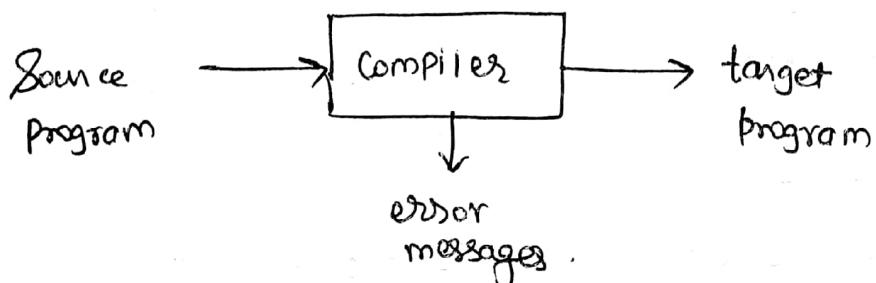


Introduction to Compilers:

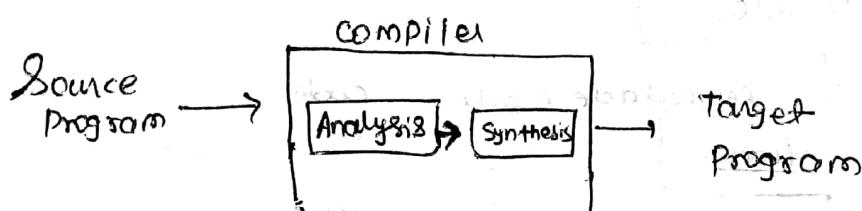
A compiler is a program that reads a program written in one language (i.e. the source program) and convert it into an equivalent program in another language i.e. the target language.



The Compiler takes a Source program such as higher level lang such as C, PASCAL, FORTRAN & converts it into low level lang. or machine level lang such as assembly lang.

Structure of Compiler:

- The compilation can be done in two parts: Analysis and Synthesis. In analysis part the source program is read and broken down into constituent pieces, and creates an intermediate representation of the source program.
- The synthesis part constructs the desired target program from the intermediate representation.



Properties of Compiler:

- The compiler itself must be bug free.
- It must generate correct machine code.
- The generated machine code must run fast.
- The compiler must be portable.

* Analysis of the Source program:

The Source program can be analyse in three phase.

1) Linear analysis:

In this type of analysis the Source String is read from left to right and grouped into tokens. e.g. Token for the lang. can be identifiers, Constants, relational operations, Keywords.

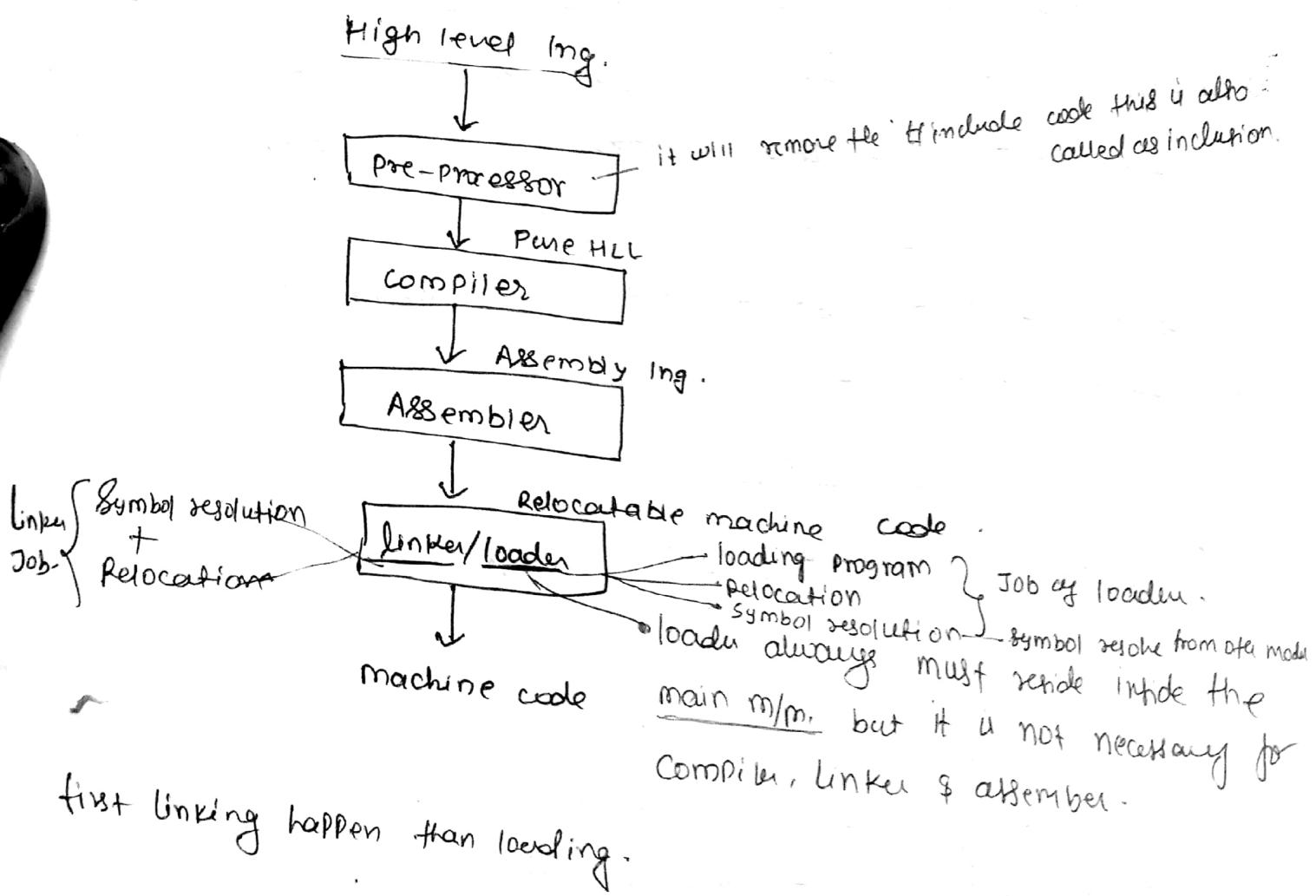
2) Hierarchical Analysis:

In this analysis, character or tokens are grouped hierarchically into nested collections for checking them Syntactically.

