

Compiler Designs

Compiler — Converts Source Code to byte Code
(High level) (Low level)
— Reports error.

Highlevel language Consists of preprocessing statements.
but it does not compile by compiler.

Highlevel \longrightarrow Assembly level \longrightarrow Relocatable language

\Downarrow
No fixed location in memory.

* Linker links all subroutines, loaders loads at fixed location.

* Compiler has 6 phases for converting purely high level \rightarrow target.

1) Lexical $\xrightarrow{\text{Analysis}}$ scans, splits as tokens.

2) Syntan \longrightarrow checks the syntan.

3) semantics

4) Intermediate Code

5) Code Optimization.

6) Code Generation

} related to language.

} Dependent on Machine,
Independent of language.

c =
a t
b * 20

Intermediate Code — (Almost 3 intermediateries)

$$T_1 = id_3 * 20$$

$$T_2 = id_2 + T_1$$

$$id_4 = T_2$$

Address Code

Code optimization —

$$T_1 = id_3 * 20$$

$$id_4 = id_2 + T_1$$

Code Generation —

Instruction	Source	destination
movf	id_3	R_1
mulF	20	R_1
movf	id_2	R_2
addf	R_2	R_1

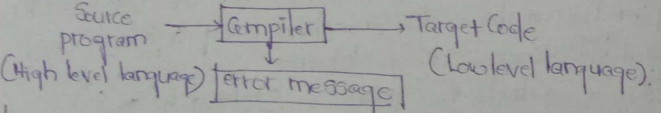
movf R_1 id_1

These 6 phases are connecting to error handler.

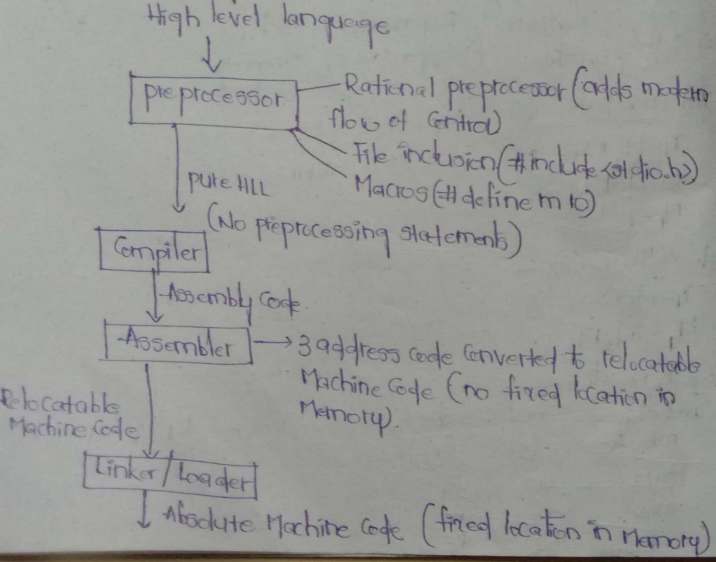
Explain about phases
 Lexical Analysis is Scanner / Linear Analysis
 Syntax " " Hierarchical Analysis

Position = Initial + Final * 20

Compiler: Converts high level language into low level language & also reports the error.

