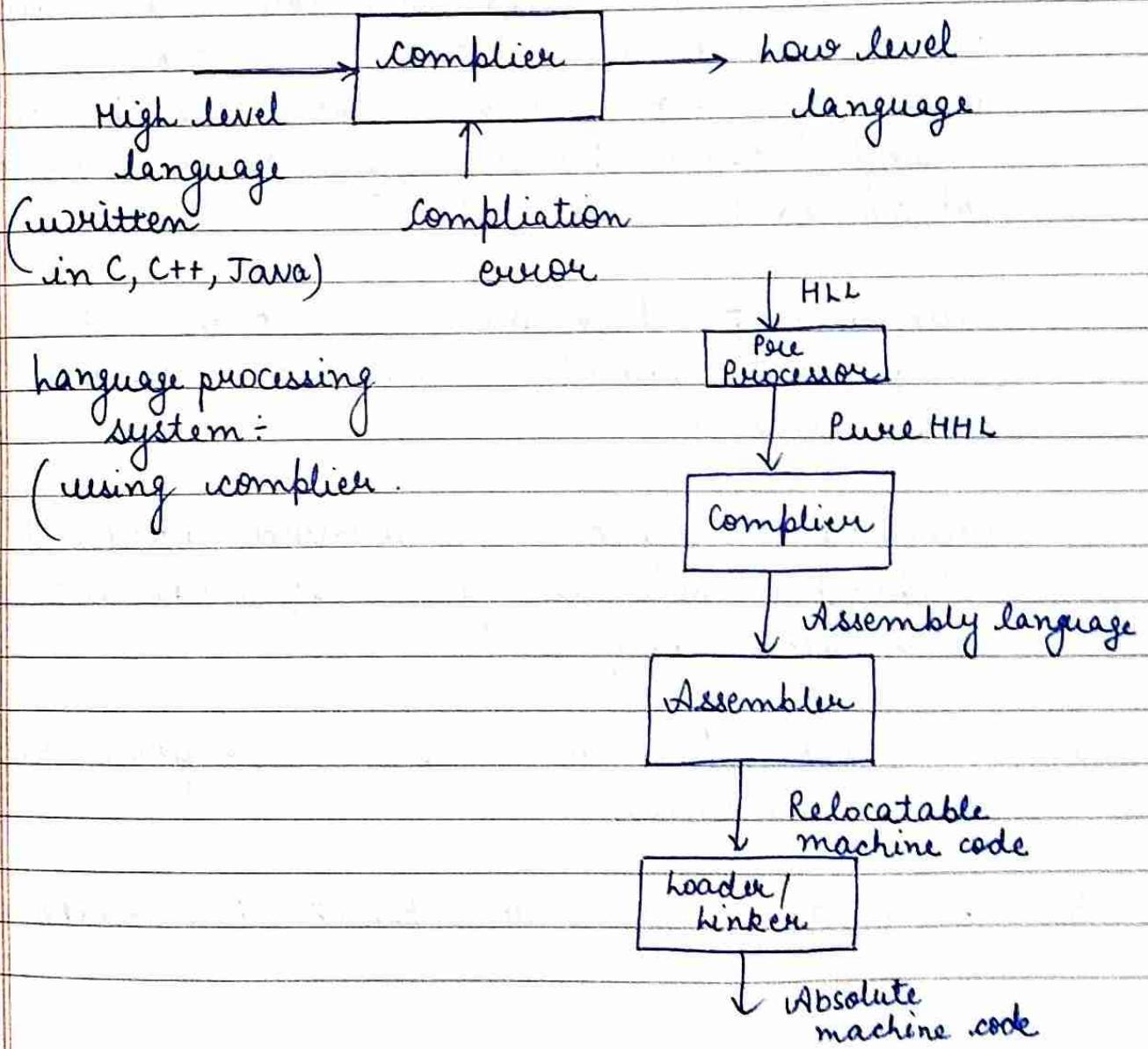


Compiler design:

A compiler is a special program that translates a program written in high-level programming language into machine code or bytecode i.e. understandable by computer.

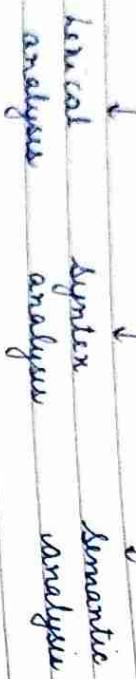
An important role of the compiler is to report the errors in the source program that is detected during translation process.

The output generated by compiler is called machine code / object code / target code / 0's / 1's.



Analysis of source program

The structure of the source program can be analyzed in 3 phases:



1. Lexical analysis:

It is the lowest level. It is also called linear analysis or scanning. In this, the source programs are viewed as sequence of characters called tokens. A token is a group of strings that make sense. It can be single character or sequence of characters. Tokens can be classified as:

1. Identifiers: These are names chosen to represent variables, functions, classes, date items etc.

2. Keywords: They are the reserved words of any programming language. Keywords can't be used as identifiers.

3. Operations: They are used to perform operations on operands.

4. Separators: These are punctuation marks used

e.g. → (,), ;, /, {, }, .
to group together a sequence of tokens that have a single meaning.

5. Literals: Some languages support literals which denote direct values

6. Comments: non-executable part of program which is used to increase readability of a program.

Eg.: A statement:

Position = Initial + rate * 60 we

grouped into following characters:

- a) Position → identifier
- b) Initial → identifier
- c) + → operator
- d) = → Assignment symbol
- e) rate → identifier
- f) * → multiplication sign (operator)
- g) 60 → number

2. Syntax analysis:

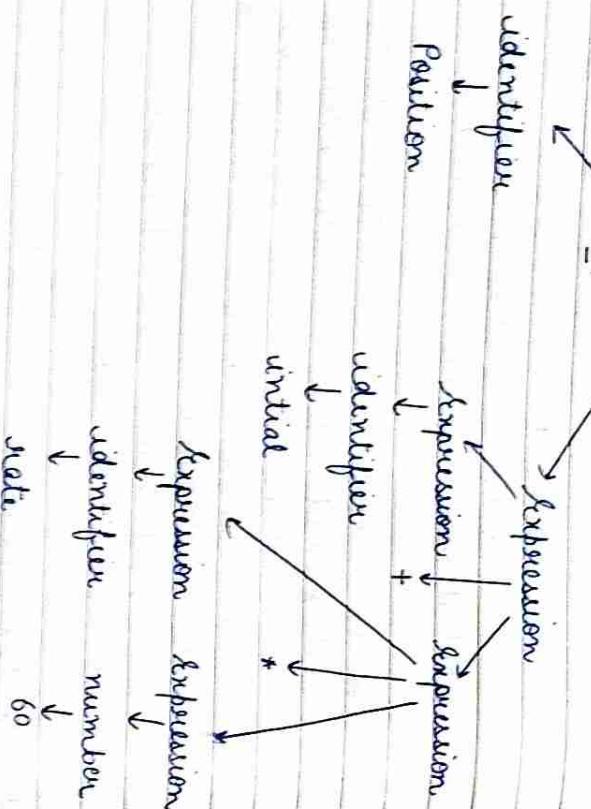
It describes the way in which program statements are constructed from tokens. This is defined using context free grammar & represented using parser tree. It involves grouping of source program into grammatical phases that are

would by the compiler to synthesize output.

→ A Syntax tree is a tree in which internal nodes are the operators & the external nodes are operands. This analysis shows an error if syntax is incorrect.

E.g.: Position = initial + rate * 60.

Parse tree → Assignment

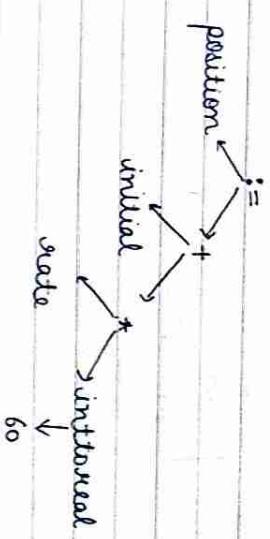


b. Abstract Syntax Tree: Compressed representation of parse tree.



a)

lexical analysis of parse tree.



In statement, position = initial + rate * 60, the identifiers position, initial & rate are declared as real numbers, but 60 is a integer. Therefore 60 is converted to type real. Thus, syntax tree becomes as shown below:

E.g. → Typechecking ⇒ checking for permitted operands.

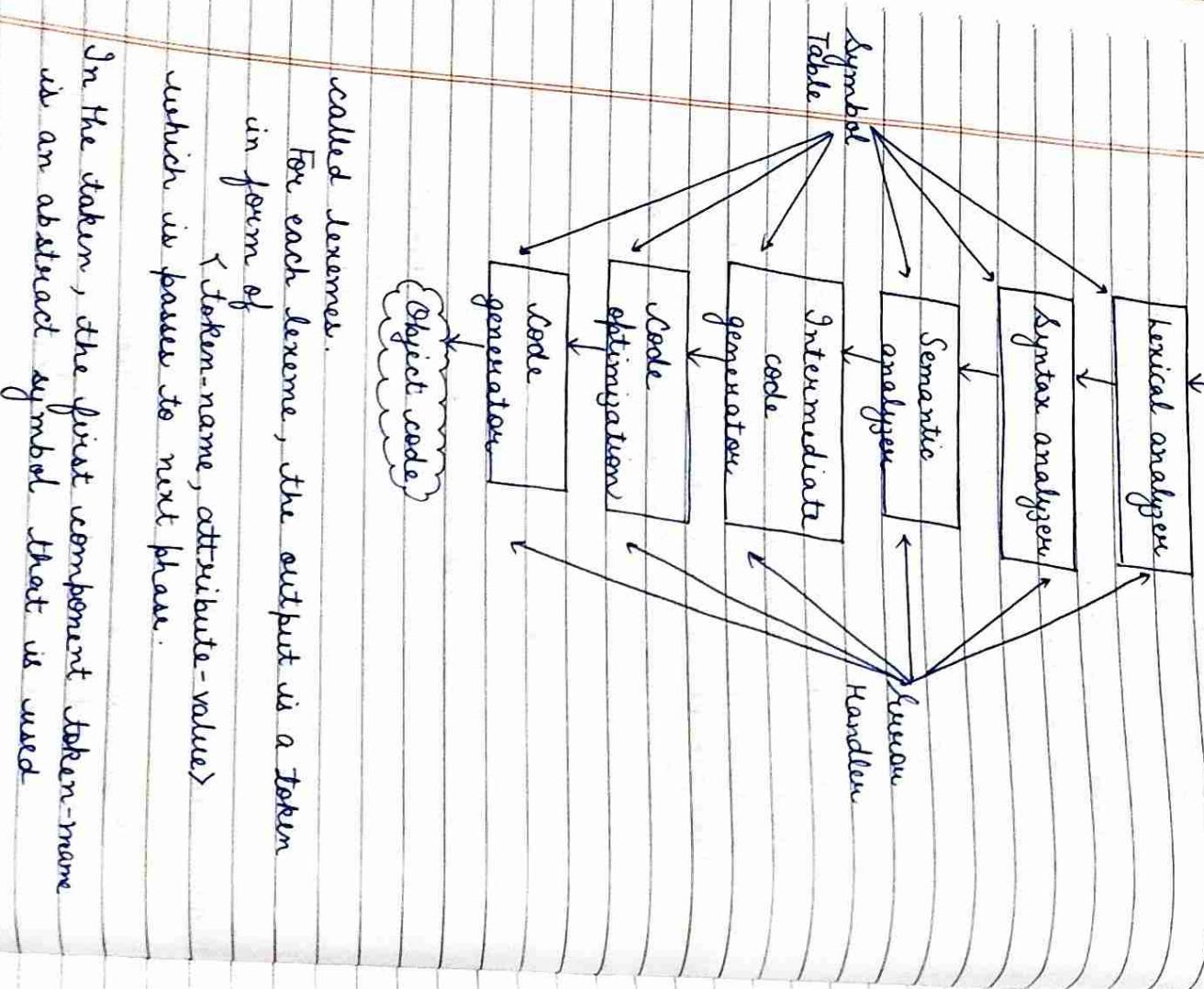
3. Semantic analysis:

This phase checks the source program for semantic errors & gathers type information for code generation phase.

It uses hierarchical structure determined by syntax level to identify operands & operators.

Source code

Done
Page



During syntax analysis & the second component attribute-value points to an entry in symbol table for this token. Information from symbol table is needed for semantic analysis & code generation.

E.g. - position = initial + rate * 60

The character in these statements would be grouped into following lexemes :

1. position is a lexeme that would be mapped into token `<id, 1>` where id is an abstract symbol for identifier & 1 points to symbol table entry for 'position'. The symbol table entry holds information about the identifier, such as name & type.

2. The assignment symbol = is a lexeme that is mapped into token `<=>`. It needs no attribute value, so second component is omitted.

called lexeme.

For each lexeme, the output is a token in form of `<token-name, attribute-value>`

which is passed to next phase.

In the token, the first component token-name is an abstract symbol that is used

- 6.
 - 5.
 - 4.
 - 3.
 - 2.
 - 1.
- 60 is the lexeme that is mapped into `<60>`.

Blanks separating the lexeme would be eliminated by lexical analyzer.

$\langle \text{id}, 1 \rangle \langle = \rangle \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle 60 \rangle$

Tokens: Token is a single character or sequence of characters treated as single entity.
→ Identifiers
→ Keywords
→ Operators
→ Separators
→ Literals/ constants

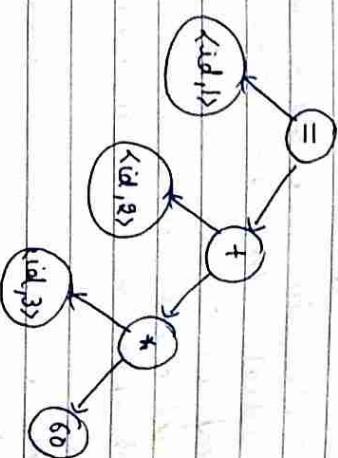
Pattern: A set of strings in input for which the same token is produced as output. This set of strings is described by rule called pattern associated with the token.

Lexeme: A lexeme is a sequence of characters in source program that is matched by pattern for the token.
e.g.: Token: Informal description: Sample lexeme →
if characters, if else characters e.g.,
Comparisons <, >, >=, <=, ==, !=
id letters followed by space, _, \$,
number any numeric constant 2.14, 3.22
literal anything but number "use"
- need by " "

b) Syntax analysis:

→ The second phase of compiler is parsing.
→ The parser uses the first component of the tokens produced by lexical analyzer to create a tree which is intermediate representation that tells about grammatical structure of token stream.
from Parse tree, syntax tree is generated.

Syntax tree:



→ The internal node is labelled with '+' has child.

($\text{id}, 3$) represents identifier node.

The node labelled with '*' makes it explicit that we must multiply $\langle \text{id}, 3 \rangle$ with 60.
The node labelled '+' indicates that we must add the result of multiplication to value of initial.

→ The root of tree, labelled with = indicate we must store the result of this addition for identifier position.

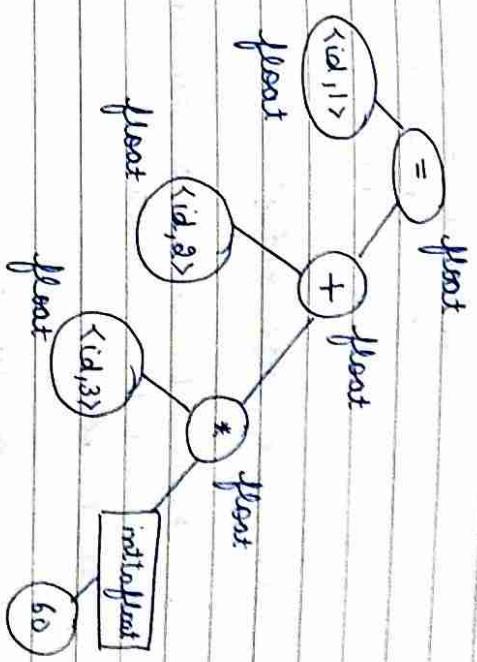
c.) Semantic analysis:

The semantic analysis uses the syntax tree & the information in symbol table to check the source program for semantic consistency with the language definition & it also gathers type of information & saves it either in syntax tree or the symbol table, for subsequent use during intermediate code generation.

→ Most important part of semantic analysis is type checking where compiler checks the each operator has matching operands.

e.g. → In many programming languages requires an array index to be integer, the compiler will report an error if floating point number is used to index an array.

d.) Intermediate code generation:



Ex: Suppose, the division program, initial state is given by itself forms an integer. The semantic analyzer discovers that the operator / is applied to floating point & integer.

So, the output of semantic analyzer adds an extra node for operator intofloat, which explicitly converts integer to floating point number.

