

Compiler → which translates the high-level code to machine-level-code

→ Source program → **Compiler** → Target program

If the transition of source program to Target program in 1 step, it becomes complex

→ so we have 6 phases

- ① Lexical analyzer
- ② Syntax analyzer
- ③ Semantic analyzer
- ④ Intermediate code generator
- ⑤ Code optimizer
- ⑥ Code generator

Lexical analyzer → is also called scanner
→ takes source program as the input

eg: "A statement in source program"

sum = sum + x;

is given to lexical analyzer

Lexical analyzer

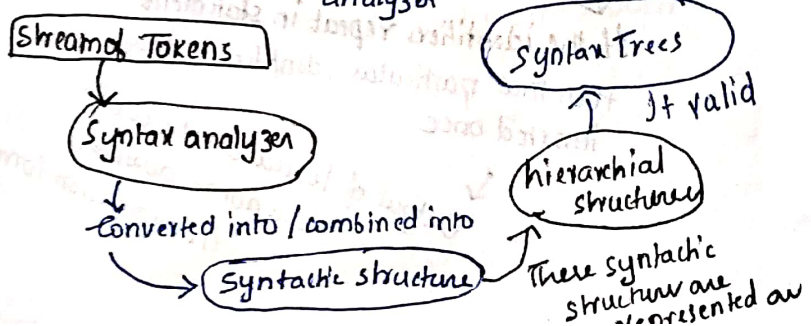
which divides all the statements into individual logical entities

- Tokens**
- Identifiers
 - Keywords
 - Operators
 - Constants
 - Punctuation marks
 - Special characters
 - Comments

So the output of lexical analyzer is stream of tokens

* In order to recognize the tokens within the source program a finite automata is implemented in the lexical analyzer

Syntax analyzer → takes stream of tokens as the input from the lexical analyzer



* syntax analyzer is also called parser.

* So syntax analyzer gives syntax trees as output

* syntactic trees are generated only when there is a valid syntactic structure

* we follow two parsing techniques to check the validity of syntactic structures

- ① Top-down parsing technique
- ② Bottom-up parsing technique (efficient)

Semantic analyzer → takes input from the Syntax analyzer: syntax trees

* functionalities of semantic analyzer

- ① Type checking → operator compatible
- ② uniqueness check → operands are present or not
- ③ Name-related check → we cannot declare more than one variable with same name under a data type.
- ④ flow control check → Case / Program / function /

In Ada, subprograms are represented using Procedure

Procedure name_of_procedure

Procedure nested_procedure

end nested_procedure

End name_of_procedure

It checks the, whether the mentioning of procedure name in end procedure is there or not

* output of semantic analyzer is Syntax tree (with type checking functionalities)

Intermediate code generator

→ Takes input as syntax trees from semantic analyzer

- * Intermediate code forms
 1. Polish Notation (Prefix) [Reverse postfix notation]
 2. Syntax tree
 3. Three Address code

* Output of Intermediate code generator is three address code, as it becomes easy to generate the code

Three address code

→ A statement which is having at most three operands. It can have any operators (any one) including the assignment operator.

eg: $x = y + z$; ✓
 $x = y$; ✓
 (no assignment) $x + y * z$ ✗

* we can convert any statement which is not in three address code to three address code using operator precedence rules & Associativity rules [manually].

eg: $x + y * z$ → But compiler uses internally a different technique "Syntax directed Translation"

$T_1 = y * z$
 $T_2 = x + T_1$ ✓

* Output of Intermediate code generator is Three address code

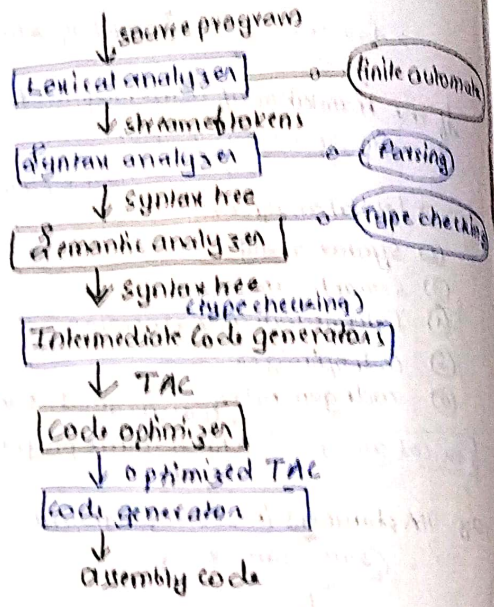
Code optimizer

→ takes Three address code as input

- * Code optimizer reduces the no. of instructions by therefore optimizing the code and occupying less space
- * There will be two code optimization techniques
 - ① machine - Independent code optimization.
 - ② machine - dependent code optimization.
- * Output of code optimizer is optimized three address code

Code generator

→ It takes optimized code as input and it produces the assembly code



► Assume a statement in the program is given to the lexical analyzer

positions = initial + rate * 60

↓ Then the lexical analyzer identifies each character & produces a stream of tokens using a finite automata

According to identifiers it stores identifiers & constants in symbol table

within the symbol table

mem. loc	ident
1	position
2	initial
3	rate
4	60

↓ + * already exists within the operator table

So the given statement is converted as

$idi := ider - ida * 60$

↓ If the identifiers repeat in statement then that particular identifier is only inserted once

↓ Output of lexical analyzer will not be available in expression form

And $id_1 := id_2 + id_3 * 60$ is given as input to syntactical analysis

which constructs context free grammar for each statement by syntax analyzer

then the CFG Generated will be

$S \rightarrow id := E$ $S \rightarrow id := E$
 $E \rightarrow E + E$ $E \rightarrow E + E / E * E / id / num$
 $E \rightarrow E * E$ (ambiguous grammar)
 $E \rightarrow id / num$

Then the syntactic structure is derived either by left most derivation or right most derivation

Therefore S, E are non terminals and deriving

$S \rightarrow id := E$
 $\rightarrow id := E + E$
 $\rightarrow id := id_2 + E$
 $\rightarrow id := id_2 + E * E$
 $\rightarrow id := id_2 + id_3 * E$
 $\rightarrow id := id_2 + id_3 * num$

Then the above syntactic structure is valid or not! : verified by two approaches
 Topdown approach & bottom up approach

Output of syntax analyzer is Syntax tree

Syntax tree: A tree is called syntax tree, if interior nodes are represented by operators, leaf nodes are represented by identifiers (or) constants

Construction of syntactic structure: syntax tree

* we convert

$id_1 := id_2 + id_3 * 60$ into postfix notation

$id_1 := id_2 + id_3 * 60$

$T_1 = id_3 * 60$

$T_2 = id_2 T_1 +$

$T_3 = id_1, T_2 :=$

$T_3 = id_1, id_2 T_1 + :=$

$T_3 = id_1 id_2 id_3 60 * + :=$ (Postfix)

and expression is evaluated using a stack

To generate tree we use 3 functions

- ① mknode (op, LP, RP) : for operators
 - ② mkleaf (id, entry) : for identifiers
 - ③ mkleaf (num, value) : for constants
- entry: is place/mem. loc where identifiers are stored.

$\therefore id_1 id_2 id_3 60 * + :=$

* $mkleaf(id, 1)$ creates a node and returns a address

$\therefore P_1 = mkleaf(id, 1)$

Similarly

$P_2 = mkleaf(id, 2)$

$P_3 = mkleaf(id, 3)$

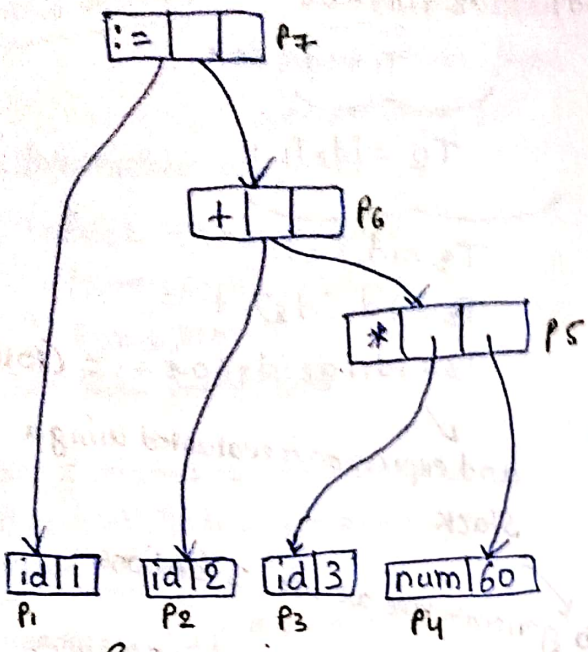
$P_4 = mkleaf(num, 60)$

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

$P_6 = mknode(+, P_2, P_5)$

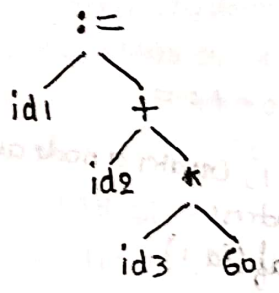
$P_7 = mknode(:=, P_1, P_6)$ Consider

we use functions in ordered postfix notation



Syntax tree.

Abstract syntax tree: A syntax tree is called abstract syntax tree, if we represent corresponding symbol at each node.



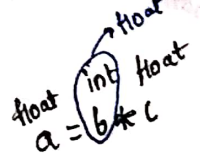
This is given as an input to semantic analyzer.

It verifies operator compatible operands are applied or not (Type checking)

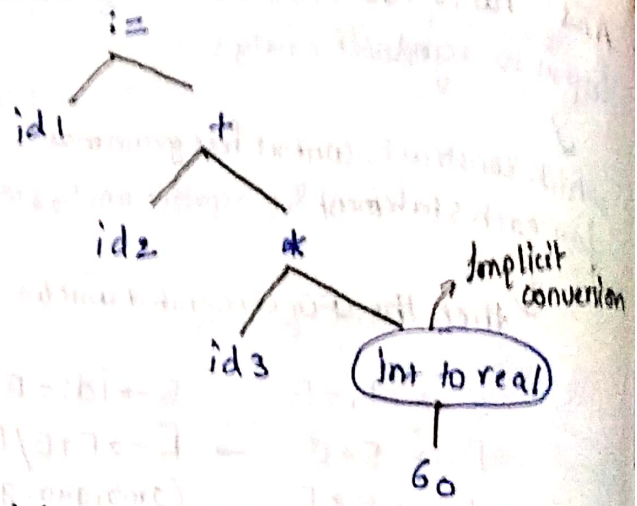
within the output of semantic analyzer, we effect the implicit conversions.

Taking the statement and reflecting implicit conversions

where lower data type calculations are converted into higher datatype automatically

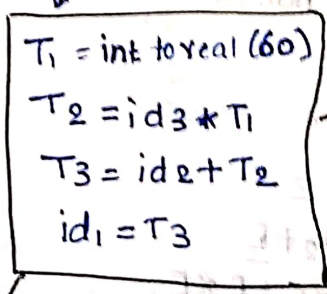


As there are no mnemonics in assembly language to perform operations on mixed mode expression



This is given as Output to Intermediate code generator.

Then it should be converted into three address code (Intermediate code)

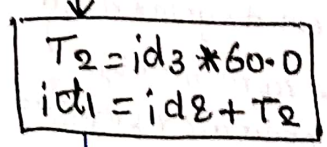


we cannot write $T_1 = id_3 * \text{int to real}$ because is itself is considered as an operator

This three address code is given as an input to the code optimizer.

Type casting has the highest precedence when compared to other operators

Then it carries out the conversions and replace with the converted value



This optimized code is given as input to the code generator.

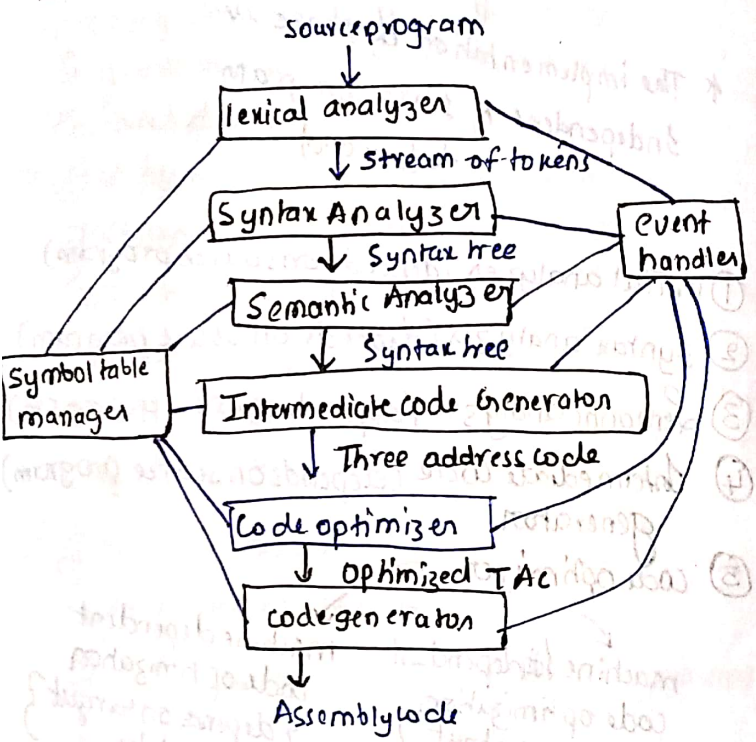
And it produces an Assembly code


```

MOVF id3, R1
MULF #60.0, R1
MOVF id2, R2
ADD R2, R1
MOVF R1, id1

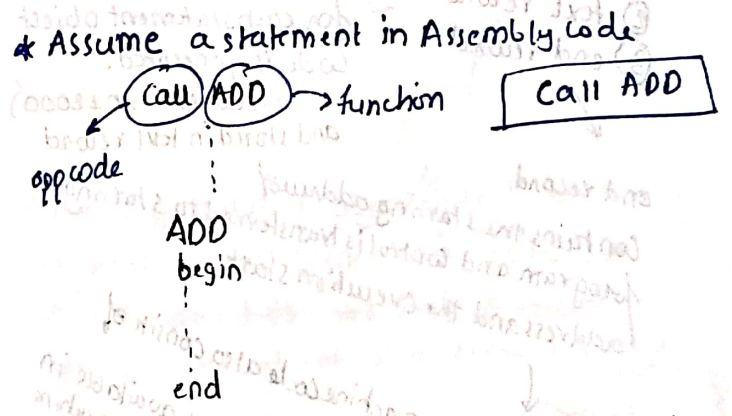
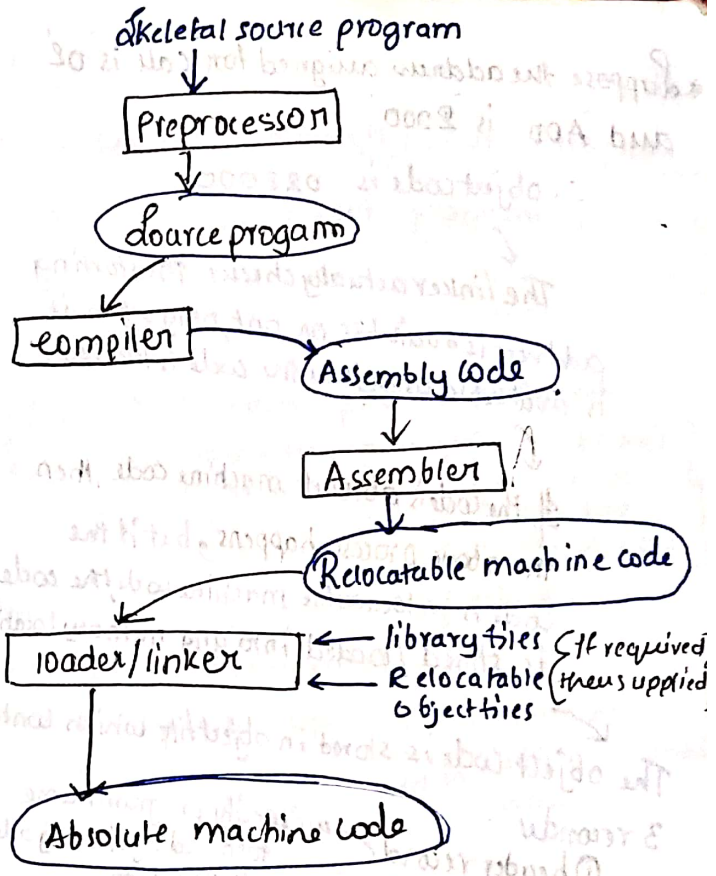
```