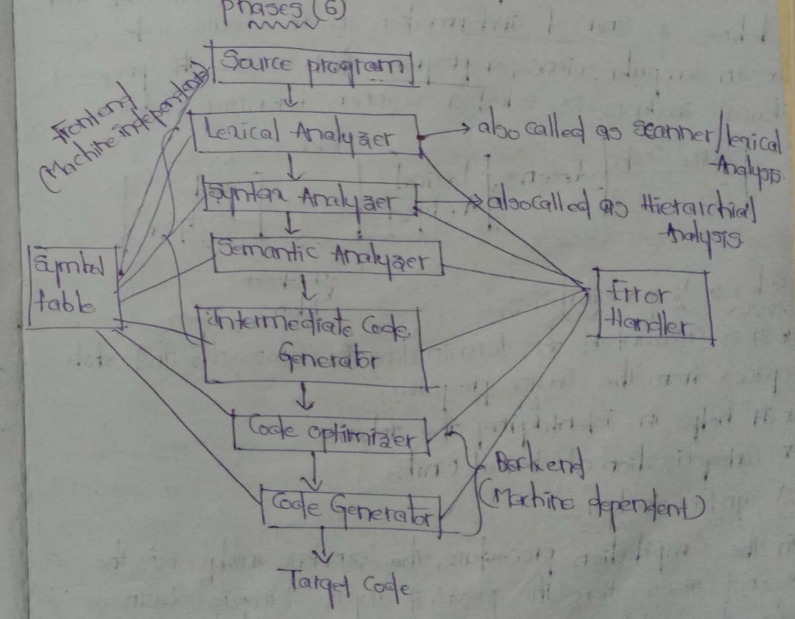


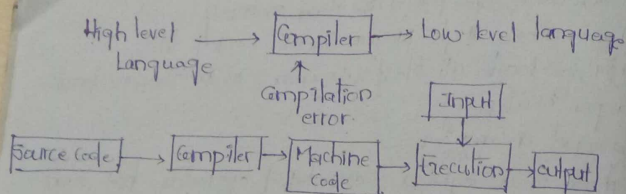
Language processing system:-



* all subroutines are placed at fixed memory.
 → Linker → allows single program from several files.
 → Loader → locates all files/subroutines at particular location or proper location.

* phases of a compiler:-
 * A phase is a logically interrelated operation that takes program in one representation & produce output in another representation (Target program) → HLL
 * there are 2 parts of compiler 1) Analysis (Machine Independent, language dependent)
 2) Synthesis (Machine dependent, language independent).



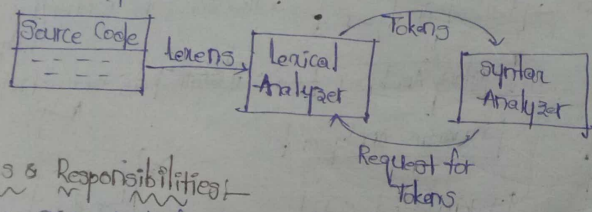


* Lexical Analysis:-

Lexical Analysis (or) Lexical Analyzer is the initial stage or phase of the compiler. This phase scans the source code and transforms the input program into a series of a token.

* A token is basically the arrangements of character that defines a unit of information in the source code.

* In Computer science, a program that executes the process of lexical analysis is called a scanner, tokenizer, or lexer.



Roles & Responsibilities:-

- * It is accountable for terminating the comments and white spaces from the source program.
- * It helps in identifying the tokens.
- * Categorization of lexical units.

2) Syntax Analysis

In the compilation procedure, the syntax analysis is the second stage. Here the provided input string is scanned for the validation of the structure of the standard

grammar. Basically, in the second phase, it analyses the syntactical structure and inspects if the given input is correct or not in terms of programming syntax.

* It is also known as parsing in a compiler.

* Roles and Responsibilities:-

- Note syntax errors.
- Helps in building a parse tree.
- Acquire tokens from the lexical analyzer.
- Scan the syntax errors, if any.

3) Semantic Analysis:-

In the process of compilation, semantic analysis is the third phase. It scans whether the parse tree follows the guidelines of language. It also helps in keeping track of identifiers and expressions. In simple words, we can say that a semantic analyzer defines the validity of the parse tree, and the annotated syntax tree comes to an output.

Roles and Responsibilities:-

- * saving collected data to symbol tables or syntax trees.
- * It notifies semantic errors.

* Scanning for semantic errors.

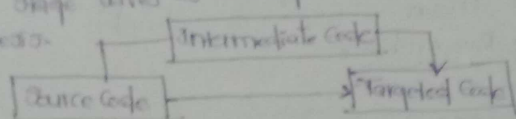
4) Intermediate Code Generation

The parse tree is semantically confirmed; now, an intermediate code generator develops three address codes. A middle level language code generated by a compiler at the time of the translation of a source program into the object code is known as intermediate code (or) text.

- * A code that is neither high-level nor machine code, but a middle-level code is an intermediate code.

Lexical Analysis code to machine code later.

- * We can translate this as a bridge or way from analysis to synthesis.