3) Semantic Analysis:

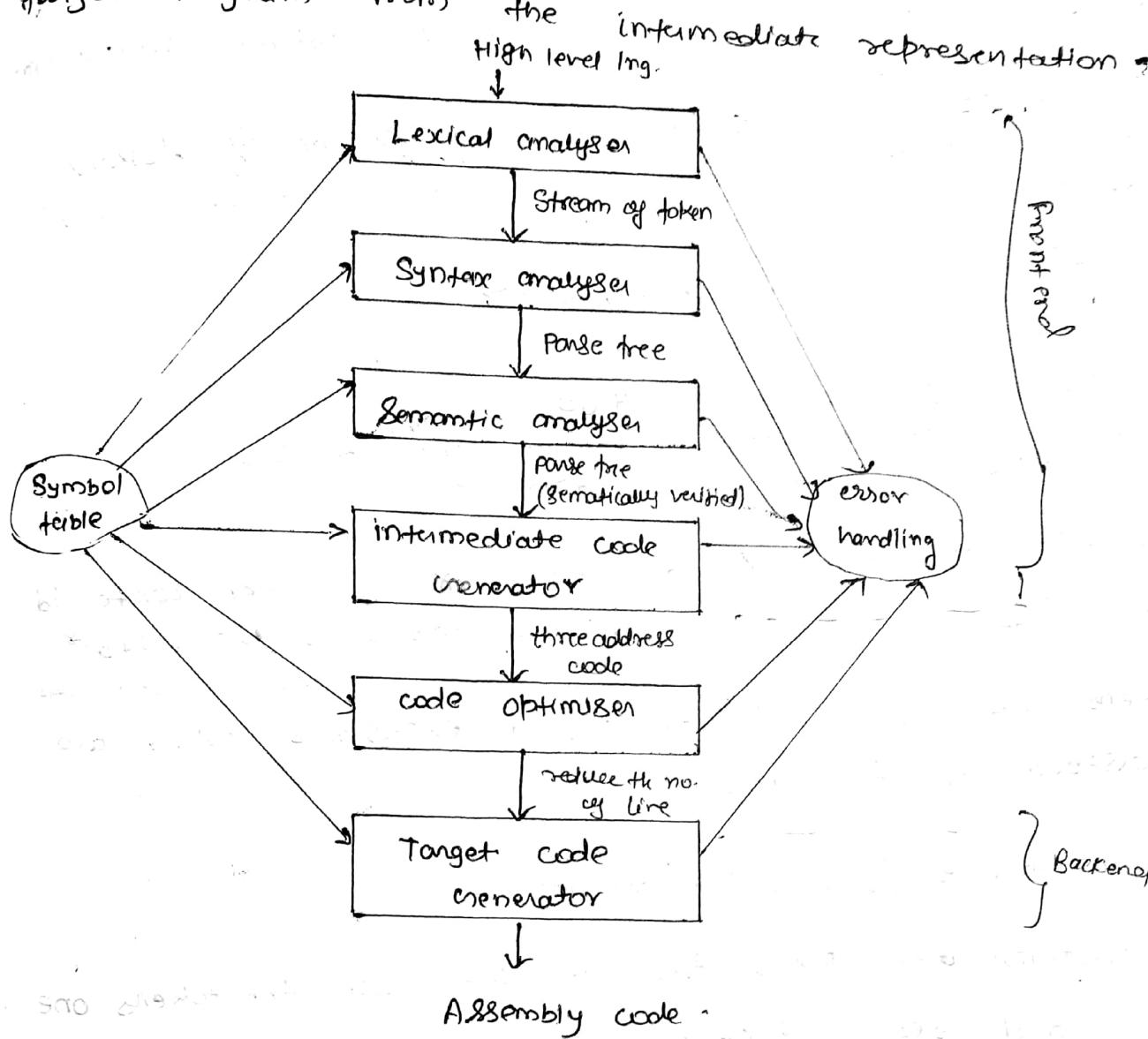
This kind of analysis ensure the correctness of meaning of the program.

Language Processing System:



Phase of a Compiler

- Basically two phases of Compilers, ~~Analyse~~ Analysis & Synthesis phase.
- Analysis Phase creates one intermediate representation from the given source code. Synthesis Phase creates an equivalent target program from the intermediate representation.



Symbol table :

- it is a data structure being used and maintained by the compiler, consist of all the identifiers name along with their types.

error handling

- Each phase can encounter errors after detecting an error, a phase must somehow deal with that error. Syntax & Semantic analysis Phases usually detect/handle a large fraction. The lexical phase can detect errors where the character remaining in

int a=5;
 a = 5; int b;
 {name value attribute}

Insert
look up

linklist
Hash table

the up don't form any token of the lang.

lex

* Lexical Analyzer: (lex tool)

- it reads the program and convert it into tokens. The char sequence forming a token is called the lexeme for the token. token for the lang. can be identifier, constant & relation op..
keyword, separators (, ;) etc.
- it convert stream of lexeme into a stream of tokens.
- non tokens (comments, Preprocessor directive, macro, blanks, tabs, newline etc.).

<token name, attribute, value>

e.g. POSITION \equiv initial + rate * 60

Here Position is the lexeme that could be mapped into token $<\text{id}_1>$

$<\underline{\text{id}_1}><= ><\underline{\text{id}_2}><+><\underline{\text{id}_3}><*><60>$

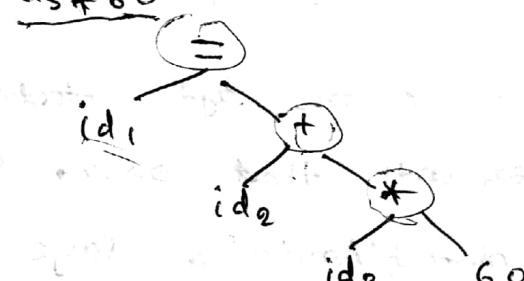
- it alphabet followed by alphabets or digit then write id (identifier) otherwise same. regular expression for id = $l(l+d)^*$ letter followed by letter or digit.
- Tokens are defined by the regular expressions which are understood by the lexical analyzer.

* Syntax Analyzer / (Parser): (tool yacc)

- It constructs the parse tree. it takes all the tokens one by one and uses context free grammar to construct a parse tree.

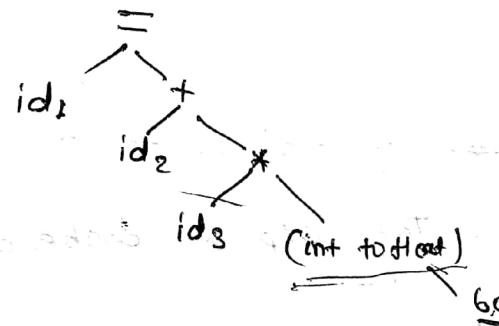
In a parse tree / Syntax tree all the internal nodes represents an operation & children of the nodes represents argument of the operation.

e.g. $\text{id}_1 = \text{id}_2 + \text{id}_3 * 60$



* Semantic Analyzer:

- it verifies the parse tree, whether it's meaningful or not.
- it further more produces a verified parse tree
- it may be classified into category type checking, type conversion
- it checks whether the program satisfies semantic constraints or not? e.g. like array index only represented using int.
 $\Rightarrow \text{array}[u] \vee \text{array}[x]$
- The language specification may permit some type conversion like int to float, float to int etc.



* Intermediate code generator:

- it generates intermediate code that is a form which can be readily executed by machine.

eg three address code etc. ^{imp. memory line code} intermediate code is converted to machine language using the last two phases which are platform dependent.

- it is easy to produce and easy to translate into the target program.

eg. temp1 = Int to real(60)

temp2 = id3 * temp1

temp3 = id2 + temp2

id1 = temp3

T

* Code optimizer:

- it attempts to improve the intermediate code so that it consumes fewer resources and produces more speed.
- The meaning of the code being transformed is not altered.
- Optimization can be categorized into two types machine dependent and machine independent.

eg. $\text{temp1} = \underline{\text{id}_3 * 60}$,

$\underline{\text{id}_1 = \text{id}_2 + \text{temp1}}$

* Code generator: (LANCE ^{TOOLS})

- its main purpose is to write a code that the machine can understand. The O/P is depends on the type of assembler.
- i.e. it converts the code into Assembly lang. i.e. Low level lang.

eg.