logical phases of compiler



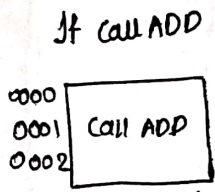
Symbol table manager: It is a datastructure which is used to store the user defined information like identifiers and constants.

- > If Identifier is a variable, then it stores the information like name, address, value, type, lifetime and scope.
- > If identifier is an array, it stores information like name of array, type of array, index type, range of array.
- > If identifier is a function, it stores the information like name of the function, no. of parameters, return type of the function.



- * for each statement address is allocated and object code is generated
- * If the starting address is defined then it is taken as starting address otherwise it takes '0' as starting address.

- > Call ADD is second statement, then the assembler calculates the length of instructions on memory locations of call ADD and takes the next address and assigns it to the second statement.



then the second statement starts occupying from 0003
 Assembler also generates the object code for each statement.

Suppose the address assigned for call is 02 and Add is 2000

∴ object code is 022000

The linker actually checks the starting address is available or not, only when it is available only then the code is loaded.

If the code is absolute machine code, then the above process happens, but if the code is relocatable machine code, the code is stored/loaded into any memory location.

The object code is stored in object file which contains 3 records

- ① header record → we specify program name followed by starting address of program
 - ② text record → for each statement object code is generated.
 - ③ end record → (eg: call Add → 022000) and stored in text record
- end record contains the starting address of program and control is transferred to starting address and the execution starts

Relocatable machine code also consist of extra record

- ④ modification record → not available in absolute machine code.

If the specified starting address is unavailable, then some other free memory locations are used to store the code and this information is stored in modification record.

loader & linker actually can check which memory locations are free

Grouping of phases

front end & backend pass

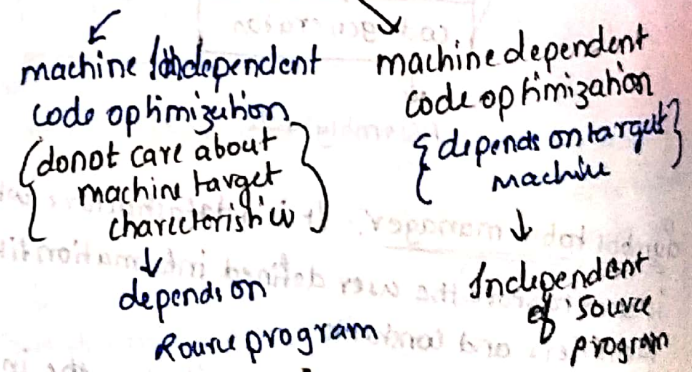
* The Implementation / Phase depends on the source program & independent of Target program

front end

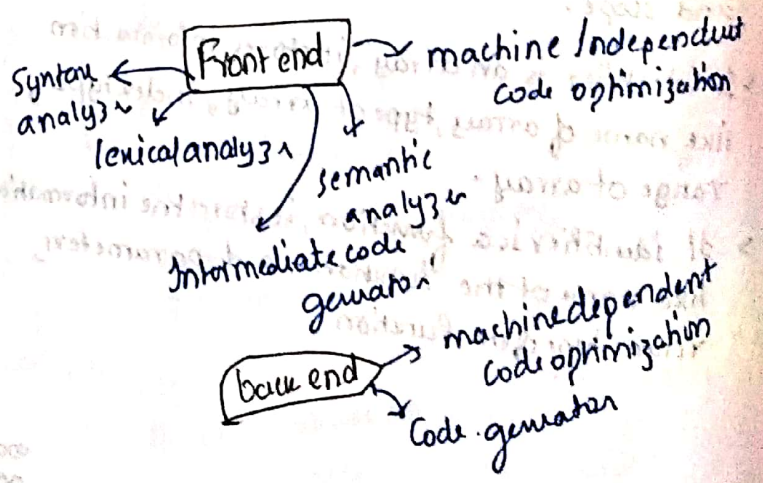
* The implementation of phase is independent of source program

back end

- ① lexical analyzer (depends on source program)
- ② Syntax analyzer (depends on source program)
- ③ Semantic analyzer (depends on source program)
- ④ Intermediate code generator (depends on source program)
- ⑤ code optimizer

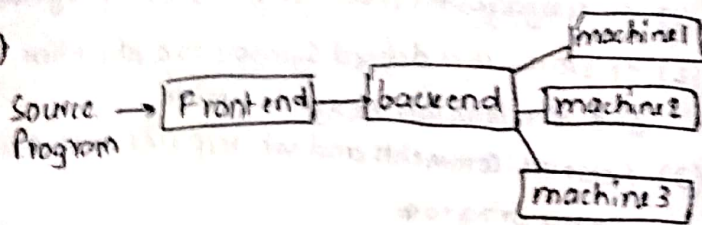


⑥ Code generation (depends on target machine)



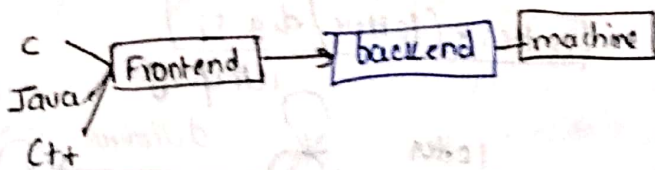
Advantages of grouping phases

(1)



- > Supposing a new programming language is being developed. Source program is given as input to front-end of the compiler design.
- > But by designing compiler; using grouping of phases. front end maybe similar to another already developed programming languages and you can include its API.
- > And the backend can be taken good care for designing it.
- > If there are no grouping of phases into front end & backend, we need to design compiler (all phases) for heterogeneous machines.

(2) By grouping of phases, for similar programming languages you can use the same front end.



and sometimes you can use the same backend too.
And this is very tedious to develop some compiler for similar programming languages.

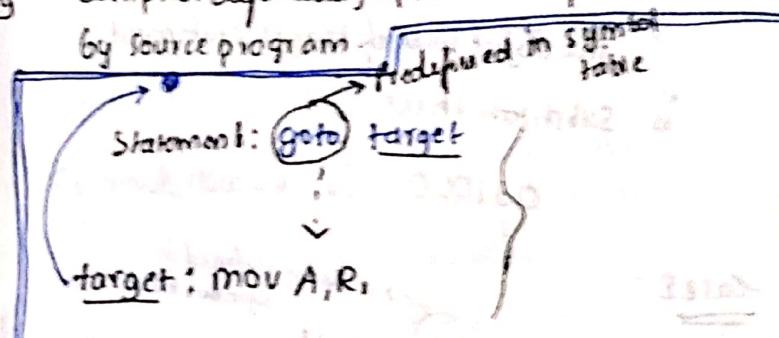
Pass

* It is a physical scan over the source program i.e., number of times a source program is scanned.

(0n)

* grouping one or more phases is called Pass
 { can be explained using legacy facts (128 KB)
 & compiler needs to be loaded into memory
 compiler size > ram size, therefore they grouped one or more phases into pass }

> Implementation of the phases, if whatever it contains one pass or multi pass compiler depending upon features provided by source program.



> Suppose goto has an object code is 03
 > Target statement comes after goto, so compiler doesn't know the object code of target statement yet. stores target in symbol table.
 e.g.: 1000 is assigned Address for target

∴ object code = 031000

> If target statement is placed before goto statement (forward reference) then object code = 03 []
 ↓ after slot is assigned in this case
 As target comes before goto, it places target statement in symbol table

target	2000

X Wrong { vice versa }

Case 1:

target statement: `mov A, R1`;

goto statement: `goto target`;

→ already defined in symbol table

In here target statement is before the goto statement.

Consider object code of statement `mov A, R1` is 1000

and it is stored in the symbol table next statement is

`goto target`;

and object code of above statement is 03

∴ the object code of target statement is substituted here

031000 (backward reference)

Case 2:

goto statement: `goto target`

target statement: `mov A, R1`

→ already defined in symbol table

In here target statement is after the goto statement

Consider object code of goto statement is 03

The object code of goto statement is 03

and object code of target is not yet computed

hence a slot is created 03

(Forward reference)

← after target object code is calculated the value is put in slot

Role of lexical analyzer

- (1) It recognizes the tokens in the source program
- (2) It stores user defined symbols like identifiers and constants into the symbol table
- (3) Removes comments and whitespaces within the source program.
- (4) Keeps track of the line number to associate with every message
- (5) Identifies the lexical errors

Token It describes class or category of a input state
ex: identifiers, constants, keywords, operators, comments etc.

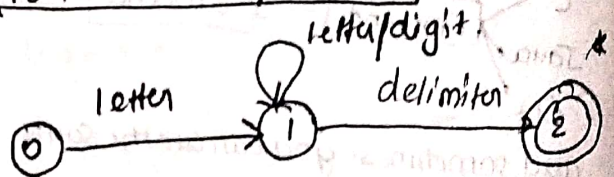
Pattern: A set of rules that describes a token

Note: Patterns are represented using regular expressions.

lexeme: Sequence of characteristics that is a form a input string which matches with the pattern of a token

Regular expression for identifier

$\text{letter} \cdot (\text{letter} / \text{digit})^*$



* when actually reading the identifier an extra character (as for =) and we perform retract operation.

* Each state is considered as a procedure to write the code.

Code for the state 0:

1. $c = \text{getchar}();$
2. IF LETTER(c) then
3. goto state 1; Else FAIL(c);

This does not generate error immediately, but sends control to the next finite automata

Code for the state 1:

1. $c = \text{getchar}();$
2. If LETTER(c) OR DIGIT(c) then
3. goto state 1;
4. Else if DELIMITER(c) then
5. goto state 2;
6. Else Fail(c);

Fail() function does not immediately generate the lexical error but transfers the control to the next finite automata.

Code for the state 2:

1. RETRACT();
2. return (id, Install());

Do the RETRACT() function eliminates extra space from the identifier

it is used to store the recognized identifier within the symbol table.

ex: If $\text{sum} = \text{sum} + x;$

In symbol table and location is allotted to it

sum	

Symbol table

then Install function returns
return (id, 1);

TOKEN	Value	Code
BEGIN	1	-
END	2	-
IF	3	-
THEN	4	-
ELSE	5	-
identifier	6	Pointer to symbol table
constant	7	Pointer to symbol table
<	8	1
<=	8	2
=	8	3
<>	8	4
>	8	5
>=	8	6

* Syntax Analyzer identifies the Tokens using their corresponding values.

* when lexical analyzer identifies the token it returns the equivalent integer values to the syntax analyzer

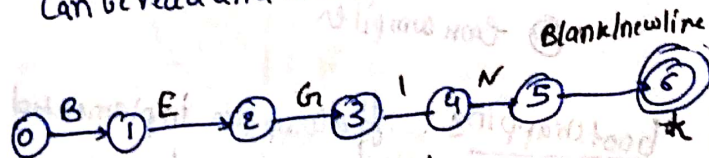
* The values may vary for any token.

