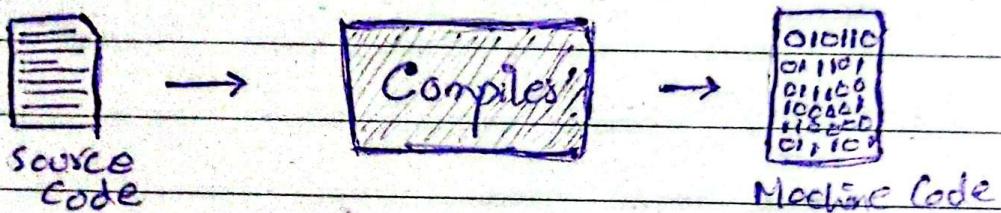


Date \_\_\_\_\_

Day \_\_\_\_\_

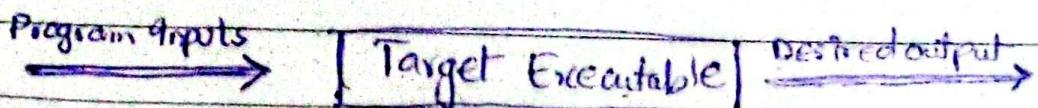
## WEEK # 01

- Programming languages — Notations for describing computations to people and to machines.
- A language understood by a machine is called Machine language (ML), Machine Code, Object code
  - ↳ Consist of streams of 0's & 1's
- The program that carries out the translation activity is called compiler



\* A compiler is a program that translates source code written in one language into the target code of another language \*

- Executables — Files that actually do something by carrying out a set of instructions. Eg: .exe file



↳ Initially there is a compile phase followed by the run/execute phase

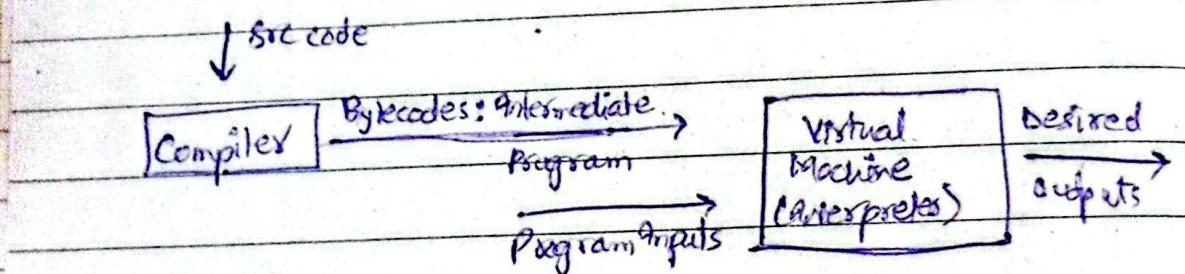
## • Interpreter

↳ A program that doesn't produce the target executable

→ It reads the code line by line & directly executes the code by using given inputs



