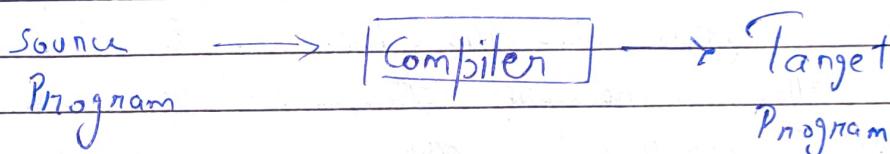


Date Jan 12<sup>th</sup>, 2023

# Compiler Design

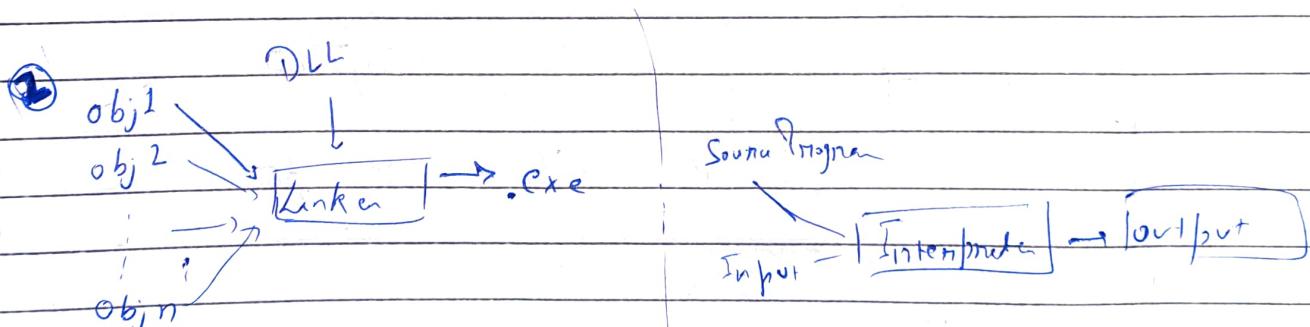
## Introduction

- \* Compiler is a program that translates a source program (a program written in a source language) & translate it to the equivalent program (target program, program written in a target language).



• C/C++  
• Java

• obj  
• asm



③ .exe → [Loader] → [memory]

\*) The machine language program

① .asm → [Assembler] → .o

Date .....

Assembler → The compiler may produce an assembly language program as its output. Because it is easier to produce to output & its easier to debug. The assembly lang. program is then processed by the programs called assembler that produces the relocatable object code as its output.

Linker - Large Programs are often compiled in pieces. So the relocatable code may have to be linked together with other relocatable object files, library files into the code that actually runs on the machine. The linker is also responsible for this task. It gives the ex.e

Loader → The loader finally puts the executable file into memory for execution.

## Structure of the Compiler

It is divided into two parts phases.

i) Analysis (front-end)

ii) Synthesis (Back-end)

