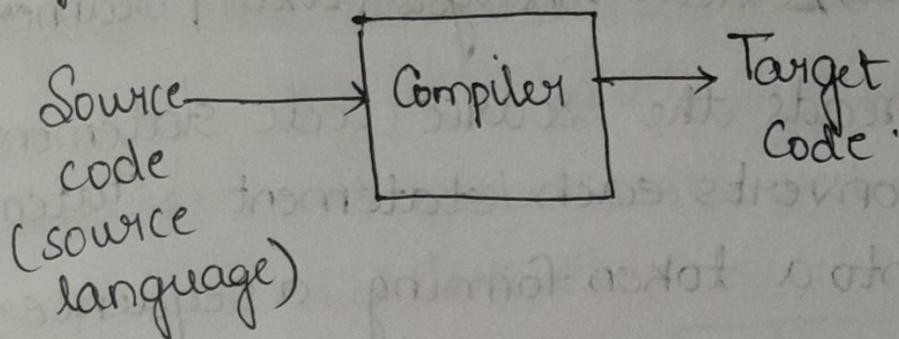
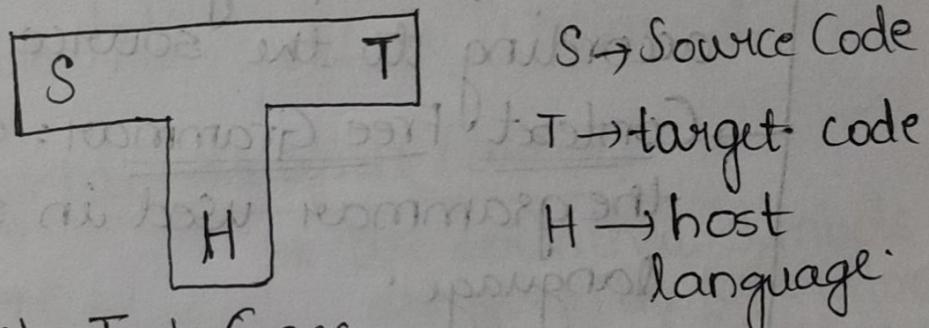


05/02/2024

- Compiler converts high-level language to low language.



- High-level language: human readable.
- Low-level language: machine understandable.
- Opcode: combination of 0's & 1's.
- Represented as:



- Compilers (C): TurboC, gcc
- Host language for C compiler is C language.
- Using the same host language as the programming language for a compiler is called Bootstrapping.
- Most of compilers are implemented using C.
- Pass: group of phases.
- Phase: consists of a specified functionality which takes a I/P and gives a O/P.

→ O/P of final phase is the target code.

* Basic phases of a compiler:

Phase-1 → Lexical Analysis (Lexer / Scanner)

- reads the source code statement
- converts each statement into a token forming a sequence.
- removes all the extra spaces and new line characters.
- removes the comments.

token:
 $\langle \text{name}, \text{loc} \rangle$

Phase-2: → Syntax Analysis (Parser) mandatory optional

- checks if the program is correct according to the source language.
- Context Free Grammar: specifies the grammar used in source language.
- O/P of this phase is a Syntax or parse tree.

Phase-1: Lexical Analysis

- If it is a reserved word then token consists only $\langle \text{name} \rangle$.
- If it is not a reserved word then token consists. $\langle \text{name}, \text{loc} \rangle$

Eg: variables / identifiers.

Eg: int a;

$\langle \text{int} \rangle \quad \langle \text{id}, 1 \rangle \quad \langle ; \rangle$

→ symbol table is the data structure where we store the pre-defined symbols.

* We store all symbols, reserved words, special characters etc.

Name	Type	descr
int	keyword	integer

→ If the symbol is not present in the symbol table, we create a token and update the symbol table. No repetition is allowed.

Phase-2: Syntax Analysis (Parser)

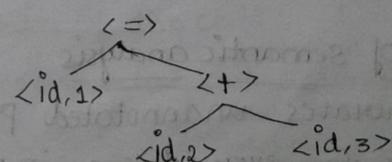
Ex: $a = b + c$ * all terminal symbols are given by operators.

= Phase-1: $\langle \text{id}, 1 \rangle \Rightarrow \langle \text{id}, 2 \rangle + \langle \text{id}, 3 \rangle ; \langle ; \rangle$

⇒ We draw the parse tree to these tokens.

⇒ Types: bottom-up; top-down

⇒ Leaf nodes are variables and operators would be the internal nodes.



Phase-3: Semantic Analysis:

- Meaning of the statement.
- Checks if the given tokens generate meaningful statements.
- for this it checks the type of identifier and verifies if further operations are meaningful (or) not.
- performs type conversion for the required types to be modified.

Interview questions strongly typed languages: The type of the variable is explicitly declared.

Eg: C, C++, Java.

Weakly typed languages: Type of the variable is known while usage in the code.

Eg: Python.

Explicit type conversion: C

Output of semantic analysis

generates unannotated parse tree.
i.e. for every token in parse tree

we add a "type" tag.

