Complexity Theory:** Turing Machines are crucial for understanding the time and space complexity of algorithms. Computational complexity classes, such as P and NP, are defined in terms of the resources (time and space) required by a Turing Machine to solve a problem.

- **Computational Theory:** Turing Machines provide the theoretical underpinning for the Church-Turing thesis, asserting that any effectively computable function can be computed by a Turing Machine. This thesis has profound implications for the nature and limits of computation.

- **Programming Language Design:** Turing Machines influence the design and study of programming languages. Concepts such as computability, decidability, and the halting problem are essential considerations in language design and compiler construction.

- **Artificial Intelligence (AI) and Computation Theory**: Theoretical concepts derived from Turing Machines, such as computability and decidability, play a role in understanding the possibilities and limitations of artificial intelligence algorithms and models.

- **Cryptography**: The study of computational complexity, influenced by Turing Machines, is relevant to the field of cryptography. Problems that are hard for a Turing Machine to solve (like factoring large numbers) form the basis for secure cryptographic algorithms.

- **Database Query Processing**: Theoretical aspects related to the efficiency of query processing in databases, such as the complexity of query evaluation, draw upon concepts from Turing Machines and computational complexity theory.

- **Theoretical Computer Science Research**: ``Turing Machines are foundational in theoretical computer science research. The study of complexity classes, decidability, and other aspects of computation is ongoing and influences the development of new algorithms and models.

- **Formal Verification**: Concepts from Turing Machines are used in formal verification methods to ensure the correctness of hardware and software systems. Formal verification involves mathematically proving that a system behaves as intended.

- **Automated Theorem Proving**: Turing Machines contribute to the study of automated theorem proving, a field that aims to develop algorithms for proving mathematical theorems using computational methods.

- **Understanding the Limits of Computation**: Turing Machines are crucial for understanding what can and cannot be computed algorithmically. The study of undecidable problems, such as the halting problem, helps define the boundaries of what is achievable through computation.

## IV.  UNDECIDABLE

Undecidable problems are problems for which no algorithm can always correctly determine the answer for all possible inputs. One of the most famous undecidable problems is the Halting Problem, introduced by Alan Turing

- **Halting Problem:** Given a description of an arbitrary computer program and an input, determine whether the program will eventually halt (terminate) or continue running indefinitely. Alan Turing proved that there is no general algorithm to solve the Halting Problem

- **Hilbert's Tenth Problem**: Hilbert's Tenth Problem asks for an algorithm to determine whether a given Diophantine equation (polynomial equation with integer coefficients) has integer solutions. This problem was proven to be undecidable by Yuri Matiyasevich, building on previous work by Julia Robinson, Martin Davis, and Hilary Putnam

In computer science and mathematics, the concept of an undecidable problem refers to a type of problem for which there is no algorithm (step-by-step procedure) that can determine a correct solution for all

possible inputs. In simpler terms, it's a problem where there's no general method to find the right answer every time. Now, let's imagine a "Super Puzzle Solver" that claims to solve any puzzle thrown at it. However, we're talking about puzzles that are so tricky, that there's no guaranteed way to solve them. It's like having a magical friend who says, "I can solve any puzzle you give me, no matter how complicated!"

But here's the catch: If the puzzle is of a certain type, known as undecidable, our Super Puzzle Solver might run into a problem. There's no one-size-fits-all solution for these puzzles. It's like trying to come up with a recipe that works for every dish in the world – it's just not possible.

To illustrate, consider the "Halting Problem," a famous example of an undecidable problem. It asks whether a given program will eventually stop running or continue indefinitely. The Super Puzzle Solver might struggle with this one because there's no algorithm that can predict whether any arbitrary program will halt or go on forever. It's like trying to predict the future actions of a mischievous cat – you never know if it will suddenly decide to keep running around.

So, while our Super Puzzle Solver is indeed impressive for many puzzles, it can't guarantee success for undecidable problems. These are the puzzles that push the limits of what can be solved algorithmically, making them particularly challenging and intriguing in the world of computer science and mathematics.

## *1.2 Essential Foundation of FSM (Define Terms )*

**Symbols**:    Symbols are entities or individual objects, which can be any letter, alphabet, or any picture. An atomic unit, such as a digit, character, lower-case letter, etc. Sometimes word. [Formal language does not deal with the "meaning" of the symbols.

Example: 1, a, b, #, 0, 2, 3 so on

**Alphabet**: The alphabet is a finite set of symbols, which are the basic elements used to construct strings. In the context of automata theory, an alphabet is typically denoted by Σ.

∑ = {a, b} ,    ∑ = {A, B, C, D} , ∑ = {0, 1, 2} ,∑ = {0, 1, ....., 5} ,∑ = {#, β, Δ}

**What is power of Alphabet ?**

In the context of Finite State Machines (FSMs), the term "alphabet" refers to the set of symbols or inputs that the machine can recognize and respond to. The power of the alphabet in an FSM lies in its ability to define the possible inputs and transitions between states. Let's break down the concept with an example.

Consider a simple finite state machine that models a door with two states: "closed" and "open." The possible inputs or alphabet for this FSM might consist of two symbols: {Push, Pull}. The door can transition between states based on these inputs. Here's a basic representation:

State : Closed
↓
Input : Push - - →state : open

State : open
↓
Input :  Pull  - - → state : closed

In this example:

**The alphabet is {Push, Pull}.**

**The states are {Closed, Open}.**

**The transitions are defined by the inputs Push and Pull.**

**The power of the alphabet becomes apparent in how it allows the FSM to respond to different inputs and change its state accordingly**. For instance:

If the door is in the "Closed" state, and the input is "Push," it transitions to the "Open" state.

If the door is in the "Open" state, and the input is "Pull," it transitions back to the "Closed" state.

The alphabet defines the language that the FSM understands, and it plays a crucial role in determining the behavior of the system. It's worth noting that the power of the alphabet in FSMs extends to more complex scenarios where there can be multiple states, inputs, and transitions, allowing FSMs to model a wide range of systems and processes.


**String**: A string is a finite sequence of symbols from an alphabet. Strings are the input and output of automata.

If $\sum$ = {a, b}, various strings that can be generated from $\sum$ are {ab, aa, aaa, bb, bbb, ba, aba.....}.

A string with zero occurrences of symbols is known as an empty string. It is represented by ε.


Some special sets of strings: Σ*  All strings of symbols from

$$\Sigma+ = \Sigma* -\{\varepsilon\}$$

Example: Σ = {0,1}

  Σ* = {**ε**, 0, 1, 00, 01, 10, 11, 000, 001,…}

  Σ+ = {0, 1, 00, 01, 10, 11, 000, 001,…}

**A language**: A set of strings from some alphabet (finite or infinite). In other words, Any subset L of Σ*

Some Special languages

{}    -    The empty set/language, containing no string.

{ε}    -    A language containing one string, the empty string.

Examples

Σ = {0,1}

L = {x | x is in Σ* and x contains an even number of 0"s}    Σ = {0000, 10101010, 1100,0011 ...}

L = {x | x is in Σ* and x forms a finite length real number} = {0, 1.5, 9.326,…}

Σ = {a, b, c,…, z, A, B,…, Z}

L = {x | x is in Σ* and x is a Pascal reserved word}

= {BEGIN, END, IF,…}

Σ = {Pascal reserved words} U { (, ), ., :, ;,…} U {Legal Pascal identifiers} L = {x | x is in Σ* and x is a syntactically correct Pascal program}

Σ = {English words}

L = {x | x is in Σ* and x is a syntactically correct English sentence}

In Alphabet Σ= {0,1}

**Finite Sets** :

L1 = set of all strings of length 2    =    {00,01,10,11}

L2 = set of all string of length 3     = {000, 001, 010,011,101,110,111}

**Infinite Sets**:

L3 = set of all strings that begin with 0

={0, 00, 000, 010,0110, 0111, …. }

Explain concatenation of languages?


**Power of Σ**

The power set of a set is the set of all possible subsets of that set. If we consider ε as an element in a set, then the power set of {ε} would be {{}, {ε}}. If ε is an empty string, then the power set of {ε} is {{}}.

In formal language theory and FSMs, the use of epsilon often signifies the possibility of a transition without consuming any input. Therefore, if ε is an element in a set, the power set includes both the set without ε

and the set with ε. If ε represents the empty string, then the power set includes only the set without ε (i.e., the set of all possible subsets of the empty set, which is {{}}).

For instance, ε ={0, 1}

$\Sigma^0$=set of all strings of length 0 : $\Sigma^0 = \{ \in \}$

$\Sigma^1$=set of all strings of length 1 : $\Sigma^1 = \{ 0,1 \}$

$\Sigma^2$=set of all strings of length 2 : $\Sigma^2 = \{ 00,01,10,11\}$

$\Sigma^3$=set of all strings of length 3 : $\Sigma^3 = \{ 000, 001, 010, \ 011, 100, 101, 110, 111 \}$

$\Sigma^n$=set of all strings of length n

Cardinality- number of elements in a set    $\Sigma^n = 2^n$

$$\Sigma^* = \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \Sigma^0 \dots \dots \dots \dots \dots$$

$$=\{\in\} \cup \{0.1\} \cup \{00,01,10,11\} \cup \dots \dots \dots \dots \dots.$$

=set of all possible strings of all lengths over {0,1}

**Concatenation of Language**:

The concatenation of languages is a fundamental operation in formal language theory. If you have two languages, say $L_1$ AND $L_2$ Their concatenation, is denoted as $L_1$. $L_2$. is a new language consisting of all possible concatenations of strings where the first part comes from $L_1$ and the second part comes from $L_2$.

It is defined as:

L1·L2= {xy |x∈ $L_1$ ,y∈ $L_2$ }

Here, xy denotes the concatenation of strings x and y.

For example, let's say L1={ab,c} and L2={12,34}, then the concatenation L1·L2 would be:

L1·L2={ab12, ab34, c12, c34}


## 1.3  Problem

The problem of deciding whether a given string is part of a particular language is essentially the task of determining if a string belongs to a set of strings that adheres to the rules defined by a specific type of automaton or grammar.

*The problem is to decide whether a given string s belongs to a particular language L.* This means determining if the string s satisfies the rules defined by the language L.

Example: Recognizing Palindromes: Let's consider a simple example where the language is defined as the set of palindromes over the alphabet {0,1} {0,1}. A palindrome is a string that reads the same forwards as backward.

Language L: The set of palindromes over {0,1} {0,1}.

Example String s: 101101101101.

Automaton or Grammar: For this example, you can use a finite state machine (FSM) or a regular grammar to define the language of palindromes. The automaton or grammar would have rules specifying how to generate or recognize palindromes.

Deciding Membership: To decide whether a given string s (e.g., 101101101101) is part of the language, you would follow the rules of the automaton or grammar. For palindromes, you might use a state machine to check if the string reads the same backward as forward.

In the case of 101101101101, you would follow the transitions of the state machine, reading the string from left to right. If, after processing the entire string, you end up in an accepting state, then the string is in the language. Otherwise, it's not.

If the state machine accepts 101101101101, then 101101101101 is in the language of palindromes. If it doesn't accept, then 101101101101 is not a palindrome.

This process of deciding whether a string is part of a language is fundamental to automata theory. Different types of automata (finite state machines, pushdown automata, Turing machines, etc.) and grammars (regular grammars, context-free grammars, etc.) define different types of languages, and the problem of deciding membership is a central concept in understanding and working with these formal languages

## 1.4  FINITE STATE MACHINE

A Finite State Automaton (FSA), also known as a Finite State Machine (FSM), is a computational model used to represent and analyze systems that can exist in a finite number of states. These states are connected by transitions, and the machine moves from one state to another based on inputs. FSAs are widely used in various areas of computer science, including formal language theory, compiler design, and modeling of systems with discrete behavior.

**the key components of a Finite State Automaton**:

- **States**: A finite set of distinct states represents the possible conditions or situations of the system. Each state represents a snapshot of the system at a specific moment.

- **Alphabet**: An alphabet is a finite set of symbols that the automaton recognizes as inputs. These symbols trigger state transitions when encountered.

- **Transitions**: Transitions define how the automaton moves from one state to another in response to input symbols. Each transition is associated with a specific input symbol and leads to a new state.

- **Start State**: One state is designated as the initial or start state. This is the state the automaton begins in before processing any input.

- **Accepting States (Optional)**: Some states may be designated as accepting states. When the automaton reaches an accepting state after processing an input sequence, it signifies that the input sequence is recognized or accepted.

As we know, Finite State Automata are used in various applications, including:

- Lexical Analysis in Compilers: Recognizing and tokenizing the structure of programming languages.

- Pattern Matching: Searching for patterns within strings or sequences.

- Network Protocols: Modeling the behavior of communication protocols.

- Digital Circuit Design: Describing the control logic of digital systems.

- Natural Language Processing: Analyzing and parsing language structures

**Formal Definition of FA**

A finite automaton is a collection of 5-tuple (Q, ∑, δ, q0, F), where:

Q: a finite set of states

∑: a finite set of the input symbol

q0: initial state

F: final state

δ: Transition function

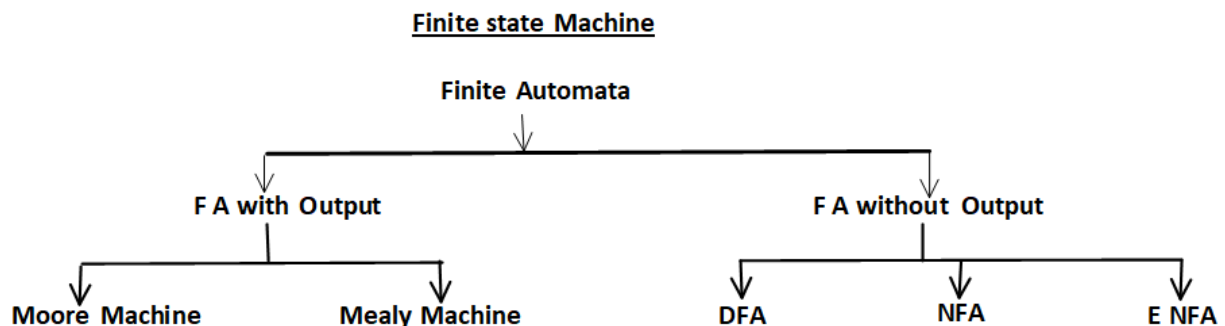Finite-state machines classified as broadly as below



Figure 1.3 : classification of Finite State machine

**FA With Output:** A Finite Automaton (FA) with output, is a type of finite-state machine that not only processes input symbols but also produces output based on its current state and using either input or without input. This can be broadly classified as **Moore Machine** and **Mealy Machine**.

**A Moore machine** is a type of finite state machine (FSM) where each state is associated with a specific output. In other words, the output of a Moore machine depends only on its current state and not on the inputs that caused the transition between states.

Imagine a traffic light that can be either in a "Green" state, "Yellow" state, or "Red" state. A Moore machine associated with this traffic light would display a specific color (output) depending on its current state. The color displayed is solely determined by the state of the traffic light and doesn't depend on any particular event or action.
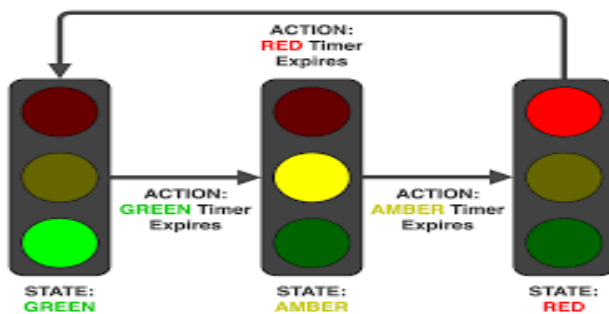


Figure 1.4: Traffic signals states and Transition

- **States**: "Green," "Yellow," "Red."

- **Inputs**: The passage of time (no specific input needed).

- **Outputs**: The color displayed: "Green," "Yellow," "Red."

In simple terms, a Moore machine is like a system that displays a certain output (like a color on a traffic light) based solely on its current state. It doesn't care about what caused it to change states; it just knows what to display given its current condition. The above example of the traffic light helps illustrate how Moore machines work by associating specific outputs with each state.

**A Mealy machine** is another type of finite state machine (FSM) where the output depends not only on the current state of the machine but also on the inputs that caused the transition between states. In contrast to a Moore machine, which associates outputs with states, a Mealy machine produces outputs based on both the current state and the inputs.

a turnstile that controls access to a building. The turnstile's behavior (output) depends on both its current state and whether you insert a valid card (input).

Figure 1.6 : Turnstile machine

- **States**: "Locked" and "Unlocked."

- **Inputs**: "Insert Card" and "Push Turnstile."

- **Outputs**: "Access Denied" and "Access Granted."

**FA Without Output**: In the context of Finite State Automata (FSAs), the term "without output" refers to the fact that the automaton is not producing any meaningful output or result as it transitions between states. The focus is primarily on modeling the sequence of states that the system goes through in response to inputs, rather than generating specific outputs.

## 1.5 A Deterministic Finite Automaton (DFA)

DFA is like a simple machine that reads input and follows specific rules to decide whether a given input is accepted or rejected. It's often used in computer science to recognize patterns or check if a sequence of symbols belongs to a certain set.

- DFA refers to deterministic finite automata. Deterministic refers to the uniqueness of the computation. The finite automata are called deterministic finite automata if the machine reads an input string one symbol at a time.

- In DFA, there is only one path for specific input from the current state to the next state.

- DFA does not accept the null move, i.e., the DFA cannot change state without any input character.

- DFA can contain multiple final states.

In the following diagram, we can see that from state q0 for input a, there is only one path that is going to q1. Similarly, from q0, there is only one path for input b going to q2.

A DFA is a collection of 5-tuples same as we described in the definition of FA.

Q: finite set of states

∑: finite set of the input symbol
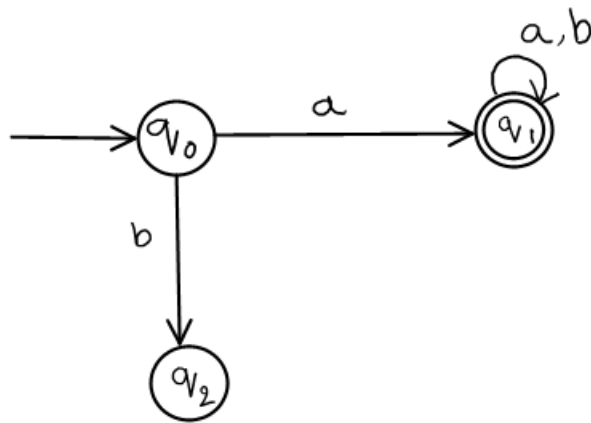
q0: initial state

F: final state

δ: Transition function



Figure 1.8 DFA

Transition function can be defined as: δ: Q x ∑→Q

The Graphical representation of DFA :

A DFA can be represented by digraphs called state diagram. In which:

- The **state** is represented by **vertices**.