## ▼ Java Combines Compilation & Interpretation



**Portability:** Bytecode compiled on one machine can be interpreted on another one

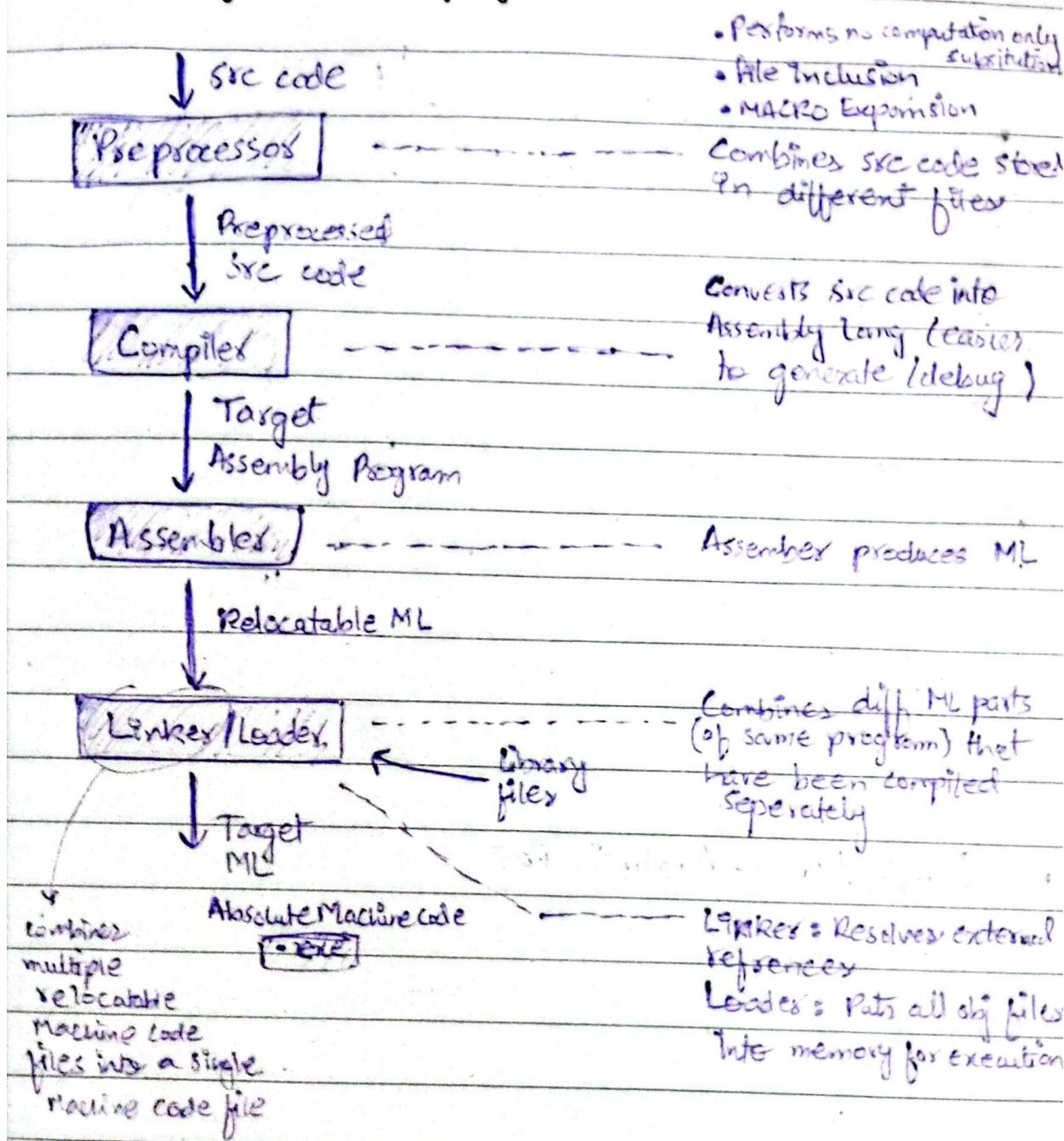
**Faster Execution thru Just-In-Time Compilers (JIT)**

- In compiler based language compiler & source code not needed to run .exe whereas interpreter requires both
- in compiler if multiple error it shows them all whereas in interpreter it only shows 1st one
- Compiles low level language → Fast
- Interpreter takes less space
- In Interpreter debugging easier

