

Compiler Design

Week 0

- Theories of Compiler Design used in NLP, Machine Translation, Compiler Information Retrieval, Semantic Web Technology, Query Language Processing.
- procedure for identifying syntax errors

"Compiler Basics"

Compiler is a translator that converts the high-level language into the machine language.

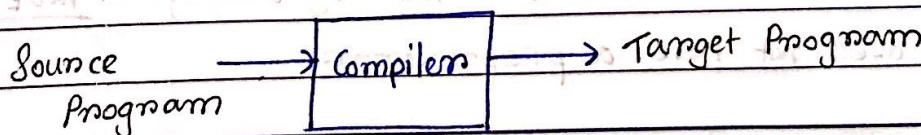
Compiler is used to show errors to the programmer.

A compiler is a computer program that translates computer code written in one programming language (the source language) into another language. (the target language.)

- Programming problems are easier to solve in high-level languages.
- Solutions are usually more efficient (faster, smaller) when written in Machine Language.

Compilers are the bridges Tools to translate programs written in HLL to efficient executable code.

Another important role of the compiler is to report any errors in the source program that it detects during the translation process. (compiler errors)



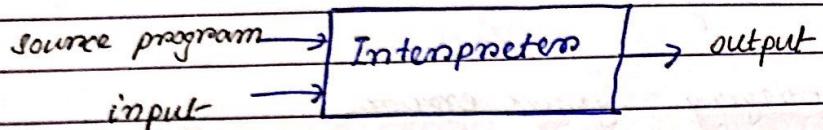
(Compiler is a Software program)

Q. What is an interpreter?

Ans:

An interpreter requires no pre-computation

it directly takes the source program and runs the input on it producing the desired output. (line by line / instruction by instruction)



Q. Difference betⁿ Interpreter and Compiler.

* The interpreter takes less memory (no .exe file needed)

* The interpreter is invoked each time the CPU has to run any piece of instruction (as no files.exe generated)

Compilers & interpreters do similar jobs, but there are differences:

To run a program you've written, it must first be translated into machine code so the computer can read it. This is what compilers and interpreters do.

However, Compilers convert the code all at once, save it, then run it. ↗ (C, C++) Java

Whereas, interpreters translate the code one line at a time, as it is run. ↗ (Python, Ruby) JavaScript

* Interpreter is slower

* The interpreter converts HLL to Assembly Language \Rightarrow FALSE

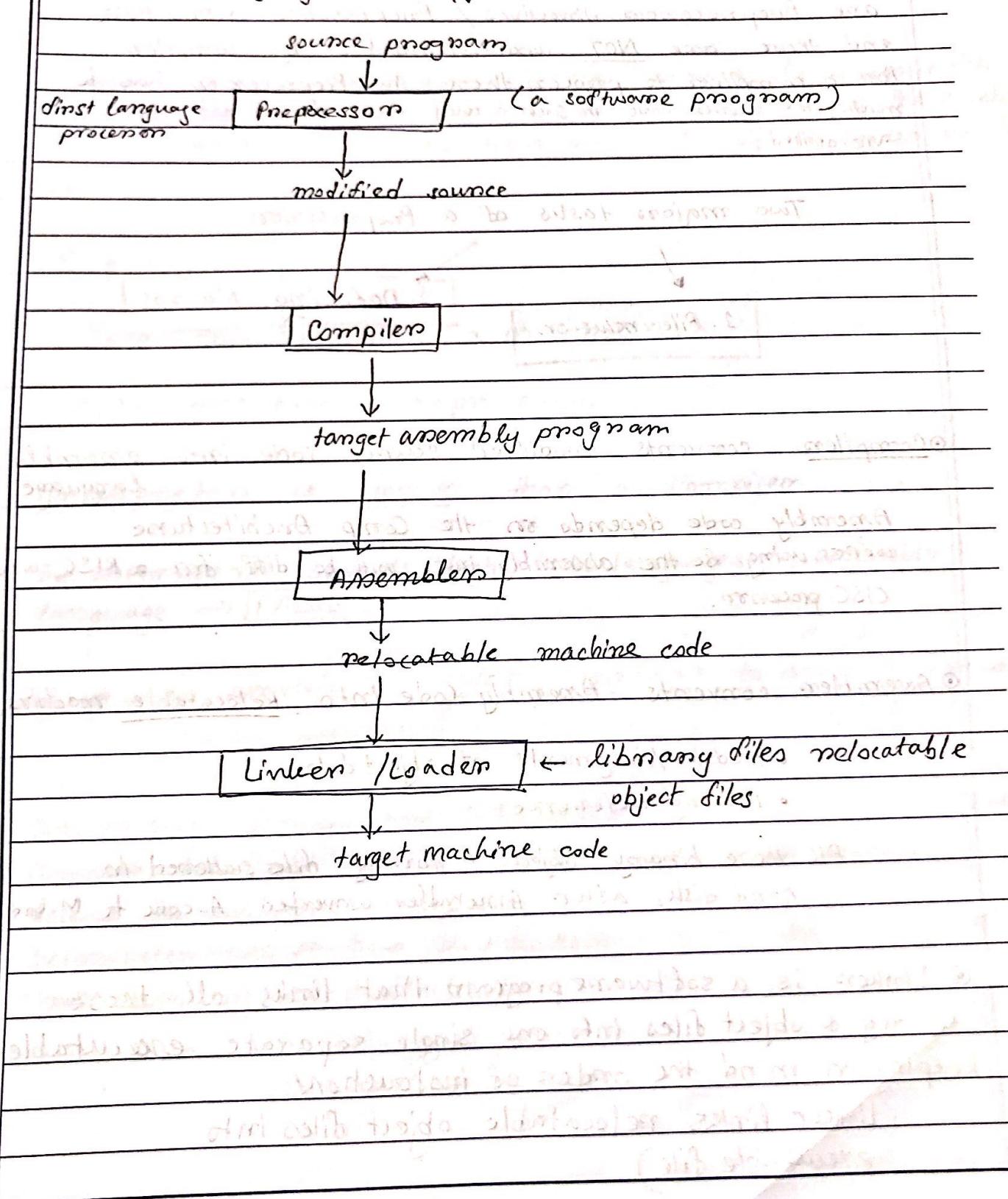
* The interpreter

* The compiler converts HLL to Assembly Language \Rightarrow TRUE

* Java does not have a preprocessor.

- A Language Processing System is a system that actually converts the source program written in one language into a target machine code consisting of 1s and 0s by incorporating different language processors such as preprocessors, compilers, assemblers, linkers and loaders.

Q Steps of a Language Processing System



(asks the Preprocessor to include the standard iostream headers file into our program)

Instructions preceded by a hash symbol are Preprocessor directives / Preprocessor instructions and these are NOT understandable by Compilers.

For a compiler to process these, the Preprocessor has to modify the source code in such a way that these instructions are omitted.

Two major tasks of a Preprocessor

1. File inclusion

2. Defining Macros

Compiler converts modified source code into assembly language.

Assembly code depends on the Comp. Architecture we're using. So the assembly inst. will be diff for a RISC and CISC processor.

Assembler converts Assembly code into Relocatable machine code.

- modules/fragments of object files
- library reference

All these binary object & library files scattered in hard disk, after Assembler converted. A. code to M. code

Linker is a software program that links all these library & object files into one single separate executable file. keeping in mind the order of instructions.

(Linker links relocatable object files into executable file)

General Summary: { * Compilers
* Compilers vs Interpreters
* Steps of LP system.

Q. Loader loads executable machine code into main memory on execution.

Loader will take the executable file produced by Linker in the hard disk and it will load them into the memory by mapping them into memory addresses, after that those segments will be loaded for execution.

Ques - 01

1. Java does not have a preprocessor.
2. Interpreter is slower than a Compiler.
3. The interpreter converts higher level languages to assembly language → **FALSE**
4. The compiler converts higher-level languages to assembly language. → **TRUE**

Compiler runs faster than Interpreter.

Compilers take less time to execute than interpreters because the run time is slower.

Interpreter doesn't have any compile time as the code never compiles, it directly runs.

"Phases of Compilers"

July 12th, 2020

Week 1, Lecture - 2

Phases of a Compiler (L8S1CT)

- Lexical Analysis
- Syntax Analysis
- Semantic Analysis
- Intermediate code generation
- Code optimization
- Code generation

built-in keyword → lexeme
(meaningful sequence
of characters)

Lexical Analyzer/scanner will find lexeme (variable, no., operator)
and generate tokens for the same.

This token stream fed as input into Syntax Analyser/Parsers.
which will generate a syntax tree based on a grammatical
structure. or CFG

Semantic Analyzer checks for inconsistency in meaning.
(array size can never be declared as a double)

ICG → (3 address code generated
called Intermediate Representation)

MICO's task → to reduce the no. of lines or to
reduce the no. of variables in 3-address code
to decrease Compilation Time
and increase Efficiency.

Code Generator → generates Assembly Instructions
which is hardware-specific.

MDCO → to reduce no. of lines + variables in assembly code.
(optional)

(machine-dependent
code optimizers)

Code Optimizer

④

Symbol Table: a table that contains ^ all the variables used in the software program for Compiler.

index/pointers to

It includes variables and also information regarding the variables such as the name, the type, the scope, the storage allocation for that variable.

All the phases of compiler will interact with Symbol Table to install or modify updates regarding the variables.

Symbol Table contains: • Variables • Procedures (names, no. and types of argument used)

- 1. Records the identifiers used in the source program
- 2. It's a data structure with collection of records.

#

Analysis - Synthesis model of compilation

(all phases of compiler broadly classified into 2 major categories)

Analysis

Synthesis

Front-end of compiler is written as it links the compiler to the software part of the machine

back-end of compiler as it links the compiler to the machine's hardware part

Breaks up the source program

into constituents → constructs the target program

[Scanners, Parsers, Semantic Analyzer]

[Code optimizer, Code generators]

Analysis Phase

Synthesis Phase

(Scan the entire input statement characters by characters to find a pattern)

contains Regular Expression
id → letter (letter / digit)*

Lexical Analysis divides the program text into lexemes or tokens.

- Tokens are the "words" of the programming language.

- Lexeme are meaningful sequences of characters from source character stream.

token \Rightarrow < token-name, attribute-value >

Lexical analysers also updates Symbol Table entries, removes the white spaces, removes the comments.

④ Syntax → understand sentence structure

⑤ Semantic analysis → performs type checking, looks for meaning and performs type conversion.

⑥ ICG • 3 address code which are assembly-like instructions

why? \Rightarrow (maxⁿ 3 addresses (variables/digits) per line)

- At most 1 operator on the RHS of equal sign
- All the 3 instructions are always executed line by line (in order)

⑦ Code Optimisation

Code Optimization → task is to reduce no. of lines and variables to improve the time and space efficiency of the generated code.

Peephole Global optimization (takes more

optimization (entire code) time)

(small fragment of code)

Compiler source will be first off compilation.

Q1

Code Generation

- takes as input the intermediate code representation
- generates assembly code that is hardware specific

Q2

Compilers

- * Each phase can encounter errors (detects & reports it along the way).

- lexical errors • syntax errors • semantic errors

- * Should be able to proceed & process the rest of the program after an error detected.

- * Should be able to link the errors with the source program.

Quiz - 02

1. Code Generation is NOT part of Analysis Part. False
2. Operators such as "-", "+", "=" get a token name. True
3. Structure of a sentence is checked during Syntax Analysis part.
4. Code Optimizer reduces the no. of lines and variables.

Lexical Analysis

July 12th, 2020

(Tokenize the text of the source program)

Lexical Analysis is the first phase of Compiler Design
Also called Scanner.

It removes the whitespace or comments.

Lexical Analyzer breaks the syntaxes into a series of tokens, by removing any whitespace or comments in the source code.

18 / C / x / 1 / 2 / = → assignment operator → assignment
 keyword ↓ operator no. → relational operator → relational operator
 identifier → relational operator → assignment operators

else → even / odd → relational operator → relational operator
 keyword → identifier → relational operator → assignment operators

Lexical analyzers also connected to Symbol Table which is a data structure. It inserts the Identifier & procedure/function tokens and the class object and instance tokens to the symbol table with their attributes.

Q) Lexical Analyzers

(key func's)

- Scan input
- remove White Space, New Lines, tabs
- identify tokens
- create Symbol Table
- insert tokens into ST
- generate errors
- send tokens to parser

- Parsen**

 - ✓ perform syntax analysis (validating syntaxes)
 - ✓ actions dictated by token order
 - ✓ update symbol table entries
 - ✓ create abstract representation of the source (using parse tree)
 - ✓ generate errors

- Q. What are major terms for Lexical Analysis?

- A **Token** → a classification for a common set of strings
example: <Identifier>, <number>, etc.

- Pattern** → the rules which characterize the set of strings for a token

- **Lexeme** → • actual sequence of characters that matches pattern and is classified by a token.
(Identifiers: x, count, name, etc.)
sum.

lexeme 1 2 3

token <number>

- No specific pattern for defining keywords (if, const)

(Each of the word has diff patterns, so no common common pattern
(if, const, do, while, etc) in general)

literal (token) - any seq. of characters or sentence

within the double quotation marks

Is called a literal.

Language & Regular Expressions.

* A Regular Expression is a set of Rules / Techniques for Constructing Sequences of Symbols (Strings) from an alphabet.

Let Σ be an alphabet, r a Regular Expression, then $L(r)$ is the Language that's characterized by the Rules of r .

Precedence is important in Parsing regular expressions

1. "*" has the highest precedence,
2. Concatenation as middle precedence
3. "|" has lowest precedence
4. Use parentheses to override these rules.

* Equality vs. Equivalence

We say two n.e. as equivalent, not equal !!

• Algebraic Law of R.E.

$$R(ST) = (RS)T = RST$$

1. Concatenation is associative

$$R(ST) = (RS)T = RST$$

2. "|" is associative

$$R | (S|T) = (R|S)|T = R | S | T$$

3. " " is commutative

$$R|S = S|R$$

4. Concatenation is distributive over |

$$R(S|T) = RS | RT$$

$$(R|S)T = RT | ST$$

5. ϵ is identity for concat $\epsilon R = R\epsilon = R$

6. * is idempotent $(R^*)^* = R^*$

7. Relation betn * (* and ϵ).

$$R^* = (R|\epsilon)^*$$

$$(\& \text{ word}) \rightarrow id = \mathcal{L}(\text{id})^*$$

important

identifiers = letter (letter | digit)*
ID

regular expression
of identifiers

letter = a | b | c | ... | z

For defining id's r.e.,
the first character
must be a letter

$x^+ = x(x^*)$

(\rightarrow one or more)

num \rightarrow digit⁺

(\rightarrow zero or one)

Q. Lexical analyzers scan away blanks, new lines, tabs
Can we define tokens for these?

ws \rightarrow delim⁺

whitespace consists one or
more delimiters.

blank \rightarrow blank

tab \rightarrow tab

newline \rightarrow newline

delim \rightarrow blank | tab | newline

Ans: No token is returned to Parser

Lexical Analyzers does NOT send any whitespace token to
the Parser, it just scans away blanks, new lines, tabs

Automata is an idealized machine to recognize / identify patterns.