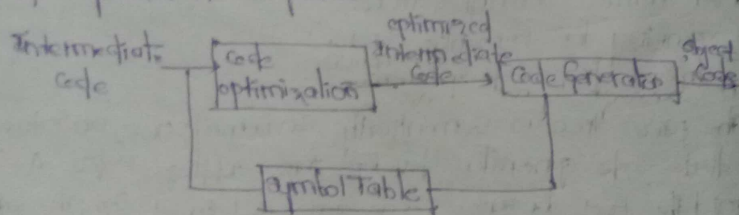


Roles and Responsibilities

- * Helps in maintaining the priority ordering of the source language.
- * Translate the intermediate code into the machine code.
- * Having operands of instructions.

Code optimizer

Now coming to a phase that is totally optional, and it is code optimization. It is used to enhance the intermediate code. This way, the output of the program is able to run fast and consume less space. To improve the speed of the program, it eliminates the unnecessary stuff of the code, and organizes the sequence of statements.



Roles and Responsibilities

- * Remove the unused variables and unreachable code.
- * Enhance runtime and execution of the program.
- * Produce streamlined code from the intermediate expression.

Code Generator

The final stage of the compilation process is the code generation process. In this phase, it tries to acquire the intermediate code as input which is fully optimized and map it to the machine code or language. Later, the code generator helps in translating the intermediate code into the machine code.

Roles and Responsibilities

- * Translate the intermediate code to target machine code.
- * select the allocate memory, opcod and registers.

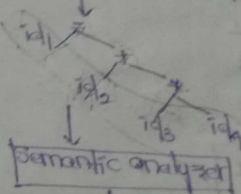
Ex

$$\text{position} = \text{Initial} + \text{Final Rate} \times 60$$

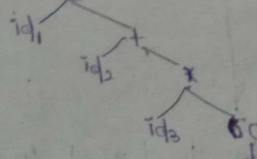
Lexical Analyzer

$$id_1 = id_2 + id_3 \times id_4$$

Syntax analyzer



Semantic analyzer



$$id_1 = id_2 + id_3 \times id_4$$

$$= id_2 + id_3 \times id_4$$

$$= id_2 + id_3 \times id_4$$

$$= id_2 + id_3 \times id_4$$

initial

Intermediate code generator

temp1 = int literal (60)
temp2 = id3 * temp1
temp3 = id2 + temp2
id1 = temp3

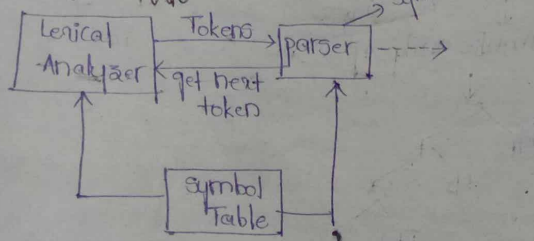
Code optimization

temp1 = id3 * 0.0
id1 = id2 + temp1

Code generation

Movf id3, R2
Mulf #0.0, R2
Movf id2, R1
Addf R2, R1
Movf R1, id1

* Lexical Analysis



parser sends requests for tokens to Lexical Analyser then send the tokens both are combinedly working but

both of them are separated for simplification & attached to symbol Tables.

* Lexical Analyser forms tokens based on Transition state.

* Transition diagrams are stored in Lexical Analyser.

* a = b || c

based on pattern it identifies the tokens whether it is identifier, keyword, etc.

* It is identified as lexeme.

* every part is identified by matching with existing pattern.

Tokens

Token is a sequence of characters that can be treated as a single logical entity. they are:

identifiers, keywords, operators, special symbols, constants

patterns

A set of strings in input for which the same token is produced as output.

* This set of strings is described by a rule called pattern associated with Token.

Lexeme

A lexeme is a sequence of characters in source program that is matched by the pattern for a Token.

ex: int a
Data type Identifier \Rightarrow existed patterns } Lexeme

in a

Not Data type Identifier \Rightarrow Not Matched.

