



Theory of Computation & Compiler Design

CONTENTS

01. Symbol	01
02. Alphabets	01
03. String	01
04. Language	02
05. Power of Sigma	02
06. Finite Automata Model	03
07. State Finite Automata Model	03
08. DFA Model	04
09. How To Construct A DFA	06
10. Minimization of DFA	07
11. NFA Figure Construction	09
12. NFA to DFA Construction	10
13. Regular Expression	12
14. Regular Expression To NFA	14
15. Context-Free Grammar (CFG)	20
16. Language Processing System	23
17. Structure of a Compiler	24

Symbol

Defination

A symbol (also called a character) is the smallest building block in a system, which can be any alphabet, letter, number, or picture used to represent information or construct larger structures like strings or languages.

Examples of Symbols:

- i. Alphabets (a to z): a,b,.....,z
- ii. Digits (0 to 9): 0,1,2,.....,9

Why uppercase letters aren't symbol?

Alphabets

Defination

An alphabet (also called an input in some contexts) is a finite set of distinct symbols, typically denoted by Σ (Greek letter sigma), used to construct strings and define languages. These symbols can be letters, digits, or special characters, and they serve as the basic building blocks for representing information in computational systems. The symbols in the alphabet are elements of the set.

Examples of Alphabets:

1. $\Sigma=\{a,b,c\}$

This is an alphabet containing three symbols: a,b and c.

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

This is an alphabet containing three symbols: 0 and 1.

3. $\Sigma=\{a,b,c,0,1\}$

This is an alphabet containing three symbols: a,b,c,0 and 1.

4. $\Sigma=\{a,a,b,c\}$

An alphabet must contain distinct symbols. In this case, adding a twice is redundant, and the set should be reduced to:

$$\Sigma=\{a,b,c\}$$

This is an alphabet containing three symbols: a,b and c.

String

Defination

A string is a finite sequence of symbols (or characters) taken from an alphabet Σ . It can include any combination of symbols from the alphabet, and its length can vary from zero (an empty string) to any finite number of symbols.

Notation:

1. A string is typically denoted by a sequence of characters inside quotation marks, like $w = \text{"abc"}$
2. If $\Sigma=\{a,b,c\}$, then $w=\text{"abc"}$ is a string made from those symbols.

Key Points:

1. Finite Sequence: A string is a sequence of symbols, and the length of the string is finite.
2. From an Alphabet: The symbols in a string come from a predefined alphabet Σ .
3. Empty String: The empty string, denoted by ϵ , is a string that contains no symbols. It has a length of zero.
4. Length of a String: The length of a string w , denoted as $|w|$, is the number of symbols it contains. For example, $|\text{"abc"}|=3$.

Language

Defination

A language is a set of strings formed from an alphabet that satisfies specific rules or properties.

Languages are two types: 1. Finite language 2. Infinite Language

A language is often denoted by L

Examples:

$$\Sigma = \{a, b\}$$

$L1 = \{W \text{ is a set of string where the length of each string must be } 2\}$

$$= \{aa, ab, ba, bb\}$$

$L2 = \{W \text{ is a set of string where the length of each string must be } 3\}$

$$= \{aaa, aab, bab, bba, abb, aba, bbb, baa\}$$

$L3 = \{W \text{ is a set of string where each string starts with 'a'}\}$

$$= \{a, ab, aa, aaa, \dots\}$$

Here, L1 and L2 is called the finite language and L3 is infinite language.

Power of Sigma

Suppose, $\Sigma = \{a, b\}$

$\Sigma^1 = \text{Set of all string over this } \Sigma \text{ of the length is '1'}$
 $= \{a, b\}$

$\Sigma^2 = \text{Set of all string over this } \Sigma \text{ of the length is '2'}$
 $= \{aa, ab, ba, bb\}$

Note: $\Sigma^0 = \text{empty string where the length of string is } 0$

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

Σ^* is called mother of all string.

Finite Automata Model

Defination

A Finite Automata Model (FA Model) is a mathematical model used to recognize patterns in input, such as strings.

Types

1. Deterministic Finite Automaton (DFA)
2. Nondeterministic Finite Automaton (NFA)

Why FA Model is used?

To give a finite representation of an infinite language

Example:

$L(M) = \{W \text{ is a set of string where each string starts with 'a'}\}$ and $\Sigma = \{a,b\}$
 $W = \{a, ab, aa, aaa, \dots\}$



Figure: State Transition Diagram

This diagram accepts only those strings that start with 'a' and rejects those that start with 'b'. It is impossible to allocate infinite memory space for a language set. Therefore, a finite memory slot is used. For this reason, Finite Automata is implemented, as it requires very little memory and can efficiently determine whether a string is accepted or not.