Phase-4: Intermediate Code Generation

↳ Eg: source code ; target code
 \Rightarrow high-level language \Rightarrow low level language

↳ intermediate code \Rightarrow 3-address code

$c = a + b * d;$
 \Rightarrow 3-address code:

$$\begin{aligned} t_1 &= b * d \\ t_2 &= a + t_1 \\ c &= t_2 \end{aligned}$$

LD R0, b

LD R1, d

MUL R1, R0

LD R2, a

ADD R1, R2

ST C, R2

to the right side of the assignment operator, we have atmost 1 operator and atmost 2 operands.

3-address code helps in generating the assembly language code easily.

Phase-5: Code Generation : Target Code

Phase-6: Code Optimization (optional)

↳ checks to reduce the size or time of execution of the code.

↳ Machine dependent code optimization.

Phase - 1, 2, 3 : Analysis Phases.

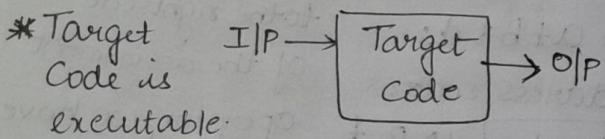
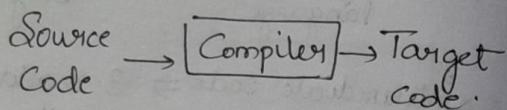
Phase - 4, 5: Synthetic Phases.

10/02/2024

UNIT - 1

* Language Processors:

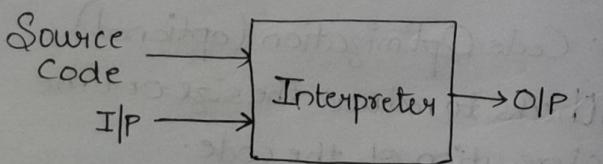
→ Compiler:



→ Compiler is machine dependent; we must execute the executable code on the same machine.

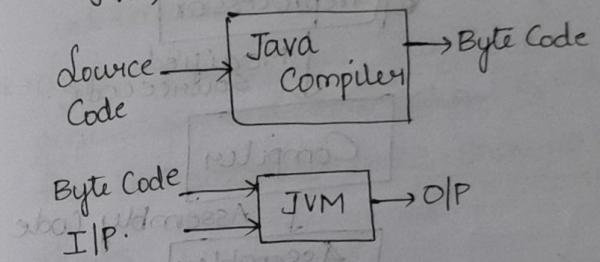
→ Interpreter:

* Converts each statement of the code one at a time, to target code.



* Compiler is faster than interpreter due to source code being

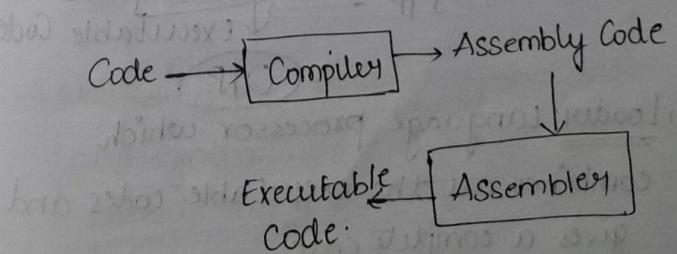
converted to target code.
→ Portability: The target code / byte code which is executable on any other machine; is called language portability
Ex: Java (runs on any system having JVM).



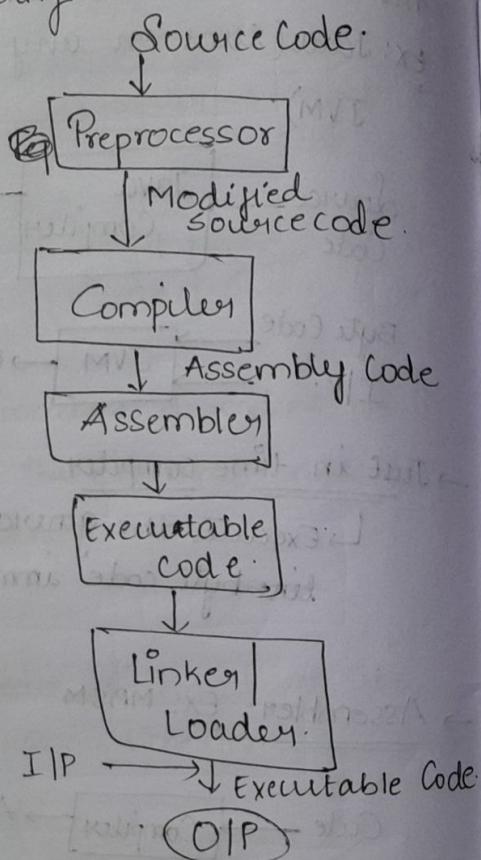
→ Just in time compiler:

↳ Executes the generated single line byte code immediately.

→ Assembler: Ex: MASM



→ Macros are replaced into the body of the code by the process of unwrapping.



→ Loader: Language processor which combines all the executable codes and gives a complete code.

→ Linker: used to link executable codes which are used to get the O/P.