Token	Lexemes	Pattern
Constant	constant	constant
if	if	if
Relation	=, <, >, <=, >=, <!, >!	<or <= or > or >= or <!= or >!=

id pi letter followed by tokens and digits.
 num 3.14 any numeric constants

* $E = M | C | * | /$ = 7 tokens
 $C = \text{"Hello"}$ \Rightarrow 3 tokens

* Error Recovery — $E \rightarrow fi$ error
 * If is panic mode recovery.

Simple strategy.

- * in panic mode error recovery we delete successive characters from the remaining input until the lexical analyser can find a well formed token / suitable pattern token.
- * other possible error recovery actions:—
 - i) Inserting a missing character. ($int \rightarrow it$)
 - ii) Replacing an incorrect character by correct character. ($int \rightarrow sint$)
 - iii) Transposing two adjacent characters. ($if \rightarrow fi$)
 - iv) Delete number of characters ($E \rightarrow int // \rightarrow int$)

Rule
 $(letter + \text{underscore}) (letter + digit + alphanumeric character)^*$
 (this is regular expression for identifiers)

Start \rightarrow (90) $\xrightarrow{\text{letter, ' - '}}$ (91) $\xrightarrow{\text{Final}}$ (96)

* Lexical-Analyser removes spaces.

Lex Tool

* Lex Tool is a lexical compiler
 No. of Lex Tools \rightarrow flex depends on regular expressions.
 Initially we write lex program that is compiled by lex tool and gets output as C program, c file.
 Extension $\rightarrow .l$, produces $\rightarrow \text{lex.yy.c}$
 c file is compiled by C compiler and gets object file.

object file produces stream of tokens.

- * there is a wide range of Tools for construction of lexical analyser the majority of these tools are based on regular expressions.
- * one of the traditional tool of that kind is lex and another tool is flex steps
- i) A specification of the lexical analyser is prepared by creating a program lex.l in the lex language.
- ii) the lex.l is run through the lex compiler to produce a C program lex.yy.c
- iii) The program lex.yy.c consists of a tabular representation of a transition diagram constructed from the regular expression of lex.l
- iv) Finally lex.yy.c is run through the C compiler to produce an object program a.out which is the lexical analyser that transforms an input stream into a sequence of tokens.

$\text{lex.l} \rightarrow \boxed{\text{lex Compiler}} \rightarrow \text{lex.yy.c}$
 $\text{lex.yy.c} \rightarrow \boxed{\text{C Compiler}} \rightarrow \text{a.out}$

Input stream $\rightarrow \boxed{\text{a.out}} \rightarrow$ stream of Tokens

- * lex specifications—
- * A lex program consists of three parts.
 - i) Declaration
 $\% \%$ \rightarrow used to separate the parts.
 - ii) Translation Rules
 $\% \%$
 - iii) Auxiliary procedures

i) Declaration:-

The declaration section includes declarations of variables, constants, and regular definitions.

ii) Translation Rules:-

The Translation Rules of a lex program (or) statements of the form $P_1/P_2/P_3/\dots/P_n$

P_1 (action 1)

P_2 (action 2)

P_3 (action 3)

P_n (action n)

(Based on regular expressions we write action in C language)

iii) Auxiliary procedures:-

Whenever Auxiliary procedures are needed by the actions and these procedures can be compiled separately and loaded with lexical analyser.

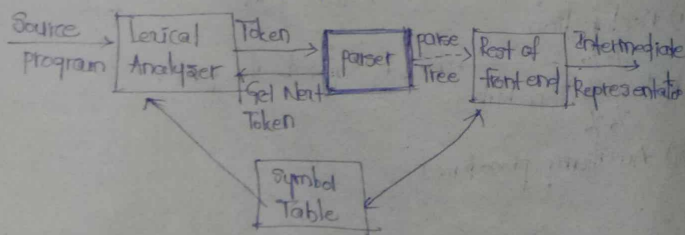
* These procedures are already pre-compiled.

* Here we use only Context free Grammar.

* Right linear, left linear, Ambiguity, Remove left Recursion.

* Syntax Analysis:- (ii)

* The parser obtains a string of tokens from lexical analyser and check whether the string can be generated by the grammar for source language.



In parser it starts i/p from left \rightarrow right only one value.

