# UNIT-I INTRODUCTION TO COMPILERS & LEXICAL ANALYSIS

Introduction - Translators - Compilation and Interpretation - Language processors - The phases of Compiler - Lexical Analysis - Role of Lexical Analyzer - Input Buffering - Specification of Tokens - Recognition of Tokens - Finite Automata - Regular Expressions to Automata NFA, DFA - Minimizing DFA - Language for specifying Lexical Analyzers Lex tool.
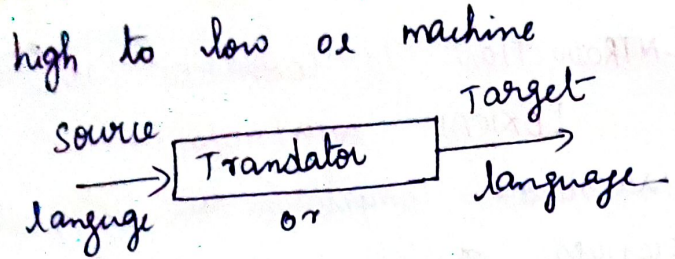
## Introduction:

Compilers are basically translators. how the source program is compiled with the help of various phases of compiler.
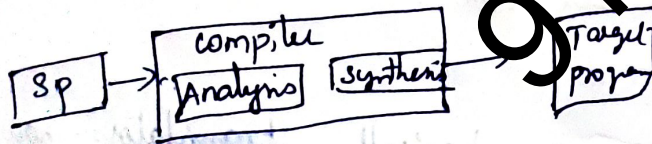
## Translators:-

Translators is one kind of program that take some code as input and converts into another form.

low level language or assembly language or high level language

high to low or machine

source $\longrightarrow$ [ Translator ] $\longrightarrow$ Target
language        or        language

Translators ( Compilers and assemblers
high to machine        assembly to mac

Analysis and Synthesis Model:
$\downarrow$        $\hookrightarrow$ intermediate form, converted
is read & broken into Constituent    into equivalent
pieces, intermediate code is created    target
                                        program

[ Sp ] $\rightarrow$ [ Compiler [ Analysis ] [ Synthesis ] ] $\rightarrow$ [ Target prog ]

Execution of program:
    $\rightarrow$ Only compiler program is not sufficient
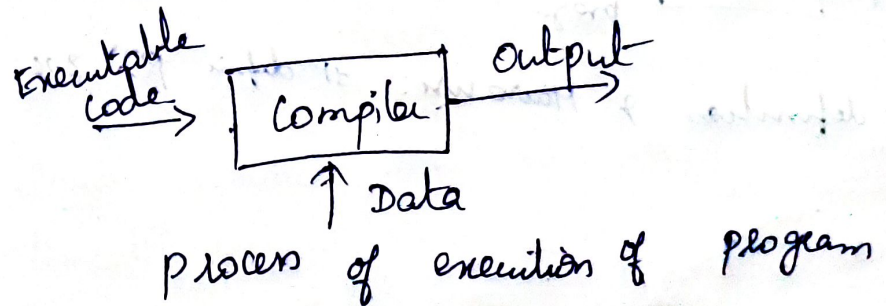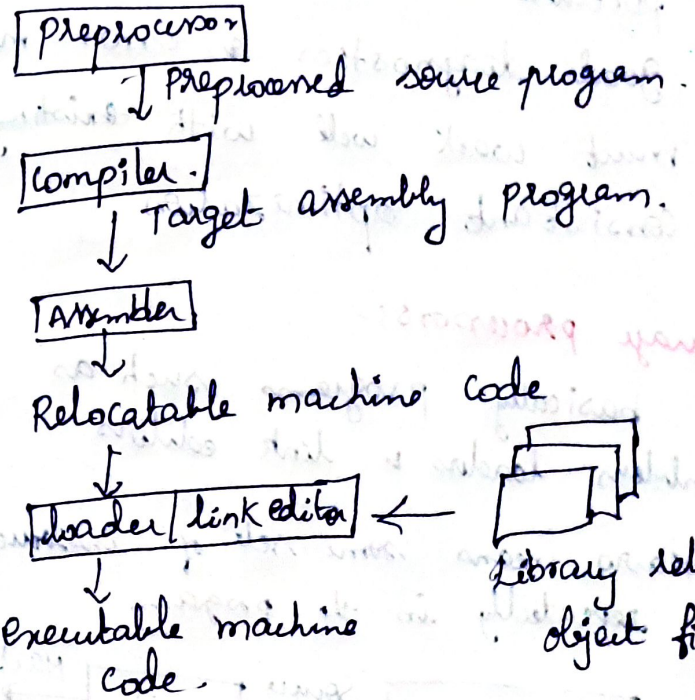compiler source to highlevel language, target-
assembly code as input, relocatable machine
code as output.

    The task of loader, relocation of object code
allocation of load time address which exists in
memory & placement of load time address,

and data in memory at proper location
link editor link several files of
relocate object modules to resolve the mutual
reference. These files may be library files.
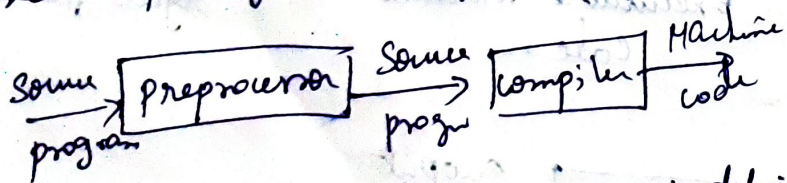
        Skeleton source program.
                $\downarrow$
        [ Preprocessor ]
                $\downarrow$ Preprocessed source program.
        [ Compiler ]
                $\downarrow$ Target assembly program.
        [ Assembler ]
                $\downarrow$
        Relocatable machine code
                $\downarrow$
        [ loader / link editor ] $\leftarrow$  Library relocatable
                $\downarrow$                    object file
        Executable machine
        code.

Executable
code $\rightarrow$ [ Compiler ] $\rightarrow$ Output
                $\uparrow$ Data
process of execution of program

properties of Compilers:

1) bug-free.
2) Correct machine Code
3) generate machine code run fast.
4) Compilation time is ∝ program size.
5) portable
6) good diagnostics & error message
7) must work well with existing debuggers.
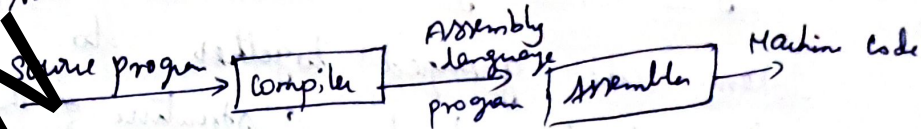8) Consistant optimization

Language processors:-

basically programs such as preprocessors assemblers loaders & link editors

* Macro means some set of instructions, used repeatedly in the program.

Source → [Preprocessor] → Source prog → [Compiler] → Machine code

Macro definitions & Macro use. # define PI 3.14.

Assemblers

Source program → [Compiler] → Assembly language program → [Assembler] → Machine Code

MOV a, R1
MUL #5, R,
ADD #, R1
MOV R1, b;

Binary language, Machine code.
relocatable machine Cod.
two passes → one pass input program
end of second pass is relocatable machine Cod.

Loaders & link editors:-

relocatable machine Code is read and the relocatable addresses are altered.

The phases of Compiler:-

1) Lexical Analysis (Scanning).
Complete source code is scanned and broken up into group of strings called tokens.

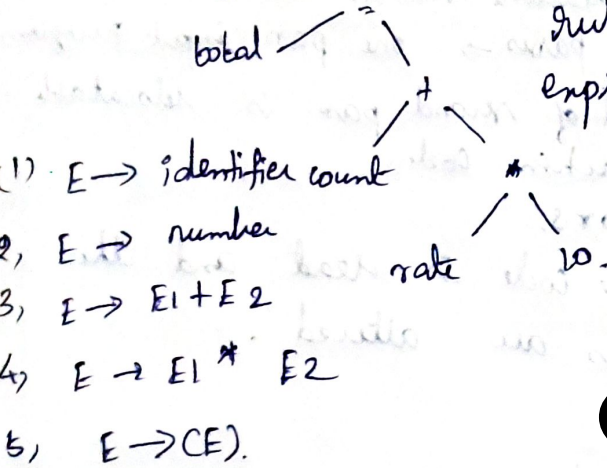A token is a sequence of characters having a Collective meaning.

total = count + rate * 10.

$id_1 = id_2 + id_3 * $ Constant no.

blank characters are eliminated during lexical analysis.

# Syntax analysis:- parsing.

tokens are grouped together to form a hierarchical structure ( structure of the source string). Called parse tree or syntax tree.

rules are usually expressed by context free grammar

total

$$E$$

(1) $E \rightarrow$ identifier count

2, $E \rightarrow$ number

3, $E \rightarrow E_1 + E_2$      rate    10.

4, $E \rightarrow E_1 * E_2$

5, $E \rightarrow (E)$.

## Semantic Analysis:- determines the meaning of the source string ( matching of parenthesis or if else statements or performing arithmetic expressions. After these phases, intermediate code as generated.

## Intermediate Code Generation:- Code can be easily converted to target code. these address code, quadruple, triple, posix

$t_1 := $ int to float (10)      $t_3 := $ count + $t_2$

$t_2 := $ rate $\times t_1$      total $:= t_3$

order of operations devised by three address code

$t := $ int_to_float(10)

total $:= 63$

## 5 Code optimization.

faster executing code or less consumption of memory, optimizing the code, overall running time of target program can be improved.

## 6, Code generation

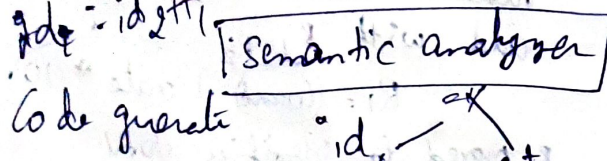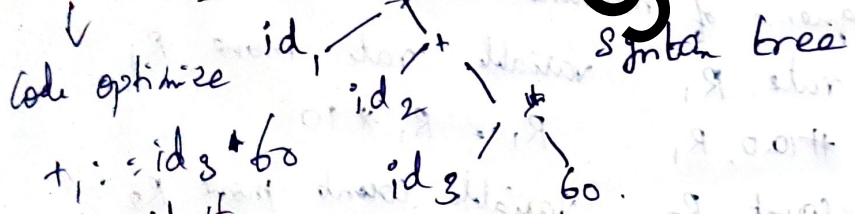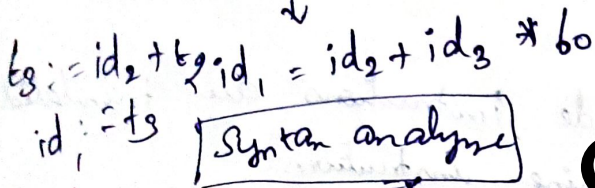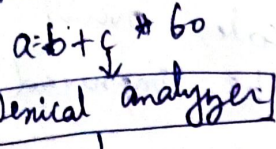intermediate code instructions are translated into sequence of machine instructions.

```
MOV rate, R₁        variable rate move R₁
MUL #10.0, R₁       R₁ = R₁ * 10
MOV count, R₂       variable count move R₂
ADD R₂, R₁          add with R₁ R₂
MOV R₁, total
```

$R_1 = $ (count) + (rate * 10)

$R_1$ moved to identifier total.

## Symbol table Management:-

store identifiers (variables) used in the program.

Stores information about attributes of each identifier.

