

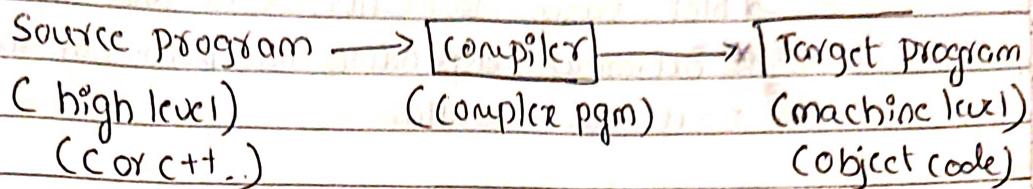
13/01/2018

Chapter-01 - Compilers

What are compilers?

Compilers are computer programs that translate one language to another. (high level to machine level)

Example: gcc, Turbo C compiler ...



Basic organization and operation of compiler should be known.

Prerequisites: Automata theory, Data structures, Discrete Maths
Machine architecture, assembly language.

Why compilers? History [John von Neumann]

Computers perform desired computations, which require writing a sequence of codes or programs.

Initially programs were written in machine language.

* The lowest-level programming language understood only by computers.

* Machine language consists of entire numbers.

Example: C7 06 0000 0002 - // Represent instruction

to move number 2 to location

0000 on Intel 8x86 processors.

Writing such programs is extremely difficult and time consuming, so it is replaced by

Assembly language is in which instructions and memory locations are given symbolic forms.

Example : mov x, 2

An assembler translates the symbolic codes, and locations into corresponding numeric codes of memory language.

Ass language Improves speed and accuracy.

Disadvantages :-

- * Not easy to write, read and understand.
 - * Extremely dependent on particular m/c for which it was written.

To overcome above there was a need to write programs in such way that it is independent of any particular machine and capable of itself being translated by a program into executable code.

(2) Easy to write, read and understand.

Example $x=2$

* FORTRAN

* Compiler — Noam Chomsky (natural language)

- * Classification of languages according to complexity of grammars
 - * Algorithms

Chomsky hierarchy (structure of Natural language)

Made easier - compiler, partial automation.

-Classification acc. complexity of grammar & power of algorithm for cognit.

Type O unrestricted

Recursively Enumerable Turing machine

Type I

Context Sensitive

Context sensitive

Linear Bounded

Type 2

Context Free

CEFL

Pushdown

Type 3

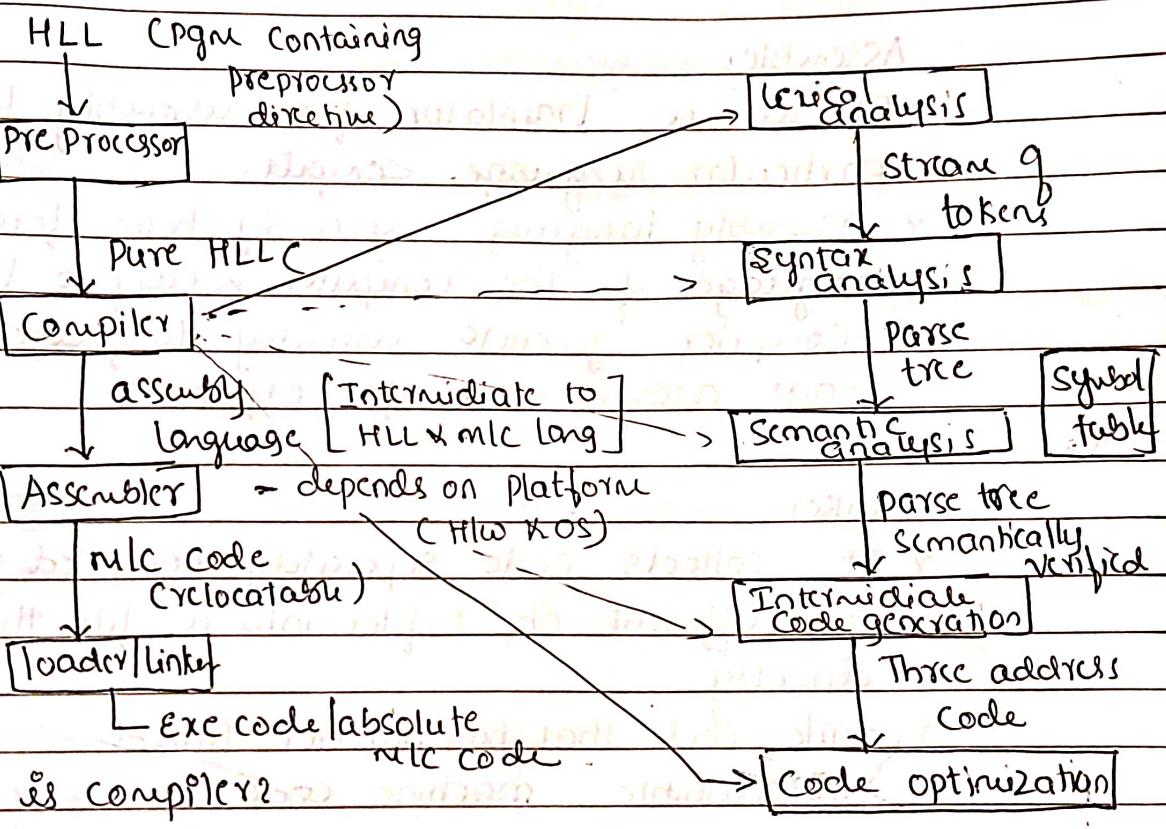
Regular

Regular language

NFA, DFA

Type 2 - Context Free Grammar useful for programming language (standard way for representing structure of programming language)

- * Parsing, Recognition of Context Free Language (1960-70)
- * Finite automata and Regular Expressions (Type 3)
- * Symbolic methods for expressing the structure of words or tokens of a programming language.
- * Code optimization (Code improvement) techniques.



What is compiler?

List of compilers?

Phases of compiler

History of compiler

Error Handler

Literal table

Phases of Compiler C Prog

Programs related to uscd together with compiler.

Interpreters:

- * It is a language translator, which executes the source program immediately rather than generating object code.
- * Any language can be either interpreted or compiled.
- * choice depends on language & circumstances under which translation occurs. Ex: BASIC, LISP
- * Compilers are preferred when speed is required.
- * Hybrid translators exist [blw interpreter & compiler]

Assembler:

- * It is a translator for assembly language of a particular language computer.
- * Assembly language is a symbolic form of machine language of the computer & easy to translate.
- * Compiler generate assembly language and then rely on assembler to get object code.

Linker:

- * It collects code separately compiled or assembled in different object files into a file that is executable directly.
- * mlc code that has not yet linked executable machine code.
- * Linker also connects an object program to the code for standard library functions & resources supplied by OS [memory allocators, I/O & O/P devices].

