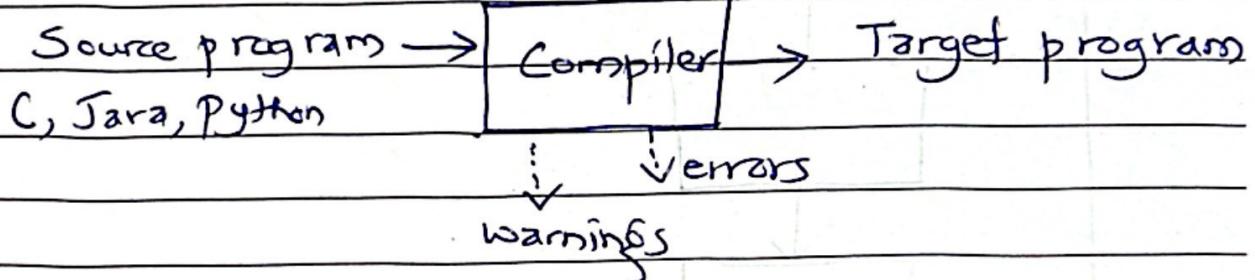


COMPILER DESIGN

Compiler:

It is a software/program that converts a program written in High level language (source language) to a low level language (Object / Target language).

- It also reports errors present in source programs -



Types of compilers.

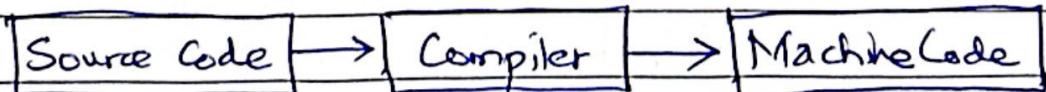
1). Single pass compilers

It is type of compiler that processes the source code only once.