It is a data structure used to store the information about identifiers.

store & retrieve data from that record efficiently

## Error detection & handling

errors are reported to error handlers, compilation can proceed, syntax analysis phase, syntax error.

$a := b + c * 60$
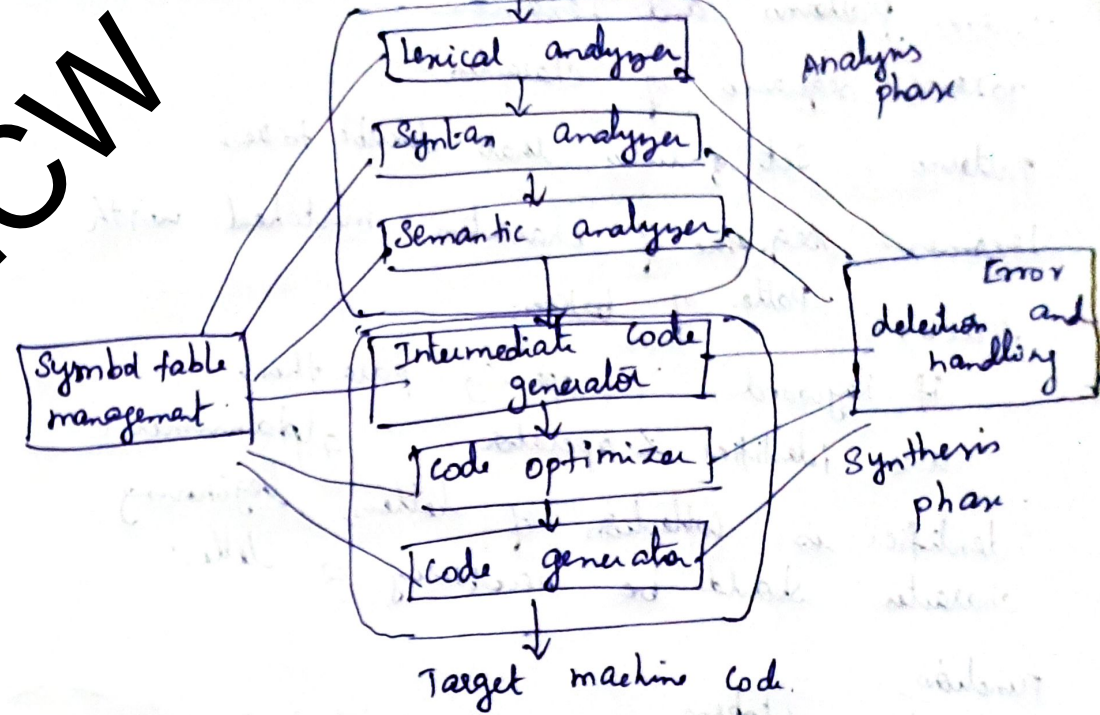
$\boxed{\text{Lexical analyzer}}$

TOKENS

$t_3 := id_2 + t_2 \quad id_1 = id_2 + id_3 * 60$
$id_1 := t_3$

$\boxed{\text{Syntax analyzer}}$  Syntax tree

Code optimize

$id_1$
$id_2$ +
$id_3$ * 
$id_3$ 60.

$t_1 := id_3 * 60$
$id_1 := id_2 + t_1$

$\boxed{\text{Semantic analyzer}}$  Semantic tree

Code generate

$id_1$
$id_2$ +
$id_3$ * int to float
$id_3$ 60.

MOV F $id_3, R_2$
MUL F #60.0, $R_2$
MOV F $id_2, R_1$
ADD F $R_2, R_1$
MOV F $R_1, id_1$

intermediate code  ic
$t_1 := $ int to float (60)
$t_2 := id_3 * t_1$

---

source program.

$\boxed{\text{Lexical analyzer}}$  Analysis phase

$\boxed{\text{Syntax analyzer}}$

$\boxed{\text{Semantic analyzer}}$

$\boxed{\text{Symbol table management}}$  $\boxed{\text{Intermediate code generator}}$  $\boxed{\text{Error detection and handling}}$

$\boxed{\text{Code optimizer}}$  Synthesis phase
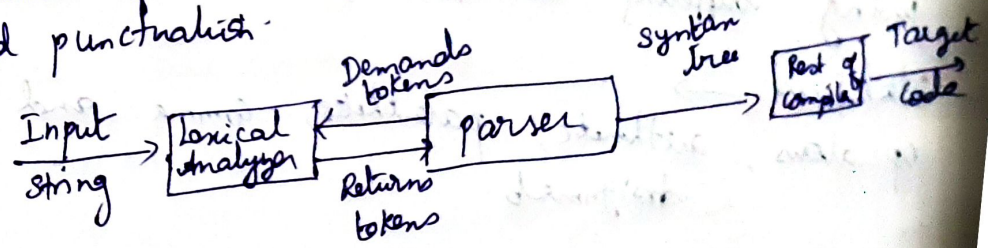
$\boxed{\text{Code generator}}$

Target machine code.

## Lexical Analysis

Role of Lexical Analyzer (first phase of compiler)

Reads the input from source program left to right one character at a time and generates the sequences of tokens (identifiers, keyword operators and punctuation.

Input string → $\boxed{\text{Lexical Analyzer}}$ ⇄ Demands tokens / Returns tokens $\boxed{\text{parser}}$ → Syntax tree $\boxed{\text{Rest of compiler}}$ → Target code

Tokens, Patterns and Lexemes:

Tokens: sequence of character

Patterns: Set of rules, that describe tokens

Lexemes: sequence of characters, matched with pattern of token.

if(a<b)

if -keyword `(` opening parenthesis

a ← identifier < operator alphanumeric

identifier is collection of letters, beginning character should be necessarily a letter.

Functions
1) Stream of tokens

Issues of lexical Analyzer.
1) Lexical Analysis and Syntax Analysis are separated out, Reduces burden of on parsing phase, using buffering techniques for efficient scan
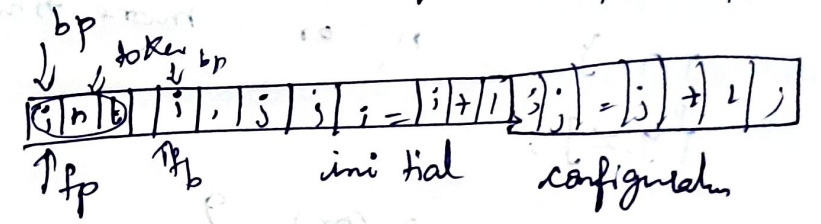
identifiers →
operators, arithmetic, parenthesis, comma and Assignment.

Keyword: some special words, int, void are keywords

Input Buffering:-

is used to identify the lexeme correctly using two pointer method

scan left to right   begin_ptr (bp)
                     forward-ptr (fp)



one buffer scheme: lexeme is very long, crosses buffer boundary, buffer has been refilled, overwriting the first part of lexeme.

Two buffer scheme: sentinel. identify the end of buffer

Regular expression:-

[abc] = a, b or c.

[^abc] = any character except a, b, c

[a-z] = a to z, [A-Z] = A to Z

[0-9] . 0 to 9.

[ ] ? → 0 or 1 tims

[ ] + occurs 1 or more

[ ] * " 0 or more.

{n} = occurs n times

{n,} n or more times

{y,z} :

1) mobile no's 8 (or)

[8 9] [0-9]{9}.

2) upper case, contains lower case chars,
one digit in between.

[A-Z][a-z]*[0-9][a-z]

\d [0-9]

\J [^0-9]

\w [a-z, A-Z, 0-9]

[^\w]

---

Mail ID. abc123 @ gmail.com

[a-z A-Z0-9_\-\.] + [@] [a-z] + [\.] (a-z){2,3}

## Recognition of Tokens:

Token representation

| Token type | Token value | [token attribut |
|---|---|---|

↓ category.

↓ information regarding

1) Symbol table is maintained
2) identifier & constants (pointer to symbol table)

| Token | code | value |
|---|---|---|
| if | 1 | |
| else | 2 | |
| while | 3 | |
| for | 4 | |
| identifier | 5 | ptr to ST |
| constant | 6 | ptr to ST |
| < <= > >= != = | 7 | 1,2,3,4,5 |
| ( ) | 8 | 1, 2 |
| + - | 9 | 1, 2 |
| = | 10 | |

steps to recognize tokens

1) stores input in input buffer

2) regular expression is built for corresponding token

3) N Deterministic finite automata is built.

4.

Finite Automata:
5 tuples $(Q, \Sigma, \delta, q_0, F)$.

$Q \rightarrow$ finite set of states, non empty.
$\Sigma \rightarrow$ input alphabet (finite set of i/ps.
$\delta -$ transition function, next state can be determined
$q_0 -$ initial state $q_0 \in Q$
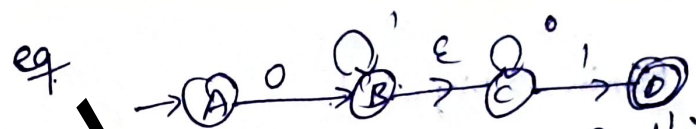$F \rightarrow$ final states.

Types of automata $\rightarrow$

Finite automata

| DFA | | NFA |

$q_1 = \delta(q_0, a)$. Only for current
next state       current input.   i/p current
state                              state.

NFA with $\varepsilon$.
$\quad \hookrightarrow$ empty symbol.

Regular NFA with 5 tuples.
$\{Q, \Sigma, q_0, F, \delta\}$.
where $\delta: Q\Sigma \rightarrow 2^Q$

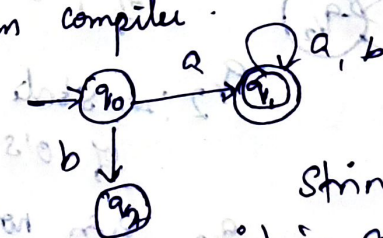$\varepsilon$-NF $\{Q, \Sigma, q_0, F, \delta\}$   $\delta: Q \times \Sigma \cup \varepsilon \rightarrow 2^Q$

eg

$\rightarrow (A) \xrightarrow{0} (B) \xrightarrow{\varepsilon} (C) \xrightarrow{1} ((D))$

$\rightarrow$ B is not seeing anything, but it goes to C.
here every state on $\varepsilon$ goes to it self.
A on $\varepsilon$ goes to A.

Deterministic finite Automata.

* if the m/c is read on i/p string one
symbol at a time.   DFA $\xrightarrow[b]{a}$

DFA uniqueness of the computation.
DFA only one path for specific i/p from the
current state to the next state.  $(A) \xrightarrow{0} B$
DFA not accept null move, does not change
state without any i/p character.
        DFA can have multiple final states. lexical
Analysis in compiler.

$\rightarrow (q_0) \xrightarrow{a} (q_1) \circlearrowleft a, b$   current i/p  next
                                        state sym    state
          $\downarrow b$              $\delta: Q \times \Sigma \rightarrow Q$
        $(q_2)$
                          String reaching final state
                          it is acceptance.

$\{ \phi \} \Rightarrow \rightarrow ((q_1))$
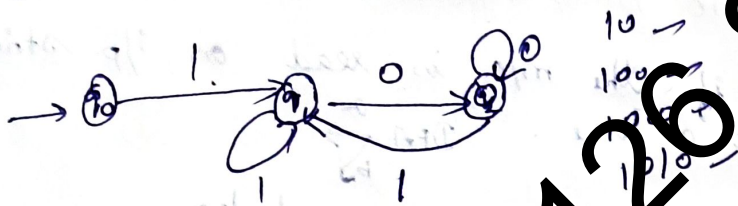$L\{\varepsilon\} = ((q_1))$

no/. of state = Maxim string + 1.