∴ for an identifier it returns (6, 1);

Finite automata for begin

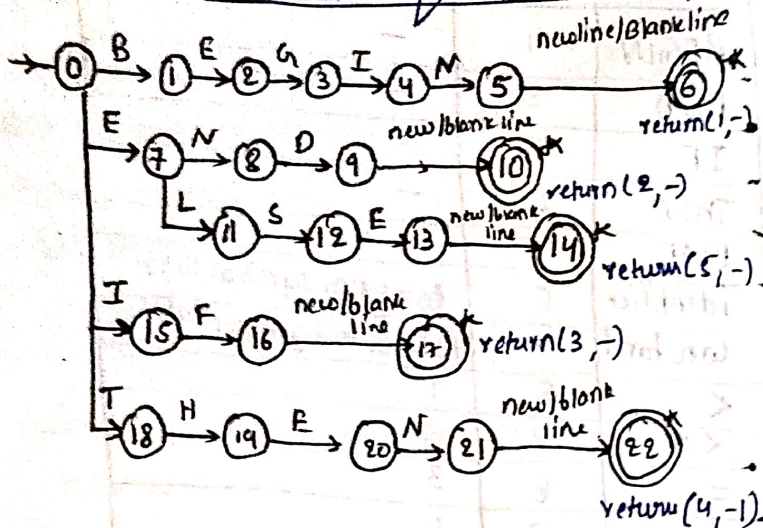
- > state 0: B
- > state 1: E
- > state 2: G
- > state 3: I
- > state 4: N

after state 4 any other character (if present) can be read and makes it as an identifier



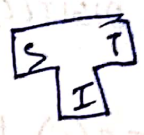
If the control does not reach final state then it calls fail function and it moves control to next finite automata
* → represents retract operation

Similarly drawing the finite automata for others (keywords)



- ① source program.
- ② Implementation language
- ③ Target language

* Characteristics of compiler represented in the



(or) can also be represented like this

C S I T

> First we design the compiler for subset features of L (language - source), as it is hard to develop a compiler to contain features of language, so we take common features

C S A A

A → machine

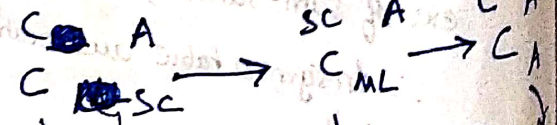
> Now we design compiler for language L

$C \xrightarrow{L} C \xrightarrow{S} C \xrightarrow{A} C \xrightarrow{A} C \rightarrow \text{object code}$

ex: C A → machine language

we take S C C

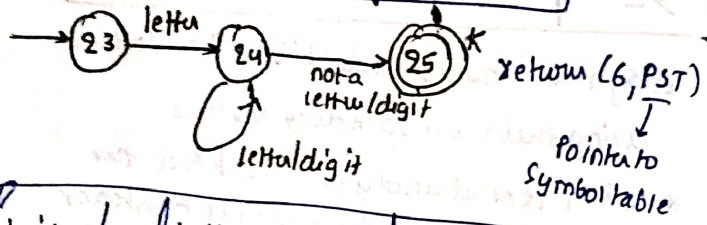
subset features etc



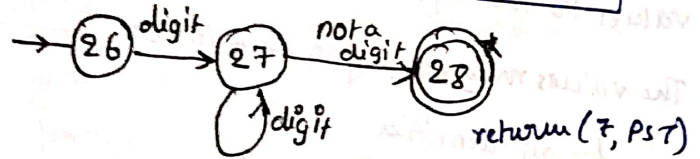
(Compiler designed for machine language for machine A for subset features of C)

(Compiler designed for machine language for machine A for subset features of C)

Similarly finite automata for identifiers



Similarly finite automata for constants



- Input Buffering
- lexical analysis in grammar
- Regular definition

Bootstrapping:-

Consider a language L that don't have any compiler available in any machine.

So we use 2 methods

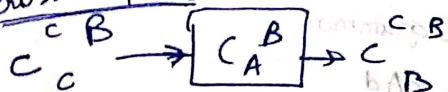
- ① Bootstrapping
- ② Cross compiler

bootstrapping: If a compiler implemented in its own language is called bootstrapping.

Cross compiler :- A compiler written on one machine produces object code for another machine.

Compiler designed for machine language for machine A for subset features of C

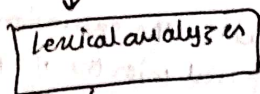
cross compiling



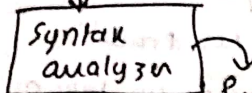
II Phase of compiler

> Imagine a statement is given as an input to lexical analyzer

sum = num + x ;



↓ stream of tokens



↓
Syntax analyzer combines stream of tokens into syntactic structures.

In order to convert the stream of tokens into syntactic structures it uses automata
(CA) CFA is defined

let s is a statement

$$S \rightarrow id = E$$

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E * E$$

$$E \rightarrow E / E$$

$$E \rightarrow id$$

for each program statement a CFA is defined.

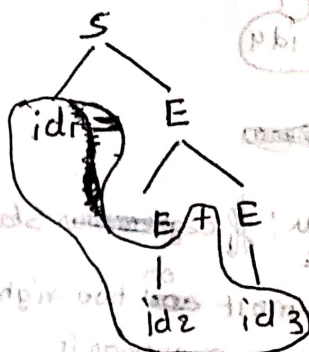
using the CFA we must derive below result
(id₁ = id₂ + id₃)
using leftmost & rightmost derivation

$$S \rightarrow id_1 = E \quad (S \rightarrow id = E) \quad \equiv$$

$$\rightarrow id_1 = E + E \quad (E \rightarrow E + E)$$

$$\rightarrow id_1 = id_2 + E \quad (E \rightarrow id)$$

$$\rightarrow id_1 = id_2 + id_3 \quad (E \rightarrow id)$$



$$id_1 = id_2 + id_3 \quad \&id_4$$

$$S \rightarrow id_1 = E \quad (\text{left-most})$$

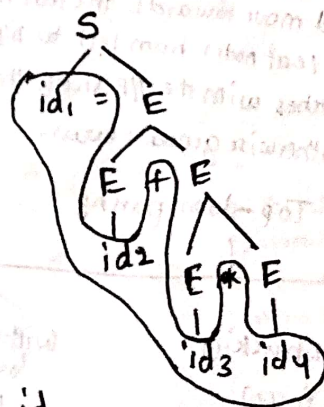
$$\rightarrow id_1 = E + E$$

$$\rightarrow id_1 = id_2 + E$$

$$\rightarrow id_1 = id_2 + E * E$$

$$\rightarrow id_1 = id_2 + id_3 * E$$

$$\rightarrow id_1 = id_2 + id_3 \&id_4$$



$$S \rightarrow id = E$$

$$S \rightarrow id = E + E$$

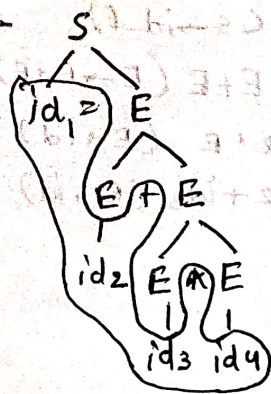
$$S \rightarrow id = E + E * E$$

$$S \rightarrow id = E + E * id_4$$

$$S \rightarrow id = E + id_3 * id_4$$

$$S \rightarrow id = id_2 + id_3 \&id_4$$

$$S \rightarrow id_1 = id_2 + id_3 \&id_4$$



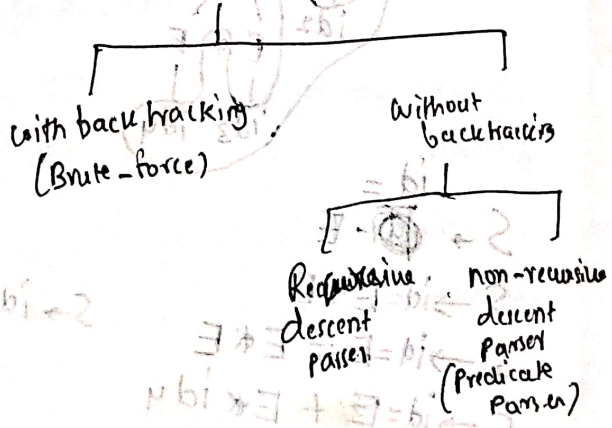
Ambiguous grammar:

contains two left-most most parse trees, then the grammar is called Ambiguous grammar

If there exist more than one parse trees then it is called Ambiguous grammar

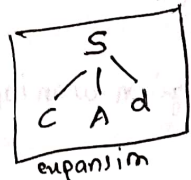
① Top-down parsing - It constructs the parse tree. It takes the starting non-terminal as root node and move towards the leaf nodes. Scan the leaf nodes from left to right. If it matches with the i/p string then it accepts otherwise generates error.

Top-down parsing



Continue the grammar

$S \rightarrow cAd$
 $A \rightarrow ab/a$

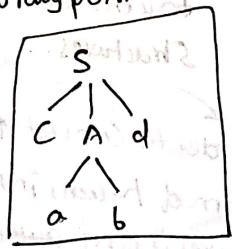


within the S we take an input pointer (pointing to c) and store it in temporary before expansion $\{IP = TP\}$

'c' in expansion is matched with i/p pointer $\{cAd\}$

> next symbol within the expansion is A and compared with the i/p pointer $\{cAd\}$

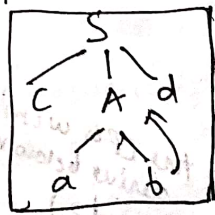
and A has two productions. before expansion we take an input pointer (pointing to a) and store it in temporary pointer.



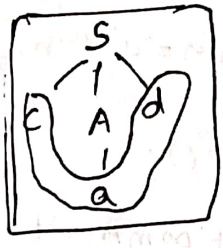
'a' in expansion matches with the i/p pointer \therefore i/p pointer is incremented

> and next time the match fails, hence we backtrack to the A (Previous). When we backtrack to A we restore the input pointer $\{IP = TP\}$

> expanding the a with second production



> then the next symbol is compared against the 'ip pointer' 'a'



and it's matching with the 'ip pointer' so it accepts the syntactic structure.

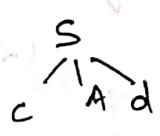
ex:-
 ⇒ Taking the string that needs to be validated **cad**

CFG = { cabd, cad }

> It starts with 's' starting symbol as root node in parse tree that needs to be validated

at before expansion we have **cad**
 { IP = TP }
 input pointer in an temporary pointer.

∴ expanding s using backtracking method



> 'input pointer' is matched with the 'c' in expansion

∴ 'ip pointer' is incremented

> so now TP points to A
 then for A { IP = TP } **cad**

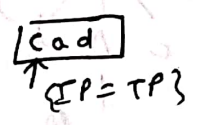
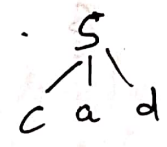
A has another productions



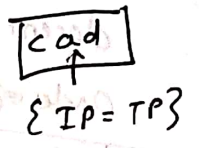
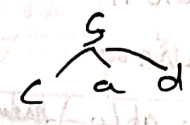
> It backtracks to root node, if there is no alternate production and ip pointer is restored

∴ ip pointer points to 'c'

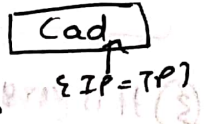
∴ Another production of s = **cad**



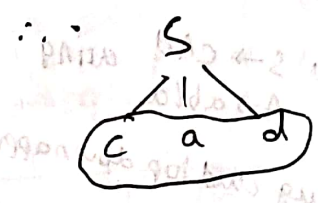
'c' matches with 'ip' ip pointer ++



'a' matches with 'ip' ip pointer ++



'd' matches with 'ip'

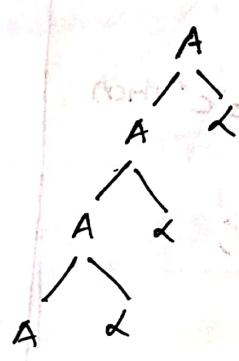


∴ it accepts

Disadvantages of backtracking :-

① Consider

$A \rightarrow A \ \& \ / \ B$



we only go to second production if we find unmatched result.

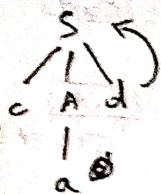
Infinite loop will occur in ~~backtracking~~ with backtracking and topdown parsing in left recursion.

So after eliminating left recursion

$A \rightarrow BA'$
 $A' \rightarrow \epsilon A' \mid B$

(reduced flat)

② $S \rightarrow CAD$
 $A \rightarrow a|ab$
 string: cabd



It should accept cabd, but it doesn't because it is indifferent order of productions

~~It (order $A \rightarrow ab|a$) then it accept both cad and cabd maybe true maybe not~~

③ It is very time taking process

⇒ For the grammar $S \rightarrow CAD$ using $A \rightarrow ab|a$

backtracking and top down approach code

each non terminal is implemented using a procedure. $\{S, A\}$ string: cabd

Procedure S ()

```

begin
  if input-pointer = 'c' then
    begin
      ADVANCE ( );
    end
  if .A ( ) then
    if input-pointer = 'd' then
      begin
        ADVANCE ( );
        return true;
      end
    end
  return false;
end
  
```

Backtrace method & recursive descent parser would not predict exact right hand side of production for a non terminal having more than one production while parsing a string but not - recursive descent parser will predict.

Procedure A ()

```

begin
  isave := input-pointer
  if input-pointer = 'a' then
    begin
      ADVANCE ( );
    end
  if input-pointer = 'b' then
    begin
      ADVANCE ( );
      return true;
    end
  end
  input-pointer := isave;
  if input-pointer = 'a' then
    begin
      ADVANCE ( );
      return true;
    end
  else return false;
end
  
```

Eliminate left recursion from the following grammar

$E \rightarrow E+T/T$
 $T \rightarrow T * F / F$
 $F \rightarrow (E) / id$

Sol:-

while $A \rightarrow A \wedge B$
 \downarrow
 $A \rightarrow \beta A'$
 $A' \rightarrow \alpha A' / \epsilon$

$E \rightarrow E+T/T$

$E \rightarrow TE'$
 $E' \rightarrow +TE' / \epsilon$

$(\langle A, B, \gamma \rangle) \in (N \cup \{ \epsilon \})^*$

//

$$T \rightarrow T * F / F$$



$$\begin{aligned} T &\rightarrow FT' \\ T' &\rightarrow *FT' / \epsilon \end{aligned}$$

after eliminating the left recursion

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' / \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' / \epsilon \\ F &\rightarrow (E) / id \end{aligned}$$

eliminate left recursion from the following grammar

$b_{exp} \rightarrow b_{exp} \text{ or } b_{term} / b_{term}$
 $b_{term} \rightarrow b_{term} \text{ and } b_{factor} / b_{factor}$
 $b_{factor} \rightarrow \text{not } b_{factor} / \text{true} / \text{false}$

eg: $b_{exp} \rightarrow b_{exp} \text{ or } b_{term} / b_{term}$

$$\begin{aligned} b_{exp} &\rightarrow b_{term} b_{exp}' \\ b_{exp}' &\rightarrow \text{or } b_{term} b_{exp}' / \epsilon \end{aligned}$$

$b_{term} \rightarrow b_{term} \text{ and } b_{factor} / b_{factor}$

$$\begin{aligned} b_{term} &\rightarrow b_{factor} b_{term}' \\ b_{term}' &\rightarrow \text{and } b_{factor} b_{term}' / \epsilon \end{aligned}$$

after eliminating

$$\begin{aligned} b_{exp} &\rightarrow b_{term} b_{exp}' \\ b_{exp}' &\rightarrow \text{or } b_{term} b_{exp}' / \epsilon \\ b_{term} &\rightarrow b_{factor} b_{term}' \\ b_{term}' &\rightarrow \text{and } b_{factor} b_{term}' / \epsilon \end{aligned}$$

$b_{factor} \rightarrow \text{not } b_{factor} / \text{true} / \text{false}$

multiple left recursions

$$A \rightarrow A \alpha_1 / A \alpha_2 / A \alpha_3 / \dots / A \alpha_n / \beta_1 / \beta_2 / \beta_3$$

$$A \rightarrow \beta_1 A' / \beta_2 A' / \dots / \beta_n A'$$

$$A' \rightarrow \alpha_1 A' / \alpha_2 A' / \dots / \alpha_n A' / \epsilon$$

eliminate left recursion

$$S \rightarrow A a / b$$

$$A \rightarrow A c / s d / e$$

$$A \rightarrow A c / s d / e$$

$$\begin{aligned} A &\rightarrow s d A' \\ A' &\rightarrow c A' / \epsilon \end{aligned} \quad \left| \quad \begin{aligned} A &\rightarrow e A' \\ A' &\rightarrow c A' / \epsilon \end{aligned}$$

$$S \rightarrow A a / b$$

$$A \rightarrow s d A' / e A'$$

$$A' \rightarrow c A' / \epsilon$$

Substituting the sd in $A \rightarrow A a / b$
it becomes $S \rightarrow S d a / b$

$$\begin{aligned} S &\rightarrow b S' \\ S' &\rightarrow d a S' / \epsilon \end{aligned}$$

To avoid the confusion, we use the algorithm.