States Finite Automata Model

States mean memory, which is denoted by a circle (\bigcirc) = {A, B, C}

1. Initial State

- i. Start scanning or traveling.
- ii. A DFA (Deterministic Finite Automaton) and an NFA (Nondeterministic Finite Automaton) both contain one initial state.
- iii. This arrow signifies the starting point of the automaton. Additionally, there is no label (value) on the edge leading to the initial state

2. Final State

- i. Stop scanning or traveling.
- ii. Denoted by double circle (\odot)
- iii. DFA (Deterministic Finite Automaton) and NFA (Nondeterministic Finite Automaton), there can be one or more final states.

3. Trap State (or Dead State)

- i. A trap state (or dead state) always has a self-loop for all input symbols.
- ii. We cannot go to the final state from a trap state
- iii. It is not mandatory to use a trap state in a DFA or NFA
- iii. When a rejected string is processed, the automaton eventually reaches the trap state (also called the dead state).

4. Normal State

- i. A state that handles input symbols and is not a final or trap state.

DFA Model

A DFA is a type of Finite Automaton (FA) that consists of five tuples:

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

Example:

Let $\Sigma = \{a, b\}$

Define the language $L(M)$ as the set of all strings that start with the letter 'a'.
Now, consider the state transition diagram for $L(M)$:

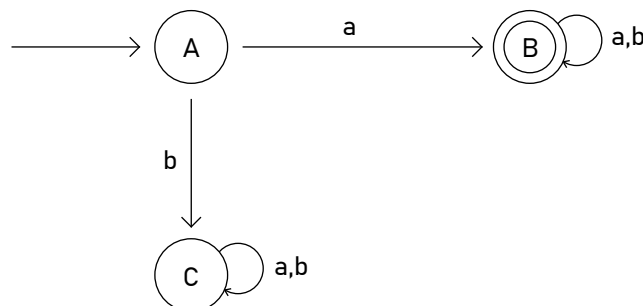


Fig: DFA Diagram

State Transition Diagram for $L(M)$

The five components of the DFA are described as follows:

Q - The set of all states: $Q = \{A, B, C\}$

Σ - The set of input alphabets: $\Sigma = \{a, b\}$

q_0 - The start (initial) state: $q_0 = A$

F - The set of final states: $F = \{B\}$

δ - The transition function

The transition function δ is defined as:

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

The Transition Function (δ) in a DFA can be represented in two ways:

1. Transition Table (Tabular Form)
2. Using Equation

1. Transition Table (Tabular Form)

This method uses a table where rows represent states and columns represent input symbols. Each cell shows the next state for a given state-input pair.

State ↓ \ Input →	a	b
A	B	C
B	B	B
C	C	C

Fig: DFA Transition Table

2. Using Equation

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

$$A \times a \rightarrow B$$

$$A \times b \rightarrow C$$

$$B \times a \rightarrow B$$

$$B \times b \rightarrow B$$

$$C \times a \rightarrow C$$

$$C \times b \rightarrow C$$

Characteristics of DFA:

1. For each input symbol, there is exactly one transition to a state.
2. A DFA has only one initial state.
3. A DFA can have more than one final state.

How To Construct A DFA

Problem 1

Construct a DFA which accepts the set of all strings over the $\Sigma=\{a,b\}$ where each string starts with an 'a'.

Solution

Alphabet: $\Sigma=\{a,b\}$

Language: $L(M)=\{a,ac,ab,aaa,aab,aba,abb,\dots\}$



Fig: DFA Diagram

Problem 2

$L(M)=\{w \mid w \text{ starts with 'aa'}\}$

Solution

Alphabet: $\Sigma=\{a,b\}$

Language: $L(M)=\{aa,aaa,aab,aaab,aabb,\dots\}$



Fig: DFA Diagram

Problem 3

Construct a DFA which accepts the set of all strings over $\Sigma=\{0,1\}$ where each string starts with 10

Solution

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

Language: $L(M)=\{10,100,101,1001,1011,10001,\dots\}$



Fig: DFA Diagram

Minimization of DFA

Example 01

Given the language:

$L = \{w \mid w \text{ has an odd number of a's and ends with 'a' or 'b'}\}$



Fig: DFA Diagram

Steps for Minimizing a DFA

1. Remove Unreachable States:

Identify and eliminate any states that cannot be reached from the initial state.