Some sort of type conversion is done by compiler.

e.g. → If we add, 1 floating point & 1 integer, the compiler may automatically convert integer to float.

In the process of translating a source program, a compiler may construct a one or more intermediate representations, which can have a variety of forms. Syntax trees are form of intermediate representations which are commonly used during syntax & semantics analysis.

After that, the compiler generate an explicit low-level or machine like representation which is like a program for abstract machine. The generated code is machine independent.

I.g. → In our example, we use 3-address code which consists of sequence of instructions with 3 operands per

instruction.

I.g. → $t1 = \text{int} \text{to} \text{float}(60)$
 $t2 = \text{id}3 * t1$
 $t3 = \text{id}2 + t2$
 $\text{id}1 = t3$.

Properties of intermediate representation:

1. It should be simple & easy to produce.
2. It should be easy to translated into target code.

c) Code optimization:

The machine-independent code optimization phase attempts to improve the intermediate code so that the better target code will be generated.

It is done to make the execution faster, short code or the target code consumes less power. Thus, it increases the efficiency.

I.g.:
LDF R2, id3
MULF R2, #60.0.
LDF R1, id2

Code generation:

It is the last phase who takes intermediate representation of source program as input & maps it into the target language. The target code can be



Machine Assembly
code code.

→ If target language is machine code, registers or memory locations are used for each of variables used by source program. As, it is converted into sequences of machine instructions to perform task.

If the target language is assembly language, it generates assembly code as a target output.

I.g.: Code generated:

$t1 = \text{id}3 * 60.0$

The 4 lines of code is converted in 3 lines of code.

The optimized code of above example is,

ADD F R1, R2
STF w1, R1.

The F in each instruction refers to floating point number.

Explanation:

1. Load the contents of address w_1 into register R_2 , then multiply it with floating point constant number 60.0.
indicates that it is intermediate constant.

2. The 3rd instruction moves w_1 's content into register R_1 . & 4th adds the contents of R_1 & R_2 . The resulted value is stored in register R_1 .

3. Finally, the value in register R_1 is stored in address of w_1 .
So, the code correctly implements statement, position = initial + scale * 60.

Symbol Table

Symbol table is the data structure containing a record for each variable name, with fields for attributes of the name.

They attributes may provide information about storage allocated for a name, its type,

its scope (where in program its value may be used), & in case of functions such things as number & type of arguments, the method of passing arguments (call by value, call by reference) & type returned.

The data structure stores records in such a way it is easy to find the record very quickly & retrieve the data for the record. The symbol table is used by every phase of compilation to get the information.

Error detection & reporting:

→ Each phase can encounter error. However, after detecting an error, a phase must somehow deal with that error, so that compilation can proceed, allowing further errors in source program to be detected.

Grouping of phases:

The process of compilation is split up into following phases:

- Analysis phase: It mainly performs 4 actions:

- lexical analysis,

- b.) Syntax analysis.
 c.) Semantic analysis.
 d.) Intermediate code generation.

The analysis part breaks up the source program into constituent pieces & impose a grammatical structure on them. Then this structure is used to generate intermediate representation.

If analysis part detects that source program is either syntactically ill formed or semantically unsound, then it must provide informative message, so that user can take corrective action.

It also collects info about source program & stores it in data structure called symbol table, which is passed along with intermediate representation to the synthesis part.

Symbol table + Intermediate code → Synthesis phase.

Synthesis phase:

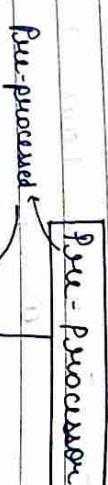
It mainly performs 2 actions:-

- a.) Code optimization.
 b.) Code generation.

The synthesis part constructs the desired

target program from intermediate representation & info. in symbol table. It also performs optimization so that execution takes less time, lesser space & thus provides efficiency.

Cousins of compiler:

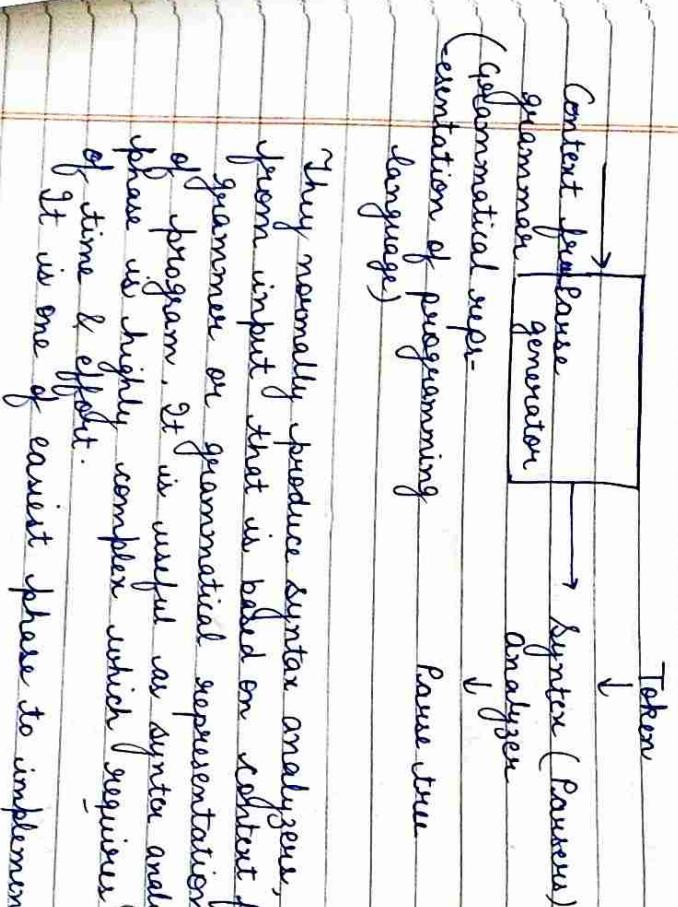


Semantic construction tools:

Complex writers use software development tools & more specialized tools for implementing various phases of compiler.

Some commonly used construction tools are:

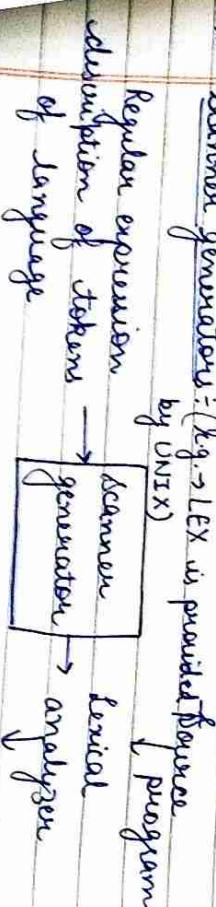
1. Parser generator:



They normally produce syntax analyzers, from input that is based on context free grammar or grammatical representations of program. It is useful as syntax analysis phase is highly complex which requires lot of time & effort.

It is one of easiest phase to implement.

2. Scanner generators: (e.g. LEX is provided by UNIX)



3. Syntax directed Translation engine:

Input: Parse tree
Output: Intermediate code.

It generates intermediate code with three address format from input consists of a parse tree.

These engines have routines (functions) to traverse (visit) each node in parse tree & produces intermediate code for that node. The basic concept is that one or more translations are related or associated with each node of parse tree & each translation is represented because of translation at its neighbor nodes in the tree.

4. Automatic code generators:

Input: Intermediate code.
Output: Machine language.
It generates the machine language for a target machine. Each operation of intermediate

It generates lexical analyzer from the input consists of specifications of regular expression which is description of token of a language. It generates finite automata (DFA & N DFA) to recognize regular expression.

Code is translated using collection of rules & then is taken as input by the code generator.

The rules must include sufficient detail that we can handle the different possible access methods for data.

In this process, template matching is used.

5. Data flow analysis engine:

It is used in code optimization. It is a part of the code optimization that gathers the information i.e. the values that flow from one part of a program to another.

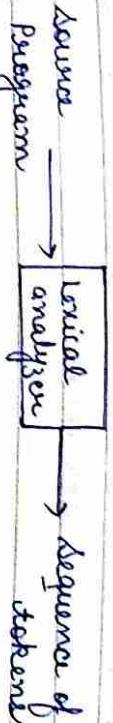
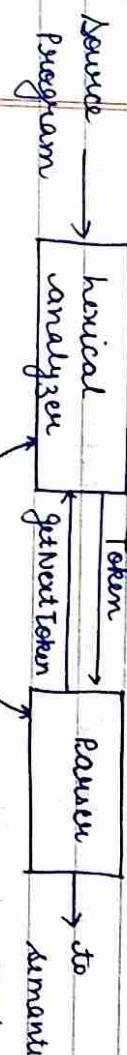
Native compiler & cross compiler

Native

1. Native compiler are complete that generate code for same platform on which it runs.

Cross

1. Cross compiler is a compiler that generate executable code for a platform other than one on which compiler is running.



When a lexeme is discovered, it needs the lexeme into the symbol table. The entries in the symbol table are needed by the other phases of compilation.

1. The main task of lexical analyzer is to read the input characters of source program, group them into lexemes to produce an output a sequence of tokens.

Role of lexical analyzer:

1. a) The main task of lexical analyzer is to read the input characters of source program, group them into lexemes to produce an output a sequence of tokens.

e.g. → Turbo C or GCC
is native compiler.

CC is that compiler running on windows produce object code for MAC or etc.

Sig's Keil.

1. It is used in code optimization. It is a part of the code optimization that gathers the information i.e. the values that flow from one part of a program to another.
2. It is dependent on machine.
3. It is used to build programs for same system, if OS is installed.
4. It generates executable file like .exe.

The interaction is implemented by sharing the power call the lexical analyzer. getnextToken command causes the lexical analyzer to read from its input until it can identify

Lexical analysis: A character sequence that can't be scanned into any valid token is a lexical error.

the next lexeme & produce a token, which it returns to parser.

It removes all the comments & whitespace (blank, newline, tab or other characters that are used to separate tokens)

It also generates error messages like exceeding length, unmatched string, illegal character. For instance, lexical analyzer may keep the track of number of newline characters seen, so it can associate a line number with each error message.

Attributes for token:

A token is associated with different pieces of information i.e. each token has different properties or attributes.

(e.g. → Information about identifier is its lexeme, type, location where it is first found etc. is kept in symbol table.)

Input buffering:

The lexical analyzer scans at input at a time from left to right. It makes use of two pointers: begin ptr (bp) & forward pointer (fp). To keep track of the input scanned.

Initially, both the pointers point to the first

character of the input string.

int a = 1;

i	n	t	=	1	;
↑	↑	↑			

ptr
bp

The forward ptr moves ahead to search for end of lexeme. As soon as the blank space is encountered it indicates end of lexeme.

In above example as soon as ptr (fp) encounters a blank space the lexeme "int" is identified.

The fp ptr moves ahead, when it encounters blank space, then both bp & fp are set at next token. The input character is then read from secondary storage, but reading this way from secondary storage is costly.

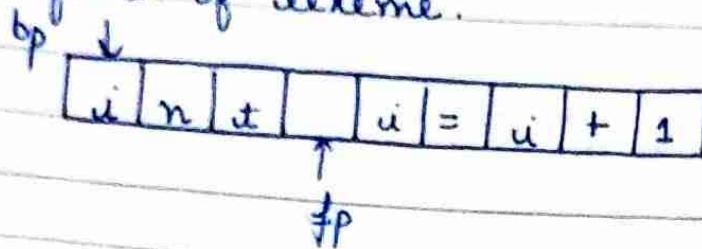
Hence, buffering technique is used. A block of data is first read into a buffer, & then by lexical analyzer. There are 2 methods used in this context.



1. One buffer scheme:

In this scheme, only one buffer is used to store the input string but in this, the problem is that if lexeme is very long, then it crosses the buffer boundary,

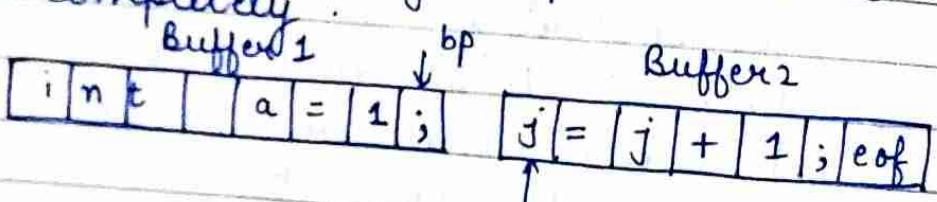
To scan rest of the lexeme the buffer has to be refilled, that makes overwriting of the first of lexeme.



2. Two buffer scheme:

To overcome the problem of one buffer scheme, in this two buffers are used to store the input string. The first & second buffer are scanned alternately. When end of current buffer is reached then the other buffer is filled.

- The only problem with this method if length of lexeme is longer than length of buffer, then scanning input can't be scanned completely.



- Initially both are pointing to ^{fp} first characters of first buffer. fp moves forward in search of lexeme. As soon as, blank character is recognized, string b/w fp & bp is identified, ^{as} token
- The "end of" character is placed at end of each buffer to identify the first & second buffer.
- So, both the buffers can be filled alternatively until the end of input program & stream of token is identified.

Specification of tokens:

To specify a token, we use Regular expressions which generate Regular languages.

Alphabet: An alphabet is a finite non-empty

set of symbols. Symbols can be letters, numbers or any special character.

E.g. $\Sigma = \{0, 1\}$ \rightarrow binary alphabet

$\Sigma = \{a, b, c, \dots, z\}$ \rightarrow lower case letters alphabet.

Strings: Strings are finite set of symbols generated from alphabet.

E.g. $\Sigma = \{0, 1\}$

String $\rightarrow 0, 01, 100, 11, 1\dots$ etc.

The length of string is total number of characters in the string.

Language: Set of strings generated from the alphabet is called language. It is collection of string.

E.g. $\Sigma = \{a, b\}$

L_1 = Set of all strings of atleast length two.

$L_1 = \{aa, bb, aaa, aaaabb, \dots\}$

So, L_1 is infinite language.

So, languages can be finite or infinite.

Regular Expression:

- 1.) Regular expression represents patterns of string of characters.
- 2.) R.E also have a special character called meta-characters. E.g. $\rightarrow *, +$ etc.

Rules for forming Regular expression:

1. Let 'a' be a character in string S. The regular expression is 'a'. It matches with character 'a' by writing $L(a) = \{a\}$
2. Any terminal symbol (element in Σ), λ & \emptyset are regular expressions.
3. The Union of two regular expression is regular expression.
4. The iteration or closure of regular expression is R^* , which is also a regular expression.
 $R^* = \{\lambda, a, aa, aaa, \dots\}$
 $R^+ = \{a, aa, aaa, \dots\}$.
5. If the alphabet of language is null i.e. $\Sigma = \emptyset$ then $L(\emptyset) = \{\}$

Operations on regular expression:

1. Choice among alternates : (Union)

Let r and s be two regular expressions. Then, union of these two regular expression represented by r/s is also a regular expression.

$$L(r/s) = L(r) \cup L(s)$$

Let us suppose, $L(r) = \{a, b\}$

$$L(s) = \{c, d\}$$

$$L(r/s) = \{a, b, c, d\}$$

2. Concatenation : Concatenation of 2 regular

expressions is also a r.e. represented by

e.g.

Two regular expressions r & s ,
concatenation is written as rs

$$L(rs) = L(r)L(s) \quad [\text{character in } L(r) \\ \text{followed by char. in } L(s).]$$

$$L(r) = \{a\}$$

$$L(s) = \{b\}$$

$$L(rs) = \{ab\}$$

c.) Repetition: Repetitive operation of a regular expression is called Kleene.

Closure: It is represented by r^* .

$$\text{e.g.} \rightarrow L(r) = \{a\}, \quad L(s) = \{b\}$$

1.) $r^* \rightarrow$ It represents the language containing the zero or more occurrences of symbol from $L(r)$.

$$L(r^*) = \{\epsilon, a, aa, aaa, aaaa, \dots\}$$

2.) $(rs)^*$

$$L((rs)^*) = \{\epsilon, ab, abab, abababab, \dots\}$$

3.) $(r|ss^*)^*$

$$L(r|ss^*)^* = \{\epsilon, a, bb, aaa, aaabbabb, \dots\}$$

LEX :-

1. lex is a computer program that generates lexical analyzer. It is commonly used with YACC parse generator.
2. It transforms an input into sequence of tokens.
3. It reads the input stream & produces source code as output through implementing the lexical analyzer in C program.

Structure of a Lex File :-

A lex program is separated into three sections by % % delimiters. The format of LEX source is as follows :-

{ definitions }

% %

{ rules }

% %

{ user subroutines }

- The definitions section defines macros & import header files written in C. It is also possible to write any C code here including declarations of constant, variables & regular definitions.
- Rules define the statement of the form $p_1 \{ \text{action } 1 \} p_2 \{ \text{action } 2 \} \dots p_n \{ \text{action } n \}$, where p_i describes the regular expressions & action_i describes actions what action the lexical analyzer should take when pattern p_i matches a lexeme.
- User subroutines are auxiliary procedure needed by the actions. The subroutine can be loaded with lexical analyzer & compiled separately. It contains C statements & functions which are called by rules in rule section.