Algorithm to eliminate left recursion

step 1: Arrange non terminals in some order such as A_1, A_2, \dots, A_n

step 2: for $i=1$ to n do
 begin
 for $j=1$ to $(i-1)$ do
 begin if $j < i$ then
 replace each production $A_i \rightarrow A_j \delta$
 with
 $A_i \rightarrow \delta_1 \delta / \delta_2 \delta / \dots / \delta_k$ where
 $A_j \rightarrow \delta_1 / \delta_2 / \dots / \delta_k$ are current
 A_j productions

end

eliminate immediate left recursion among A_i productions

end

end

Example $S \rightarrow Aa/b$
 $A \rightarrow Ac/sd/e$

Sol: rename non terminals with $A_1 \dots A_n$

$S \rightarrow A_1$
 $\downarrow \quad \downarrow$
 $A_1 \quad A_2$

\therefore Production:

$A_1 \rightarrow A_2 a / b$

$A_2 \rightarrow A_2 c / A_1 d / e$

do $j < i$

$\therefore A_2 \rightarrow A_1 d$ (satisfy the condition)

\therefore replace

~~$A_1 \rightarrow A_1 d$~~
 \Rightarrow replacing the A_j productions

$A_2 \rightarrow A_1 d$

$A_2 \rightarrow A_2 a d / b d$

\therefore production are

$A_1 \rightarrow A_2 a / b$
 $A_2 \rightarrow A_2 c / A_2 a d / b d / e$

$A_2 \rightarrow A_2 c$	$A_2 \rightarrow A_2 a d$
$A_2 \rightarrow A_2'$	$A_2 \rightarrow A_2'$
$A_2' \rightarrow c A_2'$	$A_2' \rightarrow a d A_2'$

\therefore eliminating immediate left recursion

$A_2 \rightarrow b d A_2'$

$A_2' \rightarrow a d A_2' / e$

$A_2 \rightarrow e A_2'$

$A_2' \rightarrow a d A_2' / e$

$A_2 \rightarrow A_2 c$ elimination

$A_2 \rightarrow b d A_2'$

$A_2' \rightarrow c A_2' / e$

$A_2 \rightarrow e A_2'$

$A_2' \rightarrow c A_2' / e$

\therefore production are:

$A_1 \rightarrow A_2 a / b$
 $A_2 \rightarrow b d A_2' / e$
 $A_2' \rightarrow a d A_2' / c A_2' / e$

\Rightarrow check ambiguity and write unambiguous

$R \rightarrow R + R / R \cdot R / R^* / a / b$

\circ deriving $a b a$

LMD

$R \rightarrow R + R$

$\rightarrow a + R$

$\rightarrow a + R \cdot R$

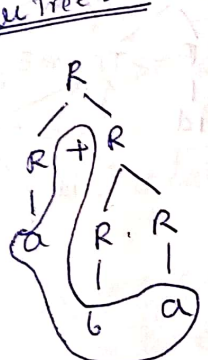
$\rightarrow a + b \cdot R$
 $\rightarrow a + b a$

LMD 2

$R \rightarrow R \cdot R$
 $R \rightarrow R + R \cdot R$
 $R \rightarrow a + R \cdot R$
 $R \rightarrow a + b \cdot R$
 $R \rightarrow a + b a$

∴ ambiguous, as there exist two left most derivations

Parse Tree LMP 1



Parse Tree LMD 2



∴ The unambiguous grammar is...

Operator precedence rule

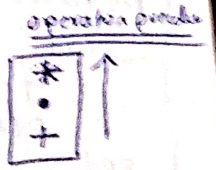
- ① Left Associativity → we have to define in left recursion
- ② Right Associativity → we have to define in right recursion

* the highest precedence operators are defined away from the starting non-terminal

* the lowest precedence operators are defined nearest to the starting non-terminal

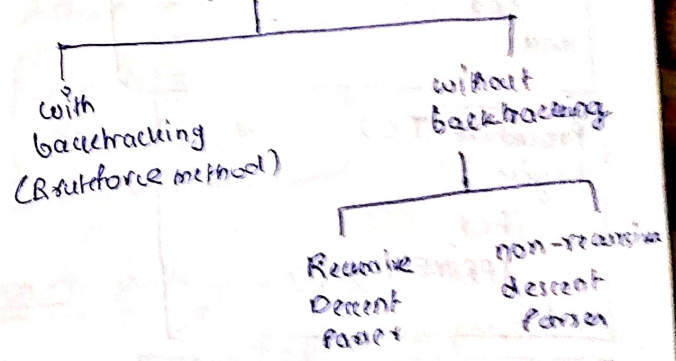
∴ modifying the grammar unambiguous grammar

$R \rightarrow R + T / T$
 $T \rightarrow \cdot T \cdot E / E$
 $E \rightarrow E * a / b$



* The below grammar was ambiguous, as the operator precedence cannot be defined clearly.

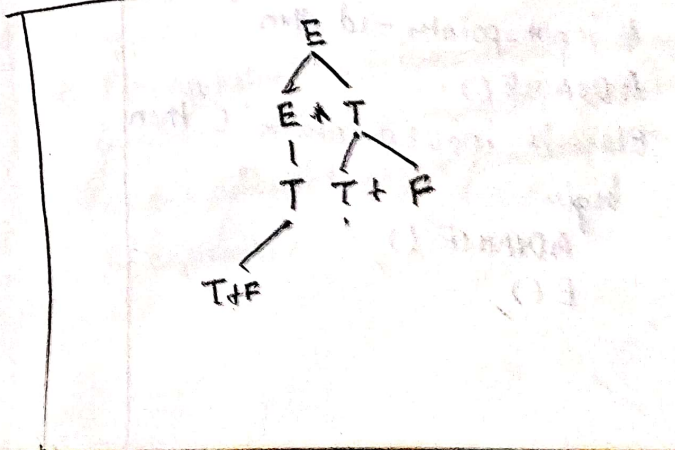
Top-down parsing



→ In recursive descent parser without backtracking we use one procedure for each non-terminal.

$E \rightarrow TE'$
 $E' \rightarrow +TE' / \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' / \epsilon$
 $F \rightarrow id / (E)$

$E \rightarrow EAT | T$
 $T \rightarrow T + F | F$
 $F \rightarrow id$




```

Procedure E()
begin
  T();
  EPRIME();
end

```

```

Procedure EPRIME()
begin
  if input-pointer = '+' then
    begin
      ADVANCE();
      T();
      EPRIME();
    end
  end
end

```

```

Procedure T()
begin
  F();
  TPRIME();
end

```

```

Procedure TPRIME()
begin
  if input-pointer = '*' then
    begin
      ADVANCE();
      F();
      TPRIME();
    end
  end
end

```

```

Procedure F()
begin
  if input-pointer = 'id' then
    ADVANCE();
  else if input-pointer = 'c' then
    begin
      ADVANCE();
      E();
    end
  end
end

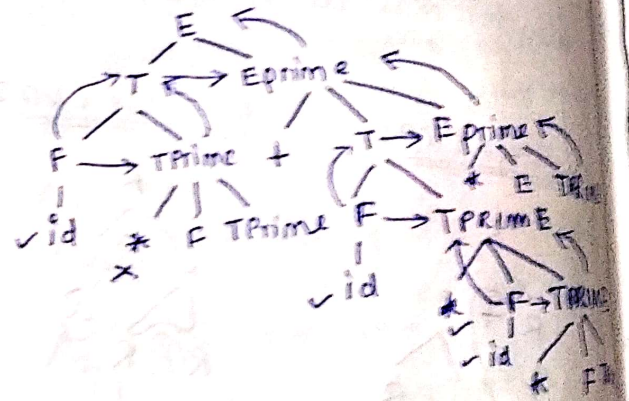
```

```

if input-pointer = ')' then
  ADVANCE();
else
  ERROR;
end
else
  ERROR;
end

```

> we can generate id + id * id from above grammar



* Consider a Production

$A \rightarrow abcd \mid abce$

\therefore Code is

```

Procedure A()
begin
  if input-pointer = 'a' then
    Advance();
  else if input-pointer = 'b' then
    Advance();
  else if input-pointer = 'c' then
    Advance();
  if input-pointer = 'd' then
    Advance();

```



```

else if ip = 'a' then
  advance c)
  - If input-pointer = 'b' then
    advance c)
  If input-pointer = 'c' then
    advance c)
  If input-pointer = 'e' then
    advance
  else
    error;
end
else
  error;
end

```

As we check the same condition twice, we reduce the grammar to

$A \rightarrow abcd / abce$
 \downarrow
 $A \rightarrow abcA'$
 $A' \rightarrow d / e$

while using recursive descent parser

↳ if a production is in form $A \rightarrow \alpha\beta / \alpha\gamma$ then it is called left factoring
 → After eliminating left factoring the productions are

$A \rightarrow \alpha A'$
 $A' \rightarrow \beta / \gamma$

→ eliminate left factoring

$C \rightarrow ics | icts / ses / a$
 $C \rightarrow b$

$S \rightarrow ics S' / a$
 $S' \rightarrow es$
 $C \rightarrow b$

After eliminating

$S \rightarrow ics S' / a$
 $S' \rightarrow es / e$
 $C \rightarrow b$

Non-recursive descent parser

↳ It is also called as the predictive parser

Suppose input-pointer is pointing to 'a' in input string and productions are

$X \rightarrow b / a$
 \bullet In backtracking and recursive descent parser the first production is expanded first (b) but in non-recursive descent parser, as the ip pointing to a and 'a' is available in 2nd production, it expands 2nd production first
 \bullet Parsing table is a 2-dimensional array which is used by non-recursive descent parser to predict which production to be expanded

$ab+ac$
 $= a(b+c)$

Parsing table

		columns - Terminal	
		a	b
row	non terminal X	$X \rightarrow a\phi$	$X \rightarrow b\alpha$

So X is expanded with
 $X \rightarrow a\phi$

→ for predictive parser there are some methods or prerequisites.

- ① Block diagram
- ② FIRST values
- ③ Follow Rule
- ④ Algorithm predictive parser table

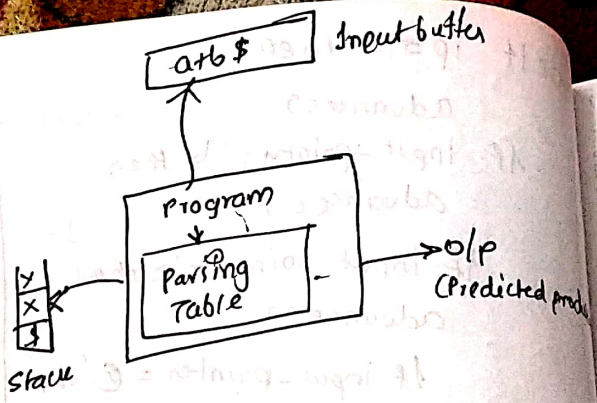
7 Block diagram for predictive parser / stack implementation

• Input buffer: - within input buffer it will contain input string & ending with ϕ

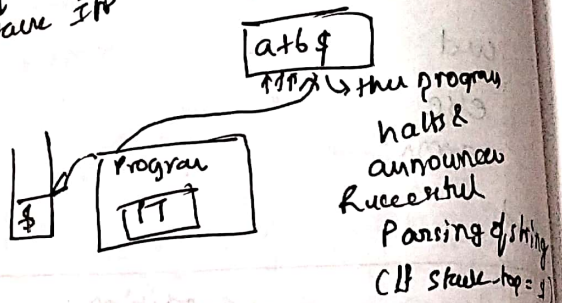
• stack: - within the stack, we push & pop the grammar symbols (VNT) and without any symbols, stack defaultly contains ϕ .

* Predictive parser actions are controlled by the program.

• Parsing table → used to predict which production needs to be expanded



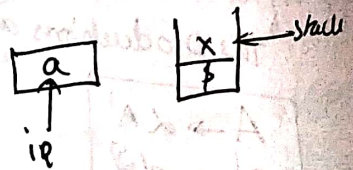
Case (i)
 ex: $X = a = \phi$
 \downarrow stack \uparrow IP



if (stack-top = ϕ) && (ip = '\$') then halt
if (stack-top = 'a') && (ip = 'a') then pop & increment input pointer so that program points next symbol

Case (ii)
 ex: $X = a \neq \phi$
 then pop() && ip++;

Case (iii)
 ex: $X \neq a \neq \phi$

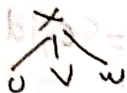


∴ we use parsing table to decide which production to expand.