- The arc **labeled** with an input character show the **transitions**.

- The **initial state** is marked with an **arrow**.

- The **final state** is denoted by **a double circle**.

1. Design a FA with ∑ = {0, 1} accepts those string which starts with 1 and ends  with 0

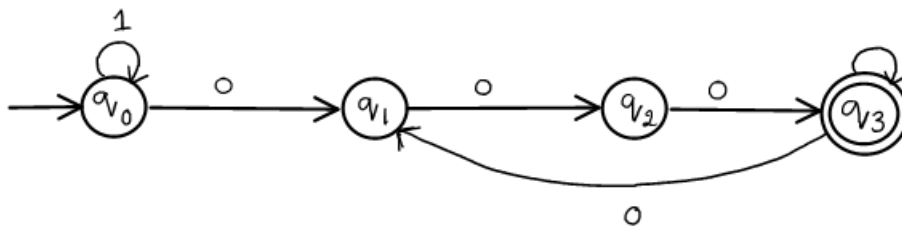2. Design a FA with ∑ = {0, 1} accepts the only input 101.
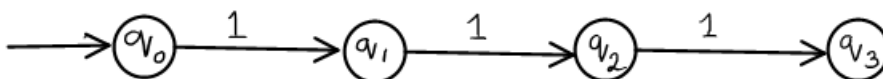


3. Design FA with ∑ = {0, 1} accepts the set of all strings with three consecutive 0's.

The strings that will be generated for this particular language are 000, 0001, 1000, 10001, .... in which 0 always appears in a clump of 3. The transition graph is as follows:
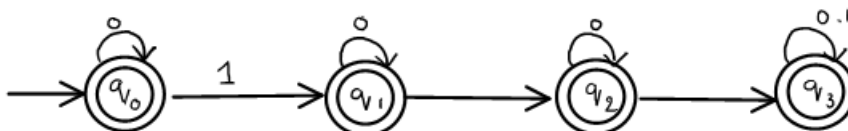


4. Design a DFA L(M) = {w | w ε {0, 1}*} and W is a string that does not contain consecutive 1's

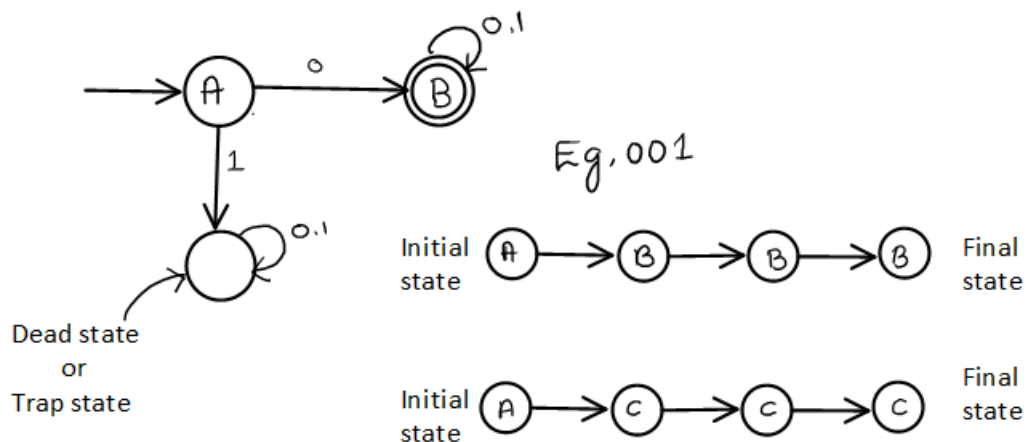When three consecutve 1's occur the DFA will be :



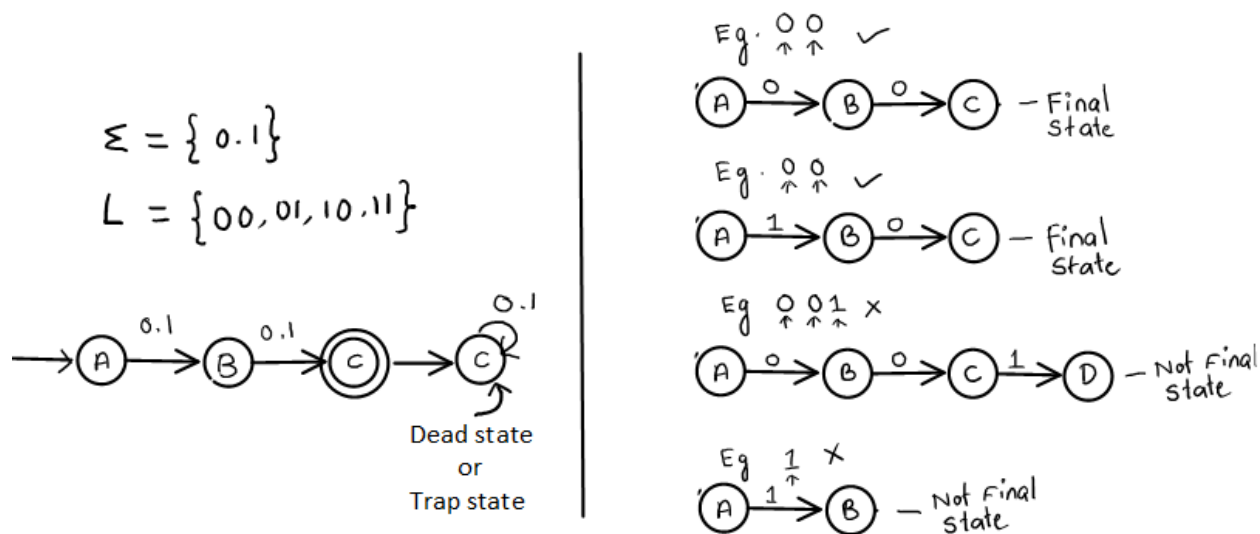Here two consecutive 1's or single 1 is acceptable , hence



The stages $q_0$ , $q_1$ , $q_2$ are the final states. The DFA will generate the strings that
Do not contain consecutive 1's like 10, 110, 101 ..............etc

5. Design DFA the set of all strings start with '0'

$$= \{\ 0,\ 00,\ 01,\ 000,\ 010,\ 011,\ 0000........\}$$

Eg. 001

Initial state — Final state

Initial state — Final state

Dead state
or
Trap state

6. Design DFA that accepts set of all strings over {0, 1} of length 2

$$\Sigma = \{0.1\}$$

$$L = \{00, 01, 10, 11\}$$

Dead state
or
Trap state

Eg. 00 ✓ — Final State

Eg. 00 ✓ — Final State

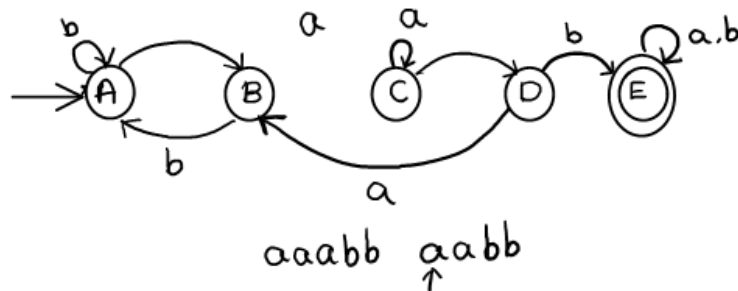Eg. 001 ✗ — Not Final state

Eg. 1 ✗ — Not Final state

7. Construct DFA that accepts any string {a, b} that doesnot contain aabb in it

$$\Sigma = \{a, b\}$$
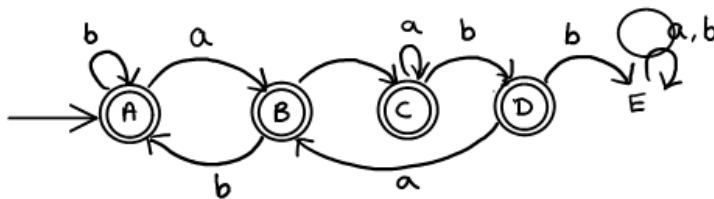
Try to design a simpler problem

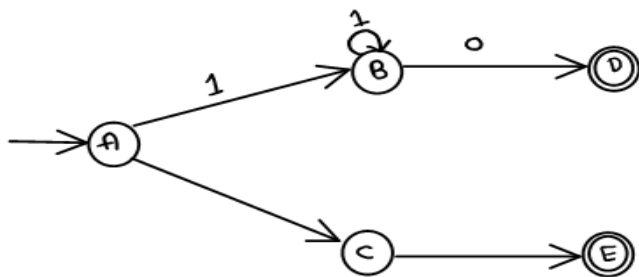Let us construct a DFA that accepts all strings over (a, b) that contains the string aa bb in it

Flip the states
- Make the final state Into non final state
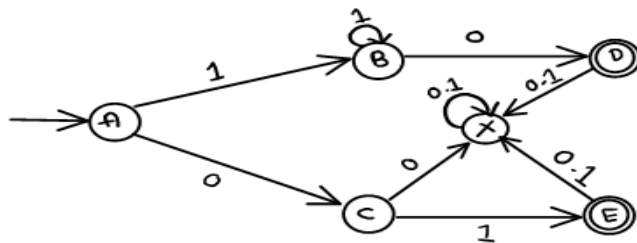
    and
- Make the non final state Into final states



aaabb aabb



8. How to figure out what following DFA recognizes ?

10

01
one binary digit 'd'

L = {Accepes the string 01 or a string of atleast
     One '1' followed by a 'o'}



10

01
one binary digit 'd'

L = {Acceptes the string 01 or a string of atleast
     One '1' followed by a 'o'}

Eg :- 001 , 010 , 011 , 1101 , 1100

**9.** Design a FA with ∑ = {0, 1} accepts the strings with an even number of 0's followed by single 1.
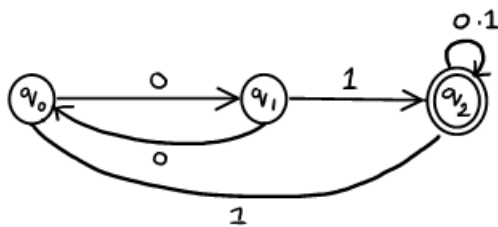
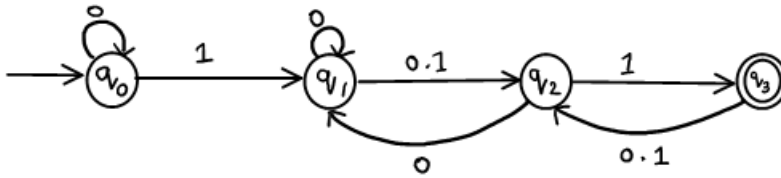The DFA Can be shown by a transition
diagram as

**Transition Table**: The transition table is basically a tabular representation of the transition function. It takes two arguments (a state and a symbol) and returns a state (the "next state").

A transition table is represented by the following things:

- Columns correspond to input symbols.

- Rows correspond to states.

- Entries correspond to the next state.

- The start state is denoted by an arrow with no source.

- The accept state is denoted by a star



| Present State | Next state For Input | Next state Of Input |
|---|---|---|
| $q_0$ | $q_1$ | $q_2$ |
| $q_1$ | $q_0$ | $q_2$ |
| *$q_2$ | $q_2$ | $q_2$ |

| Present State | Next state For Input (0) | Next state of Input (1) |
|---|---|---|
| $q_0$ | $q_0$ | $q_1$ |
| $q_1$ | $q_1, q_2$ | $q_2$ |
| $q_2$ | $q_1$ | $q_3$ |
| *$q_3$ | $q_2$ | $q_2$ |

Example: Vending Machine for Snacks



**States**: Ready, Dispensing

**Alphabet**: Buttons for different snacks

**Rules**:

• If in the "Ready" state and a snack button is pressed,

  **transition to "Dispensing**."

• If in the "Dispensing" state, transition back to "Ready" after

**dispensing**.

　　　Accept or Reject:

Acceptance when transitioning to "Dispensing."

**Figure 1.7: Vending machine of Snacks**



Another example will be Elevator Control :

**States**: Ground Floor, First Floor, Second Floor

**Alphabet**: Buttons for floor selection

**Rules**:

If on the "Ground Floor" and the button for the "First Floor"

　　is pressed, transition to "First Floor."

If on the "First Floor" and the button for the "Second Floor" is

　　pressed, transition to "Second Floor."

Accept or Reject: No specific acceptance or rejection; it's a system for floor navigation.

**DFA Conclusion:**

- In DFA, In given current state we know what the next state to be

- It has only one unique next state

- It is simple and easy to design

**Regular Language** : A language is said to be REGULAR Language if and only if some Finite State Machine recognizes it. In general, The language is not regular in case the language not recognized by FSM and requires memory(it cannot store FSM or Count strings).

Eg:- **abbab abbab** – In this language required to store the repetitive and memory required that's the reason we cannot store it and therefore it is not reqular language.

$a^n b^n$ -aaa bbb or aa bb – it means no of aa should have similar bb then we have to store n. which cannot stored by FSM.

UNION              -      A U B ={x/ x$\epsilon$A or x$\epsilon$B}

CONCATENATION    -      A **o** B ={xy/x$\epsilon$A or y$\epsilon$B}

STAR               -      A * ={$x_1 x_2 x_3$ ... ... ... $x_k$|k ≥ 0 and each x; ∈ A}

Eg      A={pq,r} ,B{ t,uv}

A U B    =    {pq, r, t, uv}

 A **o** B ={pqt, pquv, rt,ruv}

A* ={$\epsilon$, pq, r, pqr, rpq, pqpq, rr, pq, pqpqpq rrr, ... ... ... }

Theorem 1: The class of regular language is closed under UNION

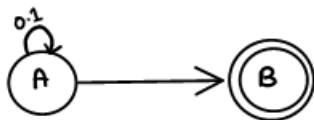Theorem 2: The class of regular language is closed under CONCATENATION

## 1. 6. Non-Deterministic Finite Automata (NFA)

A Non-deterministic Finite Automaton (NFA) is like a flexible version of a Deterministic Finite Automaton (DFA). While a DFA has precisely defined rules for each combination of state and input, an NFA can have multiple possible transitions for a given state and input. It allows for a bit more "guesswork" in the decision-making process.

- NFA stands for non-deterministic finite automata. It is easy to construct an NFA than DFA for a given regular language.

- The finite automata are called NFA when there exist many paths for specific input from the current state to the next state.

- Every NFA is not DFA, but each NFA can be translated into DFA.

- NFA is defined in the same way as DFA but with the following two exceptions, it contains multiple next states, and it contains ε transition



3 STATES -> $A^1$ -> A,B,C, AB,AC,BC,ABC, $\emptyset$ -> $2^3$

NFA also has five states same as DFA, but with different transition function, as shown follows

δ: Q x ∑ → $2^Q$

WHERE

Q: finite set of states

∑: finite set of the input symbol

q0: initial state

F: final state

δ: Transition function

## Non-deterministic finite Automata

**NON- DETERMINISM**

>> In NFA, given the current state there could be multiple next states.

>> The next state may be chosen at random.

>> All the next states may be chosen in Parallel.

Simple Calculator: Evaluating arithmetic expressions.

States: Start, Number, Operator, Accept

Alphabet: Digits, +, -, *, /

Rules:

• Start reading a number.

• Transition to the operator state after reading a number.

• Allow multiple paths for different sequences of numbers and operators.

Accept or Reject:

Accept if it reaches the "Accept" state, indicating a valid expression; otherwise, reject.

1. Design NFA that all strings that end with 0

NFA- Example-1

0.1



L=(set of all strings that end with o)

Eg :- 100



>> If there is anyway run the machine that end in any set of states out of which at least one state in a final state the NFA accepts.



×

2. Design NFA that set of all string start with 0

L= ( set of all strings that start with 0 )
    = ( 0, 00, 01, 000............)



Eg:001



Eg :101

 Dead configuration

3. Design NFA accept set of all string of length 2

Design NFA accept set of all string of length 2

Eg : 101

$\rightarrow \text{(A)} \xrightarrow{1} \phi$  Dead configuration

>> construct a NFA that accepts sets of all strings over {0, 1} of length 2

$\Sigma$ = {0, 1}
L = {00, 01, 10, 11}

$\rightarrow \text{(A)} \xrightarrow{0,1} \text{(B)} \xrightarrow{0,1} \text{(C)}$

Eg:00 ✓

$\rightarrow \text{(A)} \xrightarrow{0} \text{(B)} \rightarrow \text{(C)}$
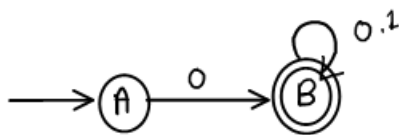
Eg:001 ✗

$\rightarrow \text{(A)} \xrightarrow{0} \text{(B)} \xrightarrow{0} \text{(C)} \xrightarrow{1} \phi$

4. Design NFA that accept set of all string that ends with '1'

01, 001, 0001, 0*1, 1, 101, 1101

$\rightarrow \text{(A)} \xrightarrow{0,1} \xrightarrow{1} \text{(O)}$

5. Design NFA that all string input contains '0'

$\rightarrow \text{(A)} \xrightarrow{0,1} \xrightarrow{0} \text{(B)} \xrightarrow{0,1}$

2. Design NFA that all string starts with '0'

$\rightarrow \text{(A)} \xrightarrow{1} \text{(B)} \xrightarrow{0} \text{(C)} \xrightarrow{0,1}$

3.  Design NFA that all string that contains '01'



4.  Design NFA that all string that ends with 11



5.  Design NFA with $\sum$={0,1} accepts all strings with 01



Transition table

| Present State | Next state For Input (0) | Next state Of Input (1) |
|---|---|---|
| $-q_0$ | $q_1$ | $\varepsilon$ |
| $q_1$ | $\varepsilon$ | $q_2$ |
| $* q_2$ | $q_2$ | $q_2$ |

Design NFA with $\sum$ = {0, 1} and accept all string of length atleast 2.

6. Design NFA with $\Sigma=\{0,1\}$ and accept all string of length at least 2



Transition table

| Present State | Next state For Input (0) | Next state of Input (1) |
|---|---|---|
| $-q_0$ | $q_1$ | $q_1$ |
| $q_1$ | $q_2$ | $q_2$ |
| $* q_2$ | $\varepsilon$ | $\varepsilon$ |

# 2. NFA TO DFA Conversion

Converting a Non-Deterministic Finite Automaton (NFA) to a Deterministic Finite Automaton (DFA) involves systematically exploring all possible states and transitions in the NFA to create an equivalent DFA

Every DFA is on NFA but vice versa but there is an equivalent DFA for NFA

$$DFA = Q \times \Sigma \longrightarrow Q$$

$$NFA \quad \partial = Q \times \Sigma \longrightarrow 2^a$$

$$NFA \cong DFA$$

**Q: Why do you need of Conversion from NFA to DFA ?**

Converting a Non-Deterministic Finite Automaton (NFA) to a Deterministic Finite Automaton (DFA) is often done for practical reasons in the field of computer science and formal language theory. Here are some reasons why this conversion is important in the real world:

- **Implementation of Lexical Analyzers (Compilers):** Compilers often use DFAs to perform lexical analysis, which involves tokenizing source code into meaningful units. Regular expressions, which are commonly used in lexical analysis, can be more efficiently implemented with DFAs.