~~mov A~~ ~~6000~~ LDF R₂, id₃

mov

MULF R₂, R₂#60.0

LDF R₁, id₂

ADD F R₁, R₂

STF id₁, R₁

(method of optimization)

eg. Find the total no. of tokens.

~~tokens~~

int main()

{

1/2 variable

int a,b;

a=10;

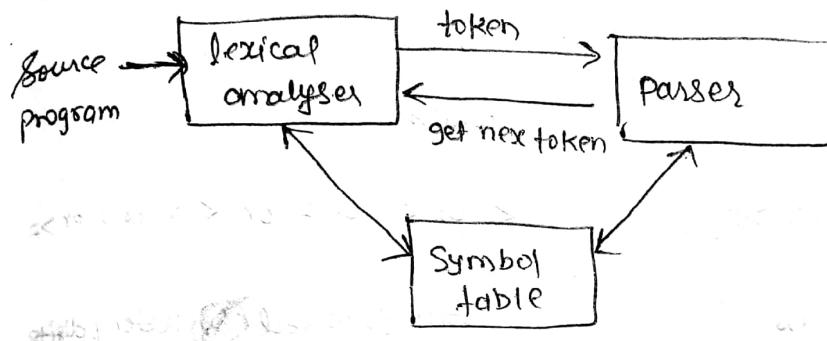
return 0;

}

tokens are.

"int main() int a", "b";", "a", "=", "10", ",", "return 0", ";" }"

- The Role of the Lexical Analyzer:
- its main task is to read the input character and produce as O/P a sequence of token that the parser uses for syntax analysis.
- upon receiving a "get next token" command from the parser, the lexical analyzer reads ip characters until it can identify the next token.



- lexical analyzer is a part of the compiler so that it may also perform certain secondary tasks at the user interface.
- one of the such tasks is skipping out from the source program comments and white space in the form of blank, tab, and newline characters.
- if any error will present then LA will correlate that error with source file and line number.

Issues in lexical analysis:

- simpler design
- compiler efficiency is improved
- compiler portability is enhanced.

Tokens, Patterns, and lexemes:

- A token is a set of string over the source alphabet.
- A pattern is a rule that describes that set.
(Set of string is described by a rule called a pattern.)

- A lexeme is a sequence of characters matching that pattern.

e.g. in Pascal for the Stmt -

Const Pi = 3.1416

The substring Pi is a lexeme for the token Identifier.

Token

Sample lexemes

informal description of pattern

if

if

if

while

while

while

relation relop

<, <=, =, >, >=

< or <= or = or < > or >=

fd

Count, Sum, i, j, pi

letter followed by letter & digits

Num

0, 12, 3.141 ...

Any numeric constant

e.g. printf("Sum=%d\n", total);

- printf and total are lexemes matching the pattern for token id,
- "Sum=%d\n" is lexeme matching literal.
- Regular expressions are used to specify lexeme patterns.

Input Buffering:

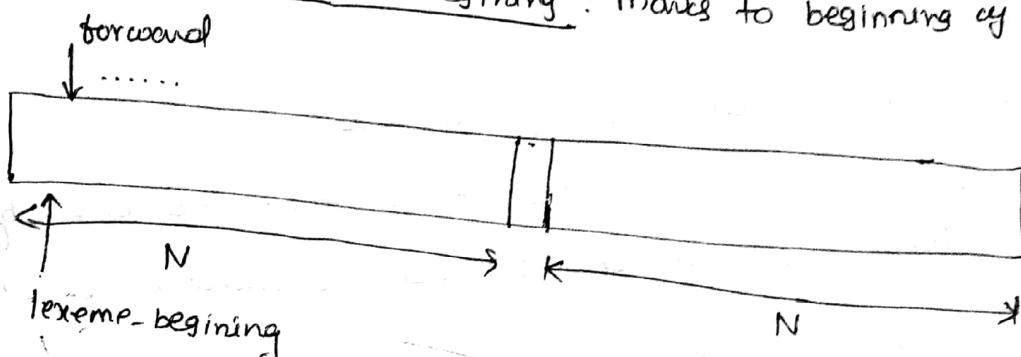
Speed of the lexical analysis is concern.

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

- We have two buffer I/P scheme that is useful when look ahead is necessary
 - Buffer Pair
 - Sentinels

A identifier would start with a letter followed by zero or more letters, digits and underscores.

- reading the source program can be split up by a token buffer scheme. each buffer is of size N and N usually the size of disk (N is no. of char. eg 3024 or 4096)
- If lesser than N character remaining in the i/p file with a special character ~~eof~~ marks the end of the source file.
- There are two pair: lexeme-beginning: marks to beginning of ~~current lexem~~ ~~half~~



Algo

if (forward at end of first half) then begin

 : reload second half

 : forward = forward + 1
end

else if (forward at end of second half) then
beginning in

 : reload first half

 : move forward to beginning of first half
end

else

forward = forward + 1

eg

abc = pqr * xyz

 ↓
 forward

: a : b : c : = :	p : : q : r : * : x : y : z : eof
-------------------	-----------------------------------

 ↑
 lexeme

beginning abc → identifier (= when forward reaches at = then lexeme.beginning token from file).

: a : b : c : = :	p : : q : r : * : x : y : z : eof
-------------------	-----------------------------------

Pqr → Identifier

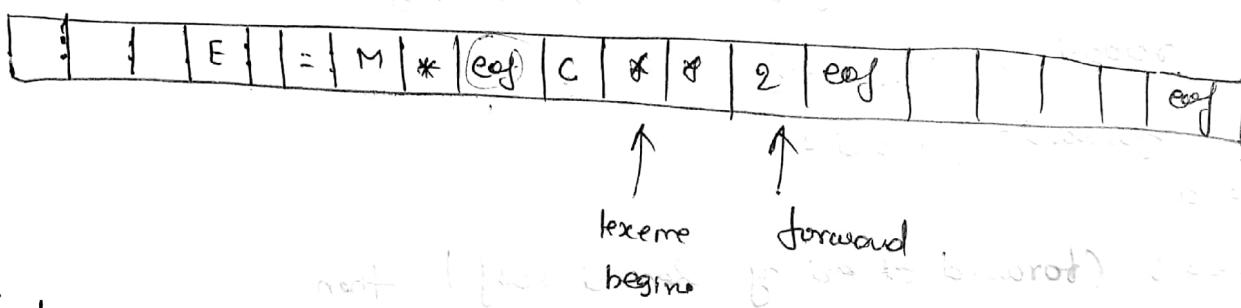
Sentinels:

- In the previous scheme each time when the forward pointer is moved, a check is done to ensure that one half of the buffer has not moved off. If it is done then other half must be reloaded.
- Therefore at the ends of the buffer require two tests for each advancement of the forward pointer.

Test 1: For end of buffer

Test 2: To determine what character is read.

- The usage of sentinel reduces two test to one by extending each buffer half to hold a sentinel character at the end.
- The sentinel is a special character that can't be part of the source program.