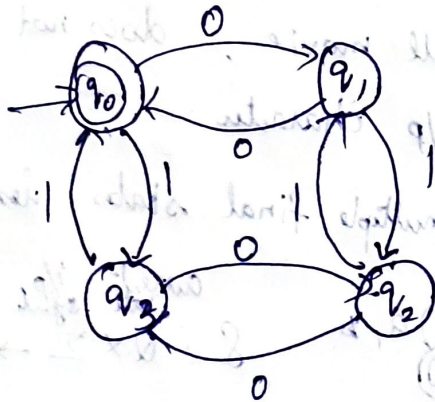q.    DFA with  $\Sigma\{0,1\}$.
          Starts with & ends with 0.

$L = \{10, 100, 1000, 1010, 1100, 1110 \ldots \}$

Min length = 2,      No/. of states = 2+1 = 3



2) FA with  $\Sigma = \{0,1\}$  even no/. of 1's &
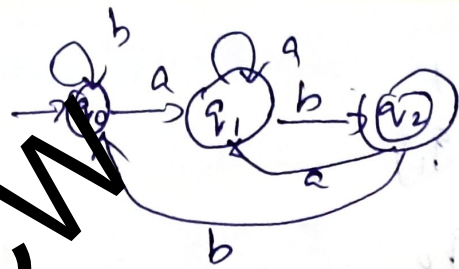$L = \{001, 110, 100, 110, \ldots\}$  even no/. of 0's



even no/. of 1's,
$q_0$: even no/. of 0's.

$q_1$ = State of odd no/.
     of 0's, even no/
     of 1's

$q_2$ :→ state of odd no.
       of 0's & 1's

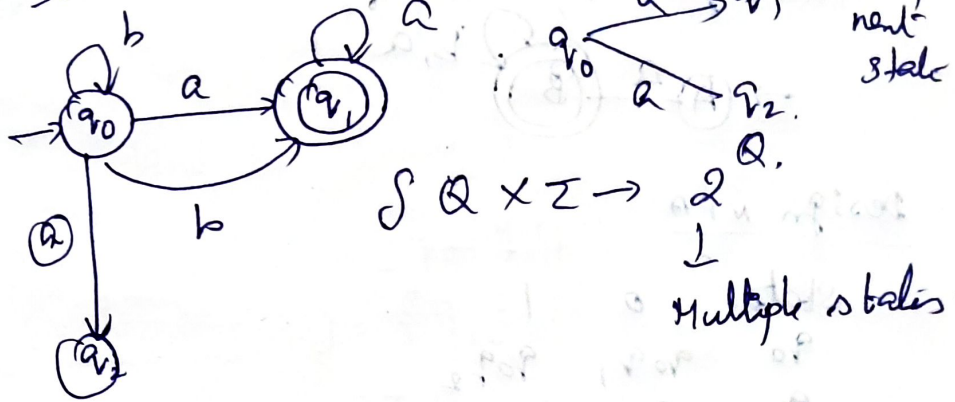$q_3$. even no/. 1 of 0's
     odd no/ of 1's

---

ab, baab, aab, bbaab



NFA Non Deterministic finite automata.

NFA is not DFA,  each NFA translated to DFA

NFA is defined in the same way as DFA
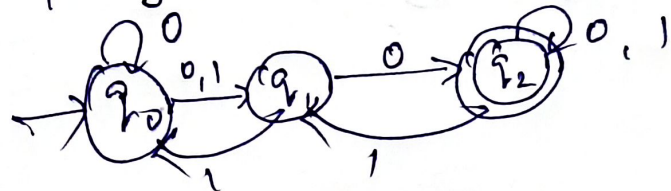but. It contains multiple next state.
      It contains $\in$ transitions.



q.



$\delta$  $Q \times \Sigma \to$  $2^Q$

$\underline{\hspace{1cm}}$
Multiple states.

q.   $Q = \{q_0, q_1, q_2\}$
     $\Sigma = \{0,1\}$,   $q_0 = \{q_0\}$     $F = \{q_2\}$

Soln

| State | 0 | 1 |
|---|---|---|
| → $q_0$ | $\{q_0, q_1\}$ | $\{q_1\}$ |
| $q_1$ | $q_2$ | $q_0$ |
| $(q_2)$ | $q_2$ | $\{q_2, q_1\}$ |

**Example**

NFA starts with 'a' $\Sigma = \{a, b\}$

$L = \{a, ab, abb, aba, ab\, abab..\}$

Min Length = 1, no. of states 3.



**Design NFA**

| State | 0 | 1 |
|---|---|---|
| $q_0$ | $q_0 q_1$ | $q_0 q_2$ |
| $q_1$ | $q_3$ | $\varepsilon$ |
| $q_2$ | $q_1 q_3$ | $q_3$ |
| • $q_3$ | $q_3$ | $q_3$ |

**NFA with $\varepsilon$ to NFA without $\varepsilon$.**



$\{Q, \Sigma, \delta, q_0, F\}$

$Q = \{q_0, q_1\}$ $\Sigma = \{0, 1, \varepsilon\}$ $q_0 = q_0$, $F = \varepsilon$,

**Transition table**

| State\Input | 0 | 1 | $\varepsilon$ |
|---|---|---|---|
| $q_0$ | $\{q_0\}$ | - | $\{q_1\}$ |
| $q_1$ | - | - | $\{q_1\}$ |

**Step : 2 :-** find $\varepsilon$- closure.

$\varepsilon$-closure $(q_0) = \{q_0, q_1\}$

$\varepsilon$-closure $(q_1) = \{q_1\}$

**Step 3** : processing states $q_0, q_1$.

$\underline{q_0}$ $\delta'(q_0, 0) = \varepsilon\text{-closure}(\delta(q_0, q_1), 0)$

$= \varepsilon\text{-closure}(q_0)$

$= q_0, q_1,$

$\delta'(q_0, 1) = \varepsilon\text{-closure}(\delta(q_0, q_1), 1)$

$$= \varepsilon \text{ closure } (q_1)$$
$$= \{q_1\}.$$

$$\delta'(q_1, 0) = \varepsilon - \text{closure } (\delta(q_1, 0)).$$
$$= \varepsilon - \text{closure } (\phi)$$
$$= \phi.$$

$$\delta'(q_1, 1) = \varepsilon - \text{closure } (\delta(q_1, 1))$$
$$= \varepsilon - \text{closure } \{q_1\}$$
$$= \{q_1\}$$

## ε NFA to NFA



|  | 0 | 1 |
|---|---|---|
| → A | {A, B, c} | {B, c} |
| B | {c} | {B, c} |
| * c | {c} | {c} |

$$\varepsilon^* \quad 0 \quad \varepsilon^*$$

| A | A → A | A | $A_{B_c}$ |
|---|---|---|---|
|  | B — $\phi$ | — |  |
| C | — C | C |  |

$$A \quad \varepsilon^* \quad 1 \quad \varepsilon^*$$

| A | A | $\phi$ |  |
|---|---|---|---|
|  | B | B — B, C |  |
| C | C → C |  |  |

B
| | $\varepsilon^*$ | 0 | $\varepsilon^*$ |
|---|---|---|---|
| B | — $\phi$ — | | |
| C | — C — | C. | |

| | $\varepsilon^*$ | 1 | $\varepsilon^*$ |
|---|---|---|---|
| B — B | | B | B, c |
| C. | | c | c. |

C
| | $\varepsilon^*$ | 0 | $\varepsilon^*$ |
|---|---|---|---|
| C | C | C | |

| | $\varepsilon^*$ | 1 | $\varepsilon^*$ |
|---|---|---|---|
| C — C | | C | C |

$\varepsilon$-closure ($\varepsilon^*$) - All the states that can be reached from a particular state only by seeing $\varepsilon$ symbol.



## NFA with ε to DFA.



| ε-NFA | a | b | c | ε. |
|---|---|---|---|---|
| → $q_0$ | {$q_0$} | $\phi$ | $\phi$ | {$q_1$} |
| $q_1$ | $\phi$ | {$q_1$} | $\phi$ | {$q_2$}. |
| * $q_2$ | $\phi$ | {$q_2$} | {$q_2$} | $\phi$ |

Set of all states that reach particular state, include that state

$\varepsilon$ closure $(q_0) = \{q_0, q_1, q_2\}$

$\varepsilon$ closure $(q_1) = \{q_1, q_2\}$

$\varepsilon$ closure $(q_2) = \{q_2\}$.

| DFA $\delta_D$ | a | b | c |
|---|---|---|---|
| $\rightarrow q_0, q_1, q_2$ | $[q_0, q_1, q_2]$ | $\{q_1, q_2\}$ | $\{q_2\}$ |
| $[q_1, q_2]$ | $\phi$ | $\{q_1, q_2\}$ | $\{q_2\}$ |
| $q_2$ | $\phi$ | $\phi$ | $q_2$ |

Start state → $\varepsilon$ closure of $\{q_0, q_1, q_2\}$

$q_0 \delta a \cup q_1 \delta a \cup q_2 \delta a$

$\delta_D (\{q_0, q_1, q_2\}, a)$

$= \varepsilon$-closure $(\delta(q_0, q_1, q_2), a)$

$= \varepsilon$-closure $(\delta(q_0, a) \cup \delta(q_1, a) \cup \delta(q_2, a))$

$= \varepsilon$-closure $(\{q_0\} \cup \phi \cup \phi)$

$= \varepsilon$-closure $(q_0)$

$= q_0, q_1, q_2$

Regular Expression to $\varepsilon$-NFA (Thomson Construction)

to $\varepsilon$-NFA.
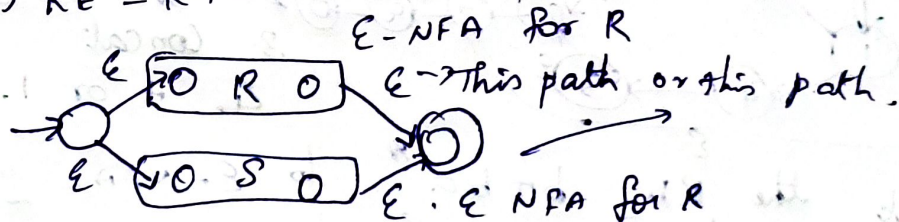
Basis:

1. $RE = \varepsilon$.

2. $RE = \phi$.

3. $RE = a$

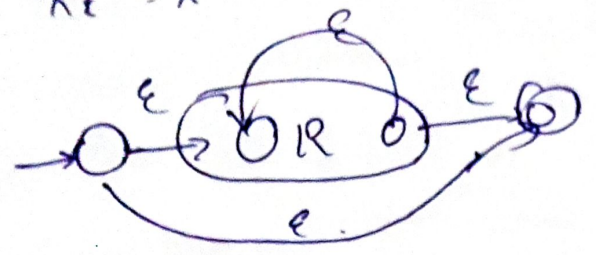Induction: → Regular Expression.

1) $RE = R + S$ or $R|S$.

$\varepsilon$-NFA for R

$\varepsilon$ → This path or this path.

$\varepsilon$ - $\varepsilon$ NFA for R

$RE = (0+1)^* | (01)$

    R    S

2) $RE = RS$ [Concatenation]

→ identify element.

3, $RE = R^*$



4, $RE = (R)$.

$\varepsilon$ NFA for $R$ = $\varepsilon$-NFA for $(R)$

Convert the $RE_\varepsilon$ $(0+1)^* 10$ to $\varepsilon$-NFA.