Loader:

- * Compiler, assembler or linker will produce data that is not yet ready to execute & fixed.
- * Principal memory references were made relative to

undetermined starting location - anywhere in memory

- relocatable

- * Loader resolve all relocatable addresses relative to given start address
- * Use of loader makes executable code more flexible (occurs behind scene)

Preprocessors

- separate program, called by compiler before actual translation begins.
- Preprocessor can delete comments, include other files and perform macro substitutions

Editors

Debugger

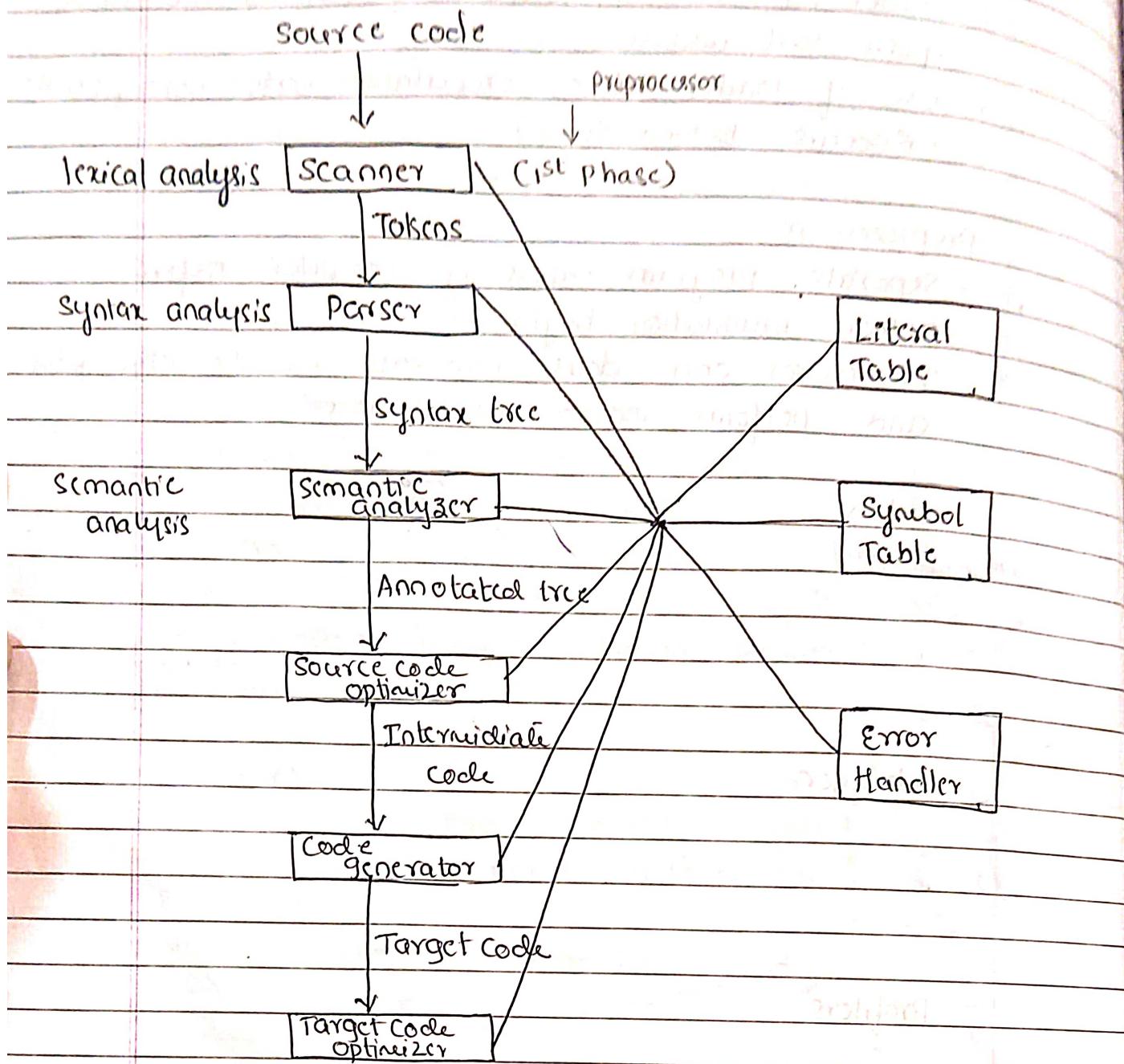
Profilers

Project Managers

Difference b/w interpreter & compiler

The Translation Process [Phases of Computer Translation]

- Compiler consists of number of Phases - performing distinct logical operations.



The Scanner: Reads the source program (stream of characters)
(lexical analysis) converts it into tokens (meaningful unit)
removes whitespace or comments

Ex: $a[\text{index}] = 4 + 2$

a - identifier , [- left bracket, index - identifier,] - right bracket

= - assignment, 4 - number, + - Plus sign, 2 - number

$$\text{id}[\text{id}] = \text{id} + \text{id}$$

nine + nine

135

1350

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow \text{id.} \end{aligned}$$

$$x = a + b * c$$

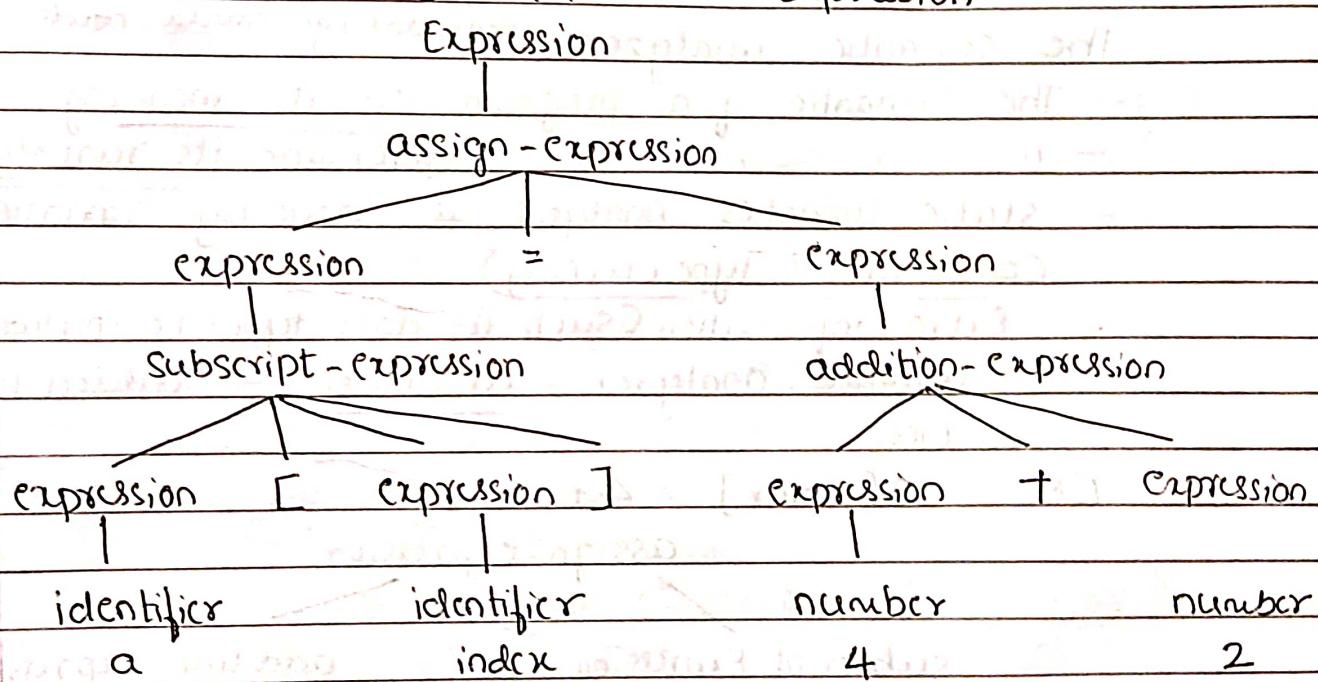
x - identifier, = - assignment, a - identifier, + - plus sign
 b - identifier, * - multiplication sign, c - identifier.