Example of LEX file:

```

/* Definition section */
%{
#include <stdio.h>
%}
%o

/* Rules section */
/* [0-9]+ matches string of one or more digits */

```

```
[0-9]+ {  
    printf ("Saw an integer: %s\n",  
           yytext);  
    /* Ignores all other characters */  
    %%  
/* C code section */  
int main (void)  
{  
    /* Call the lexer, then quit. */  
    yylex();  
    return 0;  
}
```

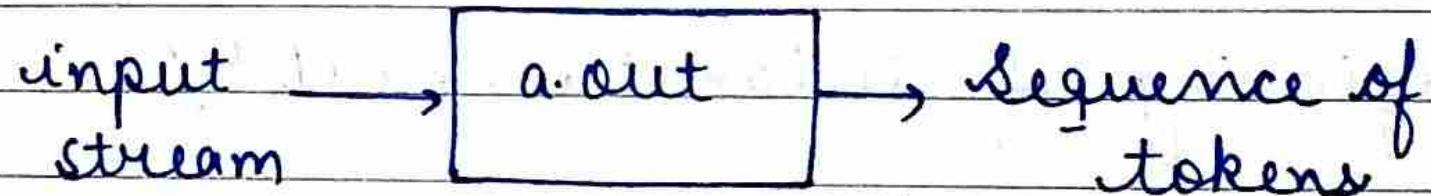
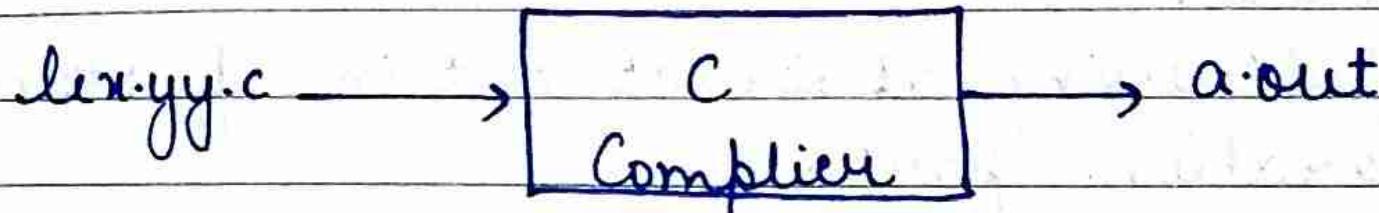
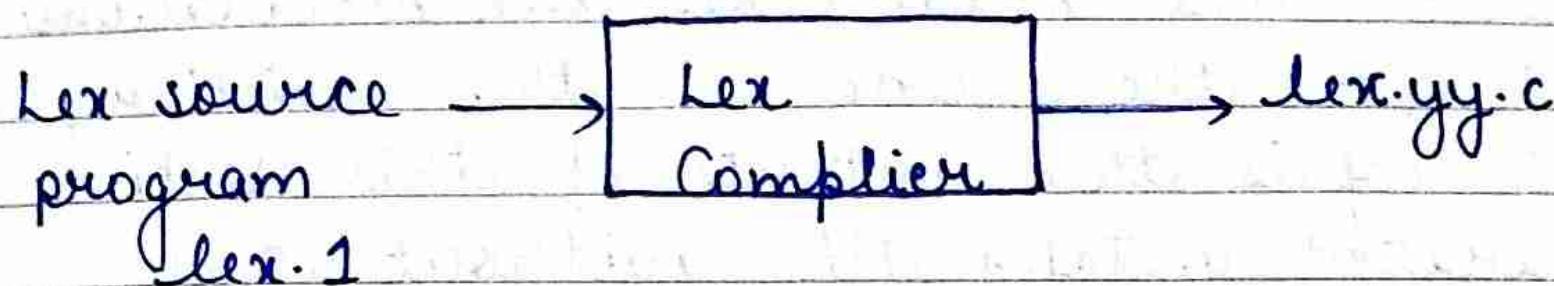
Given input: abc 123z.1g & * 2gh6

Output: Saw an integer: 123
Saw an integer: 2
Saw an integer: 6

Function of LEX:

1. Firstly, lexical analyzer creates a program `lex.1` in lex language. The lex compiler runs the `lex.1` program & produce a C program `lex.yy.c`.
2. Finally C compiler runs the `lex.yy.c` program.

- & produces a object program a.out.
3. a.out is lexical analyzer that transforms an input into sequence of token.

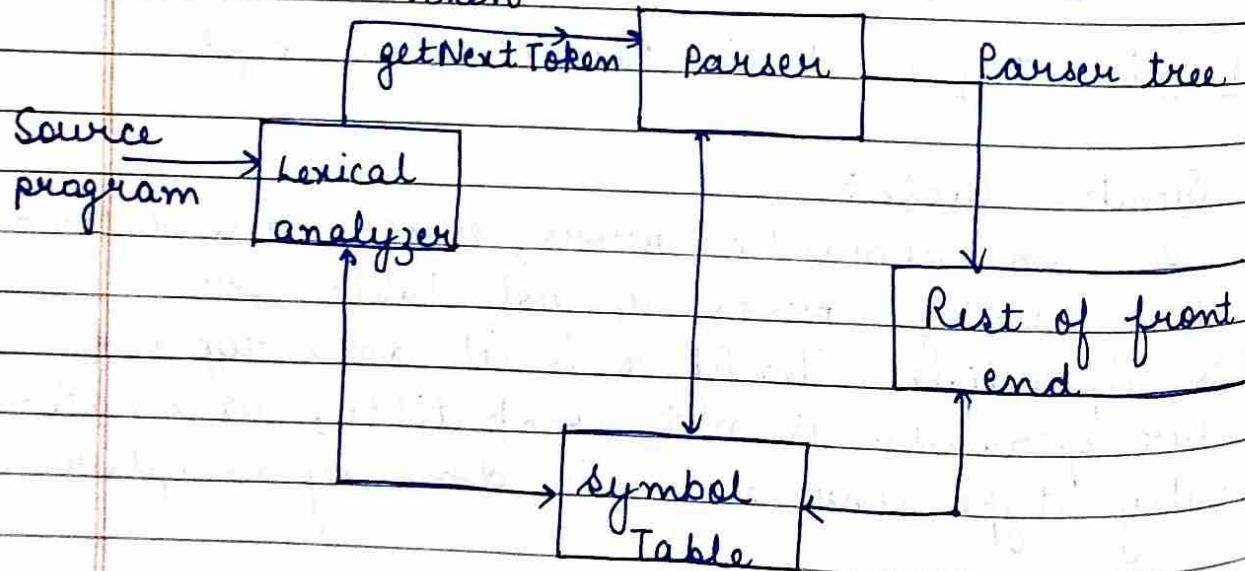


Role of parser:

In syntax analysis phase, a compiler verifies whether or not the tokens generated by the lexical analyzer are grouped according to the syntactic rules of the language. This is done by parser.

The parser obtain string of tokens from lexical analyzer & verifies that the string can be the grammar ~~from~~ for ~~the lexical analyzer~~ source language. It ~~decides~~ detects & reports any syntax error & produce a parse tree from which intermediate code can be generated.

Token



Context-free grammar:

The syntax of any programming language is described by context-free grammar. It is represented by tuples (V, S, Σ, P) where V is set of non-terminals, Σ is set of terminals, S is start symbol &

P is set of production in the form,
 $V \rightarrow (V + \Sigma)^*$

Parse tree:

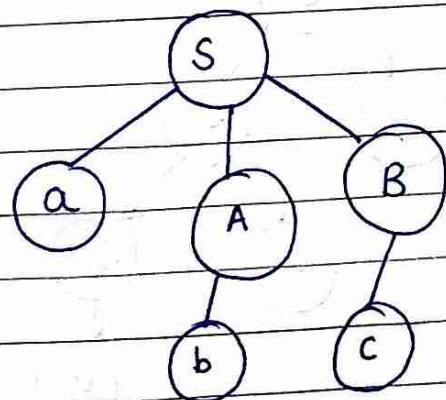
Parse tree or derivation tree is an ordered, rooted tree that represents the syntactic structure of a string according to some context-free grammar.

1. The starting symbol of a grammar must be used as root of the parse tree
2. Leaves or external nodes represent terminals.
3. Each interior symbol represents an production in the grammar.

e.g.: grammar: $S \rightarrow aAB$

$$A \rightarrow b$$

$$B \rightarrow c$$



A parse tree for a CFG $G = (V, \Sigma, P, S)$ is a tree satisfying these conditions:

1. Root has label S , S is start symbol.
2. Each vertex of parse tree has a label which can be variable (V), terminal (Σ) or ϵ (null).

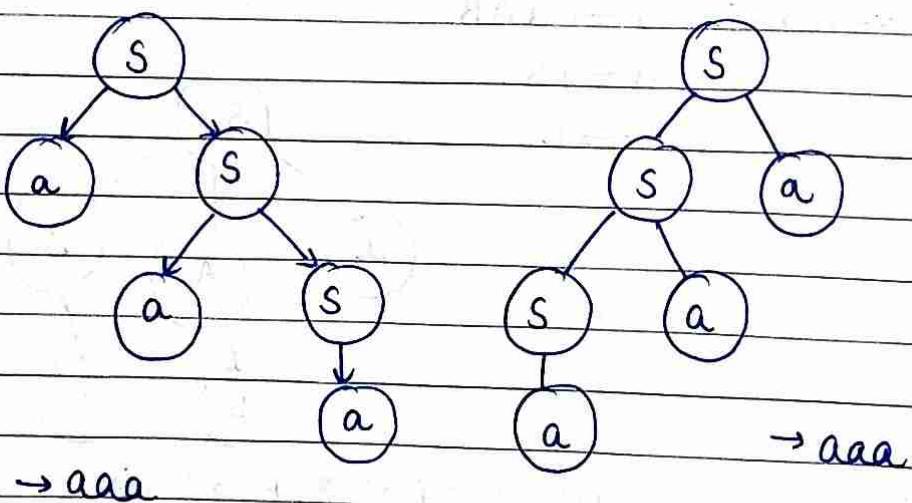
3. If $A \rightarrow C_1 C_2 \dots C_n$ is a production, then $C_1, C_2 \dots C_n$ are children of node labelled with A .
4. Leaf nodes are terminal (Σ) & interior nodes are (V) variable.

Yield: Yield of parse tree is concatenation of labels of the leaves in the left to right ordering.

Ambiguity: A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

E.g.: $S \rightarrow aS \mid Sa \mid a$

For string $\rightarrow aaa$



Parsing: The process of deriving the string from the given grammar is known as parsing (derivation).

Depending on how derivation is done, we have 2 kind of parsers:

1. Top-Down Parsing:

Top-down parsing attempts to build the parse-tree from root-to-leaf. The top down parser starts from the start symbol & proceed to the string. It follows left most derivation. In left most derivation, the left most non-terminal in each sentential is always chosen.

Top-down parsers

Recursive
Descent parsing

Non-recursive
descent parsing (LL)

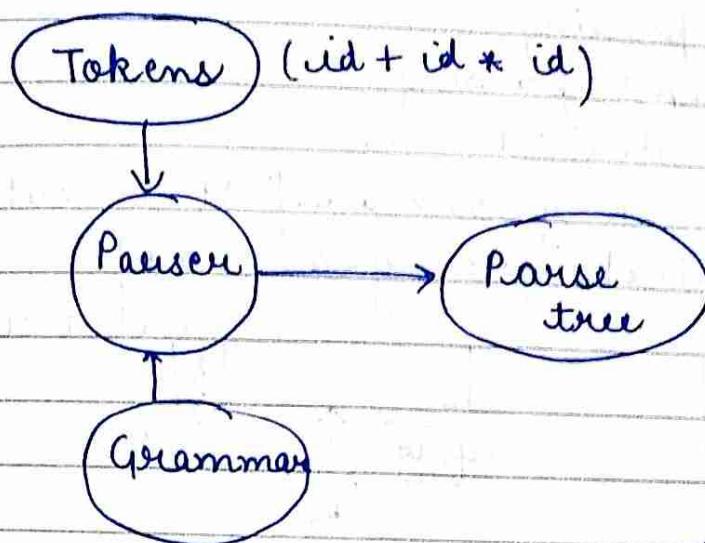
2. Bottom-Up parsing:

It attempts to build the parse-tree from leaf-to-root. It starts from string & proceed to the start symbol. It follows right-most derivation. In right-most derivation, the right most non-terminal in each sentential is always chosen.

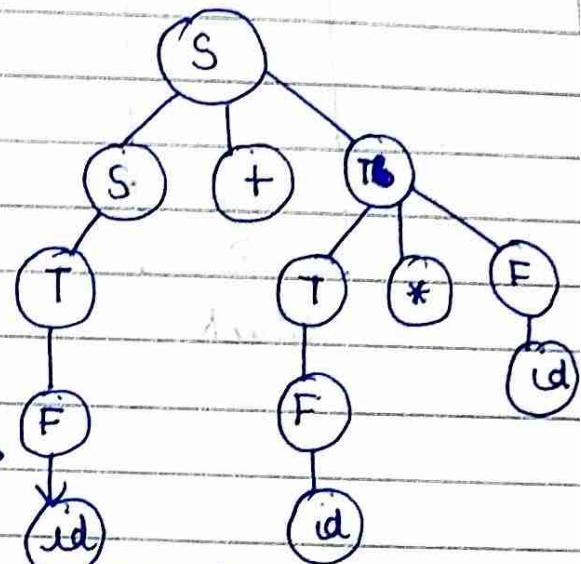
Bottom-Up

Shift reduce (LR)
LR(0) SLR CLR LALR
Operator precedence

Syntax analysis %

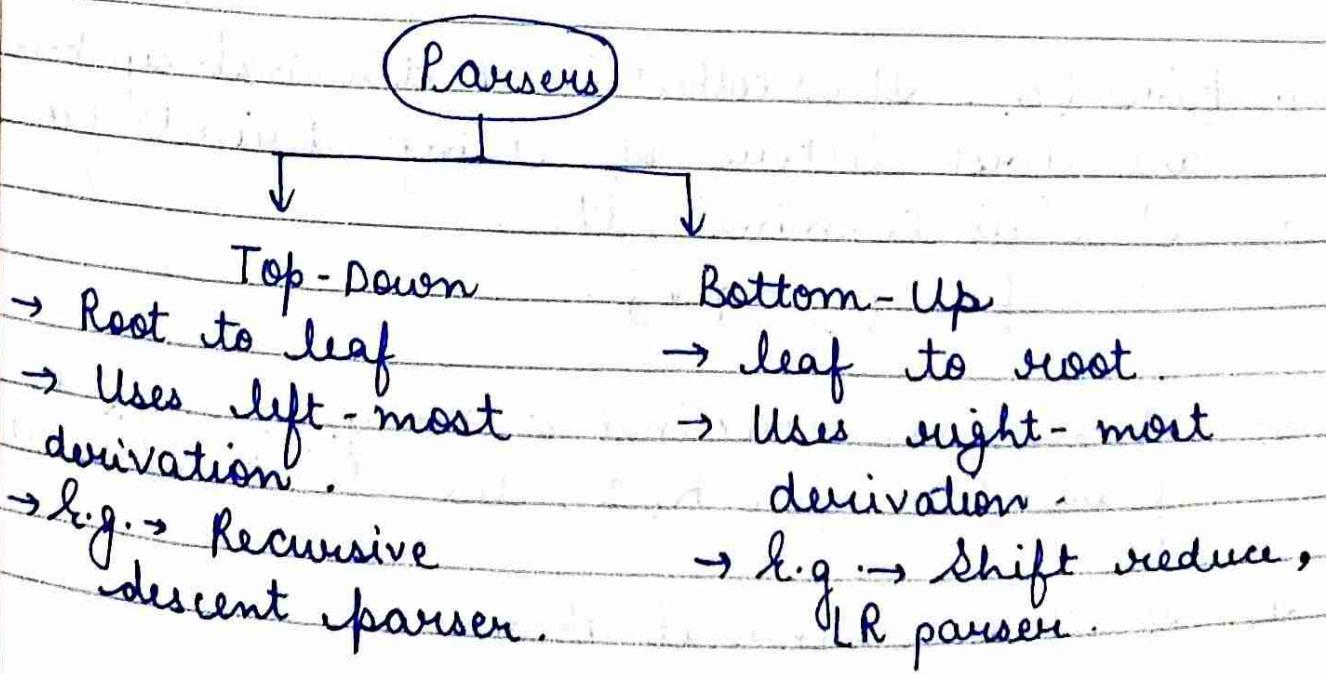


$$\begin{aligned} S &\rightarrow S + T / T \\ T &\rightarrow T * F / F \\ F &\rightarrow id \end{aligned}$$



[Top to down & left] →
to right.