Finite State Automata

finite numbers of states

FA is a program / recognizer that takes a string as an input and determines whether it's a sentence of the language or not.

one start state
many final states

R.E. defines the pattern of a sentence on a string.

each state labelled with a state name

directed edges (labelled with symbols)

→ DFA (deterministic)

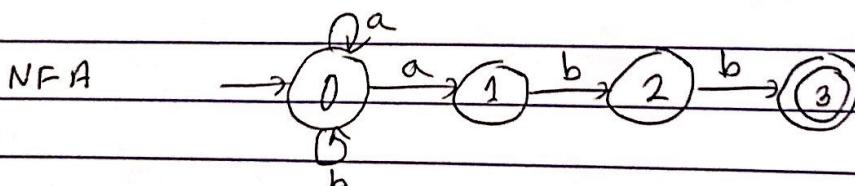
two types → NFA (non-deterministic)

The regular expression is compiled into a recognizer by constructing a generalized transition diagram called a finite Automaton.

A DFA is an NFA with the following restrictions:

• ϵ moves are NOT allowed. • for every state, there's one and only one path from an input symbol

$(a|b)^*abb$



DFA

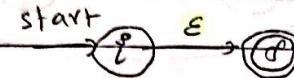
draw later

(NFA → we can move to two or more states with same input symbol and it accepts ϵ)

Both DFA and NFA are "recognizers" of regular sets.
But time-space trade space exists.

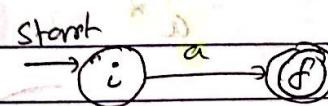
DFA's are faster recognizers (can be much bigger too)
drawback
so more storage space

Converting RE to NFAs.
(Thompson's Construction)



rule 1 • ϵ star empty string is a n.e. denoting $\{\epsilon\}$

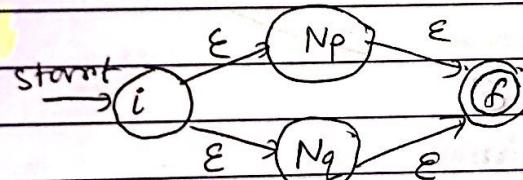
rule 2 • a is a n.e. denoting $\{a\}$ for any a in Σ



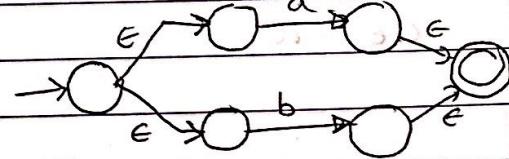
If P and Q are n.e. with NFAs N_P, N_Q :

rule 3 #

P | Q (union)



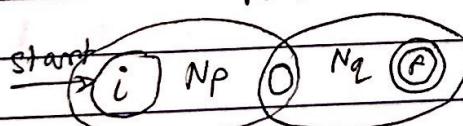
example
a|b



rule 4 #

PQ (concatenation)

P.Q

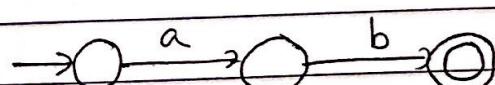


✓ the last symbol of N_p is state
first symbol of N_q state

a.b

example

a b

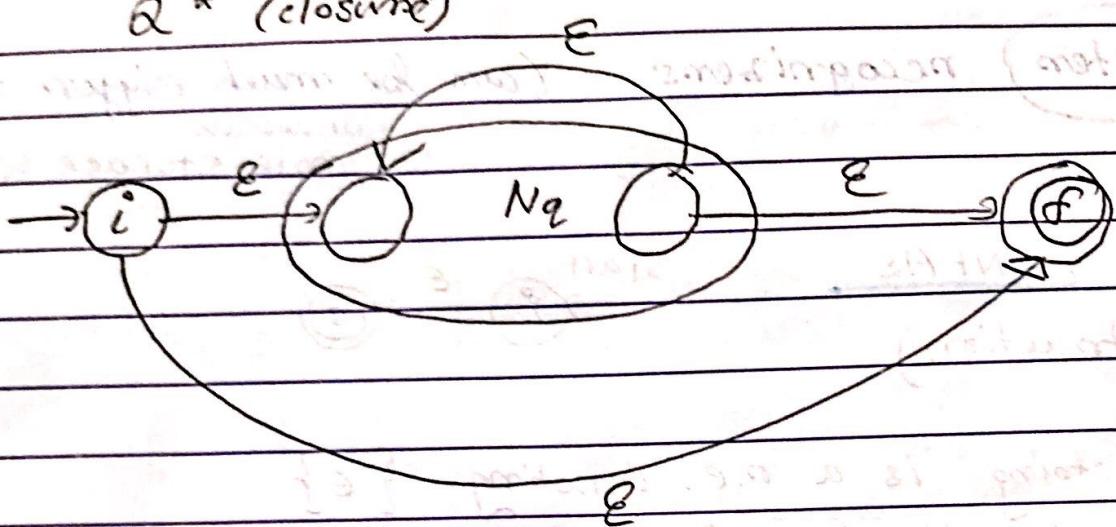


last state of a
first state of b

rule 5

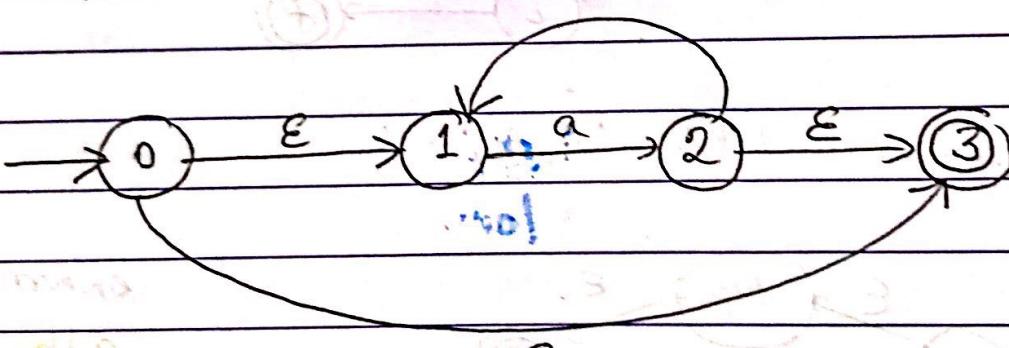
If Q is a n.e. with an NFA N_Q then Q^* has state ϵ which is start state and final state.

Q^* (closure)



example

a^*



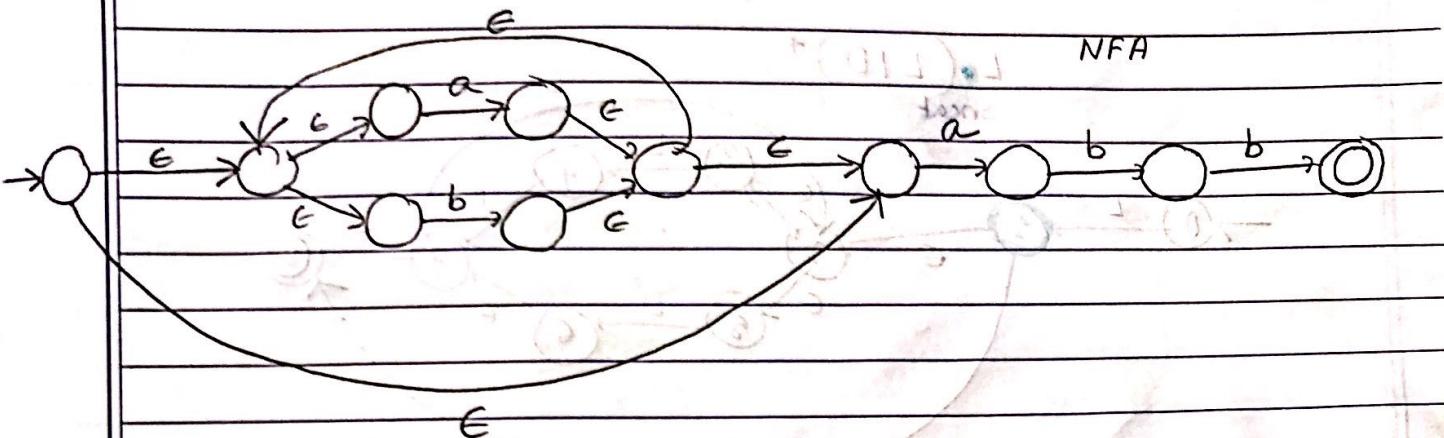
$$a\epsilon = a$$

#

Example 1

 $(a|b)^*abb$

Using Thompson's C. rule



H.W.

#

Example 2

 $(ab^*c) \mid (a(b)c^*)$

* → highest

• → middle

| → lowest

Precedence

of operators

Precedence of Operators

* → highest

• → medium

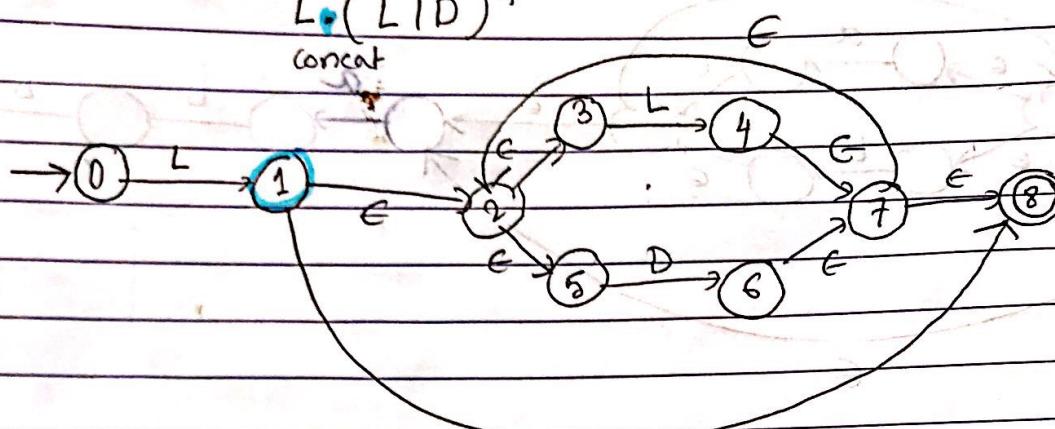
| → lowest

Example
3

id \rightarrow letter (letter | digit)*

$L \bullet (L \cup D)^*$

concat



Quiz - 03

1. Lexical Analysis does not generate errors \rightarrow False
2. Token and lexeme mean the same thing \rightarrow False
3. We cannot use decimal point in RE as we only use digits \rightarrow False

Converting RE to DFA

RE \rightarrow NFA \rightarrow DFA

(subset construction in 331 Automata)

• direct method RE \rightarrow DFA

(don't have to create NFA)

First, augment the R.E.

adding # at the end

$r \rightarrow (r) \#$

Then, we create syntax tree from this augmented r.e.

Input: A regular expression r Output: A DFA D that recognizes $L(r)$ (the language defined by r.e. r)

Syntax Tree

- leaves
- internal nodes

(alphabet or ϵ)

also

(operators)

• , *, |

concatenation star-node, or-node
cat-node $(a/b)^* a bb$

* Add # at end of r.e. to augment it.

$((a/b)^* a bb) \#$ To construct ST

1 2 3 4 5

• Consider precedence of operators

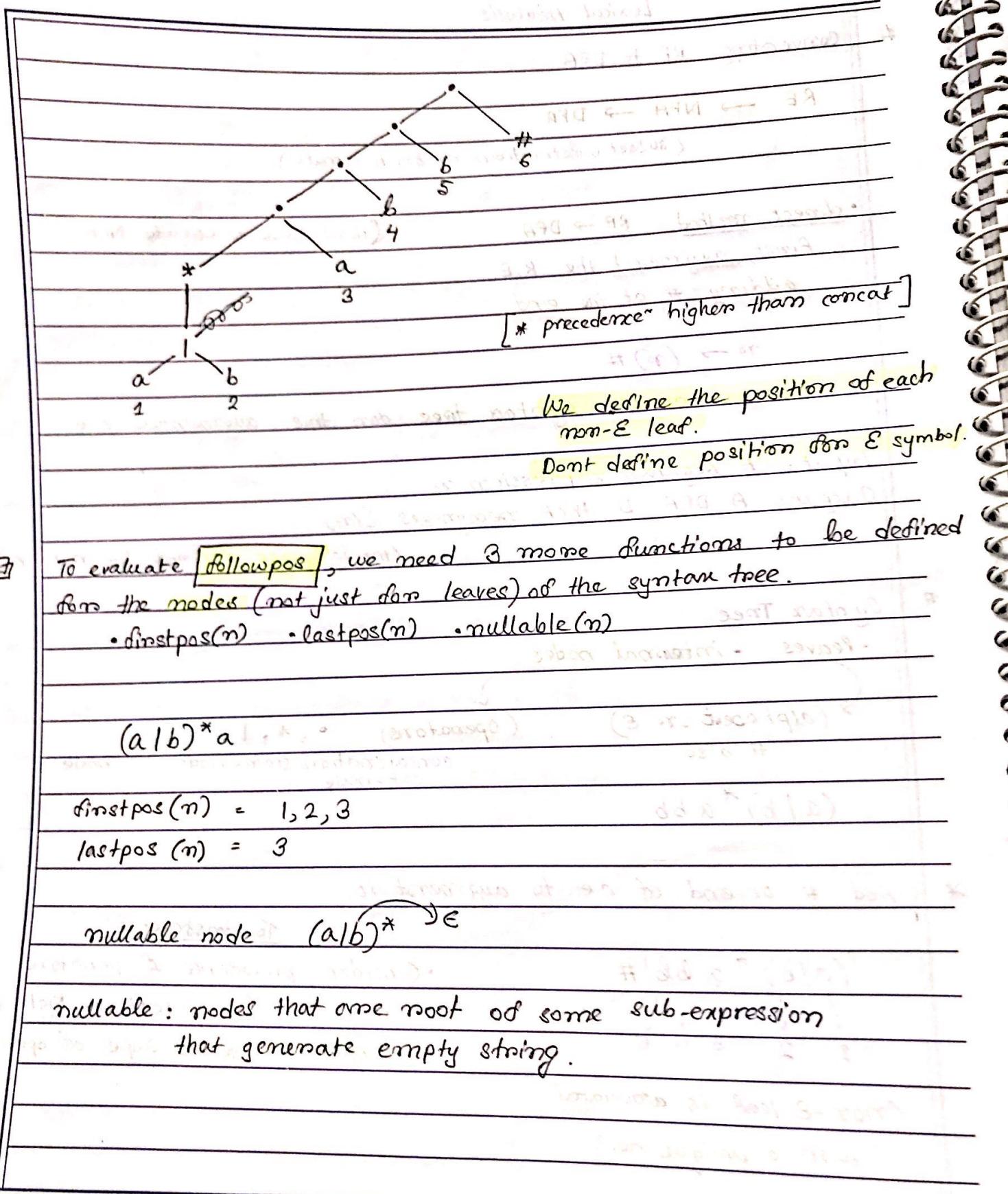
• Scan r.e. from Left to Right (L to R)