$x \mid \alpha$
 $x \rightarrow uvw$

Consider
 $x \rightarrow uvw$ production

\therefore In above case (iii)
 Path u, v, w in
 order into the stack



First Rules:-

① If x is a terminal then
 $FIRST(x) = \{x\}$

② If $x \rightarrow \alpha$ is a production, where
 x is a non-terminal

α is a terminal and
 $\alpha \in (T \cup \{ \epsilon \})^*$

$FIRST(x) = \{\alpha\}$

and if $x \rightarrow \epsilon$ is a

Production then $FIRST(x) = \{ \epsilon \}$

③ $x \rightarrow y_1 y_2 \dots y_k$ where y_i
 may be terminal or non-terminal

For each y_i production

$FIRST(x) = FIRST(y_i)$

If $FIRST(y_i) = \epsilon$ then

$FIRST(x) = (FIRST(y_1) - \epsilon) \cup FIRST(y_2)$

If $FIRST(y_2) = \epsilon$ then

$FIRST(x) = (FIRST(y_1) - \epsilon) \cup$
 $(FIRST(y_2) - \epsilon) \cup$
 $(FIRST(y_3))$

and so on

If $y_1 y_2 \dots y_k \Rightarrow \epsilon$ then also x include
 ϵ into $FIRST(x)$

ex: for 3rd rule

consider a grammar

$S \rightarrow AB$

$A \rightarrow a$

$B \rightarrow b$

Non terminal = $\{S, A, B\}$

$FIRST(A) = \{a\}$

$FIRST(B) = \{b\}$

$FIRST(S) =$

as rule ③

$FIRST(S) = FIRST(A)$
 $= \{a\}$

ex: $S \rightarrow AB$

$FIRST(A) = FIRST(a)$ (rule ①)

$= \{a, \epsilon\}$

$FIRST(B) = FIRST(b)$

$= \{b\}$

$FIRST(S) = FIRST(A) - \epsilon \cup FIRST(B)$

$= \{a, \epsilon\} - \epsilon \cup \{b\}$

$= \{a\} \cup \{b\}$

$= \{a, b\}$

ex: $S \rightarrow AB\epsilon$

$A \rightarrow a/\epsilon$

$B \rightarrow b/\epsilon$

$\therefore \text{FIRST}(A) = \{a, \epsilon\}$

$\text{FIRST}(B) = \{b, \epsilon\}$

$\text{FIRST}(S) = \text{FIRST}(A)$

~~$(\{a, \epsilon\} - \epsilon) \cup (\text{FIRST}(B) - \epsilon)$~~
 ~~$(\{a, \epsilon\} \cup \{b, \epsilon\})$~~

@

$\text{FIRST}(S) = \text{FIRST}(A) - \epsilon \cup \text{FIRST}(B) - \epsilon \cup \epsilon$

$\therefore \text{FIRST}(S) = (\{a, \epsilon\} - \epsilon) \cup (\{b, \epsilon\} - \epsilon) \cup \epsilon$

$\text{FIRST}(S) = \{a, b, \epsilon\}$

compute FIRST for following grammar

$E \rightarrow TE'$

$E' \rightarrow +TE' / \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' / \epsilon$

$F \rightarrow (E) / id$

~~$\text{FIRST}(E) = (\text{FIRST}(T) \cup \text{FIRST}(E'))$~~

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

$\text{FIRST}(+) = \{+\}$

$\text{FIRST}(*) = \{*\}$

$\text{FIRST}(\epsilon) = \{\epsilon\}$

$\text{FIRST}(\epsilon) = \{\epsilon\}$

$\text{FIRST}(id) = \{id\}$

$\text{FIRST}(E) = \{(, id\}$

$E \rightarrow TE'$

$E \rightarrow FT'E'$

$E \rightarrow \underline{(E)T'E'} / \underline{idT'E'}$
 \hookrightarrow so include FIRST terminal in FIRST(E)

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

$E' \rightarrow +TE' / E' \rightarrow \underline{\epsilon}$ (value)
 \hookrightarrow so include in FIRST of E'

$\text{FIRST}(T) = \{(, id\}$

$T \rightarrow FT'$

$T \rightarrow \underline{(E)T'} / \underline{idT'}$

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

$T' \rightarrow *FT' / \underline{\epsilon}$

include in FIRST(T')

$\text{FIRST}(F) = \{(, id\}$

$F \rightarrow \underline{(E)} / \underline{id}$

compute first

$S \rightarrow ACB / cB / Ba$

$A \rightarrow da / Bc$

$B \rightarrow g / \epsilon$

$C \rightarrow h / \epsilon$

- $FIRST(L) = \{h, \epsilon\}$
 $FIRST(B) = \{g, \epsilon\}$
 $FIRST(A) = \{d, g, \epsilon, h\}$
 $FIRST(S) = \{d, g, \epsilon, h, a, b, \bullet\}$

Follow rules

Given a non-terminal A the set Follow(A) consisting of terminals. Apply the following rules until nothing can be added to any follow set

1. If A is start symbol, then ϕ is in follow(A)
2. If there is any production $B \rightarrow aAy$ then $FIRST(y) - \{\epsilon\}$ is in follow(A)
3. If there is a production $B \rightarrow aAy$ or $B \rightarrow A$ such that a is in $FIRST(y)$ then everything of Follow(B) is in Follow(A).

- ex:1
- $E \rightarrow TA$
 - $A \rightarrow +TA / \epsilon$
 - $T \rightarrow FC$
 - $C \rightarrow *FC / \epsilon$
 - $F \rightarrow (B) / id$

$FIRST(E) = FIRST(T) = FIRST(F) = \{\epsilon, id\}$
 $FIRST(A) = \{+, \epsilon\}$ $FIRST(C) = \{*, \epsilon\}$

$\therefore Follow(E) = Follow(A) = \{\$, \epsilon\}$
 $Follow(T) = Follow(C) = \{\$, \epsilon, +\}$
 $Follow(F) = \{+, *, \epsilon, \$\}$

moves made by predictive parser for above grammar (id*id*id) for above grammar

Stack	input symbol	output
\$E	id*id*id\$	
\$AT	id*id*id\$	$E \rightarrow TA$
\$ACE	id*id*id\$	$T \rightarrow FC$
\$ACid	id*id*id\$	$F \rightarrow id$
\$AC	+id*id\$.
\$A	+id*id\$	$C \rightarrow \epsilon$
\$AT+	+id*id\$	$A \rightarrow +TA$
\$AT	id*id\$.
\$ACF	id*id\$	$T \rightarrow FC$
\$ACid	id*id\$	$F \rightarrow id$
\$AC	*id\$.
\$ACF*	*id\$	$C \rightarrow *FC$
\$ACF	id\$.
\$ACid	id\$	$F \rightarrow id$
\$Ac	\$	$C \rightarrow \epsilon$
\$A	\$	$A \rightarrow \epsilon$
\$	\$	

Algorithm: Construct the predictive parsing table

Input: Grammar G

Output: Predictive parsing table

Step (1): For each production $A \rightarrow \alpha$ do
Step (2) and step (3)

Step (2): For each terminal a in $FIRST(\alpha)$ add
 $A \rightarrow \alpha$ into $M[A, a]$.

Step (3): If ϵ is in $FIRST(\alpha)$ then add
 $A \rightarrow \alpha$ into $M[A, b]$ for each
terminal b in $FOLLOW(A)$.

If ϵ is in $FIRST(\alpha)$ and $\$$ is
in $FOLLOW(A)$ then add $A \rightarrow \alpha$
into $M[A, \$]$

Step (4): make undefined entries ϵ non-
error.

ex: for step (2)

$S \rightarrow AB$
 $A \rightarrow a|b$
 $B \rightarrow b$

	a	b
S	$S \rightarrow AB$	$S \rightarrow AB$
A	$A \rightarrow a$	$A \rightarrow b$

$S \rightarrow AB$ $A \rightarrow a$

$A \rightarrow b$

$$FIRST(AB) = FIRST(aB) \cup FIRST(bB) \\ = \{a, b\}$$

ex: for step (3) - (i)

$S \rightarrow Ac$

$A \rightarrow \epsilon$

$FOLLOW(C) = \{c\}$

	c
A	$A \rightarrow \epsilon$

ex: for step (3) - (ii)

$S \rightarrow Ac | A$

$A \rightarrow \epsilon$

$A \rightarrow \epsilon$

$A \rightarrow \alpha$

$FIRST(\epsilon) = \{\epsilon\}$

$FOLLOW(A) = \{c\}$

Constructing predictive parsing table for
grammar

$E \rightarrow TE'$

$E' \rightarrow +TE' | \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' | \epsilon$

$F \rightarrow (E) | id$

sol: $E \rightarrow TE'$: ($A \rightarrow \alpha$) form

$$FIRST(\alpha) = FIRST(TE')$$

$$= FIRST(FT'E')$$

$$= FIRST((E)TE')$$

$$FIRST(idTE')$$

$$= \{c, id\}$$

$$M[E, c] \quad M[E, id] \rightarrow \text{to}$$

these entries were needed to place $E \rightarrow TE'$
production

$E' \rightarrow +TE' | \epsilon$: ($A \rightarrow \alpha$)

$$FIRST(+TE') = \{+\}$$

$$FIRST(\epsilon) = \{\epsilon\}; FOLLOW(E') = \{c\}$$

$(E' \rightarrow +TE')$ is added in $M[E', c]$

$(E' \rightarrow \epsilon)$ is added in $M[E', c]$
 $M[E', c]$

$$T \rightarrow FT' \quad (A \rightarrow \alpha)$$

$$\text{FIRST}(FT') = \text{FIRST}(ET') \cup \text{FIRST}(idT')$$

$$= \{ \epsilon, id \}$$

$(T \rightarrow FT')$ production is added in
 $M[A, C], M[A, id]$

$$T' \rightarrow \epsilon FT' / \epsilon \quad (A \rightarrow \alpha)$$

$$\text{FIRST}(\epsilon FT') = \{ \epsilon \}$$

$$\text{FIRST}(\epsilon) = \{ \epsilon \}$$

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

follow(T)
 as ϵ is there after T'

$T' \rightarrow \epsilon FT' / \epsilon$ production is added
 with $M[T', +], M[T',)],$
 $M[T', \$]$

$$F \rightarrow (E)$$

$$\text{FIRST}((E)) \rightarrow \{ (\}$$

$$F \rightarrow id$$

$$\text{FIRST}(id) \rightarrow \{ id \}$$

	id	+	*	^)	\$
E		$E \rightarrow TE'$			$E \rightarrow TE'$	
E'			$E \rightarrow TE'$		$E' \rightarrow +$	$E' \rightarrow +$
T		$T \rightarrow FT'$			$T \rightarrow FT'$	
T'			$T' \rightarrow \epsilon FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$				$F \rightarrow (E)$	

The initial configuration of the predictive parser:

Stack I/P Buffer
 $\$ S$ $W \$$

The acceptance configuration

Stack I/p buffer
 $\$$ $\$$

Stack I/P buffer o/p

$\$ E$ $id + id + id \$$

$\$ E'$ $id + id * id \$$ $E \rightarrow TE'$

$\$ E' T'$ $id + id * id \$$ $T \rightarrow FT'$

$\$ E'$ $+ id * id * id \$$

$\$ E'$ $+ id * id \$$ $T' \rightarrow \epsilon$

$\$ E' T'$ $+ id * id \$$

$\$ E' T'$ $id * id \$$

$\$ E' T'$ $id * id \$$ $T \rightarrow FT'$

$\$ E' T' E$ $id * id \$$ $F \rightarrow id$

$\$ E' T' id$ $id * id \$$

$\$ E' T' id$ $id * id \$$

$\$ E' T' F$ $id * id \$$ $T' \rightarrow \epsilon FT'$

$\$ E' T' F$ $id \$$

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

$\$ E' T' F$ $id \$$

$\$ E' T' F$ $id \$$

$\$ E' T'$ $\$$

$\$ E'$ $\$$ $T' \rightarrow \epsilon$

$\$$ $\$$ $E' \rightarrow \epsilon$

$\$$ $\$$ (Accept)

$S \rightarrow ictss'$

$FIRST(ictss') = \{i\}$

This production is added to $M[S, i]$

$S \rightarrow a$

$FIRST(a) = \{a\}$

This production is added to $M[S, a]$

$S' \rightarrow \epsilon$

$FOLLOW(S') = FOLLOW(S) = \{\$, \epsilon\}$

This production is:

$FOLLOW(S) = FIRST(S') = \{\epsilon, e\}$
 \downarrow
 ϵ is ignored

$S' \rightarrow es/\epsilon$

$FOLLOW(S) = FOLLOW(S')$

$M[S', \$]$ and $M[S', e]$ we add the production $S' \rightarrow \epsilon$

$M[S', e]$ we add production $S' \rightarrow es$

	i	a	e	b	ϵ	$\$$
S	$S \rightarrow ictss'$	$S \rightarrow a$				
S'			$S' \rightarrow e$ $S' \rightarrow es$			$S' \rightarrow \epsilon$
C				$C \rightarrow b$		

Follow () rules

- $FOLLOW(S') = FOLLOW(S) = \{\$, \epsilon\}$
- $FOLLOW(S) = FIRST(S') = \{\epsilon, e\}$
- $FOLLOW(S) = FIRST(S')$

LL(1)

* A grammar

Follow gives

• LL

left to

Scan

input

Pro

①

②

③

LL(1) Grammar

* A grammar for which parsing table with no multiplied defined entries is called LL(1) grammar

Following are the properties to check the given grammar LL(1) or not.

• LL(k) → no. of lookahead symbols to take parsing decisions
 left to right scanning of input strings
 left most derivation

Properties

① no ambiguous, left recursion, left factoring is LL(1) grammar

② A parsing table with no multiple definition entries is LL(1) grammar

③ For each production $A \rightarrow \alpha / \beta$

should satisfy following conditions.