E.g. Early versions of Pascal compilers

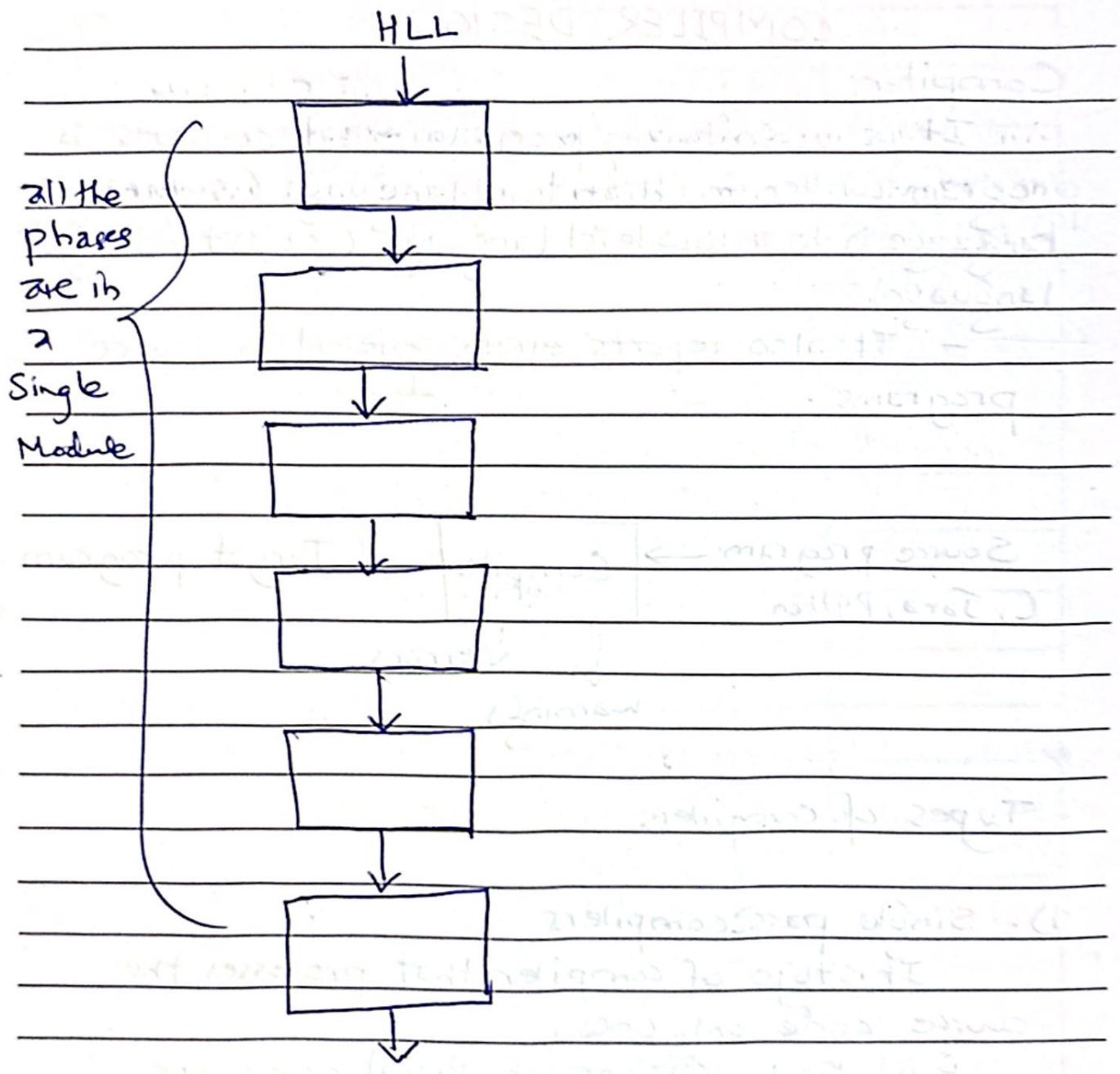
Advantage: faster compilation

Disadvantage: Limited optimization capabilities



SMTWTHFS

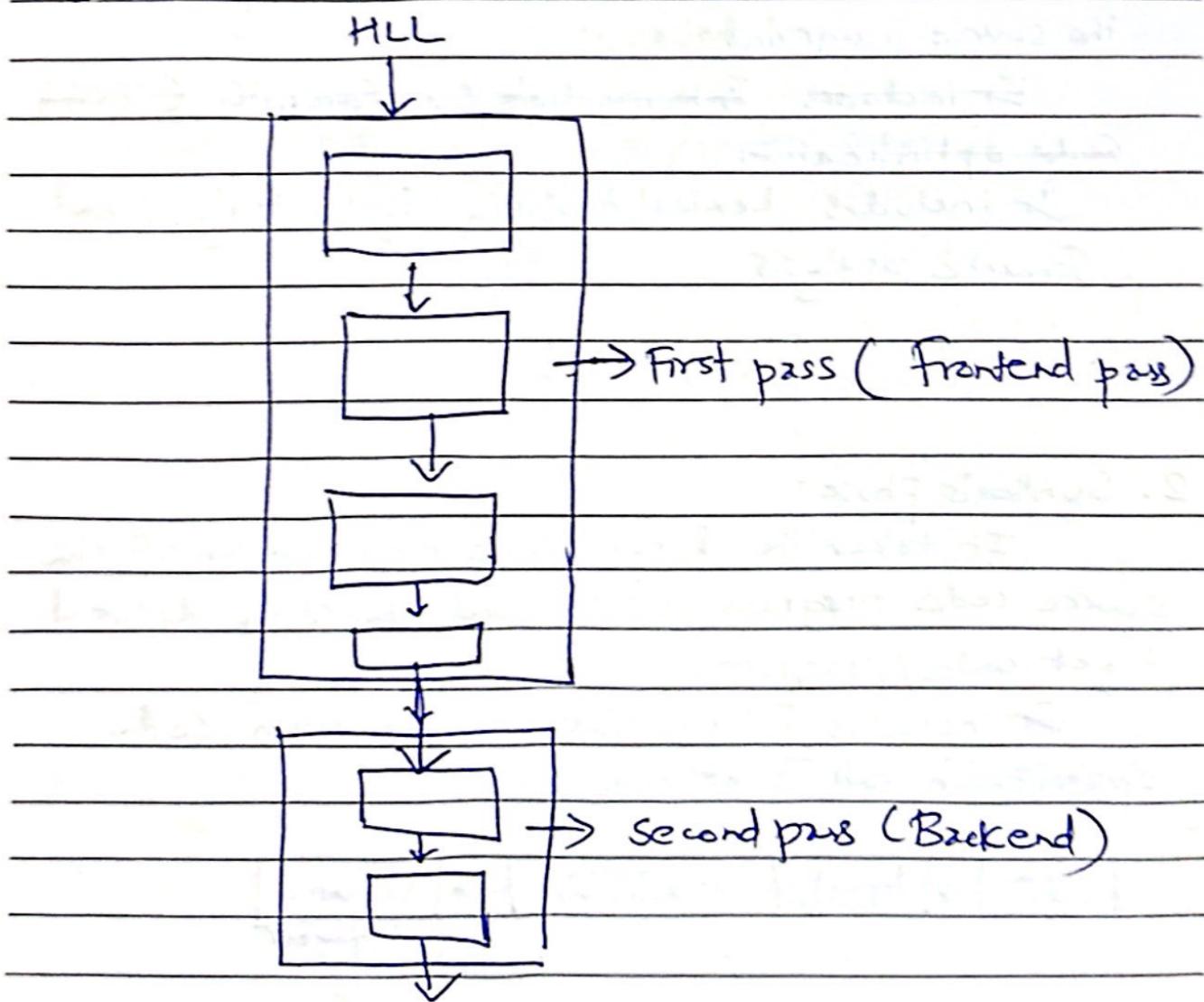
Date



KHWOPA

2) Multipass Compiler

It is a type of compiler that processes source code multiple times (to convert HLL \rightarrow Low level language) i.e. to convert source code to target/object code.



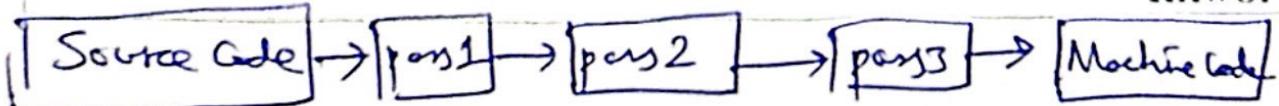
Low level language

Advantage: Better optimization and Error Handling

Disadvantage: Slower compilation

E.g.: Modern C++ Compilers (GCC, Clang)

KHWOPA



Compilation Process

1. Analysis Phase:

It breaks down the source code / program into small parts (tokens, patterns, lexeme) and checks syntax and semantics. Creates an intermediate representation of the source program.

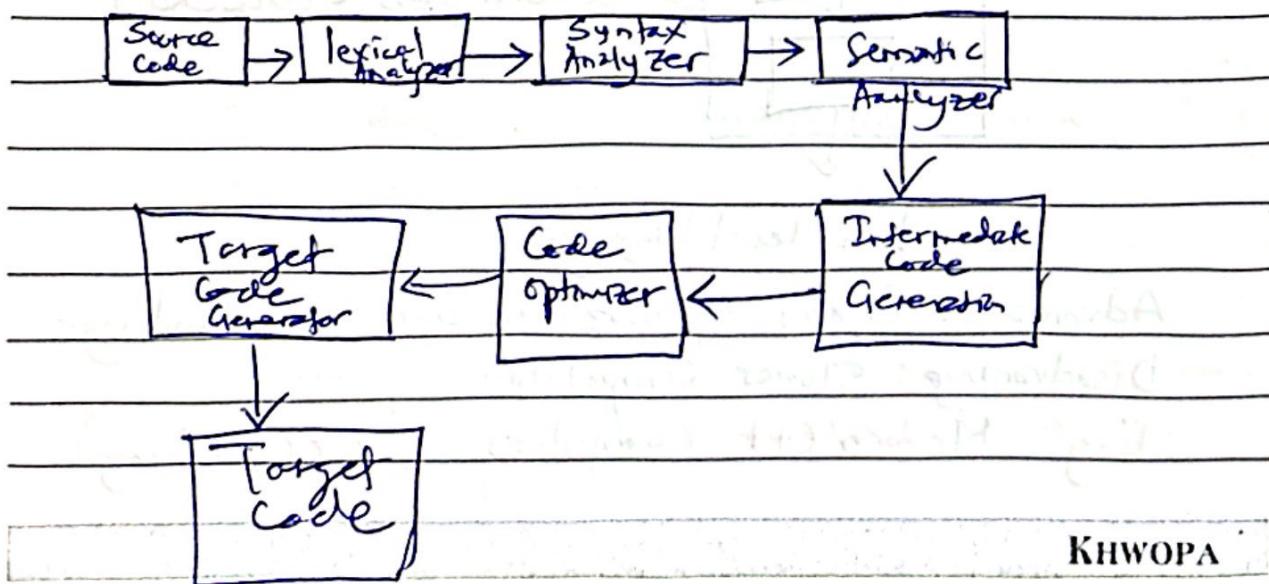
It includes Intermediate Code Generation (ICG), Code optimization

It includes Lexical Analysis, Syntax Analysis and Semantic analysis

2. Synthesis Phase:

It takes the Intermediate representation of the source code program as input and creates the desired target code / program.

