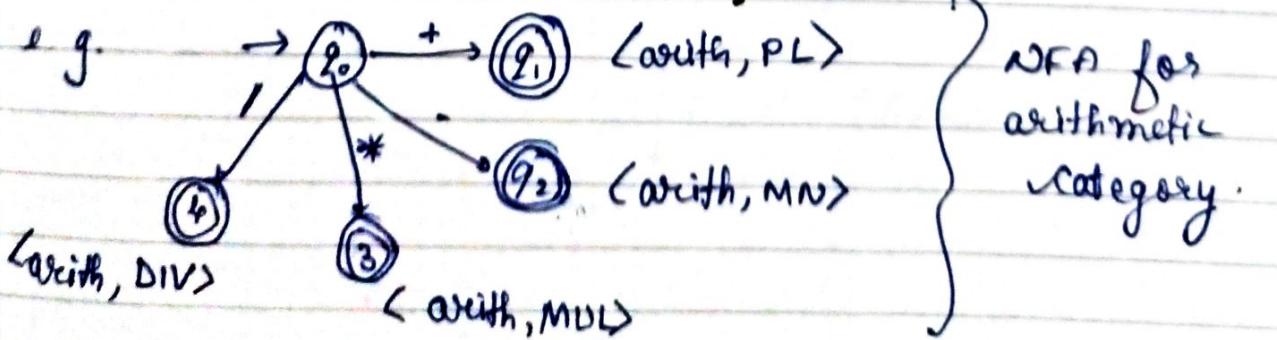


→ Implement NFA for following

CATEGORY	VALID LEXEMES	TOKEN ISSUED
① Keywords		
② Variables		
③ Numbers		
④ Arithmetic	+,-,*,/	e.g. <arith, PL> <sup>Plus</sup>
⑤ Relational Operators		

Therefore, each lexeme belongs to a ~~sector~~ category which is its name.

↓  
token issued



→ Design NFA

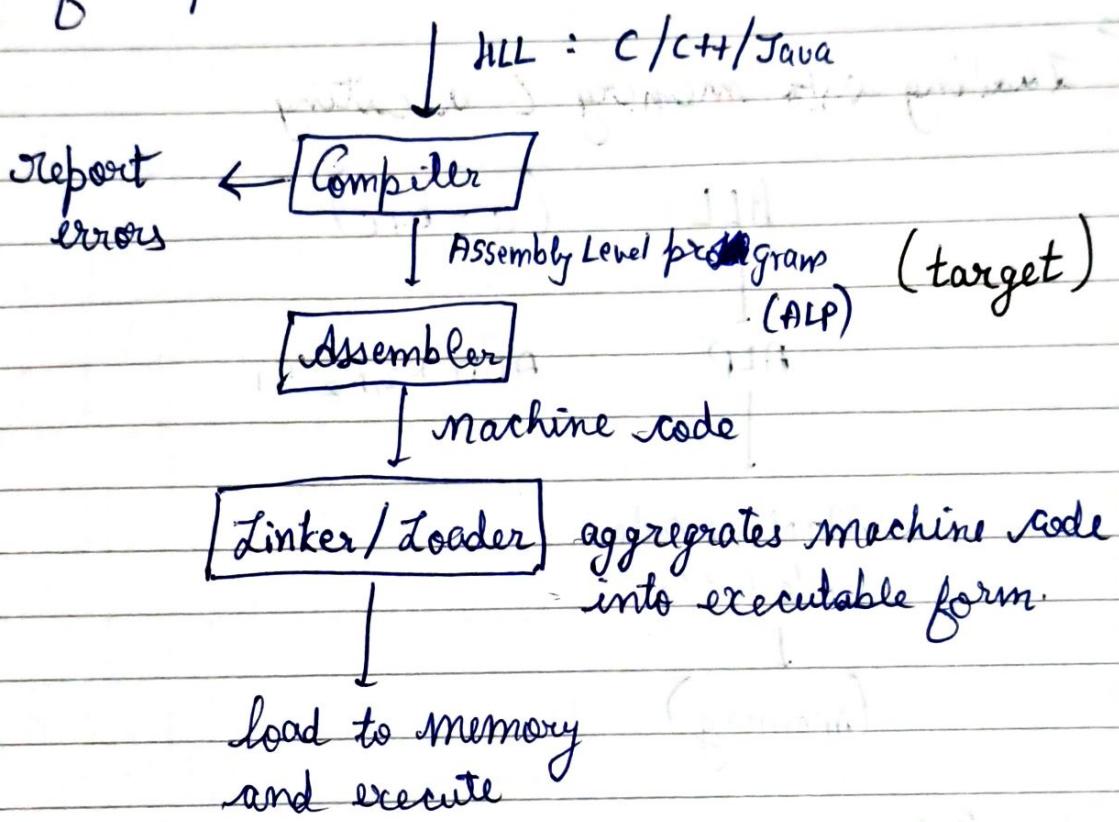
- ① User enters string
- ② Call function NFA
- ③ Return back a signal to main