If the yield of the parse tree isn't matched with the input, it means there is some error.



First & Follow:

1. First (A): It is collection of terminal symbols which are first letters of strings derived from A.
 - a) If x is terminal, then
 $\text{first}(x) = \{x\}$.
 - b) If x is non-terminal & $x \rightarrow y_1, y_2 \dots y_K$ is the production for $K \geq 1$, the $\text{first}(x)$ will be
 - b1) If y is terminal then,

$$\text{First}(x) = \text{First}(y_1, y_2, y_3) = \{y_1\}$$

b2) If y_1 is non-terminal then

If y_1 doesn't derive to an empty string i.e. if $\text{first}(y_1)$ doesn't contain ϵ then

$$\text{First}(x) = \text{First}(y_1, y_2, y_3) = \text{first}(y_1)$$

b3) If y_1 contains ϵ then

$$\text{First}(x) = \text{First}(y_1, y_2, y_3) =$$

$$\text{First}(y_1) - \{\epsilon\} \cup \text{First}(y_2, y_3)$$

Similarly, $\text{First}(y_2, y_3) = \text{first}(y_2)$ if it doesn't contain ϵ

If $\text{first}(y_2)$ contains ϵ , then

$$\text{First}(y_2, y_3) = \text{First}(y_2) - \{\epsilon\} \cup \text{First}(y_3)$$

Similarly, this method will be repeated for further grammar symbols i.e. for $y_4, y_5 \dots y_k$

c.) If x is non-terminal & $x \rightarrow ax$ then $\text{First}(x) = a$

d.) If $x \rightarrow \epsilon$, then $\text{first}(x) \rightarrow \epsilon$.

e.g. 1. $S \rightarrow ABC \mid ghi \mid jkl$

$$A \rightarrow a \mid b \mid c$$

$$B \rightarrow b$$

$$D \rightarrow d$$

$$\text{First}(D) \rightarrow d$$

$$\text{First}(B) \rightarrow b$$

$$\text{First}(A) \rightarrow a, b, c$$

$$\text{First}(S) \rightarrow a, b, c, g, k$$

$$(iii) \quad S \rightarrow ABC$$

$$A \rightarrow a/b/\epsilon$$

$$B \rightarrow c/d/\epsilon$$

$$C \rightarrow e/f/\epsilon$$

$\text{First}(c) \rightarrow e, f, \epsilon$

$\text{First}(B) \rightarrow c, d, \epsilon$

$\text{First}(A) \rightarrow a, b, \epsilon$

$\text{First}(S) \rightarrow \text{First}(ABC)$

$\text{First}(A)$ contains ϵ , So, $\text{First}(S) \rightarrow \text{First}(A) - \{\epsilon\} \cup \text{First}(BC)$

$\text{First}(S) \rightarrow a, b \cup \text{First}(BC)$

$\text{First}(B)$ contains ϵ , So, $\text{First}(BC) = \text{First}(B) - \{\epsilon\} \cup \text{First}(C)$

$$\begin{aligned} \text{First}(BC) &= (c, d) \cup (e, f, \epsilon) \\ &= c, d, e, f, \epsilon \end{aligned}$$

So, $\text{First}(S) \rightarrow a, b, c, d, e, f, \epsilon$.

$$(iii) \quad E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

$$F \rightarrow id | (E)$$

$\text{First}(F) \rightarrow id, ($

$\text{First}(T') \rightarrow \epsilon \mid +.$

$\text{First}(T) = \text{First}(F) = \text{id}, (.$

$\text{First}(E') \rightarrow *, \epsilon.$

$\text{First}(E) \rightarrow \text{First}(TE') \rightarrow \text{First}(T) - \text{id}, (.$
 [$\text{First}(T)$ doesn't contain $\epsilon.$]

iv) $S \rightarrow (L) \mid a$
 $L \rightarrow SL'$

$L' \rightarrow , SL' \mid \epsilon$

$\text{First}(S) \rightarrow (, a$

$\text{First}(L) \rightarrow \text{First}(S) \rightarrow (, a$

$\text{First}(L') \rightarrow , , \epsilon.$

v) $S \rightarrow ACB \mid CbB \mid Ba$

$A \rightarrow da \mid BC$

$B \rightarrow g \mid \epsilon$

$C \rightarrow h \mid \epsilon$

$\text{First}(C) \rightarrow h, \epsilon$

$\text{First}(B) \rightarrow g, \epsilon$

$\text{First}(A) \rightarrow d, \text{First}(BC)$

$\text{First}(B)$ contains $\epsilon,$ So, $\text{First}(BC) \rightarrow \text{First}(B) -$

$\{\epsilon\} \cup \text{First}(C)$

$\rightarrow g, h, \epsilon.$

First(A) $\rightarrow d, g, h, \epsilon$.

First(S) \rightarrow First(ACB), First(CB),
First(Ba)

First(ACB) $\rightarrow d, g, h, \epsilon$.

First(CbB) $\rightarrow h, b$.

First(Ba) $\rightarrow g, a$.

So, First(S) $\rightarrow a, b, d, g, h, \epsilon$.

Q. Follow(A):

Follow(A) is defined as collection of terminal symbols that occur directly to the right of A .

Rules:

- If S is the start symbol, Follow(S) $\rightarrow \$$
- If production is of form,

$$A \rightarrow dB\beta, \beta \neq \epsilon$$

- If First(β) contains ϵ , then

$\text{Follow}(B) = \text{First}(\beta) - \epsilon \cup \text{Follow}(A)$
 \because when β derives ϵ , then the terminal after A will follow B .

- If production is in form,
 $A \rightarrow dB$ then

→ * Follow never contains ϵ

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

e.g.

1) $S \rightarrow AaAb \mid BbBa$
 $A \rightarrow \epsilon$
 $B \rightarrow \epsilon$

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

$$\text{Follow}(A) \rightarrow a, b$$

$$\text{Follow}(B) \rightarrow b, a$$

2) $S \rightarrow ABC$

$$A \rightarrow DEF$$

$$B \rightarrow \epsilon$$

$$C \rightarrow \epsilon$$

$$D \rightarrow \epsilon$$

$$E \rightarrow \epsilon$$

$$F \rightarrow \epsilon$$

$$\text{Follow}(A) = \text{First}(BC)$$

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

$$\text{So, } \text{follow}(A) = \text{First}(C)$$

$$\text{First}(C) = \epsilon$$

$$\therefore \text{follow}(A) = \text{follow}(S) = \$$$

3) $S \rightarrow (L) \mid a$

$$L \rightarrow SL'$$

$$L' \rightarrow , SL' \mid \epsilon$$

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

$$\text{Follow}(L) \rightarrow)$$

$$\text{Follow}(L') \rightarrow)$$

4) $S \rightarrow ACB \mid CbB \mid Ba$ $C \rightarrow h \mid \epsilon$
 $A \rightarrow da \mid BC$
 $B \rightarrow g \mid \epsilon$

1.) $\text{Follow}(c) \rightarrow b$, $\text{First}(B)$, $\text{Follow}(A)$

$\text{First}(B) \rightarrow g, \epsilon$

so, $\text{Follow}(c) \rightarrow b, g$, $\text{Follow}(S)$, $\text{Follow}(d)$

2.) $\text{Follow}(A) \rightarrow \text{First}(c) \rightarrow h, \epsilon$

$\therefore \text{Follow}(A) = \text{Follow}(\text{First}(B)) \cup h$

$\text{Follow}(A) \rightarrow f, g, h$.

$\text{Follow}(S) \rightarrow \$$.

$\therefore \text{Follow}(c) \rightarrow b, g, h, \$$

3.) $\text{Follow}(B) \rightarrow a, \$$, $\text{First}(c)$

$\rightarrow a, \$, h$, $\text{Follow}(A)$

$\rightarrow a, \$, g, h$.

5.) $E \rightarrow TE'$

$E' \rightarrow +TE' | E$

$T \rightarrow FT'$

$T' \rightarrow *FT' | E$

$F \rightarrow (E) | id$

i) $\text{Follow}(E) \rightarrow \text{Follow}(E')$, $\text{Follow}(T')$, $,$, $\$$

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

$\text{Follow}(T) \rightarrow +$, $($, id

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

$\text{Follow}(F) \rightarrow \text{First}(T') \rightarrow *, \text{First}(E) \rightarrow (, \text{id}, *$

Final:

$\text{Follow}(E) \rightarrow \$,), +, (, \text{id}.$

$\text{Follow}(E') \rightarrow \$,), +, (, \text{id}$

$\text{Follow}(T) \rightarrow +, (, \text{id}$

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

$\text{Follow}(F) \rightarrow (, \text{id}, *$

6.) $E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) | \text{id}$

$\text{Follow}(E) \rightarrow \$,)$

$\text{Follow}(E) \rightarrow \$,)$

$\text{Follow}(T) \rightarrow +, \$,)$

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

$\text{Follow}(F) \rightarrow *, +, \$,).$

Parsing

Top-Down Parsing

Predictive Parsing

LALR(k) Parsing

left recursion

left factoring.

Bottom-Up

Shift left reduce

LALR(k) Parsing

Predictive Parsing:

1. Construct the parsing table from the given grammar.
2. Apply the predictive parsing algorithm to build the parse tree.

Also to construct parsing table,

For each production $A \rightarrow \alpha$ its following steps:

1. For each terminal a in $\text{first}(\alpha)$, add $A \rightarrow \alpha$ to Table $[A, a]$.
2. If ϵ is in $\text{first}(\alpha)$, then add $A \rightarrow \alpha$ to Table $[A, b]$ for each terminal b in $\text{Follow}(A)$.
3. If ϵ in $\text{first}(\alpha)$ & $\$$ in $\text{Follow}(\alpha)$, then add $A \rightarrow \alpha$ to Table $[A, \$]$.

E.g: Grammar: $E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) | id$

First

$E \quad \{(, id\}$

$E' \quad \{+, \epsilon\}$

$T \quad \{(, id\}$

$T' \quad \{\ast, \epsilon\}$

$F \quad \{(, id\}$

Follow

$\{ \), \$ \}$

$\{ \), \$ \}$

$\{ +, \), \$ \}$

$\{ +, \), \$ \}$

$\{ \ast, +, \), \$ \}$

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

Apply the predictive parsing algo to construct parse tree.

Set up ipt to point to first symbol of input string $w\$$.

repeat

 if $\text{Top}(\text{Stack})$ is terminal or $\$$ then

 if $\text{Top}(\text{Stack}) == \text{current}(\text{ipt})$ then

$\text{Pop}(\text{Stack})$ & $\text{ipt}++$.

 else

 Null.

 else

 if $\text{Table}[x, a] = x \rightarrow y_1 y_2 y_3$ then

 begin

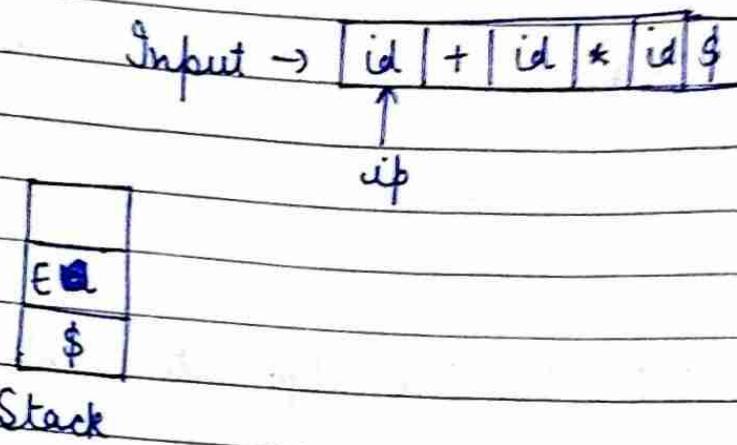
$\text{Pop}(\text{Stack})$;

 Push $y_1 y_2 \dots y_k$ onto stack with y_1 on top.

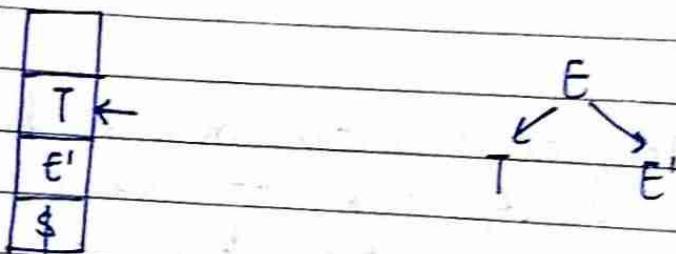
Output the production $x \rightarrow y_1 y_2 \dots y_k$

1

end
else
Null
until ($\text{Top}(\text{Stack}) = \$$ (i.e. Stack become empty)).



(i) E is non-terminal,
Table $[E, \text{id}] = E \rightarrow TE'$



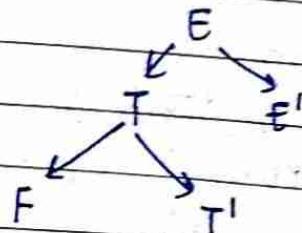
(ii) T is non-terminal

Table $[T, \text{id}] = T \rightarrow FT'$

Pop (T)

Push (F, T')

F
T'
E'
\$



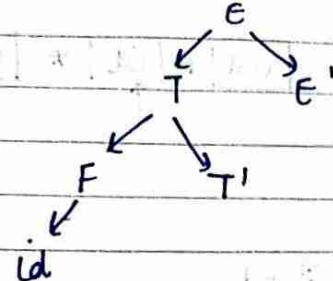
(iii) F is non-terminal,

Table $[F, id] = F \rightarrow id$.

Pop(F)

Push(id)

id
T'
E'
\$



T'
E'
\$

(iv) Top(Stack) = terminal

Top(Stack) = id, Yes

Pop(id)

up++

id	+	id	*	id	\$
----	---	----	---	----	----

(v) T' is non-terminal

Table $[T', +] = T' \rightarrow E$

Pop(T')

Push(E)

E	T	E'
E'		
\$	F	T'

(vi) E is non-terminal but null

Table $[E, +] = E \rightarrow \emptyset$

so, simply Null.

(vii) E' is non-terminal

Table $[E', +] \rightarrow +TE'$

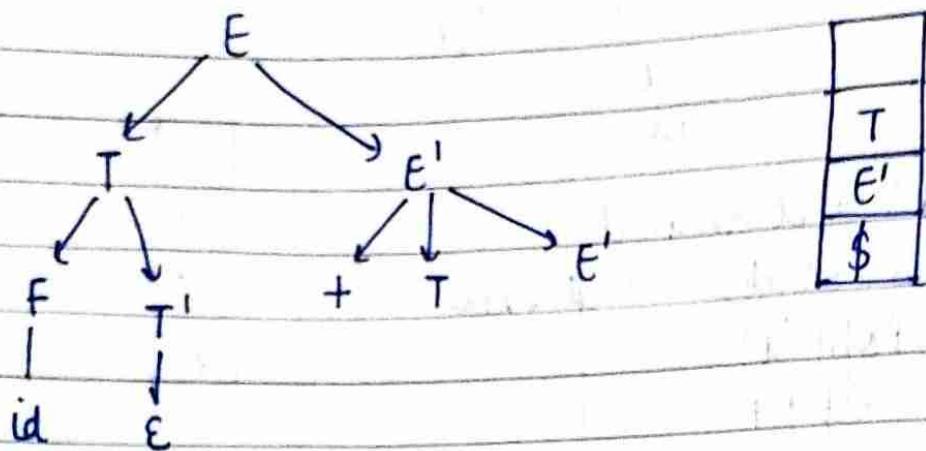
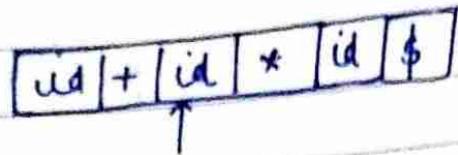
Pop(E')

Push(+TE')

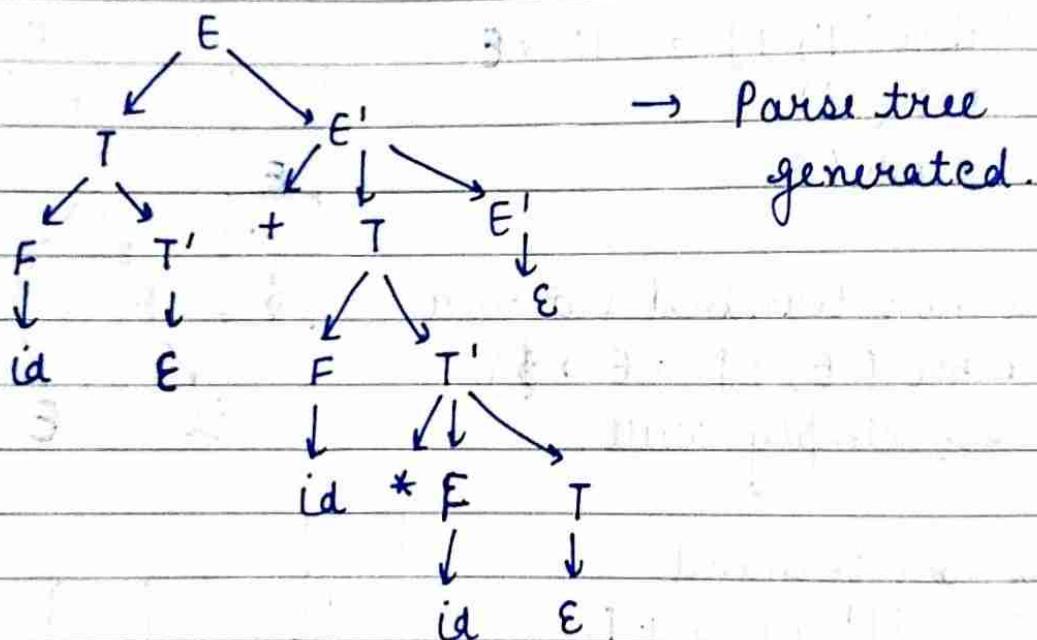
E'
\$

+
T
E'
\$

(viii) + is non-terminal
 $+ \leftarrow$ current (up), Yes
 Pop (+)
 $\text{up}++$



$$T \rightarrow FT'$$



When $\text{Top}(\text{Stack}) = \$ = \text{input}$, the parse tree is completed, parser halts & accepts the strings.