(i) $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) \neq \{\epsilon\}$

ex:

$S \rightarrow xy / xw$

$S \rightarrow xy$

$\text{FIRST}(xy) = \{x\}$

$M[S, x]$

$S \rightarrow xw$

$\text{FIRST}(xw) = \{x\}$

$M[S, x]$

(ii) $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) \neq \{\epsilon\}$

ex: $S \rightarrow A / B$

$A \rightarrow \epsilon$

$B \rightarrow \epsilon$

$S \rightarrow A$

$\text{FIRST}(A) = \{\epsilon\}$

$\text{Follow}(S) = \{\epsilon\}$

$S \rightarrow B$

$\text{FIRST}(B) = \{\epsilon\}$

$\text{Follow}(S) = \{\epsilon\}$

almost one production can derive ϵ

(iii) $\text{FIRST}(\alpha) \cap \text{Follow}(A) = \phi$

ex:

$S' \rightarrow \epsilon / e$

$\text{FIRST}(\epsilon) = \{\epsilon\}$

$S' \rightarrow e$

$\text{Follow}(S') = \{\$, \epsilon\}$

$\text{FIRST}(\alpha) \cap \text{Follow}(S')$

$\{\epsilon\} \cap \{\$, \epsilon\} = \{\epsilon\}$

(iv) $\text{Follow}(A) \cap \text{FIRST}(B) = \phi$

ex: $S' \rightarrow e / es$

$\text{FIRST}(es)$

• Check following grammar is LL(1) grammar

$S \rightarrow AaAb / BbBa$

$A \rightarrow \epsilon$

$B \rightarrow \epsilon$

Sol: $S \rightarrow AaAb / BbBa$

$\text{FIRST}(AaAb) = \text{FIRST}(A)$

$= \{a\}$

$\text{FIRST}(BbBa) = \text{FIRST}(B)$

$= \{b\}$

$\text{FIRST}(A) \cap \text{FIRST}(B) \neq \{a\}$

$\text{FIRST}(A) \cap \text{FIRST}(B) = \{\phi\}$

∴ given grammar is LL(1) grammar

• Construct predictive parsing table for

$$S \rightarrow aAbB \mid bAbB \mid \epsilon$$

$$A \rightarrow s$$

$$B \rightarrow s$$

① $S \rightarrow aAbB \mid bAbB \mid \epsilon$

$$S \rightarrow aAbB :$$

$$\text{FIRST}(S) = \text{FIRST}(aAbB) = \{a\} \quad M[S, a]$$

$$S \rightarrow bAbB$$

$$\text{FIRST}(bAbB) = \{b\} \quad M[S, b]$$

$$S \rightarrow \epsilon$$

$$\text{Follow}(S) = \{a, b, \$\}$$

$$A \rightarrow s$$

$$\text{Follow}(A) = \text{Follow}(S)$$

$$S \rightarrow aAbB$$

$$\text{Follow}(A) = \text{FIRST}(b\beta)$$

$$= \{b\}$$

$$\text{Follow}(A) = \text{FIRST}(a\beta)$$

$$= \{a\}$$

$$B \rightarrow s$$

$$\text{Follow}(B) = \text{Follow}(S)$$

$$S \rightarrow bAbB$$

$$M[S, a], M[S, b], M[\$, \$]$$

② $A \rightarrow s$

$$\text{FIRST}(A) = \text{FIRST}(S) = \{a, b, \$\}$$

$$= \text{FIRST}(aAbB) + \text{FIRST}(bAbB) + \text{FIRST}(\epsilon)$$

$$= \{a\} + \{b\} + \text{Follow}(B)$$

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

③ $B \rightarrow s$

$$\text{FIRST}(B) = \text{FIRST}(S) = \{a, b, \$\}$$

$$M[B, a]$$

$$M[B, b]$$

$$M[B, \$]$$

	a	b	\$
S	$S \rightarrow aAbB$ $S \rightarrow \epsilon$	$S \rightarrow bAbB$ $S \rightarrow \epsilon$	$S \rightarrow \epsilon$
A	$A \rightarrow s$	$A \rightarrow s$	$A \rightarrow s$
B	$B \rightarrow s$	$B \rightarrow s$	$B \rightarrow s$

∴ This is not LL(1) grammar as it is possible (C1).

• Construct LL(1) parsing table for given grammar & also show parsing action for input

$$S \rightarrow aAdc \mid BcAe$$

$$A \rightarrow b/e$$

$$B \rightarrow Cf/d$$

$$C \rightarrow fe$$

$$\text{FIRST}(S) = \text{FIRST}(aAdc) \cup \text{FIRST}(BcAe) = \{a, d, f\}$$

$$\text{FIRST}(A) = \text{FIRST}(b) \cup \text{FIRST}(e) = \{b, e\}$$

$$\text{FIRST}(B) = \{d, f\}$$

$$\text{FIRST}(C) = \{f\}$$

$$\rightarrow \text{Follow}(S) = \{\$, \}$$

$$\text{Follow}(A) = \{c, e\}$$

$$\text{Follow}(B) = \{c\}$$

$$\text{Follow}(C) = \{f\}$$

Parsing table for LL(1) parsing table

Non-terminal	input symbol						
	a	b	c	d	e	f	\$
S	$S \rightarrow aAdc$			$S \rightarrow BcAe$			$S \rightarrow \epsilon$
A		$A \rightarrow b$	$A \rightarrow e$		$A \rightarrow \epsilon$		
B				$B \rightarrow d$			$B \rightarrow \epsilon$
C							$C \rightarrow f$

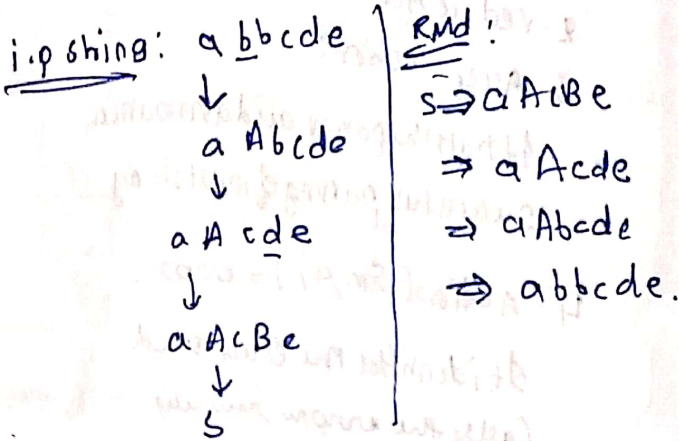
questions:

- Every SLR grammar is unambiguous but not every unambiguous grammar is SLR explain
- yacc calculator program with own production

UNIT - II

Bottom-up parsing :- constructs the parse tree starting from the leaf node (i/p string) working up towards root node (starting non terminal).

ex: $S \rightarrow aAcBe$
 $A \rightarrow Ab|b$
 $B \rightarrow d$

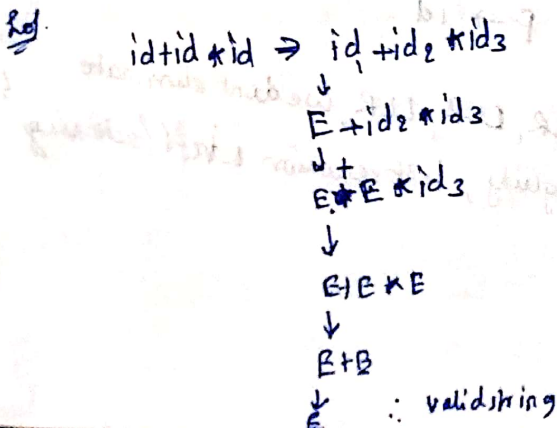


Reduction: A substring matching with the right hand side of the production, replace it with the left hand side of production is called reduction.

Handle:- A substring matching with the right hand side of production replacing such substrings with left hand side of production. eventually leads to starting non terminal is called handle.

- find appropriate handles to reduce given i/p string 'id+id*id'

- $E \rightarrow E+E$
- $E \rightarrow E * E$
- $E \rightarrow (E)$
- $E \rightarrow id$



Right sentential form	Handle	Production
id ₁ + id ₂ * id ₃	id ₁	$E \rightarrow id$
E + id ₂ * id ₃	id ₂	$E \rightarrow id$
E + E * id ₃	id ₃	$E \rightarrow id$
E + E * E	E * E	$E \rightarrow E * E$
E + E	E + E	$E \rightarrow E + E$
E		

Stack	I/P Buffer	Action
\$	w\$	
\$	id ₁ + id ₂ * id ₃ \$	Shift
\$id ₁	+ id ₂ * id ₃ \$	Reduce $E \rightarrow id$
\$E	+ id ₂ * id ₃ \$	Shift
\$E +	id ₂ * id ₃ \$	Shift
\$E + id ₂	* id ₃ \$	Reduce $E \rightarrow id$
\$E + E	* id ₃ \$	Shift
\$E + E *	id ₃ \$	Shift
\$E + E * id ₃	\$	Reduce $E \rightarrow id$
\$E + E * E	\$	Reduce $E * E \rightarrow E$
\$E + E	\$	Reduce $E \rightarrow E + E$
\$E	\$	Accept

Shift: Shift the next input symbol into stack

Reduce:- The right end of string to be reduced must be on top of stack

Locate the left end of string within stack and replace the string with non terminal

Accept: Successful completion of parsing

Error: Discove a syntax error & call an error recovery routine

Conflicts during shift-reduce parsing

- There are context free grammars for which shift-reduce grammar cannot be used. We cannot decide whether to use shift or reduce eg: dangling if-else grammar

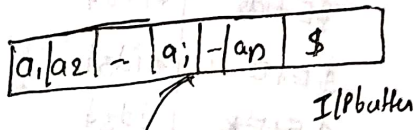
LR parsers :-

- (1) SLR (Simple LR)
- (2) CLR (Canonical LR)
- (3) LALR (Look Ahead LR)

Stack:-

$S_0 \ x_1 \ S_1 \ x_2 \ S_2 \ \dots \ x_m \ S_m$
 $x_1, x_2 \ \dots \ x_m \rightarrow$ Standard for
 Grammar symbol

$S_0 \ S_1 \ \dots \ S_m \rightarrow$ Standard for
 State symbol



Input buffer



Parsing Table	
Action	GOTO

Terminal Action	Statement terminal
S, R, E	State symbol

Configuration of LR Parser:-

$(S_0 \ x_1 \ S_1 \ x_2 \ S_2 \ \dots \ x_m \ S_m, a_i, a_{i+1} \ \dots \ a_n \ \$)$

$ACTIONS[S_m, a_i] = \text{shift } S$

$(S_0 \ x_1 \ S_1 \ x_2 \ S_2 \ \dots \ x_m \ S_m, a_i, a_{i+1} \ \dots \ a_n \ \$)$

$S = GOTO[S_m, a_i]$

- 1. Shift Action
- 2. reduce Action
- 3. Accept Action:

It is the parser and announces successful parsing of input string.

4. $ACTIONS[S_m, a_i] = \text{error}$:

It identifies the error and calls the error recovery routine ().

Check the following ~~grammar~~ string

id + id + id Valid or not for the following given grammar & Parsing table

- (i) $E \rightarrow E + T$
- (ii) $E \rightarrow T$
- (iii) $T \rightarrow T * F$
- (iv) $T \rightarrow F$
- (v) $F \rightarrow (E)$
- (vi) $F \rightarrow id$

• For SLR, CLR, LALR, we don't eliminate ambiguity, left recursion & left factoring

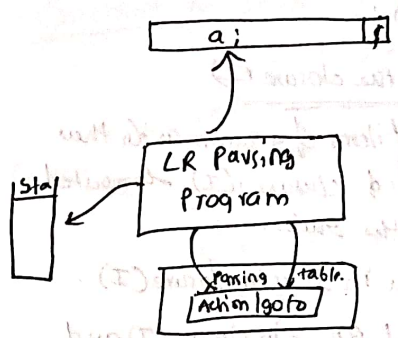
Parsing table

	id	t	*	C)	\$	E	T	F
0	S5			S4			1	2	3
1		S6				acc			
2		r2	S1		r2	r2			
3		r4	r4		r4	r4			
4	S5			S4			8	2	3
5		r6	r6		r6	r6			
6	S5			S4				9	3
7	S5			S4					10
8		S6			S11				
9		r1	S7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			

* empty cells are errors

Stack	LLP buffer	Action
0	idtid*ids	Shift
0 id \$	tid*ids	reduce F → id
0 F \$	tid*ids	reduce T → F
0 T \$	tid*ids	reduce E → T
0 E \$	tid*ids	Shift
0 E t \$	idtid*	Shift
0 E t d \$	*ids	reduce F → id
0 E t \$	*ids	reduce T → F
0 E t t \$	*ids	Shift
0 E t t \$	*ids	Shift
0 E t t \$	id	
0 E t t \$	*ids	

⇒



Prerequisites for LR parse

(i) Item: placing the dot in right hand side of the production in any position
 each item describes how much production we have seen at a given point in the parsing process.

$A \rightarrow XYZ$
 $A \rightarrow \cdot XYZ$ → expecting a string derived from X
 $A \rightarrow X \cdot YZ$ → expecting a string derived from Y
 $A \rightarrow XY \cdot Z$ → expecting a string derived from Z

eg: $S \rightarrow Aa$
 $A \rightarrow b$
 $\therefore S \rightarrow \cdot Aa$
 $S \rightarrow A \cdot a$
 $S \rightarrow Aa \cdot$

(ii) Augmented grammar : (G')

Adding the production $S' \rightarrow s$ to the actual productions

$S' \rightarrow s$
 $S \rightarrow Aa$
 $A \rightarrow b$

Purpose of adding this production is to halt the parsing process & announce the successful parsing of the input string

000
 (ii) closure (I) :

Steps to compute the closure (I)

- If I is a set of items of grammar G then the set of items of closure(I) computed from I by the rules.
 - ① every item in I is in closure(I).
 - ② If $A \rightarrow \alpha \cdot B \beta$ is in closure(I) and $B \rightarrow \gamma$ is a production in Grammar G then add $B \rightarrow \cdot \gamma$ to I if it is not already in I.

eg: $S \rightarrow Aa$ I: $S \rightarrow \cdot Aa$
 $A \rightarrow b$ $A \rightarrow \cdot b$

Closure(I)
 • closure($S \rightarrow \cdot Aa$)
 • $A \rightarrow \cdot b$

Compute closure of

- $E' \rightarrow \cdot E$
- $E \rightarrow E+T$
- $E \rightarrow T$
- $T \rightarrow T * F$
- $T \rightarrow F$
- $F \rightarrow (E)$
- $F \rightarrow id$

Sol: $E' \rightarrow \cdot E$ I: $E' \rightarrow \cdot E$
 $E \rightarrow \cdot E+T$
 $F \rightarrow \cdot T$

Closure(I)
 • closure($E' \rightarrow \cdot E$)
 $E \rightarrow \cdot E+T$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$

$F \rightarrow id$

Procedure closure(I)

```

begin
  repeat
    for each item  $A \rightarrow \alpha \cdot B \beta$  is in closure(I)
    and
      each production  $B \rightarrow \gamma$  is in grammar G
      such that  $B \rightarrow \cdot \gamma$  is not in I
    do add  $B \rightarrow \cdot \gamma$  to I
  until no more set of items added to I
  return I
end.
    
```

goto(I, x) is defined as closure of all set of items of $A \rightarrow \alpha x \cdot \beta$ such that $A \rightarrow \alpha \cdot x \beta$ is in I.