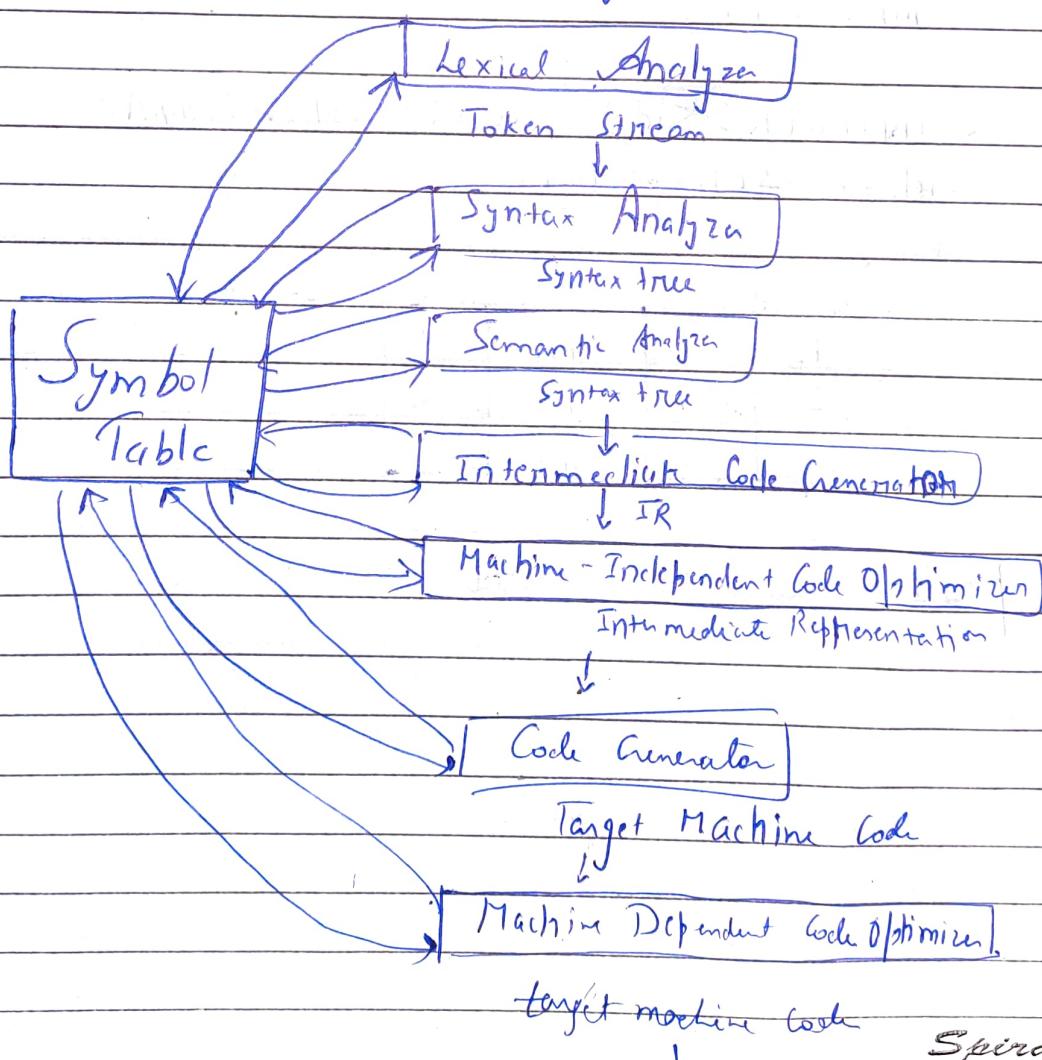
The Analysis breaks up the source program into constituent pieces & imposes the grammatical structure on them

Date .....

It then uses the structure to create an intermediate representation (IR) of source program. The Analysis also collects the information about the source program & stores it in a data structure called a symbol table, which is passed along with IR to the Synthesis part.

Synthesis part constructs the desired target program from the IR & information in the symbol table.

### Phases of a Compiler (fig 1.6)



Date .....

## Lexical Analysis

- \* It <sup>is known</sup> done by the program called lexical analysis  
Lexical Analysis or scanning

The first phase of a compiler is called a scanner.  
It reads the stream of characters & groups the characters into meaningful sequences called lexemes.

For each lexeme, the lexical analysis produces as output a token of the form

<token-name, attribute value>

Ex →

int a,b,c ;