* Top-Down parsing Techniques \Rightarrow Root \rightarrow yield (leaf Nodes)

* Bottom-up parsing Techniques \Rightarrow Yield \rightarrow Root (leaf Node)

* A parser for programming language always make a single left to right scan over the i/p. Everytime it take one at a time. A parser can be constructed for any grammar. The methods commonly used in compilers are classified as

i) Top-Down parsing

ii) Bottom-up parsing

* Syntax analyser also have error handling.

Assignment \rightarrow Error handling Techniques in syntax analysis or in a parser.

Context free Grammar:-

Definition.

* Derivations.

$G: S \rightarrow SS + / SS * / a$

$G(\{S\}, \{+, *, a\}, P, S)$

$w = aa + a *$

$S \rightarrow SS *$

$\rightarrow SS + S *$

$\rightarrow SS / S +$

$\rightarrow a aa + a *$

$\Rightarrow aS / Sa / a$

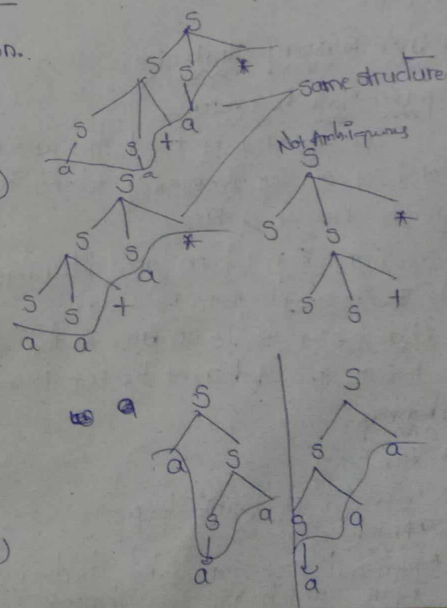
$w = aS$

asa

aaa

(Different structures)

(Ambiguous)



Closure properties for all

Assignment

* Syntax Analyzer or parser:-

Analyzes the syntactical structure and checks if the given input is in the correct syntax of the programming language or not.

* Syntax Error Handling:-

The error handler in a parser has goals that are simple to state but challenging to realize:-

* Report the presence of errors clearly and accurately.

* Recover from each error quickly enough to detect subsequent errors.

* Add minimal overhead to the processing of correct programs.

* Error Recovery Strategies:-

i) Panic-Mode Recovery:-

once an error is found, the parser intends to find a designated set of synchronizing tokens by discarding input symbols one at a time.

* Synchronizing tokens are delimiters, semicolon or } whose role in source program is clear.

* When parser finds an error in the statement, it ignores the rest of the statement by not processing the i/p.

Advantage:-

* Simplicity

* Never get into infinite loop.

Disadvantage:-

* Additional errors cannot be checked as some of the input symbols will be skipped.

ii) Phrase Level Recovery:-

* When a parser finds an error, it tries to take corrective measures so that the rest of inputs of statement allow the parser to parse ahead. The corrections may be

* Replacing a prefix by some string.

* Replacing Gamma by Semicolon.

* Deleting extraneous Semicolon.

* Inserting missing Semicolon.

Advantage:-

* It can correct any input string.

Disadvantage:-

* It is difficult to cope up with actual error if it has occurred before the point of detection.

iii) Error productions:-

* The use of the error production method can be incorporated if the user is aware of common mistakes that are encountered in grammar in conjunction with errors that produce erroneous constructs.

Ex:- Write 5r instead of 5*

Advantage:-

* If this is used then, during parsing appropriate error messages can be generated and parsing can be continued.

Disadvantage:-

* The disadvantage is it is difficult to maintain.

iv) Global Correction:-

* The parser considers the program in hand as a whole and tries to figure out what the program is intended to do and tries to find out a closest match for it, which is

error-free.

* When an erroneous i/p statement X is fed, it creates a parse tree for some closest error-free statement Y .

Advantage:-

* This may allow the parser to make minimal changes in the source code.

Disadvantage:-

* Due to the complexity (time and space) of this strategy, it has not been implemented in practice yet.