Date \_\_\_\_\_

Day \_\_\_\_\_

## • Language Processing System:

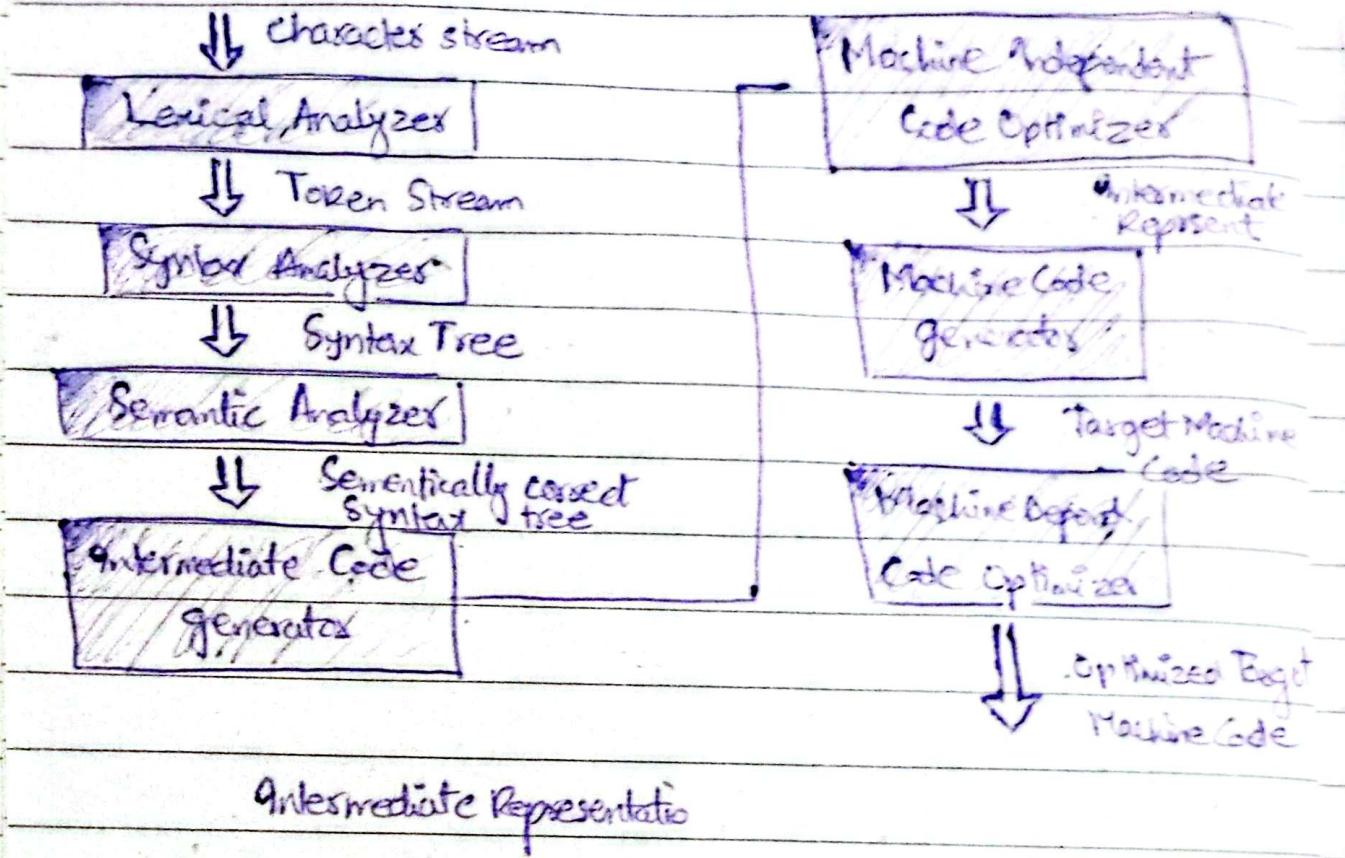


These all phase don't run sequentially

Date \_\_\_\_\_

Day \_\_\_\_\_

## ► Structure of A Compiler



Intermediate Representation

## ► Compilation - Analysis Part

→ Phases Involved

- (1) Lexical Analysis
- (2) Syntax Analysis
- (3) Semantic Analysis

↳ Determines the operations implied by the source program which are recorded in a tree structure called Syntax Tree.