I: $A \rightarrow \alpha \cdot x \beta$
 $goto(I, x) \rightarrow$ dot is moved to right of x

\therefore I: $A \rightarrow \alpha \cdot x \beta$

Closure($A \rightarrow \alpha x \cdot \beta$)

eg: I: $E' \rightarrow \cdot E$
 $\downarrow \downarrow$
 closure($E' \rightarrow \cdot E$)
 $E \rightarrow \cdot E+T$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot id$

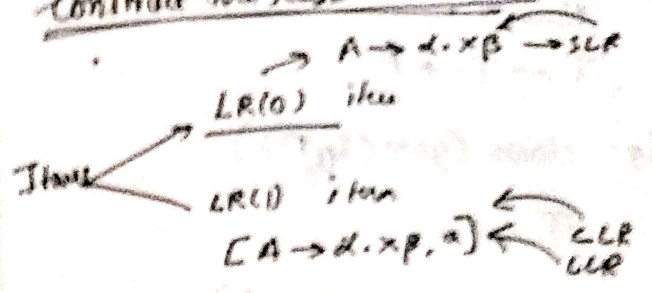
closure (goto (I, E))
 $E' \rightarrow E$
 $E \rightarrow E \cdot TT$
 closure (goto (I, T))
 $E \rightarrow T$
 $T \rightarrow T \cdot RF$
 closure (goto (I, C))
 $F \rightarrow (\cdot E)$
 $E \rightarrow \cdot E \cdot TT$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T \cdot RF$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot id$ — halt —

Source of Item Construction:

Procedure ITEM(G')

begin
 $c = \{ \text{closure}(\{ E' \rightarrow \cdot E \}) \}$
 repeat
 for each set of items in c &
 each grammar symbol x
 such that (I, x) is not empty &
 $\text{goto}(I, x)$ is not c .
 Do add $\text{goto}(I, x)$ to c
 until no more sets of items added
 to c
 end

Construct the set of LR(0)



eg: $F \rightarrow E \cdot TT$
 $E \rightarrow T$
 $T \rightarrow T \cdot RF$
 $T \rightarrow F$
 $F \rightarrow (E)$
 $F \rightarrow id$

def: augmented grammar is (G')

$E' \rightarrow E$
 $E \rightarrow E \cdot TT$ — — — — — start and Producers

Initial Itemset I_0 : $\text{closure}(E' \rightarrow \cdot E)$

$E \rightarrow \cdot E \cdot TT$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T \cdot RF$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot id$

$c = \{ I_0, I_1, I_2, I_3, I_4, I_5, I_6, I_7, I_8, I_9, I_{10}, I_{11} \}$
 $I_1: \text{closure}(\text{goto}(I_0, E))$
 $E' \rightarrow E \cdot$
 $E \rightarrow E \cdot \cdot TT$

$I_2: \text{closure}(\text{goto}(I_0, T))$

$E \rightarrow T$
 $T \rightarrow T * F$

$I_3: \text{closure}(\text{goto}(I_0, F))$

$T \rightarrow F \cdot$

$I_4: \text{closure}(\text{goto}(I_0, \epsilon))$

$F \rightarrow (\cdot E)$

$E \rightarrow \cdot E + T$

$E \rightarrow \cdot T$

$T \rightarrow \cdot T * F$

$T \rightarrow \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot \text{id} \quad \text{--- halt ---}$

$I_5: \text{closure}(\text{goto}(I_0, \text{id}))$

$F \rightarrow \text{id} \cdot$

$I_6: \text{closure}(\text{goto}(I_1, +))$

$E \rightarrow E + \cdot T$

$T \rightarrow \cdot T * F$

$T \rightarrow \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot \text{id}$

$I_7: \text{closure}(\text{goto}(I_2, *))$

$T \rightarrow T * \cdot F$

$F \rightarrow \cdot (E)$

$F \rightarrow \cdot \text{id}$

$I_8: \text{closure}(\text{goto}(I_4, E))$

$F \rightarrow (E \cdot)$

$E \rightarrow E + \cdot T$

~~$I_9: \text{closure}(\text{goto}(I_4, T))$~~

~~$E \rightarrow T \cdot$~~

~~$T \rightarrow \epsilon T \cdot * F$~~ = I_2

$I_{10}: \text{closure}(\text{goto}(I_4, F))$

~~$F \rightarrow \text{id} \cdot$~~

= I_3

~~$I_{11}: \text{closure}(\text{goto}(I_4, \epsilon))$~~

= I_4

~~$I_{12}: \text{closure}(\text{goto}(I_4, \text{id}))$~~

= I_5

$I_{13}: \text{closure}(\text{goto}(I_6, T))$

$E \rightarrow E + T \cdot$

$T \rightarrow T \cdot * F$

~~$I_{14}: \text{closure}(\text{goto}(I_6, F))$~~

$\Rightarrow I_3$

~~$I_{15}: \text{closure}(\text{goto}(I_6, \epsilon))$~~

$\Rightarrow I_4$

~~I_{10} closure (goto (I_6 , id))~~
 $\Rightarrow I_5$

I_{10} closure (goto (I_7 , E))
 $T \rightarrow T * F.$

~~I_{11} closure (goto (I_7 , ϵ))~~
 $\Rightarrow I_{10}$

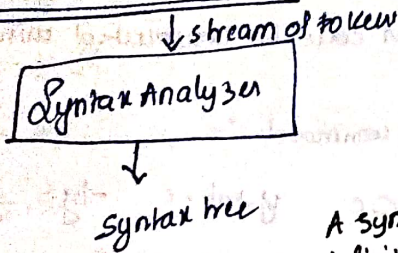
~~I_{11} closure (goto (I_7 , id))~~
 $\Rightarrow I_5$

I_{11} closure (goto (I_8 , ϵ))
 $F \rightarrow (E).$

~~I_{12} closure (goto (I_8 , +))~~
 $\Rightarrow I_6$

UNIT - III

⇒ Syntax directed translation :



- $E \rightarrow E + E$
- $E \rightarrow E * E$
- $E \rightarrow (E)$
- $E \rightarrow id$

Ambiguous grammar

A syntax directed definition (SDD) is a context-free grammar together with attributes & rules

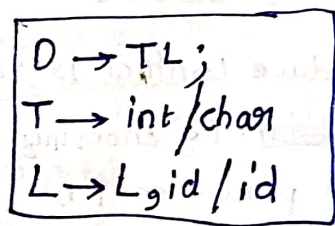
Attributes associated with grammar symbols & rules associated with productions.

∴ unambiguous grammar is

- $E \rightarrow E + T / T$
- $T \rightarrow T * F / F$
- $F \rightarrow (E) / id$

∴ D is used to declare the datatype of variable

Datatype → list of variables



L → list of variables

grammar for declaring variables of particular datatype

This grammar don't explain how the expression is evaluated

so we need to add some additional information in form of semantic rules.

we have 2 methods to associate semantic rules to the context free grammar

- (1) Syntax directed definition
 - (2) Translation schema
- CFG + Semantic rule

evaluation of arithmetic expressions
(using syntax directed definition)

- $E \rightarrow E + F$
- $E \rightarrow T$
- $T \rightarrow T * F$
- $T \rightarrow F$
- $F \rightarrow (E)$
- $F \rightarrow digit.$

- we associate the attributes with the terminal & non-terminals using the (·) operator
- attributes are used to hold some information. So we associate val with non-terminals for value.
- And attribute needs to be associated based on the type.

- eg: • type (datatype)
- string (string)

- If terminal is predefined we don't associate any attribute
- eg: (+, -, *, /, (,), , - . etc)

$E \rightarrow E + T$

$$E.val = E_1.val + T.val$$

$E \rightarrow T$

$$E.val = T.val$$

$T \rightarrow T_1 * F$

$$T.val = T_1.val * F.val$$

$T \rightarrow F$

$$T.val = F.val$$

$F \rightarrow (E)$

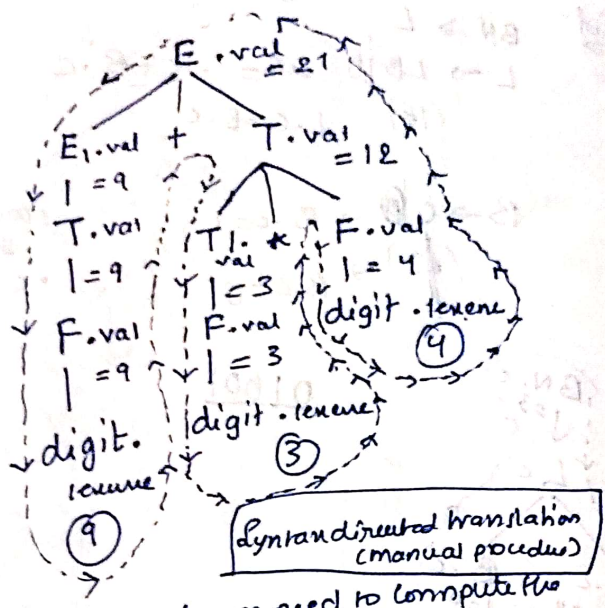
$$F.val = E.val$$

$F \rightarrow digit$

$$F.val = digit$$

$9 + 3 * 4$

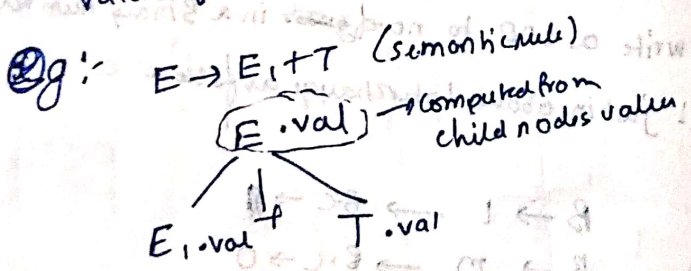
digit + digit * digit



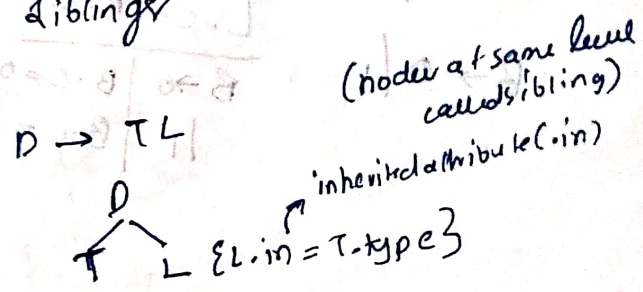
- * At each node we need to compute the attribute values
- * we are going to use post order traversal
left → right → root

→ There are two types of attributes

① Synthesized attribute :-
In a parse tree attribute value at a node computed from attribute values of child nodes



② Inherited attribute : In a parse tree at a node attribute value computed from attribute value of parent / or siblings



S-attributed definition :

• A syntax directed definition synthesized attributes (exclusively) is called s-attributed definition.

$$x + ty * 3$$

$$\underbrace{\hspace{2cm}}_{T_1}$$

$$x + T_1$$

$$\boxed{T_2 = x + T_1} \quad \checkmark$$

Different types of 3-address code statements :

① Assignment ~~operation~~ Statement : $A = B \text{ op } C$

② Assignment Instruction : $A = \text{op } B$

unary minus
~ (negation)

Type cast

int float

$$a = b$$

$$\hookrightarrow a = (\text{int}) b$$

③ Copy statement :

$$A = B$$

④ unconditional Jump statement :-

goto label ; / Jump

⑤ Conditional Jump statement

If A (relational operator) B .

goto L

⑥ Procedure call statement :-

Param n call (P, n)

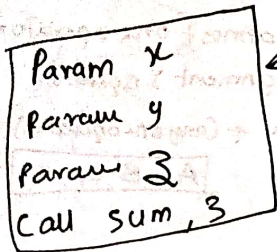
Procedure name

no. of parameters
Passing in
call statement

↳ each parameter is represented
by param

eg: $\text{Sum}(x, y, z)$

This procedure call is converted into the three address code



$\text{Sum}(x, y, z)$ is converted as like this

① Indexed Assignment Statement

$$x = y[i]; \quad x[i] = y$$

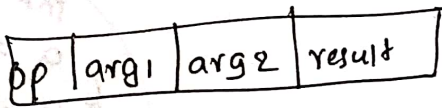
② Address Assignment pointer Assignment

$$x = \&y$$

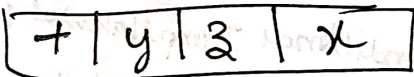
$$*x = y$$

Record structure : To store the Three address code

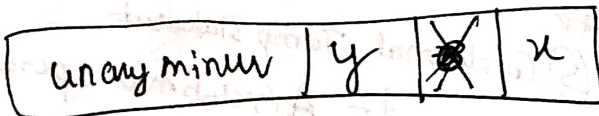
① Quadruple record structure



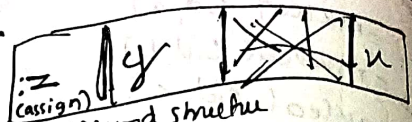
eg: 1 $x = y + z$



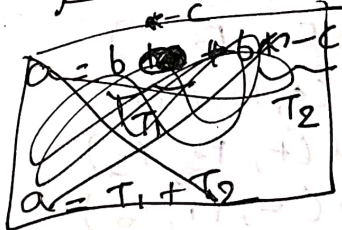
eg: 2 $x = y$



eg: 3 $x = y$ (if no op then we have to assign)



② Triple record structure
③ Indirect Triple record structure
: Representing the following in Quadruple, Triple, Indirect Triple



$$① a = b * c + b * c$$

$$T1 = -c$$

$$T2 = b * T1$$

$$T3 = -c$$

$$T4 = b * T3$$

$$T5 = T2 + T4$$

$$a = T5$$

one record \rightarrow and memory is assigned to each record

id	op	Arg1	Arg2	result
0	uminus	c	--	T1
1	*	b	T1	T2
2	uminus	c	--	T3
3	*	b	T3	T4
4	+	T2	T4	T5
5	assign (:=)	T5	--	a

In actual implementation the pointer to symbol table is written instead of names

disadvantage
Structure

Disadvantage of Quadruple Record Structure

Structure:

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

Lexical analyze

each identifier stored in symbol table

while translating into three address code

In the quadruple record

the symbols are pointers to symbol tables & some are not temporary variables
eg: T₁, T₂

And these temporary variables also need to store inside the symbol table

So it occupies huge memory

To avoid this we use Triple record structure

Instead of storing in temporary variable & also memory location we use triple records

Triple Record structure for $a = b * -c + b * -c$

mem loc	op	Arg1	Arg2
(0)	unary m	c	--
(1)	*	b	(0)
(2)	unary m	c	--
(3)	*	b	(T ₃) (2)
(4)	+	(1)	(3)
(5)	assign	a	(4)

Disadvantage of Triple Records

- During execution the assigned memory locations for records may not be free.
- So we face runtime overhead to change the memory location during execution & otherwise we get undesired results

In order to overcome this issue we go for the Indirect Triple record structure.