• Parentheses define scope of operation

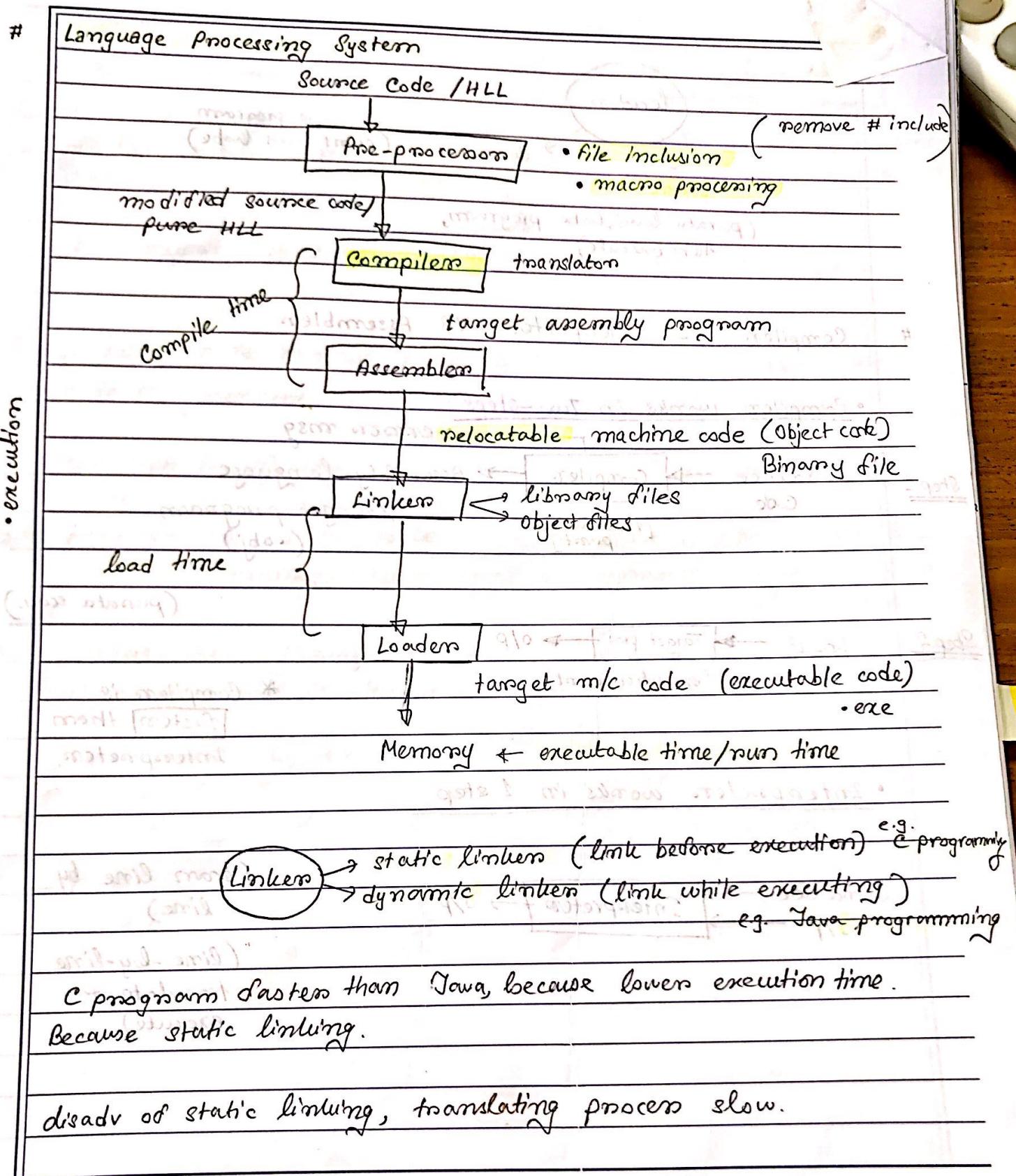
(non- ϵ leaf is associated

with a unique no.)

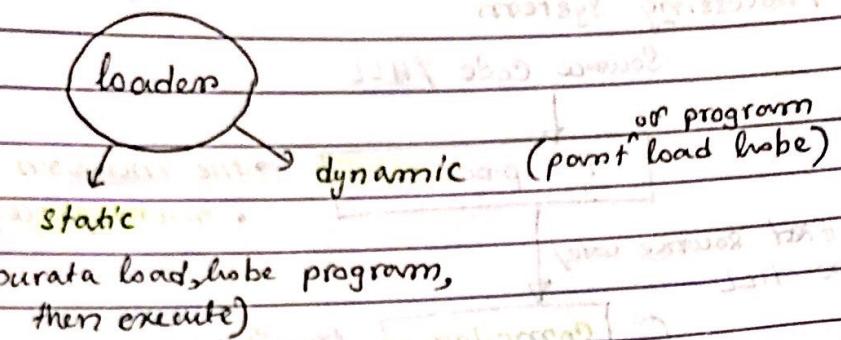
abcd, ϵ , #
numbers these



Revision of Lecture 1

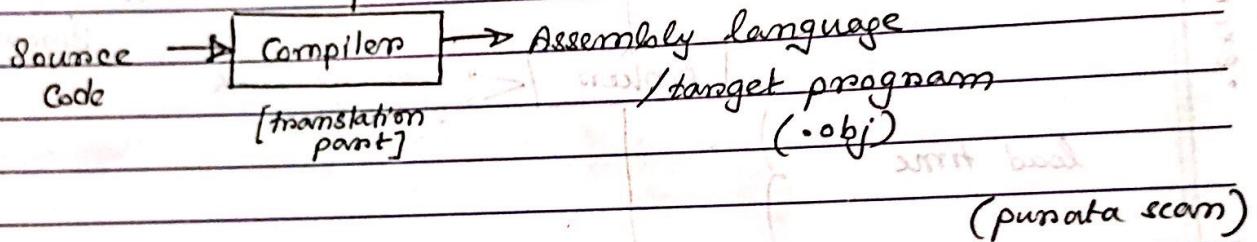


1. Approach To Executing A



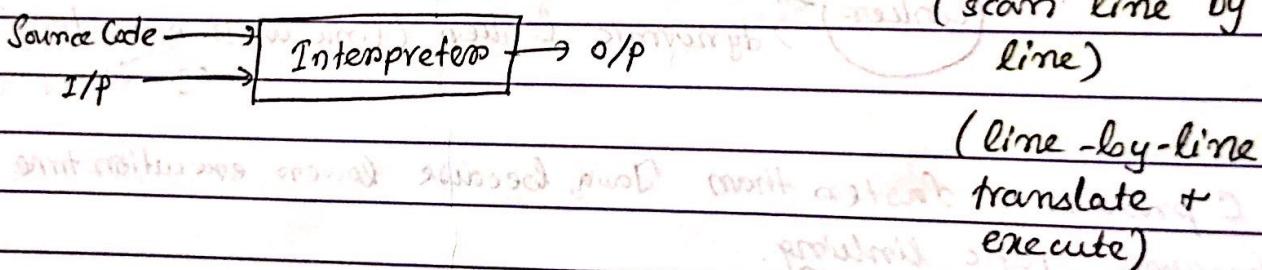
Compiler VS Interpreter VS Assembler

- Compiler works in Two steps



* Compiler is faster than Interpreter.

- Interpreter works in 1 step



Note: Compiler generates machine code, while interpreter translates & executes directly.

Q. Differences between Compiler & Interpreter.

Compiler

Interpreter

1. Compiler scans the whole program and translate into machine code.
1. Interpreter scans one line at a time and translate into machine code.
2. Intermediate object is generated.
2. Intermediate object code is not generated.
3. It takes more memory since object code is generated.
3. It takes less memory.
4. Compiler is faster.
4. Interpreter is slower.
5. Program need to be compiled everytime.
5. Every time higher level program is converted into lower level program.
6. Errors are displayed after entire program is checked.
6. Errors are displayed for every instruction interpreted.
For example: C, C++
so efficient error diagnosis.

for example: Basic,
Python

The output of assembly language to machine code.



Java is a compiled language.

Python is an interpreted language.

High-level language - If a program contains `#define` or `#include` directives such as `#include` or `#define` it is called HLL. They are closer to humans but far from machines.

The (#) tags are called preprocessor directives.

They direct the preprocessor about what to do.

Preprocessor - The preprocessor removes all the `#include` directives by including the files called file inclusion and all the `#define` directives using macro expansion.

Assembly Language - It's neither in binary form nor high level. It is an intermediate state that is a combination of machine instructions and some other useful data needed for execution.

Compilers and Interpreters do similar jobs, but they're different in the way they read input.

Compiler reads the inputs in one go,

The Compiler in one go reads the inputs, does the processing and executes the source code whereas the Interpreter does the same line by line.

Compiler scans the entire program and translates it as a whole into machine code whereas an Interpreter translates the program one statement at a time.

Compilers FASTER than Interpreters.

Linker loads a variety of object files into a single file to make it executable.

The loader loads it in memory and executes it.

- Frontend is platform-independent

L8S1/CT

Frontend, Backend/Synthesis
Analysers

Instead of producing a target program as a translation, an interpreter appears to directly execute the operations specified in the source program on inputs supplied by the user.

The machine-language target program produced by a compiler is usually much faster than an interpreter at mapping inputs to outputs.

Interpreter gives better error diagnostics than a compiler, because it executes the source program statement by statement.

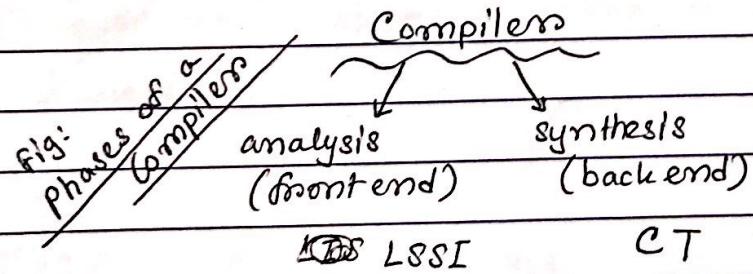
Assembly language is easier to produce as output and is easier to debug.

Compiler produces assembly-language program as its output.

The assembly language is then processed by a program called an Assembler that produces relocatable machine code as its output.

Linker resolves external memory addresses, where the code in one file may refer to a location in another file.

Loader puts together all of the executable object files into memory for execution.



- Compilation process operates as a sequence of phases.

L8S1 CT

Symbol Table stores info about the entire source program, is used by all phases of the compiler.

Symbol Table helps the compiler to function smoothly by finding the identifiers quickly.

Phases of a Compiler

Analysis Phase

An intermediate representation is created from the given source code.

Synthesis Phase

Equivalent target program is created from the intermediate representation.

1. Lexical Analyzer
2. Syntax Analyzer
3. Semantic Analyzer
4. Intermediate Code Generator
5. Code Optimizer
6. Code Generator

The analysis phase consists of three main steps: lexical analysis, syntax analysis, and semantic analysis. These steps are typically performed sequentially, although some compilers may use a more parallelized approach.

The synthesis phase consists of two main steps: code generation and optimization. The goal of the synthesis phase is to produce efficient and portable target code.

The analysis phase is responsible for determining the meaning of the source code, while the synthesis phase is responsible for generating the target code.

The analysis phase is typically implemented using parser generators like YACC or ANTLR, while the synthesis phase is often implemented using code generation libraries like LLVM or Clang.

The analysis phase is responsible for identifying errors in the source code, while the synthesis phase is responsible for generating error messages and performing code optimization.

The analysis phase is typically implemented using parser generators like YACC or ANTLR, while the synthesis phase is often implemented using code generation libraries like LLVM or Clang.

- Q. Why does a lexical analyzer use buffers rather than input?
- Q. Lexical error → initializing a variable name starting with a number.
- Q. DFA is a state pattern recognizer.
- Q. Is string 'abbc' accepted by n.e. "a b | b*c" \Rightarrow NO
- Q. Syntax Analyzer deals with the grammatical structure of a programming language.
- Line 1. and 2 (both)
 - Compiler phases are broadly classified into analysis & synthesis phases.
 - One buffer cannot hold all input characters, it needs two.
 - One buffer has less capacity than the other.
 - To increase space complexity

Q. C++ does not use an interpreter.

Q. Syntax error → int[] arr = new int [3.5]

Q. Syntax → missing ;

Q. True → $a a^* (b|c) = a^* a (c|b)$

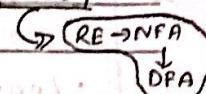
Q. Which of the following is a n.e. of all the binary numbers divisible by 8

Q. $\Sigma \{0,1\}$ is described by n.e. $(0+1)^* 0(0+1)^* 0(0+1)^*$

Set of all strings containing atleast two 0's

Converting RE to DFA directly

331



RE → DFA (direct method)

Step 1: RE augmentation

(adding additional symbol hash #)

 $r \rightarrow (r)\#$ this symbol is concatenated with whole r.e.

Step 2:

Then, we create a syntax tree for this augmented regular expression.

#

Construction of DFA from RE

1.

Input: A regular expression r .

2.

Output: A DFA D that recognizes $L(r)$. (the language defined by r.e. r)

Method:

① Construct Syntax Tree from augmented RE $(r)\#$

② Construct the nullables, firstpos, lastpos and followpos

③ Construct D states: set of states of D D trans: transition table for D



- Q. Intermediate Code Generator produces machine-dependent intermediate code → FALSE
- Q. Interpreter better does error diagnosis.
- A. Preprocessor is used to expand macros into Source Language statements.
- Q. $R^+ = (R^*)^*$ valid
- $R^* = (R \mid \epsilon)^*$ valid
- $R^+ = R + UE$ valid
- ~~E.R.E = E~~ invalid

Q. Set of all strings over $\Sigma = \{0, 1\}$ that begins & ends with 0 is:

$$0(0+1)^+0$$

[3.8]

(a)

$$((I+D)D(I+D))^* (I+D)$$

→ 10 part forming pattern can't be 00

Ques

Draw

Syntax Tree has two types of nodes.

- Leaves contain: Alphabet symbols or ϵ , # symbol
(non- ϵ leaf associated with unique position no.)

- Internal nodes contain: Operators.
(cat-node, on-node, star-node)

$((a|b)^*abb) \#$

1

2

3

4

5

6

unique position no.

- When constructing Syntax Tree, consider Precedence of the Operations.

- Scan the n.e. from Left to Right

- Parentheses (brackets) define the scope of the operation.

[* precedence higher than concat]

[define the

position of each of
the non- ϵ leaf.

Don't define position
form ϵ -symbol]

a
1
b
2

fig: Syntax Tree

for $((a|b)^*abb) \#$

Nullable: nodes that are the root of some sub-expression that generate empty string.

If n is a leaf labeled by ϵ then

$$\text{nullable}(n) = \text{true}$$

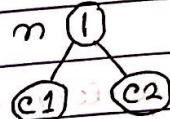
(leaf node)

If n is a leaf labeled with position i

$$\text{nullable}(n) = \text{false}$$

If n is an or-node (\sqcup) with children c_1 and c_2 .

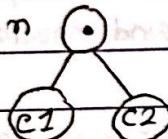
$$\text{nullable}(n) = \text{nullable}(c_1) \text{ or } \text{nullable}(c_2)$$



If n is a cat-node (\bullet) with children c_1 and c_2 .

$$\text{nullable}(n) = \text{nullable}(c_1) \text{ and } \text{nullable}(c_2)$$

(both)



If n is a star-node (*) with children c_1 .

$$\text{nullable}(n) = \text{true}$$

n (*) nullable because it
applies to strings combining both it shows it will produce the
empty string.

$$*(\text{dis})$$

Rules for firstpos (See slide) }

Lecture 1, Part 2

Rules for lastpos (See slide) }

Summary table

- ④ To evaluate followpos, we need 3 more functions to be defined for the nodes (not just for leaves) of the Syntax Tree.
- **firstpos(n)** - the set of the positions of the first symbols of strings generated by the sub-expression rooted by n.
 - **lastpos(n)** - the set of the positions of the last symbols of strings generated by the sub-expression rooted by n.
 - **nullable(n)** - true if the empty string is a member of strings generated by the sub-expression rooted by n false otherwise.