- **Efficient Pattern Matching**: DFAs are generally faster than NFAs when it comes to pattern matching. Many string-matching algorithms rely on DFAs for efficient searching and matching operations

- **Parsing and Syntax Analysis**: DFAs are used in parsing techniques for syntax analysis of programming languages. Some parsing algorithms, like LL(1)( **Left-to-Right, Leftmost derivation with 1 symbol lookahead**) and LR(1)( **Left-to-Right, Rightmost derivation with 1 symbol lookahead**.), are designed to work with DFAs, making the conversion from NFA to DFA an essential step in language processing.

- **Network Protocol Analysis**: DFA-based models are used to analyze and validate network protocols. DFAs can efficiently recognize and process patterns in network traffic

  .

- **Regular Expression Matching in Text Editors**: Text editors and search engines often use DFAs for efficient pattern matching. Converting regular expressions to DFAs allows for faster searching and matching in large text documents.

- **Modeling and Verification**: DFAs are used in formal methods for modeling and verifying systems. They can represent state machines in a deterministic manner, making it easier to analyze and reason about system behaviour

  .

- **Optimizing Finite State Machines**: In embedded systems and hardware design, DFAs are used to model and optimize finite state machines. DFAs can be implemented more efficiently in hardware than NFAs.

- **Database Query Optimization**: Query optimizers in databases may use DFAs to efficiently process regular expressions and pattern-matching queries.

In summary, the conversion of NFA to DFA is crucial in various applications where efficient pattern matching, parsing, and recognition of languages play a significant role. DFAs offer advantages in terms of simplicity, determinism, and efficiency, making them well-suited for practical implementations in real-world systems and applications.

**Q: Convert NFA to DFA L={set of all strings over (0,1) that starts with '0' }**

$\Sigma = \{ 0, 1\}$

NFA



| | 0 | 1 |
|---|---|---|
| A | B | $\phi$ |
| B | B | B |

DFA



| | 0 | 1 |
|---|---|---|
| A | B | C |
| B | B | B |
| C | C | C |

C - Dead state/Trap state

**Q: Convert NFA to DFA L={set of all strings over (0,1) that ends with '1' (SUBSET CONSTRUCTION METHOD)**

L = { set of all strings over (0,1) that ends with '1' }

$\Sigma$ ={ 0, 1}

NFA



| | 0 | 1 |
|---|---|---|
| A | {A} | {A.B} |
| B | B | B |

DFA

AB - single state



| | 0 | 1 |
|---|---|---|
| A | {A} | {AB} |
| AB | {A} | {AB} |

Find the equivalent DFA for NFA given by

M = [ {A, B, C} , (a, b) $\delta$ , A{c} ]   where $\delta$ is given by

|   | 0 | 1 |
|---|---|---|
| A | A,B | C |
| B | A | B. |
| Ⓒ | – | A,B |



|   | 0 | 1 |
|---|---|---|
| A | AB | C |
| B | AB | BC |
| Ⓑⓒ | A | AB |
| Ⓒ | D | AB |
| D | D | D |



Q : Using below diagram convert NFA to DFA

given below is the NFA for a Language

L = { set of all strings over (0.1) that ends with '0.1' }
construct the equivalent DFA



|   | 0 | 1 |
|---|---|---|
| –A | A,B | C |
| B | $\phi$ | C |
| Ⓒ | $\phi$ | $\phi$ |

Given below is the NFA for a language

L = { set of all strings over (0,1) that ends with '01'}

   construct its equivalent DFA

NFA



| | 0 | 1 |
|---|---|---|
| $\rightarrow$A | A,B | C |
| B | $\phi$ | C |
| Ⓒ | $\phi$ | $\phi$ |

DFA



1101
1110

| | 0 | 1 |
|---|---|---|
| $\rightarrow$A | A,B | A |
| AB | AB | AC |
| AC | AB | A |

Q. Design an NFA for a language that accepts all strings over { 0, 1 }in which The second last symbol is always '1'

   Then covert it to its equivalent DFA

NFA



| | 0 | 1 |
|---|---|---|
| $\rightarrow$A | A,B | C |
| B | C | C |
| Ⓒ | $\phi$ | $\phi$ |

Eg:- 1010
110
1101010

:

Q . Design an NFA for language that accepts all strings over { 0, 1 } in which the Second last symbol is always '1'

Then convert it to equivalent DFA

NFA



| | O | 1 |
|---|---|---|
| −A | A,B | C |
| B | C | C |
| Ⓒ | φ | φ |

Eg:- 1010
110
1101010

DFA



| | O | 1 |
|---|---|---|
| →A | A | AB |
| AB | AC | ABC |
| (AC) | A | AB |
| (ABC) | AC | ABC |

# 3. Minimization of DFA

Minimization of a DFA (Deterministic Finite Automaton) is a process that involves reducing the number of states in the DFA while preserving its ability to recognize the same language. The goal is to create a smaller DFA that is equivalent to the original DFA in terms of the language it recognizes.

**DFA minimization process:**

- **Equivalent States Identification**: Begin by identifying pairs of states that are equivalent. Two states are considered equivalent if, for every input symbol, they lead to equivalent states and both are either accepting or non-accepting.

  States are equivalent if, upon reading any string, they lead to states that are also equivalent.

- **State Merging**: Merge the equivalent states into a single state. This process reduces the number of states in the DFA. The transitions from the merged states are determined by the transitions from the original states.

- **Update Transitions**: Update the transitions of the remaining states to account for the merged states. For each input symbol, the transition from the merged state is the state resulting from merging the transitions of the original states.

- **Repeat**: Repeat the process of identifying equivalent states and merging them until no more merging can be done. This ensures that the DFA is minimized to its smallest possible form.

- **Result**: The final minimized DFA is obtained, which has the smallest number of states while still recognizing the same language as the original DFA.

Minimizing a DFA is important for various reasons, including simplifying the representation of the language, improving efficiency in terms of both time and space, and ensuring that the DFA is as concise as possible while maintaining its language recognition capabilities. The minimization process is a key concept in automata theory and formal language theory.

## DFA

5 states                                    4 states

Equivalent

Two states 'A' and 'B' are said to be equivalent if

$$\partial(A, X) \longrightarrow F$$
and
$$\partial(B, X) \longrightarrow F$$

Or

$$\partial(A, X) \longleftrightarrow F \quad \text{\{where 'x' is any input string\}}$$
$$\partial(B.K) \longleftrightarrow F$$

Types of Equivalent

If $|X| = 0$, then A and B are said to be 0 equivalent

If $|X| = 1$, then A and B are said to be 1 equivalent

If $|X| = 2$, then A and B are said to be 2 equivalent

…

..
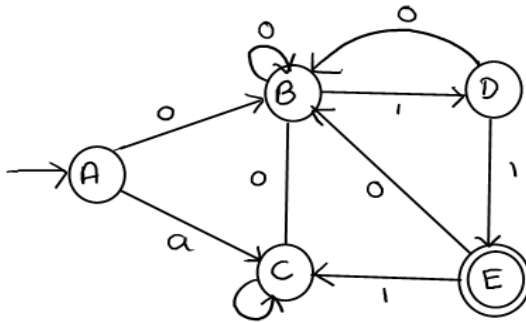
If $|X| = n$, then A and B are said to be n equivalent

**Why Minimization of DFA?**

- **Simplicity and Understandability**: Minimizing a DFA results in a simpler and more compact representation. The reduced number of states makes it easier for humans to understand and analyze the structure of the automaton. This is especially important for manual inspection and comprehension.

- **Efficiency**: A minimized DFA often leads to improved efficiency in terms of both time and space. With fewer states, the DFA requires less memory to represent, leading to reduced storage requirements. Additionally, processing and recognizing strings may be faster with a minimized DFA, as there are fewer states to traverse.

- **Optimization of Resources**: In practical applications, such as compilers, regular expression matching, and other language recognition systems, minimizing DFAs is crucial for optimizing resource usage. Smaller DFAs consume less memory and computational resources, making them more suitable for embedded systems and environments with limited resources.

- **Equivalence**: The minimized DFA is equivalent to the original DFA in terms of the language it recognizes. This means that both DFAs accept the same set of strings. Minimization ensures that no states are redundant, and the DFA is as small as possible while still recognizing the language.

- **Algorithmic Efficiency**: Minimization algorithms, such as Hopcroft's algorithm or the Myhill-Nerode theorem, provide systematic and algorithmic approaches to optimize DFAs. These algorithms are well-studied and efficient, contributing to the development of automated tools for DFA minimization.

- **Language Complexity Analysis**: Minimizing a DFA can provide insights into the complexity of the language it recognizes. The size of the minimized DFA is related to the complexity of the regular language it represents, and this information is valuable in the study of formal languages and automata theory.

In summary, the minimization of DFAs is a crucial step in the design and optimization of systems involving pattern matching, language recognition, and automata-based algorithms. It contributes to the efficiency, simplicity, and overall effectiveness of these systems in various computational applications

A. **Minimize the given following DFA**

Step 1 : construction state transition table

|  | O | 1 |
|---|---|---|
| →A | B | C |
| B | B | D |
| C | B | C |
| D | B | E |
| Ⓔ | B | C |

**STEP 2: Compute Equivalence**

Equivalence 0: {A,B,C,D}    {E}

Equivalence 1: {A,B,C} {D}, {E} →A,B , A,C (EQUIVALENCE)    C,D →BOT DESTINATION IN DIFFERENT SET

Equivalence 2: {A, C} ,{B},{D},{E}

Equivalence 3: {A,C}, {B},{D}, {E}

Equivalence 2 and 3 are same set. We should stop the prcess.

**STEP3 : Draw minimized DFA**

| | 0 | 1 |
|---|---|---|
| →A | B | C |
| B | B | D |
| C | B | C |
| D | B | E |
| Ⓔ | B | C |

**B.   MINIMIZE THE DFA FROM BELOW TRANSITION TABLE**

| | 0 | 1 |
|---|---|---|
| $-q_0$ | $q_1$ | $q_5$ |
| $q_1$ | $q_6$ | $q_2$ |
| $q_2$ | $q_0$ | $q_2$ |
| $q_3$ | $q_2$ | $q_6$ |
| $q_4$ | $q_7$ | $q_5$ |
| $q_5$ | $q_2$ | $q_0$ |
| $q_6$ | $q_6$ | $q_4$ |
| $q_7$ | $q_6$ | $q_2$ |

STEP 1: Compute equivalences

Q:
DFA
one

Design   Minimize
when   more   than
final states

O Equivalence

$$\{q_0, q_1, q_2, q_3, q_4, q_5, q_6, q_7\} \{q_2\}$$

I – Equivalence

$$\{q_0, q_4, q_6\}$$

$$\{q_1, q_7\}$$

$$\{q_3 q_5\} \quad \{q_2\}$$

2 - Equivalence

$\{q_0, q_4\}$ $\{q_6\}$ $\{q_1, q_7\}$ $\{q_3, q_5\}$ $\{q_2\}$

3 - Equivalence

$\{q_0, q_4\}$ $\{q_6\}$ $\{q_1, q_7\}$ $\{q_3, q_5\}$ $\{q_2\}$

Step 2 :- Design DFA using above equivalence

|  | 0 | 1 |
|---|---|---|
| $(q_0, q_4)$ | $\{q_1, q_7\}$ | $\{q_3, q_5\}$ |
| $q_6$ | $\{q_6\}$ | $\{q_0, q_4\}$ |
| $(q_1, q_7)$ | $\{q_6\}$ | $\{q_2\}$ |
| $(q_3, q_5)$ | $\{q_2\}$ | $\{q_6\}$ |
| $(q_2)$ | $\{q_0, q_4\}$ | $\{q_2\}$ |

|  | 0 | 1 |
|---|---|---|
| $-q_0$ | $q_1$ | $q_5$ |
| $q_1$ | $q_6$ | $q_2$ |
| $(q_2)$ | $q_0$ | $q_2$ |
| $q_3$ | $q_2$ | $q_6$ |
| $q_4$ | $q_7$ | $q_5$ |
| $q_5$ | $q_2$ | $q_0$ |
| $q_6$ | $q_6$ | $q_4$ |
| $q_7$ | $q_6$ | $q_2$ |

Minimize the following DFA



| | O | I |
|---|---|---|
| A | B | C |
| B | A | D |
| C | E | F |
| D | E | F |
| E | E | F |
| F | F | F |

0 - Equivalence - {A,B,F} {C,D,E}
1 - Equivalence - {A,B} {F} {C,D,E}



| | O | I |
|---|---|---|
| {A,B} | {A,B} | {C,D,E} |
| {F} | {F} | {F} |
| {C,D,E} | {C,D,E} | {F} |

**Q: Design DFA when unreachable state**



A state is said to be Unreachable if there Is no way it can be reached from the initial state

**Step 1: remove unreachable state F and write Transition table excluding F**

Step 1 :- Remove unreachable sate F and write transition
table excluding F



|   | 0 | 1 |
|---|---|---|
| - A | B | C |
| ⑧ | D | E |
| © | E | D |
| D | G | G |
| E | G | G |
| F | G | G |

0 - Equivalence : {A,D,E} {B,C,G}

1 - Equivalence : {A, DE} {B,C} {G}

2 - Equivalence : {A} {D,E} {B, C} {G}

3 - Equivalence : {A}, {D,E} {B,C} {G}

Step 2 :- construct transition table

|   | 0 | 1 |
|---|---|---|
| -{A} | {B,C} | {B,C} |
| {D,E} | {G} | {G} |
| {B,C} | {D,E} | {D,E} |
| {G} |  | {G} |

# 4. Epsilon-NFA

An ε-NFA, or epsilon-Nondeterministic Finite Automaton, is a type of finite automaton that extends the capabilities of a nondeterministic finite automaton (NFA) by allowing transitions labeled with the empty string, denoted as ε (epsilon). In regular NFAs, transitions are made based on input symbols, but in ε-NFAs, transitions can also be made without consuming any input symbol.

The transition function of an ε-NFA includes ε-transitions, which means the automaton can move from one state to another without reading any input symbol. This non-deterministic behavior enhances the expressive power of the automaton and makes it more flexible in recognizing certain types of languages.

Formally, an ε-NFA is defined by a 5-tuple (Q, Σ, δ, $q_o$, F), where:

Q is a finite set of states.

Σ is a finite set of input symbols (the alphabet).

δ is the transition function, which maps Q × (Σ ∪ {ε}) to 2^Q (the power set of Q).

$q_o$ is the initial state.

F is the set of accepting states.
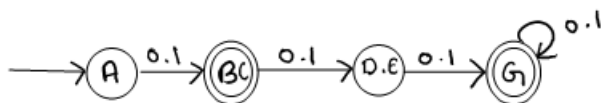
The key difference from a regular NFA is that ε can be used as an input symbol in transitions. This allows the automaton to transition to new states without consuming any input, giving it more flexibility in recognizing languages.

The transition function δ is defined for ε-NFAs as follows:

δ: Q × (Σ ∪ {ε}) → 2^Q

In practice, ε-NFAs are often used in the theory of formal languages and automata to describe and analyze certain types of languages. They are particularly useful in the conversion process between different types of automata, such as converting regular expressions to NFAs or DFAs.



**Applicability of Epsilon-NFA:** An ε-NFA (epsilon-Nondeterministic Finite Automaton) is a type of non-deterministic finite automaton where transitions can be taken without consuming any input symbol,

represented by the epsilon (ε) symbol. The use of ε-NFA is primarily found in theoretical computer science and formal language theory. Some key applications include:

- **Regular Expression to NFA Conversion**: ε-NFAs are often used in the conversion process from regular expressions to nondeterministic finite automata (NFAs). This conversion helps in understanding and implementing regular expression matching algorithms.
- **Theoretical Computations**: ε-NFAs are used in theoretical computer science to study the properties of regular languages. They provide a formalism for understanding the structure of regular languages and their relationships with other types of formal languages.
- **Language Recognition**: In the context of formal language theory, ε-NFAs are employed to recognize and define regular languages. They can be used to model and recognize patterns in strings that conform to a certain regular language.
- **Automata Theory**: ε-NFAs are an essential concept in automata theory, helping researchers and students understand the relationship between different types of automata and formal languages.

In real-time applications, the direct use of ε-NFAs may be less common. However, the concepts and theoretical insights gained from studying ε-NFAs can have practical implications in various areas of computer science and software engineering. For example:

**Compiler Design**: Techniques from formal language theory, including ε-NFAs, are applied in compiler design for lexical analysis and pattern matching.

**Text Processing and Search Algorithms**: Regular expression engines often use the concepts of NFAs, including ε-NFAs, to efficiently search for patterns in text.

**Parsing and Syntax Analysis**: Concepts related to NFAs can be applied in parsing and syntax analysis in the development of compilers and interpreters.

**String Matching Algorithms**: Algorithms that involve pattern matching, such as those used in text processing or search engines, may draw inspiration from the theoretical foundations provided by ε-NFAs.