It includes Intermediate Code Generation, Code optimization and Target code



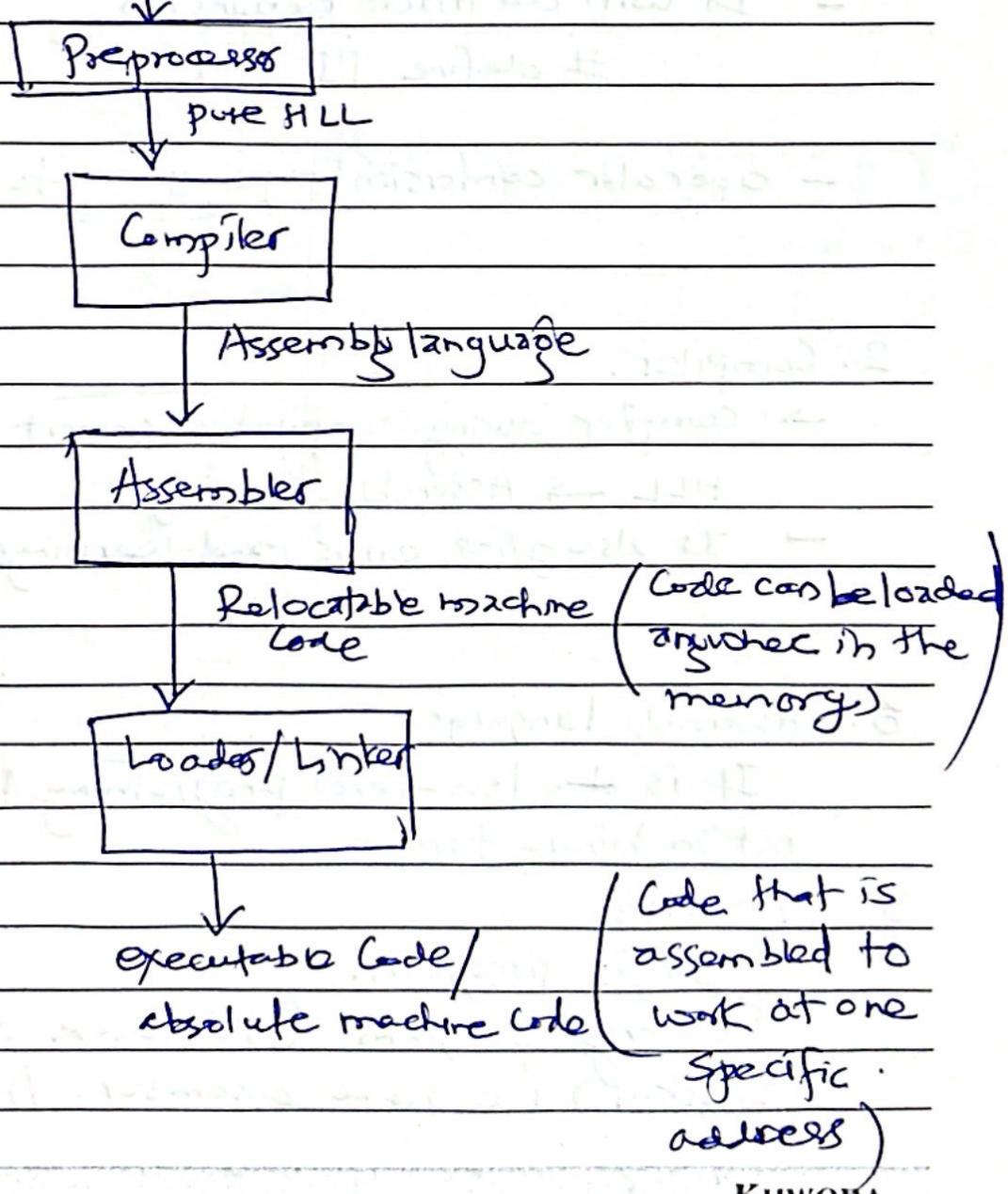
KHWOPA

Language Processing Systems.

In Compiler Design we write the programs in High level language which is easy for us to understand.

These programs are fed into a series of tools and components to get the desired code that can be used by the machine code. This is called language.

Input = High Level Language / Source Code



1. Pre-processor

In Preprocessing, HLL Converted to Pure HLL

It handles macros, file inclusion and conditional compilation.

- Preprocessor removes the preprocessor directives (`#include <stdio.h>`) and will add the respective files i.e. file inclusion.

- It will do macro expansion

`#define PI 3.14`

- Operator conversion e.g. `a--` to `a=a-1;`

2. Compiler.

→ Compiler during compilation convert pure HLL → Assembly language.

→ It also gives errors and warnings.

3. Assembly language

It is ~~not~~ low-level programming language. It is not in binary form.

3. Assembler

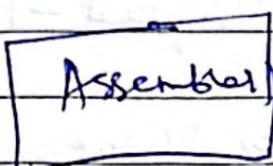
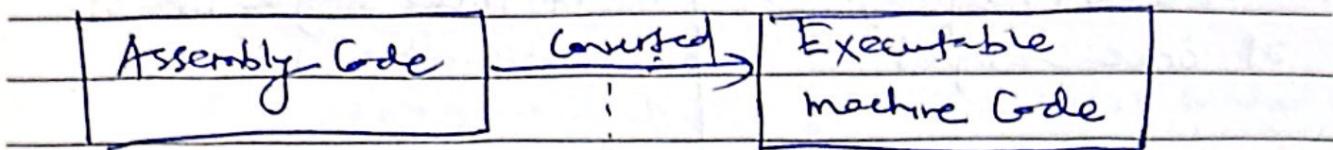
It is program.

for every program (Hardware or Operating System) we have assembler. An Assembler

SMTWTHFS

Date

for one platform will not work for another platform



Can be loaded
at any time &
can be run.

4. Loader & Linker.

Converts relocatable Code → Absolute ~~Machine~~
Machine Code

A Linker links different object files into a single file / executable file. and then loader loads that executable file in memory and executes it.

KHWOPA

Compiler

1. Takes entire program at once as input

2) Speed - High

3) Generates intermediate object code. So, memory requirement is more.

4) Errors:

All errors are displayed together

5. Error detection is difficult

6) Compilers are larger in size.

7) C, C++, Scala

Interpreter

1. It takes single line of code at a time.

2. Speed - low

3) Memory requirement is less because no intermediate code created.