To determine followpos, we need to determine firstpos, lastpos and nullable.

$(a|b)^* a$

$\hookrightarrow \{aa, ab, ba, aa\}$
 $\hookrightarrow \{aa, ba, \epsilon a\}$

firstpos : 1 on 2 on 3

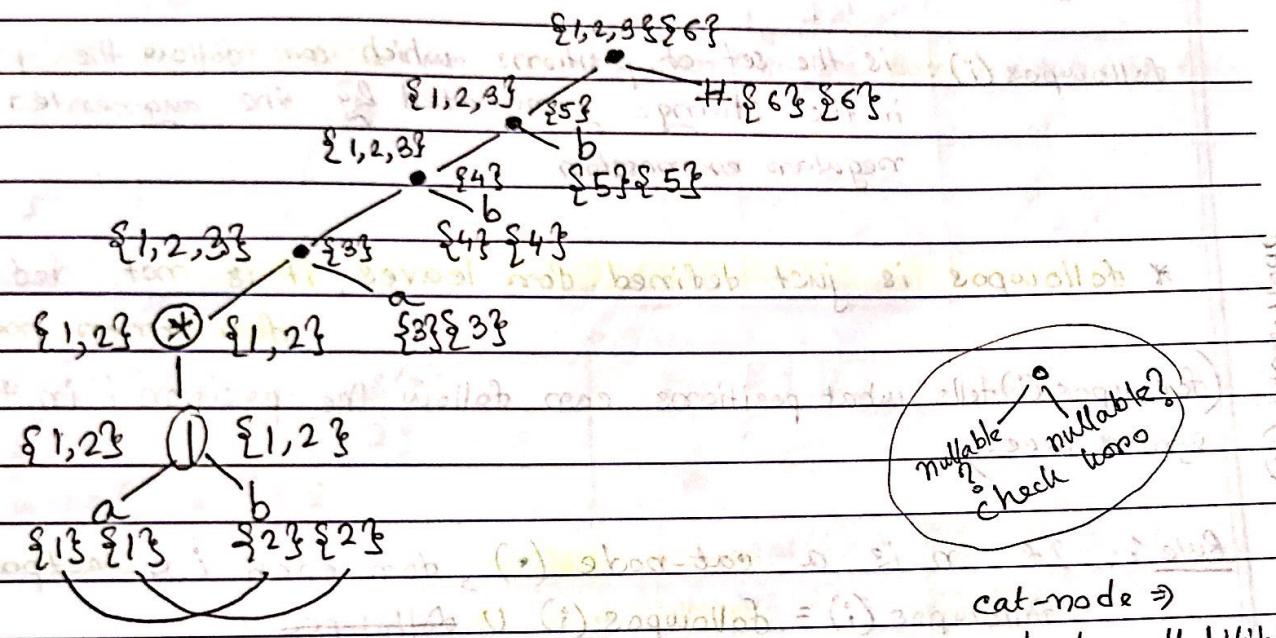
lastpos : 3

nullable node is the node n that produces empty string.

nullable node : $(a|b)^* a$

firstpos and lastpos example (Lecture 1, Part 3)

- first, find leaves $\Rightarrow \text{firstpos}()$ and $\text{lastpos}()$ ← 38
- Traverse the tree in a bottom-up approach.



cat-node \Rightarrow
check nullability
and
 c_1, c_2

* Lecture 2, Part 1 *

Lexical Analysis

RE \rightarrow DFA (direct method)

followpos function returns the possible next positions of the designated alphabet of n.e.

followpos(i): is the set of positions which can follow the position i in the strings generated by the augmented regular expression.

* **followpos** is just defined for leaves, it is not defined for inner nodes.

(**followpos(i)** tells what positions can follow the position i in the syntax tree.)

Rule 1: If n is a cat-node (\cdot), for each $i \in \text{lastpos}(c_1)$

$$\text{followpos}(i) = \text{followpos}(i) \cup \text{followpos}_{\text{firstpos}}(c_2)$$

Rule 2 If n is a (*) star-node, for each $i \in \text{lastpos}(n)$

$$\text{followpos}(i) = \text{followpos}(i) \cup \text{firstpos}(n)$$

- After computing firstpos and lastpos for each node, followpos of each position can be computed by making depth-first traversal of the Syntax Tree.

After constructing ST, determine firstpos and lastpos of leaf nodes and interior nodes.

We'll determine followpos only for (*) start-node and (•) cat-node.

Traverse the ST using depth-first traversal technique.

on-node doesn't define any next position so we don't consider this for followpos.

Start node

$$i \in \{1, 2\}, \text{firstpos}(*) = \{1, 2\}$$

$$\text{followpos}(1) = \{1, 2\}$$

$$\text{followpos}(2) = \{1, 2\}$$

Catnode

$$i \in \{1, 2\}, \text{firstpos}(c2) = \{3\}$$

$$\text{followpos}(1) = \{1, 2, 3\}$$

$$\text{followpos}(2) = \{1, 2, 3\}$$

$\{1, 2, 3, 4, 5\}$

$\{1, 2, 3, 4, 5\}$

$A = \{1, 2, 3, 4, 5\}$

$i \in \{3\}, \text{firstpos}(b) = \{4\}$

$$\text{followpos}(3) = \{4\}$$

so $\{1, 2, 3\}$

$\{1, 2, 3\}$

$i \in \{3\}, \text{firstpos}(b) = \{4\}$

$$\text{followpos}(3) = \{4\}$$

for $i \in \{4\}, \text{firstpos}(b) = \{5\}$

$$\text{followpos}(4) = \{5\}$$

$$\text{followpos}(5) = \{6\}$$

start tool 2 - #

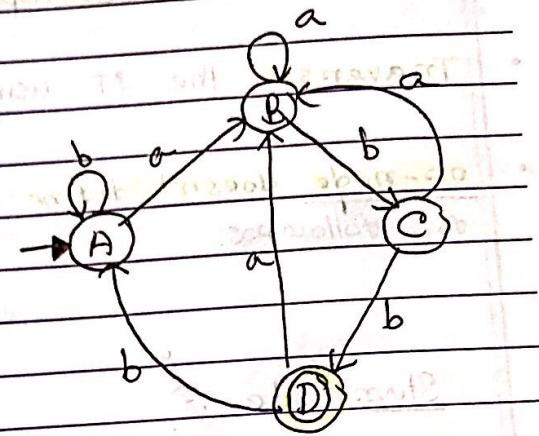
use these two to construct DFA

- ST with firstpos and lastpos
- followpos Table

In DFA construction procedure, we start with firstpos of root.

root:

Dstate	a	b
$\{1, 2, 3\} = \emptyset$	$\{1, 2, 3, 4\}$	$\{1, 2, 3\}$
$\equiv B$	$\equiv A$	
$\{1, 2, 3, 4\}$	$\{1, 3\}$	$\{2, 4\}$
B	C	D
$\{1, 2\}$	$\{1, 2\}$	$\{2, 3\}$
$\{2\}$	$\{2\}$	$\{3\}$
$\{1, 2, 3\}$	$\{1, 2, 3\}$	$\{1, 2, 3\}$

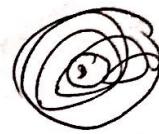


root
firstpos

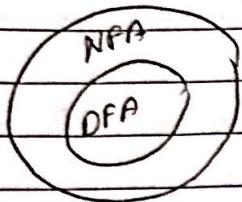
Dstate	a	b
$\{1, 2, 3\} = \emptyset$	$\{1, 2, 3, 4\} = B$	$\{1, 2, 3\} = A$
$B = \{1, 2, 3, 4\}$	$\{1, 2, 3, 4\} = B$	$\{1, 2, 3, 5\} = C$
$C = \{1, 2, 3, 5\}$	$\{1, 2, 3, 4\} = B$	$\{1, 2, 3, 6\} = D$
$D = \{1, 2, 3, 6\}$	$\{1, 2, 3, 4\} = B$	$\{1, 2, 3\} = A$

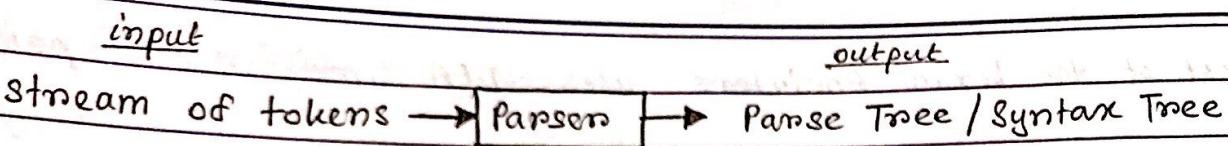
- 6 last state

DFA
recognize pattern generated by RegEx.



Most of the Lexical Analyzers uses DFA simulation for pattern/token recognition.





operators (internal node)	
operand (leaf node)	

lexical error → misspelling of keywords (printf, return)

syntax error → ",) missing

Parsing

Week 4, Lecture - 01
(SRF)

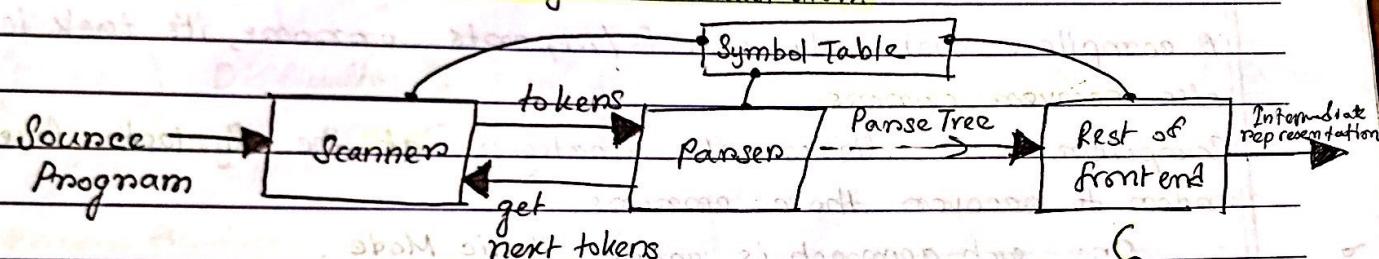
Parsing / Syntax Analysis

When a lexical analyzer processes these lexemes into tokens, these lexemes are then taken as inputs by Parser.

The Syntax Analyser / Parser will apply a grammatical structure on those tokens in order to generate a Syntax Tree. This ST follows a basic CFG.

The ST is passed to the Semantic Analyser which will check for any meaningful inconsistency within the tree.

Syntax tree is used to guide translation.



Scanners & Parser work simultaneously side by side in order to generate a

Parse Tree from a given stream of tokens.

When a parser takes a stream of tokens as inputs and generates a Parse Tree, it causes

the func to get next tokens to the Scanners in order to generate more tokens so as to generate more Syntax Trees.

- collecting token info
- type checking

- Intermediate code generation

Both Scanners & Parser take help from Symbol Table, by updating and modifying the values and info of variables.

bottom up parser -

goal

bottom up parser -

against above slide point (where?)

diff. approach to semantcs

1. bottom up

Lexical error - misspelt keyword / typos

Syntactic error - extra hyphens or extra closing bracket.

Semantic error - type errors, undefined IDs.

Logical errors cannot be handled by Compiler, rather these are handled by the Programmers themselves.

Compiler cannot detect logical Errors.

Most of the error-detection happens during Syntactic Analysis level. Adequate error reporting is a major task to do. Inadequate recovery may introduce more errors.

A compiler not only detects/reports errors, its task is to also recover errors.

Compiler uses different strategies at the Syntactic level in order to recover those errors.

One such approach is called Panic Mode.

Panic Mode Recovery Approach discards tokens until it finds a 'synchronizing' token.

When Syntax Analyzer finds one or a set of tokens that doesn't fall into a grammatical structure (CFG), the Syntax Analyzer will keep on deleting those tokens one by one until it finds a synchronizing token.

Synchronizing token could mean the end of something.

It could mean the end of a class, end of a func, end of a statement and so on

? end ;

- Adv:
- simple to implement
 - does not go into infinite loop
 - commonly used method

- Disadv:
- may skip over large sections of source with some errors.

2#

Phrase Level Recovery

works by inserting, deleting, or modifying tokens within a code segment.

Compiler connects the program by deleting or inserting tokens so it can proceed to parse from where it was.

(Adv): does not get into infinite loop

... constantly inserting tokens and never scanning the actual source.

• Generally used for error-repairing compilers.

(Disadv):

Difficulty: Point of error detection might be much later than the point of occurrence.

3#

Error productions

• Augment the CFG with "Error Productions"; now the CFG accepts anything.

Used with:

- LR (Bottom-up parsing)
- Parser Generators.

Error productions are prod. rules that are written within a CFG.

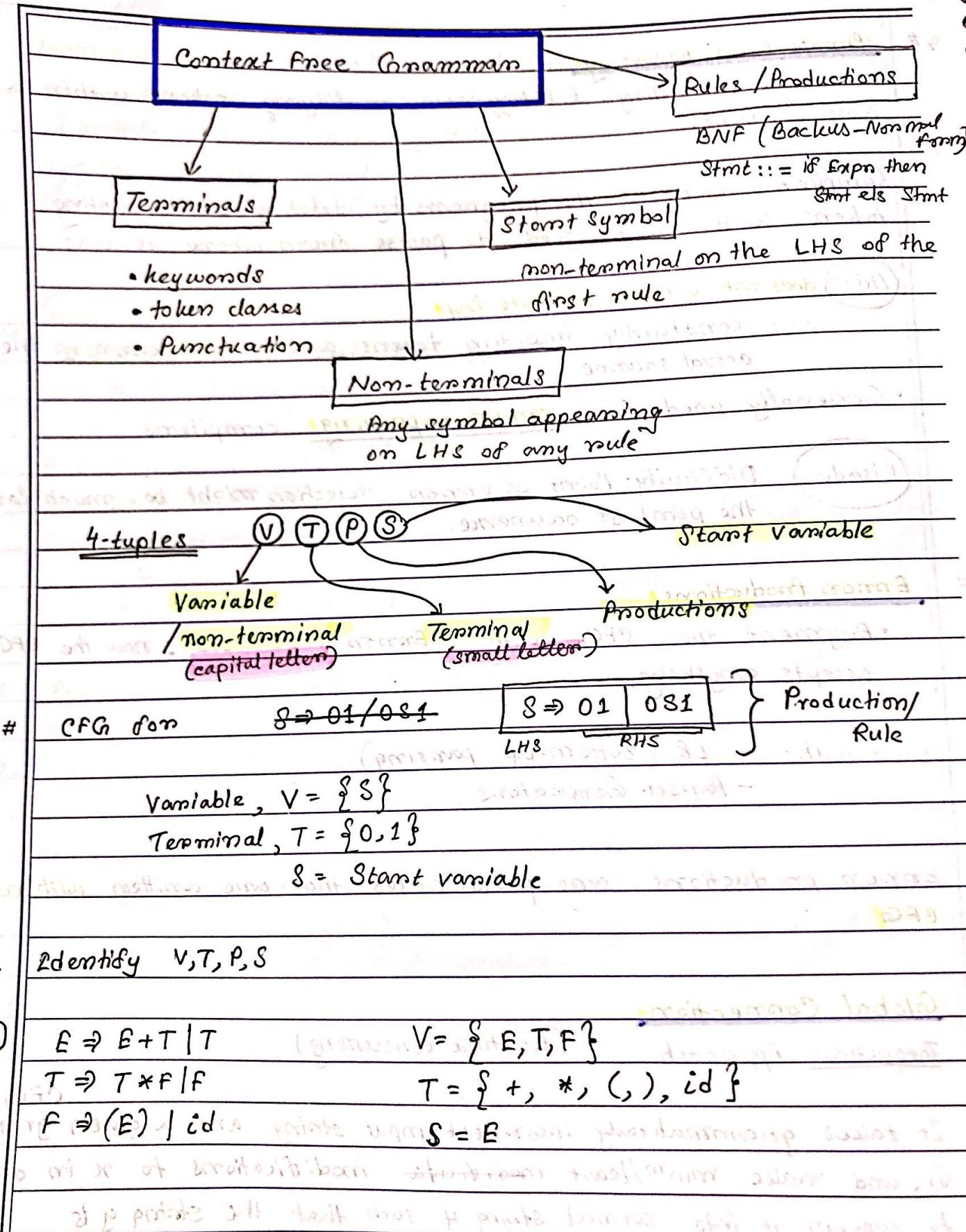
4#

Global Correction

Theoretical Approach (too time consuming)

It takes grammatically incorrect input string and a given grammar G_1 , and make min^m/least modifications to x in order to convert it into correct string y such that the string y is acceptable by grammar G_1 .