While ε-NFAs themselves may not be directly implemented in real-time systems, the theoretical understanding gained from their study contributes to the development of efficient algorithms and tools in various areas of computer science

Conversion of E- NFA to NFA

Convert the following E- NFA to its equivalent NFA



State    E* input E*
E- closure (e-) all the states
that can be reached from a
particular state only by
seeing the E symbol

|   | 0 | 1 |
|---|---|---|
| →A | {A,B,c} | {B.c} |
| B | {c} | {B.c} |
| C |   |   |

$$A \quad \begin{matrix} E^* & O & E^* \\ A - A - A \\ B - \phi - \overset{*}{c} \\ C - C - C \end{matrix}$$

$$A \quad \begin{matrix} E^* & O & E^* \\ A - \phi - \\ B - B - B.C \\ C - C - C \end{matrix}$$

$$A \quad \begin{matrix} E^* & O & E^* \\ B - B - BC \\ C - C - C \end{matrix}$$

$$A \quad \begin{matrix} E^* & O & E^* \\ - & - \\ C - C - C \end{matrix}$$

|   | {A,B,c} | {B,c} |
|---|---|---|
| →Ⓐ | {A.B.c} | {B.c} |
| Ⓑ | {c} | {B.c} |
| Ⓒ | {c} | {c} |

**Q: Design NFA**

Convert the following - E- NFA to its equivalent NFA



| | $\overset{*}{E}$ | O | $\overset{*}{E}$ |
|---|---|---|---|
| A | A | B | C |
| B | B C | $\phi$ C | $\phi$ C |
| C | C | C | C |
| D | D | $\phi$ | $\phi$ |

| | $\overset{*}{E}$ | I | $\overset{*}{E}$ |
|---|---|---|---|
| A | A | $\phi$ | $\phi$ |
| B | B C | B D | B C D |
| C | C | D | D |
| D | D | $\phi$ | $\phi$ |

NFA

| | O | I |
|---|---|---|
| →A | B,C | $\phi$ |
| B | C | B,C,D |
| C | C | D |
| D | $\phi$ | $\phi$ |

Computed NFA for above



**Q: Convert EPSILAONT TO NFA**

Conversion of E- NFA to NFA - examples (part-D)

Conversion the following E- NFA to its equivalent NFA



| | $\overset{*}{E}$ | a | $\overset{*}{E}$ |
|---|---|---|---|
| P | P Q R | φ P φ | φ P Q R |
| Q | Q | P | P Q R |
| R | R | φ | — |

| | $\overset{*}{E}$ | b | $\overset{*}{E}$ |
|---|---|---|---|
| P | P Q R | Q R φ | Q R — |
| Q | Q | R | R |
| R | R | φ | — |

| | $\overset{*}{E}$ | c | $\overset{*}{E}$ |
|---|---|---|---|
| P | P Q R | R Q Q | R Q — |
| Q | Q | Q | Q |
| R | R | Q | — |

| | a | b | c |
|---|---|---|---|
| -P | {P,Q,R} | {Q,R} | {Q,R} |
| Q | {P,Q,R} | {R} | {Q} |
| ®R | φ | φ | φ |

# 5. Compiler Design:

Compiler design is a crucial field in computer science that focuses on creating software systems known as **compilers**. Compiler design is a complex and challenging field that requires a deep understanding of computer architecture, programming languages, and optimization techniques. Efficient compilers play a crucial role in software development by translating high-level abstractions into machine-executable code, enabling the execution of diverse applications on various computing platforms.

**Compiler** : A compiler is a specialized software tool that translates the entire source code of a high-level programming language into an equivalent machine code or an intermediate code. The primary purpose of a compiler is to enable the execution of a program on a computer by converting it from a human-readable and portable form into a form that can be executed by a computer's central processing unit (CPU).



**Interpreter**: An interpreter is a type of language processor that reads and executes a program or script line by line, without the need for a separate compilation step. Instead of translating the entire source code into machine code before execution, an interpreter processes the source code one statement at a time, interpreting and executing it directly.

**Q: Explain the differences between compiler and Interpreter ?**

| Feature | Compiler | Interpreter |
|---|---|---|
| Execution Approach | Translates the entire code before execution, producing an independent executable file or intermediate code. | Interprets the code line by line or statement by statement in real-time without producing a separate executable. |
| Output | Generates an executable file or intermediate code. | Does not generate a separate executable; executes source code directly. |
| Efficiency | Often results in more efficient code due to optimization during compilation. | May have slower execution compared to compiled code, as it interprets code in real-time. |
| Debugging | Debugging is typically done on the source code level, which can be more challenging. | Allows for easier debugging, often providing interactive debugging tools. |

| | | |
|---|---|---|
| **Memory Usage** | May consume more memory due to standalone executables or larger compiled code. | Generally has lower memory overhead as it doesn't produce a separate executable. |
| **Portability** | Compiled code may be platform-dependent. | The interpreter itself is platform-independent, but the source code may need an interpreter for each platform. |
| **Examples** | GCC (GNU Compiler Collection), Clang. | Python, Ruby, JavaScript. |
| **Execution Speed** | Can result in faster execution once compiled. | Tends to have slower execution compared to compiled code. |
| **Modification & Execution** | Code needs to be recompiled after modifications, and the new executable is run. | Modifications can be directly applied to the source code, and changes take effect immediately during execution. |

## 10.1 Language Processing System:

The computer functions as an intelligent synergy of both hardware and software. Hardware, essentially mechanical equipment, operates by following instructions provided by the relevant software. In this context, hardware interprets instructions as electronic charges, akin to the binary language employed in software programming, which consists solely of 0s and 1s. To elaborate, writing instructions directly in binary format—represented by a sequence of 0s and 1s—would be a cumbersome and intricate task for computer programmers. As a more user-friendly alternative, programmers write programs in high-level languages that are easier for them to understand and remember.

*These programs are then processed through a series of devices and components within the operating system (OS) to generate the desired code usable by the machine. This comprehensive process is known as* a **language processing system.**

To Generate the machine code following stages

- **Pre-processor**: A source program may be divided into modules stored in separate files. The task of collecting the source program is entrusted to a separate program called the pre-processor. It may also expand macros into source language statements.
- **Compiler** : A compiler is a program that takes a source program as input and produces an assembly language program as output.
- **Assembler**: Assembler is a program that converts assembly language programs into machine language programs. It produces re-locatable machine code as its output.
- **Loader and link-editor** :
  a. The re-locatable machine code has to be linked together with other re-locatable object files and library files into the code that actually runs on the machine.
  b. The linker resolves external memory addresses, where the code in one file may refer to a location in another file.
  c. The loader puts together the entire executable object files into memory for execution.

**Figure: A language-processing System**

```
              Source program
                    │
                    ▼
            ┌─────────────────┐
            │  Pre - Processor │
            └─────────────────┘
                    │   Source
                    ▼   Program
            ┌─────────────────┐
            │    Compiler     │
            └─────────────────┘
                    │  Target assembly program
                    ▼
            ┌─────────────────┐
            │    Assembler    │
            └─────────────────┘
                    │  Re lecatable machine program
                    ▼
            ┌─────────────────────┐         Library, Relocatable
            │  Loader link editor │ ◄──────  – Object files
            └─────────────────────┘
                    │
                    ▼
            ┌─────────────────┐
            │  Target machine │
            └─────────────────┘
```

## 10.2 Phases of Compiler :

- **Lexical analysis**: It scans the source code and converts the source code into tokens. Here, input is source code and output is stream of tokens.
- **Syntax Analysis**: After converting the source code into tokens, the next phase comes i.e. syntax analysis. It checks the grammatical mistakes of the source code. To verify the syntax of the source code the language must be defined by CFG. It takes streams of input as input and generates a parse tree.
- **Semantic Analysis**: It verifies the meaning of every sentence by performing type check. The syntax analyzer just verifies whether the operator is operating on a required number of operands or not.
- **Intermediate Code Generator:** The source code is converted into intermediate representation to make the code generation process simple and easy.
- **Code Optimization**: This process includes in reducing the number of instructions without affecting the outcome of the source program.
- Target Code Generator: The optimization source code should convert into assembly code.

let's walk through each phase of compilation with a simple example. We'll use a basic code snippet for illustration:

int main() {

    int a = 5;

    return a * 2;

}

**1. Lexical Analysis:**

Tokens:

[INT, MAIN, OPEN_PAREN, CLOSE_PAREN, OPEN_BRACE, INT, IDEN(a), ASSIGN, NUM(5), SEMICOLON, RETURN, IDEN(a), STAR, NUM(2), SEMICOLON, CLOSE_BRACE]

**2.Syntax analysis :**

Syntax Tree:

   - Function: main

     - Declaration: int

- Body:

- Assignment: a = 5

- Return: a * 2

**3. Semantic Analysis:**

Semantic Analysis:

- Check: 'main' is declared as a function.

- Check: 'a' is declared as an integer variable.

- Check: Return statement has an integer expression.

**4.Intermediate Code Generation:**

1. allocate space for variables

2. a = 5

3. t1 = a * 2

4. return t1

**5. Code Optimization**

Optimized Intermediate Code:

1. allocate space for variables

2. t1 = 10

3. return t1

**6. Target Code Generation**

Assembly Code:

```
MOV eax, 10      ; store 10 in register eax

RET              ; return from the function
```

In this example, the code snippet declares a simple main function with an integer variable a. It assigns a value to a and returns the result of a * 2. Each phase of the compiler processes the source code, transforming it into a more refined representation until the final assembly code is generated

# Assignments and Practices

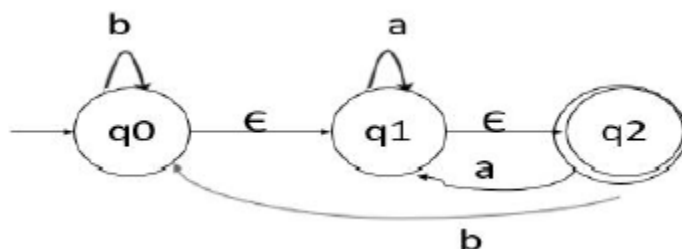1. Consider the below finite automata and check whether the strings are accepted or not

| States (Q) | Input Alphabtes | |
|---|---|---|
| | 0 | 1 |
| →q0 | q1 | q3 |
| q1 | q0 | q2 |
| (q2) | q3 | q1 |
| q3 | q2 | q0 |

1. 0001        2.1010        3.1001        4. 0101

2. Define alphabets, strings and Languages?
3. Compare DFA V/S NFA
4. Construct DFA for given NFA

| | Next state | |
|---|---|---|
| | 0 | 1 |
| → q0 | q0,q1 | q0 |
| q1 | q2 | q1 |
| q2 | q3 | q3 |
| (q3) | - | q2 |

5. Convert the following NFA with ε moves to DFA.



6. Design DFSM to accept the each of the following language
   i) $L = \{w \in \{a, b\}^* ; w$ has all strings that ends with sub string abb $\}$
   ii) $L = \{w;$ where $|w| \bmod 3 = 0$ where $\Sigma = \{a\}\}$
   iii) $L = \{w \in \{a, b\}^*$ every a region in w is of even length.$\}$

7. Construct an equivalent DFA from the following given NFA using subset construction method
8.

9. Construct the following NFA with epsilon -moves construct equivalent DFA



Fig.Q.2(b)

https://nptel.ac.in/courses/106/106/106106049/#

https://nptel.ac.in/courses/106/104/106104123/#

www.jflap.org

# 2 Regular Expressions and Languages

## 2.1 Regular Expressions

Regular expressions, often abbreviated as regex or regexp, are powerful and concise sequences of characters used to define search patterns. They are widely employed in various programming languages, text editors, and tools for string manipulation, pattern matching, and data validation. Here are some fundamental concepts related to regular expressions:

▪ **Character Classes:**

    ▪ [ ]: Defines a character class. Matches any one of the characters inside the brackets. For example, [aeiou] matches any vowel.

▪ **Quantifiers:**

    ▪ **\*:** Matches zero or more occurrences of the preceding character or group.

    ▪ +: Matches one or more occurrences of the preceding character or group.

    ▪ ?: Matches zero or one occurrence of the preceding character or group.

▪ **Anchors:**

    ▪ ^: Anchors the regex at the beginning of a line.

    ▪ $: Anchors the regex at the end of a line.

▪ **Wildcards:**

    ▪ . (dot): Matches any single character except a newline.

▪ **Escape Characters:**

    • \: Escapes a metacharacter, allowing it to be treated as a literal character.

▪ **Grouping and Capturing:**

    • ( ): Groups characters together. Also used for capturing substrings.

▪ **Alternation:**

    • | (pipe): Represents alternation, allowing the matching of either the expression before or after it.

▪ **Quantifiers:**

    • {n}: Matches exactly n occurrences of the preceding character or group.

    • {n,}: Matches n or more occurrences.

Regular Expression: b{2,}

- {n,m}: Matches between n and m occurrences.

- **Character Escapes:**

  - \d: Matches any digit (0-9).

  - \w: Matches any word character (alphanumeric + underscore).

  - \s: Matches any whitespace character (space, tab, newline).

- **Negation:**

  - [^ ]: Negates a character class, matching any character not inside the brackets.

- **Anchors:**

  - \b: Matches a word boundary.

  - \B: Matches a non-word boundary.

- **Modifiers:**

  - i: Case-insensitive matching.

  - g: Global matching (find all matches rather than stopping after the first).

Regular expressions provide a concise and expressive way to describe complex patterns in strings, making them a versatile tool for tasks like text searching, data validation, and text manipulation. They are supported in many programming languages, including Python, JavaScript, Java, and others, often through libraries or native language support.

Regular expressions (regex or regexp) can be understood in terms of Finite State Machines (FSMs), which are theoretical models used to represent and recognize patterns in strings. **Regular expressions and FSMs are closely related because regular expressions define patterns that can be recognized by FSMs.**

- The language accepted by finite automata can be easily described by simple expressions called Regular Expressions. It is the most effective way to represent any language.
- The languages accepted by some regular expressions are referred to as Regular languages.
- A regular expression can also be described as a sequence of patterns that defines a string.
- Regular expressions are used to match character combinations in strings. The string searching algorithm used this pattern to find the operations on a string.

In the context of regular expressions and FSMs:

**Alphabet** ($\Sigma$): The alphabet consists of all possible symbols that can appear in the input string. In regular expressions, these symbols are the characters or tokens that make up the strings you want to match.

**States (Q):** States represent different stages in the recognition process. Each state corresponds to a specific part of the pattern that has been recognized so far.

**Transitions (δ):** Transitions describe the changes of state based on input symbols. In the context of regular expressions, transitions represent the progression through the pattern.

**Accepting States (F):** Accepting states indicate successful recognition of the pattern. In the case of regular expressions, an accepting state is reached when the entire input string has been matched according to the pattern.
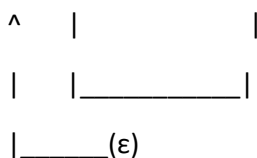
Now, let's consider a simple example to illustrate the connection between regular expressions and FSMs:

For instance:

**In a regular expression, x\* means zero or more occurrence of x. It can generate {e, x, xx, xxx, xxxx, .....}**

In a regular expression, "x\*" means zero or more occurrences of 'x'. It can generate the set {ε, x, xx, xxx, xxxx, ...}, where ε represents an empty string. Here's how you can interpret it in terms of a Finite State Machine (FSM):

Start State ---(x)---> Accepting State

```
    ^     |              |
    |     |_____|
    |_____(ε)
```

Additionally, the epsilon transition (ε) allows the machine to move from the start state to the accepting state without consuming any 'x', allowing for an empty string match.

the pattern "x\*" means zero or more occurrences of the character 'x'. Here are some sample examples:

Pattern: a\*      Matches: "", "a", "aa", "aaa", ...

Pattern: \d\* (zero or more digits)     Matches: "", "0", "123", "987654", ...

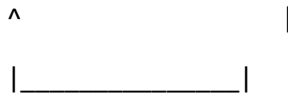Pattern: [A-Za-z]\* (zero or more letters) Matches: "", "abc", "XYZ", "Hello", ...

Pattern: .\* (zero or more of any character)     Matches: "", "x", "xy", "xyz", "abc", "123", ...


**In a regular expression, x+ means one or more occurrences of x. It can generate {x, xx, xxx, xxxx, .....}**

In terms of a Finite State Machine (FSM), the regular expression "x+" (one or more occurrences of 'x') can be represented by a simple state machine with two states: the start state and the accepting state. The transition from the start state to the accepting state occurs on each occurrence of the character 'x'. Once in the accepting state, the machine remains in that state for any subsequent occurrences of 'x'. This reflects the idea of matching one or more occurrences of 'x'.

Here is a simple illustration of the FSM for the regular expression "x+":


Start State ----(x)----> Accepting State

   ^                |

   |_____|


This FSM will accept sequences like "x", "xx", "xxx", "xxxx", and so on, but it won't accept the empty string.

In a regular expression, the notation "x+" means one or more occurrences of the character or pattern 'x'. Here's an example of how it works:

If 'x' is a single character, like a letter or a digit, then the regular expression "x+" would match one or more occurrences of that character 'x'. For example:

Pattern: a+

Matches: "a", "aa", "aaa", "aaaa", ...

Pattern: 1+

Matches: "1", "11", "111", "1111", …

If 'x' is a more complex pattern, the regular expression will match one or more occurrences of that entire pattern. For example:

Pattern: \d+

Matches: "0", "123", "987654", …

Pattern: [A-Za-z]+

Matches: "abc", "XYZ", "Hello", …

**In a regular expression, ?: Matches zero or one occurrence of the preceding character or group.**

In terms of a Finite State Machine (FSM), the regular expression? (question mark) indicates zero or one occurrence of the preceding character or group.

Here's a simple FSM representation:

Start State ---(ε)---> Accepting State ---(?)---> Optional State

   ^              |                |

   |_____|_____|

in a regular expression denotes zero or one occurrence of the preceding character or group. Here are some sample examples:

Pattern: a?    Matches: "", "a"

Pattern: \d? (zero or one digit) Matches: "", "1", "9", ...

Pattern: [A-Za-z]? (zero or one letter) Matches: "", "A", "b", ...

Pattern: x? Matches: "", "x"

**{n}: Matches exactly n occurrences of the preceding character or group.**

In terms of a Finite State Machine (FSM), the regular expression {n} indicates exactly n occurrences of the preceding character or group. Representing this in an FSM is a bit more complex, as it involves having n transitions from the start state to an intermediate state and from the intermediate state to an accepting state.

A sample illustration is

Start State ---(x)---> Intermediate State ---(x)---> ... ---(x)---> Accepting State

For example, let's say you have the pattern a{3}, which means exactly three occurrences of the character 'a':

This FSM will accept sequences like "aaa" and only "aaa."

Here are some sample examples:

Pattern: a{3}    Matches: "aaa"

Pattern: \d{2} (exactly two digits) Matches: "12", "34", ...

Pattern: [A-Za-z]{4} (exactly four letters) Matches: "abcd", "WXYZ", ...

Pattern: x{0} (exactly zero occurrences of 'x') Matches: ""

**{n,}: Matches n or more occurrences.**

In terms of a Finite State Machine (FSM), the regular expression {n,} indicates matching n or more occurrences of the preceding character or group. Representing this in an FSM involves having at least n transitions from the start state to an intermediate state and then having optional transitions (0 or more) from the intermediate state back to itself.

Here's a simplified illustration:

Start State ---(x)---> Intermediate State ---(x)---> ... ---(x)---> Intermediate State ---(x)---> Accepting State

```
    ^                                                                              |

    |------------------------------------------------------------------------------ |
```

For example, let's say you have the pattern a{2,} which means matching two or more occurrences of the character 'a':

This FSM will accept sequences like "aa", "aaa", "aaaa", and so on.

ome sample examples of the regular expression {n,}, which matches n or more occurrences of the preceding character or group:

Pattern: a{2,}

Matches: "aa", "aaa", "aaaa", ...