2. Construct the State Transition Table:

Create a table representing state transitions for each input symbol.

3. Mark Initial and Final States:

Clearly indicate which states are initial and final.

4. Find Equivalent States:

Identify and merge states that belong to the same equivalence class.

Now, proceed with minimizing the DFA of Example 01 using these steps.

Step 1: Remove unreachable states

Remove those states which are unreachable from the initial state. No states are removed here because states 1B, 2A, and 2B are reachable from the initial state.



Fig: DFA Diagram

2. Construct the State Transition Table:

State ↓ \ Input →	a	b
1A	2A	1B
1B	2A	1B
2A	1A	2B
2B	1A	2B

Fig: DFA Transition Table

3. Mark Initial and Final States:

Clearly indicate which states are initial and final.

State ↓ \ Input →	a	b
→ 1A	2A	1B
1B	2A	1B
2A	1A	2B
* 2B	1A	2B

4. Find Equivalent States:

Identify and merge states that belong to the same equivalence class.

$$\begin{aligned} 0\text{-equivalent} &= [\text{Set of non-final states}] [\text{Set of final state}] \\ &= [1A \ 1B \ 2A] [2B] \end{aligned}$$

$$1\text{-equivalent} = [1A \ 1B] [2A] [2B]$$

$$2\text{-equivalent} = [1A \ 1B] [2A] [2B]$$

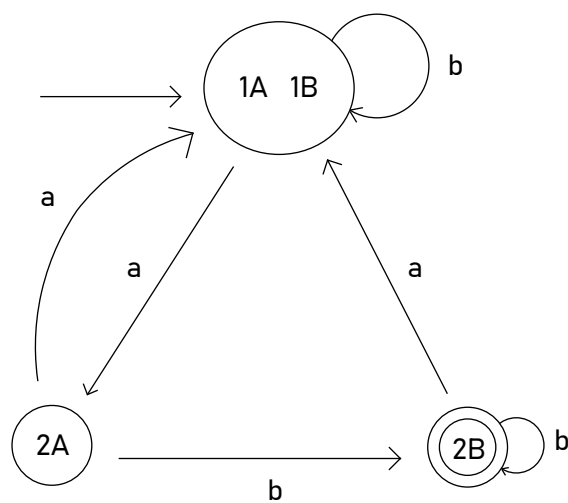


Fig: Minimized DFA Diagram

NFA Figure Construction

Rules for Each Input

1. No Transition:

1.1 Imagine you're walking on a path, but suddenly, the road ends. You have no way to move forward.

1.2 In an NFA, if there is no transition for a given input, the machine cannot continue from that state.

2. Single Transition:

2.1 Imagine you're at a door, and there is only one way to go forward.

2.2 In an NFA (or DFA), if a state has exactly one transition for a symbol, it moves to only one specific next state.

3. Multiple Transitions:

3.1 Imagine you're at a junction where you can go in two or more different directions at the same time.

3.2 In an NFA, a state can have multiple paths for the same input symbol, meaning the machine can follow multiple options at once.

Example 01

$$L(M) = \{ w \mid w \text{ starts with 'a'} \}$$

$$\Sigma = \{a, b\}$$

$$W = \{ a, ab, aba, abb, \dots \}$$



Fig: NFA Diagram

Example 02

$$L(M) = \{ w \mid w \text{ is a string that ends with 'a'} \}$$

$$\Sigma = \{a, b\}$$

$$W = \{ a, ba, aa, aba, bba, abba, \dots \}$$



Fig: NFA Diagram

NFA to DFA Construction

Language

$$L(M) = \{ w \mid w \text{ ends with } ab \} \text{ where } \Sigma = \{a, b\}$$

Solution

1. Construct the NFA

Create an NFA for the given language. The transition diagram consists of states A, B, and C, with the following transitions:



Fig: NFA Diagram

2. State Transition Table for NFA

State ↓ \ Input →	a	b
A	{A, B}	{A}
B	{ }	{C}
C	{ }	{ }

Fig: NFA Transition Table

3. Construct the DFA Transition Table

State ↓ \ Input →	a	b
A	[A B]	[A]
B	[A B]	[A C]
C	[A B]	[A]

Fig: DFA Transition Table

5. Construct the DFA Figure

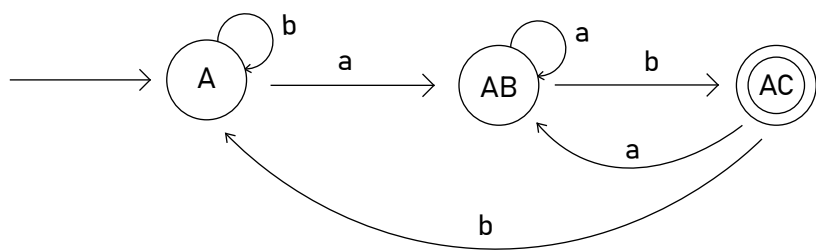


Fig: DFA Diagram

Regular Expression

A regular expression is a method or representation technique used to match patterns on strings. It is widely used in compiler design, is easy to understand, and has efficient implementation.

