# Lexical Analysis

→ The scanning / lexical analysis phase of a compiler performs the task of reading a stream of characters as an input and produce a sequence of tokens such as names, keywords, numbers etc for syntax analyzer.

→ It discards the white space and comments between the tokens and also keep tracks of line number.

→ lexical analyzer correlate error message with source program.

## Role of lexical Analyzer

→ lexical analyzer is the first phase of a compiler.

→ Its main task is to read the input characters and produce as output a sequence of tokens that parser uses for syntax analysis.
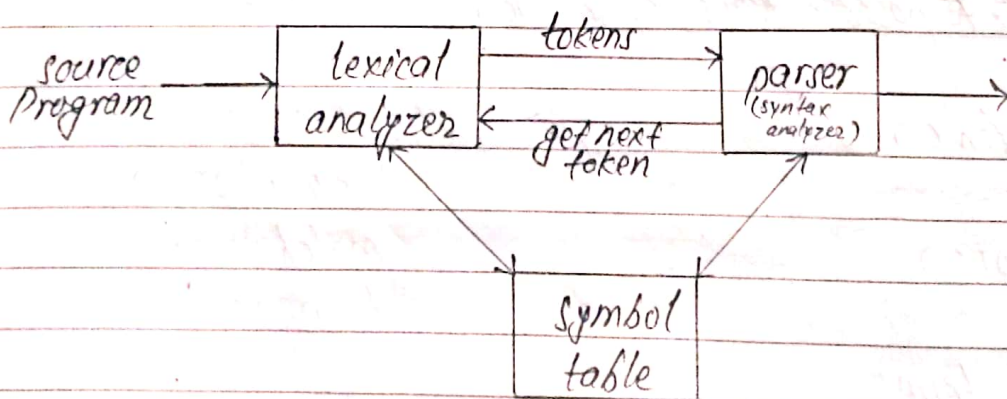
fig: lexical Analyzer

→ As in the figure, upon receiving a "get next token" command from the parser the lexical analyzer read input characters until it can identify the next token.

→ It removes the comments from the source program.

→ It keeps track of line numbers while scanning the new line characters. These line numbers are used by the error handler

to print the error message.


## Tokens, Patterns, Lexemes

→ _Token_ is a sequence of characters that can be treated as a single logical entity.
  e.g. Identifiers, keywords, operators, constants etc.

✦ e.g. of non-tokens:
  Comments, preprocessor directive, macros, blanks, tabs, newline etc.

→ _Pattern_: A set of string in the input for which the same token is produced as output. This set of string is described by a rule called a pattern associated with the token.
  e.g. The pattern for the Pascal identifier token, id is:
  $$id \rightarrow letter\ (letter\ |\ digit)^*$$  i.e. letter followed by letter & digits.


→ A _lexeme_ is a sequence of characters in the source program that is matched by the pattern for a token.
  e.g. The pattern for the RELOP token contains six lexemes
  $$=,\ <=,\ <>,\ >,\ <,\ >=$$
  → In SQL
  ↳ != (In C)


# Input : x = x * (acc + 123)

| token | lexemes | token | lexemes | token | lexemes |
|---|---|---|---|---|---|
| identifier | x | star | * | plus | + |
| equal | = | left-paren | ( | integer | 123 |
| identifier | x | identifier | acc | right-paren | ) |

# Attribute of tokens

When a token represents more than one lexeme, lexical analyzer must provide additional information about the particular lexeme. This additional information is called as the attribute of the token.

→ Attributes are used to distinguish different lexemes in a token.

Some attributes:
- <id, attr>   where attr is pointer to the symbol table.
- <assign-op, >   no attribute is needed (if there is only one assignment operator)
- <num, val>   where val is the actual value of the number.

→ Token type and its attribute uniquely identifies a lexeme.

e.g.  dest = source + 5
Tokens:  < id, pointer to symbol_table entry for dest>
         < assign-op,>
         <num, integer val 5
         < id, pointer to symbol_table entry for source>
         < add-op, >
         <num, integer val 5>

E = M * C ** 2
    < id, pointer to symbol-table entry for E >
    < assign-op,>
    < id, pointer to symbol-table entry for M>
    <mult-op, >
    < id, pointer to symbol-table entry for C>
    <exp-op, >
    < num, integer value 2 >

# Lexical Error

→ During the lexical analysis phase this type of error can be detected.

→ lexical error is a sequence of characters that doesnot match the pattern of any token. lexical phase error is found during the execution of the program.

→ lexical phase error can be:

- spelling error
- Exceeding length of identifier or numeric constant
- To remove the character that should be present
- Appearance of illegical characters.
- To replace a character with an incorrect character
- ~~Transposition of two characters~~

E.g.

```
void main()
{

    int x=10, Y=20;
    char *a;
    a=&x;
    x = 5xab;
}
```

In this code, 5xab is neither a number nor an identifier. so this code will show the lexical error.

⇒ Possible error recovery actions are:

- Panic mode recovery - deleting successive characters ~~un~~ until a well formed token is formed.
- Inserting a missing character.
- Replacing a missing character by a correct character.
- Transposing two adjacent character.
- Deleting an extraneous character.

# General approaches to the implementation of a lexical analyzer

There are three general approaches to the implementation of a lexical analyzer.

1. Use a lexical-analyzer generator like lex or flex, to produce lexical analyzer from a regular-expression based specification. The generator provides routines for reading & buffering the input. (easiest to implement; least efficient).

2. Write the lexical analyzer in high level programming language like c, using the I/O facilities of that language to read and buffering the input. (Intermediate in ease, efficiency).

3. Write the lexical analyzer in assembly language and explicitly manage the input and buffering. (Hardest to implement, but most efficient)

## Lookahead and Buffering

→ Many times, a scanner has to look ahead several ~~times~~ characters from the current character in order to recognize the token.

For e.g. 'int' is keyword in c, while the term 'inp' may be a variable name. When the character 'i' is encountered, the scanner cannot decide whether it is a keyword or a variable name until it reads two more characters.

→ In order to efficiently move back and forth in the input stream, input buffering is used.

# Input Buffering

→ lexical analysis needs to look ahead several characters before a match can be announced.

→ we have two buffer input scheme that is useful when look ahead is necessary
- Buffer Pair
- Sentinels

## Buffer Pair (2N Buffering)

In this technique buffer is divided into two halves with N-characters each, where N is number of characters on the disk block like 1024 or 4096. Rather than reading character by character from file we ~~read N~~ read N input character at once. If there are fewer than N character in input eof marker is placed. There are two pointers : lexeme pointer and forward pointer.
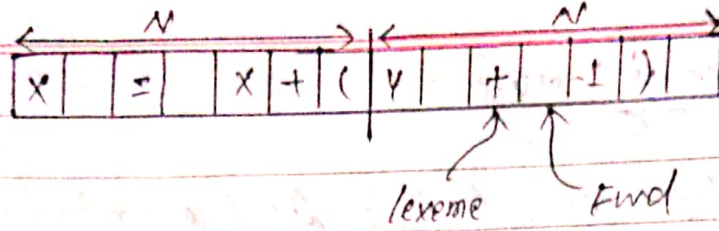
- Lexeme pointer points to the start of the lexeme while forward pointer scans the i/p buffer for lexeme.