* Structure of the compiler:

- 2 parts
 - Analysis: Analyse each part of source code.
 - Intermediate representation code.
 - Synthesis: Code and O/P

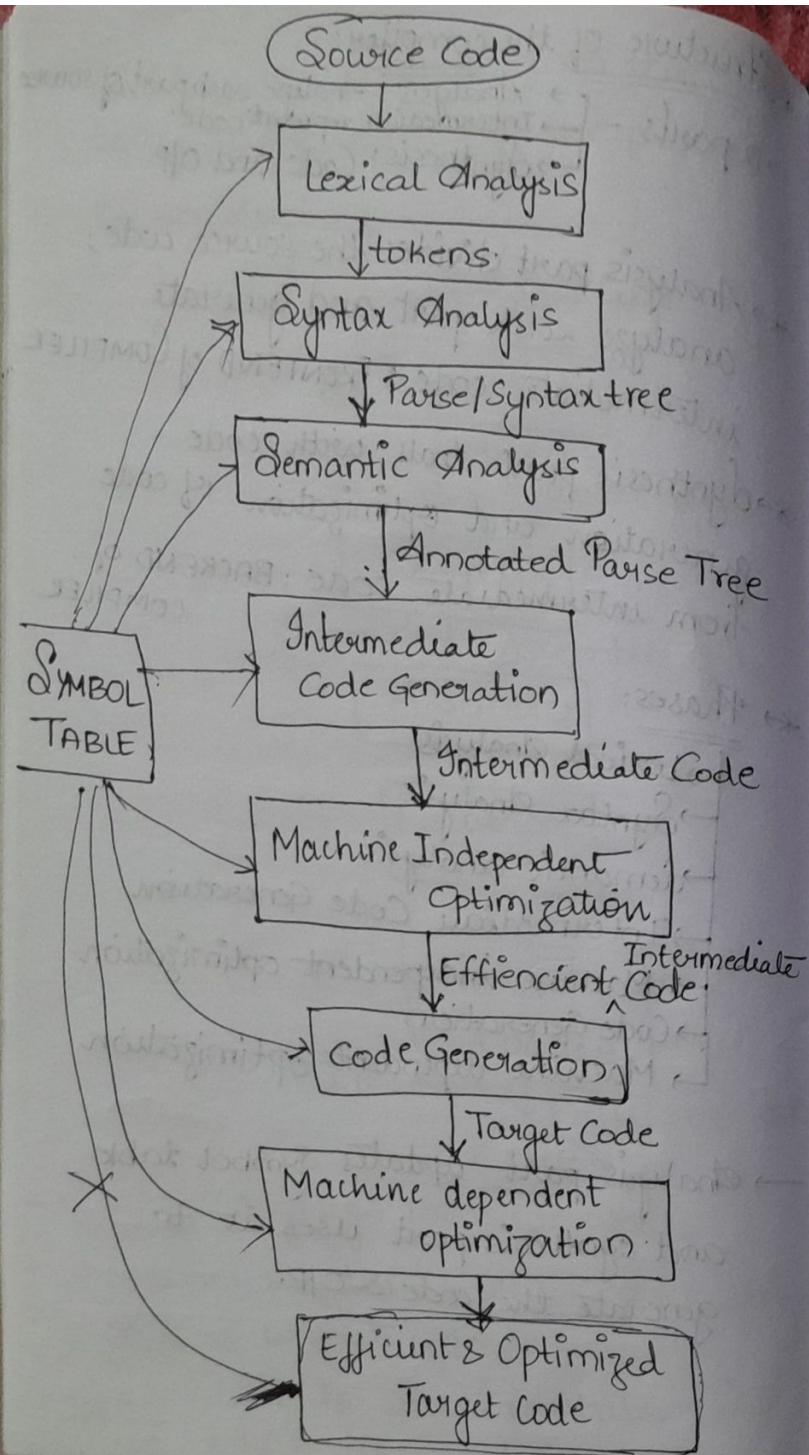
* Analysis part divides the source code; analyze each part and generate intermediate code : FRONTEND of COMPILER

* synthesis part deals with code generation and optimization of code from intermediate code : BACKEND of COMPILER

* Phases:

- Lexical Analysis
- Syntax Analysis
- Semantic analysis
- Intermediate Code Generation
- Machine independent optimization
- Code Generation
- Machine dependent optimization

→ Analysis part updates symbol table and synthesis part uses it to generate the code & O/P.



*Ex:

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

① Lexical Analysis:

→ Converts into legal set of lexemes.

→ lexeme: set of symbols that can match with one of pre-defined pattern

→ lexemes: position, =, initial, +, rate, *, 60

→ format of token: <token name, attribute value>

Ex: identifier: name of token = ID
number: name of token = num

Symbol Table:

attribute value (pointer)	Name	Type	Scope	Default
---------------------------	------	------	-------	---------

Ex:

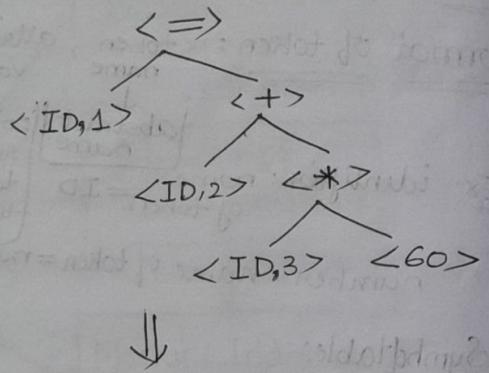
→ <ID, 1> <= > <ID, 2> <+> <ID, 3> <*>

② Syntax Analysis:

→ If any one of the lexeme is not matched to the given patterns is called lexical error.