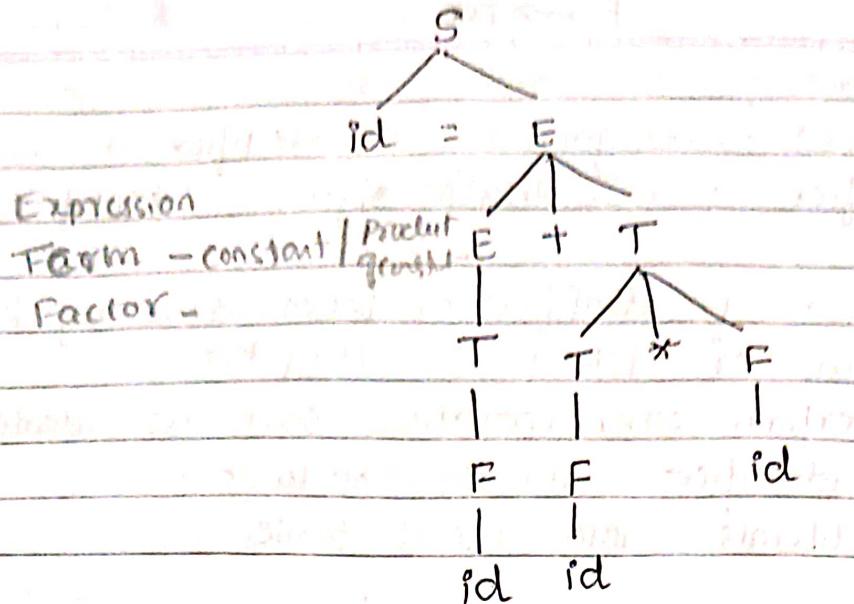
- * Identification of identifiers or tokens done by Regular Expressions or tokens ex $L(L+d)^*$
- * Scanner performs other operations such as identification of
 - entering identifiers into symbol table
 - " literals into literal table

The Parser (Syntax analysis)

- Receives source code in the form of tokens and perform Syntax analysis - which determine structure of program
- Performing grammatical analysis as in natural language
- Determine structural elements of program and their relationship.
- The Parser produces parse tree or syntax tree.

Ex: $a[\text{index}] = 4 + 2$ - Expression



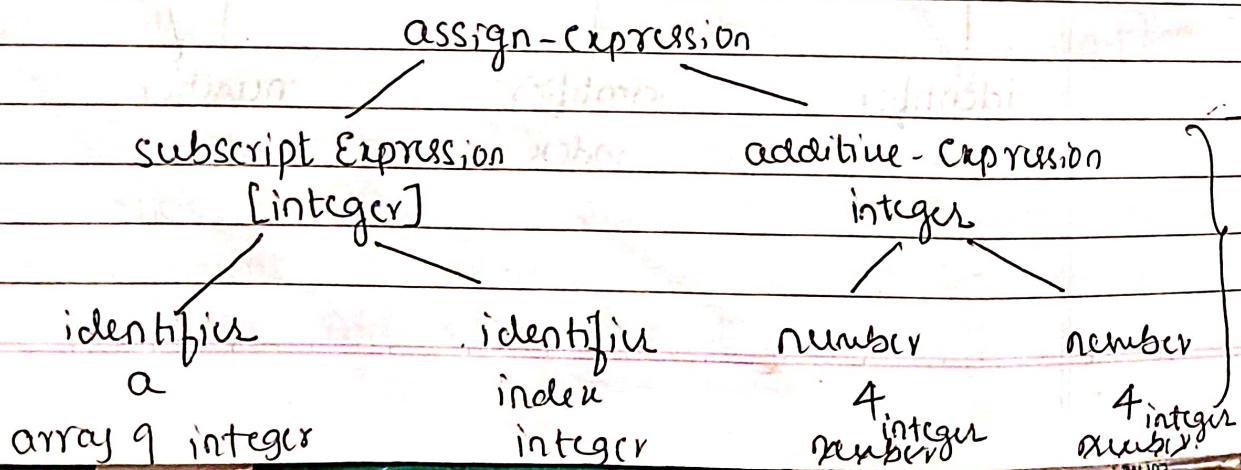


- Internal nodes of parse tree are labeled by names of structures. they represent
- leaves of parse tree represent sequence of tokens from input.
- Parse tree is inefficient in its representation
- So Parser generate Syntax tree [Condensation of information contained in Parser tree]. - abstract Syntax tree

The semantic analyzer: (Redundant info. removed, delete nodes of exp)

- The semantic of a program are its meaning.
- " " " determine its runtime behaviour.
- static semantics Analysis is done by semantic analysis (declarations, type checking)
- Extra information (such as data type) Computed by semantic analyzer - attributes - added to parse tree

Ex: $a[\text{index}] = 4 + 2$



prior information gathered before analysis

a - array of integer values with subscript from
subrange of integers

index - integer value.

semantic analyzer

- Annotate syntax tree with types of all sub-expressions
- precision, & checks that assignment makes
- declares a type mismatch error
- Result of annotated tree shown in Figure
- static semantics, dynamic semantics [salary 10% to 20%]

The Source Code optimizer:

- Compiler includes code improvement or optimization steps after semantic analysis.
- Code improvement based on source code
- Code optimization depends on compilers.
- Example the expression $4+2$ can be precomputed by compiler to return 6. [constant folding]
- Can be performed directly on (annotated) syntax tree by collapsing right-hand subtree of root.
- Three-address code (standard), P-code (Pascal compilers)

6

integer

Ex : $t = 4 + 2$

$a[i] = t$

- Intermediate code [code representation / intermediate between source and object code]

The Code generator

- Receives intermediate code [IR] and generates code for target machine.
- Properties of target machine become important
- Representation of data & instruction play major role [how many bytes or words variable of int & local data type occupy in memory]

prior information gathered before analysis

- array of integer values with subscript from
subrange of integers
- index - integer value.

semantic analyzer

- Annotate syntax tree with types of all sub-expressions.
 - checks that assignment makes sense
- If mismatch declare a type mismatch error
- Result of annotated tree shown in Figure
- static semantics, dynamic semantics [salary 10% to 20%]

int a, c;
float b;
c = a + b;

The Source code optimizer:

- Compiler includes code improvement or optimization steps after semantic analysis.
- Code improvement based on source code
- Code optimization depends on compilers.
- Example the expression $4+2$ can be precomputed by compiler to return 6. [constant folding]
- Can be performed directly on (annotated) syntax tree by collapsing right-hand subtree of root.