Pattern: \d{3,} (three or more digits)

Matches: "123", "4567", ...
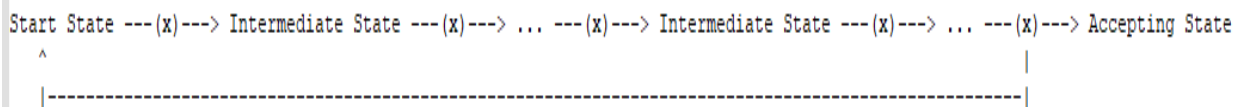
Pattern: [A-Za-z]{4,} (four or more letters)

Matches: "abcd", "WXYZ", "lmno", ...
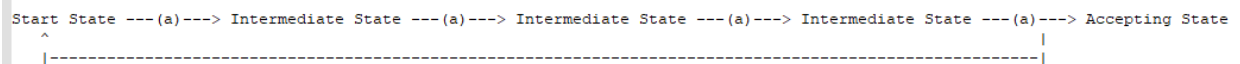
Pattern: x{0,} (zero or more occurrences of 'x')

Matches: "", "x", "xx", "xxx", ...

**{n,m}: Matches between n and m occurrences**

In terms of a Finite State Machine (FSM), the regular expression {n,m} indicates matching between n and m occurrences (inclusive) of the preceding character or group. Representing this in an FSM involves having at least n transitions from the start state to an intermediate state, optional transitions (0 or more) from the intermediate state back to itself, and at most (m-n) optional transitions from the intermediate state to an accepting state.

```
Start State ---(x)---> Intermediate State ---(x)---> ... ---(x)---> Intermediate State ---(x)---> ... ---(x)---> Accepting State
    ^                                                                                              |
    |----------------------------------------------------------------------------------------------|
```

For example, let's say you have the pattern a{2,4} which means matching between 2 and 4 occurrences of the character 'a':

```
Start State ---(a)---> Intermediate State ---(a)---> Intermediate State ---(a)---> Intermediate State ---(a)---> Accepting State
    ^                                                                                         |
    |-----------------------------------------------------------------------------------------|
```

This FSM will accept sequences like "aa", "aaa", and "aaaa."

sample examples of the regular expression {n,m}, which matches between n and m occurrences (inclusive) of the preceding character or group:

Pattern: a{2,4}    Matches: "aa", "aaa", "aaaa"

Pattern: \d{3,5} (between 3 and 5 digits) Matches: "123", "4567", "98765"

Pattern: [A-Za-z]{2,3} (between 2 and 3 letters) Matches: "ab", "XYZ", "lmn"

Pattern: x{0,2} (between 0 and 2 occurrences of 'x') Matches: "", "x", "xx"

## 2.3 Operations of Regular Expressions

Regular languages, as defined in formal language theory, can be expressed using regular expressions and recognized by finite state machines (FSMs). Let's discuss some common operations on regular languages along with examples represented in FSMs.

**1. Union (L1 ∪ L2):** The union of two regular languages results in a language containing all strings that belong to either language.

Example FSM:

Let L1 be the language of strings ending with "a."

Let L2 be the language of strings starting with "b."

The FSM for L1 ∪ L2 recognizes strings that either end with "a" or start with "b."

L U M = {s | s is in L or s is in M}

LUM = {ba, baaaba, bbbbaaa.. }

**2. Concatenation (L1 · L2):** The concatenation of two regular languages results in a language containing all possible concatenations of strings from the first language with strings from the second language.

Example FSM:

Let L1 be the language of strings containing only "0" or "1."

Let L2 be the language of strings containing only "2" or "3."

The FSM for L1 · L2 recognizes strings like "01," "32," etc.

**3. Kleene Closure (L):*** The Kleene closure of a regular language results in a language containing all possible repetitions (including zero) of strings from the original language.

Example FSM:

Let L be the language of strings with an even number of "a"s.

The FSM for L* recognizes strings with any number of "a"s, including an empty string, "aa," "aaaa," etc.

**4. Intersection (L1 ∩ L2):** The intersection of two regular languages results in a language containing only strings that are common to both languages.

Example FSM:

Let L1 be the language of strings containing "ab."

Let L2 be the language of strings containing "bc."

The FSM for L1 ∩ L2 recognizes strings that contain both "ab" and "bc."

 L ∩ M = {st | s is in L and t is in M}

5. **Complementation (¬L):** The complement of a regular language results in a language containing all strings not in the original language.

Example FSM:

Let L be the language of strings containing an even number of "0"s.

The FSM for ¬L recognizes strings with an odd number of "0"s.

These operations demonstrate how regular languages can be manipulated and combined using various set-theoretic operations. FSMs provide a graphical representation of these languages and their operations, making it easier to understand and implement algorithms for recognizing and generating strings within these languages.

Regular expressions are used for representing certain sets of strings in an algebraic fashion

1) Any terminal symbol i.e. symbols $\epsilon \sum$ $\qquad$ a,b,c, $\qquad$ ^ , $\phi$

   Including ^ and $\phi$ are regular expressions.

2) The Union of two regular expressions is $\qquad$ $R_1, R_2,$ $\qquad$ ( $R_1 + R_2$)

   also a regular expresion.

3) The Concatenation of two regular expressions $\qquad$ $R_1, R_2, \rightarrow ( R_1 + R_2 )$

   is also a regular expression.

4) The iteration (or Closure) of a regular expressions $\qquad$ $R \rightarrow R^*$ $\qquad$ $a^*$ =^,a, aa, aaa,

    is also a regular expression.

5) The regular expression over $\Sigma$ are precisely those obtained recursively by the application of the above rules once or several times.


**Q: Write the regular expression for the language accepting all combinations of a's except the null string, over the set ∑ = {a}**

The regular expression has to be built for the language

L = {a, aa, aaa, ....}

This set indicates that there is no null string. So we can denote regular expression as:

R = a+

**Q: Write the regular expression for the language accepting all the string containing any number of a's and b's.**

This will give the set as L = {ε, a, aa, b, bb, ab, ba, aba, bab, .....}, any combination of a and b.

The (a + b)* shows any combination with a and b even a null string.

**Q: Write the regular expression for the language accepting all the string which are starting with 1 and ending with 0, over ∑ = {0, 1}.**

In a regular expression, the first symbol should be 1, and the last symbol should be 0. The r.e. is as follows:

R = 1 (0+1)* 0

Q: Write the regular expression for the language starting and ending with a and having any having any combination of b's in between.

The regular expression will be:

R = a b* b

**Q: Write the regular expression for the language starting with a but not having consecutive b's.**

The regular expression has to be built for the language:

L = {a, aba, aab, aba, aaa, abab, .....}

The regular expression for the above language is:

R = {a + ab}*

**Q: Write the regular expression for the language accepting all the string in which any number of a's is followed by any number of b's is followed by any number of c's.**

As we know, any number of a's means a* any number of b's means b*, any number of c's means c*. Since as given in problem statement, b's appear after a's and c's appear after b's. So the regular expression could be:

R = a* b* c*

**Q: Write the regular expression for the language over ∑ = {0} having even length of the string.**

The regular expression has to be built for the language:

L = {ε, 00, 0000, 000000, ......}

The regular expression for the above language is:

R = (00)*

**Q: Write the regular expression for the language having a string which should have atleast one 0 and alteast one 1.**

The regular expression will be:

R = [(0 + 1)* 0 (0 + 1)* 1 (0 + 1)*] + [(0 + 1)* 1 (0 + 1)* 0 (0 + 1)*]

**Q: Describe the language denoted by following regular expression**

**r.e. = (b* (aaa)* b*)***

The language can be predicted from the regular expression by finding the meaning of it. We will first split the regular expression as:

r.e. = (any combination of b's) (aaa)* (any combination of b's)

L = {The language consists of the string in which a's appear triples, there is no restriction on the number of b's}

**Q: Write the regular expression for the language L over ∑ = {0, 1} such that all the string do not contain the substring 01.**

The Language is as follows:

L = {ε, 0, 1, 00, 11, 10, 100, .....}

The regular expression for the above language is as follows:

R = (1* 0*)

**Q: Write the regular expression for the language containing the string over {0, 1} in which there are at least two occurrences of 1's between any two occurrences of 1's between any two occurrences of 0's.**

At least two 1's between two occurrences of 0's can be denoted by (0111*0)*.

Similarly, if there is no occurrence of 0's, then any number of 1's are also allowed. Hence the r.e. for required language is:

R = (1 + (0111*0))*

Q: Write the regular expression for the language containing the string in which every 0 is immediately followed by 11.

The regular expectation will be:

R = (011 + 1)*

**Q: Write regular expressions for the following language**

  **a) the set of all strings such that the number of 0's is odd**

  **b) set of all strings that do not contain 1101**

**a) Set of all strings such that the number of 0's is odd:**

To express the language of strings with an odd number of 0's, you can use the regular expression:

Regular Expression: (1*01*0)*1*

Explanation:

(1*01*0)*: Zero or more occurrences of the pattern "1", followed by "0", followed by zero or more occurrences of "1".

1*: Zero or more occurrences of "1" at the end.

This regular expression represents strings with an odd number of 0's.

**b) Set of all strings that do not contain 1101:**

To express the language of strings that do not contain the substring "1101", you can use the regular expression:

Regular Expression: (0|1|10|110)*

Explanation:

(0|1|10|110)*: Zero or more occurrences of the patterns "0", "1", "10", or "110".

This regular expression represents strings that can be formed by concatenating any combination of "0", "1", "10", or "110", and it ensures that the forbidden substring "1101" is not present in the string.

**Q: Write regular expressions for the for the following over {0,1}***

**a) the set of all strings that begin with 110**

**b) the set of all strings that contain 1011**

**c) the set of all strings that contain exactly three 1's**

**a) Set of all strings that begin with 110:**

Regular Expression: 110.*

Explanation:

110: Specifies that the string must start with "110".

.*: Represents zero or more occurrences of any character (wildcard) after "110".

This regular expression captures strings that start with the sequence "110".

**b) Set of all strings that contain 1011:**

Regular Expression: .*1011.*

Explanation:

.*: Represents zero or more occurrences of any character (wildcard) before and after the pattern.

1011: Specifies the required pattern "1011".

This regular expression captures strings that contain the substring "1011" anywhere within them.


**c) Set of all strings that contain exactly three 1's:**

Regular Expression: 0*1(0*1){2}0*

Explanation:

0*: Represents zero or more occurrences of "0".

1: Specifies the first occurrence of "1".

(0*1){2}: Specifies exactly two occurrences of the pattern "0*1", meaning two additional occurrences of "1".

0*: Represents zero or more occurrences of "0" after the three 1's.

This regular expression captures strings that contain exactly three occurrences of the digit "1".

**Q: Write regular expressions for the following languages**

**a) the set of all strings of 0s and 1s not containing 101 as a substring**

**b) the set of all strings with an equal number of 0's and 1s. such that no prefix has two more 0'S THAN 1's , nor two more 1's than 0's**

**c) the set of strings of 0's and 1's whose number of 0's is divisible by five and whose number of 1's is even**

a) Set of all strings not containing 101 as a substring:

Regular Expression: (0|1)*((?!101)(0|1))*

Explanation:

(0|1)*: Matches any combination of 0's and 1's (including an empty string).

((?!101)(0|1))*: Utilizes negative lookahead to ensure that "101" is not present as a substring.

This regular expression captures strings that do not contain "101" as a substring.

b) Set of all strings with an equal number of 0's and 1's, without a prefix having two more 0's than 1's or two more 1's than 0's:

Regular Expression: ε|(0*10*1)*0*1*

Explanation:

ε: Represents the empty string.

(0*10*1)*: Matches pairs of 0's and 1's in any order.

0*1*: Matches any additional 0's or 1's that may follow.

This regular expression captures strings with an equal number of 0's and 1's, ensuring that no prefix has two more 0's than 1's or two more 1's than 0's.

c) Set of strings with the number of 0's divisible by five and the number of 1's being even:

Regular Expression: 0*(00000)*(1*01*01*)*

Explanation:

0*: Matches any number of 0's at the beginning.

(00000)*: Matches groups of five 0's.

(1*01*01*)*: Matches any combination of 1's, 0's, and 1's in a way that ensures the number of 1's is even.

This regular expression captures strings where the number of 0's is divisible by five, and the number of 1's is even.


**Q.Give english description of the languages of the following regular expressions**

**a) (1+ε)(00*1)*0***

**b)(0*1*)*000(0+1)***

**c)(0+10)*1***

**a) (1+ε)(00*1)*0***

This regular expression describes the language of strings that either start with a "1" or are empty, followed by zero or more occurrences of the pattern "00*1," and ending with zero or more "0"s. In simpler terms:

- The string may start with a "1" or be empty.
- It can then have zero or more occurrences of "001," where "0" represents zero or more 0's.
- Finally, the string may end with zero or more "0"s.

**b)(0*1*)*000(0+1)***

This regular expression describes the language of strings that consist of any combination of 0's and 1's, including the empty string, followed by the substring "000," and ending with zero or more occurrences of 0's or 1's. In simpler terms:

- The string can consist of any number of 0's and 1's, including the empty string.
- It must then contain the substring "000."
- The string can end with zero or more occurrences of 0's or 1's.

**c)(0+10)*1***

This regular expression describes the language of strings that consist of any combination of "0" or "10," followed by zero or more occurrences of "1." In simpler terms:

- The string can have any combination of "0" or "10."
- It must end with zero or more occurrences of "1."

## Regular Expression- Examples

Describe the following sets as regular Expressions

1) {0, 1, 2}                                    0 or 1 or 2

    R = 0 + 1+ 2

2) {^, ab}

    R = ^ ab

3) {abb, a, b, bba}                            abb or a or b or bba

    R =abb + a + b + bba

4) {^, 0, 00, 000, ……….}                    closure ∞ 0

    R = 0*

5) {1, 11, 111, 1111,    …………}

    R =1*

## Identities of Regular Expressions

1) $\emptyset$+ R=R

2) $\emptyset$R + $\emptyset$R = $\emptyset$     ^

3) $\in R = R\emptyset = \emptyset$

4) $\in^* = \in$ and $\emptyset^* = \in$

5) R + R = R

6) $R^*R^* = R^*$

7) $RR^* = R^*R$

8) $(R^*)^* = R^*$          $RR^+$

9) $\in + RR^* = \in + R^*R = R^*$

10) $(PQ)^*P = P(QP)^*$

11) $(P + Q)^* = (P^*Q^*)^* = (P^* + Q^*)^*$

12) (P+Q)R =PR +QR    and
    R(P + Q) =RP +RQ

## An Example Proof using Identities of Regular Expressions

**Prove that (1+00*1) + (1+00*1) (0+10*1)* (0+10*1) is equal to 0*1(0+10*1)***

LHS = (1+00*1) + (1+00*1) (0+10*1)* (0+10*1)

= (1+00*1) [∈+ (0+10*1)* (0+10*1)]                    ∈+R*R = R*

= (1+00*1) (0+10*1)*

=(∈. (1+00*1) (0+10*1)*                    ∈.R =R

=(∈+ 00*)1 (0+10*1)*

= O*1 (0+10*1)*=RHS

Design Regular Expression for the following languages over {a,b}

1) Language accepting strings of length exactly 2

2) Language accepting strings of length atleast 2

3) Language accepting strings of length atmost 2

Soln

1) L = {aa, ab, ba, bb}                    2) $L_1$={aa, ab, ba,bb, aaa,…………..}

R= aa+ab+ba+bb                    R= (a+b)(a+b) (a+b)*

= a (a+b)+b(a+b)          = (a+b) (a+b)

3) $L_1$ = {∈, a, b, aa, ab, ba, bb}

  R=∈+a+b+aa+ab+ba+bb

= (∈+a+b) (∈+a+b)

## 2.4 Finite Automata and Regular Expression

### 2.4.1 From DFA to Regular Expression

Converting a Deterministic Finite Automaton (DFA) to regular expressions is a process that allows you to represent the same language in a more compact and human-readable form. This conversion is useful for various reasons, including simplifying the understanding of the language, facilitating further analysis and manipulation, and providing a more intuitive representation for certain applications.

### 2.4.2 Why Convert DFA to Regular Expressions:

- **Compact Representation**: Regular expressions are often more concise and easier to understand than DFAs, especially for complex languages. They provide a high-level view of the language without delving into the detailed transitions and states of the automaton.

- **Ease of Communication**: Regular expressions are widely understood and used in computer science and programming. Converting a DFA to regular expressions makes it easier to communicate the language's rules and patterns with others.
- **Simplification of Language Description**: Regular expressions allow you to describe patterns and rules more naturally and intuitively. This can be especially beneficial when dealing with complex languages that may have intricate DFA representations.
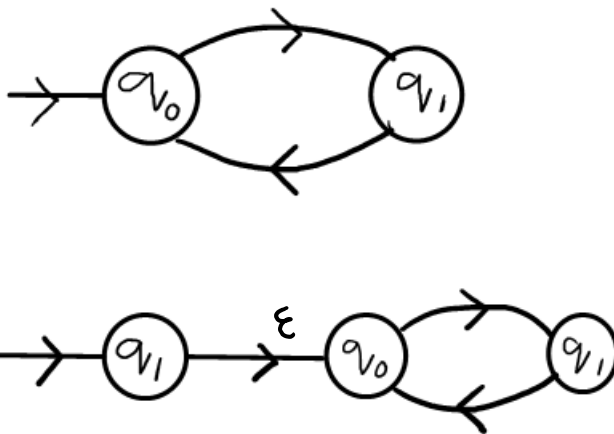
### 2.4.3 Different methods of DFA to Regular expression

There are several methods to convert a Deterministic Finite Automaton (DFA) to a regular expression. Two common methods are the state elimination method and the use of the Arden's Theorem. Let's briefly discuss both methods:

1. **State Elimination Method**: This method involves eliminating states one by one until only the initial and final states remain. At each step, you replace two states with a regular expression that represents the language accepted by the original DFA between those states.