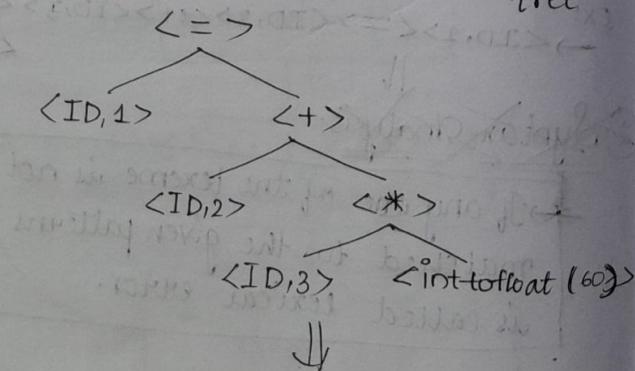
② Syntax Analysis:

Errors detected here are called syntax errors.
 → If no error, syntax tree is generated.



③ Semantic Analysis:

Annotated parse tree:



④ Intermediate Code Generation:

3-address code: ($a = b + c$)

$$t_1 = \text{int-tofloat}(60)$$

$$t_2 = \text{rate} * t_1 \quad | \quad t_2 = \langle \text{id}, 3 \rangle * t_1$$

$$t_3 = \text{initial} + t_2 \quad | \quad t_3 = \langle \text{id}, 2 \rangle + t_2$$

$$\cancel{t_4 = \langle \text{id}, 1 \rangle} = t_3$$

⑤ M10:

$$t_1 = \langle \text{id}, 3 \rangle * 60.0$$

$$\langle \text{id}, 1 \rangle = \langle \text{id}, 2 \rangle + t_1$$

⑥ Code Generation:

MOVF R0, #60.0

MULF R0, R0, ID3

ADDF R0, R0, ID2

STF ID1, R0

12/02/2024

→ Target Code is machine dependent as we convert intermediary code (machine independent) as they are just repetition of source code to target code which is a set of micro instructions. As micro instructions contains register and unique instruction set that depends on your machine (risk/cisk). Hence target code is machine dependent.

* SYMBOL TABLE:

We store info regarding constants & identifiers.

* Compiler design tools:

→ Compiler is a software
→ Lex: It is a tool that automatically generates/implements lexical analyser.

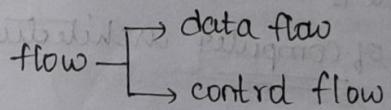
→ YACC / parser generator: automatic

parse tree generation with taking input as grammar!

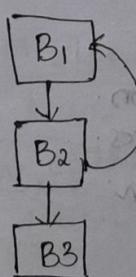
YACC: Yet Another Compiler Compiler.

→ Syntax directed translation engine.
Annotated parse-tree generation.

→ Data flow analysis engine: Tracks flow of data in 3-address code.



→ 3-address code is divided into blocks where each block has statements that execute only sequentially.



→ code generator: helps to efficiently use register.

* Applications of Compiler Technology

13/02/24

- * register: variable is allocated to the register.
- Implementation of High-level programming languages:

- * Improve the efficiency of generated code:
 - * using registers
- optimization of computer architectures:
 - ↳ Parallelism:
 - ↳ Instruction level
 - ↳ Processor level.
- Program Translation:
 - ↳ Compiled simulation.
 - ↳ Hardware synthesis.
 - ↳ Binary Translation.

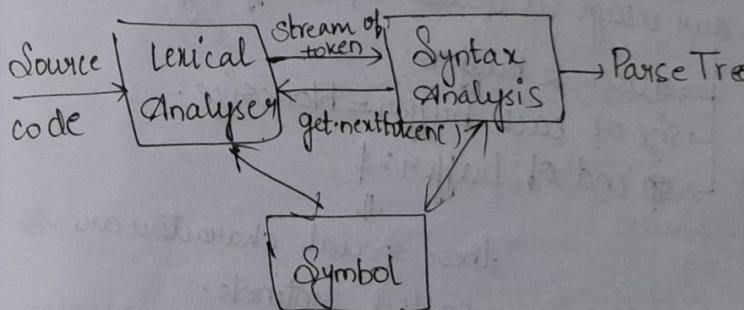
Programming language basics

- Static / dynamic distinction.

- Aliasing
- Parameter passing
- Dynamic Scope
- Explicit Access Control
- Static scope and block structures
- Environment & states.

13/02/2024

Lexical Analysis



* lexical error:

- If the lexeme is not matched with any one of the pattern, such error is called lexical error.

- Error recovery takes place in a compiler whenever an error is occurred; then it overcomes the error and compiles ~~ear~~ other statements.