RHS -
number

Ex $t = 4+2$

$a[i] = t$

- Intermediate code [code representation between source and object code]

The Code generator

- Receives intermediate code [IR] and generates code for target machine.
- Properties of target machine become important
- Representation of data & instruction play major role [how many bytes or words variable of int & float data type occupy in memory]

Target code generator.

Input : IR C Three code address

Process: generate Code \rightarrow Assembly language.

Target Code	Mov R0, index		Suppose a[5] the
	MUL R0, 2		R0 \Rightarrow 5, a \Rightarrow 1000
	MOV R1, Xa		R0 \Rightarrow 10 [integer 2 byte memory]
	ADD R1, R0		R1 \Rightarrow 1000
	MOV *R1, 6		\rightarrow R1 \leftarrow 1010 \leftarrow 6
	1000 1002 1004 1006 1008 1010 1012		
	[] [] [] [] [] [] []		
	a[0] a[1] a[2] a[3] a[4] a[5] a[6]		

Target Code Optimizer

- * improve target code
 - choose addressing modes
 - replace slow instruction by faster one
 - eliminating redundant or unnecessary operations.

Improvement for the above code

- * use shift instruction instead of multiplication [second line]
[expensive in terms of execution time]
- * use powerful addressing mode as indexing addressing.

Mov R0, index // R0 - 5

multiplicand - SHL R0 // 0101 - 5

indexing. Mov Xa[R0], 6 // 01010 - 10

01010, - 10

Major Datastructures in a compiler

- * The interaction b/w the algorithms used by compiler and data structures that support these phases is strong.
- * Ideally a compiler should compile a program in time proportional to size of program i.e $O(n)$, n - pgm size
- * Data structures needed by the phases as part of their operation and that serve to communicate information among the phases.

Tokens [Lexeme - is an actual character sequence]

- * Token is a sequence of characters that can be treated as a unit in the grammar of programming language.
- * In programming language, keywords, constants, identifiers, strings, numbers, operators and punctuations are considered as tokens. (In C six type of tokens)
Ex: int value = 100; [keyword, const, identifier, string, operator, symbol]
contains tokens int (Keyword), value (Identifier), = (operator), 100 (constant), ; (symbol)
- * The Scanner represents tokens of source language as a value of an enumerated data type.

* Example: i) Sum = 3 + 2;

Tokens: sum, =, 3, +, 2, ;

ii) if ($x_1 * x_2 < 1.0$) {

Tokens: if, C, x_1 , *, x_2 , <, 1.0, }

Table	Tolsen	Lexemes	Pattern	Co regular expression
1.1	Const	Const	Const	Const expression
	if	if	if	Notation used to distinguish between if and else
	relation	<, <=, =, >, >=	< <=, =, <, >, >=	tokens
	id	Pi, count, D2	Letter (Letter + digit)*	
	num	3.1426, 0.6	numeric constant [0-9]*	
	literal	"Core dumped"	char b/w "	

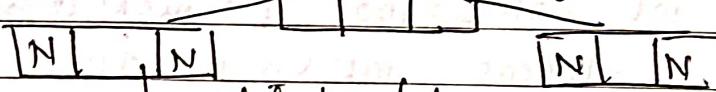
Example : Const Pi = 3.1416

Pi is a lexeme for the token "identifier"

Syntax Tree

Syntax Tree [Tree representation of the abstract syntactic structure of source code]

- * Pointer based structure, dynamically allocated
- * Syntax tree represented using single variable pointing to root node
- * Each node represent information [record] collected both by parser & semantic analyzer [linked list / array]



Symbol table ... datatype / dynamically allocated or symbol table

- * This data structure keeps information associated with identifiers: functions, variables, constants and data types.
- * Symbol table interacts with every phase of compiler.
 - Scanner, parser, or semantic analyzer
 - Semantic analyzer will add data type and other information.
 - optimization and code generation phases will use the information provided by symbol table.
 - Symbol table is accessed frequently - access operation needs to be efficient (constant-time operations)
- For this standard datastructure is hashtable O(n)

Literal table

- Stores constants and strings used in program.

- insertion & lookup are frequent.
- Literal table need not allow deletions [bcz of global data]
- Constant or string appear only once in table.
- Literal table reduces the size of a program in memory by allowing reuse of constants & strings.
- Required by code generator.

Intermediate Code

- Three-address code & P-code
- Different optimizations performed.
- Depending on above parameters code may be kept as an array of text strings, a temporary text file or as linked list of structures.

Temporary files

- As computers did not possess enough memory for an entire program to be kept in memory during compilation.
- The above problem solved using temporary files, to hold results of intermediary steps.
- Keep only enough information from earlier parts of source program to enable translation to proceed.

GATE Question on Lexical Analysis

Printf("i=%d,xi=%x",i,xi);

printf

"i=%d,xi=%x"

i

x

i

i

i

i

Q.2. In a compiler, keywords of a language are recognized during lexical analysis of the program.

Example:

$$x = a + b * c$$

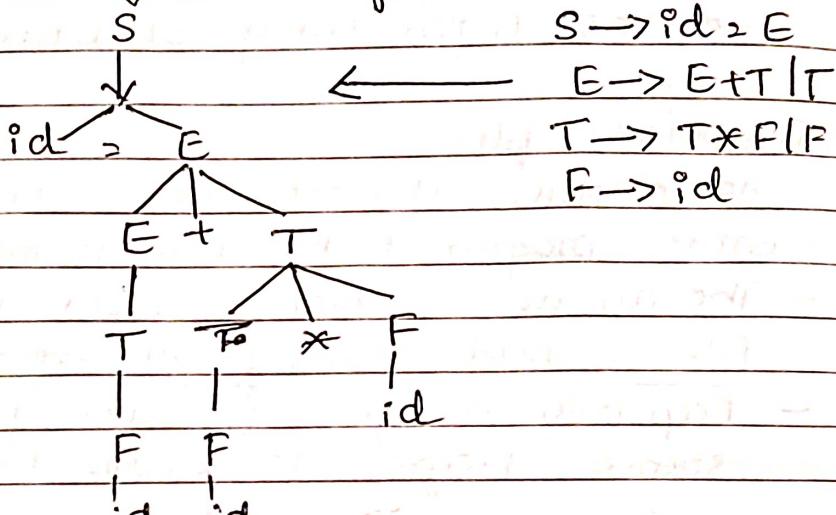
Lexical Analysis

$$\downarrow$$

$$id = id + id * id \text{ (Token)}$$

Syntax Analysis

\downarrow (Syntax tree)