Algo

forward = forward + 1

if forward = eof then

begin

: reload sentinel half

: forward = forward + 1

end

else if forward at end of second half then begin

: reload first half

: move forward to beginning of first half

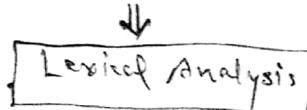
end

else terminate lexical analysis

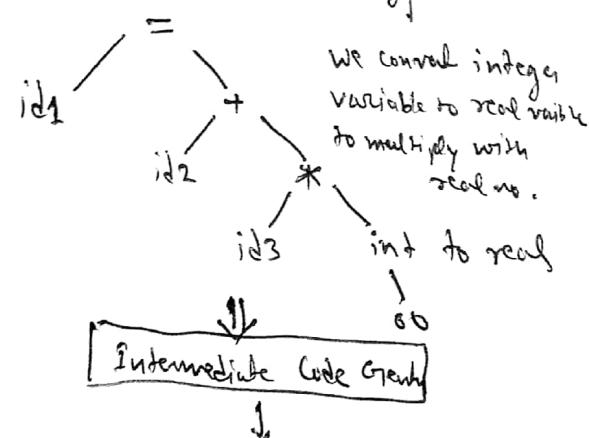
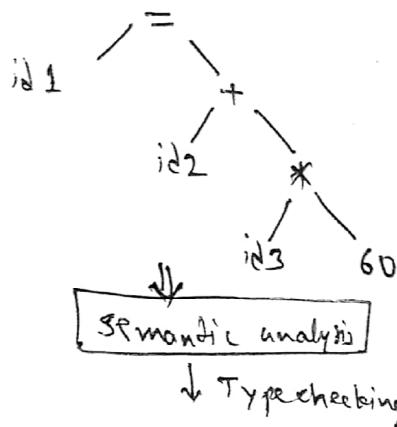
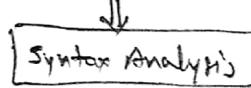
Adv.

- most of the time it performs only one test to see whether forward pointer points to an eof.
- only when it reaches the end of the buffer half or eof. it performs more tests.

$$\text{posn} = \text{initial} + \text{rate} * 60$$



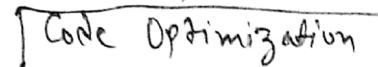
$$\text{id1} = \text{id2} + \text{id3} * 60$$



$$\text{temp1} = \text{int to real}(60)$$

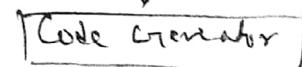
$$\text{temp2} = \text{id3} * \text{temp1}$$

$$\text{id1} = \text{id2} + \text{temp2}$$



$$\text{temp} = \text{id3} * \text{int to real}(60)$$

$$\text{id1} = \text{id2} + \text{temp}$$



MOV F id3, R1 Left Register

MUL F #60, R1

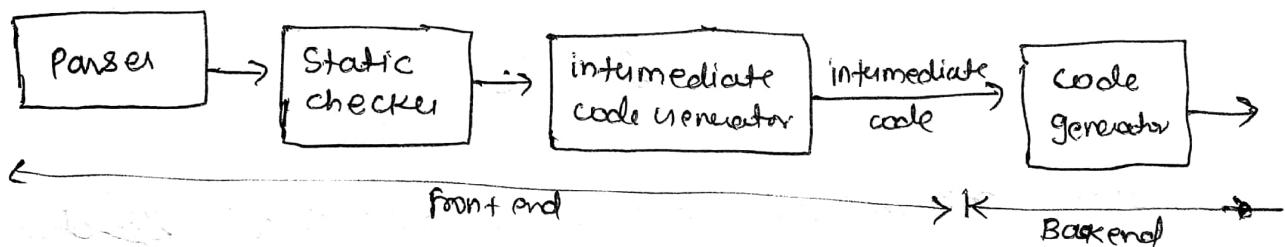
MOV F id2, R2

ADD F RL, R2

MOV F R2, id1

Intermediate code generation:

- in analysis-synthesis model of a compiler the front end translates a source program into an intermediate representation from which the back end generates target code.
- The benefits of using machine independent ~~code~~ intermediate code are:
 - Because of machine independent intermediate code, portability will be enhanced.
 - Retargeting is facilitated. (A compiler for a different machine can be created by attaching a back end for the new machine to an existing front end.)
 - It is easier to apply source code modification to improve the performance of source code by optimizing the intermediate code.



- uses register names, but has unlimited registers.
- uses control structures (like jump goto).
- uses OP-codes (eg push).

e.g. Three address code (if uses only 3 address)

$$x = y + z$$

Intermediate languages:

- Syntax tree and postfix notation are two kind of intermediate representation. and third is called three address code

Syntax tree

Production

$$E \rightarrow E_1 + T$$

$$E \rightarrow E_1 - T$$

$$E \rightarrow +$$

$$T \rightarrow (E)$$

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

$$T \rightarrow \text{num}$$

Semantic Rule.

$E\text{-node} = \text{new Node}('+' , E_1\text{-node}, T\text{-node})$

$E\text{-node} = \text{new Node}(' - ', E_1\text{-node}, T\text{-node})$

$E\text{-node} = T\text{-node}$

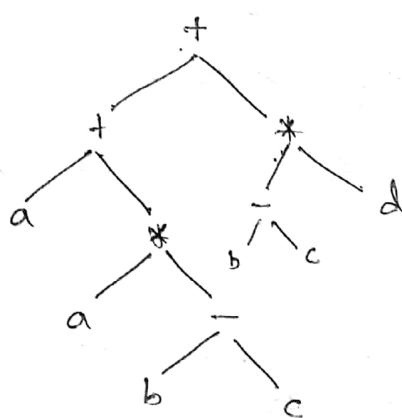
$T\text{-node} = E\text{-node}$

$T\text{-node} = \text{new Leaf}(\text{id}, \text{id.entry})$

$T\text{-node} = \text{new Leaf}(\text{num}, \text{num.val})$

$$a + a * (b - c) + (b - c) * d$$

representation of tree:

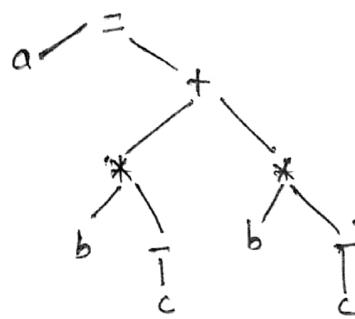
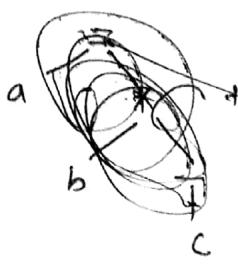


- Syntax tree
- Postfix expression
- 3 address code

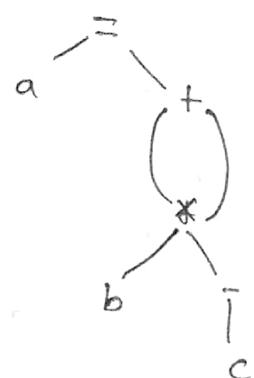
A syntax tree depicts the natural hierarchical structure of a source program.

A dag gives the same information but in a more compact way bcz common expressions are identified.

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



Syntax tree

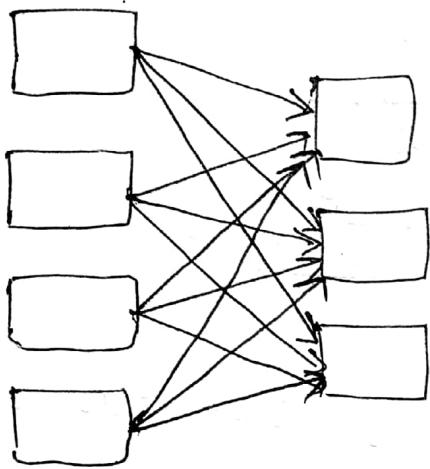


Dag (Directed Acyclic)

To representation of Syntax tree.