4) Errors:

Continues translating the program until the first error is met, in which case it stops

5. Error detection is easy

6. Smaller in size

7. Perl, python, Matlab, Ruby

Phases of a Compiler

HLL / Source Code

Lexical Analyzer

Tokens

uses a tool called lex to generate these tokens.

Syntax (parser)

Analyser

Context Free Grammar

parse tree

Semantic Analyser

parse tree

(verified semantically)

Intermediate Code Generator

TAC (Three Address Code)

Code Optimiser

optimized code

Reduce number of TAC

backend

Target Code Generator

Assembly Code / Target program

KHWOPA

Symbol Table Manager

→ Data structure being used by compiler to store all the information related to identifiers.

e.g. its types, scope, size, location, name etc.

- It helps the compiler to function smoothly by finding the identifier quickly.

→ All the phases interacts with the Symbol table manager.

	None	Type	Scope	Size
a	int	global	2 Bytes	
b	float	local	4 Bytes	

Error Handler

It is a module which takes care of the errors which are encountered during compilation and it takes care to continue the compilation process even if error is encountered.

The task of error handling process are to detect each error, report it to the user and then some recovery strategy and implement them to handle errors.

Phases of Compiler

- | | | |
|---------------------------------|---|------------------------------|
| 1) Lexical Analysis | } | Analytic phase
(Frontend) |
| 2) Syntax Analysis | | |
| 3) Semantic Analysis | | |
| 4) Intermediate Code Generation | } | Synthetic phase
(Backend) |
| 5) Code optimizer | | |
| 6) Target Code Generator | | |

1. Lexical Analyzer.

- Reads the program and Converts characters into tokens using a tool called LEX tool.
- Tokens are defined by regular Expression which are understood by lexical Analyzer
- Removes white Spaces, Comments, tabs

E.g. int a = 10;
 Tokens: int, a, =, 10, ;

2) Syntax Analyzer / PARSER

- Constructs the parse tree.
- Takes the tokens one by one and uses Context free Grammar (CFG) to construct the parse tree. If it is not possible to

construct the parse tree then input is ~~syntactically~~
 Syntactically correct and error message shown
 displayed. Using productions we can represent
 what the program actually is. The input has
 to be checked whether it is in the desired
 format or not.

- Syntax errors can be detected by it if the given grammar is not accordingly.
- Consider input as tokens and checks their syntactic consistency.

e.g. $\text{if } (a > b) \{ \dots \}$ → Valid Syntax.

3. Semantic Analyzer.

- Verifies the parse tree, whether it is meaningful or not.
- It uses parse tree and information in the symbol table to check the source program for semantic consistency with the language definition.
- type checking, depicting of the variables used, local checking, flow control checking

e.g. int a = "Hello"; → Error (Type mismatch)

4. Intermediate Code Generation.

→ Generates the intermediate code

① Three Address Code (TAC)

or

Abstract Syntax tree (AST)

5. Code optimizer

→ We get optimized code like our number of lines are reduced removing unnecessary code lines.

→ Code optimization phase attempts to improve the intermediate code so that it runs fast and consumes less resources

6. Target Code Generator

- Final phase of the Compiler which generates the target code.

Sample Code & Compiler phases example

#include <stdio.h>

```
int main() {
    int a=5, b=10;
    int sum = a+b;
    printf ("sum=%d", sum);
    return 0;
}
```

→ Lexical Analysis

Breaks Codes to tokens

int, main, a, b, sum, printf, etc.

→ Syntax Analysis :

checks grammar.

e.g. Correct use of int and printf.

→ Semantic Analysis : Checks Variables are declared before use.

→ Intermediate Code Generation

Sum = a+b;

t1 = b

t2 = a+t1;

Sum = t2;

- Code Optimization : Eliminates redundant operations.

→ Code Generation.

MOV AX, b
Add AX, a
Mov Sum, AX.

I. Lexical Analyzer

Function of Lexical Analyzer

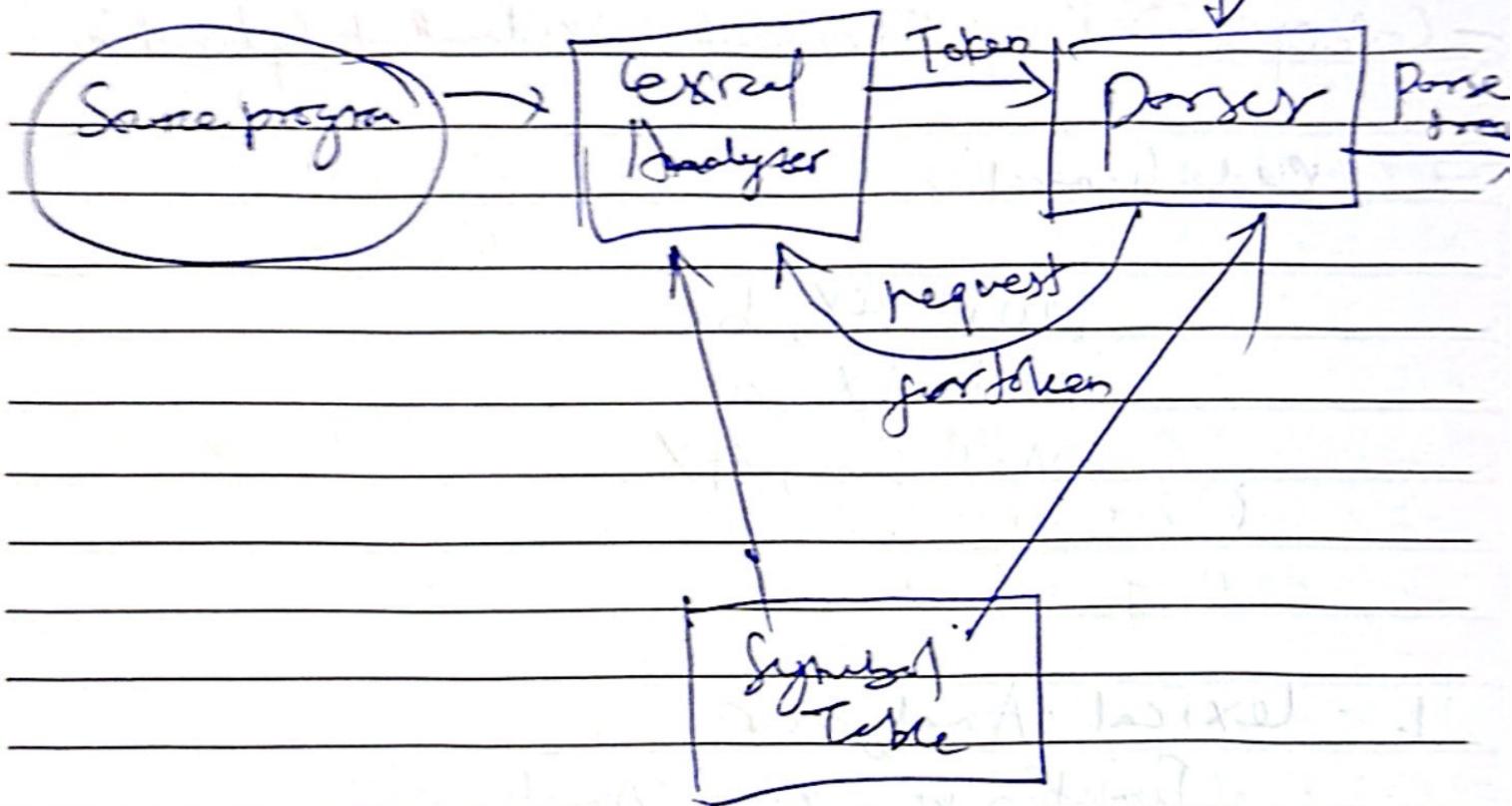
The main function of Lexical Analyzer.