- when forward pointer reaches the end of one half, second half is loaded and forward pointer points to the begining of the next half.
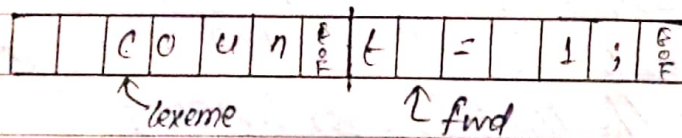
Pseudo code :

```
if (fwd at end of first half)
    reload second half;
    set fwd to point to begining of second half;
else if (fwd at end of second half)
    reload first half;
    set fwd to point to begining of first half;
else
    fwd++;
```

{fwd = forward pointer

| x | = | x | + | ( | y | + | · 1 | ) | |

lexeme   Fwd

• Sentinels

→ While using buffer pair technique, we have to check each time fwd is moved that it doesn't crosses the buffer half and when it reaches end of buffer, the other one needs to be loaded.

→ We can solve this problem by introducing a sentinel character at the end of both halves of buffer.

→ This sentinel can be any character which is not a part of source program. EOF is usually preferred as it will also indicate end of source program.

→ It signals the need for some special action (fill other buffer-half, or terminate processing).

| | c | o | u | n | EOF | t | | = | | 1 | ; | EOF |

lexeme     fwd

Pseudo code:

```
fwd++
if (*fwd == EOF)
{
      if (fwd at end of first half)
         ...
      else if (fwd at end of second half)
         ...
      else
            terminate processing
}
```

# Specification of Tokens

There are 3 specifications of tokens:
1) strings
2) Language
3) Regular expression

## strings and Language:

→ An alphabet $\Sigma$ is a finite set of symbols (characters)
   e.g. $\{0, 1\}$ is the binary alphabet

→ A string 's' is a finite sequence of symbols from $\Sigma$.
   - $|S|$ denotes the length of string s.
   - $\epsilon$ denotes the empty string, thus $|\epsilon| = 0$

→ A language is a specific set of strings over some fixed alphab $\Sigma$.
   - $\phi$ → the empty set language.
   - $\{\epsilon\}$ → language consisting of only empty string.

## Operations on strings

1. Prefix of s : A string obtained by removing zero or more trailing symbols of string s. e.g. ban is a prefix of banana.
2. suffix of s : A string formed by deleting zero or more of the leading symbols of s. e.g. nana is a suffix of banana.
3. substring of s : A string obtained by deleting a prefix and a suffix from s. e.g. nan is a substring of banana.
4. Proper prefix, suffix, or substring of s : Any non-empty string x that is a prefix, suffix, or substring of s that $s <> x$.

5. **subsequence of S** : Any string formed by deleting zero or more not necessarily contiguous symbols from s. e.g. baaa is a subsequence of banana.

## Operations on Languages

let L and M be two language then

1. union : $L \cup M = \{ s \mid s \text{ is in } L \text{ or } s \text{ is in } M \}$
2. Concatenation : $LM = \{ st \mid s \text{ is in } L \text{ and } t \text{ is in } M \}$
3. kleene closure of L : $L^* =$ "zero or more concatenation of" L.
4. Positive closure of L : $L^+ =$ "one or more concatenation of" L.

## Regular Expressions :

The regular expression over alphabet specifies a language according to the following rules :
1. $\epsilon$ is a regular expression that denotes $\{\epsilon\}$, i.e. the set containing the empty string.
2. $a \in \Sigma$ is a regular expression denoting $\{a\}$.
3. If r and s are regular expressions denoting languages $L(r)$ and $L(s)$ respectively, then
   a) $(r)|(s)$ is a R.E denoting the language $L(r) \cup L(s)$.
   b) $(r)(s)$ is a R.E. denoting the language $L(r)L(s)$.
   c) $(r)^*$ is a R.E denoting the language $(L(r))^*$.
   d) $(r)$ is a R.E denoting the language $L(r)$.

→ A language denoted by a regular expression is said to be a regular set.

## Properties of regular expression

For regular expression r, s & t

1. $r|s = s|r$     (| is commutative)
2. $r|(s|t) = (r|s)|t$     (| is associative)
3. $(rs)t = r(st)$     (Concatenation is associative)
4. $r(s|t) = rs|rt$     (Concatenation distributes over |)
5. $\epsilon r = r\epsilon = r$     ($\epsilon$ is the identity element for concatenation)
6. $r^* = (r|\epsilon)^*$     (Relation betⁿ * and $\epsilon$)
7. $r^{**} = r^*$     (* is idempotent)


## Regular Definition :

- If $\Sigma$ is an alphabet of basic symbols, then a regular definition is a sequence of definition of the form

$$d_1 \rightarrow r_1$$
$$d_2 \rightarrow r_2$$
$$\dots$$
$$d_n \rightarrow r_n$$

where, $d_i$ is a distinct name and $r_i$ is a regular expression over symbols in $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$

    ↗ Basic symbols     ↰ previously defined names

e.g. In c the RE for identifiers can be written using the regular definition as

$$letter \rightarrow A|B|\dots|z|a|b|\dots|z|\_$$
$$digit \rightarrow 0|1|2|\dots|9$$
$$identifier \rightarrow letter\,(letter\,|\,digit)^*$$

## Notational shorthands :

- This shorthand is used in certain constructs that occur frequently in regular expression.
- The following shorthands are often used:

$$r^+ = rr^*$$

(r?)   $r? = r | \epsilon$    ( zero or ~~more~~ one occurrences )

$$[a-z] = a|b|c| \cdots |z$$

e.g.

$$digit \rightarrow [0-9]$$
$$num \rightarrow digit^+ (.digit^+)? (E(+|-)? digit^+)?$$

## Recognition of tokens.

- A recognizer for a language is a program that takes a string x, and answers "yes" if x is a sentence of that language, and "no" otherwise.
- Recognition of token implies implementing a regular expression recognizer. That entails the implementation of finite automaton.
- ✓ The tokens that are specified using RE are recognized by using finite automata.
- Recognizer of tokens takes the language 'L' and the string's as input and try to verify whether $s \in L$ or not.
- There are two types of Finite Automata
  1. Deterministic Finite Automata (DFA)
  2. Non-deterministic Finite Automata (NFA)
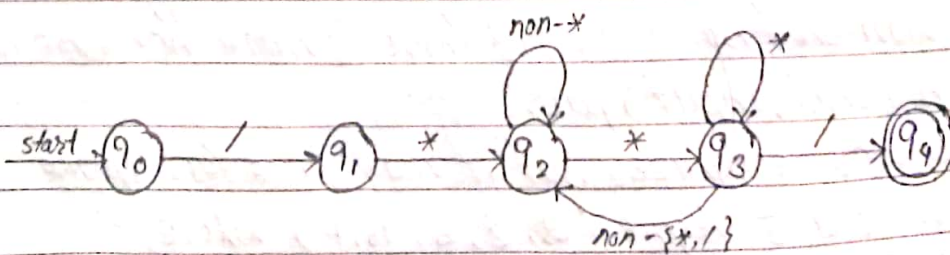
# Deterministic Finite Automaton (DFA)

→ FA is deterministic, if there is exactly one transition for each (state, input) pair.