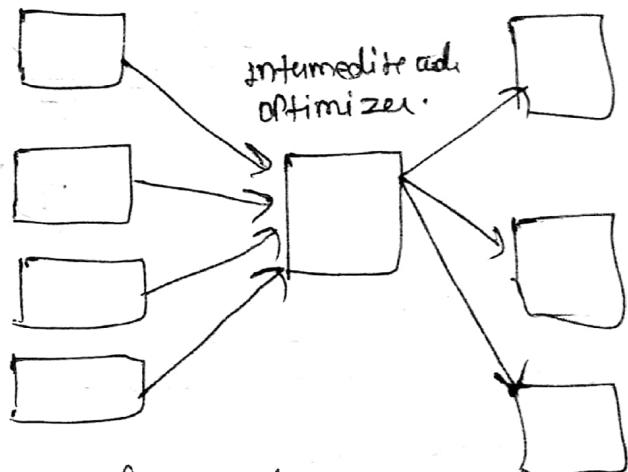
Why we need intermediate code.

4 Source
Img



3 target
machine

4 Comp
ing.



3 target
machines

4 - frontend

- + 4*3 optimizers
- + 4*3 code generators

4 frontend

1 optimizer

3 code generator

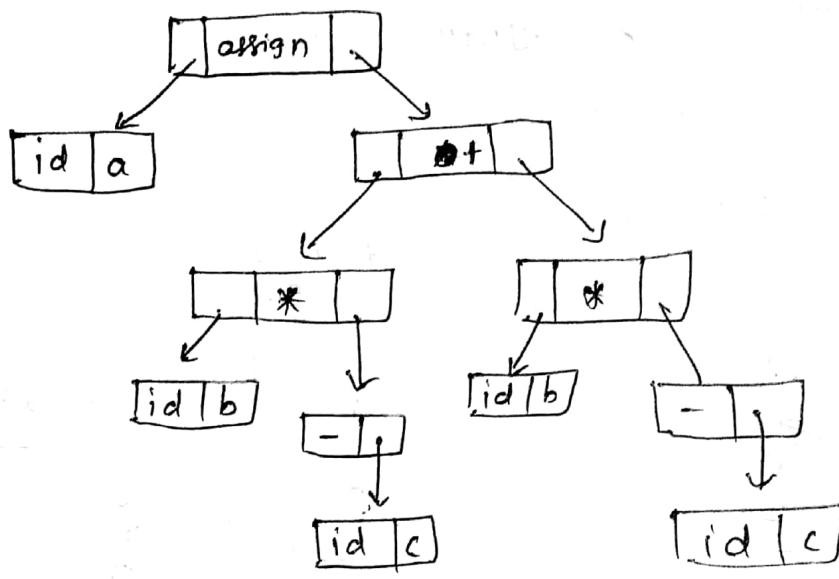
① nodes are allocated from an array of records and the index or position of the nodes serves as the pointer to the node.

All the nodes in the Syntax tree can be visited by following pointers starting from the root at pos. 10.

0	id	b
1	id	c
2	unminus	1
3	*	0 2
4	id	b
5	id	c
6	unminus(=)	5
7	*	4 6
8	+	3 7
9	id	a
10	assign(=)	9 8

left and right children of an intermediate node.

② each node is represented as a record with a field for its operator and additional fields for pointers to its children.



3 address code:

- Three address code is a sequence of statements of the general form.

$$x := y \text{ OP } z$$

- There is at most one operator (OP) on the right side of an instruction.
- At most three address are used to represent one statement. (2 for operands & one address for the result).

eg

$$x + y * z \quad (\text{Here 3 add. is used } (x, y, z))$$

$$t_1 = y * z \quad (\text{Assume this Alc to procedure of associativity}).$$

$$t_2 = x + t_1$$

$$\text{eg } a = \underline{\underline{x}} + y + z$$

$$t_1 = x + y$$

$$t_2 = t_1 + z$$

$$a = t_2$$

- Three address code is build from two concept address & instructions.

address can be: name, constant, c

Type of three address Stmt:

- Statement can have symbolic labels and there are Stmt for flow of control.

1) Assignment statements of the form $[x = y \text{ OP } z]$

OP → binary / arithmetic / logical opn.

2) Assignment instructions of the form $[x = OP y]$

OP → unary operation (+, -, <<, >> ...)

3) COPY Statement of the form $[x = y]$ when $x \leftarrow y$

4) The unconditional jump $.[goto L]$ three address Stmt with label L is the next to be executed.

5) Conditional jump such as $[IF X \text{ relOp } Y \text{ goto L}]$

6) Procedure call $[Param x \text{ and. call P return y}]$

7) Address and pointer assignment $x := &y \quad x := *y$.

8) $x = y[i] \quad x[i] = y$.

9) $x = *y \quad \text{and} \quad *y = y$.

param x
call P, n
 $y = \text{call}(P, n)$

Implementation of three address Stmt:

- A three address Statement is an abstract form of intmed - ate code.
- it can be implemented as quadruples, triples, indirect triples tree or dag.

eg. Three address code for

$$a + b * c - d / (b * c)$$

3 add. code.

$$t_1 = b * c$$

$$t_2 = a + t_1$$

$$t_3 = b * c$$

$$t_4 = d / t_3$$

$$t_5 = t_2 - t_4$$

Quadruples: it is a record structure with 4 fields which we call OP, arg1, arg2, ~~arg3~~ result.

The OP field contains an internal code for the operator. for the above eg.

OP	arg1	arg2	result
*	b	c	t_1
+	a	t_1	t_2
*	b	c	t_3
/	d	t_3	t_4
-	t_2	t_4	t_5

Exceptions:

- 1) With unary operator $x = -y$ or $x = y$ don't use arg2
- 2) Operator like `param` use neither arg2 nor result.
- 3) Conditional & unconditional Jump put target in result.

Triples: it is a record structure with three fields OP, arg1, arg2.

- one field
- arg1 & arg2 for the arguments of OP are either pointers to symbols or constant or pointers to structure.
 - It avoids the entering temporary names into the symbol table.
 - So it optimizes the space complexity.

OP	ang1	ang2
*	b	c
+	a	(0)
*	b	c
/	d	(2)
-	(1)	(3)

$b c_2 + t_3$ at index 0.
 $b c_3 + t_3$ present at index 2.

Indirect triples:

- it lists the pointer to triples rather than listing the triples themselves.

instruction

(0)
(1)
(2)
(3)
(4)
(5)
...

	OP	ang1	ang2
(0)	*	b	c
(1)	+	a	(0)
(2)	*	b	c
(3)	/	d	(2)
(4)	-	(1)	(3)

list of pointers to triples \Rightarrow Optimizing Compiler can record instruction list instead of affecting the triples themselves.