<INTEGER><id,1><COMMA><id,2><COMMA>  
<id,3><SEMI-COLON>

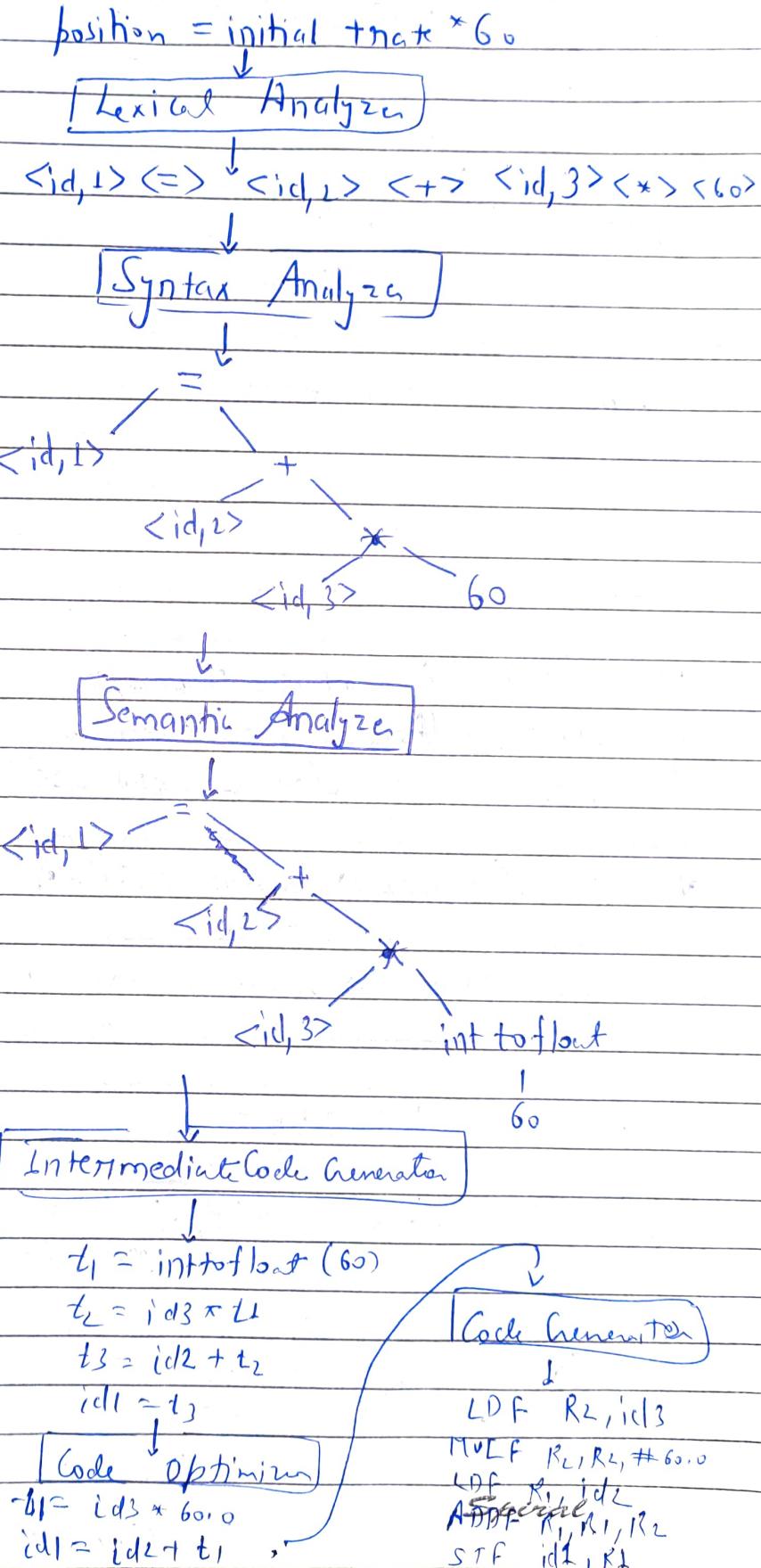
Symbol Table →

1	a	int	4
2	b	int	4
3	c	int	4

position = initial + n \* 60

# Translation of an Assignment Statement

fig 1.7



Date Jan. 16, 2023

## Chapter-3

# Lexical Analysis

### Syntax

whitespace

pattern  $\downarrow$  {action}  
pattern {action}  
|  
pattern {action}

Rule set +  
Syntax  
% %

definition  
% %  
Rules

% %  
Auxiliary func'