→ It is faster recognizer but it ~~make~~ may take more space.

→ A DFA is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$ where,
- Q is a finite set of states.
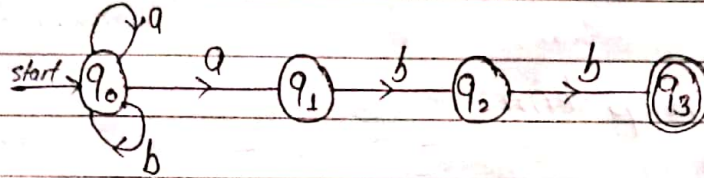- $\Sigma$ is a finite set of input alphabets
- $\delta$ is a transition function that maps $Q \times \Sigma \rightarrow Q$
- $q_0 \in Q$ is the initial state
- $F \subseteq Q$ is a set of final states.

E.g.
- DFA for R.E. $(a+b)^* abb$ :



- DFA to match c-style comments:



Note: Transition not showing from a state with any symbol is going to non-accepting (trapping state).

→ A state transition from one state to another on the path is called a move.

## ✳ Implementing DFA : Algorithm

- Algorithm to simulate a DFA(D), with start state $q_0$, that returns "yes" if the input string is accepted else return "no".

```
DFASim ( D, q0 )
{
    q = q0 ;
    c = getchar();
    while ( c != EOF )
    {
        q = move (q, c);      // transition function
        c = getchar();
    }
    if ( q is in F )
        return 'yes';
    else
        return 'No';
}
```

## ✳ Non-Deterministic Finite Automaton (NFA)

→ FA is non-deterministic, if there is more than one transition for each (state, input) pair.

→ It is slower recognizer but it may take less space.

→ An NFA is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$ where,
- Q is a finite set of states
- $\Sigma$ is a finite set of input alphabet
- $\delta$ is a transition function that maps $Q \times \Sigma \rightarrow 2^Q$
- $q_0 \in Q$ is the start state.
- $F \subseteq Q$ is the set of final states.

**E.g.**

- NFA for R.E. $(a+b)^*abb$.



→ **E-NFA :**

→ E-transitions are allowed in NFAs.

→ In other words, we can move from one state to another one without consuming any symbol.
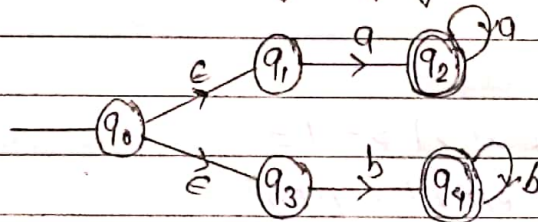
**E.g.**



fig: state machine with e-moves that is equivalent to the regular expression $aa^* + bb^*$.

✓ **Implementing NFA :**

```
q = E-closure ({q₀});
c = getchar();
while ( c != eof )
{
    q = E-closure (move (q,c));
    c = getchar();
}
if (q∩F ≠ ∅) then
        return "Yes";
else
        return "No";
```

# Grammar

stmt → if expr then stmt
| if expr then stmt else stmt
| ∈

expr → term relop term
| term

term → id
| num

Regular definitions for above grammar:

if → if
then → then
else → else
relop → < | <= | <> | > | >= | =
id → letter (letter | digit)*
num → $digit^+ (.digit^+)? (E(+|-)? digit^+)?$

Transition diagram:

relop → < | <= | <> | > | >= | =

start →(0) —<→ (1) —=→ ②  return (relop, LE)
(1) —>→ ③  return (relop, NE)
(1) —other→ ④*  return (relop, LT)
(0) —=→ ⑤  return (relop, EQ)
(0) —>→ (6) —=→ ⑦  return (relop, GE)
(6) —other→ ⑧*  return (relop, GT)

$id \rightarrow letter \ (\ letter \ / \ digit \ )^{*}$

letter or digit

start $\rightarrow$ (9) $\xrightarrow{letter}$ (10) $\xrightarrow{other}$ (11)$^{*}$ return (gettoken(),
install-id())

$id \rightarrow letter \ (\ letter \ / \ digit \ )^{*}$
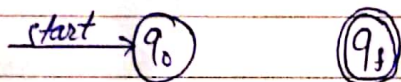
# R.E to NFA / Thomson's Construction

→ It guarantees that the resulting NFA will have exactly one final state, and one start state.
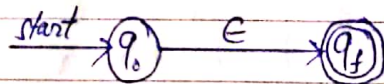
Input : RE, r, over alphabet $\Sigma$
Output : $\epsilon$-NFA accepting $L(r)$

The method consists of following steps:

i) For $\phi$, we construct the NFA as

start $\longrightarrow$ $(q_0)$      $(\!(q_f)\!)$

where $q_0$, is an initial state & $q_f$ is a final state.

ii) For $\epsilon$, we construct the NFA as

start $\longrightarrow$ $(q_0)$ $\xrightarrow{\;\epsilon\;}$ $(\!(q_f)\!)$

iii) For every $a \in \Sigma$ we construct the NFA as

start $\longrightarrow$ $(q_0)$ $\xrightarrow{\;a\;}$ $(\!(q_f)\!)$
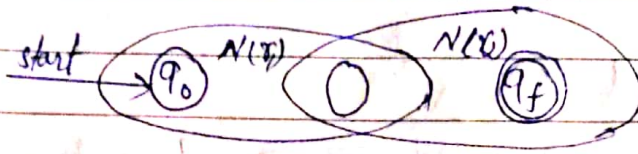
iv) If $N(r_1)$ and $N(r_2)$ are NFAs for R.E $r_1$ and $r_2$
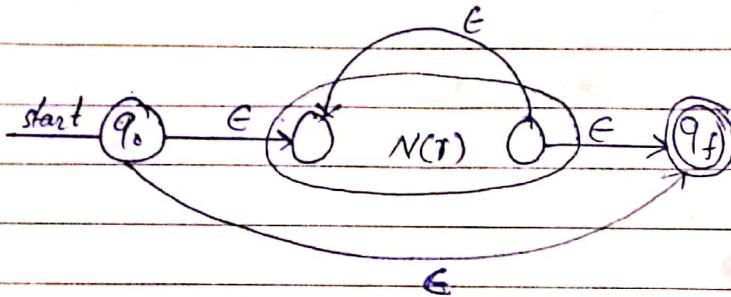
a) For R.E. $r_1 | r_2$ (i.e. $r_1 + r_2$) we construct the NFA as

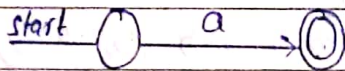b) For the R.E $r_1 r_2$, we construct the NFA as

start $\rightarrow$ $q_0$ $N(r_1)$ $\bigcirc$ $N(r_2)$ $q_f$

v) For the R.E. $r^*$, we construct the NFA as

$\epsilon$

start $q_0$ $\epsilon$ $\bigcirc$ $N(r)$ $\bigcirc$ $\epsilon$ $q_f$

$\epsilon$
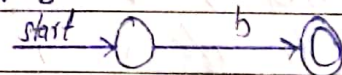
Examples

① $(a+b)^* a$ to NFA.