Formal Definition (6 Cases):

Case 1

$R = a$ or $R = b$ where $a, b \in \Sigma$

A language that accepts only one single character: either 'a' or 'b'.

If $R = a$, then the language is: $L = \{ "a" \}$



If $R = b$, then the language is: $L = \{ "b" \}$



Case 2

$R = \epsilon$

The empty string ϵ is a string with length 0 — it contains no characters.

The language is: $L = \{ \epsilon \}$

That means the NFA should accept the empty string, and only the empty string.



Case 3

$R = \emptyset$

\emptyset represents the empty set — a language that contains no strings at all, not even ϵ . So the NFA should reject everything.

Case 4

Union ($R = R_1 \cup R_2$ or $R_1 \mid R_2$ or $R_1 + R_2$)

Let's take:

$R_1 = ab$

$R_2 = bb$

Build NFA for $R_1 = ab$



Build NFA for $R_2 = bb$



Now we have to union N1 and N2



When we perform a union, we create an extra state and connect it to the initial states of both NFAs using ϵ .

Case 5

$R = R_1 \cdot R_2$



Why ϵ -transition?

Because after finishing R_1 , the machine jumps to R_2 without reading any input — and continues reading.

When we concatenate two NFA models, the second NFA remains unchanged. We make the final state of the first NFA non-final and connect them using an ϵ .

Case 6: Kleene Star

$R = R_1^*$ This means: 0 or more repetitions of R_1

Let's say

$R_1 = ab$, so $R = (ab)^*$

That means:

Accept: "", "ab", "abab", "ababab", etc.

Build NFA for $R_1 = ab$



Build NFA for $R_1 = (ab)^*$



When we apply the Kleene Star, we create an extra final state, but the previous final state remains unchanged. We connect the new final state to the NFA using ϵ -transitions. Additionally, we draw an ϵ -transition from the old final state back to the initial state.

Regular Expression To NFA

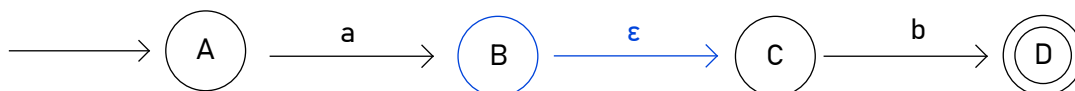
Example 01: $(ab \cup a)^*$

$a^* = 0$ or more
 $= \{\epsilon, a, aa, aaa, \dots\}$
 $a^+ = 1$ or more
 $= \{a, aa, aaa, \dots\}$

Plot NFA for each symbol. Here, symbol is a and b.



here a and b are concatenated, so **Case 5** is used here



(ab U a)

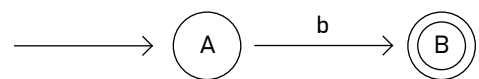
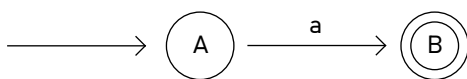


(ab U a)*

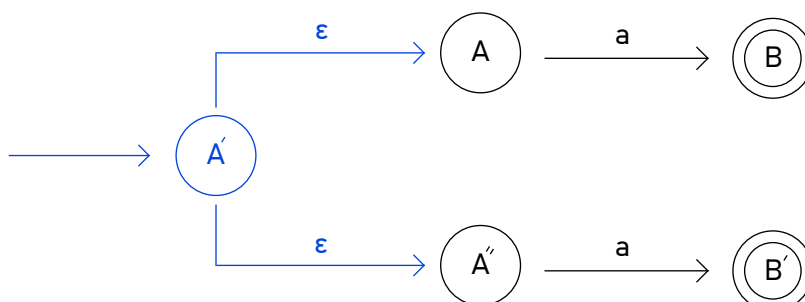


Example 02: (a U b)* aba

plot NFA for each symbol. Here, symbol is a and b.