Panic Mode Error Recovery:

- 1) Skip the symbols unless the lexeme is matched with any one of the pattern.
- 2) Transpose the symbols.

* The process of moving the data from main memory to the buffer and scanning for our usage is called input buffering.

- uses 2 buffers.
- size of each buffer = blocksize.
- end of buffer: \$

↓
these special characters are called sentinels.

Specification of patterns/tokens

- Regular Expressions
- Representation of regular languages
- Operations:
 - Union $L_1 \cup L_2$
 - Intersection $L_1 \cap L_2$
 - Concatenation $L_1 L_2$

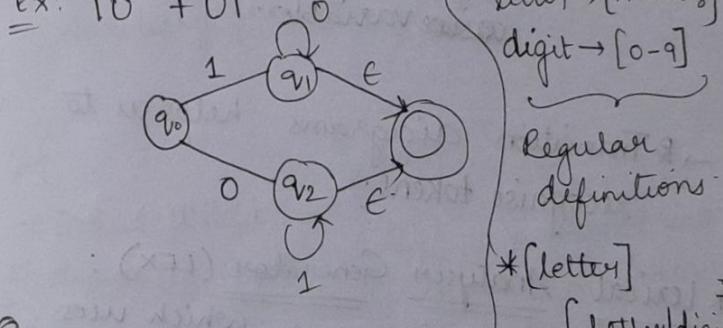
closure: L_1^* ; L_1^+
 complement $\overline{L_1}$
 reverse L_1^R

Ex: 1000110

0110 → substring

0010 → subsequence

Ex: $10^* + 01^*$



letter $\rightarrow [A-z a-z]$
 digit $\rightarrow [0-9]$

Regular definitions

$*(\text{letter})$

[letter|digit]

$\Rightarrow [l][l|d]^*$

17/02/24

→ Regular expressions are used to recognise tokens.

* The final state of the automata drawn to ~~recognise~~ recognise the given regular expression; gives the token to be sent to the parser.

→ Predefined variables:

yytext: current lexeme scanned by the lexical analyser.

character variable

yyval: holds the attribute value^{of token} which is passed to parser.

integer

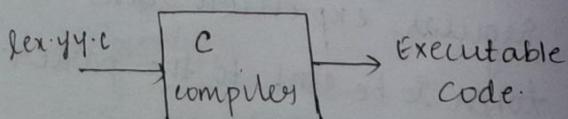
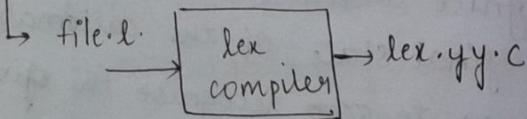
variable

→ Transition diagrams helps us to recognise tokens.

* Lexical Analyser Generator (LEX):

→ programming language ; which uses C and lexical specifications to write the code.

→ .l is the extension



Executable code → . /a.out → O/P

* LEX specifications:

Declaration section

% %

Transition Rules

% %

Auxiliary functions

Regular definitions
& variables are declared.

Set of patterns & corresponding actions

Ex:

% {

include <stdio.h>

define LT, GT, EQ, IF, THEN,
ID, NUM

% }

delim [\n \t]

ws { delim } +

digit [0 - 9]

digit { digit } { digit } * [.] ? [E | e] ? [+ | -] ?

digit { digit } { digit } * [.] ? { digit } * [E | e] ?
[+ | -] ? { digit } *

In Lex

: → any character
\: → helps in generating.

letter [a - z A - Z]

ID { letter } [{ letter } { digit }] *

% %

{ws} { no action to be performed }

IF { Return (IF); A₂; A₃ ... }
 | n no. of actions
 | → to the parser.

THEN { Return (THEN); }

< { Return (LT); }

> { Return (GT); }

= { Return (EQ); } yyval =

{ ID } { Return (ID); InstallID(); } yyval =

{ digits } { Return (NUM); InstallNUM(); }

: } any no. of patterns can be written.

% %

int main ()

{ yylex(); → func where complete lexical analysis is stored }

int InstallID() { } → Return (pointer to symbol table where the value is stored)

int InstallNUM() { } → Return (value of the number)

17/02/2024

* Syntax Analysis: (Parser)

→ stream of tokens & grammar are the I/P.

→ CFG is used $G = (V, T, P, S)$ ($V \cup T$)

$P \rightarrow E \rightarrow E + T \mid T$ → body of the production

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid id$

head of the production (only one variable)

$V = \{E, T, F\}$

$T = \{id, *, +, (,)\}$

total no. of productions = 6

order of production:

$S = E$

→ Issues in CFG:

Correction
 $A \rightarrow BA'$
 $A' \rightarrow A$

→ Left recursion: $A \rightarrow AX \mid B \mid a$.

→ Left factoring: $A \rightarrow xB_1 \mid xB_2 \mid xB_3$

→ Ambiguous: $(A \rightarrow \alpha A \mid \beta)$
 $A' \rightarrow B_1 \mid B_2 \mid \dots \mid B_n$

* Few algorithms do not produce efficient parsers.