Semantic analyzer [semantically verified]

\downarrow
Annotated Tree

(Type checking, Coercions, Call by type)

ICG

(Three-address
code)

$$t1 = b * c$$

$$t2 = a + t1$$

$$x = t2$$

\downarrow CO

Target code

generator optimized
(shorter code)

$$t1 = b * c$$

$$x = a + t1$$

\leftarrow Source code optimizer

MLC

independent

MLC

dependent

Writing code for assembly [depends on platform]

Target code (Assembly level) Generation.

need not be save

Mova → MUL R1, R2 a → R0
 ADD R2, R0 b → R1
 MOV R2, R1 C → R2

② Position = initial + rate * 60

\$ gcc -Wall filename.c -o filename

-Wall - enable all compiler's warning messages - generate better code

-o - specify output name (if not specified use .out)

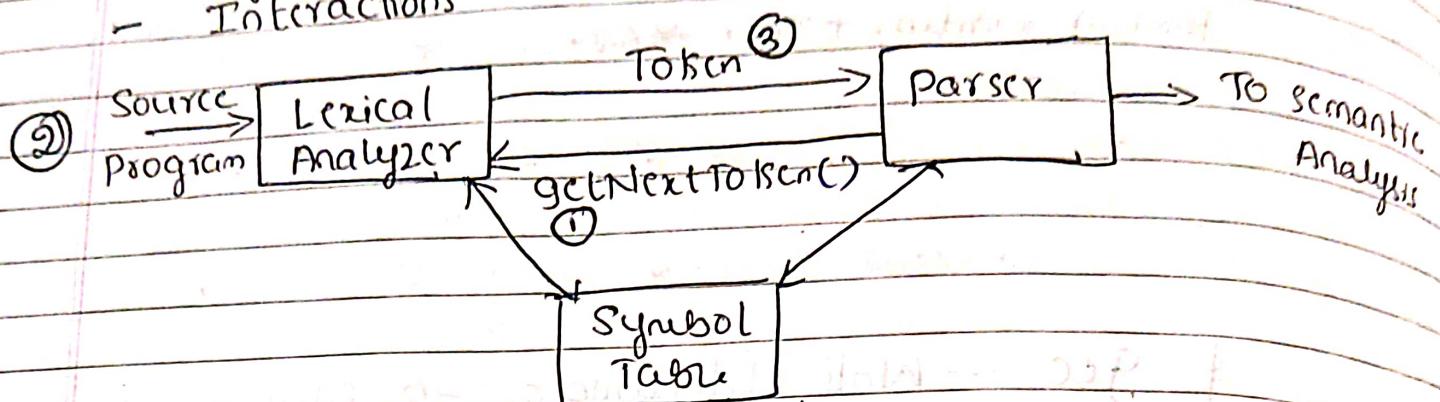
- Pre-processing { Removal of comments, macro expansion, included files, conditional compilation, Assembly, Linking }

\$ gcc -Wall -fno-strict-aliasing -O2 filename.c -o filename

→ filename.i - Assembly level instruction
→ filename.s - Assembly level instruction
→ filename.o - mle level instruction

\$ size filename.o \$ size filename

- Role of Lexical Analyzer
- Lexical analyzer when identifies a lexeme as identifier, it needs to enter that lexeme into symbol table.
- Interactions



* Other Functionalities of Lexer are

- Stripping out comments, white space (blank, newline tab or other)
- Correlating error messages generated by compiler with source program.
- associating line number with error message as lexer keep track of new line character
- source pgm copied to insert error message
- cascading of two processes -
1) scanning - 2) lexical analysis (Complex)
Simple L do not require tokenization.

Advantages of Separating lexer & parser

* Simplicity: Atleast one of the phase will be simplified
Ex: handling of whitespace & comments

Efficiency: Using buffering techniques for reading i/p can speed up compiler

Portability: I/p device specific properties can be restricted to only lexer

Tokens, Patterns & Lexemes:

Tokens: <token name, attribute value>
(optional)

main() - not a keyword w/ compiler
- same as other function ~~function~~ ~~sort()~~

Token name is abstract symbol representing lexical unit. [Keyword, identifier]

Token names are input to parser.

* Pattern: Description of the form that the lexeme of a token may take.

Ex: For a keyword - the pattern is just sequence of characters that form keyword.

For identifiers - complex pattern - $I(I+cl)^*$

* Lexeme - sequence of characters in source pgm that matches the pattern for a token.

Table 1.1 [Previous pages] shows example for above

Example: printf("Total %d\n", score)

printf, score are lexemes matching pattern for token id, "Total %d\n" is a lexeme matching literal.

Classes of token

- * one token for each keyword
- * Token for the operators either individually or in classes such as token Comparison
- * one Token representing all identifiers
- * one or more tokens representing constants [numbers and literal strings]
- * Token for each Punctuation symbol, left and right parenthesis, comma and semicolon.

Keywords

identifiers (main)

Constants

Strings

Special symbols

operators

Attributes for Tokens

- * The lexical analyzer must provide the subsequent compiler phases an additional information about particular lexeme matched.
- * So lexical analyzer returns to the parser not only a token name, but an attribute value that describes the lexeme represented by token [token name influences Parser attribute value influences translation tokens after the parser]
- * Tokens have almost one associated attribute. [attribute - may be a structure consisting several pieces of information]
- * identifier - information may be - its lexeme, its type and location at which it first found (to handle error message) - Isct in symbol table
- * Thus appropriate attribute value for an identifier is a pointer to symbol-table entry for that identifier.

Example

$E = M * C \rightarrow$

Symbol table

	id	Value
	id	Pointer to symbol table
	assigop	-
	fd	* to ST
	mult-op	-
	id	* to ST
	exp-op	-
	number	2

In practice a typical compiler instead store a character string representing

constant and use an attribute value for number, a pointer to that string

- PAGE NO. _____
DATE. _____
- * Symbol Table created by lexer and parser.
 - * Structure of symbol table is designed by compiler designer.
 - * Info - Name, type, location, scope, other attribute
 - * Operations associated with symbol table
Lookup, Insert, Modify, delete

Issues in symbol table

- 1) Format - what data structure need to be used
- 2) Access method - " Search method " " "
- 3) Location of storage - RAM or secondary memory
- 4) Scope

Lexical errors

- An input that can be rejected by lexer.