For a                                    For b

start $\bigcirc$ $\xrightarrow{a}$ $\bigcirc$          start $\bigcirc$ $\xrightarrow{b}$ $\bigcirc$

For $(a+b)$

start $\bigcirc$ $\xrightarrow{\epsilon}$ $\bigcirc$ $\xrightarrow{a}$ $\bigcirc$ $\xrightarrow{\epsilon}$ $\bigcirc$

$\xrightarrow{\epsilon}$ $\bigcirc$ $\xrightarrow{b}$ $\bigcirc$ $\xrightarrow{\epsilon}$

for $(a+b)^*$

$\epsilon$

start $\bigcirc$ $\xrightarrow{\epsilon}$ $\bigcirc$ $\xrightarrow{\epsilon}$ $\bigcirc$ $\xrightarrow{a}$ $\bigcirc$ $\xrightarrow{\epsilon}$ $\bigcirc$ $\xrightarrow{\epsilon}$ $\bigcirc$

$\xrightarrow{\epsilon}$ $\bigcirc$ $\xrightarrow{b}$ $\bigcirc$ $\xrightarrow{\epsilon}$

$\epsilon$
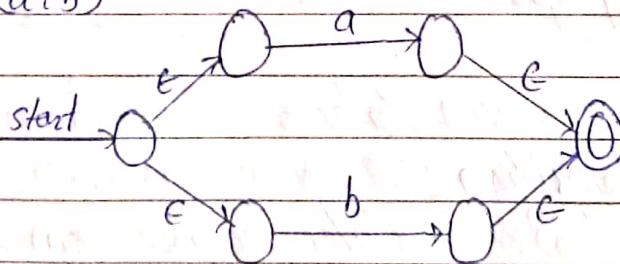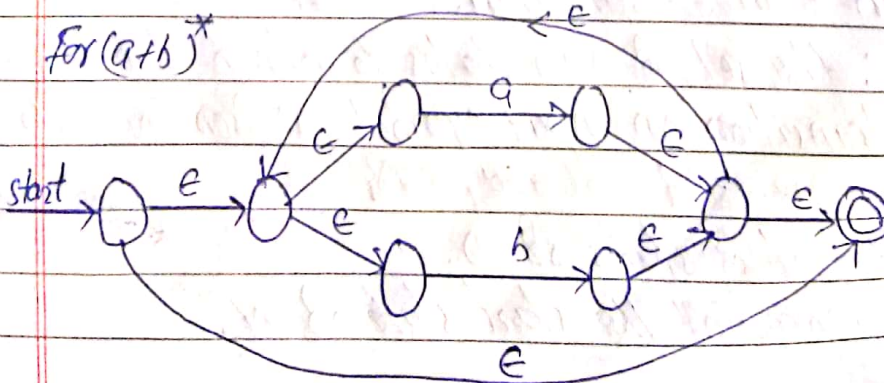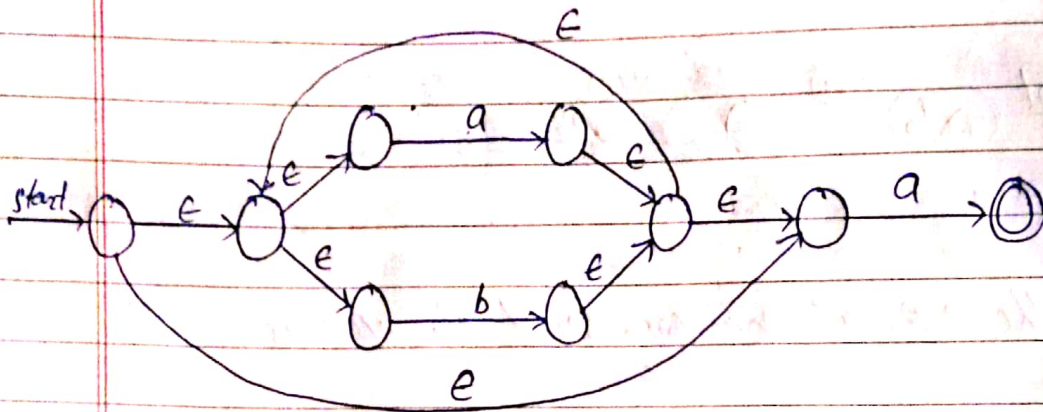
For $(a+b)^*a$



See Bnk for other Eg.

# NFA to DFA  (subset construction)

The subset construction algorithm converts an NFA into a DFA.

→ In this algorithm we use the symbol N to represent an NFA and D for representing DFA.

→ This algorithm constructs a transition table Dtran for D.

We use the following operations: (s represent an NFA state and T a set of NFA states)

- $\epsilon$-closure(s) : the set of NFA states rechable from NFA state s on $\epsilon$-transition. i.e. $\epsilon$-closure(s) = $\{s\} \cup \{t / s \xrightarrow{\epsilon} \ldots \xrightarrow{\epsilon} \ldots \xrightarrow{\epsilon} t\}$

- $\epsilon$-closure(T) : the set of NFA states rechable from NFA state s in T on $\epsilon$-transition. i.e. $\bigcup_{s \in T} \epsilon$-closure(s).

- move (T, a) : the set of NFA states to which there is a transition on input symbol 'a' from NFA states s in T. ie. $\{t / s \xrightarrow{a} t$ and $s \in T\}$

→ Dstates is the set of states of D.

→ we use so to represent the start state of N.