↳ Breaks up src code into constituent pieces while storing info in symbol table

Date \_\_\_\_\_

Day \_\_\_\_\_

## ▷ Compilation — Synthesis Part

→ Phases Involved

### ① Intermediate Code Generation

- ② " " Optimization      ↗ Both intermediate & machine code are optimized to achieve efficiency
- ③ Machine Code Generator
- ④ " " Optimizer

↳ Constructs target code from Syntax Tree & from info in Symbol Table

## ▷ Frontend & Backend

- Lexical Analysis
  - Syntax "
  - Semantic //
  - Intermediate Code Generation
  - " " Optimization }
  - Machine Code generation }
  - " " Optimization }
- } Frontend  
                  (Machine Independent)
- } Backend  
                  (Machine Dependent)

### ① Lexical Analysis

Let  $x = a + b * c \Rightarrow$  Lexical Analysis  
generates Lexems  $\Leftarrow$   
& tokens

→ Lexems are similar to words but they individually don't have any meaning

Date \_\_\_\_\_

Day \_\_\_\_\_

→ Lexical Analysis Recognizer

the tokens using regex

Eg: RegEx for identifier

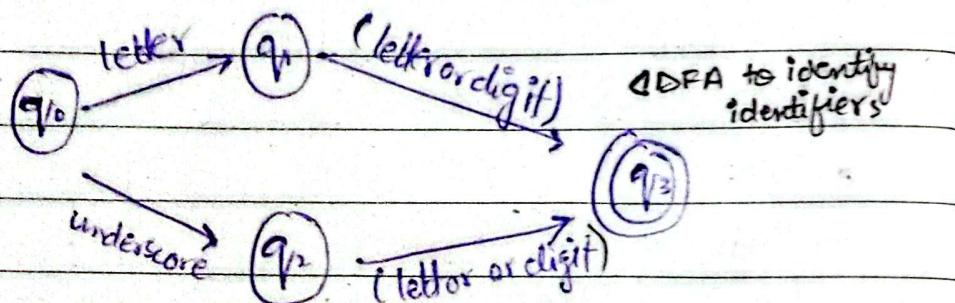
$$l(l+d)^*|_-(l+d)^*$$

l = letter

d = digit

\_ = underscore

Lexems	Token
x	identifier
=	operator
a	identifier
+	operator
b	identifier
*	operator
c	identifier



→ It reads the character stream of source code and break it up into sequences call lexems.  
(A.K.A Scanning)

→ Each lexem is represented as token

## ② Syntax Analysis

→ Also known as Parser

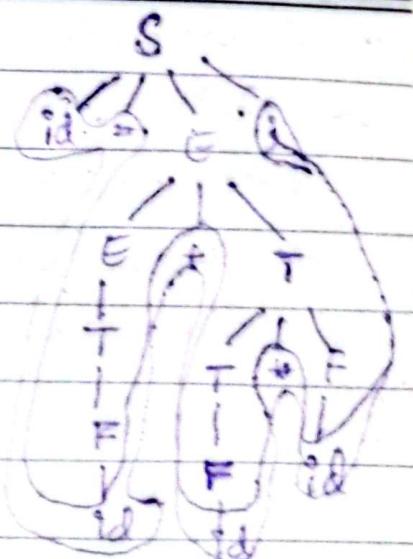
→ The stream of token is pass to this phase to identify grammatical structure of the stream.

→ Syntax analyzer depends on Type 2 or Content Free Grammars

Date \_\_\_\_\_

Day \_\_\_\_\_

- It builds a parse tree, which replaces the linear sequence of tokens with a tree structure
- The tree is built according to the rules of grammar.
- It is analyzed, augmented & transformed by later phases



### (3) Semantic Analyzer

- Takes the parse tree as input
- Performs semantic checks e.g:
  - Type checking
  - object binding
  - rejecting incorrect programs
  - issue warning
  - correctness of scope resolution
  - multiple variables declaration
- performs logical analysis of parse tree

undeclared variable  
TC

### (4) Intermediate Code Generator

- Takes the semantically verified parse tree
- Three Address Code is used [3 operands/instruction]
  - $t_0 = b * c$  → bottom to top & chain to precedence
  - $t_1 = a + t_0$  place chain w.r.t. parse tree
  - $x = t_1$