- **Identify Dead-End States**: Identify and eliminate any non-final states that have transitions only to final states.
- **State Elimination**: For each pair of states ($q_i$, $q_j$) in the remaining set, find a regular expression that represents the language accepted from $q_i$ to $q_j$, excluding eliminated states.
- **Repeat**: Repeat the process until only the initial and final states remain, and you have a regular expression representing the entire language.
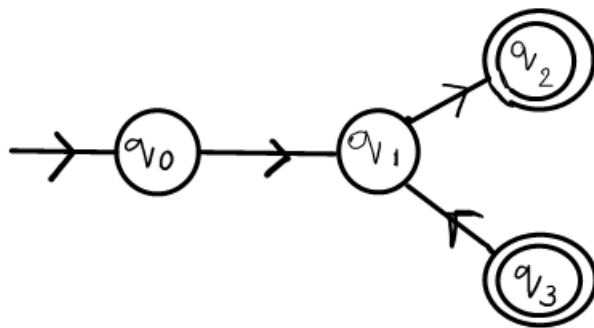
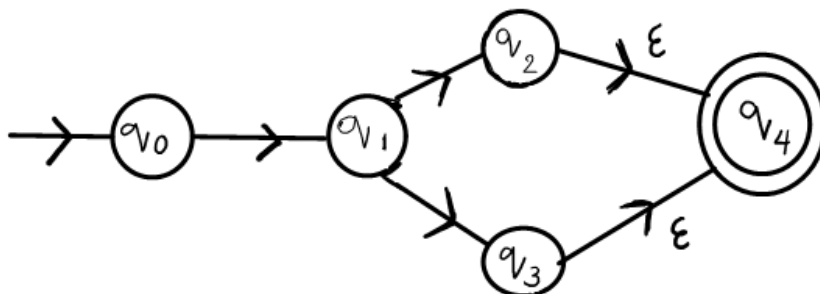**STEP 1: if there exists any incoming edge to initial state add a new initial state**



**Step 2: if there exist an outgoing edge to final state create a new final state**

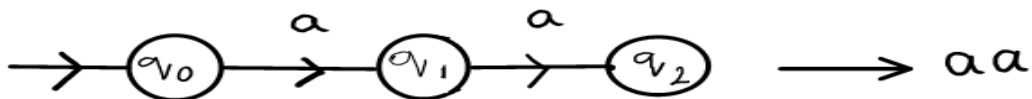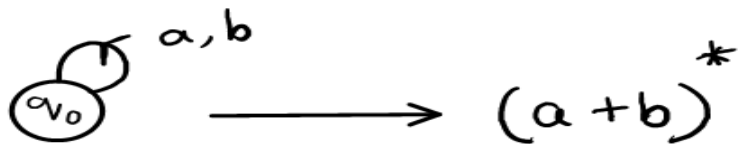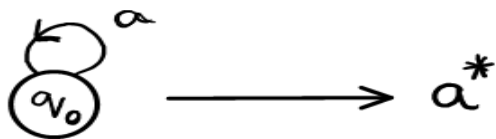**Step 3: if there exist multiple final states convert into non-final state and create a single final state**
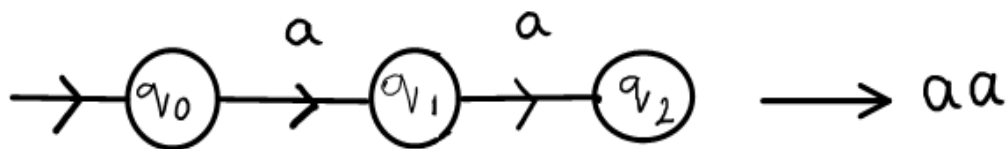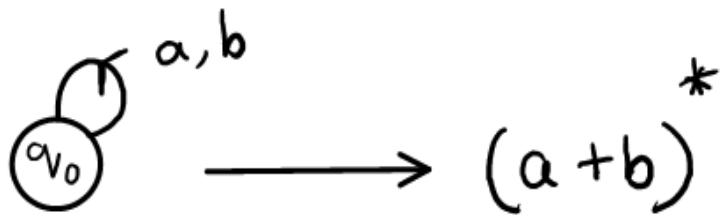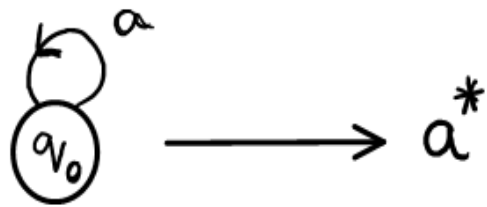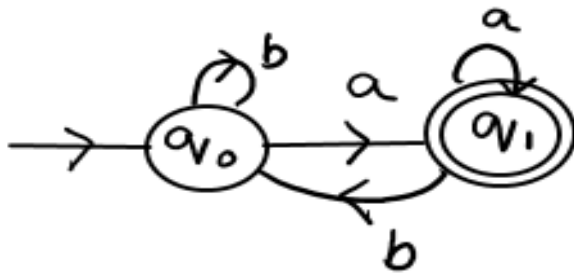


**Step 4: Eliminate all intermediate states one by one in any order**



**Step 4: Eliminate all intermediate states one by one in any order**

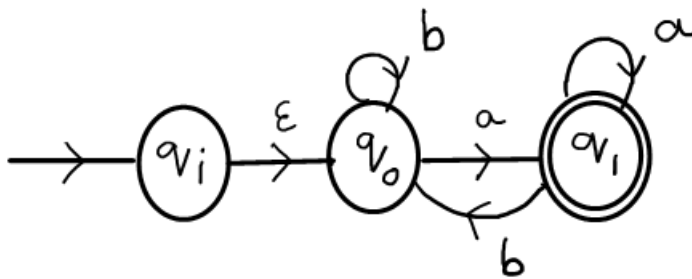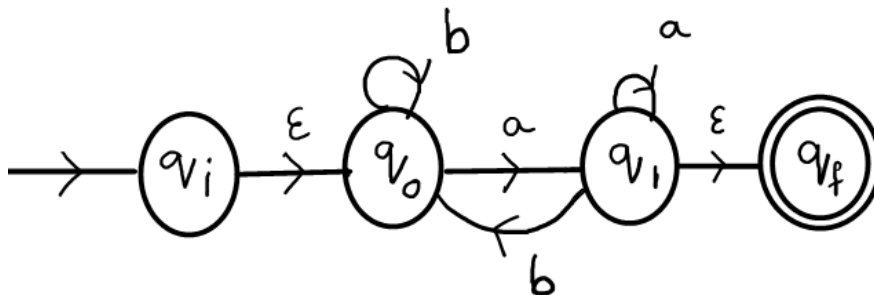**Background: converting finite automata to regular expression**

$$a^*$$

$$(a+b)^*$$

$$aa$$

$$a^*$$

$$(a+b)^*$$

$$aa$$

1. **Convert Following DFA to Regular expression**
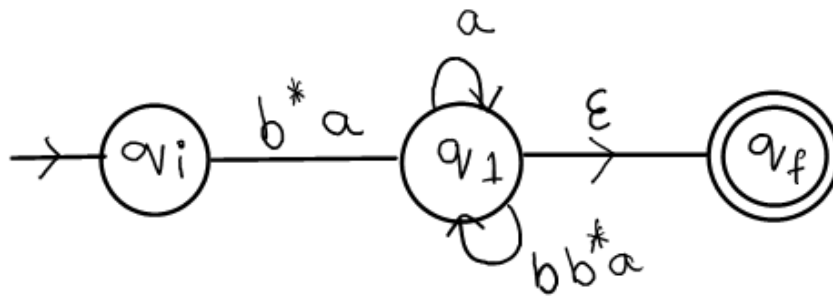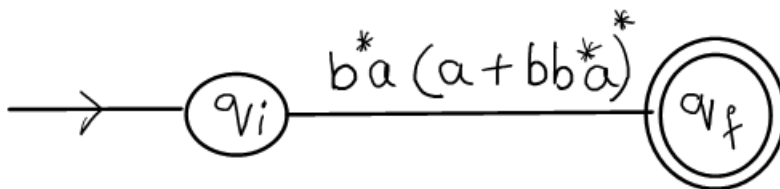
Add Initial state



Add final state



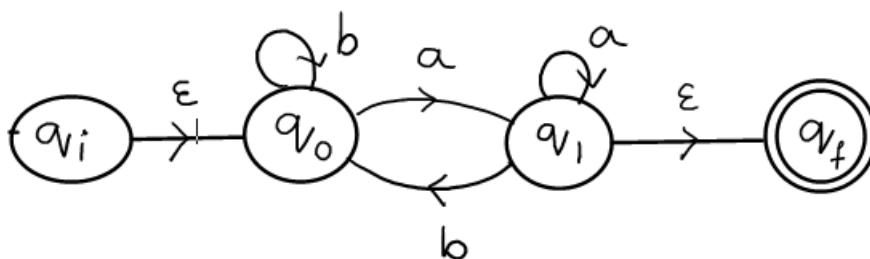## Approach 1: (from initial state)
### Eliminate state $q_0$

**Eliminate state q₁**



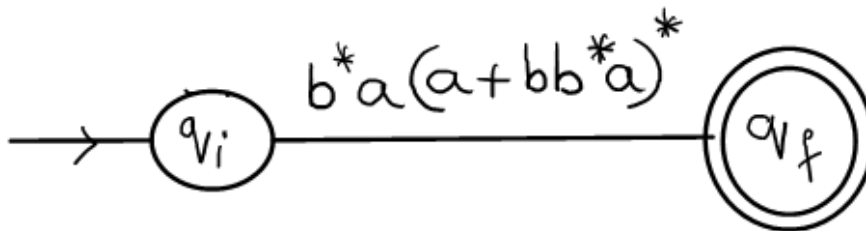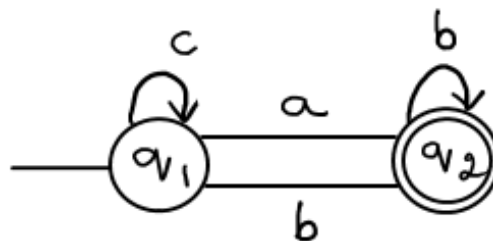# Approach 2: (from final state)
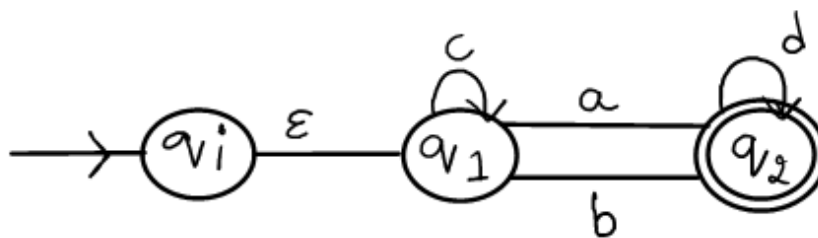
# Eliminate q1

**Eliminate state $q_0$**



2.  **Construct Regular Expression using following DFA**



**Step 1** − Initial state q1 has incoming edge. So, create a new initial state qi.



**Step 2** − Final state q2 has outgoing edge. So, create a new final state.

**Step 3** − Start eliminating intermediate states one after another.
**Eliminate q1**



Now eliminate q2.

After eliminating q2 direct path from state qi to qf having cost.

c*a(d+bc*a) * ε=**c*a(d+bc*a)***

2. **Arden's Theorem Method**:

Arden's Theorem provides a way to find the solution to a system of linear equations involving regular expressions. This method can be applied to express the regular language accepted by a DFA.

Steps:

- **Set up Equations**: For each state $q_i$, set up an equation expressing the regular expression representing the language accepted by $q_i$ in terms of the regular expressions of its neighbour's.
- **Solve Equations**: Use Arden's Theorem to solve the system of equations and find the regular expression for each state.
- **Combine Expressions**: Combine the regular expressions obtained for each state to get the overall regular expression for the language accepted by the DFA.

**1. Design or Find Regular expression for the following NFA (Using Arden's theorem)**

$$q_3 = q_2 a \rightarrow ①$$
$$q_2 = q_1 a + q_2 b + q_3 b \rightarrow ②$$
$$q_1 = E + q_1 a + q_2 b \rightarrow ③$$

$$① \rightarrow q_3 = q_2 a$$
$$= (q_1 a + q_2 b + q_3 b) a$$
$$= q_1 aa + q_2 ba + q_3 ba \rightarrow ④$$

2)

$$q_2 = q_1 a + q_2 b + q_3 b \quad \text{Putting value of } q_3 \text{ from} - ①$$
$$= q_1 a + q_2 b + (q_2 a) b$$
$$= q_1 a + q_2 b + q_2 ab$$
$$\underset{R}{q_2} = \underset{Q}{q_1 a} + \underset{R}{q_2} \underset{P}{(b + ab)}$$

$$R = Q + RP \text{ Arden's Theorem}$$
$$R = QP$$

$$q_2 = (q_1 a)(b + ab)^* \rightarrow ⑤$$

③ → $q_1 = \epsilon + q_1 a + q_2 b$

Putting value of $q_2$ from ⑤

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

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

$R = Q + RP$

$R = QP^*$

$\epsilon . R = R$

Under the equation: $\underset{R}{q_1} = \underset{Q}{\epsilon} + \underset{R}{q_1} \underset{P}{(a + a(b+ab)^*) b}$

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

$q_1 = (a + a(b+ab)^* b)^* \longrightarrow ⑥$

Final state $q_3$

$q_3 = q_2 a$

   $= q_1 a(b+ab)*a$      putting value of $q_2$ from 5

 $q_3 = (a+a(b+ab)*b)*$ $a(b+ab)*a$   putting value of $q_1$ from 6

   = Required Regular expression for the given NFA

**2. Design the Regular expression for the following DFA**



$q_1 = \epsilon + q_2 b + q_3 a \longrightarrow ①$

$q_2 = q_1 a \longrightarrow ②$

$q_3 = q_1 b \longrightarrow ③$

$q_4 = q_2 a + q_3 b + q_4 a + q_4 b \longrightarrow ④$

$①↪ q_{v_1} = E + q_2 b + q_3 a$

Putting values of $q_2$ and $q_3$ from $②$ and $③$

$q_{v_1} = E + q_1 ab + q_1 ba$

$q_1 = E + q_{v_1} (ab + ba)$

$\underbrace{q_1}_{R} = \underbrace{E}_{Q} + \underbrace{q_{v_1}}_{R} \underbrace{(ab+ba)}_{P}$

$R = Q + RP$

$R = QP^*$  Arden's Theorem

$q_1 = E (ab + ba)^*$

$E R = R$

$q_1 = (ab+ba)^*$

## 3    Design Regular expression for following when there are multiple Final states

Find the Regular Expression for the following DFA



$q_1 = E + q_{v_1} 0 \longrightarrow ①$

$q_2 = q_{v_1} 1 + q_{v_2} 1 \longrightarrow ⑪$

$q_3 = q_2 0 + q_3 0 + q_3 1 \longrightarrow ⑪⑪$

Final state $q_{v_1}$

$①↪$

$\underbrace{q_1}_{R} = \underbrace{E}_{Q} + \underbrace{q_{v_1}}_{R} \underbrace{0}_{P}$

$h = Q + RP$

$R = QP^*$  Arden's theorem

$Eh = h$

$q_1 = E . 0^*$

$q_1 = 0^* \longrightarrow ④$

Final state $q_2$

$$q_2 = q_1 1 + q_2 1$$

$$q_2 = 0^* 1 + q_2 1 \qquad \text{Putting value of } q_1 \text{ from ④}$$

$\underset{h}{\downarrow} \quad \underset{Q}{\downarrow} \quad \underset{R}{\downarrow} \underset{P}{\downarrow}$

$$q_2 = 0^* 1 (1)^* \qquad \begin{array}{l} R = Q + RP \\ R = QP^* \end{array}$$

$h =$ union of both final states

$$= 0^* + 0^* 11^*$$

$$= 0^* (\epsilon + 11^*) \qquad \epsilon + RR^* = R^*$$

$$= 0^* 1^*$$

### 2.3.4  Regular Expression to Finite Automata

The conversion of a regular expression to a finite automaton is a fundamental concept in formal language theory. The importance of this conversion lies in the fact that it provides a systematic way to represent and understand the relationship between regular languages (described by regular expressions) and the corresponding machines (finite automata)

The importance of converting regular expressions to finite automata:

1. **Equivalence**: The regular expression and finite automaton representations are equivalent in terms of expressive power; that is, they describe the same class of languages. This equivalence demonstrates that regular languages can be recognized by both regular expressions and finite automata.

2. **Understanding Language Recognition**: The conversion process helps in understanding how certain patterns and structures in regular expressions can be translated into the operational behavior of a finite automaton. It provides insights into the relationship between the syntactic and semantic aspects of regular languages.

3. **Algorithmic Recognition**: Finite automata are often used as the underlying mechanism for recognizing regular languages in computer science and automata theory. Algorithms for pattern matching, lexical analysis (e.g., in compilers), and string searching often involve finite automata. By converting regular expressions to finite automata, we can implement efficient algorithms for recognizing patterns in strings.

4. **Tool for Language Specification**: Regular expressions are widely used for specifying patterns in various contexts, such as text processing, search algorithms, and lexical analysis. Finite automata provide a

concrete and implementable representation of these patterns, allowing developers to build systems that recognize and process languages described by regular expressions.

5. **Optimization and Implementation**: The conversion process can reveal opportunities for optimizing the recognition process. Understanding the finite automaton corresponding to a regular expression can lead to the elimination of redundant states or transitions, resulting in more efficient implementations.

6. **Educational Purposes**: Learning about the conversion from regular expressions to finite automata is a common topic in courses on formal languages and automata theory. It serves as a foundational concept for understanding more complex language classes and machine models.

7. **Automation in Software Development:** Tools and libraries that automate the conversion from regular expressions to finite automata are valuable for software developers. These tools enable efficient implementation of language recognition mechanisms and are used in various applications, including text processing and validation.

In summary, the conversion of regular expressions to finite automata is a crucial step in the study of formal languages and provides a bridge between theoretical language descriptions and practical language recognition implementations in computer science and software development.

**Important points:**



1. **Convert following Regular expression to their Finite Automata**
   a) ba*b
   b) (a+b)c
   c) a(bc)*

a) b a* b :

1) b a*b                          bb, bab, baab



**b) (a+b).c**

(a+b) c



**c) a(bc)*8**

3) a (bc)*                a, abc, abcbc, abcbcbc



**6) (a | b)* (abb|a⁺b)**

$$(a|b)^* (abb|a^+b)$$



$$a^+ = \{a, aa, aaa\}$$

$$a^* = \{\epsilon, a, aa\}$$

**7) 10+(0+11)0*1**

$$10 + (0 + 11)\, 0^* 1$$



$$(0+11)\, 0^* 1$$





**8). 1 (1* 01* 01*)*.**

Step : 1

$$\left(1^*01^*01^*\right)$$





**9) 0*1+10**

## 2.3.5  Equivalence of two finite automata

Equivalence of two finite automata means that the two automata recognize the same language. In other words, they accept the same set of strings. Checking the equivalence of finite automata is a crucial task in various areas of computer science and formal language theory. Here are some reasons why one might need to check the equivalence of finite automata:

- **Language Equivalence**: The primary reason for checking equivalence is to determine whether two finite automata recognize the same language. If two automata are equivalent, it means they accept the same set of strings. This is essential for ensuring that different implementations or representations of a language are consistent.
- **Program Verification**: In the field of software engineering, finite automata and regular languages are often used to model certain behaviors or specifications of programs. Checking the equivalence of automata can be part of program verification processes to ensure that different implementations of a system behave identically in terms of the recognized language.

- **Compiler Construction**: In compiler construction, finite automata are commonly used to model lexical analyzers (scanners). It is crucial to verify that the lexical analyzers generated by different compiler tools or components are equivalent, ensuring consistent token recognition.
- **Optimization**: Checking equivalence is relevant when optimizing finite automata. For example, if you have multiple representations of the same language, you might want to check if they are equivalent and choose the most efficient one. Minimization algorithms are often applied to reduce the size of an automaton while preserving its language.
- **Automata Learning**: In the context of machine learning and automata learning, equivalence testing is used to verify if a learned automaton is equivalent to a target automaton. This is essential for validating the correctness of the learned model.
- **Formal Language Theory**: Equivalence testing is a fundamental concept in formal language theory. It helps researchers and practitioners understand the relationships between different automata models and their expressive power.
- **Security Protocols**: Finite automata are used in security protocols to model the behavior of systems and attackers. Checking the equivalence of models is crucial in security analysis to ensure that the system and the attacker's model are consistent.
- **Database Query Optimization**: In the context of query optimization, finite automata are sometimes used to represent query languages. Checking equivalence can help optimize queries by choosing equivalent but more efficient representations.