Question on predictive Parsing:

$$1.) \quad S \rightarrow (L) \mid a$$

$$L \rightarrow SL'$$

$$L' \rightarrow)SL' \mid \epsilon$$

	First	Follow
S	{(, a}	{), \$}
L	{(, a}	{)}
L'	{), ε }	{)}

	a	()	\$	
S	$S \rightarrow a$	$S \rightarrow (L)$			
L	$L \rightarrow SL'$	$L \rightarrow SL'$			
L'			$L' \rightarrow)SL'$	$L' \rightarrow \epsilon$	$L' \rightarrow \epsilon$

Here, we can see that there are 2 productions in the same cell. Hence, this grammar is not feasible for LL(1) Parser.

$$2.) \quad S \rightarrow A \mid a$$

$$A \rightarrow a$$

	First()	Follow()
S	{a}	{ \$ }
A	{a}	{ \$ }

Date _____
Page _____

a \$

S $S \rightarrow A, S \rightarrow a$

A $A \rightarrow a$

This grammar is not feasible for LL(1)
Parser.

Note: In above 2 examples, the grammar
is ambiguous grammar. So the grammar
doesn't satisfy the essential conditions.

LL(k) Parser:

- This parser parses from left to right, does a left most derivation. It looks up i^k symbol ahead to choose its next action.
- LL(1) grammar is a grammar whose parsing table has no multiply-defined entries.
- LL(1) grammar is non-ambiguous & non-left recursive.

Handle & Handle pruning:

A handle is a substring that matches the body of production i.e. ~~right~~ side of production that can be reduced to ~~left~~ - side of production. So, the reduce of handle represents one step along the reverse of a right most derivation.

If $A \rightarrow \alpha$

Then α is a handle that can be reduced to A .

E.g.: Grammar $\rightarrow S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

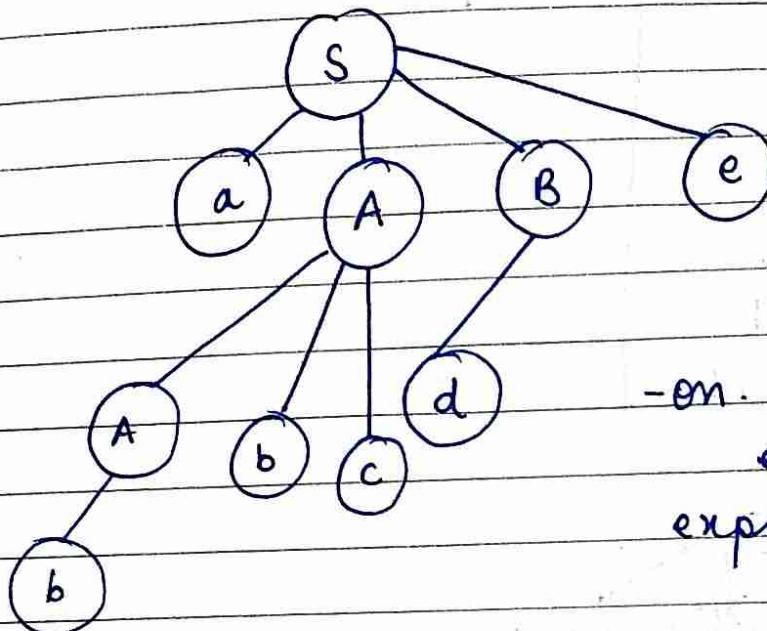
Input String $\rightarrow abbcde$.

Right sentinal form	Handle	Reduction.
\underline{abbcde}	b	$A \rightarrow b$
$a\overline{Abc}de$	Abc	$A \rightarrow Abc$

aAde
aABe
S

d
aABe

B \rightarrow d
S \rightarrow aABe



As, Bottom-up uses right-most derivation. Then firstly, B is expanded & A is expanded.

S
S \rightarrow aABe
S \rightarrow aAde
S \rightarrow aAbcde
S \rightarrow abbcde

So, handle pruning is the process of removing children of left hand side non-terminal from parse.

A right most derivation in reverse can be obtained using handle pruning.

Q. E \rightarrow E+E | E*E | (E) | id

Right Sentimental
id + id * id
E + id * id

Handle
id
id

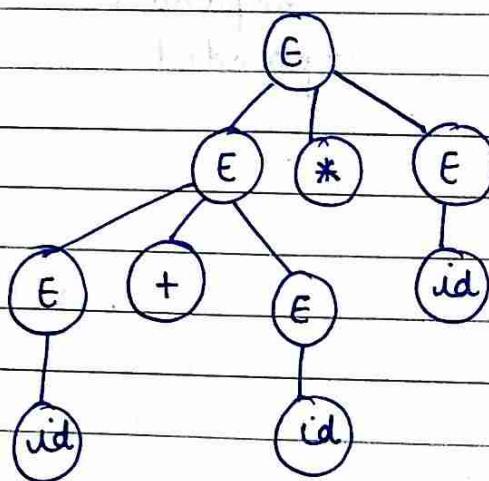
Reduction
E \rightarrow id
E \rightarrow id

$$\begin{array}{c} E + E * id \\ \underline{E * id} \\ E * E \\ \underline{E} \end{array}$$

$$\begin{array}{c} E + E \\ id \\ E * E \end{array}$$

$$\begin{array}{c} E \rightarrow E + E \\ E \rightarrow id \\ E \rightarrow E * E \end{array}$$

E	
$E \rightarrow E * E$	
$E \rightarrow E * id$	
$E \rightarrow E + E * id$	
$E \rightarrow E + id * id$	
$E \rightarrow id + id * id$	



Note:

A certain sentinel form may have many different handles.

Problems:

1. Locate a handle.
2. Decide which production to use (if there are more than two candidate productions).

Shift reduce parsing :

Shift reduce parsing is a process of reducing a string to the start symbol of grammar.

It uses a stack to hold the grammar & an input tape to hold the string.

String $\xrightarrow{\text{reduce to}}$ Start Symbol

Shift reduce parsing performs two actions:

Shift & reduce. At shift action, the current symbol in input ~~symbol~~ string is pushed to a stack. At each reduction, the symbols will be replaced by the non-terminals. The symbol is right-side of production & non-terminal is left-side of production.

$$\text{E.g. 1.) } S \rightarrow S + S$$

$$S \rightarrow S - S$$

$$S \rightarrow (S)$$

$$S \rightarrow a$$

$$W = a1 - (a2 + a3)$$

Stack	Input	Action
\$	a1 - (a2 + a3)\$	Shift
\$ a1	- (a2 + a3)\$	Reduce $S \rightarrow a$
\$ S	- (a2 + a3)\$	Shift
\$ S -	(a2 + a3)\$	Shift
\$ S - a	a2 + a3)\$	Reduce S $\rightarrow a$ Shift
\$ S - a	+ a3)\$	Shift
\$ S - a S +	a3)\$	Shift, Reduce $S \rightarrow a$

\$ S - (S + S)) \$	Reduce S → a
\$ S - (S + S)) \$	Shift
\$ S - (S + S)	\$	Reduce S → S +
\$ S - (S)	\$	Reduce S → (S)
\$ S - S	\$	Reduce S → S -
\$ S	\$	Accept

2.) $E \rightarrow E - E$

$E \rightarrow E * E$

$E \rightarrow id$

$w = id - id * id$

Stack	Input	Action
\$	id - id * id \$	Shift
\$ id	- id * id \$	Reduce E → id
\$ E	- id * id \$	Shift
\$ E -	id * id \$	Reduce Shift
\$ E - id	* id \$	Reduce E → id
\$ E - E	* id \$	Reduce E → E -
\$ E	* id \$	Shift
\$ E *	id \$	Shift
\$ E * id	\$	Reduce E → id
\$ E * E	\$	Reduce E → E * E
\$ E	\$	Accept.

3.) $S \rightarrow T L ; T \rightarrow \text{int} / \text{float} ; L \rightarrow L, id / id$
 $w = \text{int } id, id$

Stack	Input	Action
\$	int id, id \$	Shift
\$ int	id, id \$	Reduce T → int

\$T	id , id \$	Shift
\$T id	, id \$	Reduce L \rightarrow id
\$TL	, id \$	Shift
\$TL,	id \$	Shift
\$TL, id	\$	Reduce L \rightarrow h, id
\$TL	\$	Reduce S \rightarrow TL
\$S	\$	Accept.

4. $S \rightarrow OSO \mid ISI \mid a$ $w = 10201$

Stack	Input	Action
\$	10201\$	Shift
\$1	0201\$	Shift
\$10	201\$	Shift
\$102	01\$	Reduce S \rightarrow a
\$10S	01\$	Shift
\$1OSO	1\$	Reduce S \rightarrow OSO
\$IS	1\$	Shift
\$ISI	\$	Reduce S \rightarrow ISI
\$S	\$	Accept.

Operator Precedence:

An operator precedence parser is bottom-up parser that interprets an operator grammar (grammar having restriction that null production is not allowed or no 2 adjacent non-terminals in its right-hand side). Ambiguous grammars are not allowed in any parser.

except operator precedence.

Rule 1 → This parser relies on three precedence relations : \prec , \equiv , \succ

- 1) $a \prec b$ This means a "yields precedence to" b
- 2) $a \succ b$ This means a "takes precedence over" b.
- 3) $a \equiv b$ This means a "has same precedence as" b.

Rule 2 :- 1.) An identifier is always given the higher precedence than any other symbol.
2.) \$ symbol is always given lowest precedence.

Rule 3 :- 1.) If two operators have same precedence, then we go by checking their associativity.

Algo steps:-

1. Insert \$ symbol at beginning & end of input string.
2. Insert precedence operator b/w every two symbols of the string by referring the operator operator precedence table.
3. Start scanning string from L.H.S. in forward direction until \succ is encountered.
4. Keep a pointer on that.
5. Start scanning the string from R.H.S. in backward direction until \prec is encountered.
6. Keep a pointer on that location.

7. everything that lies in middle of 1 & 2 forms the handle.

8. Replace the handle with head of the respective production rule.

Repeat step 3 to 8, until the start symbol is reduced.

E.g.:

$$1.) E \rightarrow E+E \mid E * E \mid id \quad @ \rightarrow \quad w = id * id + id$$

	id	+	*	\$
id	▷	▷	▷	
f	▷	▷	▷	
+	◁	▷	◁	▷
*	◁	▷	▷	▷
\$	◁	◁	◁	

\$ id * id + id \$

\$ <. id .> * <. id .> + <. id .> \$

\$ E * E + E \$

\$ E + E \$

\$ E \$ → \$ \$

$$2.) S \rightarrow (L) \mid a \quad w = (a, (a, a))$$

	a	()	,	\$
a	▷	▷	▷	▷	▷
(◁	▷	▷	▷	▷
)	◁	▷	▷	▷	▷
,	◁	◁	▷	▷	▷
\$	◁	◁	◁	◁	

$\$(a, (a, a)) \$$

$\$ \leftarrow (\underline{a}, \underline{(\underline{a})}, \underline{(\underline{a})}) \rightarrow \$$

$\$ \leftarrow (S, \underline{a} \rightarrow \$((\underline{a})) \rightarrow \$$

$\$ \leftarrow (S \rightarrow (S, \underline{a} \rightarrow \$) \rightarrow \$$

$\$ \leftarrow (\underline{S} \rightarrow (S \rightarrow \$) \rightarrow \$$

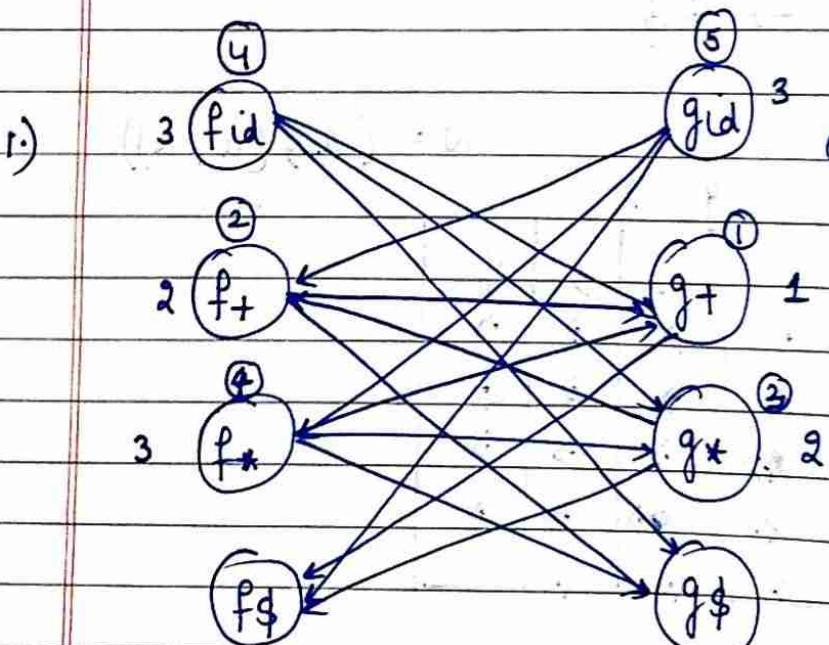
$\$ \leftarrow (\underline{S} \rightarrow (S \rightarrow \$) \rightarrow \$$

$\$ \leftarrow (\underline{L} \rightarrow (\underline{S} \rightarrow \$) \rightarrow \$$

$\$ \leftarrow (\underline{L} \rightarrow \$ \rightarrow \$$

$\$ \leftarrow (L) \rightarrow \$$

$\$ \leftarrow \$$



Q: Find operator
precedence operat
functions.

Resulting precedence functions are:

	id	+	*	\$
f	4	2	4	0
g	5	1	2+1	0

→ [refer to longest path from a symbol to \$ sign].

left factoring:

If R.H.S of more than one production starts with the same symbol, then such a grammar is called grammar with common prefixes. This kind of grammar creates a problem for top down parsers because they can't decide which production rule must be chosen to parse the string. So, we use left factoring.

Left factoring is a process by which the grammar with common prefixes is transformed to make it useful for top down parser.

How?