* Top-Down parsing Techniques:-

* Backtracking is also called Brute Force Method.

i) Brute-Force Method (Back-Tracking):-

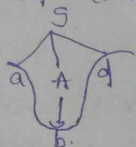
G:- $S \rightarrow aAd/aB$

$A \rightarrow b/c$

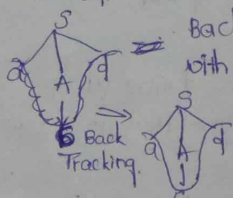
$B \rightarrow ccd/ddc$

We fail to produce from first production again we start Backtracking.

$w = abd$ (Starts input from left to right)



$w = acd$

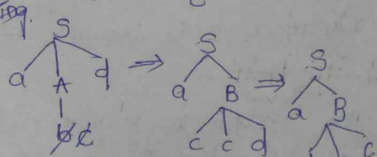


Backtracking is possible with recently used one

$w = addc$



Backtracking



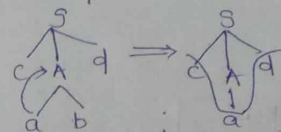
Error occurs when we fail to produce given string.

* Backtracking is possible based on recently used

G:- $S \rightarrow cAd$

$A \rightarrow ab/a$

$w = cad$



(Not match with the i/p string so we choose alternative production Rule)

ii) Recursive descent parser:- (Repeat same function)

A parser that uses a set of recursive procedures to recognize its input with no backtracking it is called a Recursive descent parser.

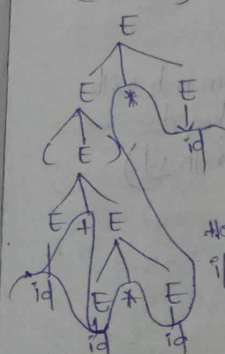
G:- $E \rightarrow E+E$

$E \rightarrow E * E$

$E \rightarrow (E)/id$

Ex:- $(2+3*4)*5$

$w = (id + id * id) * id$



Here we have different structures for same i/p & same string \Rightarrow Ambiguous.

$A \rightarrow p/pz$

$z \rightarrow q/az$

Removing left recursion:-

$G: E \rightarrow E'I/I'$

$T \rightarrow T'f/f'$

$F \rightarrow (E)/id$

$E \rightarrow TE'$

$E' \rightarrow +TE'/E \Rightarrow (z/z)/a/z$

$T \rightarrow FT'$

$T' \rightarrow *FT'/e$

$F \rightarrow (E)/id$

Rule:-

$A \rightarrow A\alpha/\beta(i)$
 $A \rightarrow \beta A' (ii)$
 $A' \rightarrow \alpha A'/\epsilon (iii)$

①

$G_2: S \rightarrow as/sa/a$

Rules:-

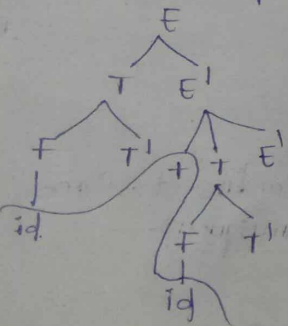
$S \rightarrow sa/as/a$ (i) $A \rightarrow A\alpha/\beta$

$S \rightarrow ass'sa'$ (ii) $A \rightarrow \beta A'$

$S' \rightarrow as'/\epsilon$ (iii) $A' \rightarrow \alpha A'/\epsilon$

First we write procedures for Non-terminals which are at left side. (from G_1)

$w = id + id$ string is stored in Buffer.



(if character is match with procedure we take, else we should replace with id)

Terminal \rightarrow i/p buffer symbol

Non-terminal \rightarrow procedure.

code for eg, ①

procedure $E()$

begin

$T()$;

$EPRIME()$;

end;

procedure $EPRIME()$

if input_symbol = '+' then

begin

$ADVANCE()$;

$T()$;

$EPRIME()$;

end;

procedure $T()$

begin

$F()$;

$TPRIME()$;

end;

procedure $TPRIME()$

if input_symbol = '*' then

begin

$ADVANCE()$;

$F()$;

$TPRIME()$

end;

procedure $F()$

if input_symbol = 'id' then

$ADVANCE()$;

else if input_symbol = 'c' then

begin

$ADVANCE()$;


```

E();
if input.symbol = ']' then
  ADVANCE();
else ERROR();
end;
else ERROR();

```

* To avoid the necessity of a Recursive language we shall also consider a Tabular implementation of recursive descent called predictive parsing. It is also called Non-recursive parsing.