Date \_\_\_\_\_

Day \_\_\_\_\_

### (5) Code Optimizer

- Takes the intermediate code as input and optimizes it.
- $t_0 = b * c ; \quad \} \text{ optimized to 2 instructions}$   
 $x = a + t_0 ; \quad \|$

### (6) Target Code Generator

- Converts the optimized intermediate code to target machine code.
- This is done thru assembly language

### (7) Machine Dependent Optimizer

- Optimizes the machine code
- It is optional activity

### (8) Error Handler

- It is a subroutine to take care of the continuation of compilation, even any error at any phase i.e error handler is responsible to continue the compilation process even any error occurs at any phase

Date \_\_\_\_\_

Day \_\_\_\_\_

## WEEK #02

### ▷ SYMBOL TABLE MANAGEMENT

- Data Structure used for storing the names of variables and their associated attributes
  - object
  - function
  - classes
- It should be designed so that the compiler can
  - ↳ find record for each variable quickly
  - ↳ Store / retrieve data for a record quickly
- Symbol table is ~~used~~ by Lexical analysis, Syntax is made & updated by

Semantic Analysis. Then used by last three phases

Lexical Analysis → Creates entries for identifiers

Syntax Analysis → Adds info regarding attributes.

Semantic Analysis → Using available info checker

Semantics & updates the symbol table

Intermediate → Available info helps in adding temporary variables information

Code Optimization → Info is used in machine dependent optimization

Target Code Gen → Generates the code using address info of identifiers

Date \_\_\_\_\_

Day \_\_\_\_\_

## → Attributes of table.

- ① Name → name of var
- ② Type → type [int, float...]
- ③ Size → How many bytes
- ④ Line of Declaration → n
- ⑤ Line of Usage → if used in multiple place this will be a linklist
- ⑥ Address → Address info of identifier

## • Non-block Structured Language (FORTRAN...)

→ Contains single instance of the variable declaration

### → Operations

- i - Insert()
- ii - Lookup()

## • Block Structured Language (C, C++, JAVA)

→ Variable declaration may happen multiple times

### → Operations

- i - Insert()
- ii - Lookup()
- iii - Set()
- iv - Reset()

defining & redefining  
scope

Date \_\_\_\_\_

Day \_\_\_\_\_

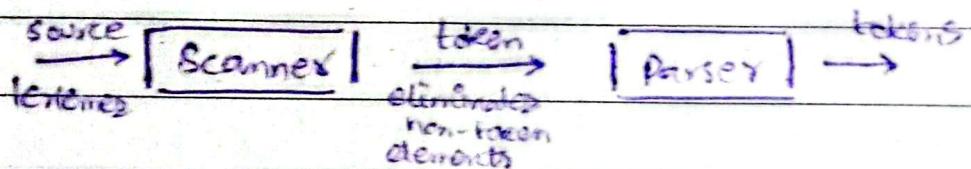
## D Parsing & Scanning

→ In compilers the recognizer is split in to two phases

① Scanner → Translate input characters to tokens

↳ Also, report lexical errors like illegal characters and illegal symbols

② Parser → Read token stream and reconstructs the derivation (parse tree)



→ These two modules are separate to maintain simplicity & separation of concerns

↳ Scanner hides details from parser (comments, whitespace, input files etc)

↳ Parser therefore has simpler stream of tokens and also due to efficiency

↳ Scanner uses simple faster design but still consumes a surprising amount of the compiler's total execution time.

Date \_\_\_\_\_

Day \_\_\_\_\_

## ▷ Tokens

- ↳ It is like a class which is assigned to each lexem (character of src code)
- ↳ Some tokens have attribute associated with it

### • Operators & Punctuations

↳ + , / , \* , ; , ( , ) , { , } , < , = ---

↳ Each of these has separate lexical class (token name)

### • Keywords

↳ while , do , int , float , switch , if , else ...