1. We make one production for each common prefixes. The common prefix may be terminal non-terminal or combination of both.
2. Rest of derivation is added by new productions.

e.g.: 1.) $S \rightarrow uEtS \mid uEtSeS \mid a$
 $E \rightarrow b$

→ $S \rightarrow uEtSS' \mid a$

$$S' \rightarrow eS \mid \epsilon.$$

$$E \rightarrow b$$

2.) $A \rightarrow aAB \mid aAC \mid aBc$

$$\rightarrow A \rightarrow aA'$$

$$A' \rightarrow AB \mid Ac \mid Bc$$

$$A' \rightarrow AA_1 \mid Bc$$

$$A_1 \rightarrow B \mid c$$

Again a left factored grammar

3.) $S \rightarrow bSSaas \mid bSSasb \mid bsb \mid a$

$$\rightarrow S \rightarrow bSS' \mid a$$

$$S' \rightarrow Saas \mid S asb \mid b$$

$$S' \rightarrow SSS \mid a \mid SAS_1 \mid b$$

$$S_1 \rightarrow aS \mid asb$$

4.) $E \rightarrow T+E \mid T$

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$

$$\rightarrow E \rightarrow TT'$$

$$T' \rightarrow +E \mid \epsilon.$$

$$T \rightarrow \text{int} T_1 \mid (E)$$

$$T_1 \rightarrow *T \mid \epsilon.$$

5.) $S \rightarrow A$

$$A \rightarrow aB \mid ac \mid ab \mid AC$$

$$B \rightarrow bBC \mid f$$

$$C \rightarrow g \mid \epsilon$$

$$D \rightarrow d \mid \epsilon.$$

$\rightarrow S \rightarrow A$

$A \rightarrow aa' | AC$

$A' \rightarrow B | C | b$

$B \rightarrow bBC | f$

$C \rightarrow g | \epsilon$

$D \rightarrow d | \epsilon$

6.) $S \rightarrow A \# 0$

$A \rightarrow Ad | aB | ac$

$C \rightarrow c$

$B \rightarrow bBc | \epsilon$

$\rightarrow S \rightarrow A \# 0$

$A \rightarrow Ad | aA'$

$A' \rightarrow B | C$

$C \rightarrow c$

$B \rightarrow bBc | \epsilon$

7.) $S \rightarrow abc | abd | ae | f$

(
 $S \rightarrow aS' | f$

(
 $S' \rightarrow bc | bd | e$

Ans: $S \rightarrow aS' | f$

$S' \rightarrow bS_1 | e$

$S_1 \rightarrow c | d$.

Left recursion : A grammar G is left recursive if it has production in the form,

$A \rightarrow A\alpha | \beta$

The above grammar is left recursive because the left of production occurring at first position is also present at right hand side of production.

It can eliminate left recursion by replacing a pair of production with

$$A \rightarrow A\alpha / \beta$$

$$\rightarrow A \rightarrow \beta A'$$

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

E.g:- 1. $A \rightarrow ABd / Aa / a$

$$B \rightarrow Be / b$$

$$\rightarrow A \rightarrow aA'$$

$$A' \rightarrow BdA' / aA' / \epsilon$$

$$B \rightarrow bB'$$

$$B' \rightarrow eB' / \epsilon$$

2. $E \rightarrow E+E / E * E / a$

$$\rightarrow E \rightarrow aE'$$

$$E' \rightarrow +EE' / *EE' / \epsilon$$

3. $E \rightarrow E+T / T$

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

$$F \rightarrow id$$

$$\rightarrow E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

$$F \rightarrow id$$

4. $S \rightarrow (L) | a$

$$L \rightarrow L, S | S.$$

$$\rightarrow S \rightarrow (L) | a$$

$$L \rightarrow SL'$$

$$L' \rightarrow , SL' | \epsilon$$

5. $S \rightarrow SOS1S | 01$

$$S \rightarrow 01S'$$

$$S' \rightarrow 0SISS' | \epsilon$$

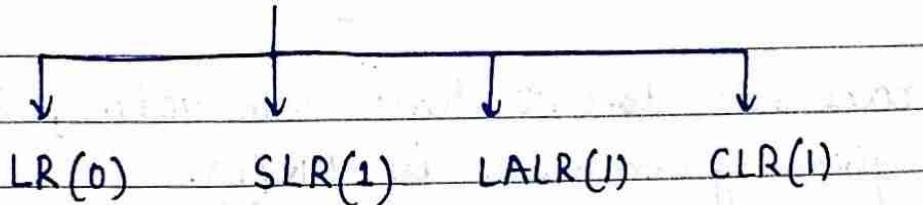
LR Parser:

LR(k) parser is type of bottom-up parser.

It is used to parse large class of grammars.

In LR(k) parsing, "L" stands for left-to-right scanning of the input. "R" stands for constructing right-most derivation. "k" stands for number of input symbols of look ahead used to make number of parsing decision.

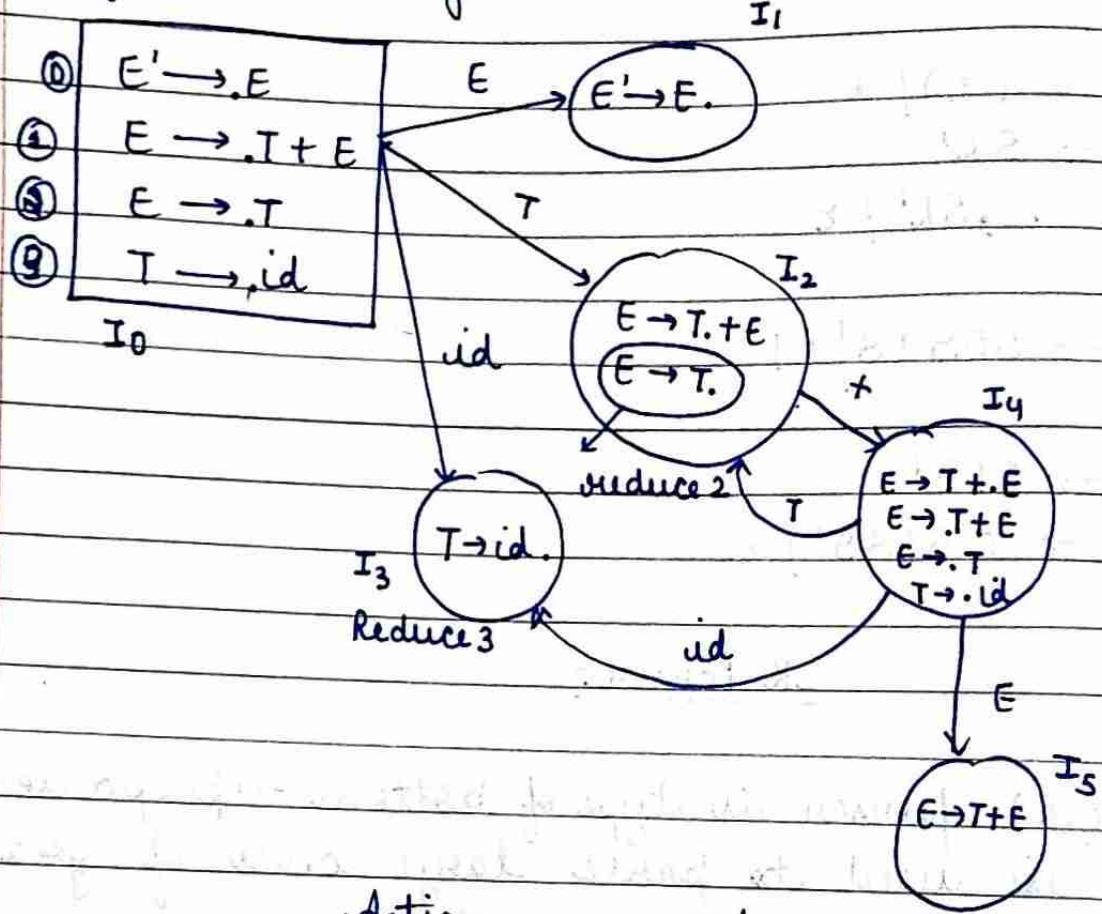
LR Parser



LR(0) 0

$$1.) \quad E \rightarrow T+E/T \\ T \rightarrow id$$

Augment the grammar :

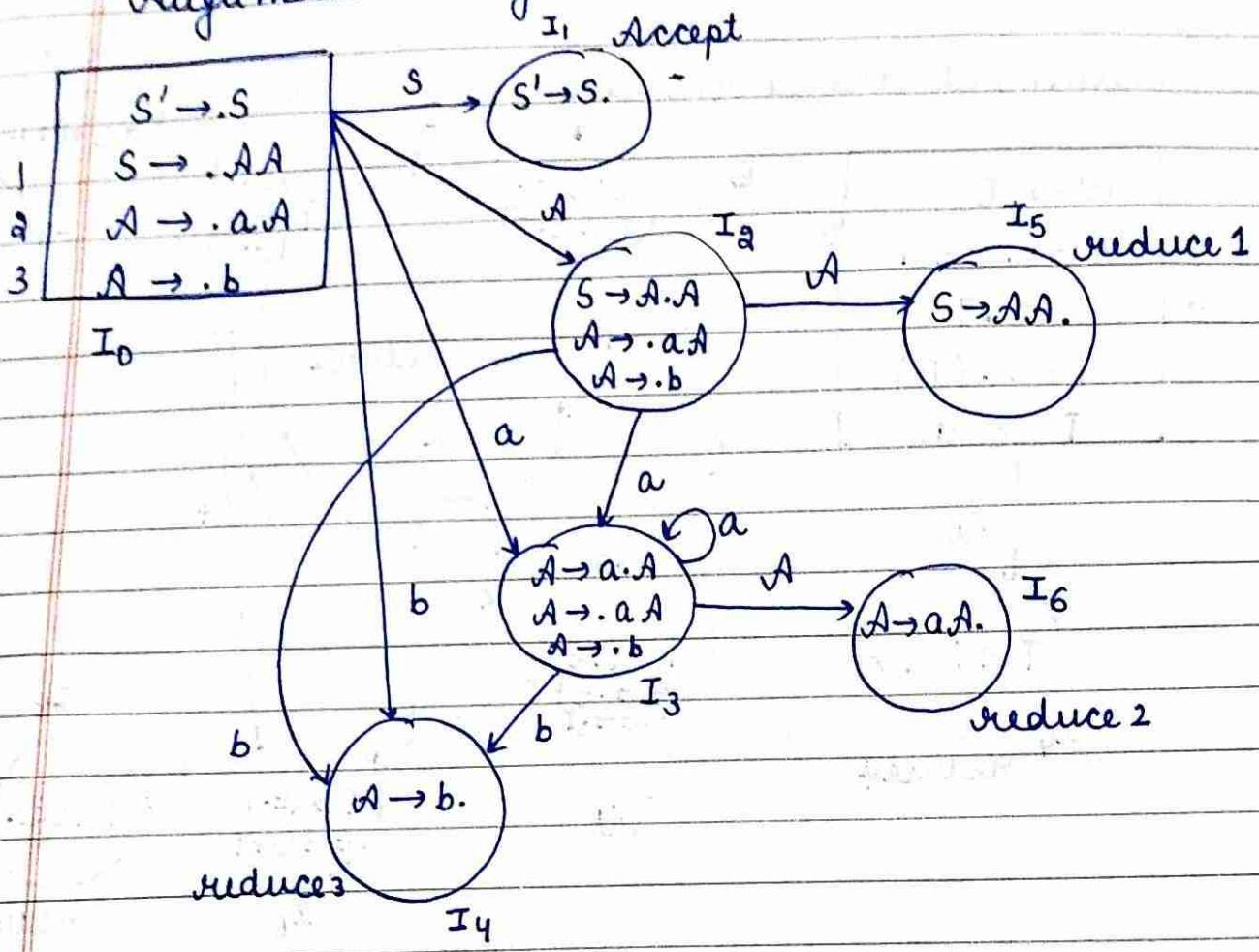


Action	goto		reduce
State	$id + \$$	$E \quad T$	
0	S_3	1 2	
1	Accept		
2	$M_2 \quad S_4/M_2 \quad M_2$		
3	$M_3 \quad M_3 \quad M_3$		
4		5 2	
5	$M_1 \quad M_1 \quad M_1$		

Each cell doesn't have one value, hence the given grammar is LR(0).

3. $S \rightarrow A.A$
 $A \rightarrow a.A / b$

Augment the grammar:

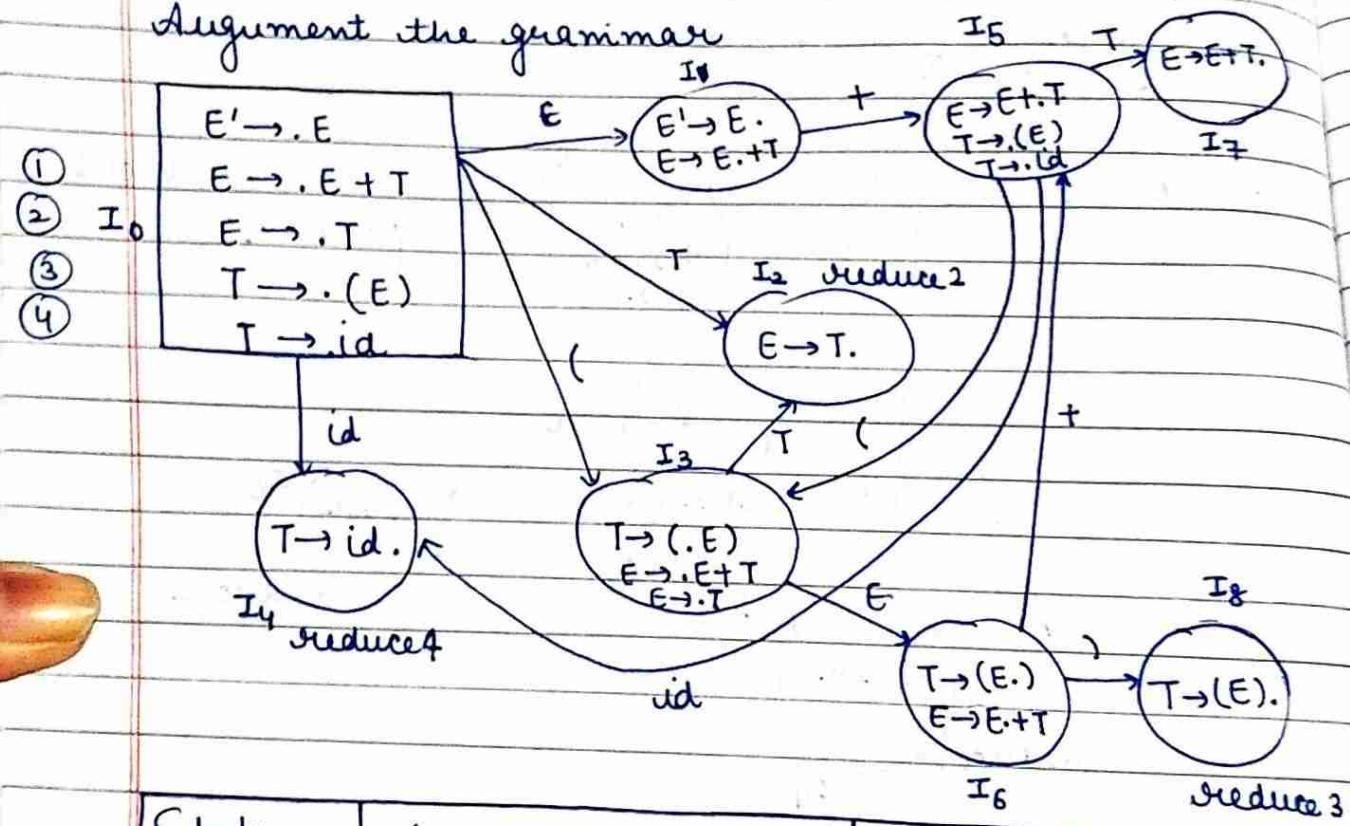


State	Action			goto	
	a	b	\$	S	A
0	S_3	S_4	-	1	2
1	-	Accept	-	-	-
2	S_3	S_4	-	5	-
3	S_3	S_4	-	6	-
4	H_3	H_3	H_3	-	-
5	H_1	H_1	H_1	-	-
6	H_2	H_2	H_2	-	-

The given grammar is LR(0).

- 3.
- $$E \rightarrow E + T$$
- $$E \rightarrow T$$
- $$T \rightarrow (E)$$
- $$T \rightarrow id$$

Augment the grammar



State	(+	id	\$	E	T
0	S_3		S_4		1	2
1		S_5				
2	g_1	g_1	g_1	g_1		
3						
4	g_1	g_1	g_1	g_1	6	2
5	S_3		S_4			
6	S_5	S_8				7
7	g_1	g_1	g_1	g_1		
8	g_1	g_1	g_1	g_1		

Action goto

Augmented grammar:

If G is a grammar with starting symbol S , then G' (augmented grammar) is a grammar with a new starting symbol S' & productions $S \rightarrow S'$. The purpose of this new starting production is to indicate the parser when it should stop parsing. The $:$ symbol indicates 2 things: the left side of $:$ has been read by a compiler & the right side of $:$ is yet to be read by a compiler.

Steps for constructing LR parsing table:

1. Write augmented grammar.
2. LR(0) collection of items to be found.
3. Defining 2 functions: goto (list of non-terminals) and action (list of terminals) in the parsing table.

Note: LR(0) economical item is a production A with dot at some position on right side of production. LR(0) items are useful to indicate how much of the input has been scanned up to a given point in process of parsing.