↳ In case of multiple NFAs, the order in which they are called is important.

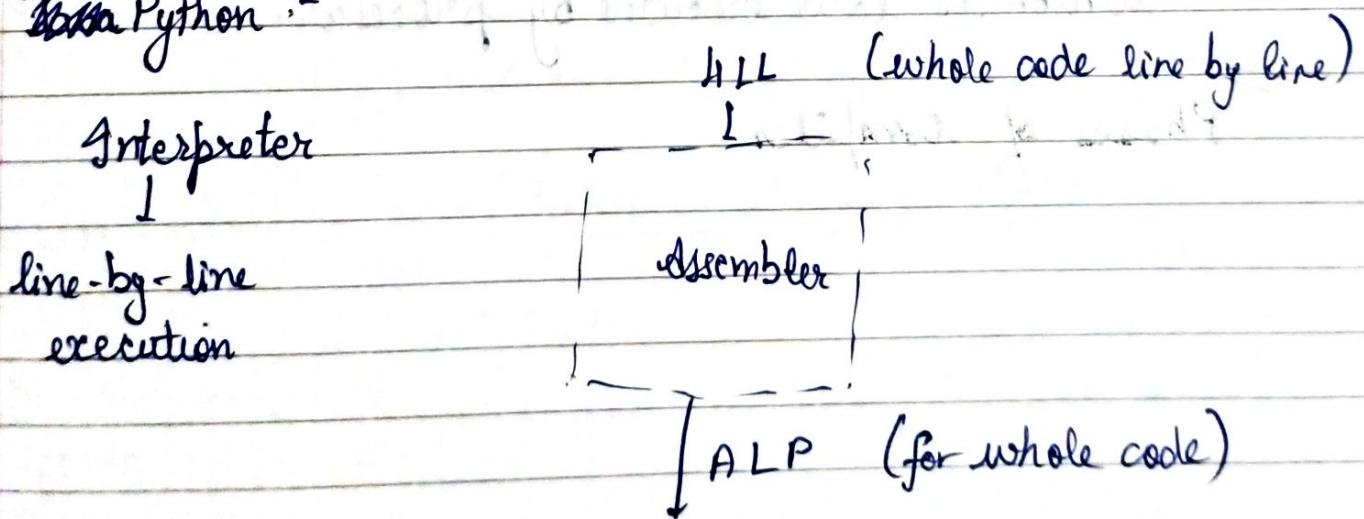
## Introduction

- Compiler : a program written in C that translates a high-level program to assembly language and reports any errors.

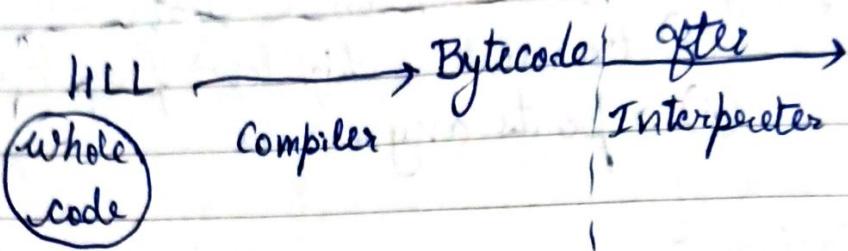
- Process of compiler :



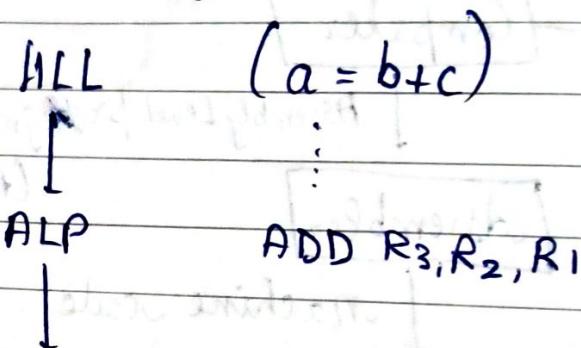
- $\Rightarrow$  ~~Java~~ Python :-



⇒ Java: Hybrid of compiler and interpreter



= Loading into memory & executing



↳ Data as well as program are stored in memory, which is later executed by processor.