## Computation of e-closure:

```
Push all states in T onto stack;
Initialize e-closure (T) to T;
while stack is not empty do begin
Pop t, the top element of stack;
for each state u with an edge from t to u labeled e do
      if u is not in e-closure (T) do begin
      add u to e-closure(T); push u onto stack;
      end
end
```
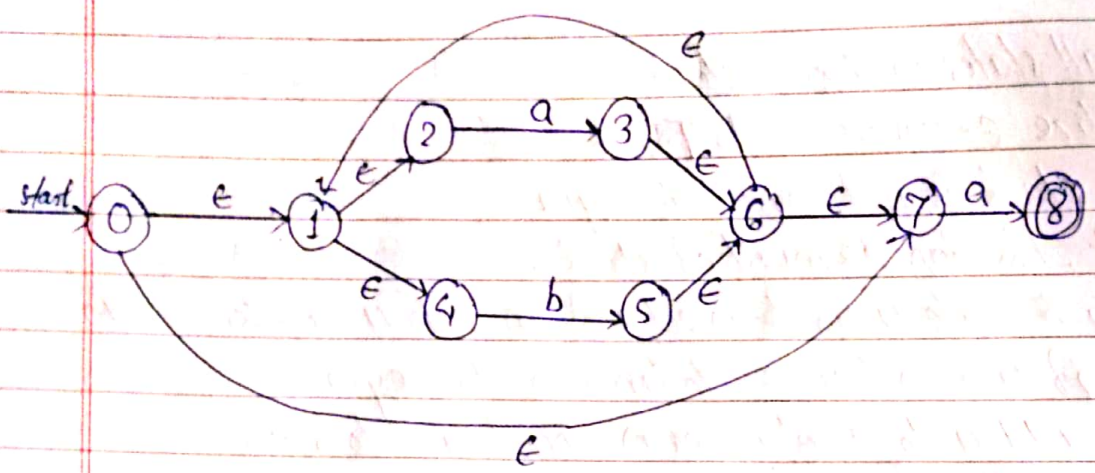
### ↙ Subset Construction Algorithm:

```
Put e-closure (S₀) as an unmarked states into Dstates.
    While there is an unmarked state T in Dstates do.
      mark T;
      for each input symbol a ∈ Σ do
          U = e-closure (move (T,a))
          if U is not in Dstates then
                add U as an unmarked state to Dstates
          end if
              Dtran[T,a] = U
      end do
    end do
```

→ The start state of DFA is e-closure (S₀)

## Example



**Sol^n**

The initial state of the NFA is 0.

Therefore, the initial state of the DFA is,

$$A = \epsilon\text{-closure}(\{0\}) = \{0,1,2,4,7\}$$

Here

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

Now,

$$Dtran[A,a] = \epsilon\text{-closure}(move(A,a))$$
$$= \epsilon\text{-closure}(\{3,8\})$$
$$= \{1,2,3,4,6,7,8\}$$
$$= B \text{ (say)}$$

// Dtran is a transition table for DFA.

$$Dtran[A,b] = \epsilon\text{-closure}(move(A,b))$$
$$= \epsilon\text{-closure}(\{5\})$$
$$= \{1,2,4,5,6,7\}$$
$$= C \text{ (say)}$$

$$Dtran[B,a] = \epsilon\text{-closure}(move(B,a))$$
$$= \epsilon\text{-closure}(\{3,8\})$$
$$= \{1,2,3,4,6,7,8\}$$
$$= B$$

Dtran [B,b] = ε-closure (move (B,b))

$\quad\quad$ = ε-closure ({5})

$\quad\quad$ = {1, 2, 4, 5, 6, 7}

$\quad\quad$ = C

| | a | b |
|---|---|---|
| { } | | |
| { } | | |
| { } | | |

Dtran [C,a] = ε-closure (move (C,a))

$\quad\quad$ = ε-closure ({3,8})

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

$\quad\quad$ = B

Dtran [C, b] = ε-closure (move (C,b))

$\quad\quad$ = ε-closure ({5})

$\quad\quad$ = {1, 2, 3, 4, 5, 6, 7}

$\quad\quad$ = C

Now the equivalent DFA is



$(Q, \Sigma, \delta, q_0, F) \rightarrow$ ε-NFA

$(Q', \Sigma, \delta^*, q_0', F') \rightarrow$ DFA $\quad\quad \rightarrow$ ε-NFA to DFA

$\delta^*$ is defined as:

$\quad \delta^*(q,a) = $ ε-closure $(\delta(q,a))$

# $\varepsilon$-NFA to NFA

$(Q, \Sigma, \delta, q_0, F) \rightarrow \varepsilon\text{-NFA}$

$(Q, \Sigma, \delta^*, q_0, F) \rightarrow \text{NFA}$

$\delta^*(q, a) = \varepsilon\text{-closure}(\delta(\varepsilon\text{-closure}(q), a))$

$$\begin{cases} \text{Initial state} \rightarrow \text{Same as of NFA} \\ \delta^*(q_0, a) = \{q_3, q_4\} \end{cases}$$

## ✗ Conversion of R.E Directly into DFA

- Syntax tree based reduction to DFA / using followposition base reduction

**Important states:**
    The state $\delta$ in $\varepsilon$-NFA is an important state if it has no null transition.

**Augmented R.E:**
    $\varepsilon$-NFA created from RE has exactly one accepting state and accepting state is not important state since there is no transition so by adding special symbol # on the RE at the rightmost position, we can make the accepting state as an important state that has transition on #.
    The RE (r)# is called the augmented regular expression of the regular expression r.
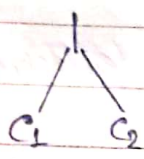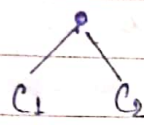
## Procedure:

1. Augment the given regular expression by concatenating it with special symbol # i.e. $r \rightarrow (r)\#$

2. Create the syntax tree of this augmented regular expression. In this tree, all operators will be inner nodes and all the alphabet, symbols including # will be leaves.

3. Numbered each leaves.

4. Traverse the tree to construct nullable, firstpos, lastpos and followpos.

5. Finally construct the DFA from the followpos.

To evaluate followpos, we need three functions to defined the nodes of the syntax tree.

- firstpos (n): The set of the position of the first symbol of strings generated by the sub-expression rooted by n.

- lastpos (n): The set of the position of the last symbol of strings generated by the sub-expression rooted by n.

- nullable(n): true if the empty string is a member of strings generated by the sub-expression rooted by n, false otherwise.
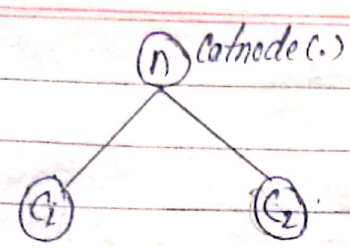
### Rules for creating nullable, firstpos & lastpos

| node n | nullable(n) | firstpos(n) | lastpos(n) |
|---|---|---|---|
| 1. leaf labeled ε | true | $\phi$ | $\phi$ |
| 2. leaf labeled with position i | False | $\{i\}$ | $\{i\}$ |

| n | nullable(n) | firstpos(n) | lastpos(n) |
|---|---|---|---|



**3.** (node with children $c_1$, $c_2$ — or node)

$nullable(c_1)$ OR $nullable(c_2)$ — firstpos($c_1$) $\cup$ firstpos($c_2$) — lastpos($c_1$) $\cup$ lastpos($c_2$)



**4.** (cat node with children $c_1$, $c_2$)

$nullable(c_1)$ AND $nullable(c_2)$

if ($nullable(c_1)$ = TRUE) firstpos($c_1$) $\cup$ firstpos($c_2$) else firstpos($c_1$)

if ($nullable(c_2)$ = TRUE) lastpos($c_1$) $\cup$ lastpos($c_2$) else lastpos($c_2$)