* SOT into three address code:

- it is constructed based on the grammar construct.
- The Synthesized attribute S.code represents the three address code for the assignment S.
- The non terminal E has 2 attributes:
 - E.place : (the name that will hold the value of E)
 - E.code : (the sequence of 3 address statement evaluating E)
- The function newtemp returns a sequence of distinct names t1, t2 ... in response to successive calls.
- hen is used to Compute the operations & exit with its Parameter.

eg

Production

Semantic Rule $S \rightarrow id : E ;$

$S \cdot code = E \cdot code \text{ || } \text{gen}(\text{top.get(id.lexeme)} \stackrel{\text{or}}{=} E \cdot \text{address})$
 $\text{gen}(id \cdot \text{place} := E \cdot \text{place})$

 $E \rightarrow E_1 + E_2$ $E \cdot \text{place} = \text{newTempC()}$ $E \cdot code = E_1 \cdot code \text{ || } E_2 \cdot code \text{ || } \text{gen}(E \cdot \text{place} = E_1 \cdot \text{place} + E_2 \cdot \text{place})$ $E \rightarrow E_1 * E_2$ $E \cdot \text{place} = \text{newTempC()}$ $E \cdot code = E_1 \cdot code \text{ || } E_2 \cdot code \text{ || } \text{gen}(E \cdot \text{place} = E_1 \cdot \text{place} * E_2 \cdot \text{place})$ $E \rightarrow -E_1$ $E \cdot \text{place} = \text{newTempC()}$ $E \cdot code = E_1 \cdot code \text{ || } \text{gen}(E \cdot \text{place} = \text{uminus } E_1 \cdot \text{place})$ $E \rightarrow (E_1)$

$E \cdot place = E_1 \cdot place \text{ (Here no need to generate the NewtempC
bez 'Parens' denote file (immediate add
which can be directly transfer).)}$

 $E \rightarrow id$ $E \cdot code \in \Sigma$

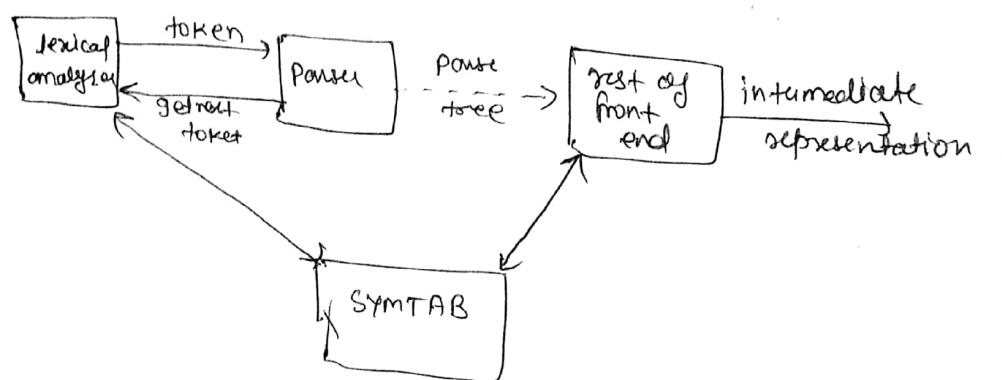
Issue in the design of a code generator:

① Input to the code generator:

- The input to the code generator consist of the intermediate representation of the source program produced by the front end together with information in the symbol table that is used to determine the run-time address of the data objects denoted by the name in the intermediate representation.

Role of the Parser:

- parser for any program take the i/p string w (which is the token) and produces as output either parse tree for w .
- if w is a valid sentence of grammar or error message indicate w is not valid sentence of given grammar.
- The goal of the parser is to determine the syntactic validity a source string is valid - tree build for the use of the next phase



Types of error can occur.

- Lexical (Such as misspelling an identifier, keyword or operator).
- Syntactic (Such as arithmetic expression with unbalance Parenthesis)
- Semantic (operator applied to incompatible operand)
- logical (infinitely recursive call)

* Error recovery Strategies:

- Panic mode recovery;
- Parse level
- Error Production
- Global Correction

→ Left recursion avoid using

$$A \rightarrow A\alpha / \beta \quad \leftarrow \text{left recr.}$$

Sol.

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' / \epsilon$$

* Left factoring:

- Procedure gives suitable grammar for converting Non-deterministic grammar to deterministic.
- 1) I/P is a grammar for predictive parsing.
 - 2) An equivalent left factored grammar.
 - 3) for each nonterminal A find the longest prefix common to two or more of its alternatives. If $\alpha \neq \epsilon$ then it is a nontrivial common prefix, replace α by the ~~the other~~ production $A \rightarrow \alpha B_1 | \alpha B_2 | \dots | \alpha B_n | \gamma$.

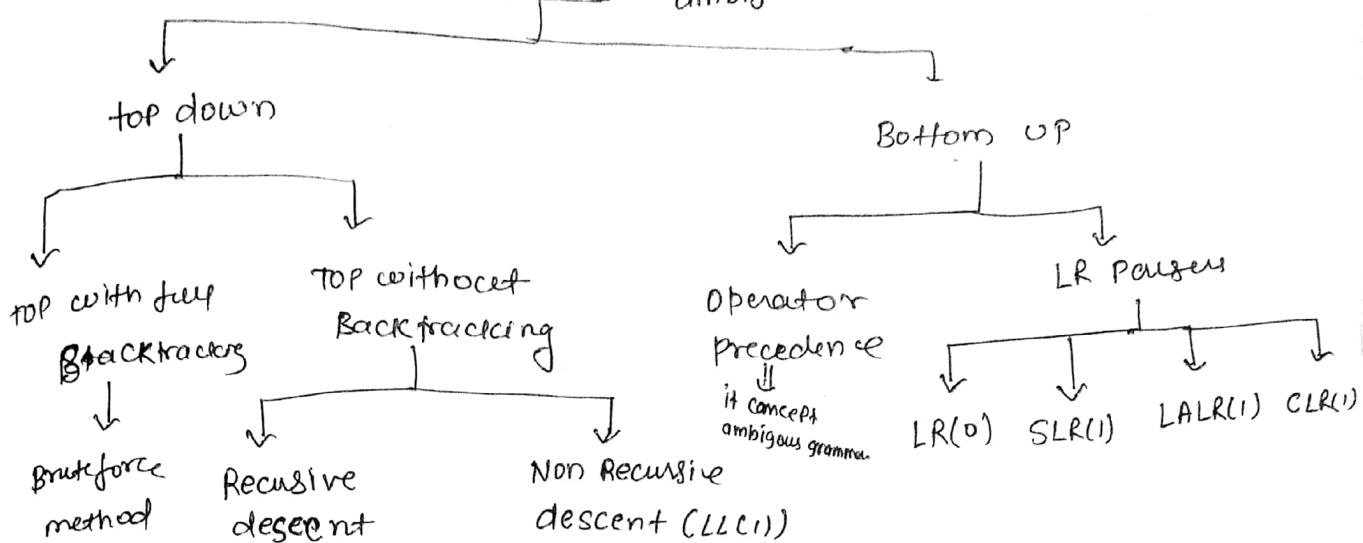
$$\begin{aligned} g: A &\rightarrow a\alpha \\ A &\rightarrow a\beta \\ A &\rightarrow a\gamma \end{aligned} \quad \left\{ \text{NDG.} \right.$$

$$\begin{aligned} A &\rightarrow a\beta \\ B &\rightarrow b\beta/b/c \end{aligned}$$

} left factored grammar

Parsing

Parsees \rightarrow For any parser grammar shouldn't be ambiguous except operator precedence.



top down Parser:

Derive from Start Symbol to last symbol.