- A CFG is a formal model consisting of 4 units



A CFG is a set of rules /productions used to generate patterns of strings.

- 2

$$S \Rightarrow (L) \mid a$$
$$L \Rightarrow L, S \mid S$$

$$V = \{S, L\}$$

$$T = \{ (,), a, , \}$$

$$S = S$$

- # Equivalence of CFGs: If two CFG's generate the same language, we say they are "equivalent".
 $G_1 \approx G_2$ whenever

$$L(G_1) = L(G_2)$$

- # In a derivation,

for Leftmost Derivation \rightarrow always expand the leftmost non-terminal (L to R)

for Rightmost Derivation \rightarrow always expand the rightmost non-terminal / variable (R to L)

- #

Parse Tree is another representation of Left Leftmost and Rightmost Derivations. Both derivations will have the same Parse Tree until and unless the grammar is ambiguous.

(The parser tree only remembers which rule to use in replacing it) ~~④~~ (See slide)

Please Note:

↓ X | ✓

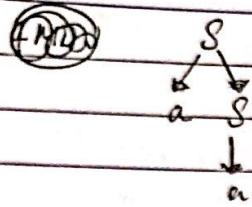
w = string
G = grammar

G.

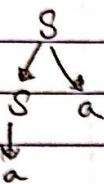
Find whether the given grammar is Ambiguous or not

①

$$S \rightarrow aS / Sa / a$$



(w = aa)



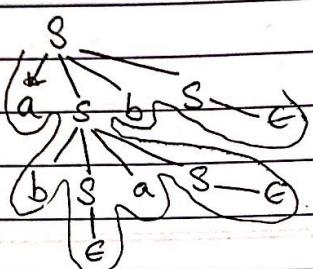
(for some G and w)
Two Parsee Trees, so
Ambiguous.

②

$$S \rightarrow aSbS \mid bSaS \mid \epsilon$$

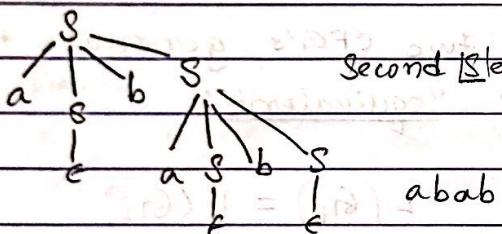
w = aabb

w = abab



abbaab
= abab

First S expand

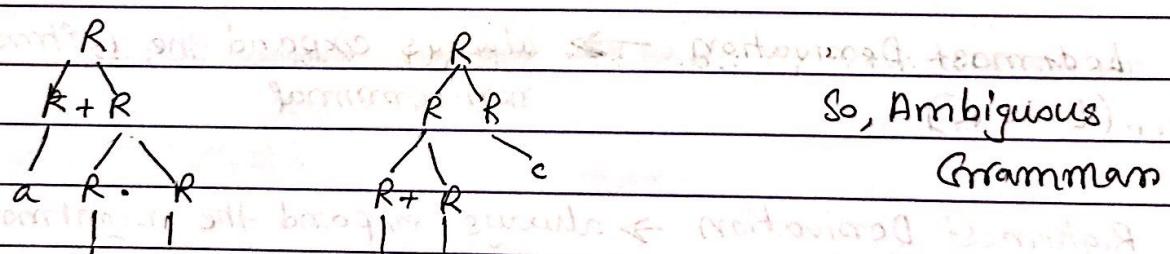


Second S expand
So, Ambiguous.

③

$$R \rightarrow R + R \mid RR \mid R^* \mid a \mid b \mid c$$

w = a+b+c



So, Ambiguous

Grammar

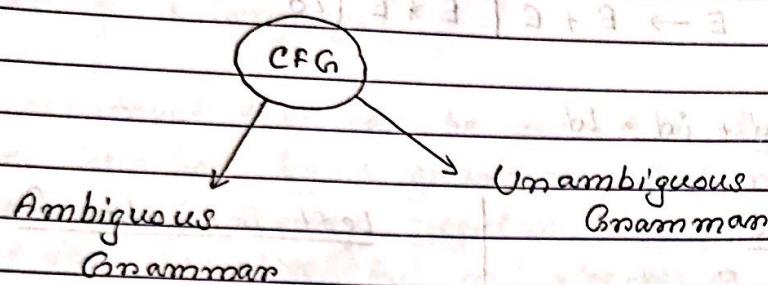
(Parser always matches terminals from left to right.)

Ambiguous Grammar (TNR)

Week 4, Lecture 2

7th August, 2020

CFG is divided into two parts based on the number of derived trees.



A grammar can have more than one parse tree generating a given string of terminals.

Such a grammar is said to be ambiguous.

Ambiguous Grammars produce either one or more or all of the below:

- more than one Parse Tree
- more than one Derivation Tree
- more than one Syntax Tree
- more than one Leftmost Derivation
- more than one Rightmost Derivation.

Most of the parsers, except one parser, will not allow ambiguous grammars.

(a * b*)

A sequence of terminals and nonterminals in
 a derivation is called "Sentential forms"

(derivation means - we're to generate a string consisting of only
 non-terminals or terminals from the body of the production rules)

Example

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

Sentence: $id + id * id$

Left Derivation 1

$$\begin{aligned} E &\rightarrow E + E \\ &\rightarrow id + E \\ &\rightarrow id + E * E \\ &\rightarrow id + id * E \\ &\rightarrow id + id * id \end{aligned}$$

Left Derivation 2

$$\begin{aligned} E &\rightarrow E * E \\ &\rightarrow E + E * E \\ &\rightarrow id + E * E \\ &\rightarrow id + id * E \\ &\rightarrow id + id * id \end{aligned}$$

different

Two L.D. while generating the same sentence of terminals.
 So by definition, this grammar is Ambiguous.

- Q. How to remove ambiguity in a grammar so as to make it Unambiguous Grammar.

Notice: all these lines that you generate in / get in the way of generating the string of terminals is known as Sentential Forms of the derivation process.]
 (id * E)

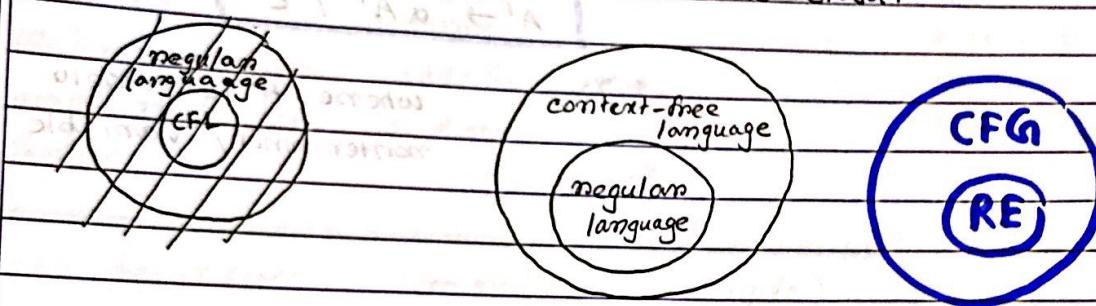
"Left Recursion and Left Factoring"

Week 05 Lecture 01

Q. What is the difference between CFG and Regular Grammar? When and why do we use CFG instead of R.G. when we're implementing a language?

Every construct that can be described by a Regular Expression can be described by a grammar, but not vice-versa.

Alternatively, every regular language is a context-free language, but not vice-versa.



Finite automata cannot count the no. of a's preceded before b.

F.A. lacks the ability of counting or memorizing, therefore, we cannot implement such languages using F.A. In these cases, we need the concepts of CFGs to implement them.

$$L = \{a^n b^n \mid n \geq 1\}$$

Left Recursive → why is it bad?

Indinite Loop

A grammar is left recursive if it can be represented in the form:

$$A \Rightarrow^+ A \alpha$$

Why is this a problem for top-down parsing?

Left Recursive
Any variable
in any of its sentential form
is calling itself to the left of the sentential form produces such a context.

Left Recursive is bad for Top-Down Parsing.
(Let the parser go into an infinite loop, hence Left Recursive is bad.)

Important

B) Converting left recursive grammar
into right recursive one.

Left Recursion Removal

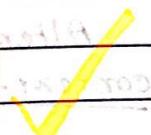
Simplest Case: "Immediate Left Recursion"

Given: $A \rightarrow A\alpha | \beta$

Transform into:

$$A \rightarrow \beta A'$$

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



where A' is a new nonterminal / variable

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

good substitution

Shorter length

[minimum length]

in better position so easy to implement that in removing it

is much faster

less time

much less

iterations left

to do your job

(program won't get stuck waiting for user input)

left parentheses

last parentheses

unlike well-known RPN in reverse Polish

left parentheses

(good structure and no memory left lost)

(but in reverse Polish lost lost)

Left and Right Recursive Grammars

Left Recursive

$$A \rightarrow A\alpha \mid \beta$$

- Causes an infinite loop.

* So we have to eliminate a left-recursive grammar.

We've to change Left Recursive grammar into the Right Recursive form.

Right Recursive

$$A \rightarrow \beta R$$

$$R \rightarrow \alpha R \mid \epsilon$$

- This never causes an infinite loop.

because parser always checks the inputs from Left to Right.

Left Recursion Elimination (example)

Grammar

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

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

(Step 1)

$$\begin{array}{|c|} \hline E \rightarrow TE' \\ E' \rightarrow +TE' \mid \epsilon \\ \hline \end{array}$$

(Step 2)

$$\begin{array}{|c|} \hline T \rightarrow FT' \\ T' \rightarrow *FT' \mid \epsilon \\ \hline \end{array}$$

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

as it is

- # Left Recursion in more than one step.
- At first, eliminate non-immediate left recursion for prod. rules
 - Then, we've to eliminate the immediate left recursion for the same prod. rule.
Do this for each and every prod. rule.

Example 1

$$S \rightarrow A^d | b$$

$$A \rightarrow Ac | Sd | e$$

left recursive, \rightarrow Yes, indirect
left recursive \rightarrow Yes, direct



$$S \rightarrow A^d | b$$

$$A \rightarrow Ac | e | \underline{Af d} | bd$$

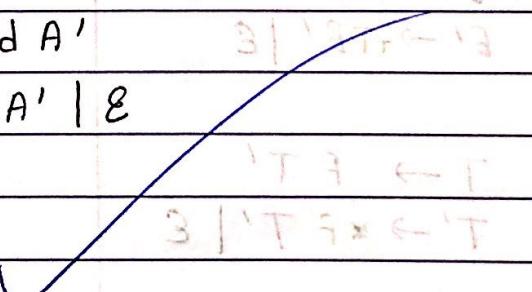


$$S \rightarrow A^d | b$$

$$A \rightarrow e A' | bd A'$$

$$A' \rightarrow f d A' | c A' | \emptyset$$

$$A \rightarrow e | bd | Af d | Ac$$



21/6/23

b11 (2) - 3

Example 2

$$S \rightarrow A\varnothing \mid b$$

$$A \rightarrow Ac \mid Sd \mid Be$$

$$B \rightarrow Ag \mid Sh \mid k$$

indirect left recursive

direct left recursive

indirect left recursive

$$\hookrightarrow S \rightarrow A\varnothing \mid b$$

$$\hookrightarrow A \rightarrow Ac \mid Be \mid A\varnothing d \mid bd$$

$$B \rightarrow Ag \mid k \mid A\varnothing h \mid bh$$

$$S \rightarrow A\varnothing \mid b$$

$$A \rightarrow BeA' \mid bdA'$$

$$A' \rightarrow cA' \mid \theta fdA' \mid \varnothing$$

$$B \rightarrow Ag \mid k \mid A\varnothing h \mid bh$$

$$B \rightarrow BeA'g \mid BeA'fh \mid k \mid bh \mid bdA'g \mid \cancel{bdA'k} \mid \cancel{bdA'bh}$$

$$B \rightarrow k \mid bh \mid (BeA'g) \mid (BeA'fh) \mid bdA'g \mid bdA'fh$$

$$B \rightarrow kB' \mid bhB' \mid bdA'gB' \mid bdA'fhB'$$

$$B' \rightarrow \varnothing \mid eA'gB' \mid \theta eA'fhB'$$

Finally,

$S \rightarrow Ad/b$

$p \rightarrow \text{Be } A' \mid \text{bd } A'$

$$A' \rightarrow c A' \mid d A' \mid e$$

$$B \rightarrow kB' \mid bhB' \mid bdA'gB' \mid bdA'ghB'$$

$$B' \rightarrow e A' g B' \mid e A' f h B' \mid e$$

1

Important

Left Factoring

Week 05
Lecture 02

to eliminate nondeterminism

- problem of limited lookaheads
- problem of generating same CPGs on doing the same operations twice or thrice (this takes more memory / time than parser)

How do we avoid these problems?

We factor the common terms out.

Factoring of common terms \rightsquigarrow Left factoring

* So to eliminate the problem of generating the same grammar string multiple times, take a common grammar string and factor it out.

Given Grammar:

$$S \rightarrow iEtS \mid iEtSe \mid a$$

$$E \rightarrow b$$

Left factored equivalent of the Grammar:

$$S \rightarrow iEtS' \mid a$$

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

$$E \rightarrow b$$

Left Factoring saves our time and improves our efficiency.

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

Left factoring removes non determinism in grammar.