**5.** (star node with child $c_1$)
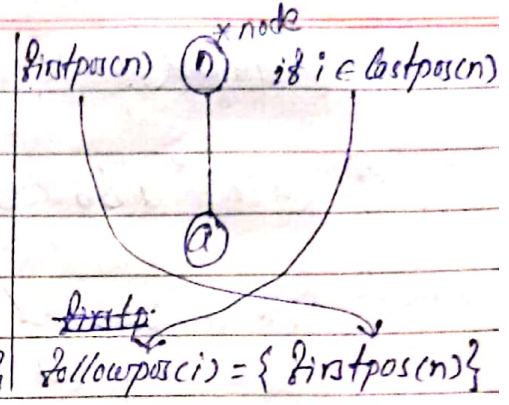
true — firstpos($c_1$) — lastpos($c_1$)

## Computation of followpos (Algorithm)

```
for each node n in the tree do
    if n is a cat-node with left child c₁ and right child c₂ then
        for each i in lastpos (c₁) do
            followpos (i) = followpos (i) ∪ firstpos (c₂)
        end do
    else if n is a star-node
        for each i in lastpos (n) do
            followpos (i) = followpos (i) ∪ firstpos (n)
        end do
    end if
end do
```

Catnode (.)

if $i \in lastpos(c_1)$
$$\longrightarrow followpos(i) = \{ firstpos(c_2) \}$$

firstpos(n) if $i \in lastpos(n)$

firstp.
$$followpos(i) = \{ firstpos(n) \}$$

## Algorithm to create DFA from RE

1. Create a syntax tree of $(r)$ #
2. Calculate the functions: nullable, firstpos, lastpos & followpos
3. Put firstpos (root) into the state of DFA as an unmarked state.
4. while ( there is unmarked state s in the states of DFA ) do
   - Mark s
   - For each input symbol $a \in \Sigma$ do
     - let $s_1, \ldots, s_n$ are positions in s and symbols in those position is a.
     - $s' \leftarrow followpos(s_1) \cup \ldots \cup followpos(s_n)$
     - move $(s, a) \leftarrow s'$
     - if ( s' is not empty and not in the states of DFA)
       put s' into the states of DFA as an unmarked state.
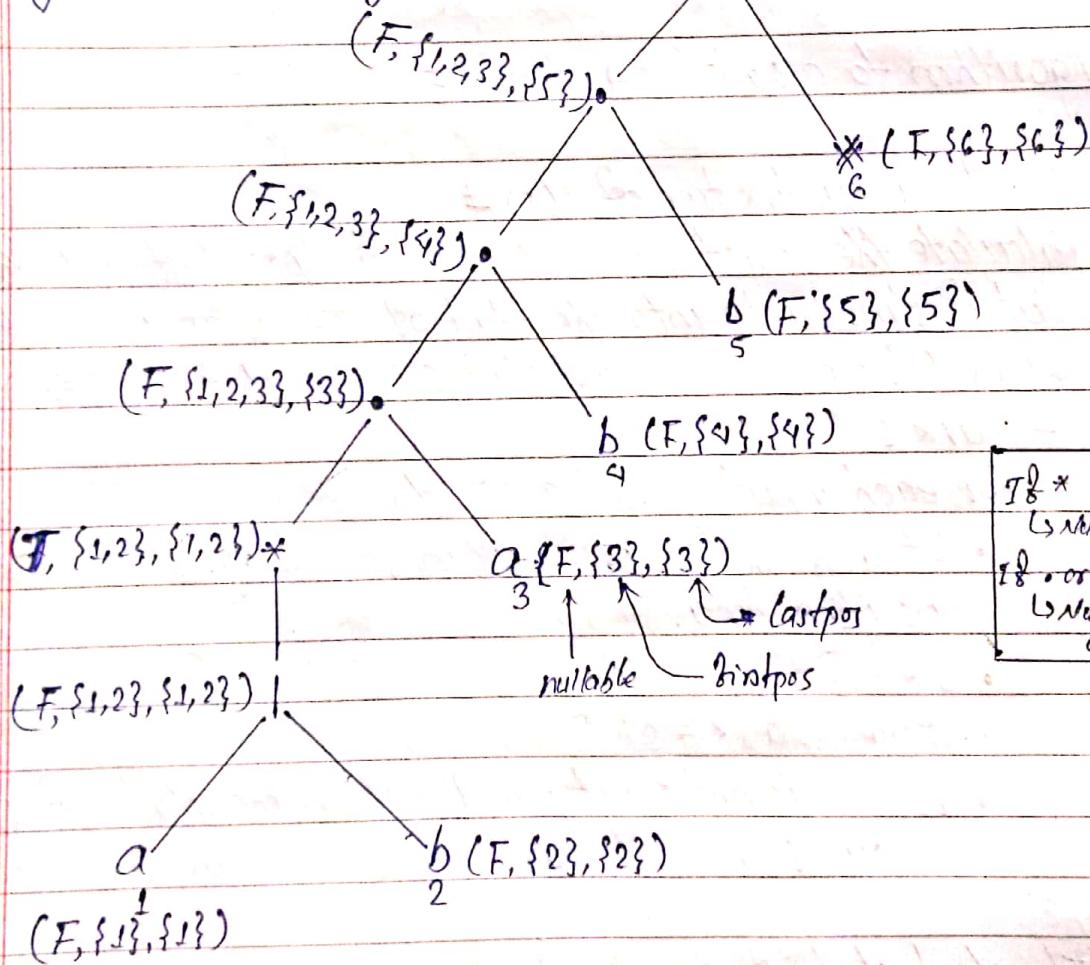
→ The start state of DFA is firstpos (root)
→ The accepting states of DFA are all states containing the position of #.

# Examples:

① $(a|b)^* abb$

The augmented regular expression of given regular expression is

$$(a|b)^* abb \#$$
$\quad\quad 1\ 2 \quad 3\ 4\ 5\ 6$

Syntax tree for augmented R.E is

$(F, \{1,2,3\}, \{6\})$

$(F, \{1,2,3\}, \{5\})$

$\ast \ (T, \{6\}, \{6\})$
6

$(F, \{1,2,3\}, \{4\})$

$b \ (F, \{5\}, \{5\})$
5

$(F, \{1,2,3\}, \{3\})$

$b \ (F, \{4\}, \{4\})$
4

$(T, \{1,2\}, \{1,2\}) \ast$

$a \ (F, \{3\}, \{3\})$
3
↑ nullable ↑ firstpos → lastpos

$(F, \{1,2\}, \{1,2\}) \ .$

$a$

$b \ (F, \{2\}, \{2\})$
2

$(F, \{1\}, \{1\})$

$T \ \ \ast$
└→ nullable is true

$T \ . \ or \ |$
└→ Nullable depends on its child

## Calculating followpos:

**for (\*) Node:**

followpos (1) = {1,2}
followpos (2) = {1,2}

**for (.) Node:**

followpos (1) = {3}
followpos (2) = {3}
followpos (3) = {4}
followpos {4} = {5}
followpos {5} = {6}

Finally,
followpos (1) = {1, 2, 3}
followpos (2) = {1, 2, 3}
followpos (3) = {4}
followpos (4) = {5}
followpos (5) = {6}.
followpos (6) = { }

Now,
1) start state of DFA = firstpos (root) = {1, 2, 3} = $S_0$

use followpos of symbol representing position in R.E to obtain the next states of DFA.