In lexical analysis phase, Lexical Analyzer converts source programs to streams of tokens.

- It is also called Scanning -

Function of Lexical Analyzer

- i.). Reads the input program character by character and it produces a stream of token (and passes all the data to the Syntax Analyzer when it demands.



- ii) Removes whitespace and comments.
- iii) Identifies keywords, operators and identifiers.
- iv) Lexical Analyzer generates error and gives the line number of the error.

E.g.:

`int x = 10;`

Tokens

Keyword : int

Identifier : x

Operator : =

Constant : 10

Punctuation : ;

if we pass `a = b + c;`

Token $id = id + id;$ ($id = \text{identifier}$)

where each id refers to its variable in

Symbol Table.

TOKENS :

A token is a sequence of characters that can be treated as a unit/single logical entity.

Typically tokens are

- keywords (for, if, while, int....)
- identifier (variable, function)
- operator (+, -, ++)
- constant
- punctuation separators , ; etc.

PATTERN:

Pattern is a rule describing all those lexemes that can

LEXEME:

It is a sequence of characters in the source program that is matched by the pattern for a token.

or

a sequence of input characters that comprises a single token is called a lexeme.

e.g., "float", "-", "723", ";".

PATTERNS

Pattern is a rule describing all the lexemes that can represent a particular token in a source language.

Now we can define pattern as:

There are some predefined rules for every lexeme to be identified as a valid token. These rules are defined by grammar rules by means of patterns.

E.g. int max (int i);

→ Lexical Analyzer first read int and find it to be valid and accepts as tokens

→ max is read by it and found to be valid function name.

→ Semicolon (, , int, ;,) and i are all valid tokens total tokens 7.

eg. int main ()

{ // 2 variables declared below

int a, b;

a = 10;

return 0;

Token = 18

}

eg. printf ("Name Age") ; Token = 5;

lexical Analysis with DFA

lexical Analyzers are Deterministic Finite Automata (DFA) to recognize patterns in source code with a finite number of states and transitions between them based on input character.

A DFA recognizing identifiers

$([a-zA-Z] \ [a-zA-Z0-9])^*$ consisting

State q_0 (start state): Transition to q_1 when seeing alphabet or _.

State q_1 (accept state): Accept letters and digits but not special symbols.

processing of Comp 1.

$C \rightarrow q_1, \circ \rightarrow q_1, \ u \rightarrow q_1, \ h \rightarrow q_1, \ t \rightarrow q_1$
 $1 \rightarrow q_1$ (Void identifier)

int year = 5 error is detected because variable names cannot start with digit.



5 - Write in

Date _____

KHWOPA

Left Recursion .

A production of Grammar is said to have left recursion if the ~~leftmost~~ leftmost variable of its RHS is same as variable of its Left hand side.

$$\text{i.e. } A \rightarrow A\alpha \mid \beta$$

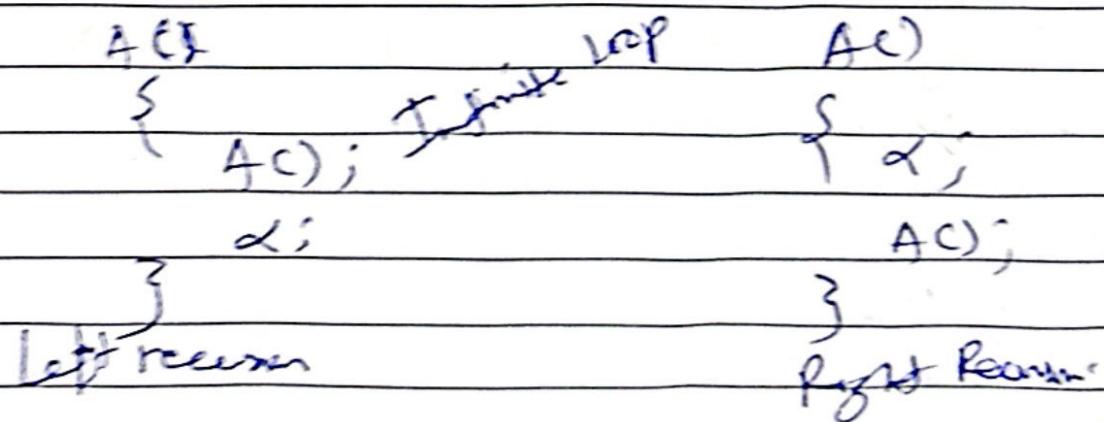
- (i) Direct left recursion $B \rightarrow B\alpha$
- (ii) Indirect left recursion $S \rightarrow A^q$
 $A \rightarrow S_m$

Why we remove Left Recursion .

Top Down Parser cannot accept the grammar having left recursion i.e. they don't allow left recursive Grammars.

→ so, we have to remove left recursion but preserve the language generated by grammar.

E.g. in term of ffn.