Steps

1) Starts from root

2) Expansion

3) LL

\rightarrow left to right scan

\rightarrow left most derivation.

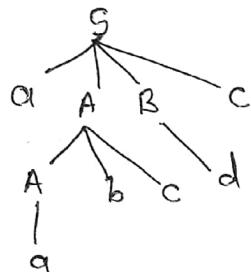
e.g.

$$S \rightarrow aABC$$

$$A \rightarrow Abc/b$$

$$B \rightarrow d$$

$$w = abcd$$



Bottom up Parser:

Starts from the word to starting production.

Steps:

1) ROOT - starts

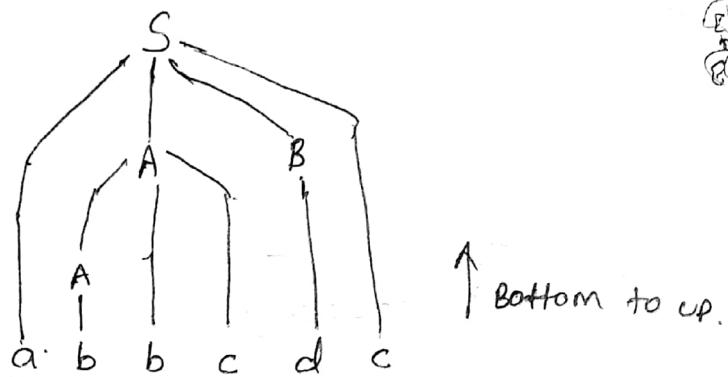
2) Reduction

3) LR

\rightarrow left right scan

\rightarrow right most derivation.

eg $S \rightarrow aABC$
 $A \rightarrow Abc/b$
 $B \rightarrow d$ w: abbcde



First and Follow :

- If we have a variable or string from that variable we try deriving all the strings then whatever is coming in the beginning, the first symbol is called as FIRST of that string.
- What is the terminal which could follow a variable in the process of derivation, is called as follow.

Algorithm FIRST(x) :

- 1) If x is terminal the FIRST(x) is $\{x\}$
- 2) If $x \rightarrow \epsilon$ is a production, then add ϵ to FIRST(x).
- 3) If x is a non-terminal and $x \rightarrow y_1 y_2 \dots y_k$ is production, then place a in FIRST(x) if for some i a is in FIRST(y_i).

Algorithm FOLLOW(x) :

- Place $\$$ in FOLLOW(S). When $S \rightarrow$ start symbol. bcz $\$$ is the right end marker of I/P buffer.
- If there is a production $A \rightarrow \alpha \beta B$ then everything in FIRST(B) except for ϵ is placed in FOLLOW(B).
- If there is a production $A \rightarrow \alpha \beta$ or production $A \rightarrow \alpha B \beta$ where FIRST(B) contains ϵ , then everything in FOLLOW(A) = FOLLOW(B).

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

	FIRST	FOLLOW
E	{id, ()}	{\$,)}
E'	{+, E'}	{\$,)}
T	{(, id)}	{+, \$,)}
T'	{*, F'}	{+, \$,)}
F	{(, id)}	{*, +, \$,)}

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

$$\text{FIRST}(E') = \{+, E'\} \quad \text{FIRST}(T') = \{*, F'\}$$

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

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

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

Q

$$S \rightarrow ABCD$$

$$\begin{array}{l}
 A \rightarrow a \\
 B \rightarrow b \\
 C \rightarrow c \\
 D \rightarrow d
 \end{array}$$

	FIRST()	FOLLOW
S	a	{\$}
A	a	{b}
B	b	{c}
C	c	{d}
D	d	{}\$}

Q

$$S \rightarrow ABCDE$$

$$\begin{array}{l}
 A \rightarrow a/\epsilon \\
 B \rightarrow b/\epsilon \\
 C \rightarrow c \\
 D \rightarrow d/\epsilon \\
 E \rightarrow e/\epsilon
 \end{array}$$

	FIRST	FOLLOW
S	{a, b, c}	{\$}
A	{a, \epsilon}	{b, c}
B	{b, \epsilon}	{c}
C	{c}	{d, e, \$}
D	{d, \epsilon}	{e, \$}
E	{e, \epsilon}	{\$}

	first	follow
$S \rightarrow Bb/Cd$	{a, b, c, d}	{\$, }
$B \rightarrow aB/E$	{a, e}	{b, }
$C \rightarrow cC/E$	{c, e}	{d, }

LL(1) Parsing table:

	first	follow
$E \rightarrow TE'$	{id, C}	{\$,)}
$E' \rightarrow +TE'/e$	{+, e}	{\$,)}
$T \rightarrow ET'$	{id, C}	{+,), \$}
$T' \rightarrow *FT'/e$	{*, e}	{+,), \$}
$F \rightarrow id / (E)$	{id, C}	{*, +,), \$}

$(E' \rightarrow e)$

Parsing table: it's aim to construct ~~top-down~~ parser.

	id	+	*	()	\$	← terminals
E	$E \rightarrow TE'$			$E \rightarrow TE'$			
E'		$E' \rightarrow TE'$			$E' \rightarrow e$	$E' \rightarrow e$	$(E \rightarrow id)$
T	$T \rightarrow FT'$			$T \rightarrow FT'$			
T'		$T' \rightarrow e$	$T' \rightarrow *FT'$		$T' \rightarrow e$	$T' \rightarrow \$$	
F	$F \rightarrow id$			$F \rightarrow ($			

variables.

- if word start with id then we can production is E so we must have to use production $E \rightarrow TE'$.
- whenever the first of any variable is E then we will place that production in follow of that variable(left hand side).
eg in the case of $E' \rightarrow e$ we will put this production in), \$
- This table is used to construct the parser tree.

g) $S \rightarrow (S)/\epsilon$

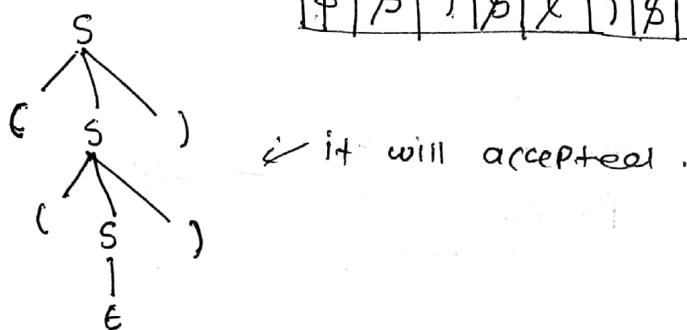
S	$\left \begin{array}{c} \text{follow} \\ \{\$,)\} \end{array} \right $	$\left \begin{array}{c} \text{first} \\ \{(, \epsilon\} \end{array} \right $
-----	--	---

LL(1) table

	()	\$
S	$S \rightarrow (S)$	$S \rightarrow \epsilon$	$S \rightarrow \epsilon$

g) Deriv the string $w = (()) \$$

pars tree



$\boxed{\$ | \$ |) | \$ | \epsilon |) | \$ | \epsilon | }$

← it will accepted.