Precedence
1. ( )
2. *
3. Concat
4. + or .

2, Convert the $RE$ $b + ba^*$ to $\varepsilon$-NFA



3, $r = a^* b$



## Regular expression to DFA

RE to $\varepsilon$ NFA to DFA.

1, $R-E \to$ NFA with $\varepsilon$.
2, Convert NFA with $\varepsilon$ to NFA without $\varepsilon$
3, Convert the obtained NFA to equivalent DFA.

$R.E \to NFA-\varepsilon \to NFA \to DFA$.

1) Design FA for given $R.E$ $10 + (0+11)0^*$

10/9/2023

LEX = {tool}
↓
→ break up an i/p streams into tokens.
→ automatically generating a Lener (scanner)

lex.l → [Lex Compiler] → lex.yy.c

lex source program,

Step 2   lex.yy.c →  [ C compiler ] → a. out } Lexical
L.A. Analyzer

Step 3:

i/p stream → [ a. out ] → seq of tokens.

i/p → [ L A ] → tokens.
↳ charactr one by one.
↳ check identifier, opers,

help of lex tool →

## STRUCTURE of LEX Program:

{ declaration } ⟹ declaration of variables.
%. %.

rules {
section} translation rules } ⟹ have the pattern { Action }
%. { i/p streams represent-
procedure %.    in Lex program.
section} auxiliary functions } ⟹ funs. can be compiled
separately.

Lex program:-
no). of vowel & constants.

%. { # include <stdio.h>
int vowels = 0;
int Cons = 0;
%. }

i/p stream
↓
grammar.

---

%. %. { a eiou AEIOU] { vowels ++; }
[a-z A-Z] { const ++; }

int yy wrap()
{
return 1;
}

main()
{
print(" enter the string at end press ^ \n")
yylex(); [ find the vowels & cons
printf(" no). of vowels = %.d \n
no). of cons = %.d \n",
vowels, cons);
}

1) session starting & ending %., { & %.}.
2) %. %.
3) two functions, main functions & yywrap functions.
yylex routine is given lex.yy.c.

$ lex x.l
↳ program name
$ cc lex.yy.c [ gcc can also used,
$ .\a. out.

yylex()

Starting point of lex from which scanning of source program starts.

yywrap()

This function is called when end of file is encountered. If yywrap returns 0 the scanner continues scanning if it returns 1 the scanner does not returns tokens.

Role of parser - Grammars - content free grammars - writing a grammar Top Down parsing - General strategies - Recursive descent predictive parser - LL(1) Parser - Shift Reduce parse - LR parse - LR(1). Item construction of SLR parsing table - Introduction on to LALR parser - Error Handling and Recovery in Syntax Analyzer - yacc tool - Design of a syntax analyzer for a sample language.

## Introduction of Syntax Analyzer:

It is a second phase in Compilation (parser) ⇒

A parsing is a process which takes the input string w and produces either a parse tree (syntactic structure). Or generates the Syntactic errors.    $a = b + 10;$

$$a \diagup = \diagdown $$
$$b \diagup + \diagdown 10$$

# Basic issues of parsing

(i) specification of syntax, (ii) Representation of input after parsing.

(i) precise & unambiguous, in detail complete.

Content free grammar.

(ii) information during parse tree has been generation, shant not be different of generating parse tree.

"A/"

# Role of parser



Source program → Lexical Analyzer ⇄ Parser → parse tree → Rest of front End → intermediate representation

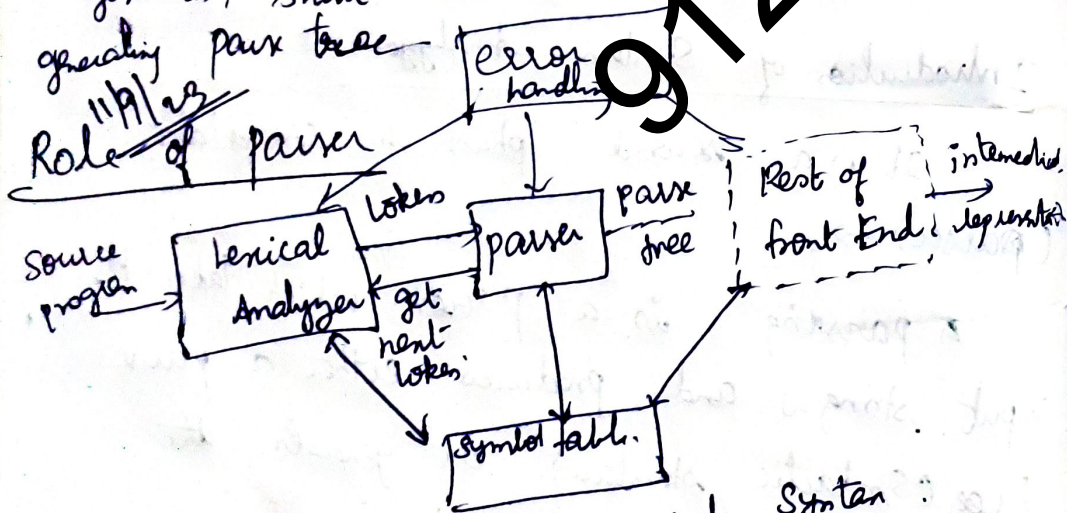Lexical Analyzer: tokens, get next token

Symbol table.

error handling

1) Syntax Analysis: check source code Syntax.
2) parsing: Break code into tree like structure.
3) Data structure: Represents syntax relationship using parse trees.
4) Error handling: Generates error messages for developers.
5) contextual Analysis: Basic semantic checks.

---

b) Intermediate Representation: Eases further Compiler phases (Processor is optimized is easy, tokens are in content free.

## Error Handling:

1. Lexical phase errors ⎡ spelling errors.
   ⎣ exceeding length of identifies
   ⎣ illegal character appears

   normally typing mistake, the wrong spelling

2) Syntax analysis phase errors ⎡ error in structure
   ⎣ Missing operators
   ⎣ unbalanced parenthesis

3) Semantic errors ⎡ Incompatible types of operands.
   ⎣ undeclared variables.
   ⎣ No matching of actual with formal parameters

## Content free grammar: → formal grammar, used to generate all possible pattern of string in a given finite language

1) V → non terminals
2) T → terminals
3) S → start symbol.
4) p → set of production rules

$G = (V, T, S, P)$ → set of production rules used to generate string of language

Production rules

non terminal $\longrightarrow (V \cup T)^*$

P → production rules which is used ttnches
replacing non- terminals symbols.

$P^o →$ aSb.   $E → E+E$
    always     $\langle S → AS.$
    forms   $S → b/a$.

1) $L = a^n b^n$   when $A \geq 1$.

$G: \{v.T S p\}$.

$v = \{S\}$,
$T = \{a,b\}$,   S is start symbol.

$P = \{ S → aSb$
       $S → ab$

1) Construct CFG for the language
having any nol. of a's are the
set $\Sigma = \{a\}$
    $\Sigma: \{a\}$
    $L = \{a, a, aa, aaa \cdots\}$
    $R.E = a^*$.

---

Production rule.

$S → aS$ —①
$S → \varepsilon$ —②
        Type List terminal
2) $S → T L T$
Type → int | float
List → List, id
List → id.;
Terminal → ;

int id, id, id ;

3) $S → a A B e$
   $A → Abc | b$
   $B → d$.

i) right most   a b b c d e
ii) left most   a b b c d e
3) parse tree

i/p $a a a a a$
⇒ $aS$   $S → aS$
⇒ $aaS$  $S → aS$
⇒ $aaaS$ $S → aS$
⇒ $aaaaS$ $S → aS$
⇒ $aaaaa S$
= $aaaaa$. $S → \varepsilon$


State
    T   L   T
    |  /|\
   int List, id
       /|\
     List , id
      |
     id.

**Soln** right most → a b b c d e

$S \to$ a A B e. ←

$S \to$ a A d e : $B \to$ d

$S \to$ a A b c d e. $A \to$ A b c

$S \to$ a b b c d e $A \to$ b

**left most** a b b c d e

$S \to$ a A B e.

$S \to$ a A b c B e. $A \to$ A b c

$S \to$ a b b c B e $A \to$ b

$S \to$ a b b c d e : $B \to$ d



a b b c d e

---

**Ambiguous grammar :-**

more than one → left most. right most parse tree

gives i/p string

If grammar is ambiguity, then it is not good for compiler construction.

must remove ambiguity

$E \to I$     $G = \{ v T P S\}$   $v = \{I, E\}$.

$E \to E + E$    $r : \{ +, *, (,), \varepsilon, 0 \cdots 9\}$

$E \to E * E$      $3 * 2 + 5$.

$E \to ( E )$

$I \to \varepsilon | 0 | 1 | 2 | \cdots 9$.

Since have two parse tree, so it is an Ambiguos grammar



$3 * 2 + 5$.    $3$   $2$

abab

$S \to a S b S$
$S \to b S a S$
$S \to \varepsilon$



right most

## Parsing Techniques:

### Top-Down parser:

The process of construction of parse tree starting from root & proceed to children is called TDP. Top → down has to be scanned.

Top down → Recursive Descent
→ Back tracking
→ Non back tracking
↓
Predictive parser
↓
LL parser

TDP internally uses left most derivation free from ambiguity - left recursion.

---

## Classification of TDP:

with back tracking → brute force technique

without back tracking → predictive parse.

### Recursive Descent parsing. LL(1)  recursive descent parser.

1) parse construct from top to Bottom.
2) i/p is read from left to right.
3) i/p is recursively passed for preparing a parse tree with or without back tracking

### Back tracking:

$W = cad,$
$S \to cad$
$A \to ab | a$



### Left recursion

$A \xrightarrow{+} A\alpha | \beta.$

Keep calling itself.

$A \to A\alpha | \beta. \Rightarrow A \to \beta A'$
$\qquad\qquad A' \to \alpha A' | \varepsilon.$

$S \to ABc$
$A \to Aa|nd|b$
$B \to Bb|c$
$c = Cc|g$

$S \to ABc$
$A \to bA'$
$A' \to aA'|dA'|\varepsilon$
$B \to cB'$
$B' \to bB'|\varepsilon$
$c \to gc'$
$c' \to ec'|\varepsilon$

$A = Aa|b$
$\boxed{\begin{array}{l} A = bA' \\ A' = aA'|\varepsilon \end{array}}$

$A = Ad|b$
$\boxed{\begin{array}{l} A = bA' \\ A' = dA'|\varepsilon \end{array}}$

$B = Bb|c$
$\boxed{\begin{array}{l} B \to cB' \\ B' \to bB'|\varepsilon \end{array}}$

$E \to E+T|T, \quad T \to T*F)|F$
$A \to ABd|Aa|a$
$B = Be|b$

**Left factoring:** two or more production are starting with same set of symbols

$A \to \alpha\beta_1 | \alpha\beta_2$
$\quad \alpha, \beta_1, \beta_2 \Rightarrow$ string

$A \to \alpha A'$
$A' \to \beta_1|\beta_2$

$S \to iEts|;Etses|a$
$\to b$

$S \to iEtSS'|a$
$S' \to es|a$
$E \to b$

$A = \alpha A'$
$A' = \beta_1/\beta_2$

$A \to AC|Aad|bd|c$
left recursion

left factoring $A \to aAB|aA|a$
$\qquad B \to bB|b$

**Removal of ambiguity.**
(i) By adding precidence.
(ii) By adding associativnes

id +id +id.