(non-deterministic grammar \rightarrow form game symbol
common prefixes
 \Rightarrow causes tree backtracking)

Q. Do left factoring on the following grammar.

$$S \rightarrow aEbS \mid aEbSe \mid x$$

$$E \rightarrow y$$

$$S \rightarrow aEbSx' \mid x$$

$$x' \rightarrow c \mid \epsilon$$

$$E \rightarrow y$$

$$S \rightarrow aBbA \mid BbA \leftarrow a$$

$$d \leftarrow \emptyset$$

→ command with the following format

$$a \mid \emptyset \leftarrow a$$

$$S \mid \emptyset \leftarrow a$$

$$d \leftarrow \emptyset$$

→ command that

invoked here with some

variables like

A A A A

Important

Top Down Parsing

Week 05, Lecture 03

- **Recursive Descent Parsing** is a form of Top Down Parsing
- Computing **First function** \rightarrow used in **Predictive Parsing**

Top Down Parsing

- Find a left-most derivation
- Build a parse tree
- Start building from the root and work down
- As we search for a derivation
(must make choices: • which rule to use • where to use it)

Two types of Top-Down Parsing

Predictive Parsing

Recursive Descent Parsing

- backtracking is needed
(if a choice of a production rule does not work, we backtrack to try other alternatives.)
- It's a general parsing technique, but not widely used.
- Not efficient

✓ no backtracking

- efficient
- needs a special form of grammars, (LL(1) grammars)

- * RDP takes a lot of time as a lot of backtracking involved.
And the prod. rules we're choosing on the basis of priority.
So, very inefficient approach.

4) Recursive Descent Parsing (Backtracking)

$$S \rightarrow Aa \mid Ce$$

$$A \rightarrow aaB \mid aaba$$

$$B \rightarrow bbb$$

$$C \rightarrow aAD$$

$$D \rightarrow bbd$$

Input: aa bbde

(priority - written or quoted before)

probable from that priority

start seeing a which

or in this language it's already generated tokens

leftmost

matching: (as non-terminal formed)

expand

backtrack to see other option

mismatch

probable岐義性

probable岐義性

probable岐義性

behavior of nondeterminism

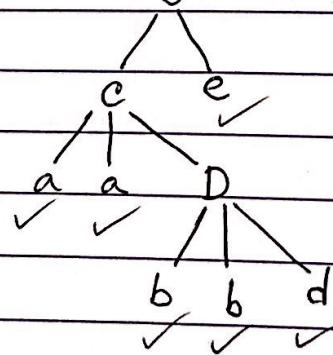
as the part of should not be matched

(continued)

S → bb, earliest occurring terminal is bb.

but ignore

Parse Tree



* First Func is needed to be used in order to do Predictive Parsing

#

First Function

Find the first terminal

$$S \rightarrow aB$$
$$B \rightarrow d$$

combo: @d

$$\text{FIRST}(S) = \{a\}$$
$$\text{FIRST}(B) = \{d\}$$

$$\text{FIRST}(\alpha) = \{a\}$$

$$S \rightarrow AB$$
$$A \rightarrow c$$
$$B \rightarrow d$$

combo: @d

$$\text{FIRST}(S) = \{c\}$$
$$\text{FIRST}(A) = \{c\}$$
$$\text{FIRST}(B) = \{d\}$$

$$\text{FIRST}(\alpha) = \{c\}$$

$$S \rightarrow AB$$
$$A \rightarrow c | \epsilon$$
$$B \rightarrow d$$

combo: cd, cd
= @d, @d

$$\text{FIRST}(S) = \{c, d\}$$
$$\text{FIRST}(A) = \{c, \epsilon\}$$
$$\text{FIRST}(B) = \{d\}$$

$$\text{FIRST}(\alpha) = \{c, d\}$$

$$S \rightarrow AB$$
$$A \rightarrow c | \epsilon$$
$$B \rightarrow d | \epsilon$$

combo: cd, ce, ed, ee
= @d, @, @, @

$$\text{FIRST}(A) = \{c, \epsilon\}$$

$$\text{FIRST}(B) = \{d, \epsilon\}$$

$$\text{FIRST}(S) = \{c, d, \epsilon\}$$

$$\text{FIRST}(\alpha) = \{c, d, \epsilon\}$$

6.

$$S \rightarrow aSe | STS$$

$$T \rightarrow RSe | Q$$

$$R \rightarrow nSn | \epsilon$$

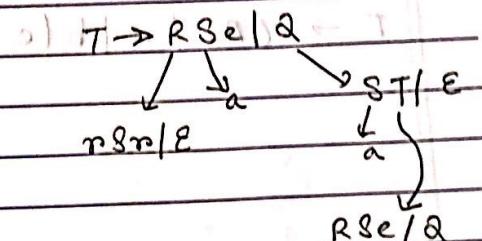
$$Q \rightarrow ST | \epsilon$$

Ans: $\text{FIRST}(S) = \{a\}$

$\text{FIRST}(T) = \{n, a, \epsilon\}$

$\text{FIRST}(R) = \{n, \epsilon\}$

$\text{FIRST}(Q) = \{a, \epsilon\}$



combo: $nSnSe, nSa, aT, \epsilon$

$$\begin{aligned} &= n, 8Sa, a, \epsilon \\ &= n, 8a, a, \epsilon \\ &= n, a, \epsilon \end{aligned}$$

$$R \rightarrow nSn | \epsilon$$

combo: n, ϵ

$$a \rightarrow ST | \epsilon$$

combo: aT, ϵ
= a, ϵ

5.

$$P \rightarrow i \mid c \mid n \mid TS$$

$$Q \rightarrow P \mid as \mid b \mid S \mid ST$$

$$R \rightarrow b \mid \epsilon$$

$$S \rightarrow c \mid Rn \mid \epsilon$$

$$T \rightarrow RSq$$

Ans:

$$\text{FIRST}(P) = \{i, c, n\}$$

$$\text{FIRST}(Q) = \{i, c, n, a, b\}$$

$$\text{FIRST}(R) = \{b, \epsilon\}$$

$$\text{FIRST}(S) = \{c, b, \epsilon\}$$

$$\text{FIRST}(T) = b, \epsilon$$

$$\text{FIRST}(S) = \{c, b, n, \epsilon\}$$

$$\text{FIRST}(A)$$

$$\text{FIRST}(T) = \{b, c, n, q\}$$

$$S \rightarrow c \mid Rn \mid \epsilon$$

$$\downarrow b \mid \epsilon$$

$$\begin{aligned} &\text{combo: } c, bn, \epsilon n, \epsilon \\ &= c, bn, n, \epsilon \\ &= c, b, n, \epsilon \end{aligned}$$

$$T \rightarrow RSq$$

$$\begin{array}{l} \downarrow b \mid \epsilon \\ \downarrow c \mid Rn \mid \epsilon \\ \downarrow \epsilon \end{array}$$

combo:

$$\begin{aligned} &b cq, b bnq, b bnq, b nq \\ &b \epsilon q, \epsilon cq, \epsilon bnq, \\ &\epsilon bnq, \epsilon \epsilon q \end{aligned}$$

$$\begin{array}{l} \Rightarrow b, b, b, b, c, b, \\ n, q \end{array}$$

$$\{i, c, n\} = (a) T2917$$

$$\{b, a, q\} = (b) T2917$$

$$\{c, b, n\} = (c) T2917$$

$S \rightarrow aSe \mid B$ $B \rightarrow bBcCf \mid C$ $C \rightarrow cCg \mid d \mid \epsilon$	$FIRST(S) = \{a, b, c, d, \epsilon\}$ $FIRST(B) = \{b, c, d, \epsilon\}$ $FIRST(C) = \{c, d, \epsilon\}$
--	--

8.	$S \rightarrow (A) \mid \epsilon$	$\text{FIRST}(S) = \{\epsilon\}$
	$A \rightarrow TE$	$\text{FIRST}(A) = \{\epsilon, a, b, c\}$
	$E \rightarrow \&TE \mid \epsilon$	$\text{FIRST}(E) = \{\&, \epsilon\}$
	$T \rightarrow (A) \mid a \mid b \mid c$	$\text{FIRST}(T) = \{\epsilon, a, b, c\}$

TNR Practice for LR(1)

Q. Given the following grammar,

$$S \rightarrow A$$

$$S \rightarrow \alpha b$$

$$A \rightarrow aAb$$

$$A \rightarrow B$$

$$B \rightarrow b$$

Compute the LR(1) items
and construct the
Parsing Table.

Ans:

Step 1: Augment the Grammar

① $S' \rightarrow S$

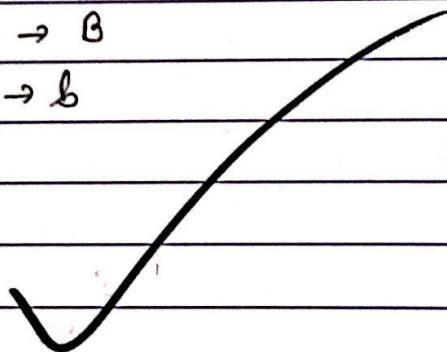
② $S \rightarrow A$

③ $S \rightarrow \alpha b$

④ $A \rightarrow aAb$

⑤ $A \rightarrow B$

⑥ $B \rightarrow b$



Step 2: Canonical collection of LR(1) items

I₀

$S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot A, \$$
 $S \rightarrow \cdot xB, \$$
 $A \rightarrow \cdot aAb, \$$
 $A \rightarrow \cdot B, \$$
 $B \rightarrow \cdot b, \$$

S

$S' \rightarrow S., \$$

I₁

A

$B \rightarrow A., \$$

I₂

b

$B \rightarrow b., \$$

B

x

$S \rightarrow A., \$$

I₃

$A \rightarrow B., \$$

I₅

$A \rightarrow a. Ab, \$$

$A \rightarrow \cdot aAb, b$

$A \rightarrow \cdot B, b$

$B \rightarrow \cdot b, b$

I₄

A

$A \rightarrow aA.b, \$$

I₈

b

B

a

$A \rightarrow a.Ab, b$

$A \rightarrow \cdot aAb, b$

$A \rightarrow \cdot B, b$

$B \rightarrow \cdot b, b$

I₉

I₁₁

$B \rightarrow b., b$

b

B

b

$A \rightarrow aA.b, b$

I₁₃

$A \rightarrow aAb., b$

b

$A \rightarrow aAb., b$

I₁₄

States	Action				Go To		
	a	b	x	\$	A	B	S
0	s4	s6	s3		2	3	52
1				Accept			1
2					r01		
3							
4	s59	s11			8	10	
5					r04		
6					r05		
7					r02		
8							
9	s9	s11			13	10	
10					r04		
11					r05		
12					r03		
13					s14		
14					r03		

final states,

$$\checkmark I_2 \Rightarrow n1 \rightarrow \$$$

$$\checkmark I_5 \Rightarrow n4 \rightarrow \$$$

$$\checkmark I_6 \Rightarrow n5 \rightarrow \$$$

$$\checkmark I_7 \Rightarrow n2 \rightarrow \$$$

$$\checkmark I_{10} \Rightarrow n4 \rightarrow b$$

$$\checkmark I_{11} \Rightarrow n5 \rightarrow b$$

$$\checkmark I_{12} \Rightarrow n3 \rightarrow \$$$

$$\checkmark I_{14} \Rightarrow n3 \rightarrow b$$

Correct

Ans 1:

(i)

	FIRST	FOLLOW
S	{a, b}	{\$}
A	{a, ε}	{a, b}
B	{b, ε}	{a, b}

(ii)

	a	b	\$
S	$S \rightarrow AaAb$	$S \rightarrow BbBa$	
A	$A \rightarrow a$ $A \rightarrow \epsilon$	$A \rightarrow \epsilon$	
B	$B \rightarrow \epsilon$	$B \rightarrow b$ $B \rightarrow \epsilon$	

(iii) No, it is not LL(1)

since there are multiple entries in a single cell

Also, since grammar is Ambiguous

it is not LL(1)

(iv) Follow set of a variable can never contain an epsilon because,

Since, error epsilon transition is for E-NFA
and not DFA.Errors cannot be handled if there is epsilon
by the parsers.