Type

- A character sequence that can not be scanned into any valid token is a lexical error.
 - Misspelling of identifiers, keyword or operators are
- Ex: if the string fi is encountered for 1st time
in a C program, then since it is not a keyword or operator, it is misspelled.

* Lexer can not tell whether fi is a misspelling of if keyword or an undeclared function identifier.

* Since fi is valid lexeme for token id, the lexer must return token id to parser & let other phases of compiler - probably parser in this case - handle it.

* When none of the patterns for tokens matches any prefix of remaining input - then lexical analyzer is unable to proceed.

— Use a simple method "Panic mode" recovery

Classification of compile time errors

- * Lexical * Syntactic * Semantic

Unmatched patterns are deleted successively from remaining input, until lexical analyzer can find well-formed token at beginning of what input is left.

Other error-recovery techniques are

- Delete one character from remaining input.
- Insert a missing character into remaining input.
- Replace a character by another character.
- Transpose two adjacent characters.

(Data structures) → Input Buffering: [Scanning Process]

* How to recognize lexeme of the input.

* We need to see one or more characters beyond next lexeme before we can be sure we have right lexeme.

* Example: we are not sure we have seen the end of an identifier until we see a character that is not a letter or digit & is not token id. In C single character operators like -, = or < could also be beginning of two-character operator like →, = or <=.

* So we have Two-buffer Scheme to handle lookaheads.

Buffer pairs:

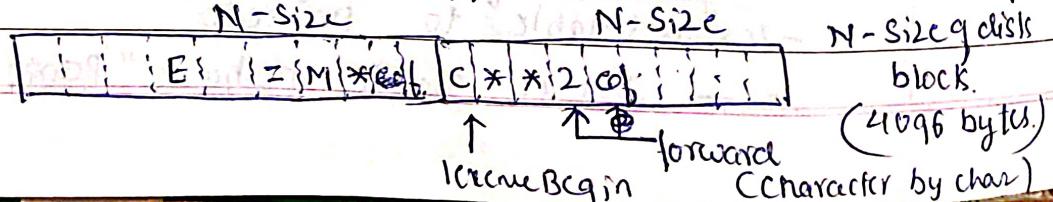
* Special buffering techniques are required, since large number of characters must be processed during compilation.

* Techniques developed to reduce overhead required to process a single input character.

Tokenizer

→ Usage of Two buffers - reloaded alternatively.

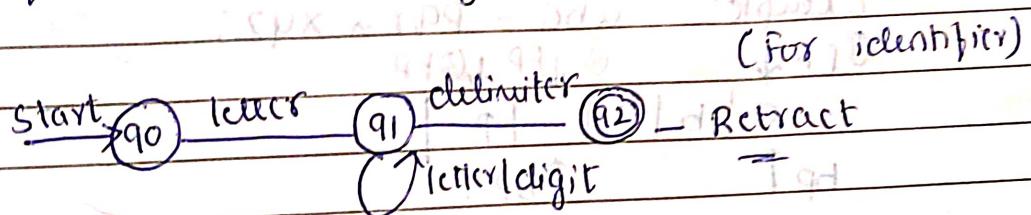
Single state



currentLexeme: set of characters b/w two pointers.

forward

- * Using system read command, read N characters into buffer, rather than one system call per character.
- * If fewer than N characters remain in input file then a special character eof needs to be inserted which is different from any other characters.
- * Two pointers to the input are:
 - lexemeBegin - marks beginning of current Lexeme
 - forward - Scans ahead until a pattern match is found.
- * Once Lexeme is determined, forward is set to character at its right end.
- * Lexeme is recorded and returned to parser.
- * lexemeBegin is set to character immediately after Lexeme found.
- * Once we reach end of one buffer, we must reload other buffer from input and move forward to begining of newly loaded buffer.



Example 2:

One i/p buffer - int i, j; 2 i/p buffer.
lexemeBegin^① int i, j; @lexemeBegin
 ② ↑ Delimit ③ --- ④
 ↑ -ptr (stops) -
 - int is lexeme
 > i is lexeme
 -> j is lexeme

Sentinels

If scheme of Input Buffering is used then we must ensure

→ each time we advance forward (we have not moved off end of the buffer). If so then must reload other buffer.

i.e. For each character read we make a test

→ one for end of buffer

→ what character read

We can combine buffer-end test for current character if we can store Sentinel character at end of buffer

Sentinel: IS a special character that can not be part of the source program and choice is the character 'eof'.

- eof is marker for the end of entire input.
- Any eof that appear other than at end of eof of a buffer means that the i/p is at an end.

Example: abc = pqr * xyz;

① FP ↓ ② FP ↓ ③ FP ↓

a | b | c | = | p |

bp ↑

↑ bp

abc → identifier

= → operator

← Reload 2nd buffer.

q | r | * | x | y | z | eof

Class Exercise

① int max(x, y)
 int x, y;
 → 25 tokens

{ /x - Find max q x & y x }

return (x > y ? x : y);

}

Is it possible to write entire 'c' program in a line? (after line)

② printf("%d Mai", x); → 8 tokens
 string literal one operator

③ xx → 2 tokens (longest match)

④ Entries of a symbol table

x a=10

b>20

* printf("%d", a+++b); No lexical error

a+++b Syntax error. (Can not

* int i, 2000;

Token for entire number → 5 tokens

longest match correct one
 a++ + b

x a++ *|+|b → No lexical error, No syntax error
 (No parsing)

* a++ |+|b → No lexical error, Syntax error.
 one opc other operand

No operator.

int main()

{

int a=10, b=20;

printf("sum is %d", a+b);

return 0;

0* a++ + b

Def Proc myproc (int A, float B) {
int D, E;

D = 0;

E = A | round(B)

if (E > 5) {

Print D

}

}

Symbol	Token	Dtype
myproc	id	Proc name
A	id	int
B	id	float
D	id	int
E	id	int

int main() {
int a, b, c; b = 10, c = 10;
a = b + c * 100;

Symbol	Token	Dtype
main	Id	funname
a	Id	int
b	Id	int
c	Id	int

PAGE NO. _____
DATE _____

Specification of Tolsons - Regular Expression (RE)

RE are important notation for specifying lexeme pattern.

- Formal Notation
- How to use RE in lexical-analyzer.
- Conversion of RE to automata to recognize tokens specified.

Terminologies

Alphabet: Finite set of symbols. Ex: letters, digits and punctuation.

Ex: {0,1} — Binary alphabet

ASCII — Alphabet used in SW systems

String: Finite sequence of symbols drawn from alphabet (sentence or word)

$|S|$ — length of string — no. of characters in a string.

Empty string — ϵ

[Compiler] — 8

Language: Finite set of strings.

Concatenation: If x and y are strings then xy denotes concatenation.

$S\epsilon = \epsilon S \rightarrow S$ [Empty string is identity under concatenation]

$$S^0 = \epsilon$$

$$S^1 = S$$

$$S^2 = SS$$

$$S^3 = SSS$$

Exponentiation of strings.