2) Here 1 and 3 represent 'a'
        2 represent 'b'
    followpos ({1,3}) = {1, 2, 3, 4} = $S_1$
            $\delta (S_0, a) = S_1$
    followpos (2) = {1, 2, 3} = $S_0$
            $\delta (S_0, b) = S_0$

3) From $S_1$ = {1, 2, 3, 4}
        1, 3 → 'a'
        2, 4 → 'b'
    followpos (1, 3) = {1, 2, 3, 4} = $S_1$
            $\delta (S_1, a) = S_1$
    followpos (2, 4) = {1, 2, 3, 5} = $S_2$
            $\delta (S_1, b) = S_2$

4) From $S_2$ = {1, 2, 3, 5}
        1, 3 → 'a'
        2, 5 → 'b'

followpos $(1,3) = \{1,2,3,4\} = S_1$
$\qquad \delta(S_2,a) = S_1$
followpos $(2,5) = \{1,2,3,6\} = S_3$
$\qquad \delta(S_2,b) = S_3$

5) From $S_3 = \{1,2,3,6\}$
$\qquad 1,3 \rightarrow 'a'$
$\qquad 2 \rightarrow 'b'$
$\qquad 6 \rightarrow '\#'$
$\quad$ followpos $(1,3) = \{1,2,3,4\} = S_1$
$\qquad\qquad \delta(S_3,a) = S_1$
$\quad$ followpos $(2) = \{1,2,3\} = S_0$
$\qquad\qquad \delta(S_3,b) = S_0$

6) Final state $= \{S_3\}$

Draw the DFA is:



② $(a+e)bc^*$

The augmented R.E is
$$(a+e)bc^*\#$$
$\qquad\quad 1 \quad\; 2\; 3\; 4$

Syntax tree for augmented R.E is



$(F, \{1, 2\}, \{4\})$

$(F, \{1, 2\}, \{2, 3\})$

$*(F, \{4\}, \{4\})$

$(F, \{1, 2\}, \{2\})$

$*(T, \{3\}, \{3\})$

$(T, \{1\}, \{1\})$   $t$

$c$ $(F, \{3\}, \{3\})$
3

$b$ $(F, \{2\}, \{2\})$
2

$a$
1

$\in (T, \phi, \phi)$
2

$(F, \{1\}, \{1\})$

Calculating followpos :

followpos (1) = { 2 }
followpos (2) = { 3, 4 }
followpos (3) = { 3, 4 }
followpos (4) = { }

Now,
   start state of DFA = firstpos (root) = {1, 2} = S₀

Use followpos of symbol representing position in R.E to obtain
the next state of DFA.

Here 1 represents 'a'
      2 represents 'b'

followpos $(1) = \{2\} = S_1$

$\qquad \delta(S_0, a) = S_1$

followpos $(2) = \{3, 4\} = S_2$

$\qquad \delta(S_0, b) = S_2$

From $S_1 = \{2\}$

$\qquad$ Here 2 represents 'b'

$\qquad$ followpos $(2) = \{3, 4\} = S_2$

$\qquad\qquad \delta(S_1, b) = S_2$

From $S_2 = \{3, 4\}$

$\qquad$ Here 3 represents 'c'

$\qquad\qquad$ 4 represents '#'

$\qquad$ followpos $(3) = \{3, 4\} = S_2$

$\qquad\qquad \delta(S_2, c) = S_2$

Accepting state $= \{S_2\}$

Now the DFA is
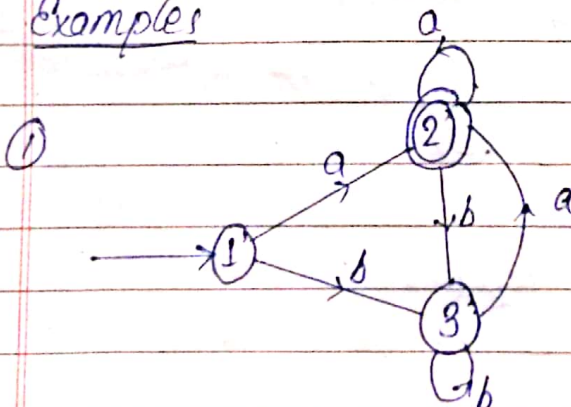
## State Minimization in DFA

→ DFA minimization refers to the task of transforming a given DFA into an equivalent DFA which has minimum number of states.

Procedure:

1. Partition the set of states into two groups: Set of accepting states and set of non-accepting states.

2. For each new group G
   - Partition G into subgroups such that states $s_1$ and $s_2$ are in the same group iff for all input symbol a, states $s_1$ and $s_2$ have transition to states in the same group.

3. Process until all the partition contains equivalent states only or have single state.

- start state of the minimized DFA is the group containing the start state of the original DFA.
- Accepting states of the minimized DFA are the groups containing the accepting states of the original DFA.

## Examples

1.



Partition the set of states as
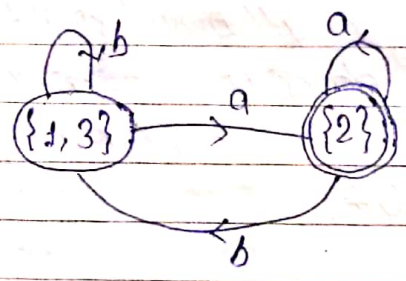$$G_1 = \{2\}$$
$$G_2 = \{1, 3\}$$

For G₂:

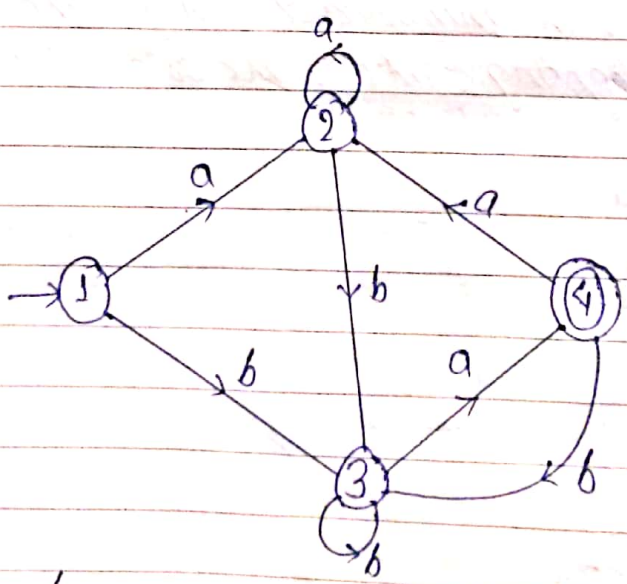| | 1 | 3 |
|---|---|---|
| a | G₁ | G₁ |
| b | G₂ | G₂ |

G₂ Cannot divided further.