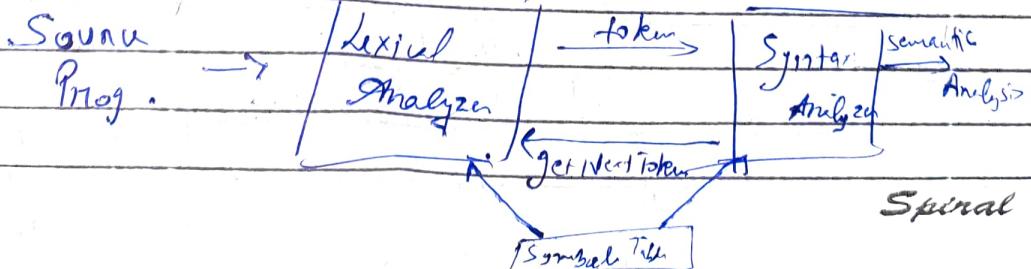
- 1) Read the input characters
- 2) Group them into lexemes.
- 3) Sequence of tokens for each lexeme.

Note  $\rightarrow$  Lexemes  $\rightarrow$  Words derived from input character stream

Token  $\rightarrow$  Lexemes mapped into token name &  
an attribute value

$$x = a + b * z$$

$$\text{Tokens} = \{ \langle \text{id}, 1 \rangle, \langle \Rightarrow, \rangle, \langle \text{id}, 2 \rangle, \langle +, \rangle, \langle \text{id}, 3 \rangle, \langle *, \rangle, \langle \text{num}, 4 \rangle \}$$



Date .....

Lexeme  $\rightarrow$  sequence of characters recognized by a pattern.

pattern  $\rightarrow$  structure that is matched by multiple strings

$[ -9-2 \text{A-Z}] [a-2 \text{A-Z} 0-9]^{*}$

Ex  $\rightarrow$  `printf("Total = %d\n", score);`

Lexeme

$<\text{id}, 1> <\text{left parenthesis}> <\text{string literal}, "Total = %d\n">$   
 $<\text{comma}> <\text{id}, 2> <\text{parenthesis}> <\text{semicolon}>$

\* We want to define 1 token for each keyword

\* Pattern for the keyword will be same as the keyword itself.

\* One token identifies all the identifiers.

\* Token for each punctuation symbol. del, (,), <, >, , <=, etc.

\* The token name influences the parsing actions.

## Handling Lexical Errors

There are 4 recovery actions that lexical analyzer takes.

- 1) Delete 1 character from remaining input.
- 2) Insert a missing character in the remaining input.
- 3) Replace a character by another character.
- 4) Transpose two adjacent characters.

Date .....

Alphabet :- finite set of symbols

String :- word made up of finite sequence of symbols drawn from the alphabet.

\*  $|s| \rightarrow$  length of string

\*  $\epsilon \rightarrow$  string of length 0.

Language :- Countable set of strings over some fixed alphabet.  
(maybe uncountable in some cases).

Prefix  $\rightarrow$  by removing 0 or more symbols from the end of string.

Total  $(n+1)$  prefixes for string of length n.

Suffix  $\rightarrow$  by removing 0 or more symbols from the starting of string.

Substring  $\rightarrow$  by removing any prefix or any suffix from the string.

Total  $\frac{n(n+1)}{2} + 1$  substrings for string of length n.  
 $(\sum_{i=1}^n i + 1)$

Proper prefix by removing at least one character but not all from the string.

Subsequence  $\rightarrow$  deleting 0 or more but not necessarily consecutive  $\downarrow$  string of  $s$ .

Total  $\rightarrow 2^n$

Date .....

Union of  $L_1$  &  $L_2$   $\rightarrow$

$$L_1 \cup L_2 = \{ s \mid s \in L_1 \text{ or } s \in L_2 \}$$

Concatenation of  $L_1, L_2$   $\rightarrow$

$$L_1 L_2 = \{ st \mid s \in L_1 \text{ and } t \in L_2 \}$$

Kleene's Star  $\rightarrow$

$$L^* = \bigcup_{i=0}^{\infty} L^i$$

Kleene's Closure  $\rightarrow$

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

Note  $\rightarrow$

$$|LUDI| = 620 \quad \left\{ \begin{array}{l} L = \{a, b, \dots, z, A, B, \dots, Z\} \\ D = \{0, 1, 2, \dots, 9\} \end{array} \right.$$

$$|LD| = 520 \quad \left\{ \begin{array}{l} D = \{0, 1, 2, \dots, 9\} \end{array} \right.$$