It is a list of pointers to the triple record structure

Indirect Triple Record Structure

Ptr	Ptr (if exists)	Statement
(17)	(13)	(0)
(18)	(14)	(1)
(21)	(15)	(2)
(22)	(16)	(3)
(23)	(17)	(4)
(24)	(18)	(5)

Address stored inside pointer will not change.

The pointer address may change but the address pointed by pointer will not change.

Disadvantage of Indirect Triple record structure

- It includes one more indirection level which is a hectic.

UNIT - IV

Type checking

4 types of errors are identified by the semantic analyzer

1. Type checks
2. Flow control checks
3. uniqueness checks
4. Name related checks

in order to leave language construct we need to have a flow of control

eg: switch construct

if in case break; is not used no case is shown, all cases are executed -- but within C#, .NET it shows an error

to conclude leave switch construct we need a flow of control to exit out of switch & transfer control to another statement

eg: B1 begin B2 begin and B2 end B1

each block should be ended by its name

for an operator to operate on operands as an operand or not

A variable can be declared of one datatype, but not more than one datatype

eg: int x; float x;

ambiguous eg: in switch each case '1' should be unique

within the Ada program ming language we use name for procedure

Procedure x
begin

end x;

checked whether these are defined or not.

even the nested blocks are supported by the Ada language

In order to perform type checking functionality a tool is used called Type checker.

It computes the datatype of language construct (It is a syntactical statement which is formed by using one or more tokens) and verifies meaning of these language construct according to the language rules.

eg: a float is a language construct.

but within the C language we cannot perform the modulus operation on the float variables.

but this is valid within the java programming language.

Type system is the main functionality that is implemented inside the type checker

Type system keeps track of datatype info of language construct and also assigns datatype to language construct. Also verifies whether it is according to the language rules or not.

In order to assign datatype to the language construct Type system uses an expression called type expression (depends on the language)

There are two types of type expressions

- basic type expression
- type construct type expression.

within the basic type expressions

- ① - boolean
- characten
- int
- real

② type_error: used to signal an error within the programming statement.

③ void: In absence of the datatype.

within the type construct type expression

* Type construct is formed by applying an operator to the other type expression

① Array (within pair)

Var A: array[1..50] of int

⇒ (In C)
int A[50];

Type expression for array variables

array(1..50, int)

② Cartesian product

T1 x T2

-- consider the formal parameters paired within the function definition

function (int a, float b, char c)

ε

These type expression is

int x float x char

③ within the c programming language we use the structure

Struct emp

```
{
  char A[10];
  int age;
}
```

3:

we use a type expression for structure

record((A x array (0..10, char) x (age x int))

This type expression uses cartesian product along with datatype and data field.

④ within the c programming language we use the pointers

```
int *p;
```

Then the type expression for the pointer is

pointer(int) | pointer(T);

```
int *P[10];
```

Then the type expression is

array (0..10, pointer(int))

⑤ functions

we use mathematical notation to represent the type expression for the functions

$D \rightarrow R$ (domain mapped to range)

eg: float sum(int a, int b)

E

3

The type expression is
int x int \rightarrow float.

write type expressions for following types

① An array of pointers to real, where the array index ranges from 1-100.

② A two dimensional array of integers (an array of arrays) whose rows are indexed from 0-9 and whose columns are indexed from -10 to 10

③ functions whose domains are functions from integers to pointer to integers and whose ranges are records consisting of an integer and character.

eg: \rightarrow expression

$P \rightarrow D; E$

declarations

This program only consists of declarations followed by expression

$D \rightarrow D; D$

$D \rightarrow id : T$

$T \rightarrow \text{char} / \text{integer} / \uparrow T / \text{array}[\text{num}] \text{ of } T$

pointers in Pascal are represented using \uparrow followed by variable in the Pascal language.

$E \rightarrow \text{literal} / \text{number} / id / E \text{ mod } E / E[E] / E^T$

$P \rightarrow D; E$ represents the program structure

So we write the semantic rules for the above grammar

$D \rightarrow id : T$

addtype(id, entry, T, type)

id stored within the symbol table.

$T \rightarrow \text{char}$

$\{ T.type = \text{char} \}$

$T \rightarrow \text{integer}$

$\{ T.type = \text{integer} \}$

$T \rightarrow \uparrow T$

$\{ T.type = \text{pointer}(T.type) \}$

$T \rightarrow \text{array}[\text{num}] \text{ of } T$ $\{ T.type = \text{array}(0..num, T.type) \}$

$E \rightarrow \text{char} \rightarrow \text{literal}$

$E \rightarrow \text{literal} \{ E.type = \text{char} \}$

$E \rightarrow \text{num} \{ E.type = \text{integer} \}$

$E \rightarrow id \{ E.type = \text{lookup}(id, entry) \}$

$E \rightarrow E_1 \text{ mod } E_2$

$\{ E.type = \text{if } E_1.type = \text{integer and } E_2.type = \text{integer then integer else type-error} \}$

$E \rightarrow E_1[E_2] \{ E.type := \text{if } E_2.type = \text{integer and } E_1.type = \text{array}(s, t) \text{ then } t_1 \text{ else type-error} \}$

Type conversion can be done in 2 ways:

(i) implicit type conversions

(ii) explicit type conversions

Equivalence of type expressions

Type expressions are built from basic types or formed by applying type constructors to other type expressions.

Two type expressions are said to be structurally equivalent if they are same basic types or formed by applying type constructors to other expressions.

They are also structurally equivalent if they are identical.

function: $equiv(s, t) : boolean$

begin
if s & t are same basic types then
return true;

else if $s = array(s_1, s_2)$ & $t = array(t_1, t_2)$
return $equiv(s_2, t_2)$

else if $s = s_1 \times s_2$ & $t = t_1 \times t_2$ then
return $equiv(s_1, t_1)$ & $equiv(s_2, t_2)$

else if $s = pointer(s_1)$ & $t = pointer(t_1)$ then
return $equiv(s_1, t_1)$.

else if $s = s_1 \rightarrow s_2$ & $t = t_1 \rightarrow t_2$ then
return $equiv(s_1, t_1)$ & $equiv(s_2, t_2)$

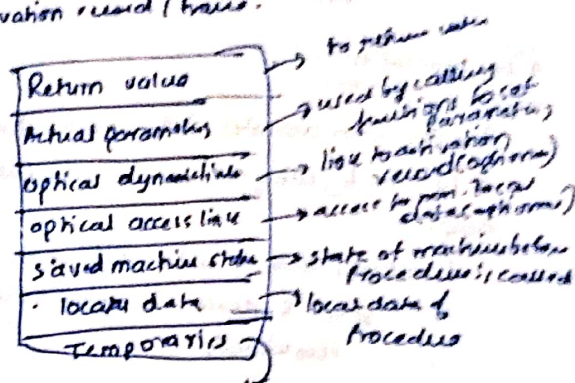
else
return false

end

Two variables are said to be name equivalent if they are associated with same type expressions.

Activation record/frame

Information needed for execution of a procedure is stored in a block of storage called an activation record/frame.



used to store values obtained during expression evaluation

Storage allocation strategies

① Static Allocation

- is done at compile time
- all data objects should be known at compile time
- no dynamic changes allowed

Limitations:

- recursive functions not allowed
- no dynamic data structure.

② Stack Allocation

manages runtime storage as a stack (control stack)

- Activation records are pushed & popped when procedure begins & ends respectively.
- local values are bound to fresh storage when activation record is pushed to stack
- are deleted when activation record is popped

Allocation tree

A tree which represents sequence of procedure calls is called Allocation tree.

eg: main() { f1(); }
f1() { f2(); }
f2() { f3(); }



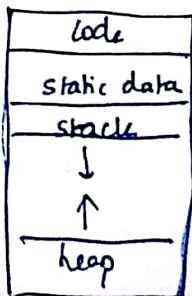
Runtime environment

Compiler obtains a block of storage from OS to run the target code.

They contain

1. generated target code
2. data objects known at compile time
3. Stack used to control execution of procedure

Run time memory is divided into



stack & heap size change dynamically

limitations:

- local values are not retained
- memory addressing can be done using pointers

Heap allocation :-

This allocation allocates & deallocates memory at runtime by using the memory area called heap.

It allocates memory when activation record starts & deallocates when activation record ends

limitations:

~~Heap~~ heap manage overhead

Parameter passing: Communication

3 ways of parameter passing

1. Call by value:

actual parameters are evaluated & values of actual parameters are passed to called procedure

Two steps

- Formal parameters are considered as local names to the called procedure so the storage of formal's is in activation record of procedure
- caller evaluates actual parameters and their r-values are stored in the storage of formal's.

eg: main () swap (int a, int b)
 { swap(a, b) { temp = x;
 } { x = y;
 { y = temp;
 }
 }

2. Call by reference: l-values of an expression are passed as actual parameters to the called procedure.

swap (&a, &b) swap (int *a, int *b)
 { { temp = *x;
 *x = *y;
 *y = temp;
 }

3. copy version 1

- This is a hybrid mechanism of call by value & call by reference.
- Also called as copy-in & copy-out

1. Before the control transfers to the called procedure, actual parameters and their r-values are stored in the storage of the formal parameters.
2. when the control returns the current r-value of formal values of formal parameters are copied into the l-value of actual parameters.

program copy (input, output)

var a: integer

procedure display (x: integer)

begin

a := x;

a := 0;

end

a := 1

display(a);

writeln('a', a);

end

Symbol table: is a data structure which consists of records of each identifier with fields for the value of the identifier

1. variable names
2. constants
3. compiler generated temporaries
4. function names
5. scope information

symbol table identifier format

Name	Information

Loop optimization :

(i) Code motion (or) loop Invariant (or) Frequency reduction

Reduction

- Reducing the number of instructions under loop or reducing the executing time of loops will speed up the program

Code motion :- an expression value that is not changing within the loop, then moving that expression to beginning before loop

```

1 limit = 100
2 while (i <= limit - 2)
3     i++;
    
```

Also called loop invariant.

And also called as frequency reduction

This technique skips the evaluation of the loop

After applying Code motion

```

1 t = limit - 2
2 while (i <= t)
3     i++;
    
```

(ii) Induction variable elimination

Induction variable : A variable which have the following form

$$I = I \pm C$$

eg: consider $C=1$ and $I=1$

∴ I value increments linearly
i.e. 1, 2, 3, 4, 5, 6, 7

eg: (3) $T_i = 4 * i$

```

1
2
3
4
i = i + 1
if i <= 20 goto (3)
    
```

i	T _i
1	4
2	8
3	12
4	16

In here T_i is also linearly incrementing within the loop.

" If two variables like above are linearly incrementing, so we replace these two variables by a single variable

∴ After applying

Induction variable elimination

```

1 Ti = 0;
2 (3) Ti = 4 + Ti
3     ...
4     if Ti < 80 goto (3)
    
```

we eliminate inside the

" And the variables can be in linear increment or linear decrement in order to perform the induction variable elimination "

(iii) Strength reduction

- Consider within the loop the same expression like

$$a * 2 = a * a$$

$$a * 2 = a + a$$

If the induction variable elimination is performed it also covers the strength reduction optimization in loops

So in place of performing multiplication, we replace with addition and we replace division operator by subtraction

replacing a high level operator with low level operator is called strength reduction

(iv) Loop unrolling

- "replicating the body of the loop in order to reduce required number of test conditions is called loop unrolling"

consider

```

1 I = 1;
2 while (I <= 100)
3     a[I] = 0;
4     I++;
    
```

replacing a variable with another variable is called copy propagation

upto how many lines the variable is going to be alive is live variable analysis

dead code elimination is elimination of unused variables and statements

```

eg:
  I = 1, 3
  while (I <= 100)
  {
    a[I] = 0;
    I++;
    a[I] = 0;
    I++;
  }
  
```

The no. of times the conditions executed is reduced to 50.

This only works when number of iterations are constant.

(V) loop jamming:

merging the two bodies of the loop, if the number of test conditions, number of indices are same.

```

eg:
  for (i=0; i<10; i++)
  {
    for (j=0; j<10; j++)
    {
      a[i][j] = 0;
    }
  }
  for (i=0; i<10; i++)
  {
    a[i][i] = 1;
    //making diagonals = 1
  }
  
```

After applying loop jamming

```

for (i=0; i<10; i++)
{
  for (j=0; j<10; j++)
  {
    a[i][j] = 0;
    a[i][i] = 1;
  }
}
  
```

Principle sources of optimization * * *

(or)
machine independent code optimization

Basic Blocks

It is a sequence of consecutive statements that in "control enters at the beginning, once entered all the statements executed sequentially without halt (or) possible branching".

Intermediate code generation

Three address code

before performing machine independent code optimization

The code optimizer generates basic blocks for three address code

Code optimizer

Algorithm Constructing basic blocks :-

Input: Three address code

Output: List of basic blocks with each three address statement exactly in one block.

Procedure:

① we determine leader, the first statement of basic block following are the rules

→ [The first statement is a leader]

→ [Any statement targeted by conditional or unconditional jump statement is a leader]

→ [Any statement that follows a conditional jump statement is a leader]

② For each leader we construct a basic block which consist of first statement as leader and all statements upto not including next leader.

→ [If there is no next leader, then it is the end of program]

Flowgraph: The relationship between basic blocks represented by a directed graph is called as a flow graph.

- Flow graph consists of nodes and edges
- each node is a basic block
- There is a distinguished node called initial block whose leader is a first statement of three address statements.

→ There is an edge from block B_1 to block B_2

(i) In block B_1 last statement is a conditional jump statement, targeting to the first statement of block B_2 , then we place an edge from block B_1 to block B_2

—(ii) Block B_2 follows block B_1 in the flow of execution and then we place an edge

(iii) If there is an edge from block B_1 to block B_2 B_1 is a predecessor of B_2 and the block B_2 is successor of block B_1