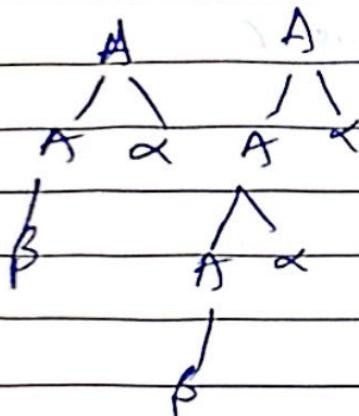


Recursion

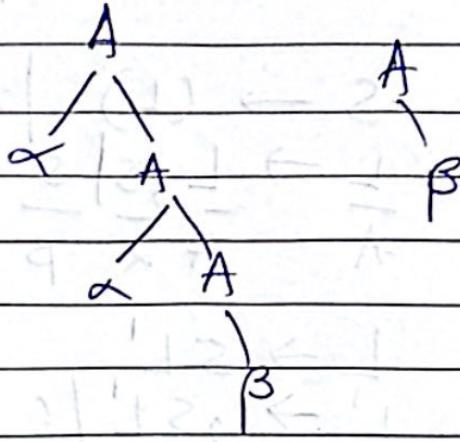
LR

RR

$$\underline{A} \rightarrow \underline{A} \alpha | \beta$$



$$A \rightarrow \alpha A | \beta$$



$$\text{yield}(I) = \beta \alpha^*$$

$$\text{Language} = \alpha^* \beta.$$

Conversion

$$\boxed{A \rightarrow \beta A' \quad - - \quad A' \rightarrow \alpha A' / \epsilon}$$

$\Rightarrow A$ will generate β followed by A'
 \Rightarrow now A' will generate α^*

e.g. $E \rightarrow E + T \quad / \quad T$
 $A \cdot \overline{A} \quad \alpha \quad \beta$

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

SMTWTHFS

Date

$$\textcircled{e_1} \quad \frac{S}{A} \rightarrow \frac{S}{A} \frac{OSIS}{\times} \frac{|0|}{\beta}$$

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' | \epsilon$$

$$S \rightarrow 0 | S'$$

$$S' \rightarrow OSIS | S' | \epsilon$$

$$\textcircled{e_2} \quad S \rightarrow (L) | x \text{ no left recr}$$

$$\frac{L}{A} \rightarrow \frac{L}{A} \frac{,S}{\times} \frac{|S}{\beta}$$

$$L \rightarrow SL'$$

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

$$\textcircled{e_3} \quad A \rightarrow A\alpha_1 | A\alpha_2 | A\alpha_3 | \dots | \beta_1 | \beta_2 | \beta_3$$

$$A \rightarrow \beta_1 A' | \beta_2 A' | \beta_3 A' | \dots$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \alpha_3 A' | \dots | \epsilon$$

Sunny

①

Grammar

Ambiguous

2 parse tree
(no parser work for it)

②

Grammar

L_1

because Top Down

parser will

not accept it

because will fall

'infinite loop -

KHWOPA

S M T W T H F S

Date

(iv) $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow id$

$$\begin{array}{l} A \rightarrow \underline{\alpha} \mid \beta \\ A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}$$

1st production

$E \rightarrow TE'$

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

Solving for 2nd production
 $T \rightarrow T * F \mid F$

$T \rightarrow FT'$

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

3rd production $F \rightarrow id$ α Final Answer ~~$E \rightarrow TE' \mid F \rightarrow id \mid \alpha$~~

$$\begin{array}{l} E \rightarrow TE' \\ E' \rightarrow *TE' \mid \epsilon \\ T \rightarrow FT' \\ T' \rightarrow *FT' \mid \epsilon \\ F \rightarrow id \end{array}$$

(v) $S \rightarrow (L) \mid a$ $L \rightarrow SL'$
 $L \rightarrow L, S \mid S$ $L \rightarrow S L' \mid \epsilon$

(vi) $S \rightarrow Aa$ $A \rightarrow Aabc \mid C$ $A \rightarrow CA'$
 $A \rightarrow SB \mid C$ $A' \rightarrow abA' \mid \epsilon$

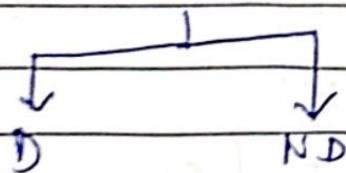
KHWOPA

LEFT FACTORING

Sometimes it is not clear which production to choose to expand a non-terminal because multiple production begin with the same terminal / non-terminal block.

This type of Grammar \rightarrow Non Deterministic Grammar or
Grammar Containing left factoring

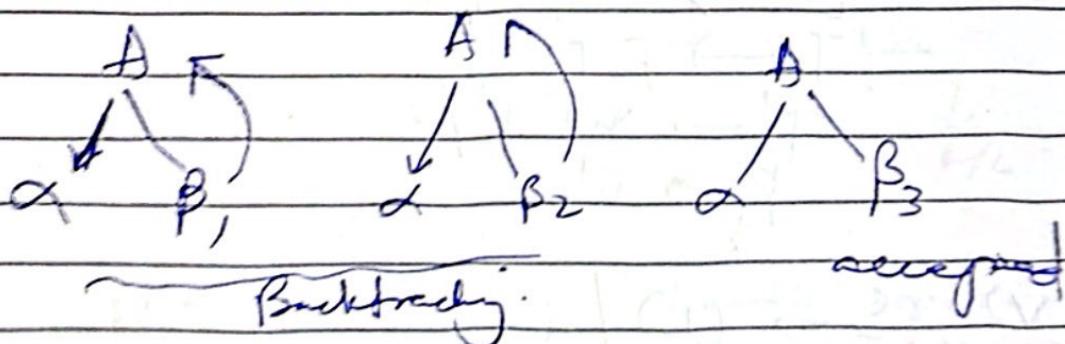
Grammar



$$A \rightarrow \underline{\alpha} \beta_1 | \underline{\alpha} \beta_2 | \underline{\alpha} \beta_3$$

common prefixes

Suppose we have to accept $\alpha\beta_3$



This is happening because of the common prefixes or (Context)

KHWOPA

One or more productions on the RHS are having something common in the prefixes (Context)

This is called Common prefix problem or NDT deterministic grammar.

Backtracking happens?

→ we are making our decision only after seeing prefix. (as in example)

i.e which production to choose without seeing full input.

If input is $\alpha \beta_3$

~~non Deterministic~~

$A \rightarrow \alpha \beta_1 / \alpha \beta_2 / \alpha \beta_3$ (NDG)

$A \rightarrow \alpha A'$

$A' \rightarrow \beta_1 / \beta_2 / \beta_3$ } we have postponed the decision making.