(a U b)



$(a \cup b)^*$



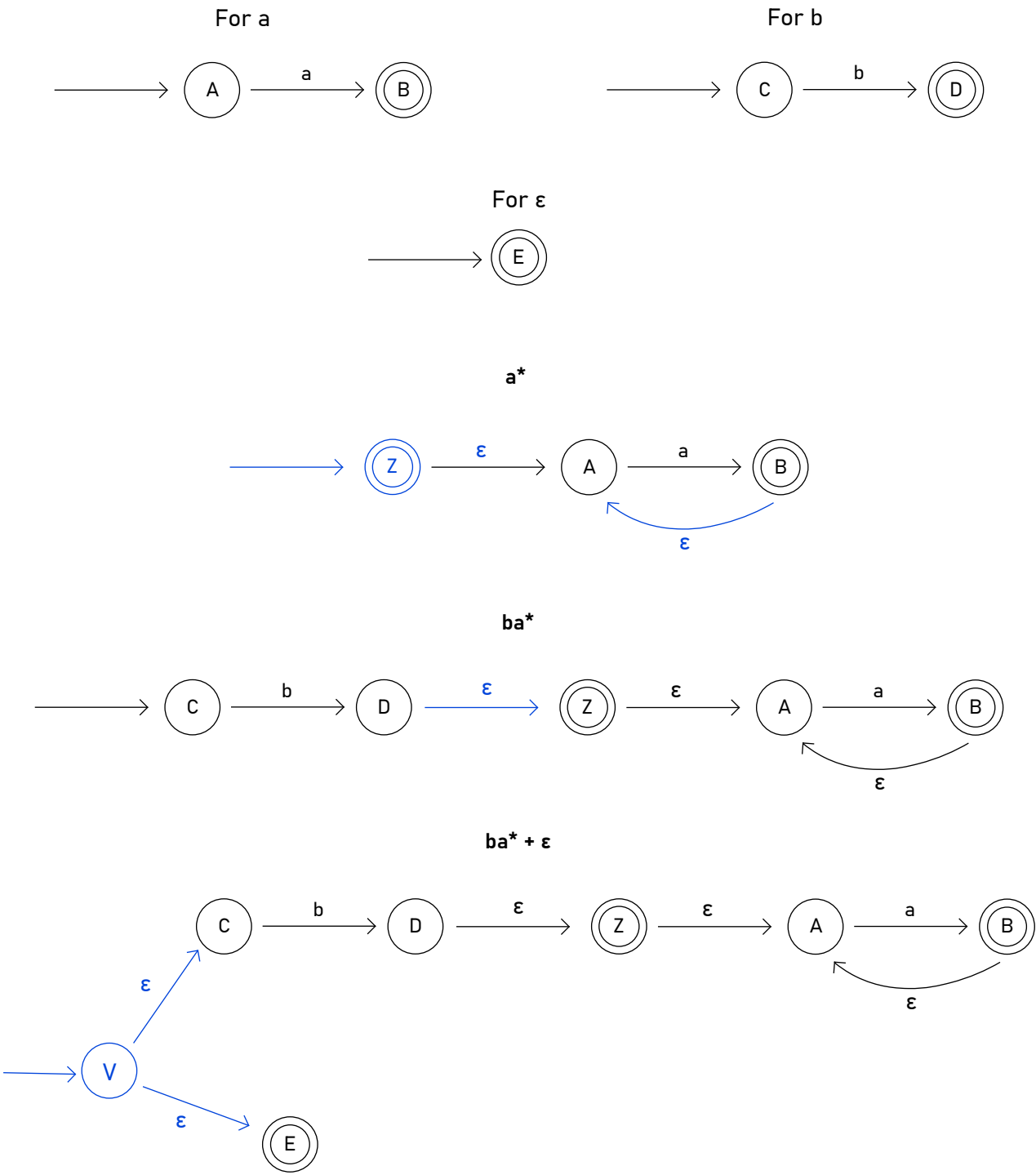
aba



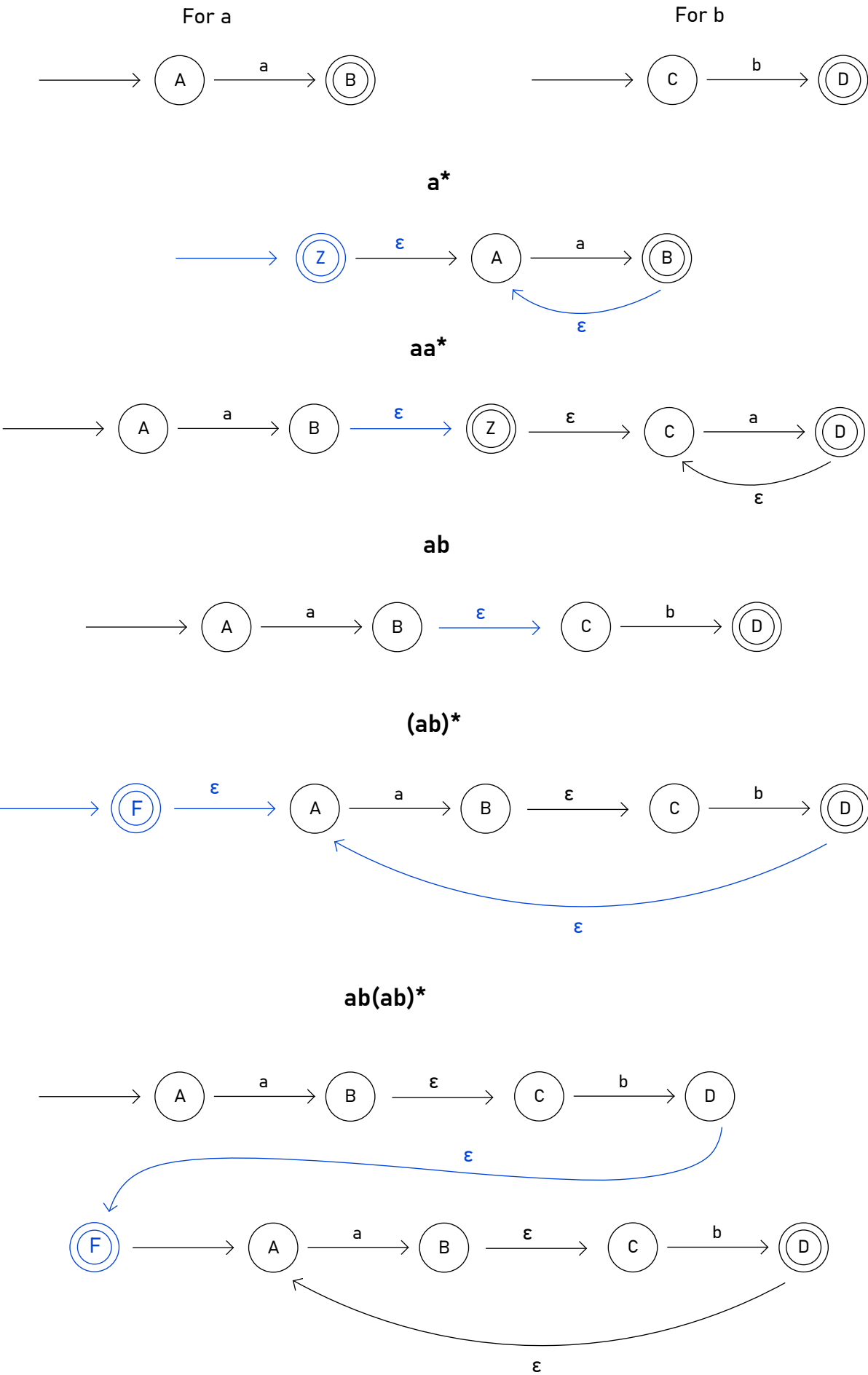
$(a \cup b)^* aba$



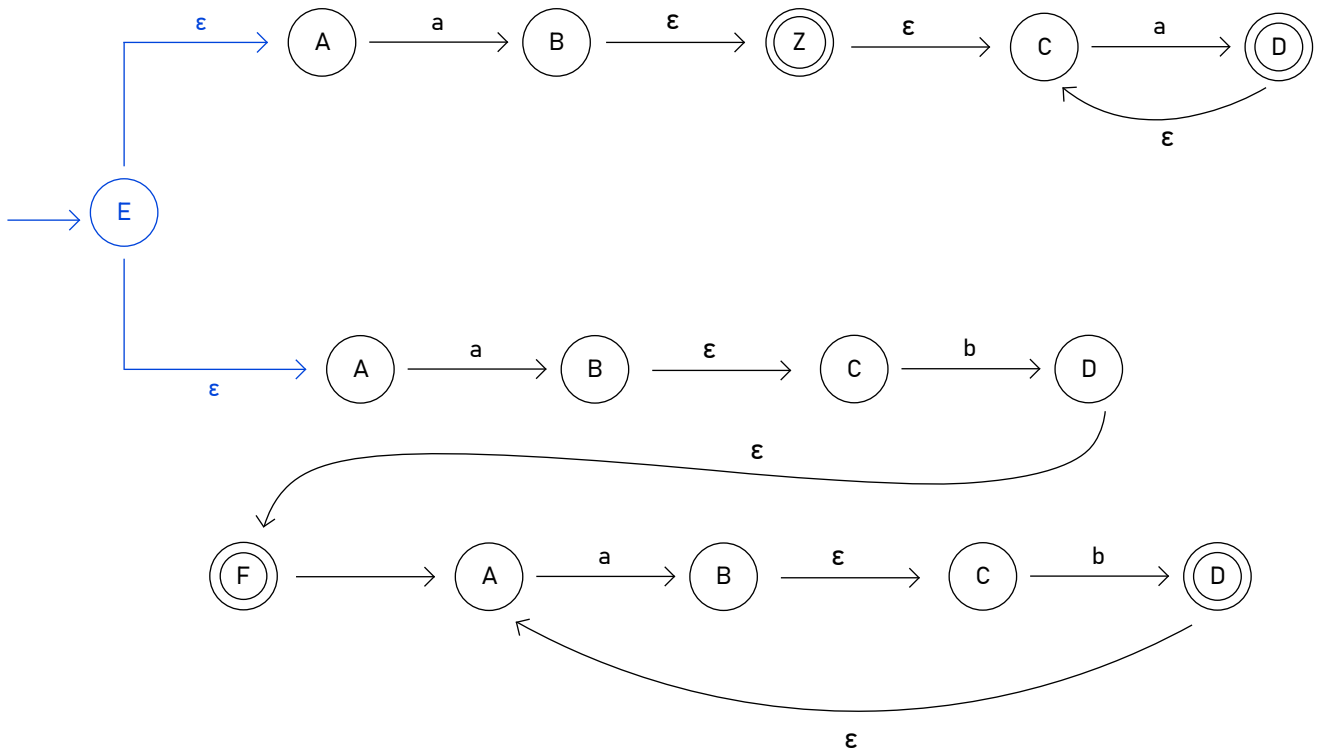
Example 03: $ba^* + \epsilon$



Example 04: $a \cup (ab)^+$
 $aa^* \cup ab(ab)^*$



$aa^* \cup ab(ab)^*$



Context-Free Grammar (CFG)

Context Free Grammar is a set of rules/production which is used to generate patterns of string.

CFG has 4-tuples. These are V, T, P, S

Description	
V	Variable / Non-Terminal
T	Set of Terminals
P	Production / rules
S	Start-Variable

Now, from a grammar, we introduce with four-tuples.

Suppose the grammar is, $S \rightarrow 01 \mid 0S1$