$L(LUDI)^*$   $\rightarrow$  Set of all strings starting with a 1-digit

$$L(a) = \{a\} \xrightarrow{\text{Reg. set}}$$

neg exp

The Unary

Order

$\Rightarrow$  Precedence  $\rightarrow$  ① The unary operator  $*$  has the highest precedence

& left associative

② Concatenation has second highest precedence & left associative

③ Union has the lowest precedence & left associative

gylext → accepted strings by the pattern  
↳ points to char.

Date .....

Ex - Parenthesizing the reg. exp.  $a b^* c$

$(a | ((b)^* c))$

## Regular Definition

letter → A | B | ... | Z | a | b | ... | z | -

digit → 0 | 1 | 2 | ... | 9

id → letter - (letter - | digit)\*

digits → digit \*

optional fraction → · digits | E

optional Exponent → (E (+) - | E) digits ) | ε

number → digits optional fraction optional exponent

In other way,

letter → [A - Z] | a - z | - | ]

digit → [0 - 9]

id → letter -

digits → digit \*

number → digits (· digits)? (E (+ -) ? digits)?

Indian Mobile Number

→ [6-9] (0-9) {3}

^ → Beginning of a line; \$ → end of the line

^ [a - z] → Beginning of a line

[^ a - z] → Negation of a - z

n<sub>1</sub> / n<sub>2</sub> → n<sub>1</sub> when followed by n<sub>2</sub>

^ [^ a c i o u ] \* \$ → Any line that does not contain any lowercase vowel. Spiral

Date .....

Q Write regular definition for the following languages:-

i) All strings of lowercase letters that contain vowels in order

ii) All strings of lowercase letters in which the letters are lexicographically correct. (in ascending order)

iii) String literal

Ans →

$\text{str} \rightarrow \text{str\_vow\_in\_order} = a^* [b-d]^* e^* [f-h]^* i^* [j-n]^* o^* [p-t]^* u^*$   
 $\text{str} \rightarrow a^* \alpha^* b^* \times i^* \times o^* \times u$   
 $\alpha \rightarrow [^a c j o u]$

$\text{str} \rightarrow \alpha^* \alpha^* \alpha^* \alpha^* \alpha^* i^* \alpha^* \alpha^* \alpha^* \alpha^* \alpha^* \alpha^*$   
 $\alpha \rightarrow [^a e i o u]$

Recommended →

~~Ans~~  
 $\text{str} \rightarrow \text{con}^* a \cdot (\text{con} | \text{al})^* e \cdot (\text{con} | \text{e})^* i \cdot (\text{con} | \text{i})^*$   
 $\text{con} \rightarrow [^a e i o u] \quad \text{or}$   
 $[b-d f-g j-n p-t v-z]$

(b)  $\text{str} \rightarrow a^* b^* c^* \dots r^* z^*$

(c)  $\times \text{ str} \rightarrow [a-z \ A-Z \ 0-9 \ -]^*$   
 $[^"]$

Spiral

Ans → `\"[""]"`

Ans →  
`\((\(\))|[""])*\"`

\* → Single line comment →

`(\//)[^\n]`

On  
`"//"[^\n]`

feb 2<sup>nd</sup>, 2023

In the process of lexical analysis, if it encounters the whitespace, it closes no return to the parser, rather it restarts the process.

Q → Lex program that copies a file, replacing each non-empty sequence of whitespace by a single blank.

Ans  
`#include <stdio.h>`

`1. }`

Date .....

%.%

[" " \t] + { fprintf (yyout, " ") ; }

\n { fprintf (yyout, "%s", yytext); }

%.%

int main() {

yyin = fopen ("input.txt", "r");

yyout = fopen ("output.txt", "w");

yylex();

fclose(yyin);

fclose(yyout);  
} return 0;