= Phases of compiler

classmate  
Date: 2023-09-05

Target:

```

LD R1, 0000 id2
LD R2, 0000 id3
ADD R3, R1, R2
ST id1, R3

```

$a = b + c$

tokens:  $\langle \text{id}, 1 \rangle$   
 $\langle \text{add}, \text{EQ} \rangle$   
 $\langle \text{id}, 2 \rangle, \langle \text{arith}, \text{P} \rangle$   
 $\langle \text{id}, 3 \rangle$

tree:

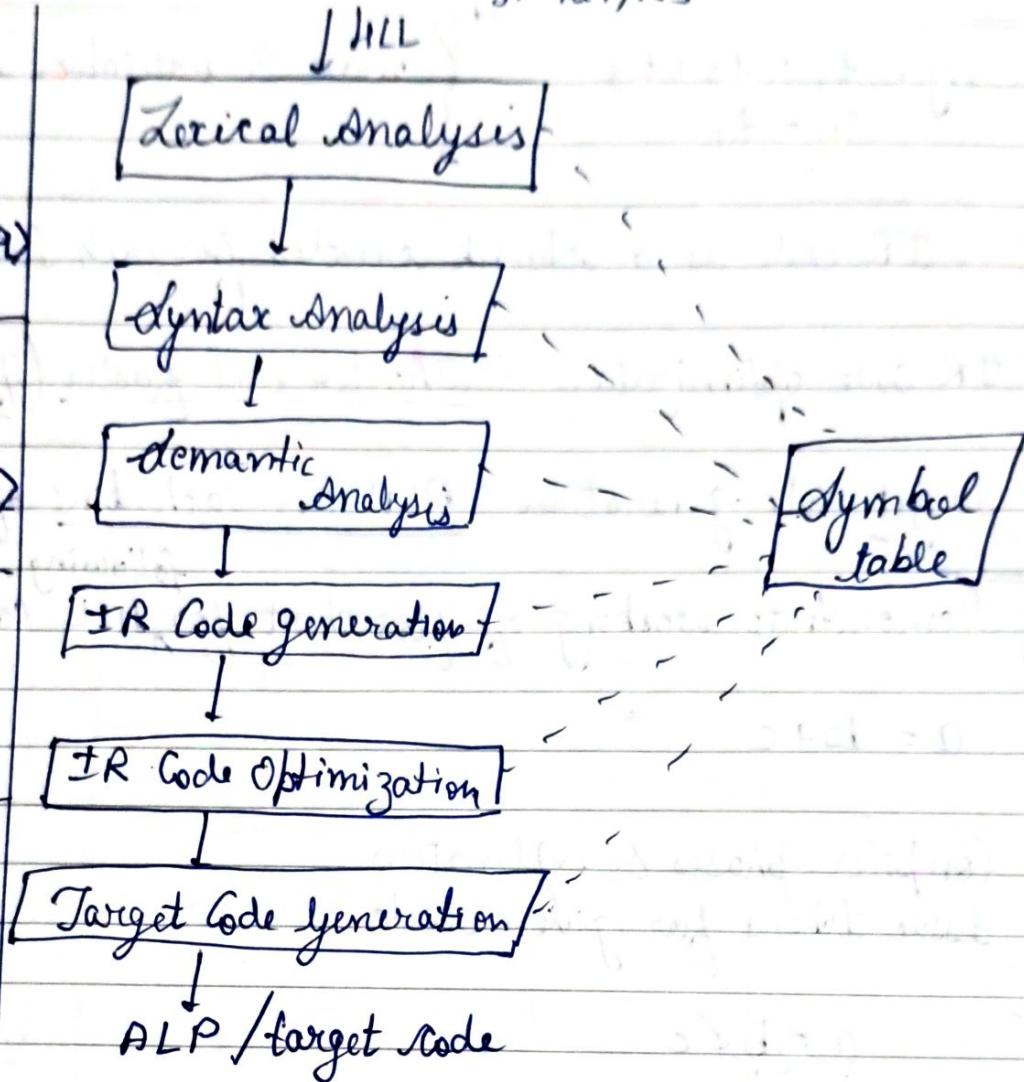
```

    E
   / \
  +   .
 / \   \
<id,1> <id,2> <id,3>