LR(0) Table:

1. If a state is going to some other state on terminal then it corresponds to a shift move.
2. If a state is going to some other state

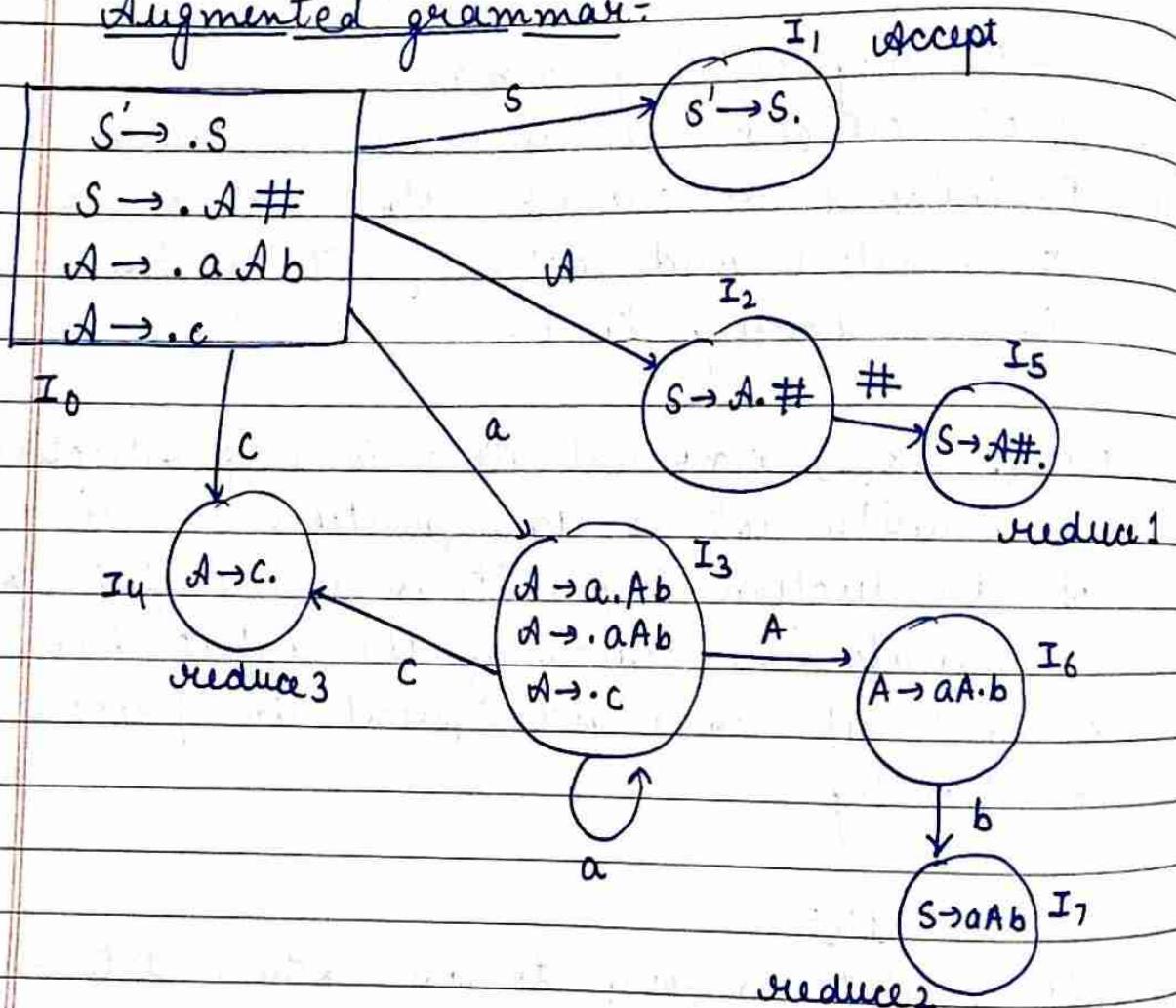
on variable then it corresponds to go to move.

3. If a state contain final item in a particular row then write the reduce node completely.

e.g.:

$$4. \quad S \rightarrow A \# \\ A \rightarrow a.A.b / c$$

Augmented grammar:



State	a	c	b	#	\$	s	A
0	S_3	S_4	S_4			1	2
1				Accept			
2					S_5		
3	S_3	S_4					6
4	M_3	M_3	M_3	M_3	M_3		
5	M_1	M_1	M_1	M_1	M_1		
6					S_7		
7	M_2	M_2	M_2	M_2	M_2		

SLR(1) = (simple LR)

1. SLR is simple LR. It is the smallest class of grammar having few number of states.
2. SLR is very easy to construct & is similar to LR parsing. The only diff. between SLR(1) and LR(0) parser is that in LR(0) parsing table, there is chance of 'shift-reduced' conflict because we are entering 'reduce' corresponding to all terminals of state but in SLR 'reduce' is entered only in "Follow" of L.H.S of production in terminating state. This is called SLR(1) collection of items

Steps :-

1. Writing augmented grammar.
2. LR(0) collection of items to be found.
3. Find Follow of L.H.S of production.
4. Defining 2 functions: 1.) go to [list of non-terminals], 2.) action [terminals].

E.g. :-

Follow

$S \rightarrow A.A$

$S \rightarrow \{ \$ \}$

$A \rightarrow aA / b$

$A \rightarrow \{ \$, ab \}$

Augmented grammar :-

Accept

1.) $S' \rightarrow S$

$S \rightarrow A.A$

$A.B \rightarrow a.A$

$A \rightarrow .b$

 I_0

$S \rightarrow S' \rightarrow S.$ I_1

 A

I_2
 $S \rightarrow A.A$
 $a \rightarrow a.A$
 $A \rightarrow .b$

 I_5 reduce 1

$S \rightarrow AA.$

 a I_3

$A \rightarrow a.A$
 $A \rightarrow .ab$
 $A \rightarrow .b$

 a A I_6

$A \rightarrow aA.$

reduce 2

 I_4

$A \rightarrow b.$

reduce 3

State	a	b	\$	S	A
0	S_3	S_4	-	1	2
1	-	Accept	-	-	-
2	S_3	S_4			5
3	S_3	S_4			6
4	H_3	H_3	H_3		
5				H_1	
6	H_2	H_2	H_2		

$$a. \quad E \rightarrow E + T \mid T$$

$$T \rightarrow TF \mid F$$

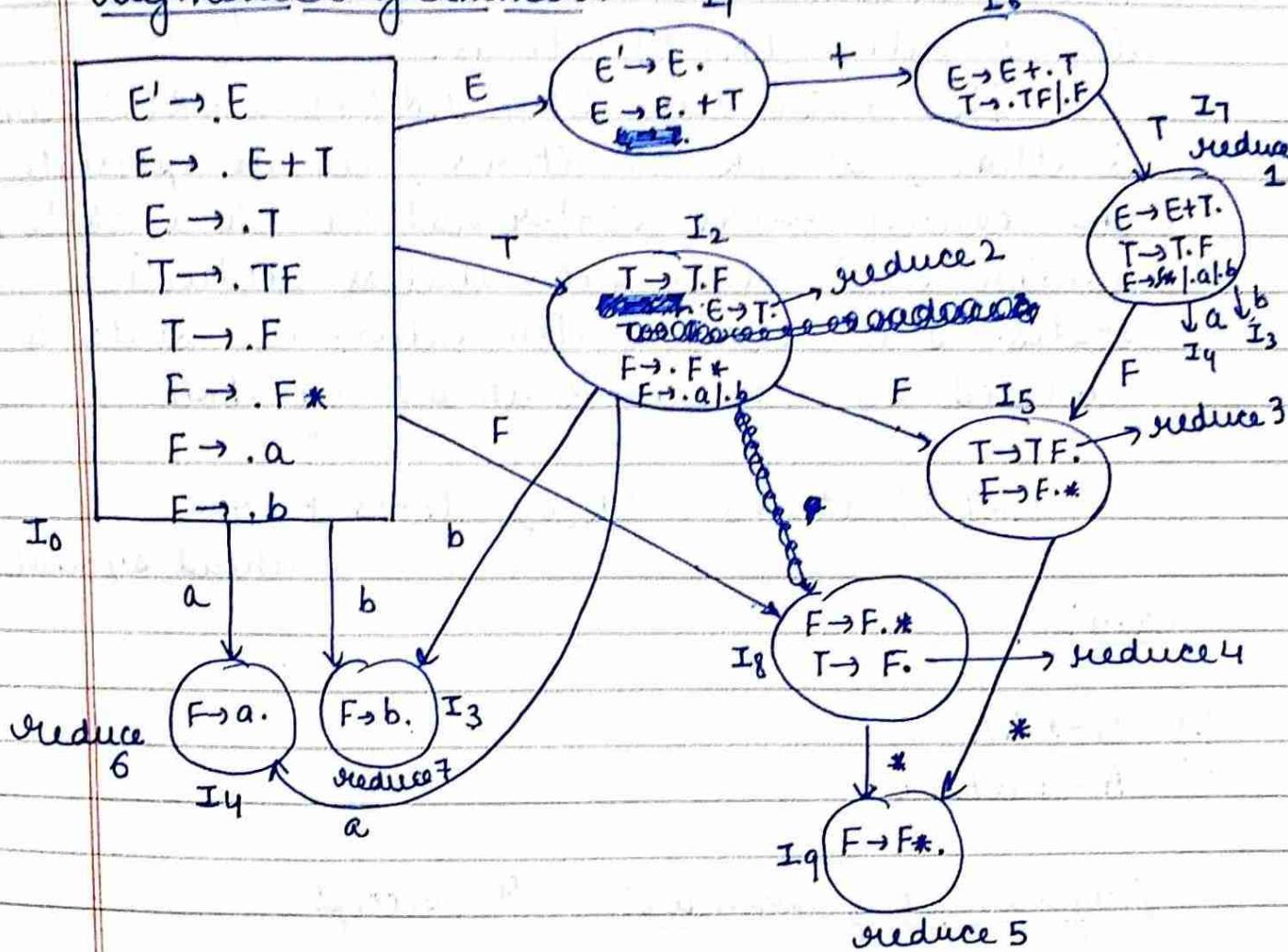
$$F \rightarrow F^* \mid a \mid b$$

$$\text{Follow: } E \rightarrow \{ +, \$, \}$$

$$F \rightarrow \{ *, \$, a, b \}$$

$$T \rightarrow \{ \$, a, b \}$$

Augmented grammar:



State	+	*	a	b	\$	E	T	F
0			S_4	S_3		1	2	8
1	S_6							
2	S_2		S_4	S_3	S_2			5
3		S_7	S_7	S_7	S_7			
4		S_6	S_6	S_6	S_6			
5		S_9	S_3	S_3	S_3			
6							7	
7	S_1	S_4	S_4	S_3	S_1			
8		S_9	S_4	S_4	S_4			5
9		S_5	S_5	S_5	S_5			

CLR(1) parser:

CLR parser stands for canonical LR parser. It is more powerful LR parser. It makes use of lookahead symbols. This methods make use of large set of items called LR(1) items.

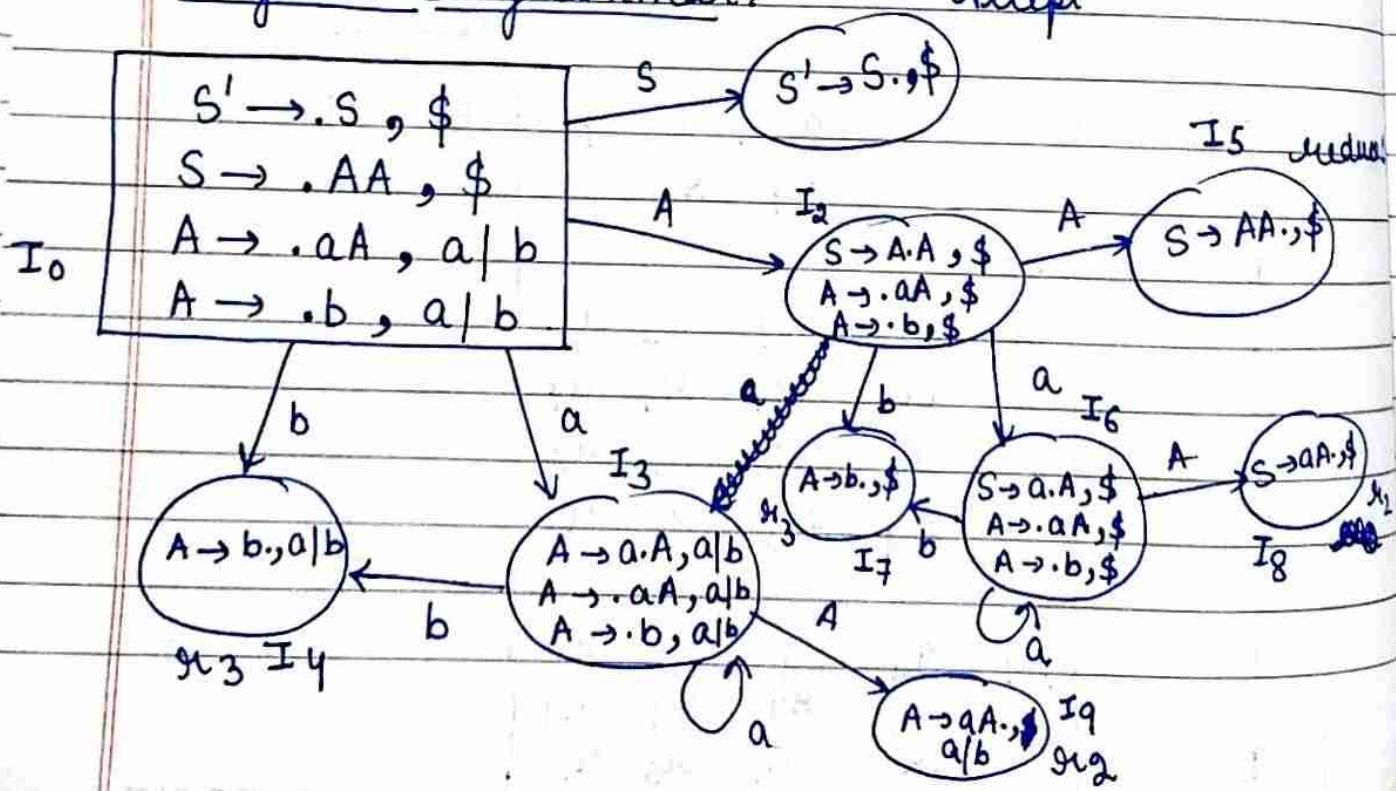
The main diff. b/w LR(0) & LR(1) items is that, in LR(1) items, it is possible to carry more information in a state, which will rule out unless reduction states. This extra information in state is carried by the lookahead symbol.

LR(1) items = LR(0) items + lookahead symbol

e.g:-

$$\begin{aligned} 1.) \quad S &\rightarrow AA \\ A &\rightarrow aA \mid b. \end{aligned}$$

Augmented grammar: I_1 Accept



State	a	b	\$	s	A
0	S_3	S_4		1	2
2	S_6	S_7			5
3	S_3	S_4			9
4	M_3	M_3			
5			M_1		
6	S_6	S_7			8
7			M_3		
8			M_2		
9	M_2	M_2			

Note: [The reductions will be written in the corresponding look ahead symbols of that state.]