Ans :
2

Step 1: Augmented Grammar

$x' \rightarrow x$

$x \rightarrow Mm$

$x \rightarrow nMo$

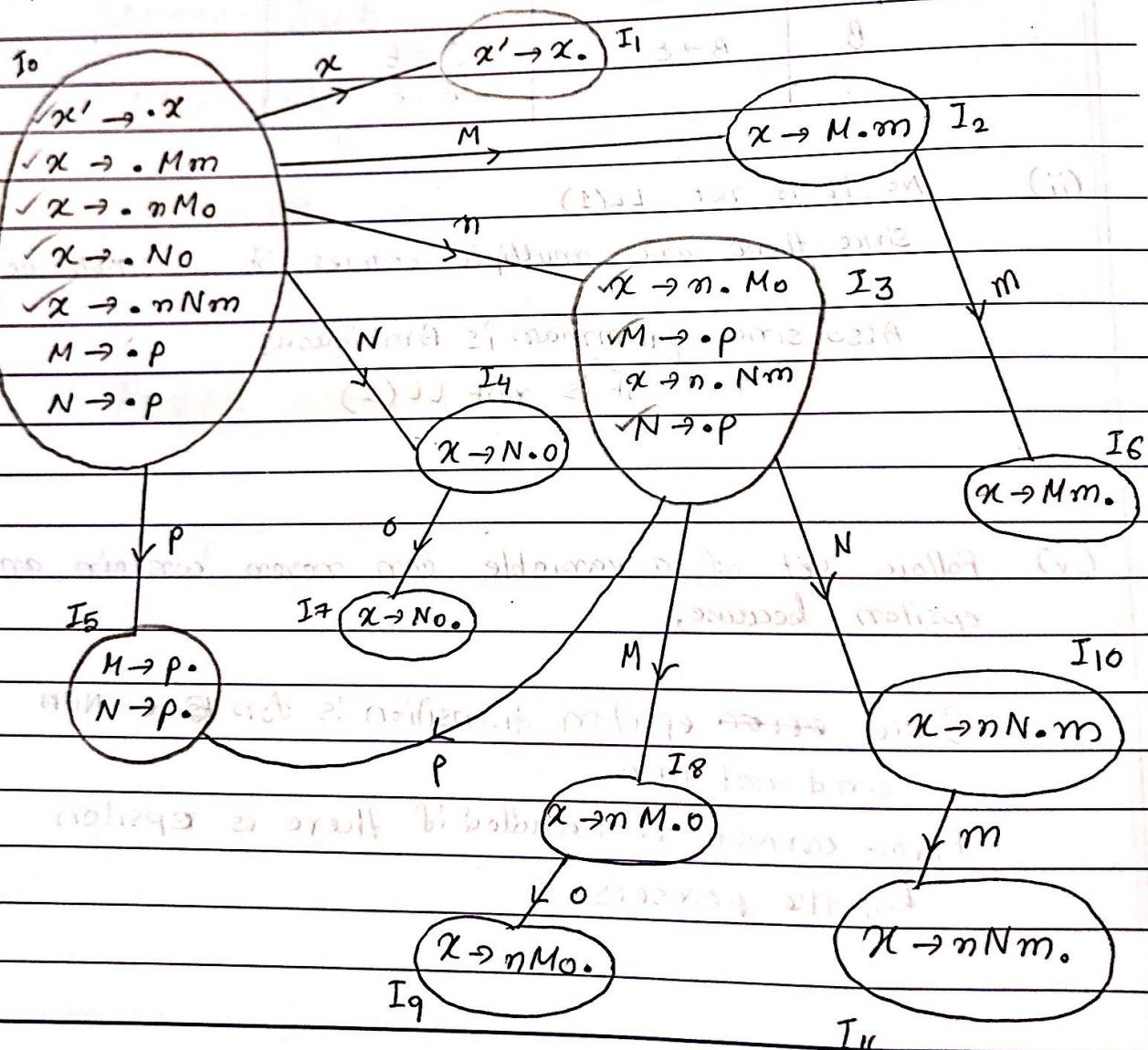
$x \rightarrow No$

$x \rightarrow nNm$

$M \rightarrow P$

$N \rightarrow P$

Step 2: Canonical collection of LR(0) items



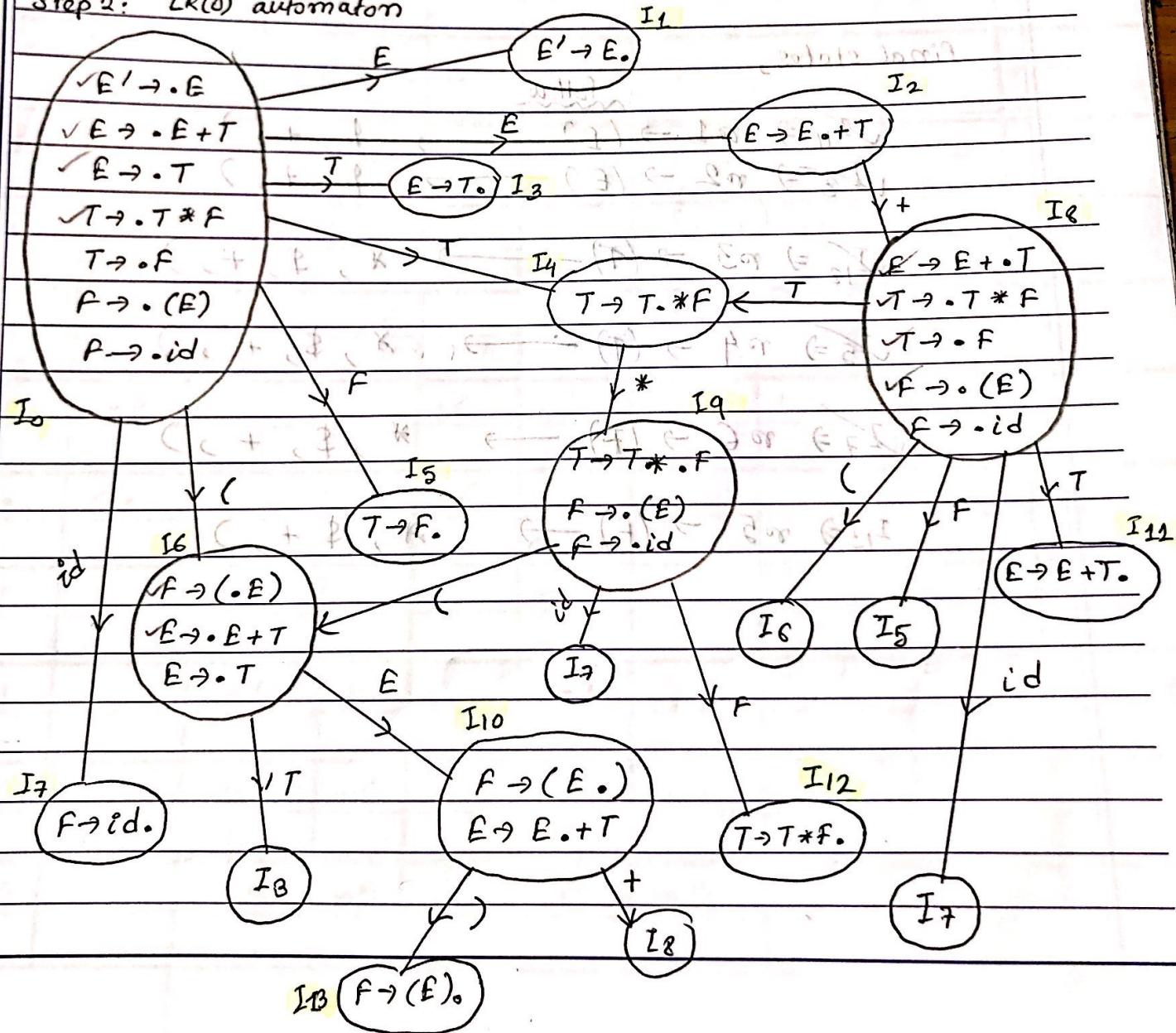
Ans 2:

Pg -

Step 1: Augment Grammar

- $E' \rightarrow E$ ①
- $E \rightarrow E + T$ ②
- $E \rightarrow T$ ③
- $T \rightarrow T * F$ ④
- $T \rightarrow F$ ⑤
- $F \rightarrow (E)$ ⑥
- $F \rightarrow id$ ⑦

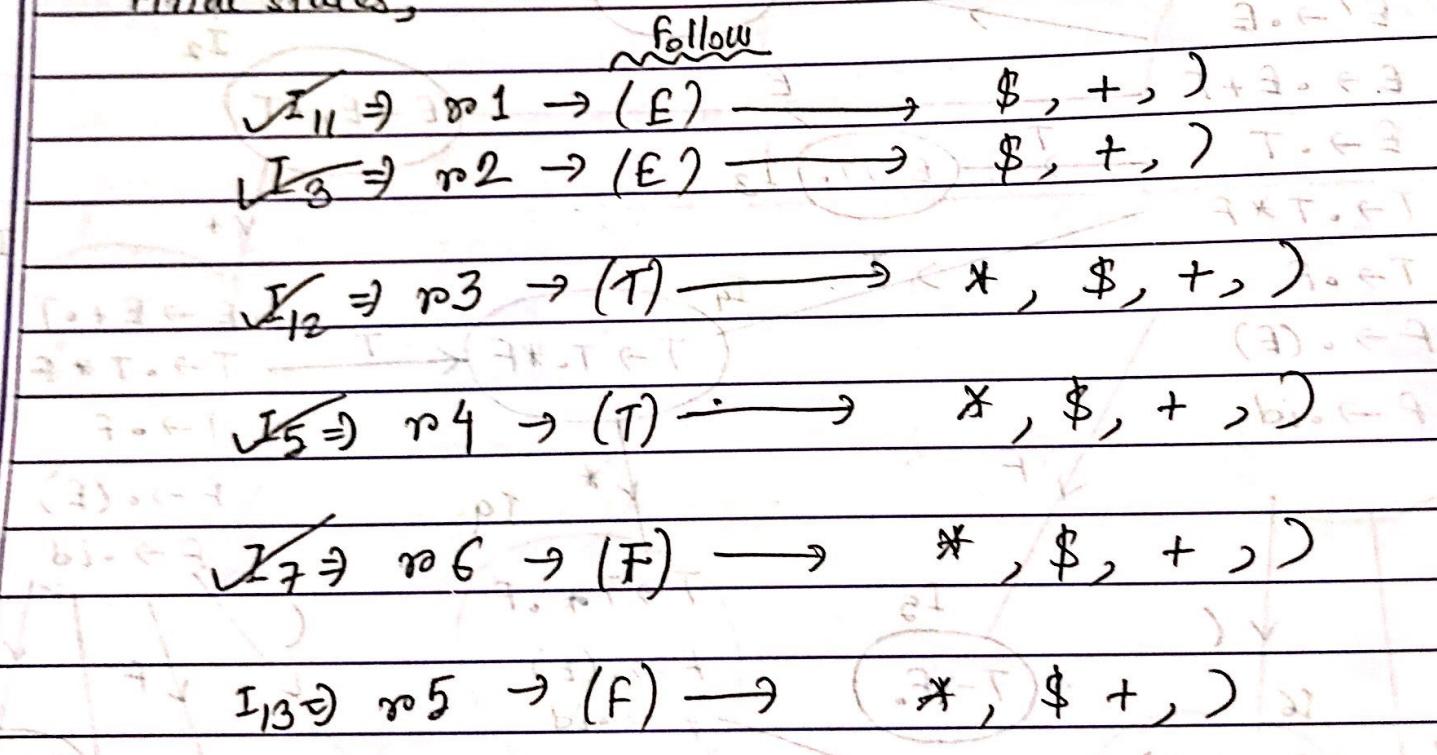
Step 2: LR(0) automaton



Variables

	FIRST	FOLLOW
E	{(, id}	{\$, +,)}
T	{(, id}	{*, \$, +,)}
F	{(, id}	{*, \$, +,)}

Final states,



	ACTION					GOTO		
i d	+	*	()	\$	E	T	F
0	s7		36			1,2	3,4	5
1					accept			
2		88						
3	r2			r2	r2			
4		s9						
5	r4	r4		r4	r4			
6						10	3	
7	r6	r6		r6	r6			
8	s7		56			4,11	5	
9	s7		36					12
10	s8		313					
11		r1		r1	r1			
12	r3	r3		r3	r3			
13	r5	r5		r5	r5			

Ans 4:

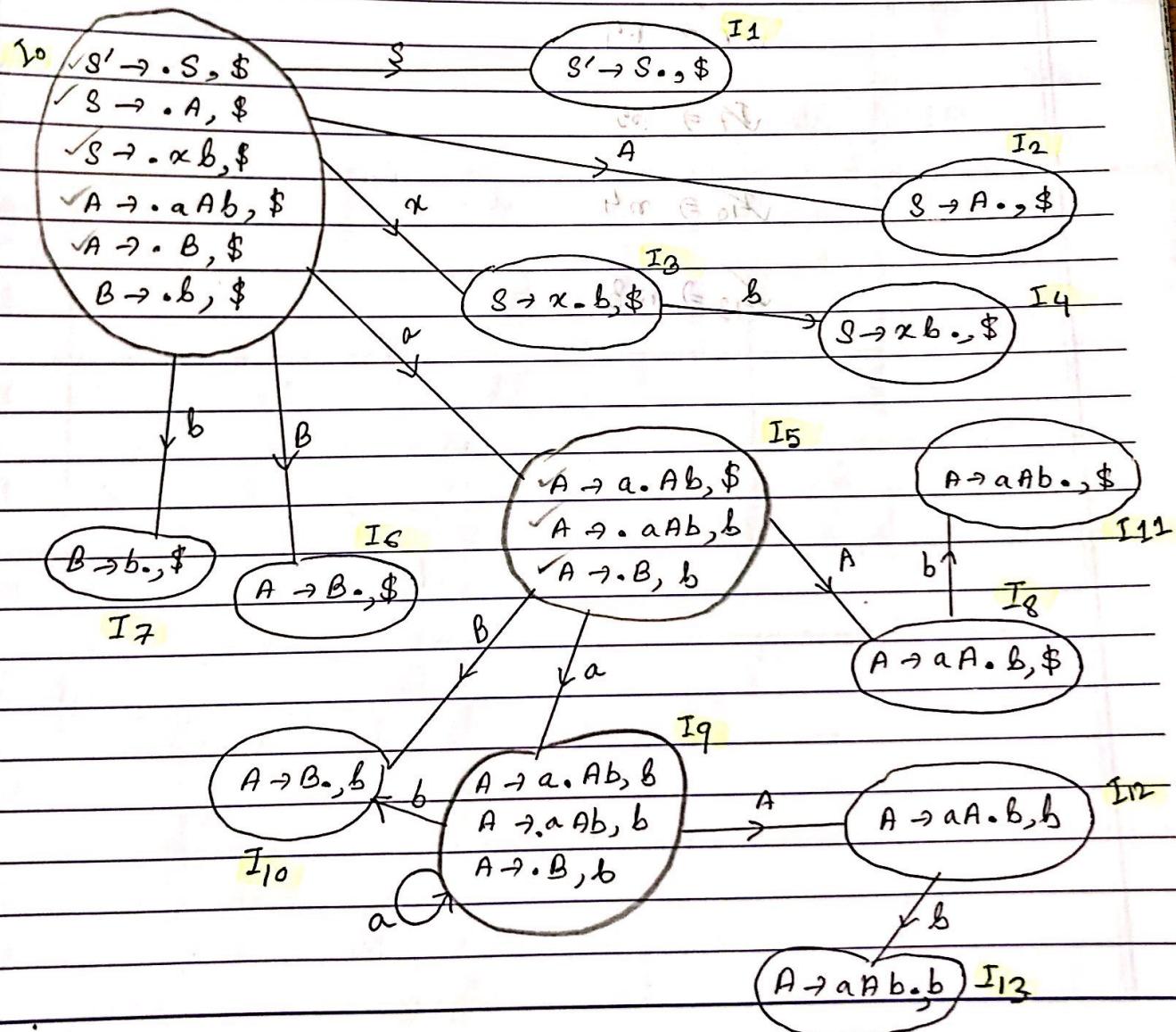
Subhi

Pg -

a) Step 1. Augment Grammar

- ① $S' \rightarrow S$
- ② $S \rightarrow A$
- ③ $S \rightarrow x b$
- ④ $A \rightarrow aAb$
- ⑤ $A \rightarrow B$
- ⑥ $B \rightarrow b$

Step 2: LR(1) items (canonical collection)



final states,

$\check{I}_2 \Rightarrow n_1$

$B \leftarrow Z$ ⑥

$A \leftarrow Z$ ⑦

$\check{I}_4 \Rightarrow n_2$

$d \times \leftarrow Z$ ⑧

$dA \leftarrow A$ ⑨

$B \leftarrow A$ ⑩

$d \leftarrow B$ ⑪

$\check{I}_1 \Rightarrow n_3$

(condition business) 2nd (1)

$\check{I}_6 \Rightarrow n_4$

$(B \leftarrow Z \leftarrow Z)$

$\check{I}_7 \Rightarrow n_5$

$\check{I}_{10} \Rightarrow n_4$

$\check{I}_{13} \Rightarrow n_3$

CLR(C) Parse Table

Subhi Bhuyan

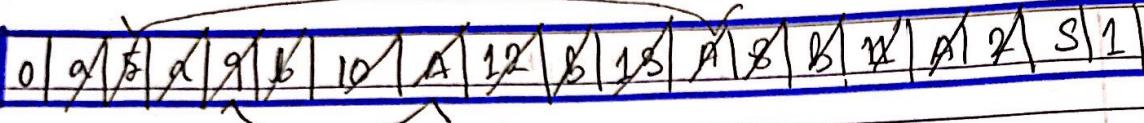
Pg -

		Action		Go To	
0	a	b	x	\$	
	s5	s7	s3	s	A
1				1	2
					6
2				accept	
				r1	
3					
				s4	
4					
				r2	
5					
				r1 ← ()	8
					10
6					
				r4	
7					
				r5	
8					
				s11	
					8s e(2, 8)
9					
	s9	s10			12
10					
				r4 ← (A, B)	
11					
				r3	
12					
				1A ← E or e(4, 8)	
					8e (6, 4)
13					
				r3	

b)

input: $\alpha \alpha b b b \$$
 $\uparrow \uparrow \uparrow \uparrow \uparrow$

Stack

 $(0, a) \Rightarrow \text{ok}$ $(1, \$) \Rightarrow \text{accept}$ $(5, a) \Rightarrow \text{ok}$ $(9, b) \Rightarrow \text{ok}$