```

tree: -- ditto --

$t_1 = \text{id}_2$   
 $t_2 = \text{id}_3$   
 $t_3 = t_1 + t_2$   
 $\text{id} = t_3$   
 $t_1 = \text{id}_2$   
 $t_2 = \text{id}_3$   
 $\text{id}_1 = t_1 + t_2$



- ① Lexical analysis: divides the input stream of characters into meaningful chunks (lexemes). For each valid lexeme, issue tokens:  $\langle \text{name}, \text{value} \rangle$
- ② Syntax analysis: checks the arrangement of tokens
- ③ Semantic analysis: type checking
  - try type conversion
    - if possible
    - or, declare type mismatch error
  - o.g. 'int' → 'float'
- ④ IR code generation
  - Intermediate representation = abstract code = 3 address code

④ Symbol table : stores info on variables

value      ↑  
type      scope  
size

e.g.  $t_1 = t_2 + t_3$

$a = t_1$

} max 3 variables in a statement

∴ IR code is a chunk created for each line of HLL.

⑤ IR code optimization : shorter and faster (if possible)

⑥ Target code generation : ALP for each line of IR code

\* Q. Demonstrate working of flowchart for HLL code

$a = b + c$

Q2 Compiler phases & explanation

Q3 Issue tokens for given code.

~~a = b + c~~

tokens M6

implemented by nfa  
↓

Lexical Analysis → meaningful chunk

Lexical categories

Valid lexemes

Tokens issued  
only kw

1. Keywords

if, while etc.  
NFA1 NFA2

must start with letter

x xy x12 ...  
single NFA

identifier order in which it comes

<id, 1>, <id, 2>, ...

3. Numbers

1 1000 ...

<num, 1>, <num, 1000>

4. Relational

>, >=, =, <, <=, <>

<relOp, GT>, <relOp, GE>,  
<relOp, EQ>, <relOp, LT>,  
<relOp, LE>, <relOp, NE>

5. Arithmetic

+, -, \*, /,

<arith, PL>, <arith, MN>  
<arith, MU L>, <arith, DIV>

6. Punctuation

, ; : #

<punc, com>, <punc, semi>,  
<punc, col>

7. Parentheses

( ) { } [ ]

<par, S-op>, <par, S-cl>,  
<par, C-op>, <par, C-cl>,  
<par, B-op>, <par, B-cl>

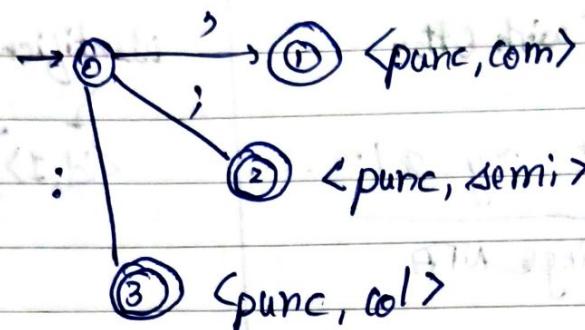
8. White space

blank, newline, tab

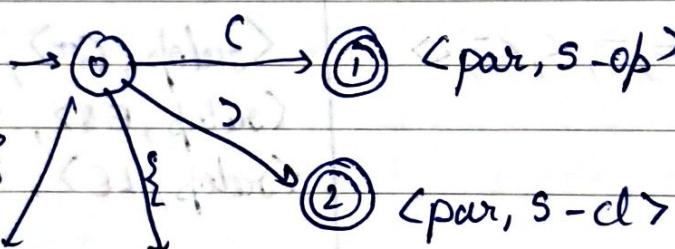
<> (empty token)

- Task is to issue tokens
- nfa issues tokens
- our task is to design nfas which issue tokens

e.g. punctuation nfa :-



e.g. parenthesis



NFA for keyword & variables

e.g.