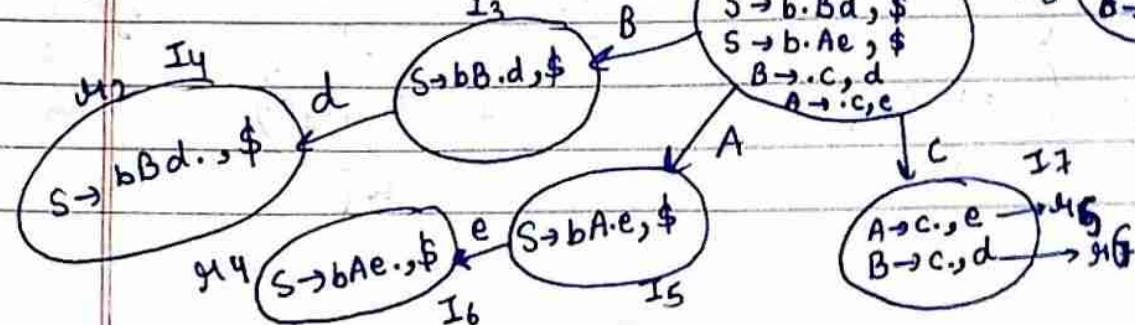
a.) $S \rightarrow aAd \mid bBd \mid aBe \mid bAe$

$$A \rightarrow c$$

$$B \rightarrow c$$

Augmented grammar

I_0	$S' \rightarrow .S, \$$
	$S \rightarrow .aAd, \$$
	$S \rightarrow .bBd, \$$
	$S \rightarrow .aBe, \$$
	$S \rightarrow .bAe, \$$



State	a	b	c	e	d	\$	S	A	B
0		s_{13}		s_{12}			1		
1	-	-	Accept	-	-	-	-	-	.
2					s_7			5	3
3						s_4			
4						s_9	s_2		
5					s_6				
6							s_4		
7					s_5	s_6			
8					s_6	s_5			
9					s_{10}				
10							s_3		
11						s_{12}			
12							s_1		
13					s_8			11	9

$$3. \quad S \rightarrow AaAb \mid BbBa$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Augmented grammar:

$$S' \rightarrow .S, \$$$

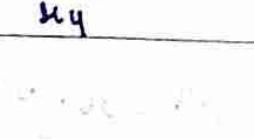
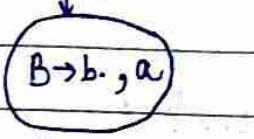
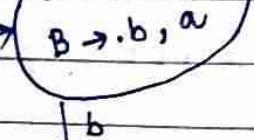
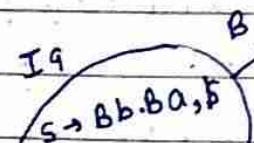
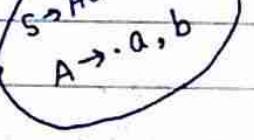
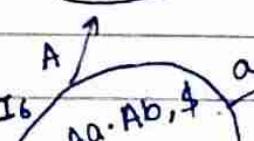
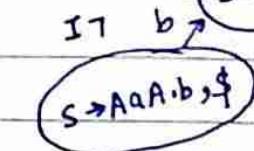
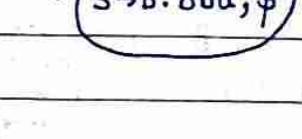
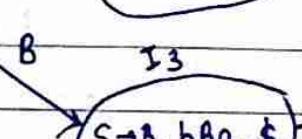
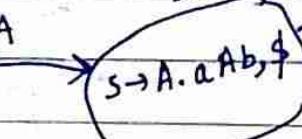
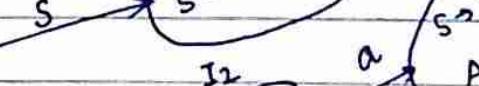
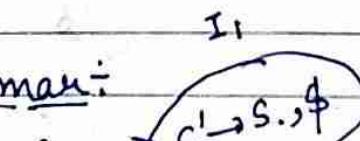
$$S \rightarrow .AaAb, \$$$

$$S \rightarrow .BbBa, \$$$

$$A \rightarrow .a, a$$

$$B \rightarrow .b, b$$

I0



State	a	b	\$	S	A	B
0	s_4		s_5	1	2	3
1	-		- Accept	-	-	-
2		s_6				
3	s_3		s_9			
4	s_3		s_9			
5			s_4			
6	s_{11}				7	
7			s_8			
8				s_1		
9						12
10	s_4					
11			s_3			
12		s_{13}				
13			s_2			

LALR(1) :

LALR is lookahead LR parser. It is powerful parser which can handle large class of grammar. The size of CLR table is quite large as compared to other parsing table. LALR reduced the size of this table. otherwise LALR is similar to CLR. The only difference is , it combines the similar states of CLR parsing table into one single state.

So, LALR also makes use of LR(1) canonical items.

$$LR(1) = LR(0) + \text{look ahead symbols.}$$

For example,

In last example, I have done in CLR. The state I_4 & I_{11} is same i.e. both are reduced 3 state but the difference is in their look aheads. So, we can combine these states like.

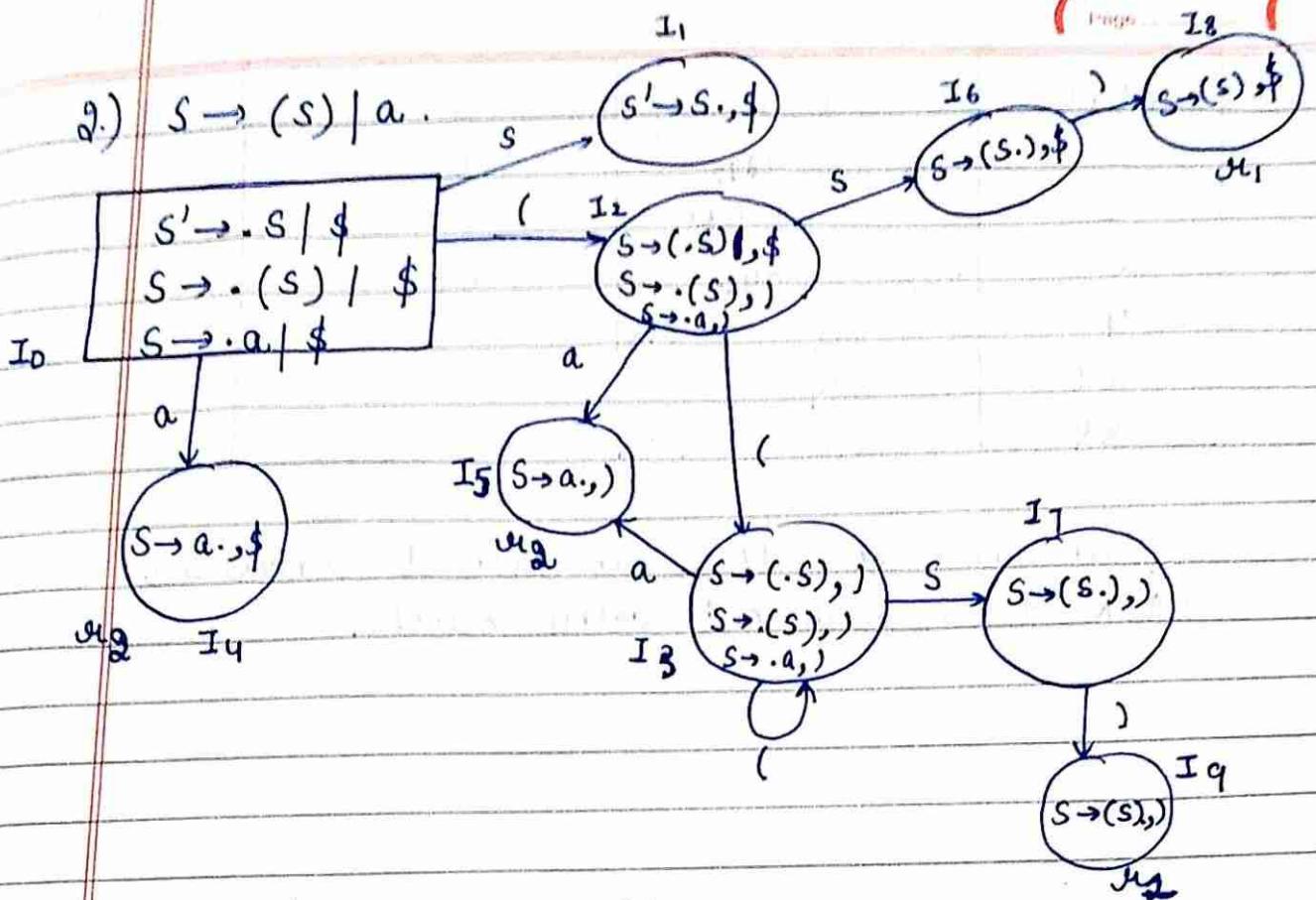
$$I_{4+11} \rightarrow A \rightarrow a, a/b.$$

Similarly, the state I_5 & I_{10} can be combined to form,

$$I_{5+10} \rightarrow B \rightarrow b, b/a.$$

Table of LAIR:

State	a	b	\$	S	A	B
0	S_{4+11}		S_{5+10}	1	2	3
1	-	Accept	-	-	-	-
2	S_6					
3			S_q			
4+11	M_3		M_3			
5+10	M_4		M_4			
6	S_{4+11}				7	
7		S_8				
8				M_1		
9			S_{5+10}			12
12	S_{13}					
13				M_2		



Now, in above example, we can see.

→ I_2 & I_3 are same except lookahead symbols. So, they can be combined.

I_{23}	$S \rightarrow (.S), \$)$
	$S \rightarrow .(S),)$
	$S \rightarrow .a,)$

→ I_4 & I_5 are both ε states. So,
 $S \rightarrow a., \$ |)$.

→ I_6 & I_7 are same,
 $S \rightarrow (S.), \$ |)$

→ I_8 & I_9 are both ε states. So, combined state = $S \rightarrow (S), \$ |)$.

State	() a \$	S
0	s_{23}	1
1	- - Accept -	-
23	s_{23}	67
45	m_2	m_2
67	s_{89}	
89	m_1	m_1

So, we can see the number of states in LALR has reduced very much.

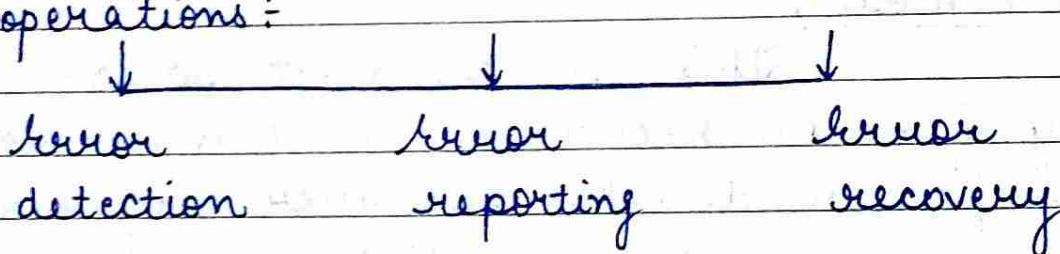
Unambiguous
grammar.

CLR LALR SLR(1)

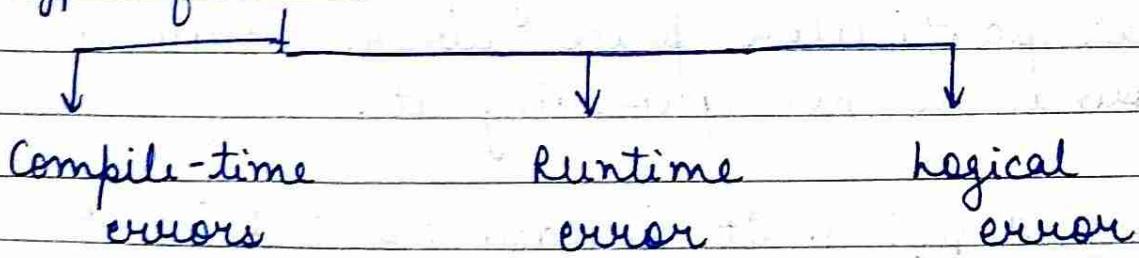
Error recovery techniques:

Error is the operation which cause an abnormal change or problems in a program.

In the phase of compilation, the errors are sent in the form of messages. This is called error handling process. Error handler performs all the functions that are supposed to performed while detecting an error in a compiler. It perform 3 operations:



Types of errors:



Compile-time errors occur while compiling a program & they occur due to incorrect syntax, like when we miss a semicolon at end of program or when we miss putting brackets at end of any loop or conditional statement etc.

Runtime errors occurs during execution of a program. These are the errors which arise due to invalid user input. For e.g. → An code

Panic mode e.g: int a, 5gh, b, fb;

//After int a, 5gh, b, fb; // parser discards
input symbol at a time

It takes input of type integer, but the user entered the character input gives an error at runtime.

Logical errors are types of errors in which the given code is unreachable or there is a presence of infinite loop.

Error recovery techniques:

1. Panic mode:

This is the most basic & easiest way of error recovery. In this, as soon as the parser detects an error anywhere in the program it immediately discards the rest of the remaining part of code from that particular place where error is found by not processing it.

Advantages: 1. It is easy to implement.

2. It never goes in an infinite loop.

Disadvantage:

1. A huge amount of input is discarded without processing that & it does not detect any more errors as it is not processing it further.

2. Statement mode / Phase level recovery:

In this mode, on discovering an error,

it performs recovery by correcting the remaining input so that remaining input statement allows parser to move further. So, the parser can easily continue its job.

In this correction includes, removing extra semicolon, fitting missing semicolons, replacing the comma with semicolon etc. This type of correction is decided by the compiler developer.

E.g. →

int a, b

// After recovery

int a, b; // Semicolon is added by compiler

Advantage: 1. It is widely used in many error repairing compiler.

Disadvantage: 1. While doing replacement in the program, it may fall into infinite loop.

3. Error production:

If a compiler designer has some knowledge of errors that are very common then these types of errors are recovered by augmented (adding extra productions) grammar for corresponding language with error productions that produce incorrect results. If error production is used

during parsing, we can generate appropriate error message to indicate error has been recognized in input.

It is very difficult to maintain, because if we change grammar, then it become necessary to change corresponding productions.

E.g.: Input string \rightarrow abcd

Grammar: $S \rightarrow A$

$A \rightarrow aA \mid bA \mid a \mid b$

$B \rightarrow cd$

The input string is not obtained by above grammar, so we need augmented grammar.

$S' \rightarrow SB$

$S \rightarrow A$

$A \rightarrow aA \mid bB \mid a \mid b$

$B \rightarrow cd$.

Now, string can be obtained.

Advantages:

1. Syntactic phase errors are generally recovered by error production.

Disadvantage:

1. It is difficult to maintain.

4. Global correction:

In this method, the parser judges or considers the full program & try to find the closest matching of this which is free from errors. When an incorrect statement is tackled, it creates a parse tree for error free statement which is closest to input. Global corrections increase time & space requirement at parsing time, so it can't be implemented practically.

Advantage: It makes very few changes in processing an incorrect input.

Disadvantage: It is theoretical concept.

5. Symbol table:

In semantic errors, errors are recovered by using ~~error~~ symbol table for corresponding identifier & if data types of two operands is not compatible, automatically type conversion is done by compiler.

Advantage: It allows basic type-conversion, which we generally do in real-life.

Disadvantage: Only implicit type conversion is possible.

YACC: (Yet another compiler compiler)

YACC is a computer program for the Unix operating system. It is an LALR(1) [lookahead, left-to-right, right most derivation producer with the 1 lookahead] parser generator. YACC was originally designed for being complemented by lex.

The input of YACC is rule or grammar & the output is C program.

Structure of YACC program:

Declarations

%.%

Rules

%.%

Auxiliary functions

%.%

1. Declarations section consists of two parts :

i.) C declarations.

ii.) YACC declarations.

C declarations contains all the declarations required for the C code in the auxilliary section. YACC copies the contents of this section in generated y.tab.c file without any modification.

The following abstract outline of structure of declaration part:

* Beginning of declaration part */

% {

// C declarations

% } .

// YACC declarations

/* End of declaration part */

% %

The YACC declarations part comprises of declarations of token (returned by lexical analyzer). The parser reads the tokens by invoking yylex() function.

2. Rules consists of 2 parts i.) the production part and ii) action part. The syntax of any programming language will be specified in form of Context Free Grammar.

A rule in YACC is of the form:

production-head : production-body { action_in_c };

% %

/* Rules section begins */

/* Rules section ends */

% %

Example of a production:

expr : expr '+' expr
| head body

expr → non-terminals

'+' → terminal

We can give names to these tokens which are defined in declaration section.

* Action part of rule consists of C statements enclose within { and }. These statements are executed when the input is matched with body of production & a reduction takes place.

E.g.

expr: DIGIT{printf ("NUM%d", pos);}

3. The auxiliary functions section contains the definitions of three mandatory functions `main()`, `yylex()`, `yyerror()` & can also contain user-defined functions. The main function must invoke `yparse()` to parse the input.

For e.g. YACC program:

```
%{  
/* Declaration part */  
#include <stdio.h>  
#include <conio.h>  
  
/* function declaration */  
void print_operator (char op)  
  
/* Variable declaration */  
int pos = 0;  
%}
```

```
/* YACC declarations */  
% token DIGIT  
% left '+'  
% left '*'  
%%
```

```
/* Rules section */
```

Start : expr '/n'

```
{ exit (1); }  
;
```

```
expr: expr '+' expr  
{ print_operator ('+') ; }  
| expr '*' expr  
{ print_operator ('*') ; }
```

```
1 DIGIT  
{ printf ("NUM %d", pos);  
};  
%
```

/* Auxiliary function section */

```
void print-operator (char c)
```

```
switch (c)
```

```
{
```

```
case '+':
```

```
    printf ("Plus");
```

```
    break;
```

```
case '*':
```

```
    printf ("MULTIPLY");
```

```
    break;
```

```
}
```

```
return;
```

```
}
```

```
yyerror (char const *s)
```

```
{
```

```
    printf ("error %s", s);
```

```
yylex ()
```

```
{
```

```
    char c;
```

```
    c = getchar();
```

```
    if (isdigit (c)) {
```

```
        pos++;
    }
```

```
    return DIGIT;
}
```

else

{

 return c;

}

}

main()

{

 yyparse();

 return 1;

}

Function :-