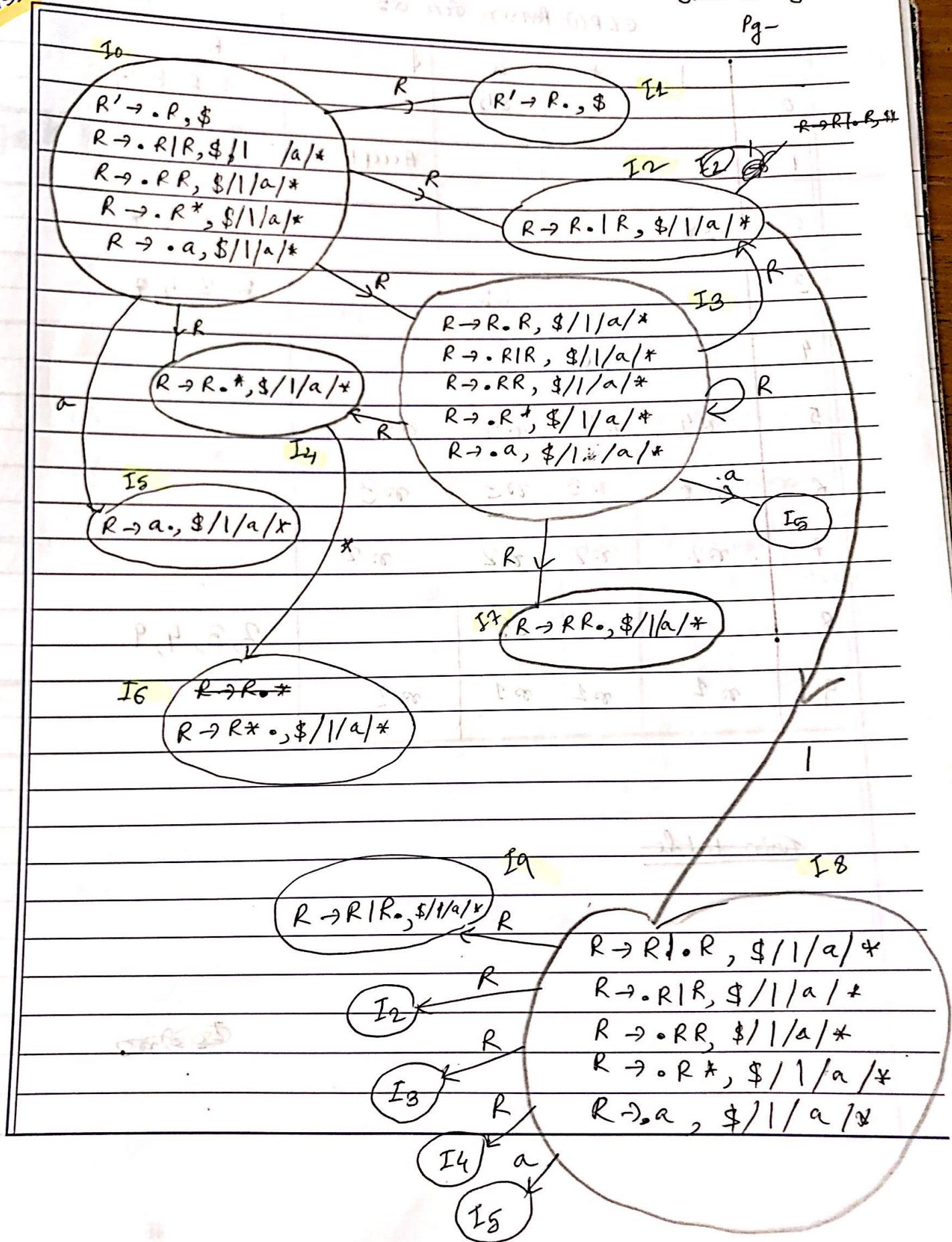
Therefore, the string
can be
parsed.

 $(10, b) \Rightarrow \text{ok}$ $(A) \rightarrow |B|$ $L = 1$ $\text{pop} = 2 L = 2$ $(12, B) \Rightarrow \text{ok}$ $\rightarrow \text{push}$ $(A, 9) \Rightarrow 12$ $(13, b) \Rightarrow \text{ok}$ $(A) \rightarrow |a A B b|$ $L = 3$ $\text{pop} = 6$ $(8, B) \Rightarrow \text{ok}$ $\rightarrow \text{push}$ $(5, A) \rightarrow \text{ok}$ $(1, \$) \Rightarrow \text{ok}$ $L = 3$ $\text{pop} = 6$ $\text{push } A$ $(2, \$) \rightarrow \text{ok}$

$S \rightarrow A$
 $L = 1$
 $\text{pop} = 2$

 $(A, 0) \Rightarrow 2$ $(S, 0) \Rightarrow 1$

Ans 5:



CLR(1) Parson for Q5

	*	a	\$	R
0		s5		1, 2, 3, 4
1		,	Accept	2, 3, 4, 5
2	s8			2, 3, 4, 5
3		s5		2, 3, 4, 7
4		s6		2, 3, 4, 8
5	n4	n4	n4	n4
6	n3	n3	n3	n3
7	n2	n2	n2	n2
8		s5		2, 3, 4, 9
9	n1	n1	n1	n1

Fair table

RE: $a^* a | a \$$

Stack: $[0] \alpha / 5 / R / 1 / 2 / 3 / A / * / \$ / R$

$(0, a) \Rightarrow S/5$

$(5, *) \Rightarrow n4 \quad L \rightarrow |a|$
 $L=1 \text{ pop} = 2$

$(R, 0) \Rightarrow 1, 2, 3, 4$

$(4, *) \Rightarrow \$\$$

$(6, a) \Rightarrow n3$

$R \rightarrow |R \neq| L=2 \text{ pop} = 4$

$(R, 2) \Rightarrow \text{blank entry}$

So, Syntax error

Hence, the given RE
cannot be parsed

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

<u>output</u>	α	match
$P \rightarrow a$	β	
$P \rightarrow (E)$	γ	
$E \rightarrow PE'$	δ	match
$P \rightarrow a$	ϵ	
$E' \rightarrow PE'$	ζ	
$P \rightarrow (E)$	η	match
$E \rightarrow PE'$	θ	
$P \rightarrow a$	ρ	
$E' \rightarrow PE'$	χ	match
$P \rightarrow a$	ψ	
$E' \rightarrow E$	γ	
$E' \rightarrow E$	φ	match
$E' \rightarrow E$	ρ	
	ϵ'	
	χ	match
	ψ	
	γ	
	φ	
	$\$$	
	α	
	β	
	γ	match

Compiler Design

Week 0

- Theories of Compiler Design used in NLP, Machine Translation, Compiler Information Retrieval, Semantic Web Technology, Query Language Processing.
- procedure for identifying syntax errors

"Compiler Basics"

Compiler is a translator that converts the high-level language into the machine language.

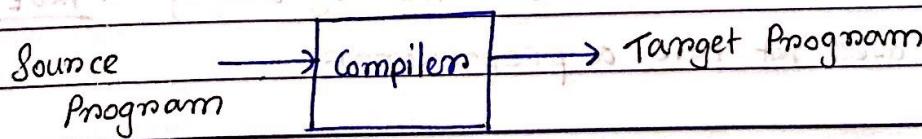
Compiler is used to show errors to the programmer.

A compiler is a computer program that translates computer code written in one programming language (the source language) into another language. (the target language.)

- Programming problems are easier to solve in high-level languages.
- Solutions are usually more efficient (faster, smaller) when written in Machine Language.

Compilers are the bridges Tools to translate programs written in HLL to efficient executable code.

Another important role of the compiler is to report any errors in the source program that it detects during the translation process. (compiler errors)



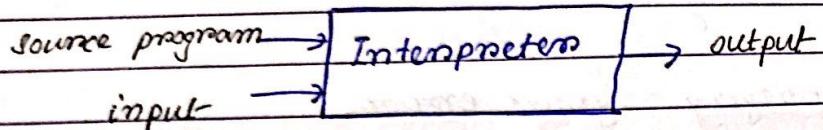
(Compiler is a Software program)

Q. What is an interpreter?

Ans:

An interpreter requires no pre-computation

it directly takes the source program and runs the input on it producing the desired output. (line by line / instruction by instruction)



Q. Difference betⁿ Interpreters and Compilers.

↳ slower

* The interpreter takes less memory (no exe file needed)

* The interpreter is invoked each time the CPU has to run any piece of instruction (as no files.exe generated)

Compilers & interpreters do similar jobs, but there are differences:

To run a program you've written, it must first be translated into machine code so the computer can read it. This is what compilers and interpreters do.

However, Compilers convert the code all at once, save it, then run it. ↳ (C, C++) Java

Whereas, interpreters translate the code one line at a time, as it is run. ↳ (Python, Ruby) JavaScript

* Interpreter is slower

* The interpreter converts HLL to Assembly Language \Rightarrow FALSE

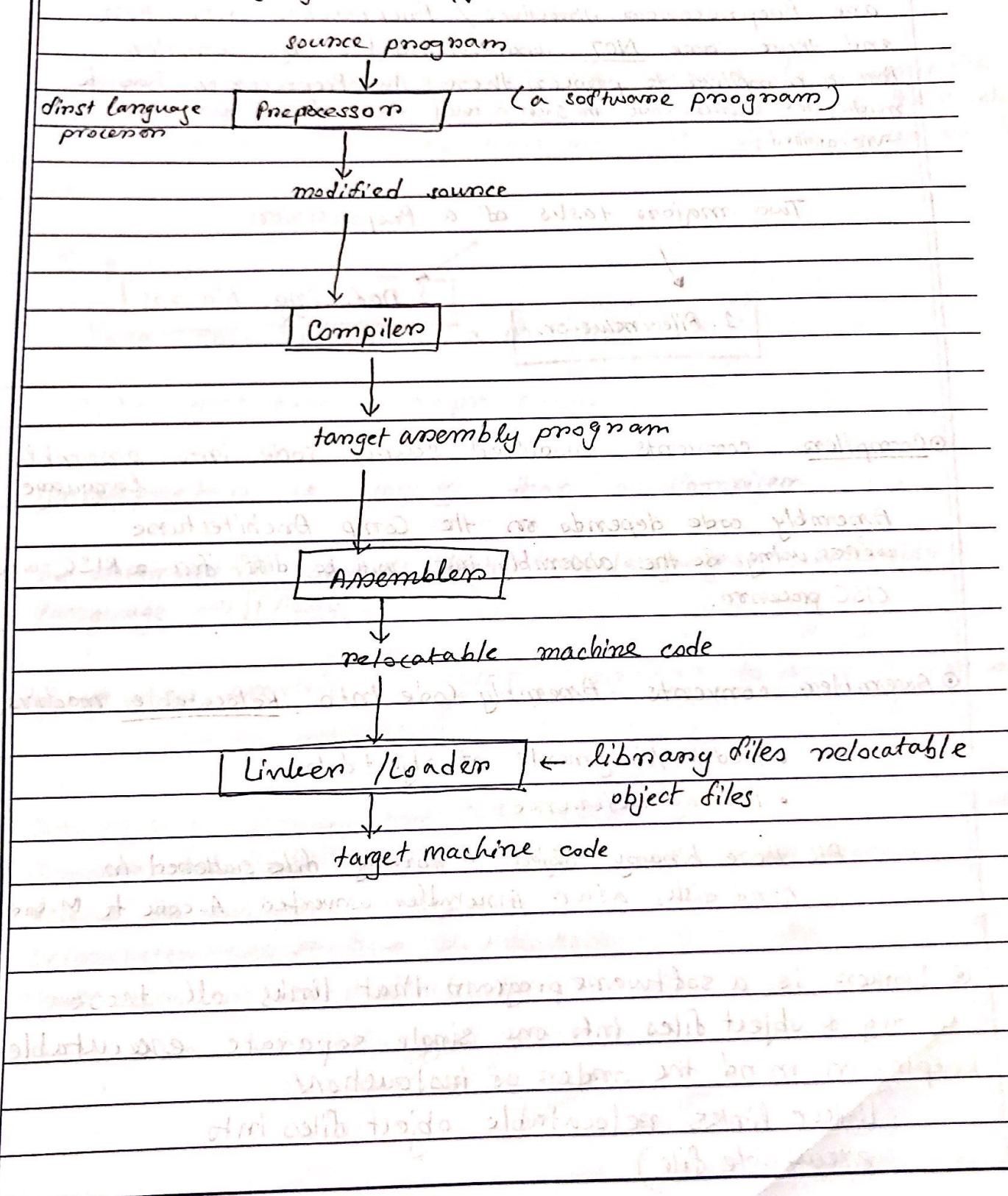
* The interpreter

* The compiler converts HLL to Assembly Language \Rightarrow TRUE

* Java does not have a preprocessor.

- A Language Processing System is a system that actually converts the source program written in one language into a target machine code consisting of 1s and 0s by incorporating different language processors such as preprocessors, compilers, assemblers, linkers and loaders.

Q Steps of a Language Processing System.



(asks the Preprocessor to include the standard iostream header file into our program)

Instructions preceded by a hash symbol are Preprocessor directives / Preprocessor instructions and these are NOT understandable by Compilers.

For a compiler to process these, the Preprocessor has to modify the source code in such a way that these instructions are omitted.

Two major tasks of a Preprocessor

1. File inclusion

2. Defining Macros

Compiler converts modified source code into assembly language.

Assembly code depends on the Comp. Architecture we're using. So the assembly inst. will be diff for a RISC and CISC processor.

Assembler converts Assembly code into Relocatable machine code.

- modules/fragments of object files
- library reference

All these binary object & library files scattered in hard disk, after Assembler converted. A. code to M. code

Linker is a software program that links all these library & object files into one single separate executable file. keeping in mind the order of instructions.

(Linker links relocatable object files into executable file)

General Summary: { * Compilers
* Compilers vs Interpreters
* Steps of LP system.

Q. Loader loads executable machine code into main memory on execution.

Loader will take the executable file produced by Linker in the hard disk and it will load them into the memory by mapping them into memory addresses, after that those segments will be loaded for execution.

Ques - 01

1. Java does not have a preprocessor.
2. C, C++ does have a preprocessor.
3. Interpreter is slower than a Compiler.
4. The interpreter converts higher level languages to assembly language → **FALSE**
5. The compiler converts higher-level languages to assembly language.

Compiler runs faster than Interpreter.

Compilers take less time to execute than interpreters because the run time is slower.

Interpreter doesn't have any compile time as the code never compiles, it directly runs.

"Phases of Compilers"

July 12th, 2020

Week 1, Lecture - 2

Phases of a Compiler (L8S1CT)

- Lexical Analysis
- Syntax Analysis
- Semantic Analysis
- Intermediate code generation
- Code optimization
- Code generation

built-in keyword → lexeme
(meaningful sequence
of characters)

Lexical Analyzer/scanner will find lexeme (variable, no., operator)
and generate tokens for the same.

This token stream fed as input into Syntax Analyser/Parsers.
which will generate a syntax tree based on a grammatical
structure. or CFG

Semantic Analyzer checks for inconsistency in meaning.
(array size can never be declared as a double)

ICG → (3 address code generated
called Intermediate Representation)

MICO's task → to reduce the no. of lines or to
reduce the no. of variables in 3-address code
to decrease Compilation Time
and increase Efficiency.

Code Generator → generates Assembly Instructions
which is hardware-specific.

MDCO → to reduce no. of lines + variables in assembly code.
(optional)

(machine-dependent
code optimizers)

Code Optimizer