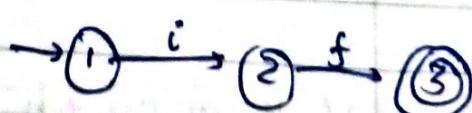
Keyword

"if"

Variable

iff, if 123

⇒ if we design nfa as



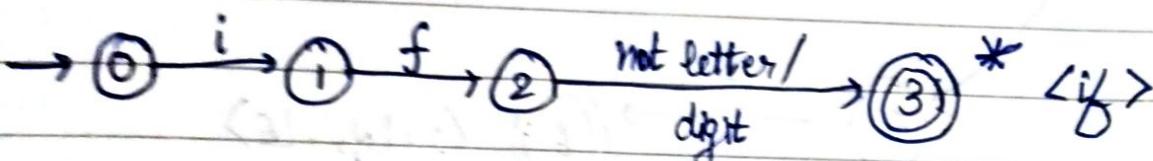
'if' will emit <if>

then even

∴ Therefore, we design nfa for each category, and when nfa accepts then break out.

\* order of keyword → variable → number → relOp is important.

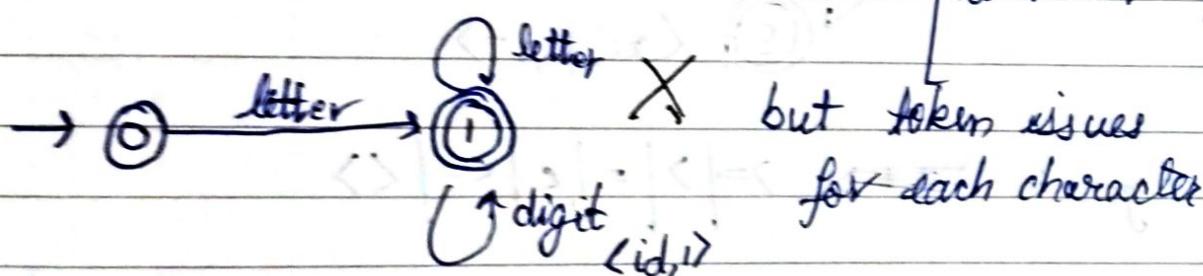
∴ but if source is 'if' we don't want keyword nfa to accept it