$E \to E+id|id$
$E \to E+T|T$
$T \to T*F|F$
$F \to id$.

$E \to E+id|id$.

$E \to E+E|E*E|id$.

id+id * id



$E \Rightarrow E+T \mid T$
$T \rightarrow T*F / F$
$F \rightarrow id.$

## Recursive Descent Parsing:

Collection of recursive procedures.

non-terminal a separate procedure is written.

### Advantages:-

1) Simple to build.
2) Constructed with the help of parse tree

### Limitations:

1) not very efficient
2) may enter infinite loop
3) not provide good error message
4) difficult to parse a string

---

## Applications of FA (need)

1) design of the lexical Analysis of a compiler.
2) Recognize the pattern by using regular Expression
3) Helpful in text editors
4) used for spell checkers.
5) use of the Mealy & Moore machine for designing the combinational & sequential circuits.

24/9/23

### Predictive parser LL(1) parser:-

top-down parsing algorithm, non-recursive type. LL(1). first L→ input is scanned from left to right.
second L→ left most derivation for i/p string.

Data structure: (i) i/p buffers → to store i/p tokens.
(ii) stack → hold left sentential form. pushed into stack in reverse order left to right

(iii) parsing table :→ row for non terminal, column for terminal $M[A, a]$.



terminal,
current i/p symbol.

(i) Elimination of left Recursion

(ii) Left Factoring.                    $E = E + T \mid T$

(iii) First & Follow functions    $T = T * F \mid F$

(iv) predictive parsing table    $F = (E) \mid id$

(5) parse the input string

## Elimination of Left recursion

$E \Rightarrow \underset{A}{E} \underset{\alpha}{\pm T} \mid \underset{B}{T}$     $T - T * F \mid F$

$E \rightarrow T E'$     $T \rightarrow F T'$

$E' \rightarrow + T E' \mid \varepsilon$     $T' \rightarrow *F T' \mid \varepsilon$

$E \rightarrow T E'$

$E' \rightarrow + T E' / \varepsilon$

$T \rightarrow F T'$

$T' \rightarrow * F T' / \varepsilon$

$F = (E) \mid id$

### FIRST

$FIRST(E) = FIRST(T) = FIRST(F)$
$= \{(, id\}$

$FIRST(E') = \{+, \varepsilon\}$

$FIRST(T') = \{*, \varepsilon\}$

### Follow :-

1, $\$ \rightarrow$ follow of Start Symbol of grammar

2, If $A \rightarrow \alpha B \beta$

3, $A \rightarrow \alpha B$

---

$Follow(E) = \{\$, )\}$

$Follow(E') = \{\$, )\}$

$Follow(T) = \{+, \$, )\}$

$Follow(T') = \{+, \$, )\}$

$Follow(F) = \{*, +, \$ )\}$

$A \rightarrow \alpha$

(1) $M[A, a] = A \rightarrow \alpha$
   $a$ is in $FIRST(\alpha)$

$M[A, b] = A \rightarrow \alpha$
   $\varepsilon$ is in $First(\alpha)$,
   $b$ is in $Follow(A)$.

### Parsing table    Row → non terminals
                     column → terminals

|     | +              | *              | (           | ) | , id           | $              |
|-----|----------------|----------------|-------------|---|----------------|----------------|
| E   |                |                | E→TE'       |   | E→TE'          |                |
| E'  | E'→+TE'        |                |             | E'→ε |             | E'→ε           |
| T   |                |                | T→FT'       |   | T→FT'          |                |
| T'  | T'→ε           | T'→*FT'        |             | T'→ε |            | T'→ε           |
| F   |                |                | F→(E)       |   | F→id           |                |

$W = id * id + id$

| Stack     | Input           | Output       |
|-----------|-----------------|--------------|
| $E        | id * id + id$   |              |
| $E'T      | id * id + id $  | E→TE'.       |
| $E'T'F    | id * id+id $    | T→FT')       |
| $E'T'id   | id * id+id $    | F→id.        |
| $E'T'     | *id + id $      |              |
| $E'T'F*   | *id +id $       | T'→*FT'      |
| $E'T'F    | id +id $        |              |

$*$ & $*$
in

$ E'T'id          id +id $.    F→id.
$ E'T'           +id $       T'→ε.
$ E'             +id $       E'→+TE'.
$ E'T +          +id $.
$ E'T $           id $        T→FT'
$ E'T' F.         id $.       F→id
$ E'T' id         id $
$ E'T'            $.
$ E'.                          $ →ε
$.                    $        E'→ε.



=>  id * id +id

## Follow:

If following the variable, you have.

Terminal → write it as it is.

Non terminal → write its First element.

last element → write Follow of LHS.

Follow set will never contains NULL.

$S→ i E t S_1 / a$
$S_1 → e S / ε$
$E → b$

First (S) = { i, a }        Follow (S) = { e, $ }
First (S_1) = { e, ε }       Follow (S_1) = { e, $ }
First ( E ) = { b }         Follow ( E ) = { t }

FOLLOW    1, Follow(S) = { $ }, 3, $A → dBβ$, then.
                        FOLLOW (B) = FIRST (B).
                                     except ε.

3, If $A → αB$  or  $A → αBβ$  where  FIRST (B)
   contains ε     (B→ε),
        FOLLOW(B) = FOLLOW(A).

$S → ABCDE$        FIRST
$A → a/ε$          FIRST (S) = { a, b, c }.
$B → b/ε$          FIRST (A) = { a, ε }
$C → c$            FIRST ( B) = { b, ε }
$D → d/ε$          FIRST (C) = { c }
$E → e/ε$.         FIRST (D) = { d, ε }
                   FIRST ( E) = { e, ε }

FOLLOW(S) = $

FOLLOW(A) = {b, c}

FOLLOW(B) = {c}

FOLLOW(C) = {d, e, $}

FOLLOW(D) = {e, $}

FOLLOW(E) = {$}

FIRST(S) — {d, g, h, ε, b, a}.

FIRST(A) — {d, g, h, ε}.

FIRST(B) — {g, ε}

FIRST(C) — {h, ε}

FOLLOW(S) — {$}

FOLLOW(A) — {h, g, $}

FOLLOW(B) — {$, a, h, g}

FOLLOW(C) — {g, $, b, h}.

## Handle pruning :

bottom up parsing to find the sub-string
that could be reduced by appropriate re-bind

2)

S → A (B) CbB | B.

A → d a | Bc

B → g / ε

C → h / ε

---

Shift Reduce parsing bottom up parsing

**Shift:** shift the next input symbol onto the top of the stack

**Reduce:** The right end of the string to be reduced must be at the top of the stack. Locate the left side of the string within the stack and decide within what non terminal to replace the string.

**Accept:** Announce successful completion of parsing.

**Error** — Discover a syntax error and call an error recovery routine.

S → cc

C → eC/d

| Stack | Input |
|-------|-------|
| $ | w $ |
| $ s | $. |

w → cdcd.

| Stack | Input | Action |
|-------|-------|--------|
| $ | cdcd $. | Shift. |
| $ c | dcd $ | Shift. |
| $ cd | cd $ | reduced by C → d. |
| $ eC | cd $ | reduced by C → cC accept |
| $ C | cd $ | shift. |
| $ Cc | d $ | shift |
| $ Ccd | $ | reduce by C → d. |

$ ccc        $    reduce by

                  c → cc

$ Cc    $    red
              3 → cc

$ s      $

## Left column

c − d − c d

```
      S
   c / \ c
  /|    /|
 c      c
 |      |
 c d   c d
```

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow id$

$id * (id + id)$

| Stack | I/p string | Action |
|---|---|---|
| $\$$ | id * (id + id) $ | shift |
| $\$ id$ | * (id + id) $ | reduced E → id |
| $\$ E$ | * (id + id) $ | shift |
| $\$ E *$ | (id + id) $ | shift |
| $\$ E * ($ | id + id) $ | shift |
| $\$ E * (id$ | + id) $ | reduce E → id |
| $\$ E * (E$ | + id) $ | shift |
| $\$ E * (E +$ | id) $ | shift reduce E → id |
| $\$ E * (E + id)$ | ) $ | shift reduce |
| $\$ E * (E + E$ | ) $ | reduce E → E + E |
| $\$ E * (E$ | ) $ | shift |
| $\$ E * (E)$ | $ | reduce E → (E) |
| $\$ E * E$ | $ | reduce E → E * E |
| $\$ E$ | $ | accept |

## Right column

4/10/23

LR parsers, non recursive, shift reduce, bottom up parse

L → left to right scanning of input stream

R → construction of right most derivation in reverse.

K → no. of lookaheads needed for derivation

$LR(0) < SLR(1) < LALR(K?) < CLR(1)$

  simple LR        lookahead LR        Canonical LR

CLR(1) is also known as LR(1)

| a | + | b | $ |



LR(0)

$S \rightarrow A A$

$A \rightarrow a A \mid b$

$S' \rightarrow \cdot S$ → augmented production

$S \rightarrow \cdot A A$ → dot in beginning.

$A \rightarrow \cdot a A / \cdot b.$ dot in end. my full production is parsed.

Left diagram item sets:

$I_0$
$S' \to .S$
$S \to .AA$
$A \to .aA \mid .b$

$I_1$: $S' \to S.$

$S \to A.A$
$A \to .aA \mid .b$

$A \to aA$ $(S \to AA)$

$A \to a.A$
$A \to .aA \mid .b$

$A \to b.$

$S \to AA.$

$\to aA.$

$S \to A.A$
$A \to .aA \mid .b$

| | action | | | goto | |
|---|---|---|---|---|---|
| | a | b | \$ | A | S |
| 0 | $s_3$ | $s_4$ | | 2 | 1 |
| 1 | | | accept | | |
| 2 | $s_3$ | $s_4$ | | 5 | |
| 3 | $s_3$ | $s_4$ | | 6 | |
| 4 | $r_3$ | $r_3$ | $r_3$ | | |
| 5 | $r_1$ | $r_1$ | $r_1$ | | |
| 6 | $r_2$ | $r_2$ | $r_2$ | | |

## SLR

$E \to E+T$  1
$E \to .T$  2
$T \to T*F$  3
$T \to F$  4
$F \to (E)$  5
$F \to id.$  6