Operations on languages

Union — $L \cup M = \{S | S \text{ is in } L \text{ or } S \text{ is in } M\}$

Concatenation — $L \cdot M = \{S | S \text{ is in } L \text{ and } t \text{ in } M\}$

Closure — Kleene closure of L $L^* = \bigcup_{i=0}^{\infty} L^i$ ($\cup \Sigma$)

Positive closure of L $L^+ = \bigcup_{i=1}^{\infty} L^i$

(Σ^*) Kleene closure denoted by L^* of a language
set of strings by concatenating L zero or more times
 L^0 - concatenation of L - ϵ

(Σ^+) positive closure - denoted by L^+ - same as Kleene closure but without term L^0 .

Ex $L = \{A..Z, a..z\}$ $D = \{0, 1..9\}$

$\Rightarrow LUD =$ set of letters and digits - 62 strings of length one.

2) $LD =$ set of strings of length two - 320 strings.

3) L^* - set of all strings of letters, including ϵ

4) L_4 - set of all 4-letter strings.

5) $L(LUD)^*$ - set of all strings of letters and digits beginning with a letter.

6) D^+ - set of all strings of one or more digits.

Regular Expressions

RE for valid identifier - $l(l+d)^*$ (Note: '-' can be included in letter), where l - letter and d - digit.
or letter $(letter | digit)^*$

RE are built recursively out of smaller regular expressions using below rules.

* Each regular expression r denotes a language $L(r)$.

Rules: To define RE over some alphabet Σ .

1) ϵ is a RE and $L(\epsilon) = \{\epsilon\}$

2) If a is a symbol of alphabet Σ then a is RE and $L(a) = \{a\}$ - language with one string of length one.

If r and s are regular expressions denoting $L(r)$ and $L(s)$ then

1. $(r)s$ is a regular expression denoting language $L(r)L(s)$
 $L(r)UL(s)$
2. $(r)s$ is a RE denoting language $L(r)L(s)$.
3. $(r)^*$ is a RE denoting $(L(r))^*$.
4. r is a RE denoting $L(r)$. - we can use parenthesis around RE without changing the language meaning.

Priority - precedence of operators. $(^* > \cdot > |)$

- 1) The unary operator $*$ has highest precedence and is left associative.
- 2) Concatenation has second highest precedence & is left associative.
- 3) $|$ has lowest precedence and is left associative.

Under above assumption

$(a)|((b)^*(c))$ can be replaced by $a|b^*c$

Let $\Sigma = \{a, b\}$ what are the languages denoted by following REs.

- 1) $a|b$
- 2) $(a|b)(a|b)$
- 3) a^*
- 4) $(a|b)^*$
- 5) $a|a^*b$
- 6) $\{a, b\}$
- 7) $\{aa, ab, ba, bb\}$
- 8) $\{\epsilon, a, a^2, a^3, a^4\}$
- 9) $\{\epsilon, a, b, aa, ab, bb, ba, bb\}$ or $(a^*b^*)^*$
- 10) $\{a, b, ab, aab, aaab\}$

A language defined by a regular expression is called regular set. If two regular expression r and s denote the same regular set, then they are equivalent i.e. $r = s$

$$\text{Ex } (a|b) = (b|a)$$

Algebraic laws for regular expression for which RE are equivalent

$$r|s = s|r \quad | \text{ is commutative}$$

$$r|(s|t) = (r|s)|t \quad | \text{ is associative}$$

$$r(st) = (rs)t \quad \text{Concatenation is associative}$$

$$r(s|t) = (rs)|rt \quad \} \quad " \text{ distributive}$$

$$(s|t)r = srltr \quad \} \text{ over } 1$$

$$\epsilon r = r\epsilon = r \quad \epsilon \text{ is identity for concatenation}$$

$$r^* = (r|\epsilon)^* \quad \epsilon \text{ is guaranteed in a closure}$$

$$r^{**} = r^* \quad * \text{ is idempotent}$$

Regular Expression Exercises.

1) Write regular expression for valid C identifiers.

$$\text{letter} \rightarrow A|B|...|z|a|b|...|z|$$

$$\text{digit} \rightarrow 0|1|...|9|$$

$$\text{id} \rightarrow \text{letter}(\text{letter}| \text{digit})^*$$

2) RE For unsigned numbers such as 5280, 0.01234
6.336E4 or 1.89E-4

$$\text{digit} \rightarrow 0|1|...|9|$$

$$\text{digits} \rightarrow \text{digit}.\text{digit}^*$$

$$\text{Optional Fraction} \rightarrow \cdot \text{digits} |\ \epsilon \quad (\cdot \text{ has to be followed by digit at least one})$$

$$\text{Optional Exponent} \rightarrow (\text{E} (+|-|\epsilon)) \text{ digits} |\ \epsilon$$

$$\text{Number} \rightarrow \text{digits} \text{ optional fraction} \text{ optional exponent}$$

Note: \ ". ^ \$ [] * + ? { } | / - special Meaning
 these symbols need to be turned off, if they need to be
 represented in character string - use " " or use \ (backslash)

Extensions of Regular Expression,

1) One or more instances - unary, postfix operator (*).

Positive closure of a RE and its language

i.e if r is a RE then $(r)^*$ denotes $(L(r))^*$.

* has same precedence & associativity as operator \times .

$$r^* = r^1 \mid \epsilon$$

$$r^+ = rr^* = r^*r \quad \left. \begin{array}{l} \text{Relate 1st cn closure} \\ \text{digit(digit)}^* - \underline{\text{digit}} \end{array} \right.$$

2) Zero or one instance - The unary, postfix operator

? - has same precedence and associativity as *, +.

i.e $r \equiv r \mid \epsilon$

$$L(r?) = L(r) \cup \{\epsilon\} \quad ; \text{Exponent.}$$

$$CE [+|-]? \text{digit}^+)?$$

3) Character class - []

A Regular Expression $a_1a_2\dots a_n$ where a_i 's are symbols of the alphabet, can be replaced by shorthand $[a_1a_2\dots a_n]$

i.e $ab1\dots 12 \equiv [a-z]$

$ab1c \equiv [a-z]b1c$

Ex: digit $\rightarrow [0-9]$

digits $\rightarrow \text{digit}^+$

number $\rightarrow \text{digit}(\text{. digits})? CE [+|-]? \text{digit}^+)?$

Describe language denoted by following RE.

1. $a(a|b)^*a$ - String of a's & b's that start & end with a.

2. $((\epsilon|a)b^*)^*$ - String of a's & b's

3. $(a|b)^*a(a|b)(a|b)$ - String of a's and b's that the character third from the last is a.

4. $a^*ba^*ba^*ba^*$ - String of a's & b's that contain three b.