∴ Given productions have left recursions:

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

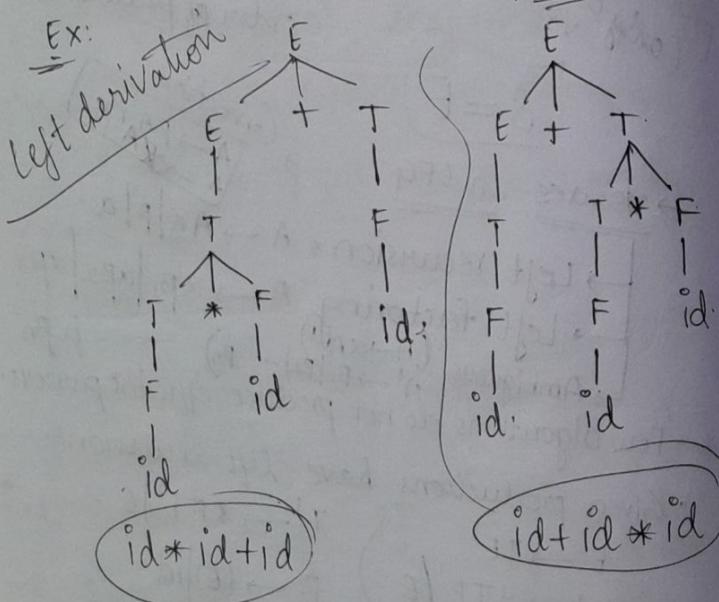
Left Factoring

Ex: $S \rightarrow iEtS \mid i^0EtSeS \mid a$
 $= S \rightarrow ietSS' \mid a$
 $S' \rightarrow eS \mid \epsilon$

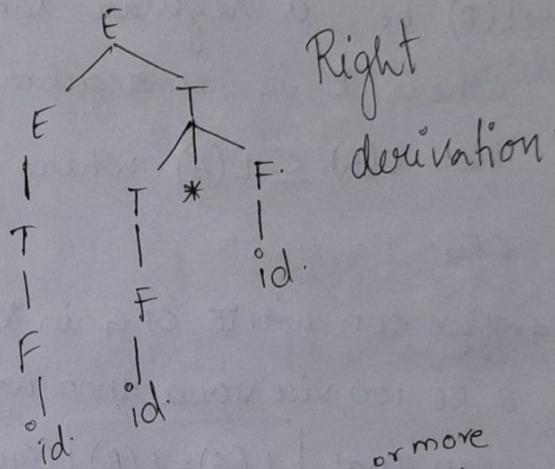
* Derivation → Left most derivation
 Right most derivation
 deriving strings from the given productions.

In any step, the replacing variable is called handle.

Parse Tree



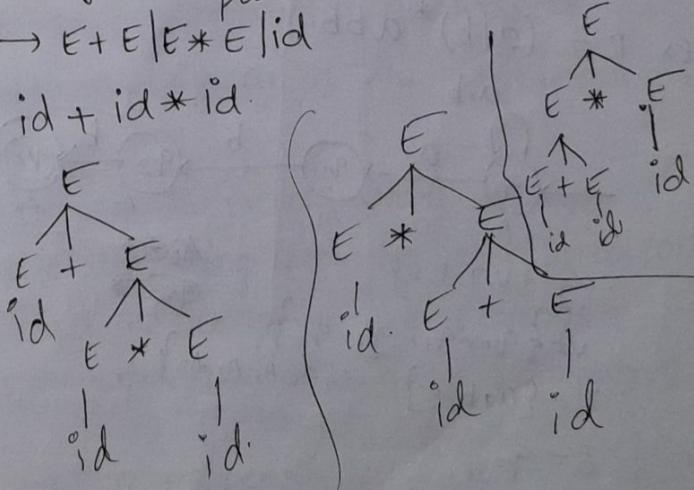
$$\begin{aligned} E = E + T &\Rightarrow E = E + T * F \Rightarrow E + F * id \\ &\Rightarrow E = E + id * id \Rightarrow E = T + id * id \\ &\Rightarrow E = F + id * id \Rightarrow E = id + id * id \end{aligned}$$



* Ambiguity: for any string z , ^{or more} ~~2~~ ≥ 2 parse trees

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

$id + id * id$



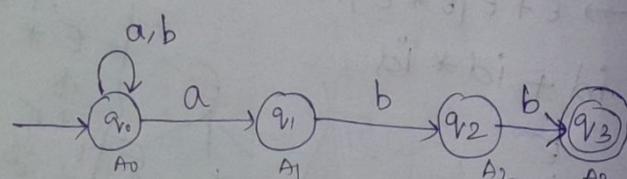
19/02/2024

* Context Free Grammar:

→ $L(R)$ be a regular language where R is a regular expression. then $L(R) \subseteq L(G)$ where $L(G)$ is CFG.