$T \to .F$
$F \to .(E)$
$F \to .id.$

goto($I_0$, id)
$I_5$: $F \to id.$

$I_0$:
$E' \to .E$
$E \to .E+T$
$E \to .T$
$T \to .T*F$
$T \to .F$
$F \to .(E)$
$F \to .id$

goto($I_1$, +)
$I_6$: $E \to E+.T$
$T \to .T*F$
$T \to .F$
$F \to .(E)$
$F \to .id.$

goto($I_4$, E)
$I_8$: $F \to (E.)$
$E \to E.+T$

goto($I_0$, E)
$I_1$: $E' \to E.$
$E \to E.+T$

goto($I_0$, T)
$I_2$: $E \to T.$
$T \to T.*F$

goto($I_2$, *)
$I_7$: $T \to T*.F$
$F \to .(E)$
$F \to .id$

goto($I_0$, F)
$I_3$: $T \to F.$

goto($I_0$, ()
$I_4$: $T \to (.E)$
$E \to .E+T$
$E \to .T$
$T \to .T*F$

goto($I_6$, T)
$I_9$: $E \to E+T.$
$T \to T.*F$

goto($I_7$, F)
$I_{10}$: $T \to T*F$

goto($I_8$, ))
$I_{11}$: $F \to (E).$

| State | Action | | | | | | goto | | |
|---|---|---|---|---|---|---|---|---|---|
| | id | + | * | ( | ) | \$ | E | T | F |
| 0 | $S_5$ | | | $S_4$ | | | 1 | 2 | 3 |
| 1 | | $S_6$ | | | | accept | | | |
| 2 | $r_2$ | $r_2$ | $S_7$ | | $r_2$ | $r_3$ | | | |
| 3 | $r_4$ | $r_4$ | $r_4$ | $r_4$ | | | | | |
| 4 | $S_5$ | | | $S_4$ | | | 8 | 2 | 3 |
| 5 | $r_6$ | | $r_6$ | $r_6$ | $r_6$ | $r_6$ | | | |
| 6 | $S_5$ | | | $S_4$ | | | | 9 | 3 |
| 7 | $S_5$ | | | $S_4$ | | | | | 10 |
| 8 | | $S_6$ | | | $S_{11}$ | | | | |
| 9 | $r_7$ | $r_7$ | $r_7$ | | $r_7$ | $r_7$ | | | |
| 10 | $r_3$ | $r_3$ | $r_3$ | | $r_3$ | 3 | | | |
| 11 | $r_5$ | $r_5$ | $r_5$ | $r_5$ | $r_5$ | | | | |

| stack | input | action |
|---|---|---|
| $0 | id * id+id $ | Shift |
| $0 id 5 | * id+id $ | reduce F→id. (1×2) |
| $0 F 3 | * id+id $ | reduce T→F (1×2) |
| $0 T 2 | * id+id $ | Shift. |
| $0 T 2 * 7 | id+id $ | Shift |
| $0 T 2 * 7 id 5 | +id $ | Reduce F→id. 1× |
| $0 T 2 * 7 F 10 | +id $. | Reduce. T→T*F 3×6 |
| | +id $ | reduce E→T 1×3 |
| $0 T 2 | +id $ | shift. |
| $0 E 1 | +id $ | Shift. |
| $0 E 1 + 6 | id $ | reduce F→id. |
| $0 E 1 + 6 id 5 | $. | reduce T→F. |
| $0 E 1 + 6 F 3. | $ | reduce E→E+T |
| $0 E 1 + 6 T 9 | $ | accept |
| $0 E 1 | $ | |

## canonical LR parsing:-

Similar to SLR parsing. Only in reduce operation.

---

CLR → canonical collection of

LR(1) items = LR(0) + look ahead.

| LR(0) | SLR(1) | CLR < LALR |
|---|---|---|
| reduce is written in full row | ⇃ reduce as written in follow of(P). | ⇃ Reduce is written only on look ahead. |

E→ BB
B→cB/d.

augument grammar along with look head items.

LR(0). E'→ . E, $→ look ahead item.
E→ .BB, $ → look ahead item
B→.cB/d, c/d → look ahead item
first of [B]

|  | c | d | $ |  | E | B |
|---|---|---|---|---|---|---|
| $r_0$ | $S_3$ | $S_4$ |  |  | 1 | 2 |
| 1 | | | accept | | | |
| 2 | $S_6$ | $S_7$ | | | | 5 |
| 3 | | | | | | 8 |
| 4 | $S_3$ | $S_4$ | | | | |
| 5 | $r_3$ | $r_3$ | | | | |
| 6 | $S_6$ | $S_7$ | $r_1$ | | | 9 |
| 7 | | | $r_3$ | | | |
| 8 | $r_2$ | $r_2$ | $r_3$ | | | |
| 9 | | | $r_2$ | | | |

## Error handling

1) important features — detect & report errors

→ Reporting errors in original source program rather than intermediate or final code

→ should not be complicated

→ should not be duplicate

→ localize the problem.

Strategies:

1) panic mode recovery

2) phrase level

3) Error productions

4) global correction

---

1) not specify ; y). EX: a+b = c;
                      d = e+f;

2) by ; delete entra ; Insert missing ;

3) good idea about common errors, find appropriate solution is stored.
   These productions detect the anticipated errors during parsing

   Ex  $E \rightarrow +E / -E / \; \ast E / /E$

4, global correction → incorrect input string required to transform x into y is

   expensive methods & not practically used

   costly in terms of time & space

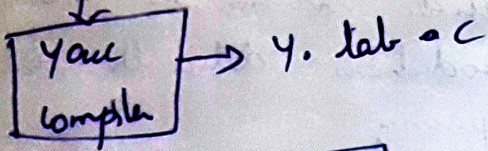YACC → automatic tool for generating the parser program

Yet Another Compiler Compiler.

lex → lexical analyzer generator

yacc → parser generator

⇓

It is a tool which generate LALR parser

[ Lex → LA →
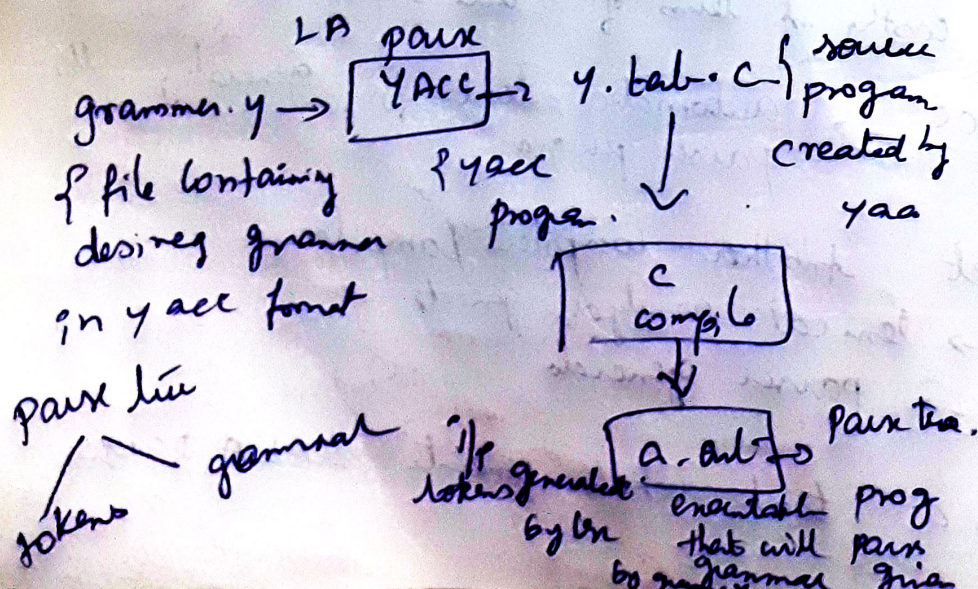[ YACC → Syntax Analysis → grammar.

Regular expression specifical

Yacc working

Step1:- yacc specification parser.y

```
┌─────────┐
│ Yacc    │ → y. tab . c
│ Compiler│
└─────────┘
```

Step  y. tab. c →
```
┌─────────┐
│ C       │ → a. out
│ compiler│
└─────────┘
```

Step:
i/p → a. out → o/p parser

tokens

LA   parse
grammar. y → ┌──────┐ → y. tab. c → source
            │ YACC │              program
            └──────┘     ↓        created by
{ file containing  { yacc          yaa
desired grammar    program.
in yacc format          ↓
                  ┌────────┐
parse lie         │ C      │
                  │ compile│
  ↗ grammar       └────────┘
 ↙                     ↓
tokens          i/p  ┌──────┐
          tokens generated │ a. out│ → Parse tree.
          by lex    executable  prog
                    that will parse
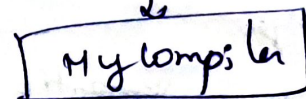                    by grammar gra
```

Syntax:
definitions { declaration of tokens
            types of values used.

Rules      { list of grammar rules ( grammar
            with semantic        rules.
                        routine.

./. /.

Supplementary code

lex / yacc.