In summary, checking the equivalence of finite automata is a fundamental and versatile concept with applications in various domains, including software engineering, formal methods, compiler construction, and security analysis. It ensures consistency and correctness in language recognition, which is essential in many computational and theoretical contexts.

## Step to identify equivalence

1) For any pair of states $\{q_i, q_j\}$ the transition for input a∈∑ is defined by $\{q_a, q_b\}$ where $\partial\{q_i, a\} = q_a$ $and$ $\partial\{q_j, a\} = q_b$

The two automata are not equivalent if for a pair $\{q_a, q_b\}$ one is INTERMEDIATE state and the other is FINAL state.

2) If Initial State is Final state of one automation, then in second automation also Initial state must be Final State for them to be equivalent.

1. **Check the two following two finite automata are equivalent or not**

Example:



A                                    B

Following are the evaluation

| Staks | c | d |
|-------|---|---|
| $\{q_1, q_4\}$ | $\{q_1, q_4\}$ <br> F's  F's | $\{q_2, q_5\}$ <br> I's  I's |
| $\{q_2, q_5\}$ | $\{q_3, q_7\}$ <br> I's  I's | $\{q_1, q_4\}$ <br> F's  F's |
| $\{q_3, q_6\}$ | $\{q_2, q_7\}$ <br> IS  IS | $\{q_3, q_6\}$ <br> IS  IS |
| $\{q_2, q_7\}$ | $\{q_3, q_6\}$ <br> IS  IS | $\{q_1, q_4\}$ <br> Fs  Fs |

A and B are equivalent

**2. Check the two following two finite automata are equivalent or not**

A                    B

Solution:

| staks | c | d |
|---|---|---|
| $\{q_1, q_4\}$ | $\{q_1, q_4\}$<br>F's  F's | $\{q_2, q_5\}$<br>I's  I's |
| $\{q_2, q_5\}$ | $\{q_3, q_7\}$<br>I's  I's | $\{q_1, q_6\}$<br>F's  I's |

A and B are not equivalent

## 2.4   Proving languages not to be regular

Regular languages have several properties that distinguish them from other classes of languages. Here are some key properties of regular languages:

- **Closure under Union (OR):** If L1 and L2 are regular languages, then their union L1∪L2 is also a regular language.
- **Closure under Intersection (AND):** If L1 and L2 are regular languages, then their intersection L1∩L2 is also a regular language.
- **Closure under Complementation (NOT):** If L is a regular language, then its complement ‾L (the set of all strings not in L) is also a regular language.
- **Closure under Concatenation:** If L1 and L2 are regular languages, then their concatenation ·L2 (the set of all strings formed by concatenating a string from L1 with a string from L2) is also a regular language.
- **Closure under Kleene Star:** If L is a regular language, then its Kleene closure L∗ (the set of all strings formed by concatenating zero or more strings from L) is also a regular language.
- **Closure under String Reversal**: If L is a regular language, then its reversal LR (the set of all strings whose reversal is in L) is also a regular language.

- **Existence of a Finite Automaton:** A language is regular if and only if there exists a finite automaton that recognizes it. This is the defining characteristic of regular languages.
- **Pumping Lemma**: Regular languages satisfy the Pumping Lemma, which provides a property used in proving that certain languages are not regular.
- **Equivalence to Regular Expressions**: A language is regular if and only if there exists a regular expression that generates it. Regular expressions and finite automata are two equivalent ways to describe regular languages.
- **Decidability**: Properties of regular languages are decidable. This means that, for any statement about regular languages, there exists an algorithm that can determine whether the statement is true or false.
- These properties make regular languages a well-defined and well-behaved class of languages, and they serve as the foundation for the study of formal languages and automata theory

### 2.4.4    Pumping Lemma:

The Pumping Lemma typically states that for any regular language, there exists a constant (the "pumping length") such that any strings in the language of length at least p can be divided into three parts, xyz, satisfying the following conditions:

1. For each $i \geq 0$, the string $xy^iz$ is in the language.
2. $|y|>0$ (i.e., y is non-empty).
3. $|xy| \leq p$ (the length of xy is at most p).

The idea is that, because the language is regular, repeating or "pumping" the middle part y should still produce strings in the language.

To use the Pumping(substring of a string repeated) Lemma(substring) to show that a language is not regular, one assumes, for the sake of contradiction, that the language is regular and then chooses a specific string in the language that violates the conditions of the Pumping Lemma. This contradiction implies that the assumption that the language is regular must be false, and therefore the language is not regular.

It's important to note that the Pumping Lemma is specific to regular languages. There are other lemmas and techniques used for proving properties of context-free languages and other language classes. The Pumping Lemma is a powerful tool, but its application requires careful reasoning and understanding of the properties of regular languages. **In layman's understanding if the substring of a string is repeated many times and if the resultant of string is available in the language then we can say that it as REGULAR.**

**STEPS :**

- Consider a language as regular
- Assume constant 'c' and select the 'w' string such that $|w| >= C$
- Divide w as XYZ
- $|Y|>0$
- $|XY| <=C$
- For $i>=0$ every string $xy^iz$ belong to L

1. **Prove that {a$^n$b$^n$ | n>=0} is not Regular**

Solution:

Let us consider the language L={ ε,ab, aabb,aaabbb,aaaabbbb,aaaaabbbbb,….}

Assume C = 6

*For instance , w =aaabbb |w| >=C*



x = aa y= ab z=bb

|y|=2 |xy| =aaab =4 < C

For instance I=0 , xy$^2$z =aaababbb – not available in L (NO. of a's are not followed b in sequence)

*For instance , w =aaabbb |w| >=C*

In this x =a y=aa z=bbb C=6

|Y| = 2 >0, |XY| = 3 < C ,

When i=2

XY$^0$Z = aaa not in our Language L

XY$^2$Z = aaaaabbb not in our Language L

{a$^n$b$^n$} is not a regular

2. Using Pumping Lemma prove that the language A={YY|Y ∈ {0,1}* } is not Regular

Solution :

Y=ε: YY=εε=ε

Y="0" YY="00"

Y="10" YY="1010"


Consider L= {00,0000,01010,0000100001.}

C = 7

For instance , w= 0000100001|w| >=C

In this i=0, x=00   z=00001 = {0000001} is not Regular

When i=1, x=0 y=0001 z=00001    |XY$^1$Z| = 00000100001 -not repeating properly so it is not regular

## 2.5   Lexical Analysis Phase and compiler design

### 2.5.4  Role of Lexical Analyzer

The lexical analyzer, also known as the lexer or scanner, is a crucial component of the compiler responsible for processing the input source code and generating a stream of tokens. Its role is primarily focused on the lexical analysis phase, which is the initial step in the compilation process.

✓ The main task of the lexical analyzer is to
  o **read the input characters** of the source program,
  o **group** them into **lexemes**, and
  o produce as **output a sequence of tokens** for the source program.
  o **stripping out comments and whitespace** (blank, newline, tab etc), that are used to separate tokens in the input.

✓ Parser invokes the lexical analyzer by ***getNextToken*** command
✓ Lexical analyzer reads the characters from input until it finds the next lexeme and produce token



Here are the key roles and functions of the lexical analyzer:

▪ **Tokenization:** The primary function of the lexical analyzer is to break the source code into a sequence of tokens. Tokens are the smallest meaningful units of the programming language, such as keywords, identifiers, literals, and operators. Tokenization simplifies the subsequent phases of the compiler by representing the source code in a more structured and manageable form.

▪ **Ignoring Whitespaces and Comments:** The lexical analyzer filters out irrelevant elements such as whitespaces and comments from the source code. These elements are important for human readability but do not contribute to the execution of the program. Removing them reduces the complexity of the subsequent phases.

▪ **Error Detection**: The lexical analyzer plays a role in detecting lexical errors in the source code. Lexical errors include issues such as invalid characters or tokens that do not conform to the language's syntax. Detecting errors early in the compilation process helps provide meaningful error messages to the programmer for debugging.

- **Symbol Table Generation:** As the lexical analyzer processes the source code, it may build a symbol table. The symbol table is a data structure that keeps track of identifiers (variable names, function names, etc.) encountered in the source code along with their attributes, such as data type or memory location. This information is crucial for later phases of the compiler, particularly in semantic analysis.

- **Pattern Matching:** The lexical analyzer uses regular expressions and finite automata to perform pattern matching on the source code. Each token type is associated with a specific pattern, and the lexical analyzer identifies and categorizes tokens based on these patterns.

- **Generating Output for Parser:** The output of the lexical analyzer is a stream of tokens, which is then passed on to the next phase of the compiler, the parser. The parser uses this token stream to build a parse tree or an abstract syntax tree, representing the syntactic structure of the program.

- **Efficiency and Optimization:** Lexical analyzers are designed to be efficient, as they are the first phase in the compilation process. Techniques such as finite automata and regular expressions are employed to optimize the scanning process and reduce the time complexity of token recognition.

**Lexical Analyzer**

- Scans the pure HLL code line by line.

- Takes Lexemes as i/p and produces Tokens.



Sample source code :

### *2.5.5 Tokens, Patterns and Lexemes*

**Lexeme**: A lexeme is a basic unit of lexical analysis and represents a sequence of characters in the source program that corresponds to a single token. Tokens are the smallest units in a programming language, and lexemes are the actual instances of these tokens in the source code.

consider the following line of code in a simple programming language:

x = 10 + y

In this line, there are several lexemes, each representing a distinct token:

x is a lexeme representing an identifier.

= is a lexeme representing the assignment operator.

10 is a lexeme representing a numeric constant.

+ is a lexeme representing the addition operator.

y is a lexeme representing another identifier.

So, in short, lexemes are the fundamental building blocks identified during lexical analysis, and they help break down the source code into meaningful units for further processing by the compiler or interpreter

**Pattern**: In lexical analysis, a "pattern" usually refers to a rule or a template that describes the structure of a token. Tokens are the basic building blocks of a programming language, and patterns help identify and classify these tokens in the source code.

Let's break it down in simpler terms:

**Token**: Think of a token as a meaningful unit in a programming language. For example, in the code int x = 5;, the tokens include the keywords int, x, =, and 5.

**Pattern**: A pattern is a set of rules that defines what a token looks like. For instance, a simple pattern for an integer constant could be "one or more digits." So, the pattern for the token 5 is "digit(s)."

Here's an example:

Consider the following line of code:

total = count * 10;

In this line, there are several tokens, each following a specific pattern:

**Identifier** (total): This follows the pattern of starting with a letter and followed by letters or digits.

**Assignment Operator** (=): This is a simple pattern of the equal sign.

**Identifier** (count): Similar to the first identifier, it follows the pattern of starting with a letter and followed by letters or digits.

**Multiplication Operator** (*): This is a single-character pattern representing the multiplication operation.

**Integer Constant** (10): This follows the pattern of one or more digits.

**Semicolon** (;): This is a single-character pattern representing the end of a statement.

So, in lexical analysis, patterns help define the rules for recognizing and extracting different types of tokens from the source code. They are crucial for breaking down the code into meaningful elements for further processing by the compiler or interpreter

**Token**: It is a pair consisting of a token name and an optional attribute value. a "token" in lexical analysis is like a building block of a programming language. It's a meaningful chunk or piece of the code that represents a fundamental unit.

- The token name as an abstract symbol represents the kind of lexical unit/lexeme(keyword/identifier, operator symbol etc)
- Processed by parser

In the statement C, printf("Total =%d\n", score)

Printf ,score -> lexmes matching the pattern for token id, "Total="%d\n" is lexme matching literal





DFA

Lexical Analyzer-Tokenization:

Tokens:

1. Keyword: int, return

2. Identifier: main, x ,a ,b ,c , print f

3. Punctuator: ( ,) ,{,,,,;}

4. Operator: =,+,*

5. Constant: 2,3,5,0

6. Literal: "The value of x is %d"

```
int main()
{
int x , a=2,b=3,c=5;
x = a + b*c;
 Print f ("The value of x is %d", x);
return 0;
}
```

Count: 39

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

One token for each keyword. The pattern for a keyword is the same as the keyword itself.

Tokens for the operators, either individually or in classes such as the token comparison

One token representing all identifiers.

One or more tokens representing constants, such as numbers and literal strings.

Tokens for each punctuation symbol, such as left and right parentheses, comma, and semicolon.

**SAMPLE CODE IN PYTHON : Construction of Symbol table**

```python
import re
# Token types
TOKEN_INT = 'INT'
TOKEN_FLOAT = 'FLOAT'
TOKEN_IDENTIFIER = 'IDENTIFIER'
TOKEN_OPERATOR = 'OPERATOR'
TOKEN_KEYWORD = 'KEYWORD'
TOKEN_SEPARATOR = 'SEPARATOR'
# Regular expressions for token recognition
regex_patterns = [
    (r'\bint\b|\bfloat\b|\bchar\b', TOKEN_KEYWORD),
    (r'\bif\b|\belse\b|\bwhile\b', TOKEN_KEYWORD),
    (r'\breturn\b', TOKEN_KEYWORD),
    (r'\b\d+(\.\d+)?\b', TOKEN_FLOAT),
    (r'\b[a-zA-Z_]\w*\b', TOKEN_IDENTIFIER),
    (r'\+|\-|\*|\/|>|=|\(', TOKEN_OPERATOR),
    (r'\)|\{|\}|\;|\,', TOKEN_SEPARATOR),
]


# Symbol table
symbol_table = {}

def getNextToken(input_string):
    if not input_string:
        return None, input_string    # Return None for an empty string
    for pattern, token_type in regex_patterns:
        regex = re.compile(pattern)
        match = regex.match(input_string)
        if match:
            value = match.group(0)
            return {'lexeme': value, 'type': token_type}, input_string[len(value):].lstrip()
    raise Exception(f'Unexpected character: {input_string[0]}')

def parse(source_code):
    global symbol_table
    tokens = []
    while source_code:
        token, source_code = getNextToken(source_code)
        tokens.append(token)
        if token['type'] == TOKEN_KEYWORD and token['lexeme'] in ['int', 'float', 'char']:
            # Handle variable declarations and update symbol table
            var_type = token['lexeme']
            var_name, source_code = getNextToken(source_code)
            if var_name['type'] == TOKEN_IDENTIFIER:
                symbol_table[var_name['lexeme']] = var_type
            else:
                raise Exception(f'Expected identifier after {var_type}, got {var_name["lexeme"]}')

        elif token['type'] == TOKEN_KEYWORD and token['lexeme'] == 'return':
            # Handle 'return' keyword (additional parsing logic can be added here)
            return_expression, source_code = getNextToken(source_code)
            if return_expression['type'] == TOKEN_IDENTIFIER:
                if return_expression['lexeme'] not in symbol_table:
                    raise Exception(f'Undeclared identifier in return statement: {return_expression["lexeme"]}')
            # Additional parsing logic for the return statement can be added here
        elif token['type'] == TOKEN_IDENTIFIER:
            # Handle variable usage (e.g., assignments) and check symbol table
            if token['lexeme'] not in symbol_table:
                raise Exception(f'Undeclared identifier: {token["lexeme"]}')
            # Additional parsing logic for assignments, etc., can be added here
    return tokens
```

```
try:
    tokens = parse(source_code)
    print("{:<15} {:<15} {:<15}".format("Lexeme", "Token Type", "Symbol Table"))
    print("="*60)
    for token in tokens:
        lexeme = token["lexeme"]
        token_type = token["type"]
        symbol_table_entry = symbol_table.get(lexeme, '-')
        print("{:<15} {:<15} {:<15}".format(lexeme, token_type, symbol_table_entry))

    print("\nSymbol Table:")
    print("{:<15} {:<15}".format("Identifier", "Type"))
    print("="*30)
    for identifier, var_type in symbol_table.items():
        print("{:<15} {:<15}".format(identifier, var_type))
except Exception as e:
    print(f'Error: {str(e)}')
```

## OUTPUT:

| Lexeme | Token Type | Symbol Table |
|--------|-----------|--------------|
| int | KEYWORD | - |
| ( | OPERATOR | - |
| ) | SEPARATOR | - |
| { | SEPARATOR | - |
| float | KEYWORD | - |
| = | OPERATOR | - |
| 3.14 | FLOAT | - |
| ; | SEPARATOR | - |
| int | KEYWORD | - |
| = | OPERATOR | - |
| 42 | FLOAT | - |
| ; | SEPARATOR | - |
| if | KEYWORD | - |
| ( | OPERATOR | - |
| x | IDENTIFIER | float |
| > | OPERATOR | - |
| 2.0 | FLOAT | - |
| ) | SEPARATOR | - |
| { | SEPARATOR | - |
| y | IDENTIFIER | int |
| = | OPERATOR | - |
| y | IDENTIFIER | int |
| * | OPERATOR | - |
| 2 | FLOAT | - |
| ; | SEPARATOR | - |
| } | SEPARATOR | - |
| else | KEYWORD | - |
| { | SEPARATOR | - |
| y | IDENTIFIER | int |
| = | OPERATOR | - |
| y | IDENTIFIER | int |
| + | OPERATOR | - |
| 1 | FLOAT | - |
| ; | SEPARATOR | - |
| } | SEPARATOR | - |
| return | KEYWORD | - |
| ; | SEPARATOR | - |
| } | SEPARATOR | - |

Symbol Table:
Identifier      Type

```
================================
main            int
x               float
y               int
```

---------------------------------------------------------------------------------------------------------------------

## 2.5.6  Lexical errors

Lexical errors, also known as lexical or scanning errors, occur during the lexical analysis phase of the compilation process when the compiler is analyzing the source code to identify and tokenize its basic components, such as keywords, identifiers, literals, and symbols. These errors are related to the structure of the code and how it conforms to the language's lexical rules. Here are some common types of lexical errors:

**Classification of Errors:**

**Identifiers that are way too long**

- Exceeding length of numeric constants.

Int i = 4567891;

Size: 2 Bytes

-32,768 to 32,767

- **Numeric constants which are ill-formed.**

Int i = 4567$91;

## 2.5.7  Attributes for token

In lexical analysis, tokens are the smallest units of meaning in a programming language. Each token represents a specific type of symbol or sequence of characters in the source code. Tokens are identified based on patterns defined by the lexical rules of the language. The attributes of tokens refer to additional information associated with each token type. Here are some common attributes of tokens along with examples:

- Token Type:

Example: int, if, while, +, =

- Lexeme: The actual sequence of characters in the source code that matches a particular token.