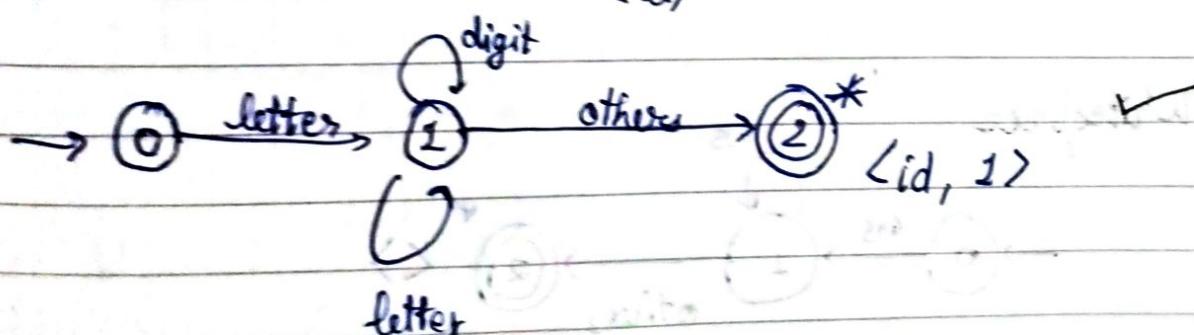
←

→ retract or  
go back

→ Variable :-



letter = a1...13/A1...12



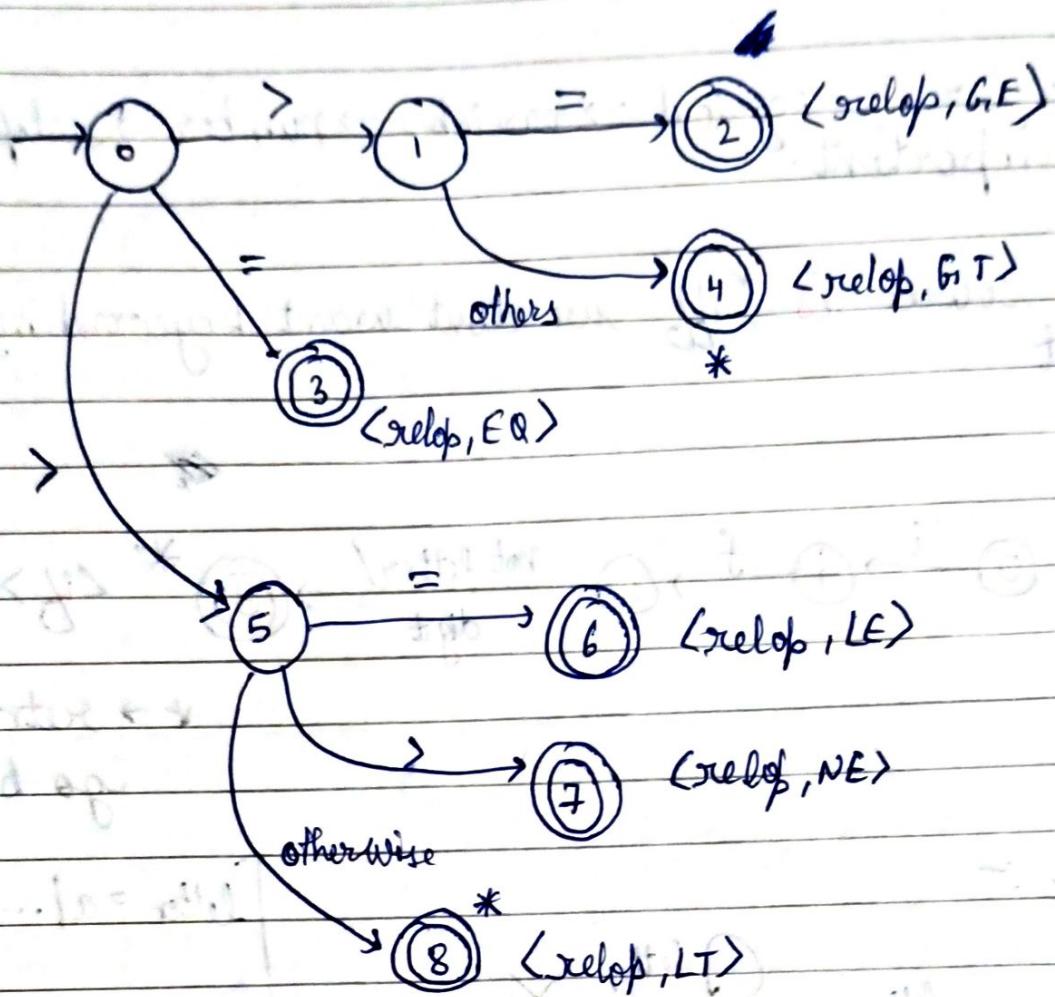
=  
Regular expression : pattern we expect

'if' : if

$\langle id, 1 \rangle$  : letter . (letter + digit)\*

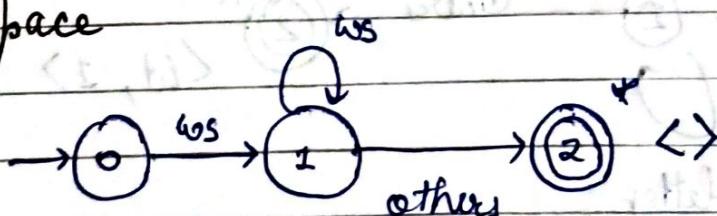
concat  
• →  $\text{letter}$   
+ → or  
\* → lots of