Here,

Equivalent states are 1 & 3.

So, the minimized ~~stat~~ DFA is



②



Accepting & non-accepting states are grouped as

G₁ = {4}

G₂ = {1, 2, 3}

For $G_2$:

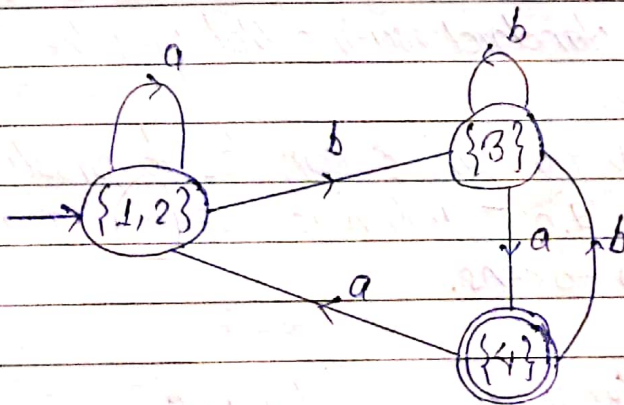| | 1 | 2 | 3 |
|---|---|---|---|
| a | $G_2$ | $G_2$ | $G_1$ |
| b | $G_2$ | $G_2$ | $G_2$ |

$G_2$ is further divided

$$G_2 = \{1, 2\}$$
$$G_3 = \{3\}$$

Also. Here, no more partitioning for $G_2$.

so the minimized DFA



## Space Time Tradeoffs : NFA vs DFA

- Given the RE $r$ and the input string $s$ to determine whether $s$ is is in $L(r)$ we can either construct NFA and test or we can construct DFA and test for $s$ after NFA is constructed from $r$.
- $\epsilon$-NFA (for NFA only constant time differs)
  - space complexity : $O(|r|)$
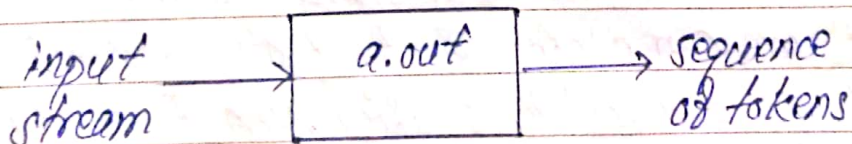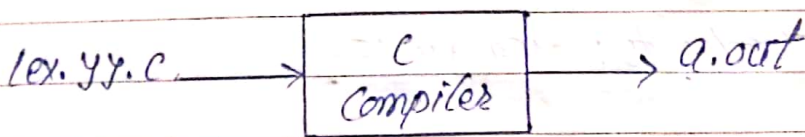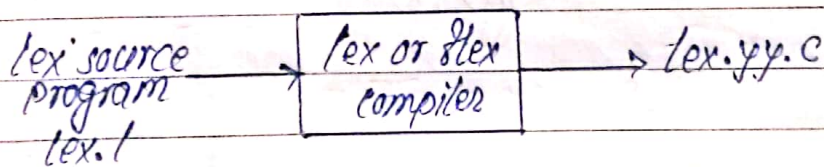  - Time complexity : $O(|r| * |s|)$
- DFA
  - space complexity : $O(2^{|r|})$ [$\epsilon$-NFA Construction & then subset const]
  - Time complexity : $O(|s|)$
- If we can create DFA from RE by avoiding transition table, then we can improve the performance

# Creating a lexical Analyzer Generator (lex)/Flex

- Systematically translate regular definitions into C source code for efficient scanning.
- Flex is a latter version of lex.
- Generated code from lex is easy to integrate in C applications by integrating.
- Firstly, specification of a lexical analyzer is prepared by creating a program lex.l in lex/Flex.
- lex.l is run into lex compiler to produce a c program lex.yy.c.
- lex.yy.c consists of the tabular representation of lex.l, together with a standard routine that uses the table to recognize lexeme.
- lex.yy.c is run through c compiler to produce the object program a.out which is a lexical analyzer that input into tokens.

```
lex source          lex or flex
program      ───→   compiler     ───→   lex.yy.c
lex.l


lex.yy.c ──────→       C         ───→   a.out
                    compiler


input    ──────→     a.out       ───→   sequence
stream                                  of tokens
```

→ A lex/flex specification consists of three parts:
1. regular definitions, C declarations in %{ %}
   %%

2. translation rules

%%

3. User-defined auxiliary procedures

The translation rules are of the form:

Pattern →
$P_1$        {action$_1$}
$P_2$        {action$_2$}
:              :
$P_n$        {action$_n$}

→ lex/Flex regular definitions are of the form:
    name definition

e.g.

Digit      [0-9]
letter     [A-Za-z]
ID         ({letter}{letter}|{digit})* or [a-z][a-z0-9]*

→ Action in lex are of the form
    Pattern  action
    where pattern must be unintended and action must be
on the same line.

Pattern                                    Action
if|then|else|for|while|do       {printf("A keyword:%s|n", yytext);}

Global function & variables

- yylex() : is the scanner function that can be invoked by the parser

- yytext : extern char* yytext; is a global char pointer holding the currently matched lexeme.

- yyleng : extern int yyleng; is a global int that contains the length of the currently matched lexeme.

Date _____
Page _____

## lex Examples

① %{
    #include <stdio.h>
    %}

Translation rules →

%%
    [0-9]+    {printf ("%s \n", yytext); }        Contains the matching lexeme
    . | \n    {  }
%%
    main()                                        Invokes the lexical analyzer
    { yylex();
    }

_____

② %{
    #include <stdio.h>
    #include <math.h>
    %}
    DIGIT    [0-9]
    ID       [a-zA-Z_][a-zA-Z0-9_]*
    OP       "+" | "-" | "*" | "/"
    %%
    {DIGIT}+    {Printf("An integer %s (%d)", yytext, atoi (yytext)); }        ascii to integer
    {DIGIT}+"."{DIGIT}+    {printf ("A float no: %s (%f)", yytext, atof (yytext); }
    {ID}        {Printf ("An identifier : %s", yytext ); }
    [\t\n]+    {  }
    .           {printf ("An unrecognized char:%s", yytext); }
    %%.
    int main (int argc, char **argv)
    {
        ++argv, --argc;

```
if (argc > 0)
    yyin = fopen (argv [1], "r");
else
    yyin = stdin;
yylex();
}
```

```
/* lex program that recognizes the identifier of c languages */
%{
#include<stdio.h>
%}
%%
^[a-z A-Z _][a-z A-Z 0-9]*    printf ("Valid Identifier");
^[^a-z A-Z _] printf ("Invalid Identifier");
.;
%%
main()
{
    yylex();
}
```

Jayanta Poudel