* predictive parsing:— (Non-recursive predictive parser)

(or) Table driven predictive parser

We use stacks & some operations, Table.

Rules:-

i) The predictive parser has an input, stack and a parsing Table & o/p.

ii) The i/p contains the string to be parsed, followed by dollar (end marker \$)

iii) The stack contains a sequence of grammar symbols (variables & terminals) preceded by dollar '\$' bottom of stack

iv) Initially the stack contains the start symbol of the grammar preceded by dollar.

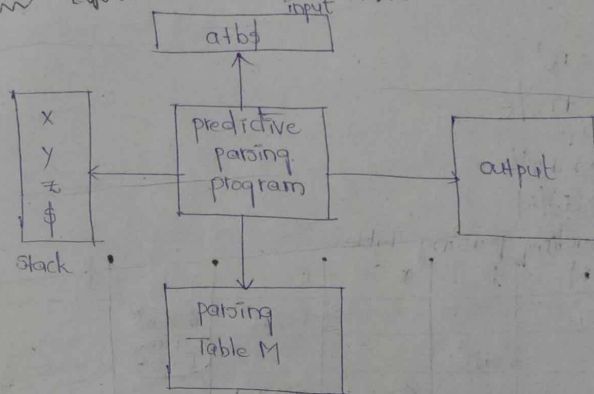
v) The parsing Table is a two dimensional array, $M[A, a]$ where A is a Non-terminal & a is a terminal. (or) the symbol '\$'.

$S \rightarrow SS/A$
 $A \rightarrow a$

	a	$\$$
S		
A		

rows \Rightarrow Non-terminals
columns \Rightarrow Terminals

Block Diagram of a predictive parser:-



Rules:-

* X the symbol on the top of the stack and a is current i/p symbol. These two symbols determine the action of the parser.

i) if $X = a = \$$ (stack & i/p is end). the parser halts and announce successful completion of parsing.

ii) if $X = a \neq \$$, (Top of the stack & i/p symbol are same). the parser pops X of the stack and advances the i/p pointer to the next i/p symbol.

iii) if X is a Non-terminal the program consult every $M[X, a]$ if X & i/p symbol is terminal & both are different.

of the parsing Table M . this entry will be either an X production of the grammar (or) an error entry.

\rightarrow if $M[X, a] = \{X \rightarrow uv\}$ the parser replaces X on top of the stack by uv .

→ if $M[x, a] = \text{error}$, the parser calls an error recovery routine.

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

D) Constructing parsing Table.

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

$E \rightarrow TE'$ i.e. $E \rightarrow TE'$
 $T \rightarrow FT'$ $\rightarrow FT'E'$
 $F \rightarrow (E) / id$ $\rightarrow (E)T'E' / idT'E'$

First of $E \rightarrow E \rightarrow TE'$ (only left side) $\Rightarrow E \rightarrow TE'$
 \uparrow
 First of E
 Follows of $E \rightarrow$ i.e. only on right side
 $F \rightarrow (E) / id$
 \uparrow
 $\Rightarrow E \rightarrow TE'$
 $\rightarrow FT'E'$
 $\rightarrow (E)T'E' / idT'E'$
 \uparrow
 $\Rightarrow E \rightarrow TE'$
 $\rightarrow FT'E'$
 $\rightarrow (E)T'E' / idT'E'$
 \uparrow
 Follows of $E = \{ (, id \}$
 $\Rightarrow E \rightarrow TE'$
 $\rightarrow FT'E'$
 $\rightarrow (E)T'E' / idT'E'$
 \uparrow
 Follows of $E = \{ (, id \}$

When there is no followed Non-Terminal at particular Non-terminal i.e.