* If we can write CFG in terms of a RE (or) vice versa then we can say that $L(G) = L(R)$ (Equivalent)

Ex: $R = (a|b)^*abb$.



$$\therefore G = \{V, T, P, S\}$$

$$V = \{q_0, q_1, q_2, q_3\}$$

$$T = \{a, b, \epsilon\} = \{A_0, A_1, A_2, A_3\}$$

$$S = q_0$$

Productions:

$$A_0 \rightarrow aA_0 | bA_0 | aA_1$$

$$A_1 \rightarrow bA_2$$

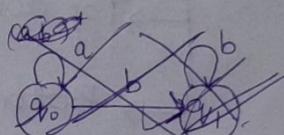
$$A_2 \rightarrow bA_3$$

$$A_3 \rightarrow \epsilon$$

→ $L(G)$ and $L(R)$ are said to be equivalent if and only if they consist of similar strings.

$$\text{Ex: } L = \{a^n b^n \mid n \geq 0\}$$

$$\begin{array}{l} A_0 \rightarrow aA_0 | bA_1 \\ A_1 \rightarrow b \\ \hline A_0 \rightarrow aA_0 | b | \epsilon \end{array}$$



* We use CFG in Syntax analysis to specify the syntax of the programming language.

* Parsing is the process where we can derive the string from the start symbol (or) to check if the given set of strings are generated from the given grammar.

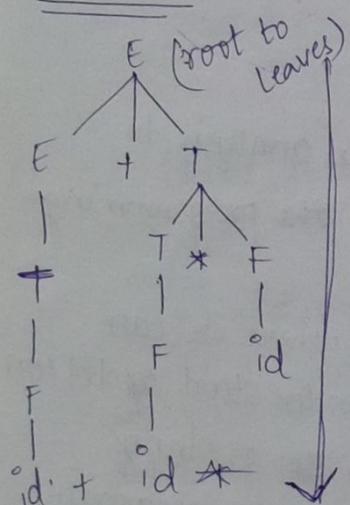
→ To get a proper string; we must keep on back track the parse tree to get a correct OIP.

→ This can be done using different parsing algorithms ~~where~~ which helps us to select the efficient production.

* Types of Parsing:

→ Top down parsing

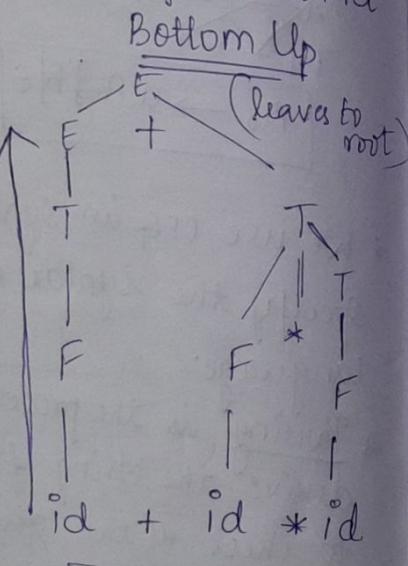
Top down



⇒ Differentiate b/w top-down & bottom-up parsing techniques (2M)

→ Bottom Up parsing

Bottom Up



* Top-Down Parsing Technique:

→ Recursive Descent Parser (RDP) : (Predictive Parser)

↳ For every non-terminal, we define a procedure.

Ex: $A \rightarrow Y$
A()

* Parser starts with calling the procedure of start symbol.

{ if (Y is ~~non-terminal~~ non-terminal)
Y();

if (Y is terminal)

terminal(Y) = current IIP
symbol
go to next symbol.

else

error

→ If above procedure works properly and reach End of string i.e. \$; the given grammar is Proper.

→ If the given grammar has left recursion; then it cannot be passed to recursive descent parser

$$\begin{array}{l} \text{Ex:} \\ \left. \begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid \text{id} \end{array} \right\} \end{array}$$

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

R DP:

```
char input[10];
for(j=0; j<n;
    if(input[j] == '$')
        printf("Success");
}
```

```
 $\Rightarrow$  char input[10];
for(j=0; j<n; j++)
{ E();
  if(input[j] == '$')
    printf("Success"); }
```

$E()$

 $\left\{ \begin{array}{l} T(); \\ E(); \end{array} \right.$
 $\left\{ \begin{array}{l} E() \\ \{ \text{if}(input[j] == '+') \\ \quad j++; T(); E(); \} \end{array} \right.$

$T()$

 $\left\{ \begin{array}{l} F(); \\ T(); \end{array} \right.$
 $\left\{ \begin{array}{l} T() \\ \{ \text{if}(input[j] == '*') \\ \quad j++; F(); T(); \} \end{array} \right.$

F()

 $\{ \text{if}(input[j] == '(') \{ j++;$
 $\quad E();$
 $\quad \text{if}(input[j] == ')') \\ \quad j++; \}$
 $\text{else if}(input[j] == 'i') \\ \quad j++; \}$
 $\}$

2d02/2024

\rightarrow LL(1) parser:

$\xrightarrow{\text{LL(1)}}$ 1 symbol is parsed at a time.
 \uparrow left most derivation is used.

Scanning of the string
 is from left
 to right

* for any G (CFG) if we can construct LL(1) parser; such G is called LL(1) grammar.

\rightarrow The grammar G must ~~be~~ not have
 left recursion and ambiguity.