```
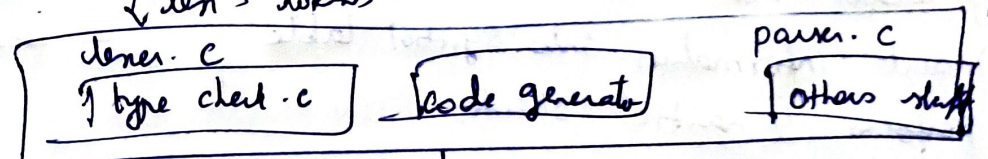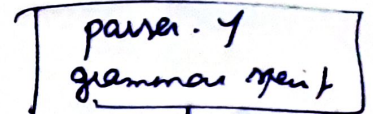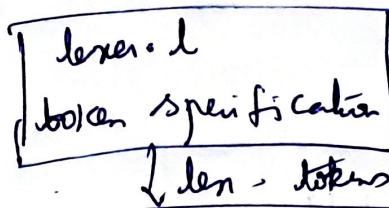┌──────────────────┐         ┌──────────────────┐
│ lexer. l         │         │ parser. y        │
│ token specification│        │ grammar specnt   │
└──────────────────┘         └──────────────────┘
        ↓ lex - tokens              ↓ yacc  parser
┌──────────────────┐                     grammar
│ lexer. c         │ ┌──────────────┐ ┌─────────────┐
│ type check. c    │ │ code generator│ │ parser. c   │
└──────────────────┘ └──────────────┘ │ others stuff│
                            ↓          └─────────────┘
                     ┌──────────────┐
                     │ My compiler  │
                     └──────────────┘
```

## Semantic rules:

1) generate code
2) insert information into symbol table
3) perform semantic check.
4) Issue error messages.

### Two notations of attaching semantic rules

1) Syntax Directed Definitions. High level specification hiding many implementation details (Attribute grammars)

2) Translation scheme. More implementation oriented. Indicate the order in which semantic rules are to be evaluated

---

Semantic rule        Production.

$E.code$   $E_1.code \,||\, T.code \,||\, '+'$    $E \rightarrow E_1 + T$

String val. attribute differetate.

We associate attributes to the grammar symbols representing the language constructs

values for attributes are computed by semantic rules associated with grammar productions

It is a generalization of context free grammar.

1) set of attributes
2) associated with semantic rules

Such formalism generates Annotated parse tree

each node of the tree is record with a field of each attribute.

$$L \rightarrow En \rightarrow \text{numerical value.}$$

$$F \rightarrow digit \quad F.val = digit.lexval$$

Numerical value of token returned by lexical analyser

$D \rightarrow TL$.    $L.inh = T.type$

$T \rightarrow int$    $T.type = integer$

$T \rightarrow float$    $T.type = float$

$L \rightarrow L_1, id$    $L_1.inh = L.inh$

$L \rightarrow id$    add.Type(id, entry, L.inh)

add Type(id, entry, L.inh)

## construction of Syntax tree:

1) mknode(op, left, right)

2) mk leaf (id, entry)

3) mk leaf(num, val)

$x * y$

$x$    $P_1 = mk leaf (id, ptr to entry x)$

$y$    $P_2 = mkleaf(id, ptr to entry y)$

$*$    $P_3 = mknode(*, P_1, P_2)$

$5$    $P_4 = mk leaf(num, 5)$

$-$    $P_5 = mknode(-, P_3, P_4)$

$z$    $P_6 = mkleaf(id, ptr to entry z)$

$+$    $P_7 = mknode(+, P_5, P_6)$



Pointer to Syntax for Z

P5

P3

P1  id  pointer to syntax for x

P2  id  pointer to syntax for y

num  5

* S attribute with synthesized attributes only

* evaluated using bottom up parse.

* purpose of stack is to keep track of values of the synthesized attributes associated with the grammar symbol, parse stack

Production  $X \rightarrow ABC$

$X \to ABC$, $X \cdot x = f(A, a, B \cdot b, C \cdot c)$

$L \to$ Attributed definition

$A \to X_1, X_2, X \dots X_n$.

$X_4$ is such that $1 \le K \le n$

$A \to X_1, X_2 \dots X_n$.

also depends upon the $X_1, X_2 \dots X_{j-1}$, to

the left of $X^n$

depends upon Inherited attribute A.

$A \to Pa$. $P \cdot in := P(A, in)$
$Q \cdot in := q(P, sy)$

---

Three address code :-

* Three address code is an abstract form of intermediate code that can be implemented as a record with the address fields.

* There are three representations used for three address code such as quadruples, triples and indirect triples.

Implementation of three address statements.

1, Quadruple
2. Triples
3. Indirect triples

Quadruple:

OP, arg1, arg2, result.

$x = y$ op $z$
$x = $ op $y$ (unary) } no arg 2
$x = y$. copy statement }

parameter $x \Rightarrow$ not arg 2 & result.

goto L unconditional Jump.

if $x$ relop $y$ goto L $\to$ conditional jump.

$a = b* - c + b * - c.$

$t_1 = $ uminus $c.$

$t_2 = b * t_1$

$t_3 = $ uminus $c$

$t_4 = b * t_3$

$t_5 = t_2 + t_4$

$a = t_5$

Quadruples structure

| | op | arg1 | arg2 | result |
|---|---|---|---|---|
| (0) | uminus | c | | $t_1$ |
| (1) | * | b | $t_1$ | $t_2$ |
| (2) | uminus | c | | $t_3$ |
| (3) | * | b | $t_3$ | $t_4$ |
| (4) | + | $t_2$ | $t_4$ | $t_5$ |
| (5) | = | $t_5$ | | a |

Triples structure :- (not using temp variables)

op, arg1, arg2

| | op | arg1 | arg2 |
|---|---|---|---|
| (0) | uminus | c | |
| (1) | * | b | (0) |
| (2) | uminus | c | |
| (3) | * | b | (2) |
| (4) | + | (1) | (3) |
| (5) | = | a | (4) |

Indirect triples

pointers to the triple structure.

| | | | op | arg1 | arg2 |
|---|---|---|---|---|---|
| (0) | 40 | (40) | uminus | c | |
| (1) | 41 | (41) | * | b | (40) |
| (2) | 42 | (42) | uminus | c | |
| (3) | 43 | (43) | * | b | 42 |
| (4) | 44 | (44) | + | (41) | (43) |
| (5) | 45 | (45) | = | (a) | (44) |

There are two ways to store two dimensional array in memory.

(i) Column major order
(ii) row - major order.

Two dimensional 3x4 array.

$$A = \begin{array}{c} 0 \\ 1 \\ 2 \end{array} \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \end{bmatrix} \; 3\times4.$$

(i) row major order of elements are stored in memory row by row )

| $a_{11}$ | $a_{12}$ | $a_{13}$ | $a_{14}$ | $a_{21}$ | $a_{22}$ | $a_{23}$ | $a_{24}$ |
|---|---|---|---|---|---|---|---|

[ first row ] [ second ro. ]

(ii) Column major order { elements are stored in memory column by column.

| $a_{11}$ | $a_{21}$ | $a_{31}$ | $a_{12}$ | $a_{22}$ | $a_{32}$ | $a_{13}$ | $a_{23}$ | $a_{33}$ | $a_{14}$ | $a_{24}$ | $a_{34}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|

1st column    2nd col.

The address of element first row & first column $A \cdot [a_{11}]$

base address of $A$ :- { base(A) }

## column major order.

address of $A[j,k] = Bax(A) + W[m(k-1) + j-1]$

$$\begin{bmatrix} 2 & 3 & 4 \\ & & \\ & & \\ & & \end{bmatrix} \begin{matrix} 1 \\ 2 \\ 3 \end{matrix}$$

$W$ = Size of element = datatype size.

$m$ = no. of rows.

## row major order

address of $A[j,k] = Base(A) + W[n(j-1) + (k-1)]$

$n$ = no. of columns.

eg

## Single dimensional array.

Single.

Arrays

two { column MO / row MO

int $A[5]$;

$A[3] = 20$;

$Add(A[3]) = 10 + 3 * 2$



$$\begin{array}{cccccc} & 0 & 1 & 2 & 3 & 4 \\ & & & & & \end{array}$$ linear.

$lo$  A  10 12 14 16 18

base add → lo.

{ add $a[i]$ = bax add + i * Size of data type}

$= lo + i * W$

$= 10 + 3 * 2$

$= 10 + 6 = 16$

---

Some compiler will start the index from 1 instead of 0.

$[1...5]$ as integer.  A



10 12 14 16 18

$A[4] = 30$

add $[A[i]] = Lo + (i-1) * W$

$= 10 + (4-1) * 2$

$= 10 + 3 * 2 = 16$.

## 2D Arrays. { In 2D Array compiler can follow row major / column major mapping}

eg  int $A[3][4]$:

$$A \begin{bmatrix} & 0 & 1 & 2 & 3 \\ 0 & 100 & 102 & 104 & 106 \\ 1 & 108 & 110 & 112 & 114 \\ 2 & 116 & 118 & 120 & 122 \end{bmatrix}$$

$3 \times 4$

Array always created is single dimension, in memory. because memory is storing in linear, & to having single integer address.

row major

| $a_{00}$ | $a_{01}$ | $a_{02}$ | $a_{03}$ | $a_{10}$ | $a_{11}$ | $a_{12}$ | $a_{13}$ | $a_{20}$ | $a_{21}$ | $a_{22}$ | $a_{23}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 100 | 102 | 104 | 106 | 108 | 110 | 112 | 114 | 116 | 118 | 120 | 122 |

$A[2][1] = [Lo + 2 * 4 + 1] * 2$

$= 1004 + a * 2$.

add $[A[i][j]] = lo + [i * n + j] * W$ → indicates 0.

$= 100 + 2 * 4 + 1 * 2$.

$= (100 + 9 * 2 = 100 + 18 = 118$.

$$A[i][j] = L_0 + [(i-1) * n + (j-1)] * w \quad \text{indicate}$$

## Column major order:

$$a_{00} \ a_{10} \ a_{20} \ | \ a_{01} \ a_{11} \ a_{21} \ | \ a_{02} \ a_{12} \ a_{22} \ | \ a_{03} \ a_{3}$$
$$100 \quad 102 \quad 104 \ | \ 106 \quad 108 \quad 110 \ 112 \quad 114 \ 116 \ 118$$

### formula

$$Add[A[i][j] = L_0 + [j * m + i] * w$$

$$Add[A[i][j] = L_0 + [(j-1) * m + (i-1)] * w .$$

$$Add(A[i][2]) = L_0 + [2 * 3 + 1] * 2$$

$$= 100 + 7 * 2 = 114$$

address of particular element in memory

UNIT-4   ## Simple Code Generator:-

1) Generates target code for a sequence of three address statements.

2) For each operator is a statement, there is a corresponding target language operator.

## Register & Address Descriptors:

### Register Descriptors:-

Keeps track of what is currently is each register.

- Initially all registers are empty.  $\boxed{R_0}$

2, Address Descriptor:- Keeps track of the location where the current value of the name can be found.

- Location may be register, a stack location or memory address. [ $\boxed{a}$

## A Code Generation Algorithm:-

For each three address statement of the form $x = y \ op \ z$ .

(i) Invoke a function get reg to determine the location L, where result of. y op z shall be stored.

getreg $\Rightarrow$ empty register & name can be stored in memory

(name not currently used)   occupied register → memory location [.L]   L.

2, Consult address descriptor for y to determine y, the current location of y. If y is not already in L, generate Mov y, L:

$y < \binom{mem}{reg}$   $x = y + z$.   $L = R_0$.

Mov y, R₀.

3, Generate the instruction OP z, L. update address descriptor of x to indicate that x is in L. If x is in register, update its descriptor, to indicate that it contains the value of x.

$x = y + z$   $L = R_0$
$R_0$.
$y + z$