↳ Each of these also have separate class

### • Identifiers

↳ All identifiers have class ID , with a value that represents its "id" in symbol table

### • Integer constants

↳ They have INT class , with the value of that constant . INT (42)

return maybe b = ifby ;

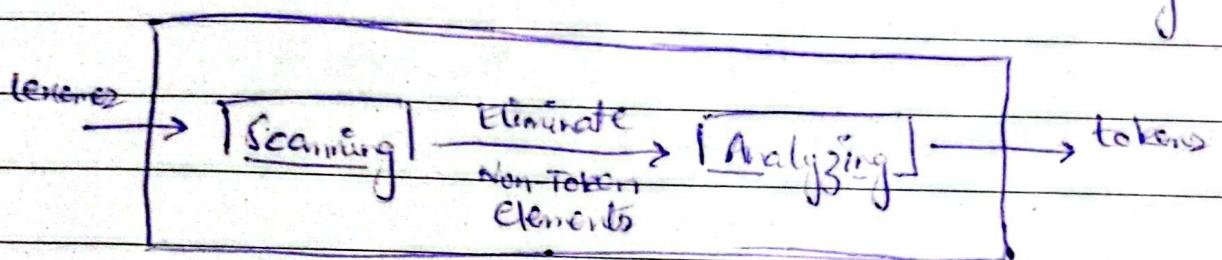
[RETURN] [ID(maybe)] [INEQ] [ID(ifby)] [colon]

Date \_\_\_\_\_

Day \_\_\_\_\_

## • Lexical Analyzer

- ↳ Scans the pure High Level Language (HLL) code line by line
- ↳ Takes lexemes as input & produces tokens or o/p
  - ↳ Uses DFA for pattern matching



C-tokens:

if :  $\rightarrow (A) \xrightarrow{?} (B) \xrightarrow{+} (C)$

identifier :  $\rightarrow (D) \xrightarrow{(a-z/A-Z/-)} (E) \xrightarrow{?} (a-z/A-Z/0-9)$

integer :  $\rightarrow (F) \xrightarrow{(+\mid=)} (G) \xrightarrow{(0-9)} (H)$

↳ Removes comments & whitespaces from pure HLL code

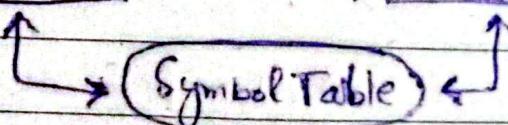
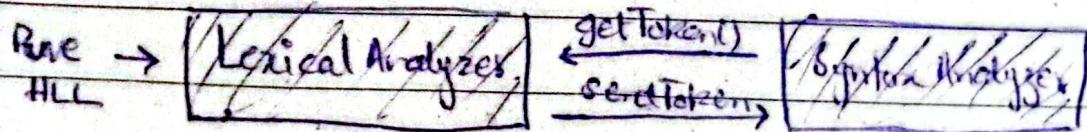
↳ ' ' space

'\t' horizontal tab

'\n' newline

'\v'

'\r' carriage return



Date \_\_\_\_\_

Day \_\_\_\_\_

→ Lexical Analyzer helps in giving error message by providing row no & col no.

→ Sabse phle Syntax Analyzer chalta ha wo get-token() call karta hai or phir LA send-token()  
→ Ags Syntax Analyzer ko esa token recv ha which was not expected by the parse tree so it reports to error handler.

→ But it doesn't stop!

→ LA uses DFA to do tokenization.

→ While doing tokenization always gives importance to longest matching

main()

{

    a a a b \* c }

    int a; b; c;

    y z x b ) ;

}

3

12

21

27

(23)

This program has  
syntax & semantic

error but lexical  
analyzer will not

stop = there is no  
lexical err

char \* f() { "string" } ;

6

in f( ) ; 200 ;

6

in /\* comment \*/ f( ) ; 200 ;

↳ comment hata kar space ajayega

Date \_\_\_\_\_

Day \_\_\_\_\_

→ ~~int i = 9; { / \* abc \* \* \* / ; abc } ; } 10 ;~~

→ ~~int i = 9; p; } 6~~

→ ~~int exp. + + + y; } 9~~

→ ~~char ch) G ("Hello ; }~~

↳ Token error (strings not closed)