g) $S \rightarrow AaAb / BbBa$

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

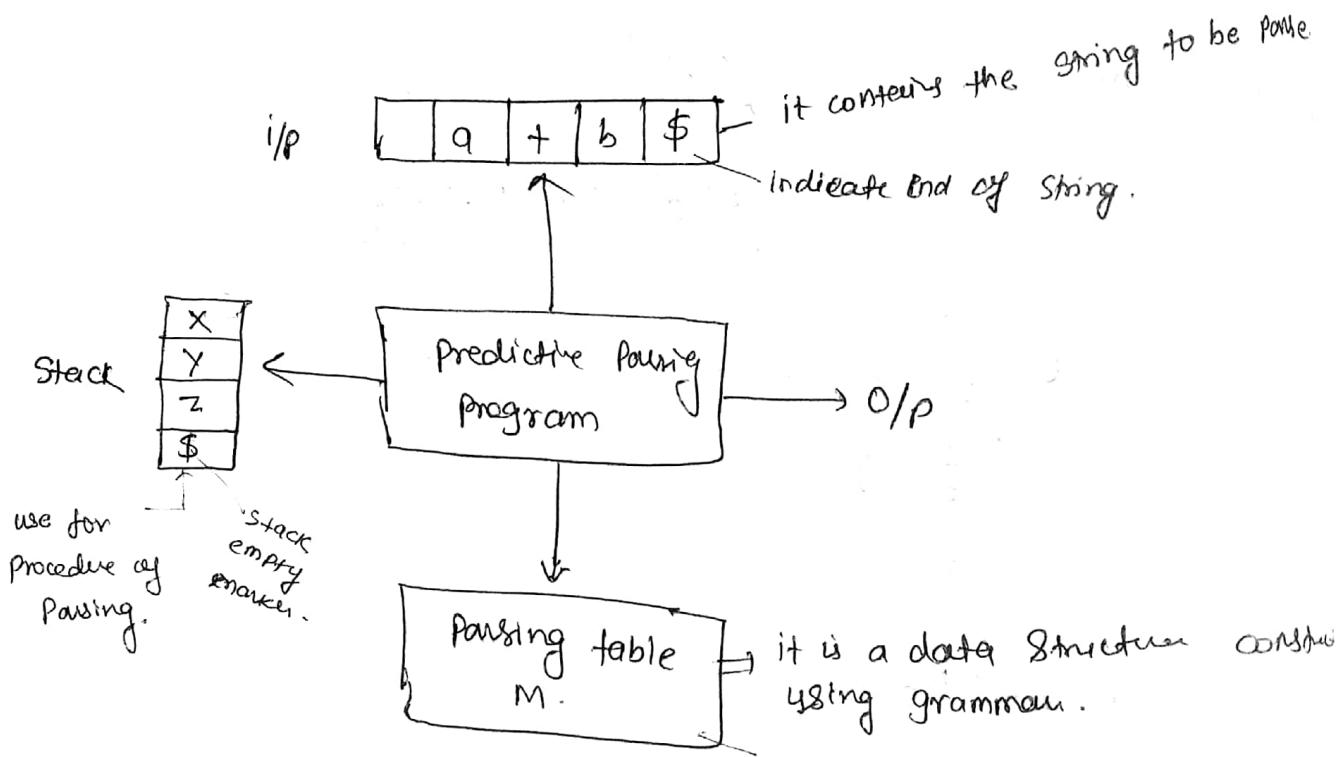
S	a	b	b	$\{\$\}$
A	$S \rightarrow AaAb$	$S \rightarrow BbBa$		
B	$A \rightarrow \epsilon$	$A \rightarrow \epsilon$		
	$B \rightarrow \epsilon$	$B \rightarrow \epsilon$		

• if two production entry will be in same cell so that grammar won't be LL(1).

Non recursive predictive parsing (LR(0))

look-ahead.

- it maintains a stack explicitly rather than implicitly via recursive calls.
- The key problem during predictive parsing is that of determining the production to be applied for a non terminal.



- Possibilities corresponding to 2 symbols.
- 1) if $X = a = \$$ the parser halts and announces successful completion.
 - 2) if $X = a \neq \$$ the parser pops a off the stack and advances the I/P pointer to the next I/P symbol.
 - 3) if X is non-terminal, the program consults entry $M[X, a]$ of the Parsing table M . This entry will be either an X -production or an error entry.

Algorithm:

set i/p to point to the first symbol of N\$
repeat

- 1) let X be the top stack symbol & a be the symbol pointed to by ip
- 2) if X is a terminal or \$ then
- 3) if $x = a$ then
- 4) pop X from the stack & advance i/p
- 5) else
- 6) error()
- 7) else /* X is a non-terminal */
 if $m[x, a] = x \rightarrow y_1 y_2 \dots y_k$ then begin
 pop X from the stack
 push $y_k y_{k-1} \dots y_1$ onto the stack with y_1 on the top.
 O/p the production $x \rightarrow y_1 y_2 \dots y_k$.
 end
 else
 error()

until $X = \$$

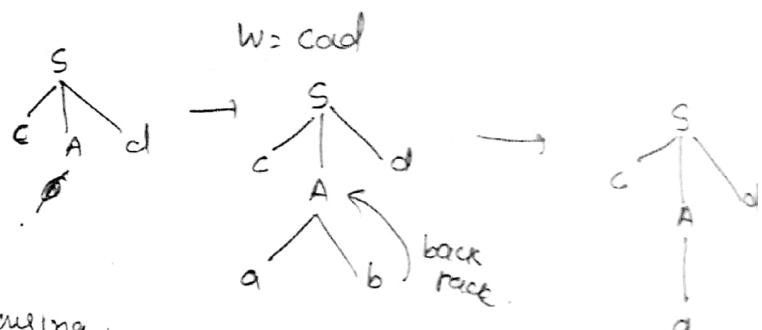
Recursive-Descent Parsing:

- It is one of the top down parsing techniques that uses a set of recursive procedures to scan its i/p.
- The Parsing method may involve backtracking, that is making repeated scan of the i/p.

Eg of backtracking.

$$S \rightarrow CAD$$

$$A \rightarrow ab/a$$



Eg of Recursive descent Parsing.

* A ~~left recursive grammar~~ can cause a recursive descent parser to go into an infinite loop.
Hence elimination of left recursion must be done before parsing.

Eg

$$\begin{aligned} E &\rightarrow iE \\ E' &\rightarrow +iE'/\epsilon \end{aligned}$$

Procedure EC()

begin
it ($i \leq l$) — looked at.

{ match(i);

} $E'()$;

end

$E'()$

{ it($odd i = 4$) }

{ match($+$) } match(i)

- $E'()$; get e

} else

} return

match(ch \neq $+$)

{ $ib(l = -t)$

$l = getch()$

} else

Prints ("error")

main()

{

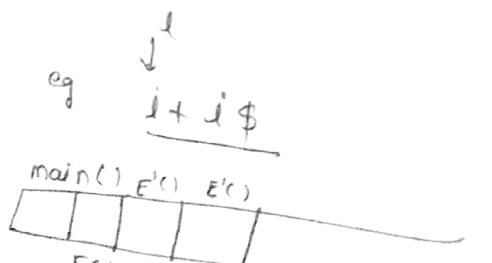
$E()$;

$ib(l = \$)$

} Prints ("Parsy Success"),

resursion
stack

provided
by OS



$$i \rightarrow E$$

Bottom up Parsing

- Constructing a Parse tree for an i/p string beginning at the leaves and going towards the root is called a bottom up parsing
- General type of bottom up parser is a Shift reduce parser