This grammar is for the language, $L(M) = \{ w \mid w \text{ is a set of string of } 0^n 1^n \text{ where } n \geq 1 \}$

$$0^1 1^1 = 01$$

$$0^2 1^2 = 0011$$

The grammar is defined with the production rule: $S \rightarrow 01 \mid 0S1$

The 4-tuples are determined as

$V = \{S\}$ (set of variables/non-terminals)

$T = \{0, 1\}$ (set of terminals)

$S = S$ (start variable)

$P = \{S \rightarrow 01 \mid 0S1\}$ (set of production rules)

Examples:

Identify V, T, P, S from the grammar G:-

$$\begin{aligned} 1) E &\Rightarrow E + T \mid T \\ T &\Rightarrow T * F \mid F \\ F &\Rightarrow (E) \mid \text{id} \end{aligned}$$

$$\begin{aligned} V &= \{E, T, F\} \\ T &= \{+, *, (,), \text{id}\} \\ P &= \{E \Rightarrow E + T, E \Rightarrow T, T \Rightarrow T * F, T \Rightarrow F, F \Rightarrow (E), F \Rightarrow \text{id}\} \\ S &= E \end{aligned}$$

Identify V, T, P, S from the grammar G:-

$$\begin{aligned} 2) S &\Rightarrow (L) \mid a \\ L &\Rightarrow L, S \mid S \end{aligned}$$

$$\begin{aligned} V &= \{S, L\} \\ T &= \{ (,), a, , \} \\ P &= \{S \Rightarrow (L), S \Rightarrow a, L \Rightarrow L, S, L \Rightarrow S\} \\ (S &= S \end{aligned}$$

USING Context-free grammar, we can match pattern of a string in two ways:-

- 1) Mathematical derivation
- 2) Parsing Tree derivation

Now, using CFG, we match the pattern of a string in two ways for the following languages.