### • Lexical Errors

- ① Identifiers that are way too long
- ② Exceeding length of numeric constant
- ③ Numeric constants which are ill-formed  
    int i = 457~~8~~92;
- ④ Illegal characters that are absent from the source code  
    char x[] = "Hello"; ↳ this is not any keyword or an identifier

### • Lexical Error Recovery

↳ Panic-mode recovery

↳ If we encounter error then skip the successive characters until delimiter is found

↳ Transpose of two adjacent characters

↳ while ( — ) ⇒ while , ill shifted

↳ Insert a missing character

↳ it sam; ⇒ int . sam;

↳ Delete an unknown / extra character

↳ intt a; ⇒ int a;

↳ Replace a character  
    itt a; ⇒ int a;

Date

# CONTEXT FREE GRAMMAR Day

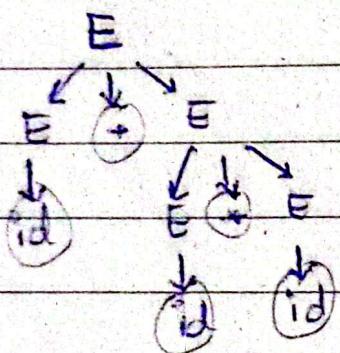
## ▷ AMBIGUOUS GRAMMAR

→ Ek grammar k lie koi es string exists koi he jiske lie ek se zyada parse tree bana sakte he to it is ambiguous.

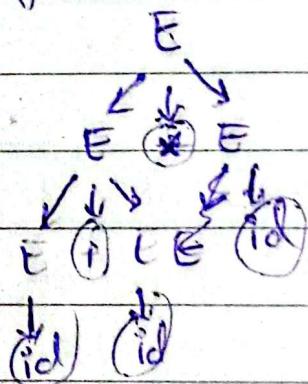
$$E \rightarrow E+E \mid E*E \mid id$$

$$W = id + id * id$$

Left Most Derivation



Right Most Derivation



## ▷ Converting Ambiguous To Unambiguous Grammar

### ① Left Recursive Grammar

$$E(S) \rightarrow S(bS/a)$$

right hand expression mein agr non terminal extreme left par ho.

### ② Right Recursive

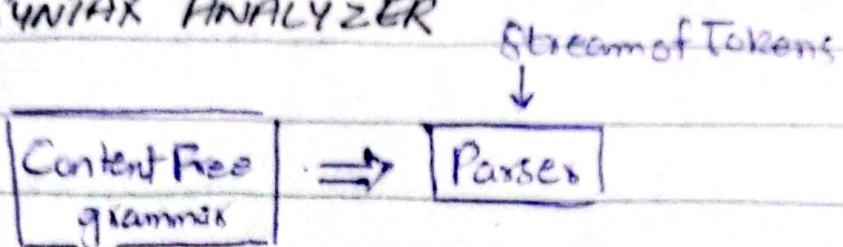
$$(S) \rightarrow a(S)/b$$

extreme right par non-terminal

Date

Day

## ► SYNTAX ANALYZER



- ## • Problems Due to Ambiguity

$$E \rightarrow E + E \quad (E \Rightarrow E) \text{ id}$$

$$W = \frac{ideid}{3+4 \times 5}$$

$3 + (4 \times 5) \rightarrow$  Precedence property violation

$$\begin{array}{c}
 E \rightarrow 35 \\
 E \times E = X \\
 \begin{array}{r}
 \overset{(+)E}{\cancel{E}} \quad \overset{id}{\cancel{E}} \\
 \overset{id}{\cancel{E}} \quad \overset{id}{\cancel{E}} \quad 5 \\
 3 \quad 4
 \end{array}
 \end{array}$$

$$\begin{array}{c} E \rightarrow 23 \checkmark \\ \swarrow \downarrow \searrow \\ E + E \\ \hline \begin{matrix} 1 & 1 \\ 3 & x \\ \hline 1 & 1 \end{matrix} \end{array}$$