Example: For the token type int, the lexeme might be int. For the token type identifier, the lexeme might be variable Name.

- Value: For tokens representing constants (e.g., numeric literals, string literals), the actual value of the constant.

Example: For the token type integer literal, the value might be the actual numeric value like 42.

- Line Number: The line number in the source code where the token was found.

Example: If a token is found on line 5 of the source code, its line number attribute would be 5.

- Position in Line: The position of the token within the line of source code.

Example: If a token is the third element on line 5, its position in line attribute would be 3.

- Token Class: A categorization of tokens into broader classes (e.g., keywords, identifiers, operators).

Example: Token class can be Keyword for tokens like if and while, or Operator for tokens like + and =.

- Scope: The scope in which an identifier is defined (used in languages with scoping rules).

Example: In a programming language with block scope, the scope might be the block in which the identifier is declared.

- Data Type: For languages with static typing, the data type is associated with a variable or constant.

Example: For an identifier representing an integer, the data type attribute might be int.

These attributes help in building the symbol table and capturing additional information needed for subsequent phases of compilation. The specifics of token attributes can vary based on the language and the design choices made by the compiler or interpreter for that language.

## 2.5.8   Input Buffering

Input buffering is a process used in computer systems and programming to efficiently handle the input data, especially in the context of lexical analysis and parsing in compilers. The primary goal of input buffering is to reduce the number of system calls or I/O operations, improving the overall performance and efficiency of the system.

Here's how input buffering works:

- **Reading Input:** When a program needs to read input, it often uses system calls or I/O functions provided by the operating system. These calls are relatively expensive in terms of processing time.

- **Buffering:** Input buffering involves reading a larger chunk of data from the input source (e.g., a file or keyboard) than is immediately needed by the program. This larger chunk is stored in a buffer—a temporary storage area in memory.

- **Processing from Buffer:** The program then processes the input data from the buffer rather than making frequent I/O operations for small amounts of data. This reduces the number of system calls, which can be a significant performance improvement, especially when dealing with large volumes of input data.

- **Efficiency:** Input buffering is particularly beneficial in situations where reading data in small chunks would result in a high overhead due to frequent system calls. By reading and processing data in larger blocks, the program can make more efficient use of resources.

**Example in the Context of Compilers:** In the context of compilers, especially during the lexical analysis phase, input buffering is commonly employed. Here's how it works in the context of reading characters from the source code:
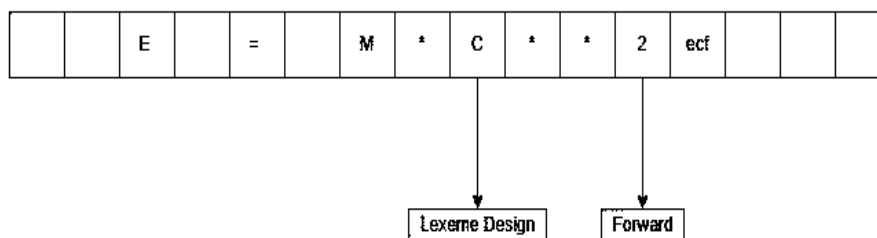
- Instead of reading one character at a time, which could lead to a large number of system calls and slow down the compilation process, a block or buffer of characters is read from the source file into memory.

- The lexical analyzer then reads characters from this buffer, and as it processes them, it advances a pointer within the buffer. When the buffer is exhausted, another block of characters is read into the buffer.

- This minimizes the overhead of reading individual characters and helps in the efficient processing of the source code.

    Specialized buffering techniques decrease the overhead required to process an input character in moving characters.

**Buffer Pairs:**

The input buffering mechanism consists of two buffers, each with a size of N characters, and they are reloaded alternately. Within this system, two pointers are employed: **lexemeBegin** and **forward**.

| | | E | | = | | M | * | C | * | * | 2 | ecf | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Lexeme Design     Forward

- LexemeBegin indicates the commencement of the current lexeme that is yet to be identified.
- Forward moves through the input buffer until it identifies a match for a specific pattern.
- Upon discovering a lexeme, lexemeBegin is adjusted to the character immediately following the newly identified lexeme, while forward is positioned at the rightmost character of that lexeme.
- The range between these two pointers encompasses the characters of the current lexeme. This approach enables efficient scanning and identification of tokens within the source code.

**Sentinels**

Sentinels In input buffering, sentinels are special markers or signals used to indicate the boundaries or endpoints of the data being processed. They help the system understand where one piece of information ends and the next one begins.

Let's imagine you're reading a book, and at the end of each chapter, there's a special symbol like "***" or "Chapter End." This symbol serves as a sentinel, telling you where one chapter ends and the next one starts. In a similar way, sentinels in input buffering help a computer system understand where one set of data or information in the input stream ends, and the next one begins.

In more technical terms, consider a scenario where you're processing a sequence of characters in a computer program. A special character, often called a sentinel, might be used to mark the end of a piece of data. For example, in strings, a null character ('\0') is commonly used as a sentinel to signify the end of the string. When the program encounters this character, it knows that it has reached the end of the current piece of data.

So, in input buffering, sentinels act like signposts, helping the system keep track of where one set of data ends and the next one starts, making it easier to read and process information.

## 2.5.9  Specification of Tokens

In the context of compilers and lexical analysis, a token is a sequence of characters in a source code file that represents a fundamental unit of meaning. Tokens are the building blocks of a programming language and are classified into different types, each serving a specific purpose in the syntax of the language. The specification of tokens involves defining the various types of tokens that can appear in the source code and the rules for recognizing them. Here's how the specification of tokens is typically done:

1  **Token Types:** Define the different types of tokens that can appear in the source code. Common token types include keywords, identifiers, literals, operators, punctuation symbols, and special symbols.
2  **Regular Expressions:** Use regular expressions to describe the patterns associated with each token type. Regular expressions are powerful tools for pattern matching and are used to define the syntactic structure of tokens. Each token type has a corresponding regular expression that specifies the allowed character patterns.
3. **Lexical Rules:** Specify lexical rules that define how tokens are formed from the input characters. Lexical rules describe the conditions under which a sequence of characters in the source code is recognized as a particular token type. These rules often include patterns defined by regular expressions.

4. **Token Attributes:** Define attributes associated with each token type. Token attributes provide additional information about the token, such as the value of a literal, the name of an identifier, or the specific operator represented by a token.

5. **Reserved Words:** Identify and list reserved words in the language. Reserved words are keywords that have special meaning in the programming language and cannot be used as identifiers.

6. **Error Handling:** Specify how lexical errors are handled. Lexical errors occur when the input does not match any of the defined token patterns. The specification should include rules for detecting and reporting lexical errors

Let's consider a simplified language with integer literals, identifiers, and basic arithmetic operators:

- **Token Types:**

  - INTEGER_LITERAL

  - IDENTIFIER

  - PLUS ( + )

- MINUS ( - )

- MULTIPLY ( * )

- DIVIDE ( / )

- **Regular Expressions:**

  - INTEGER_LITERAL: \d+

  - IDENTIFIER: [a-zA-Z][a-zA-Z0-9]*

  - PLUS: \+

  - MINUS: \-

  - MULTIPLY: \*

  - DIVIDE: \/

- **Lexical Rules:**

  - INTEGER_LITERAL: Matches one or more digits.

  - IDENTIFIER: Starts with a letter, followed by letters or digits.

  - PLUS: Matches the plus symbol.

  - MINUS: Matches the minus symbol.

  - MULTIPLY: Matches the asterisk symbol.

  - DIVIDE: Matches the forward slash symbol.

- **Token Attributes:**

  - INTEGER_LITERAL: Value of the integer.

  - IDENTIFIER: Name of the identifier.

  - PLUS, MINUS, MULTIPLY, DIVIDE: No additional attributes.

- **Reserved Words:**

  - None in this simple example.

This is a basic illustration, and in real-world scenarios, the token specification would be more comprehensive, covering a broader range of language constructs and incorporating more advanced lexical analysis techniques.

## 2.5.10 Recognition of Token

The recognition of tokens is a crucial step in the lexical analysis phase of a compiler. During this phase, the source code is scanned, and sequences of characters are identified as tokens based on the specifications provided in the token definition.

Finite Automata (FA) is a simple idealized machine that can be used to recognize patterns within input taken from a character set or alphabet (denoted as C). The primary task of an FA is to accept or reject an input based on whether the defined pattern occurs within the input.

There are two notations for representing Finite Automata. They are:

- Transition Table
- Transition Diagram

a general process for the recognition of tokens:

- **Lexical Analysis:** The process begins with the lexical analyzer (also known as the lexer or scanner) scanning the source code character by character.

- **Token Specification:** Referencing the token specification, the lexical analyzer uses regular expressions and rules to define the patterns associated with different token types.

- **Pattern Matching:** As characters are read from the source code, the lexical analyzer attempts to match them against the defined token patterns. Regular expressions play a crucial role in this step, as they describe the valid patterns for each token type.

- **Token Creation:** When a sequence of characters matches a token pattern, a token is created. The token includes the type of the token (e.g., identifier, keyword, operator), and in some cases, additional attributes such as the value of a literal or the name of an identifier.

- **Error Detection:** If a sequence of characters does not match any token pattern, the lexical analyzer may detect a lexical error. Error handling mechanisms are implemented to report and handle such errors, providing feedback to the programmer.

- **Token Stream:** The result of the recognition process is a stream of tokens, representing the fundamental units of meaning in the source code. This token stream is then passed on to the next phase of the compiler for further processing.

1. **Transition Table:**

    It is a tabular representation that lists all possible transitions for each state and input symbol combination.

    Example:

    Assume the following grammar fragment to generate a specific language

$$stmt \rightarrow \textbf{if } expr \textbf{ then } stmt \mid \textbf{if } expr \textbf{ then } stmt \textbf{ else } stmt \mid \varepsilon$$

$$expr \rightarrow term \textbf{ relop } term \mid term$$

$$term \rightarrow \textbf{id} \mid \textbf{number}$$

where the terminals if, then, else, relop, id and num generates sets of strings given by following regular definitions.

$$\textbf{if} \rightarrow \textbf{if}$$

$$\textbf{then} \rightarrow \textbf{then}$$

$$\textbf{else} \rightarrow \textbf{else}$$

$$\textbf{rebop} \rightarrow < \mid <= \mid < > \mid > \mid > =$$

$$\textbf{id} \rightarrow \textbf{letter ( letter} \mid \textbf{digit )*}$$

$$\textbf{num} \rightarrow \textbf{digits optional-fraction optional-exponent}$$

- where letter and digits are defined as - (letter → [A-Z a-z] & digit → [0-9])

- For this language, the lexical analyzer will recognize the keywords if, then, and else, as well as lexemes that match the patterns for relop, id, and number.

- To simplify matters, we make the common assumption that keywords are also reserved words: that is they cannot be used as identifiers.

- The num represents the unsigned integer and real numbers of Pascal.

- In addition, we assume lexemes are separated by white space, consisting of nonnull sequences of blanks, tabs, and newlines.

- Our lexical analyzer will strip out white space. It will do so by comparing a string against the regular definition ws, below.

- If a match for ws is found, the lexical analyzer does not return a token to the parser
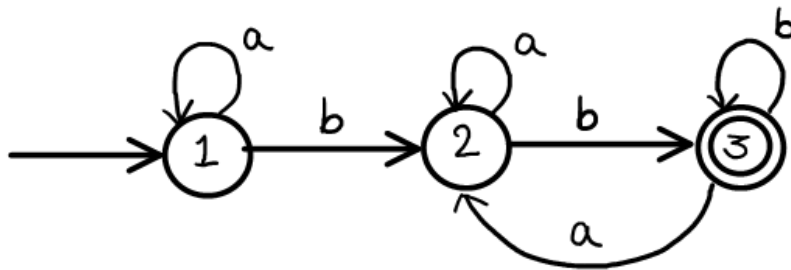
It is the following token that gets returned to the parser.

| LEXEMES | TOKEN NAME | ATTRIBUTE VALUE |
|---|---|---|
| Any ws | - | - |
| if | if | - |
| then | then | - |
| else | else | - |
| Any id | id | Pointer to table entry |
| Any number | number | Pointer to table entry |
| < | relop | LT |
| <= | relop | LE |
| = | relop | EQ |
| <> | relop | NE |

## 2. Transition Diagram

It is a directed labeled graph consisting of nodes and edges. Nodes represent states, while edges represent state transitions.
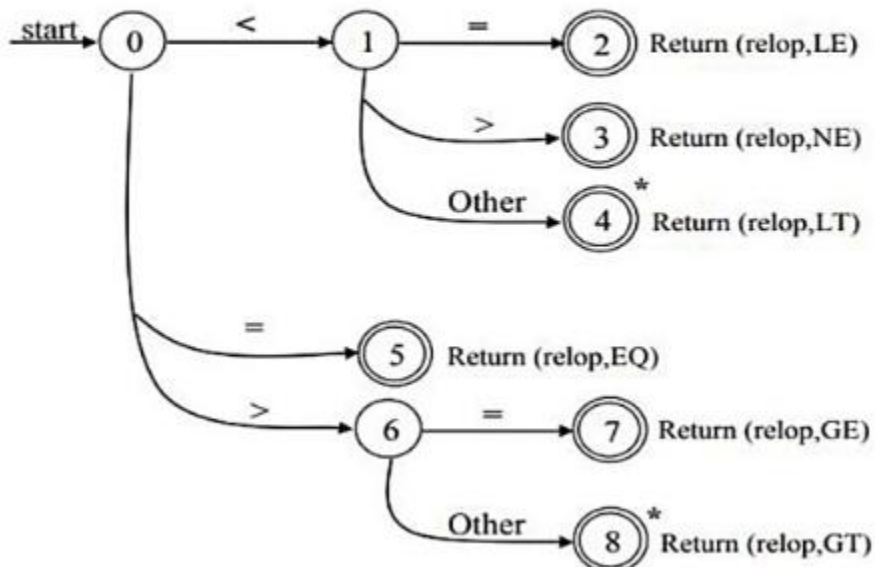
1. One state is labelled the **Start State**. It is the initial state of transition diagram where control resides when we begin to recognize a token.
2. Position is a transition diagram are drawn as circles and are called **states**.
3. The states are connected by Arrows called **edges**. Labels on edges are indicating the input characters
4. Zero or more **final states** or **Accepting states** are represented by double circle in which the tokens has been found.

the transition diagram of Finite Automata that recognizes the lexemes matching the token relop.

**Example:** A Transition Diagram for the token relation operators
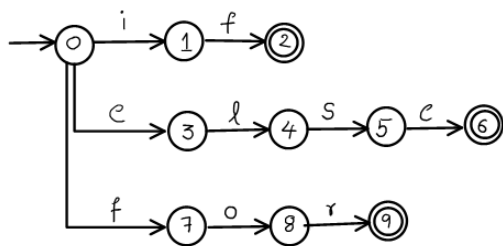
**"relop"** is shown in Figure below:



the Finite Automata Transition Diagram for the Identifiers and Keywords.

Letter⟶ a|b|........|z|      A|B|........|Z|
Digit ⟶ 0|1|........|9|
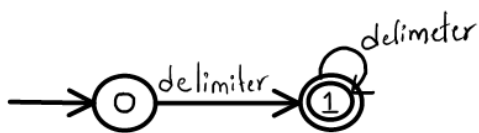Id    ⟶ letter ( letter | digit )*

abc12

**Here is the Finite Automata Transition Diagram for recognizing white spaces.**

WS $\longrightarrow$ delimiter (delimiter)*



**Write regular expression for number**

Example:- Reg. Expr for digits ( int no, floating point no)

$$\text{digit} \longrightarrow 0|1|........|9$$
$$\text{digits} \longrightarrow \text{digit (digit)*}$$
$$\text{Number} \longrightarrow \text{digits ( . digits )?  (E [+ -]? digits)}$$

123  456
123.45
123.45 E.23
123.45 E-23
123 E+23

5. **Recognition of numbers (int | flating points)**

$$\text{Number} \longrightarrow \text{Digits}^{+} \; (\text{digit}^{+}) \; ? \quad (E \; [+ -] \; ? \; \text{digit}^{+}) \; ?$$

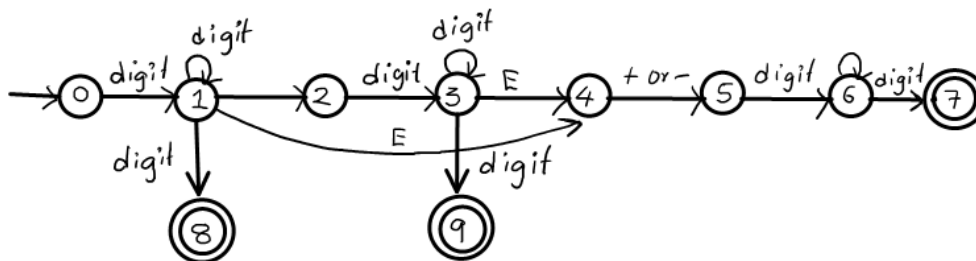123
123.45
123 E 20
123.45 E +2
123.45 E -2

**Recognize the given token in python**

```
import re

def tokenize(input_string):
    # Define regular expressions for different types of tokens
    token_patterns = [
        (r'\bif\b|\belse\b|\bwhile\b|\bfor\b', 'KEYWORD'),
        (r'\b\d+\b', 'NUMBER'),
        (r'\b[a-zA-Z_]\w*\b', 'IDENTIFIER'),
        (r'[-+*/=<>]', 'OPERATOR'),    # Updated to include comparison operators
        (r'\s', 'WHITESPACE'),
        (r'\(|\)', 'PARENTHESIS'),    # Add a pattern for parentheses
        (r':', 'COLON'),    # Add a pattern for colon
    ]

    tokens = []

    while input_string:
        for pattern, token_type in token_patterns:
            match = re.match(pattern, input_string)
            if match:
                value = match.group(0)
                tokens.append((value, token_type))
                input_string = input_string[len(value):].lstrip()
                break
        else:
            # If no match is found, raise an error or handle accordingly
            raise ValueError(f"Unable to tokenize: {input_string}")

    return tokens

# Example usage
input_string = "if x > 5 for x in range(10): print(x)"
tokens = tokenize(input_string)

for token in tokens:
    print(token)
```

**OUTPUT**:

('if', 'KEYWORD')
('x', 'IDENTIFIER')
('>', 'OPERATOR')
('5', 'NUMBER')
('for', 'KEYWORD')
('x', 'IDENTIFIER')
('in', 'IDENTIFIER')
('range', 'IDENTIFIER')
('(', 'PARENTHESIS')
('10', 'NUMBER')
(')', 'PARENTHESIS')
(':', 'COLON')
('print', 'IDENTIFIER')
('(', 'PARENTHESIS')
('x', 'IDENTIFIER')
(')', 'PARENTHESIS')