5. $(aa|bb)^*((ab|ba)(aa|bb)^*(ab|ba)(aa|bb)^*)^*$
 Strings of a's and b's that has even number of a's and b's.

Regular expression that matches "1" — "11"

Recognition of tokens

Branching statement

If (expression) then stmt
or if (expression) then stmt else stmt
or ε

stmt → if expr then stmt

| if expr then stmt else stmt

| ε

expr → term relop term

| term

term → id

| number

— tokens wrt
lexical analyzer.

Pascal code - conditional statement

- then appears explicitly after Condition.
- relop - Comparison operators - =, <>
- terminals - if, then, else, relop, id, number
- White Pattern or RE.

Chomsky hierarchy.

Hierarchy for classification of languages.

Chomsky classified languages and grammar into 4 types.

Type-3

Type-2

Type-1

Type-0

Type-3: If the grammar has all the productions of the form

$$A \rightarrow \alpha B \beta$$

Where A, B ∈ V (Non-terminal)

α, β ∈ T* (Terminal)

If variable occurring at rightmost side then

grammar said to be Right linear grammar (RLG)
i.e. $A \rightarrow \alpha B \beta$
variable rightmost side of grammar.

$A \rightarrow B\alpha | B$

variable on left side of grammar

so grammar called left linear grammar
 $A, B \in V$ (LLG)

$B, \alpha, C \in T^*$ (α, B) can be - E also

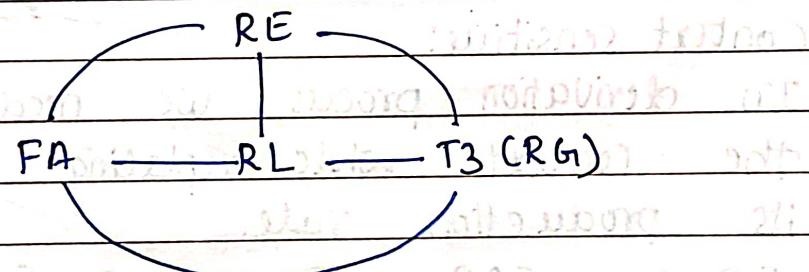
Eg: $A \rightarrow aB | a$] RLG
 $B \rightarrow ab | bBb | alb$

$A \rightarrow Bala$] LLG
 $B \rightarrow Ba | Bb | a1b$

$A \rightarrow Bala \rightarrow$ LLG Not Type 3.
 $B \rightarrow aB | a2 \rightarrow$ RLG

Note: Combination of LLG & RLG is not allowed.

Language generated by Type 3 grammar are called regular language and Regular language can be represented using Finite Automata.



Equivalent in Power.

Type 2: If all the productions in the grammar are of the form

$$A \rightarrow \alpha, A \in V$$

$$\alpha \in (VUT)^*, \text{ can be } \epsilon$$

Type 2 are called Context Free Grammar (CFG)

Regular Grammar are subpart of Context Free Grammar which can be either RLG or LLG (Type 2)

If a grammar is regular then its a context free grammar.

If a grammar is context free then it may not be regular.

Eg: $A \rightarrow \alpha A \beta \gamma \delta$

Deriving $A \rightarrow A \beta \gamma \delta \rightarrow \alpha \beta \gamma \delta$

$$\rightarrow aabb$$

Context Free vs context sensitive

Context Free: Allow a variable (non terminal) to be replaced by a corresponding production rule whenever it appears in a derivation process. The replacement occurs irrespective of context (what lies before or after the non terminal).

Eg: $A \rightarrow OAI$ - whenever A found replace it with OAI.

Context sensitive:

In derivation process we need to think about the context while replacing a variable by its production rule.

Eg: $CAB \rightarrow C OAI B$ (we can replace A by OAI whenever it is preceded by C and followed by B.)

Language generated by CFG are called Context Free Language

Machine which can be used as an acceptor for CFL is Push Down Automata (PDA)

Type 3 : RG \rightarrow FA \rightarrow RL

Type 2 : CFG \rightarrow PDA \rightarrow CFL

Type 1 : Context Sensitive Grammar - Restricted Grammar

$\alpha \rightarrow \beta$

α is $(V+T)^*$ \vee $(V+T)^*$ at least one variable

β is $(V+T)^*$, ϵ on LHS.

Restriction: Number of symbols in α should be less or equal to β .

$S \rightarrow AB$ | $AB \rightarrow AbBC$

$AB \rightarrow abc$ | $A \rightarrow bCA$

$B \rightarrow b$ | $B \rightarrow b$

$S \rightarrow E$ ('s' should not appear in RHS of the rule)

Context Sensitive grammar are used for context sensitive language and machine used is Linear bounded automata.

Type 0: Recursively Enumerable Grammar (Unrestricted)

$\alpha \rightarrow \beta$

α is $(V+T)^*$ \vee $(V+T)^*$

| One at least variable onto LHS

$B = (V+T)^* = \epsilon$

α - can not be null.

$Sab \rightarrow ba$

$A \rightarrow S$

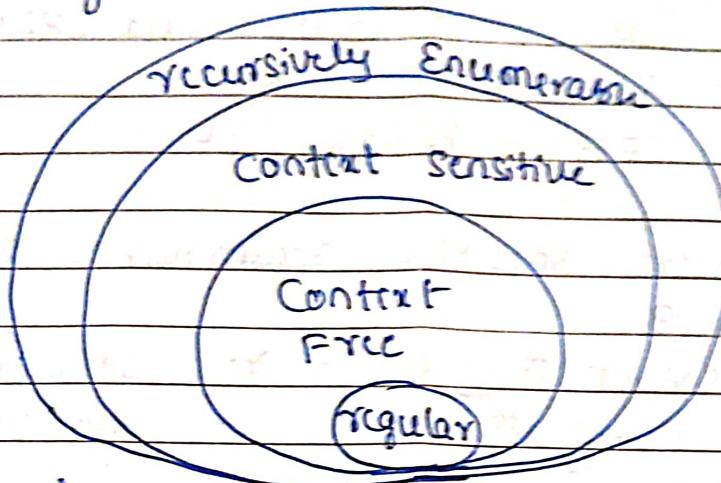
$S \rightarrow AcAB$

$Bc \rightarrow acB$

$CB \rightarrow DB$

$aD \rightarrow Db$

Recursively Enumerable Grammar for Recursively
Enumerable language & machine used is
Turing machine.



mc

Turing machine

Linear Bounded

PushDown Automata

Finite State

Grammar

Recursively Enumerable
Context sensitive

Context Free

Regular

Questions on RE. - Character class.

1. The first ten letters (upto j) in either upper or lower
[A-Ja-j]
2. LowerCase Consonants - [bcd{ghjklmnpqrstvwxyz]
3. The digits in a hexadecimal number (either upper
or lower case for digits above 9)
[0-9a-f]

4. The character that can appear at the end of
English sentence (exclamation, point)

[.?!]