$$W = 4 + 8 + 2$$

$$\underline{(4+8)+2}$$

$\Rightarrow$  Associativity  
property  
Violation

$$E \rightarrow 14 \checkmark$$

$E \rightarrow B$  X

$E \downarrow E$  and is  $E +$

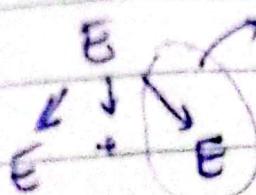
$$F + E$$

$\begin{matrix} \swarrow & \downarrow & \searrow \\ E & + & E \end{matrix}$  Same but we follow left to right associativity with plus.

Date

- Solution for Associativity

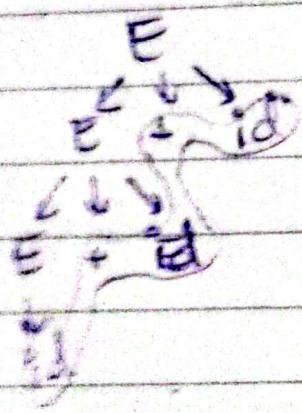
$$E \rightarrow E + E \mid id$$



this is causing the issue. Since  $+$  is L $\rightarrow$ R associative we will replace the right most terminal so the tree does not grow in that direction

$$E \rightarrow E + id \mid id$$

- We converted the grammar to left recursive to ensure L $\rightarrow$ R associativity



- For R $\rightarrow$ L associativity we convert to Right Recursive

$$E \rightarrow id + E \mid id$$

- Solution for Precedence

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

we broke this down to levels to remove violations  
ie LR associative we replace the right non-terminal

$$E \rightarrow E + T$$

$T \rightarrow T * T \rightarrow$  But  $*$  is also LR associative so we change it to  $* F$

$$E \rightarrow E + T \mid T \rightarrow \text{added this if are only event multipli cali}$$

$$T \rightarrow T * F$$

$$F \rightarrow id$$

left to precedence yields how we lower level got more priority

Date \_\_\_\_\_

low precedence  
will be on top

Day \_\_\_\_\_

Q:  $bExp \rightarrow bExp \text{ AND } bExp \mid bExp \text{ OR } bExp \mid \text{NOT } bExp \mid \text{True/F}$

$bExp \rightarrow bExp \text{ OR } bExp_1 \mid bExp_1$

$bExp_1 \rightarrow bExp_1 \text{ AND } bExp_2 \mid bExp_2$

$bExp_2 \rightarrow \text{Not } bExp_2 \mid \text{True} \mid \text{False}$

Q:  $R \rightarrow R + R \mid R \cdot R \mid R^* \mid a \mid b \mid c$

Precedence  $\rightarrow * , \circ , +$  (High  $\rightarrow$  Low)

Associativity  $\rightarrow$  Left  $\rightarrow$  Right

$R \rightarrow R + R_1 \mid R_1$

$R_1 \rightarrow R_1 \cdot R_2 \mid R_2$

$R_2 \rightarrow R^* R_2 \mid F$

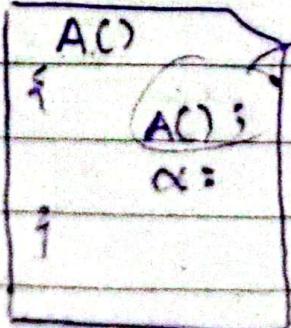
$F \rightarrow a \mid b \mid c$

## ► Problem with Left Recursion

Left Recursive

↳  $A \rightarrow A\alpha \mid \beta$

↳ If we consider this in form of code



it goes into infinite recursion.

" $\alpha$ " this statement never gets executed!

↳ But this is only a problem in Top Down Parser

Date \_\_\_\_\_

Day \_\_\_\_\_

therefore we need to convert left recursive to right recursive without changing the language

$$A \rightarrow \beta A'$$

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

here  
we can  
add termination  
condition.