$E \rightarrow TE'$

(not follows any Non-terminal)

So we should take follows of

i.e. $E' \rightarrow$ has follows of E

So we should replace Follow(E') = $\{ (, id \}$ as Follow(E)

if $E \rightarrow TE'$

(follows any terminal we should replace same thing as follows)

* The Constructive of a predictive parser is aided by two functions associated with a grammar G . These functions, FIRST and FOLLOW allow us to fill in the entries in the predictive parsing table for G , whenever possible.

1) FIRST :-

if α is any string of grammar symbols, the First(α) be the set of terminals that begins strings derived from α . if $\alpha \rightarrow \epsilon$, then ϵ is also in First(α)

To Compute FIRST(X) for all grammar symbols X , apply the following rules until no more terminals of ϵ can be added to any FIRST SET

i) if X is terminal, then First(X) is $\{X\}$
 ii) if X is non-terminal and $X \rightarrow \alpha\beta\gamma$ is a production, then add α to First(X).

if $X \rightarrow \epsilon$ is a production, then add ϵ to First(X)

iii) if $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production, then for all i such that Y_i is not the first of Y_i

$\gamma_1, \dots, \gamma_{i-1}$ are non-terminals and $FIRST(\gamma_j)$ contains ϵ for $j=1, 2, \dots, i-1$, (i.e. $\gamma_1, \gamma_2, \dots, \gamma_{i-1} \rightarrow \epsilon$) add every non- ϵ symbol in $FIRST(\gamma_i)$ to $FIRST(X)$.

FOLLOWs:-

For non terminal A , $FOLLOW(A)$ to be the set of terminal a , that can appear immediately to the right of A in some sentence. Form i.e. $S \rightarrow \alpha A \beta$ for some α and β to compute $FOLLOW(A)$. For all non-terminals A , Apply the following rules until nothing can be added to any FOLLOW sets.

- i) ϕ is in $FOLLOW(S)$, where S is the start symbol.
- ii) if there is a production $A \rightarrow \alpha \beta$, then everything in $FIRST(\beta)$ except for ϵ , is in $FOLLOW(A)$.
- iii) if there is a production $A \rightarrow \alpha \beta$ or a production $A \rightarrow \alpha \beta \gamma$ where $FIRST(\beta)$ contains ϵ , then everything in $FOLLOW(A)$ is in $FOLLOW(B)$.

* Construction of a predictive parsing Table:-

Input:- Grammar G .

Output:- parsing Table M .

Method:-

- 1) For each production $A \rightarrow \alpha$ of the grammar do the following.
 - i) For each terminal a in $FIRST(\alpha)$ add $A \rightarrow \alpha$ to $M[A, a]$.
 - ii) If ϵ is in $FIRST(\alpha)$, add $A \rightarrow \alpha$ to $M[A, b]$ for each terminal b in $FOLLOW(A)$.
 - iii) If ϵ is $FIRST(\alpha)$ and ϕ is in $FOLLOW(A)$, add $A \rightarrow \alpha$ to $M[A, \phi]$.
- 2) Make each undefined entry of M be an error.

$G:-$

$E \rightarrow TE'$
 $E \rightarrow \#TE'/\epsilon$
 $T \rightarrow Ft'$
 $T' \rightarrow *Ft'/\epsilon$
 $F \rightarrow (E)/id$

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

$FIRST(E) = \{ \epsilon, \# \}$

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

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

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

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

* Based on these values we write production rules.

* FOLLOWs:-

i) $FOLLOW(E) = \{ \phi,) \}$

$FOLLOW(E') = \{ \phi,) \} \Rightarrow \text{Rule (ii)}$

$FOLLOW(T) = \{ +, \phi,) \}$

$FOLLOW(T') = \{ +, \phi,) \}$

$FOLLOW(F) = \{ *, +, \phi,) \}$

$\rightarrow \text{Rule ii}$

Rule (i) $E \rightarrow TE'$

$F \rightarrow (E)/id$ & Rule (ii)

follow

if we get ϵ again we write $FOLLOW(E')$

$\Rightarrow FOLLOW(T)$

* Stack Implementation:-

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

γ
 $\gamma =$
 $\gamma \rightarrow \gamma$
 γ