## → Relational Operators :-



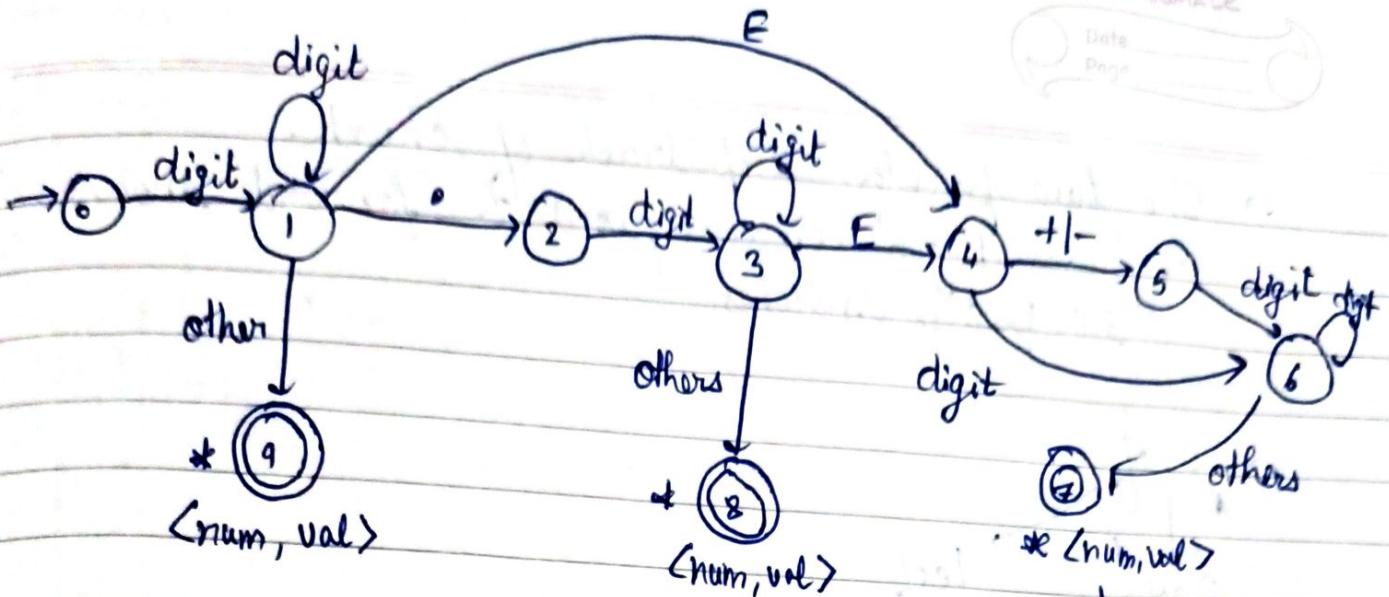
Regex : ~~>=|>|=|<|=|<=|<>~~

## → Whitespace



## → Numbers

~~1024~~



where  $\text{digit} = 0 | 1 | \dots | 9$

$$y = -3.3$$

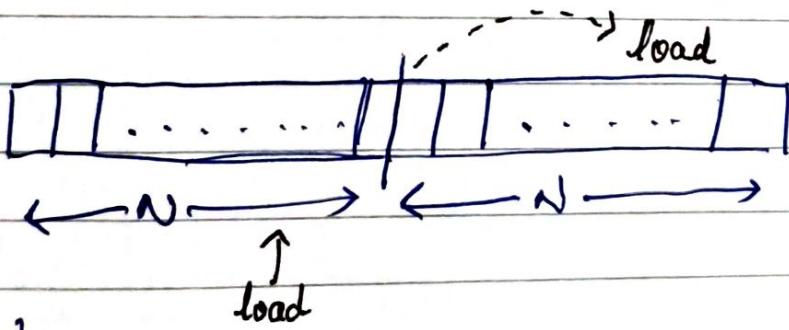
$\langle \text{id}, 1 \rangle \langle \text{relOp}, E0 \rangle \langle \text{arith}, MN \rangle, \langle \text{num}, 3.3 \rangle$