* Ex: $A \rightarrow \alpha \mid \beta \Rightarrow$ LL(1) grammar (Checking)

(i) $\text{First}(\alpha) \cap \text{First}(\beta) = \emptyset$

(ii) $\alpha \xrightarrow{*} \epsilon \Rightarrow \text{Follow}(A) \cap \text{First}(\beta) = \emptyset$

$\beta \xrightarrow{*} \epsilon \Rightarrow \text{Follow}(A) \cap \text{First}(\alpha) = \emptyset$

* FIRST SET and FOLLOW SET:

$\alpha \rightarrow \text{terminal}$: First Set only

$\alpha \rightarrow \text{Non-terminal}$: Both first set & follow set.

$\alpha \rightarrow \text{grammar symbol}$. if terminal symbols are derived from α

$$\text{First}(\alpha) = \{\text{Set of terminal symbols}\}$$

$\text{Follow}(\alpha) = \text{symbol that is following } \alpha \text{ in the body of production.}$

$\alpha \rightarrow \text{non-terminal}$; find a terminal.

First Set: $\Rightarrow \text{First set of current I/P symbol}$

\hookrightarrow if current I/P symbol is in First set (α); Such α is substituted

Followset: used when non-terminal produces ϵ .

* First Set:

Let x be grammar symbol.

1. $\text{First}(x) = \{x\}$; if x is terminal.

2. If x is non-terminal

$$x \rightarrow y_1, y_2, y_3, \dots, y_k$$

if y_i is terminal symbol

$$\text{First}(x) = \{y_i\}$$

If y_i is non-terminal having ϵ

$$\text{First}(x) = \text{First}(y_1) \cup \text{First}(y_2) \cup \dots \cup \text{First}(y_k)$$

If $x \rightarrow \epsilon$ $\text{First}(x) = \{\epsilon\}$

$$\left. \begin{array}{l} \text{EX: } E \rightarrow TE \\ E' \rightarrow +TE' \mid \epsilon \\ T \rightarrow FT \\ T' \rightarrow *FT' \mid \epsilon \\ F \rightarrow (E) \\ F \rightarrow id \end{array} \right\} \quad \begin{array}{l} \text{First}(E) = \text{First}(T) = \text{First}(F) \\ = \{\{, id\} \} \\ \text{First}(F) = \{\{, id\} \} \\ \text{First}(E') = \{+, \epsilon\} \\ \text{First}(T') = \{\ast, \epsilon\} \end{array}$$

$$\left. \begin{array}{l} \ast \rightarrow E \rightarrow ET \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid id \end{array} \right\} \quad \begin{array}{l} \text{First}(E) = \text{First}(T) \\ = \text{First}(F) \\ = \{\{, id\} \} \\ \text{First}(F) = \{\{, id\} \} \end{array}$$

$$\left. \begin{array}{l} \ast \rightarrow S \rightarrow iESS' \\ S' \rightarrow eS \mid \epsilon \\ E \rightarrow b \end{array} \right\} \quad \begin{array}{l} \text{First}(S) = i \\ \text{First}(S') = \{e, \epsilon\} \\ \text{First}(E) = \{b\} \end{array}$$

* $S \rightarrow (L) | a$ } Left recursion $A \rightarrow \beta A'$
 $L \rightarrow L, S | S$ $A' \rightarrow \alpha A' | \epsilon$

$\overbrace{S \rightarrow (L) | a}$ $\text{First}(S) = \{ (, a \}$
 $L \rightarrow \varnothing | L'$ $\text{First}(L) = \text{First}(S) = \{ (, a \}$
 $L' \rightarrow , S | \epsilon$ $\text{First}(L') = \{ , , \epsilon \}$

* $\text{First}(+EE') = \{ + \}$

* $\text{First}(ET) = \text{First}(E)$

↓
if E is having ϵ

$\text{First}(E) \cup \text{First}(T) \cup \{ \epsilon \}$
↓
if T' is having ϵ

* Follow Set

↳ for non-terminals.

① If 'S' is start symbol; then $\text{Follow}(S) = \{ \$ \}$

② Let A be any non-terminal except start symbol; $X \rightarrow \alpha A \beta B$

if B is terminal

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

if B is non-terminal
then $\text{Follow}(A) = \text{First}(B) - \{ \epsilon \}$

③ $X \rightarrow \cancel{AB}$ (or) $X \rightarrow \alpha A$
if $B \rightarrow \epsilon$
then
 $\text{Follow}(A) = \text{Follow}(X)$

$\Rightarrow \underline{\underline{EX}}:$ $E \rightarrow TE'$ $\text{Follow}(E) = \{ \$, \} \}$
 $E' \rightarrow +TE' | G$ $\text{Follow}(E') = \{ \$, \} \}$
 $T \rightarrow FT'$ $\text{Follow}(T) = (\text{First}(E) - \text{Follow}(E)) \cup \text{Follow}(E')$
 $T' \rightarrow *FT' | \epsilon$
 $F \rightarrow (E)$
 $F \rightarrow id$ $= \{ +, \$, \} \}$

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

$\text{Follow}(F) = (\text{First}(T') - \{ \epsilon \}) \cup \text{Follow}(T)$
 $= \{ *, +, \$, \} \}$

$\because E'$ has
 ϵ product