The procedure we used to convert non-Deterministic grammar \rightarrow Deterministic Grammar is called left factoring.

e.g. $S \rightarrow iEtses / iEts | a$
 $E \rightarrow b$

$S \rightarrow iEts S'$
 $\Rightarrow S' \rightarrow es | e | a$
 $E \rightarrow b$

$S \rightarrow aSSbs /$
 $aSaSb | abb / b$

$S \rightarrow aS' / b$
 $S' \rightarrow ssbs / Sasb / bsb$
 $S' \rightarrow SS'' / bb$

$S'' \rightarrow sbs / asb$

KHWOPA

S M T W T H F S

Date

$S \rightarrow bSS \cup aS \cup b \cup \underline{SS} \cup a$

$\cancel{\{ } S \rightarrow bSS' \cup a$

$\rightarrow S' \rightarrow Sas \cup sas \cup b$

$\cancel{\{ } S' \rightarrow Sas'' \cup b$

$\cancel{\{ } S'' \rightarrow as \cup sb$

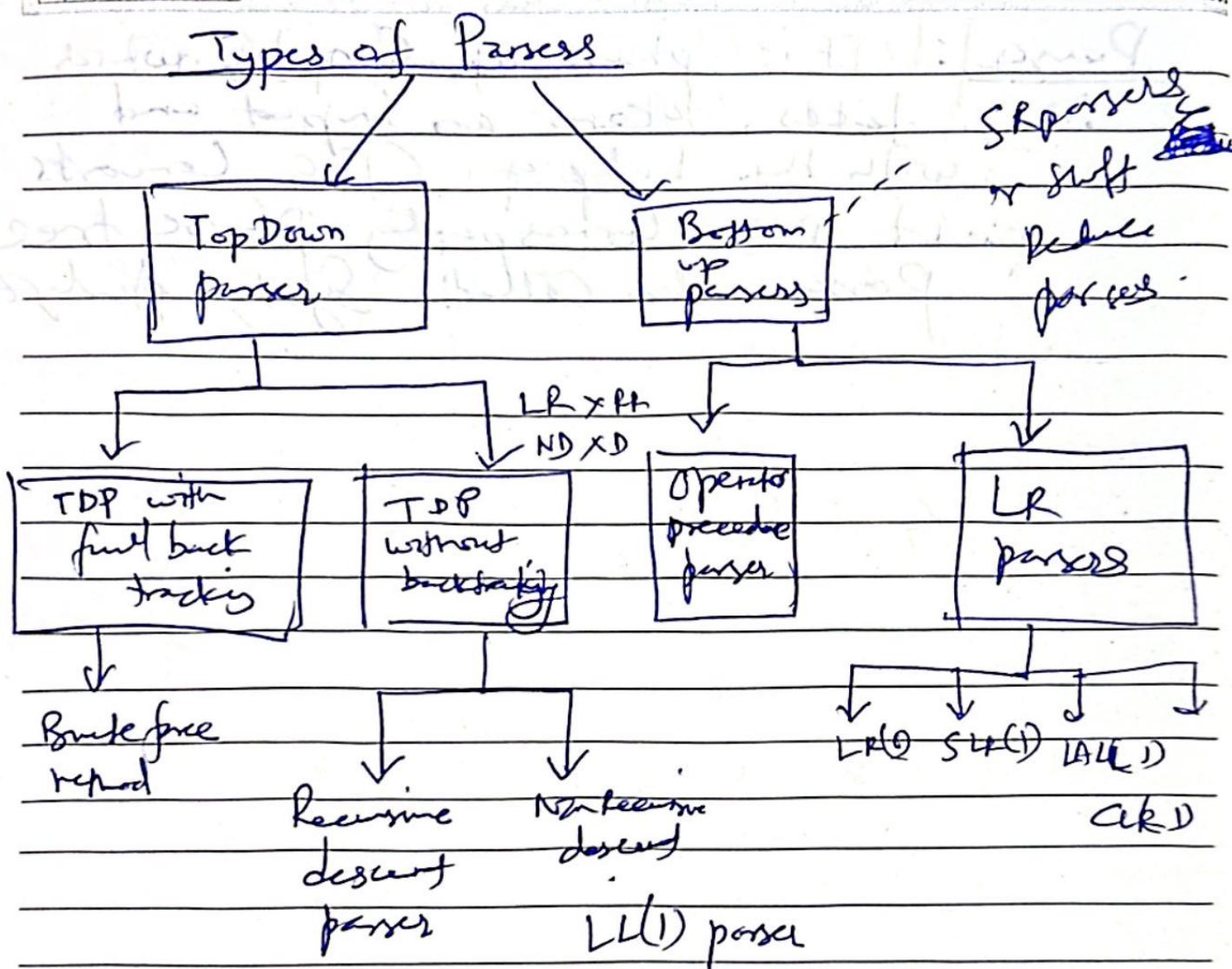
Parser

SMTWTHFS

Date

Parser: It is phase of Compiler which takes tokens as input and with the help of CFG, converts it into Corresponding parse tree.
Parser also called Syntax Analyzer.

KHWOPA



we will have some
recursion

LL(1) parser

LR parser

L - left-right (using of symbol)

L - Left most derivative

L - left-right (can)

R → R MD using
reverse of RMD.

abcd

→

KHWOPA

SMTWTHFS

Date

	TD parser	Bottom up parser
1) Variety of Grammars	X	
→ Ambiguous Grammars		X except operator precedence parser
→ Unambiguous	✓	✓
3) LR	X	✓
4) RR	✓	✓
5) ND	X	✓
6) Deterministic	✓	✓

KHWOPA

Top Down Parser

It generates parse tree for the given input String with the help of grammar production by expanding the non-terminals i.e.

it starts from the Start Symbol and ends on terminals.

- It uses LMD (Left Most Derivation)

Bottomup Parser

It generates the parse tree for the given input String with the help of grammar production by compressing the non-terminals i.e. it starts from the terminals and ends on the Start Symbol.

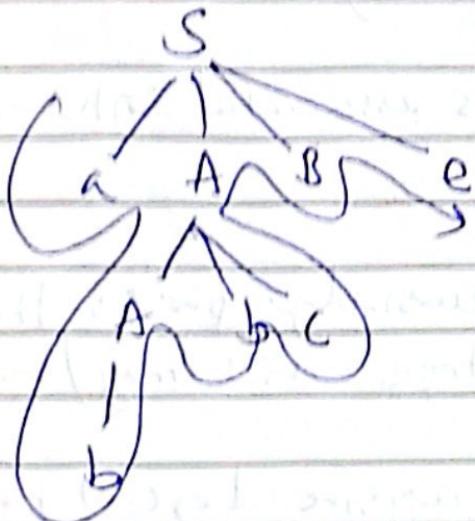
- It uses reverse of Right Most Derivation

$$S \rightarrow aABe$$

$$A \rightarrow Abc/b$$

$$B \rightarrow d$$

$$\omega \rightarrow abb_{}cde$$

TDP:LMD:

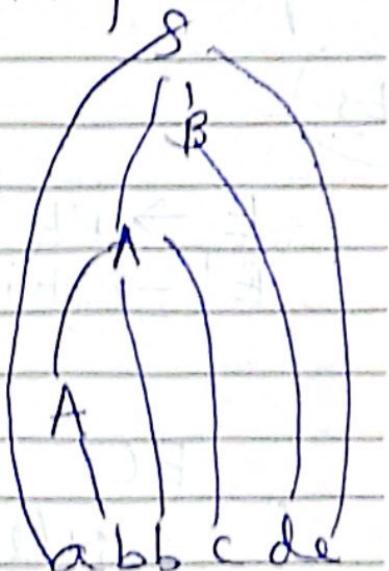
$$\begin{aligned}
 S &\rightarrow aABe \\
 &\rightarrow aAbcBe \\
 &\rightarrow abbcBe \\
 &\rightarrow abbcd_
 \end{aligned}$$

Top down parsing \rightarrow At every point we have to decide what is the next production we should use.