Example 01:

$$\begin{aligned} L &= \{a^n \text{ where } n > 0\} \\ &= \{\epsilon, a, aa, aaa, \dots, \dots\} \end{aligned}$$

$$\begin{aligned} \text{Grammar for the language is,} \\ A &\Rightarrow aA \mid \epsilon \end{aligned}$$

NOW, using the grammar, we match the pattern of string in two ways,

Mathematical Derivation: Suppose, here we match the pattern of a string which is 'aaa'

$$\begin{aligned} A &\Rightarrow aA \\ A &\Rightarrow aaA \\ A &\Rightarrow aaaA \\ A &\Rightarrow aaa\epsilon \\ A &\Rightarrow aaa \end{aligned}$$

Parsing Tree Derivation: Using this derivation technique, now we match the pattern of a string which is 'aaa'



From the above, two derivation techniques, we see that the string pattern 'aaa' is matched by two techniques / CFG techniques. So it can be said that, the grammar $A \Rightarrow aA \mid \epsilon$ is correct for the language, $L(M) = \{a^n \mid n > 0\}$

Example 02:

Match the pattern of the given string 'abbaababba' using context free grammar:-

$$\begin{aligned}
 S &\Rightarrow aSa \mid bSb \\
 S &\Rightarrow \epsilon
 \end{aligned}$$

Show both mathematical and parsing tree derivation

Mathematical: [abbaabba]

$$\begin{aligned}
 S &\Rightarrow aSa \\
 S &\Rightarrow a \text{ } bSb \text{ } a \\
 S &\Rightarrow a \text{ } b \text{ } bSb \text{ } ba \\
 S &\Rightarrow abb \text{ } aSa \text{ } bba \\
 S &\Rightarrow abba \text{ } \epsilon \text{ } abba \\
 S &\Rightarrow abbaabba
 \end{aligned}$$

Parsing tree:



Structure of a Compiler



Lexical Analyzer (Scanner)

1. Processes the pure high-level language input.
2. Converts the input into a stream of tokens.

Syntax Analyzer

1. Takes the stream of tokens as input.
2. Produces a syntax tree (also known as a parsing tree).

Semantic Analyzer

1. Verifies the syntax tree for semantic correctness.
2. Outputs a verified syntax tree or parsing tree.
3. Generates intermediate code.

Code Generator

1. Converts the intermediate code into three-address code.
2. Follows three-address code rules for this conversion.

Code Optimizer

1. Optimizes the three-address code.
2. Produces optimized code for better performance.

Target Code Generator

1. Converts the optimized code into assembly code.
2. Outputs the final assembly code

Lexical Analyzer (Scanner)

Lexeme

A lexeme is a sequence of characters in the source code that forms a single unit of meaning during the lexical analysis phase of compilation. It's essentially the raw input that the lexical analyzer (scanner) identifies and categorizes into a token. For example, in the code `int x = 5;`, the lexemes would be `int`, `x`, `=`, `5`, and `;`. Each lexeme is then classified into a token type (like keyword, identifier, operator, or literal) for further processing by the compiler.

A pattern is a rule defining the structure of a lexeme for a token type.

Example: For identifiers, the pattern is `letter (letter | digit)*`.

Starts with a letter `[a-zA-Z]`.

Followed by zero or more letters `[a-zA-Z]` or digits `[0-9]`.

Valid examples: `x`, `var`, `myVar123`.

Invalid examples: `1var` (starts with a digit), `_var` (starts with an underscore, if not allowed).

Function of the Lexical Analyzer

The Lexical Analyzer, or scanner, breaks down the source code into a sequence of tokens by:

1. Reading the input character by character.
2. Grouping characters into lexemes (meaningful units like keywords, identifiers, operators).
3. Matching lexemes against predefined patterns.
4. Assigning each lexeme to a token type (e.g., KEYWORD, IDENTIFIER, OPERATOR).
5. Removing whitespace, comments, and newline characters.
6. Passing the stream of tokens to the Syntax Analyzer.

How the Lexical Analyzer Performs Its Job

```
float position, initial, rate;  
position = initial + rate * 60
```

ID/ Number / Serial Number	Data Type	Variable / Function
01	Float	position
02	Float	pnitial
03	Float	rate

Token

Syntax : <token_name, token_number>

Note Under Construction