```
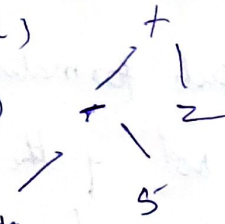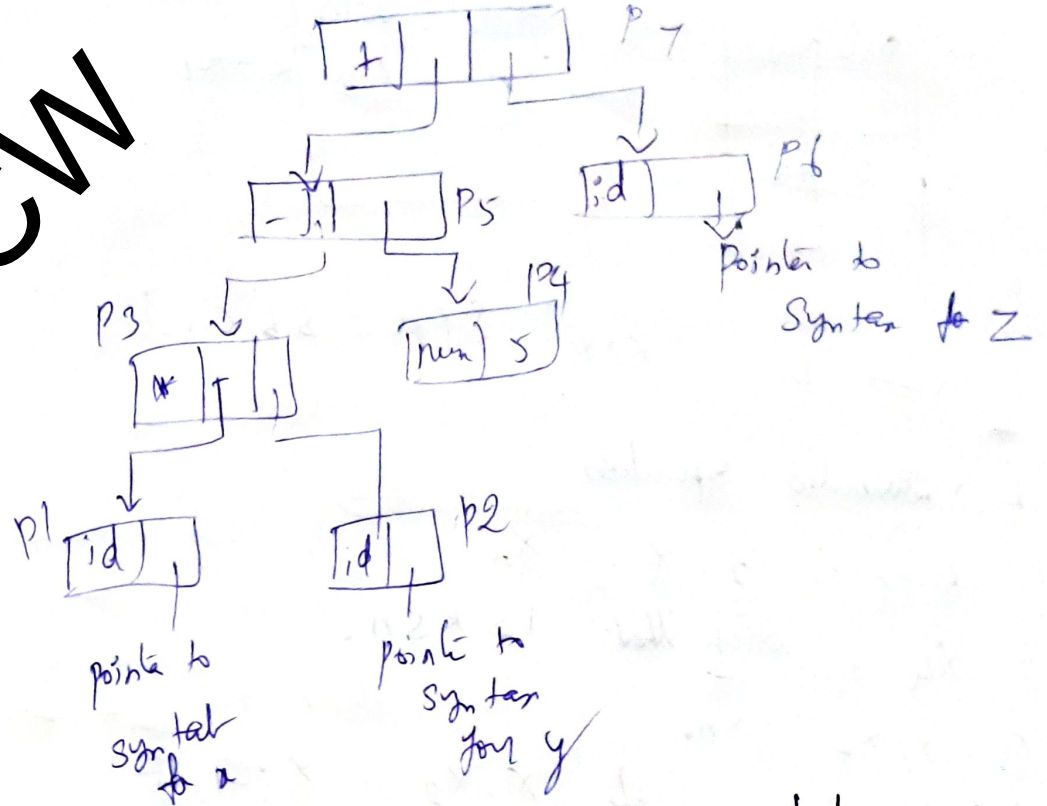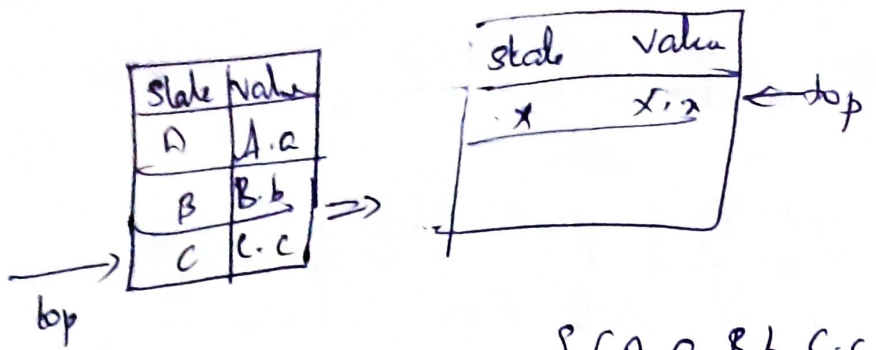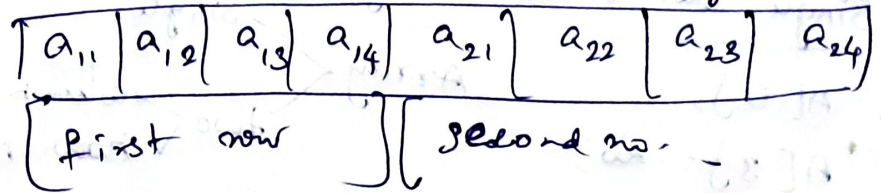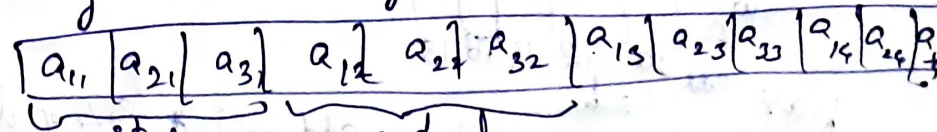Mov  y,  R₀
ADD  z, R₀
```

If y & z have no next uses and not live on emit, update the descriptors to remove y & z

x is live on emit Mov R₀, x.

---

example:-
(a-b) + (a-c) + (a-c).

Three address code sequence

$t_1 = a - b$
$t_2 = a - c$
$t_3 = t_1 + t_2$
$d = t_3 + t_2$
  ↳ live on emit

| Statements | Code generated | Register Descriptor | Address Descriptor |
|---|---|---|---|
| | | Register are empty. | t₁ in R₀ |
| $t_1 = a - b$ | Mov a, R₀ | R₀ contains t₁ | |
| | SUB b, R₀ | | |
| $t_2 = a - c$ | Mov a, R₁ | R₀ contain t₁ | t₁ in R₀ |
| | SUB c, R₁ | R₁ contain t₂ | t₂ in R₁ |
| $t_3 = t_1 + t_2$ | ADD R₁, R₀ | R₀ contains t₃ | t₂ in R₁ |
| | | R₁ contains t₂ | t₃ in R₀ |
| $d = t_3 + t_2$ | ADD R₁, R₀ | R₀ contains d | d in R₀ |
| | Mov R₀, d | | d in R₀ and memory. |

## Symbol table:-

1) Symbol table are data structures that are used by compilers to hold information about source program constructs.

→ It is used to store information about the

occurence of various entities such as,
objects, classes variable name, function etc.
It is used by both [Analysis phase] & [Synthesis phase]
Front end.                     Bad end

LA  SA  SA  IC        IC  CO  CG

The symbol table used for following purposes:

1) It is used to store the names of all entities in a structured form at one place.

2) It is used to verify if a variable has been declared.

3) It is used to determine the scope of a name.

4) It is used to implement type checking by verifying assignments & expressions in the source code are semantically correct or not.

A symbol table can either be linear or hash table.

It mains the entry for each name as
< symbol name, type, attribute >.

eg    < static, int, Raday >.
Symbol table stores an entry in this format.

Use of Symbol table:-

1) A symbol table information is used by the analysis and synthesis phases.

2) To verify that used identifiers have been defined (declared).

3) To verify that expressions and assignments are semantically content - type checking.

4) To generate intermediate or target code.

# UNIT-5   Directed acyclic graph

## DAG - optimization of Basic Block
(transformation on Basic blocks)

1) Structure preserving transformation.

2) Algebric transformation

## I. Structure preserving transformation.

1) **Common subexpression elimination**

if it Previous computed, value cannot be changed ($\times$ for computing)

$$x + y - w$$

$$
\begin{cases}
x = y+z \\
y = x-w \\
z = y+z \\
w = x-w
\end{cases}
\Rightarrow
\begin{aligned}
x &= y+z \\
y &= x-w \\
z &= y+z \\
w &= y+z
\end{aligned}
\quad
\begin{aligned}
& y+z-w \\
& y+z-y+
\end{aligned}
$$

2) **Dead code elimination.**

$$
\begin{aligned}
a &= a+2 \Rightarrow \\
b &= b+c \\
&\downarrow \\
b &= b * c \\
e &= b+2
\end{aligned}
\Rightarrow
\begin{aligned}
b &= b+c \\
b &= b * c \\
c &= b+2.
\end{aligned}
$$

3) **Renaming of temporary variables:**

$$
\begin{aligned}
t_1 &= x * y \\
t_2 &= z - t_1 \\
t_1 &= t_1 * w \\
w &= t_2 + t_1
\end{aligned}
\Rightarrow
\begin{aligned}
t_1 &= x * y \\
t_2 &= z - t_1 \\
t_3 &= t_1 * w \\
w &= t_2 + t_3.
\end{aligned}
$$

4) **Interchanging of statements.**

$$
\begin{aligned}
t_1 &= x * y \\
\begin{cases} t_2 = z - t_1 \\ t_3 = t_1 * w \end{cases} \\
\text{independent} \quad w &= t_2 + t_3.
\end{aligned}
\qquad
\begin{aligned}
t_1 &= x * y \\
t_3 &= t_1 * w \\
t_2 &= z - t_1. \\
w &= t_2 + t_3.
\end{aligned}
$$

## Algebric transformation:

Basic block, completed elimination.

$$x = x + 0. \quad \times$$
$$x = x - 0. \quad \times$$
$$a = a * 1$$
$$b = b/1 \quad \times$$

$$c = d ** 2 \Rightarrow pow(d, 2)$$
$$\Downarrow$$
$$c = d * d, \quad \text{for improving the basic blocks}$$

# Global Data flow Analysis:

→ To do code optimization & code generation.

Compiler:- Collect information about the whole program & distribute it to each block in the flow graph.

## Data flow equation:

Out [S] = gen [S] ∪ [In [S] − kill [S]]

Out [S] = Info at end of S

gen [S] − Info generated by S

in [S] − Info enters at the beginning of S

kill [S] → Info killed by S

## points & paths:-

Within B, there is a point between 2 adjacent statements B₁ → 4 points.

A path from $p_1$ to $p_n$ is a sequence of points $p_1, p_2 \dots p_n$ such that for each i; between 1 and n−1, a the

(i) $p_i$ → point preceding the statement

$p_{i+1}$ → $p_i$ → end of block.
$p_{i+1}$ beginning of the successor block

---

# Reaching definition:

A definition of variable x is a statement that assigns a value of x.

A definition of d reaches a point p, there is a path, from d do p, such that d is not killed along the path



## Data flow analysis of structural program. if.

do while.

```
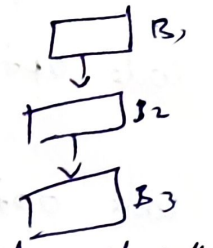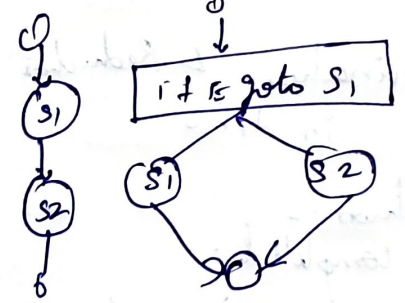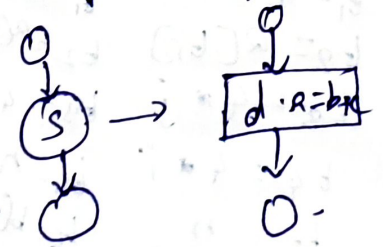a = 3
b = a+2
a = x +y
c = a +2
```



## Data flow equation for reaching definition.



gen [S] = {d}

kill [S] = $D_a$ − {d}

Out [S] = gen[S] ∪ in [S] − kill[S]]

gen[S] = gen [S₂] ∪ ( gen(S₁) − kill [S₂] ).

kill [S] = kill [S₂] ∪ ( kill[S₁] − gen[S₂] )

in [S₁] = in[S]

in [S₂] = out [S₁]

out [S] = out [S₂]

Principal sources of optimization:

code optimization: Improves the intermediate code
→ Less space & time

code optimization Techniques
1. Common subexpression elimination
2. Constant folding        * preserves the
3. copy propagation        meaning of the program
4. Dead code elimination
5. Code motion
6. Induction variable elimination & Reduction
                              in time.

Common Subexpression elimination:-
(i) if it was previously computed.
(ii) values of variables have not changed

$a = b + c$      $a = b + c$      $t_1 = 4 * i$      $t_1 = 4 * i$
$b = a - d$      $b = a - d$      $t_2 = a[t_1]$     $t_2 = a(t_1)$
$c = b + c$      $c = b + c$      $t_3 = 4 * j$      $t_3 = 4 * j$
$d = a - d$      $d = b$.         $t_4 = 4 * i$      $t_5 = n$
                                  $t_5 = n$          $t_6 = b(t_1)$
2) Constant folding:-                                $t_6 = b(t_4 + t_5)$. 

value of a expression is constant, use
the constant instead of expression

$$PI = 22/7, \quad \Rightarrow \quad 3.14,$$

Copy propagation:
$f = g$  use $g$ for $f$ after $f = g$

$x = a$              $x = a$
$y = x * b$  $\Rightarrow$  $y = a * b$
$z = x * c$          $z = a * c$

Dead Code elimination:
variable is live, if its value can be
used subsequently, otherwise it is dead

$x = a$              $y = a * b$
$y = a * b$          $z = a * c$.
$z = a * c$

Code Motion:
moves code outside a loop.

```
while (j < 10)
{  x = y + z;
   j = i + 1;
}
```
                    $x = y + z;$          loop
                    while (i < 10)        invariant
                    {  i = i + 1;         computation.
                    }

b) Induction variable Elimination & (loop control variable)
   Reduction in strength (Complex to simple)

```
i=1;                    t=4;
while (i<10)            while (t<40)
{ t = i*4;             { t=t+4;
  i=i+1;                }
}
}
```

Loop optimization:
 - Most execution time of a program is spent on loops.
 - Decreasing the number of instructions in an inner loop improves the running time of a program.

Loop optimization techniques,
1) code motion
2) Induction variable elimination & Reduction in Strength.

Loop unrolling:
   duplicates the body of the loop Multiple times, inorder to decrease the number of times the loop condition is tested

```
for (i=0; i<100; i++)        for (i=0; i<50; i++)
{ display();                 { display();
}                              display();
                             }
```

4) Loop Jamming:
   - combines the bodies of a adjacent loop that would iterate the same no. of times

```
for(i=0; i<100; i++)          for i=0; i<100; i++
  a[i]=1;                     { a[i]=1
for (i=0; i<100; i++)           b[i]=2
  b[i]=2;                     }
```

global

gen[s] = gen [s₁] ∪ gen [s₂]
kill [s] = kill [s] ∩ kill [s₂]

if    in[s₁] = in[s]
else  in [s₂]= in[s]
      Out[s] = Out[s₁] ∪ Out[s₂]

$$gen[S] = gen[S_i]$$
$$kill[S] = kill[S_i]$$
$$in[S_i] = in[S] \cup gen[S_i]$$
$$out[S] = out[S_i]$$