BUP: Bottom Up Parser

main decision we have to make is

"when to reduce the grammar"



$$RMD: \quad S \rightarrow aA_{}Be$$

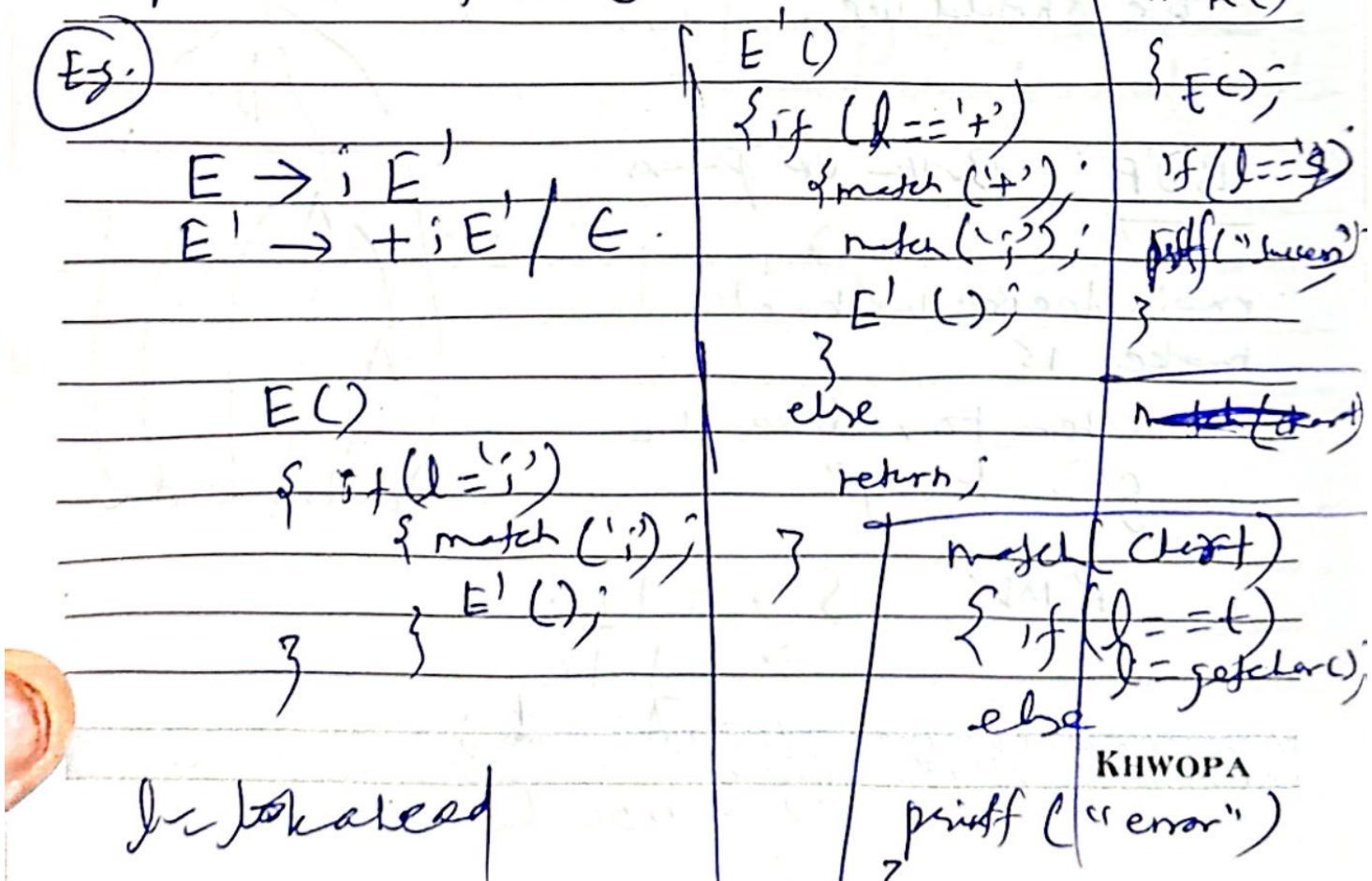
$$\rightarrow aA_{}de$$

$$\rightarrow a_{}A_{}bcd_{}e$$

$$\rightarrow abbcd_{}e$$

Recursive Descent Parser

- It is top-down parser technique that constructs the parse tree from top and the input is read from left-right.
- A procedure is associated with each non-terminal of the grammar.
- This technique recursively parses the input to make a parse tree which may/may not require backtracking.
- A form of recursive descent parser that does not require any backtracking is called predictive parsing.



Q

SMT WTHS

1 + 9 + 1 end of string.

Date

E

1 + 9 + 1 end of string.

Kiwopa

Operator Grammer

A grammar that is used to define mathematical operators is called an operator grammar or operator precedence grammar.

A grammar is said to be operator precedence grammar if it has two properties.

- (1) No R.H.S of any production is C.
- (2) No two non terminals are adjacent in P.R.H.S.

e.g. $E \rightarrow E+E \mid E * E \mid id. / operator$
gramm-

e.g. $S \rightarrow SAs \mid a$ } Not operator
 $A \rightarrow bS \mid b$ } grammar

we can say it is operator grammar

$S \rightarrow SbSbS \mid a \mid SbS'$

$A \rightarrow bSb \mid b$ not regular.



There are 3 operator precedence relations:

- 1) $a \rightarrow b$ means a has higher precedence than the term b. or "a takes precedence over b".
- 2) $a \leftarrow b$ means a 'yields precedence to' b
or a has lower precedence than term b.
- 3) $a \doteq b$ means a has same precedence as b.

Operator Precedence Parser

- Bottom up parser that interprets an operator grammar.
- This parser is only used for operator grammar.
- Ambiguous grammar are not allowed in any parser except operator precedence parser.