3  
 306  
 300.3  
~~3E+4~~  
~~3E4~~  
~~3.3E4~~

## Input Buffering

↳ 2 buffers

↳ size  $N = 4096$  bytes ~~bytes~~ i.e. 4096 characters

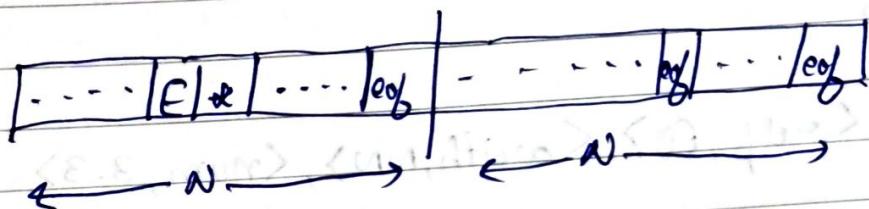
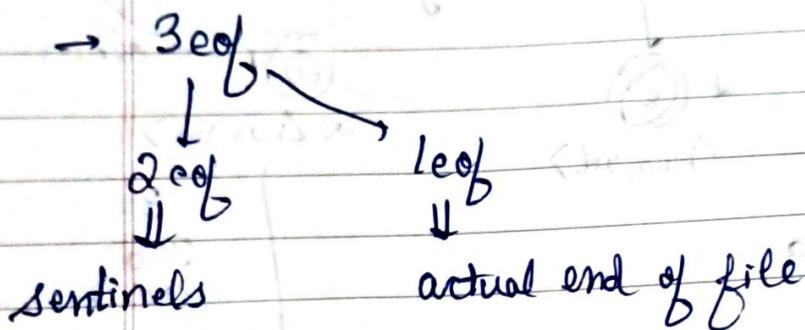


① size  $N$

② alternative loading

③ 2 pointers : `lexemeBegin`, `forward`  
`(Fix)`                    `(Moves)`

- The two pointers keep track of chunk.  
∴ 'lexemeBegin' remains fixed & 'forward' moves ahead till pattern matches

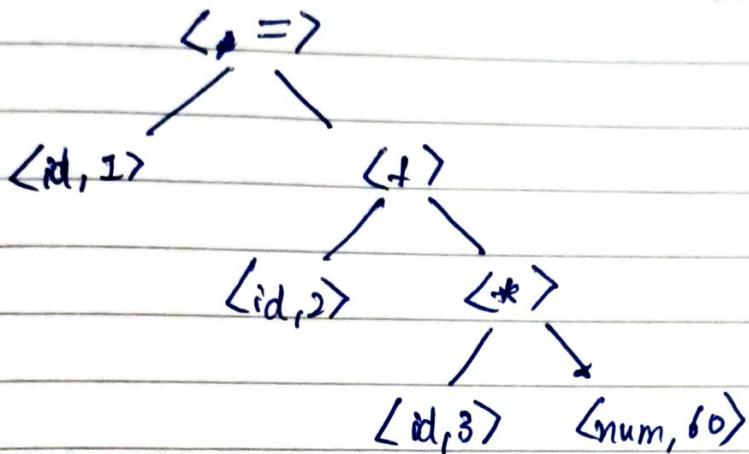


- when actual of encountered, the procedure ends.
  - when sentinel is reached, next buffer is loaded and forward continues its journey till pattern matches.

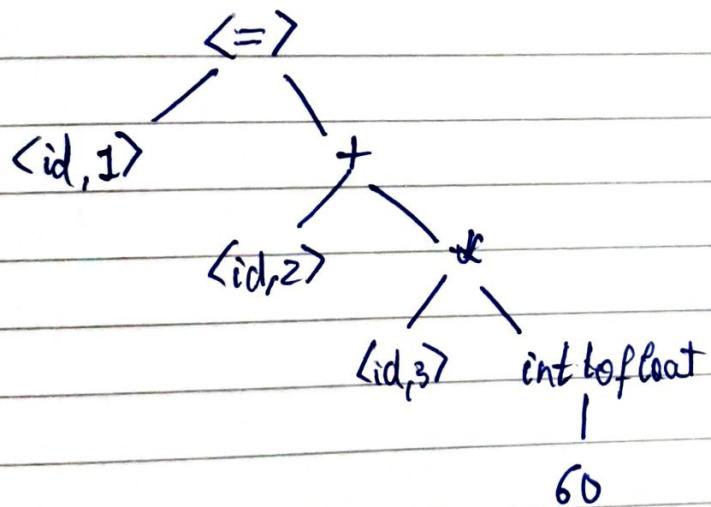
$$① pos = in + rate * 60$$

②  $\langle id, 1 \rangle, \langle \text{relOp}, = \rangle, \langle id, 2 \rangle, \langle \text{arith}, PL \rangle, \langle id, 3 \rangle, \langle \text{arith}, MUL \rangle,$   
 $\langle \text{num}, 60 \rangle$

③



④



$$t_1 = \text{inttofloat}(60)$$

$$t_2 = id3 * t1$$

$$t_3 = id2 + t2$$

$$id1 = t3$$

$$⑤ t_2 = id3 * \text{inttofloat}(60)$$

$$id1 = id2 + t2$$

⑥ LDF R<sub>1</sub>, id 3

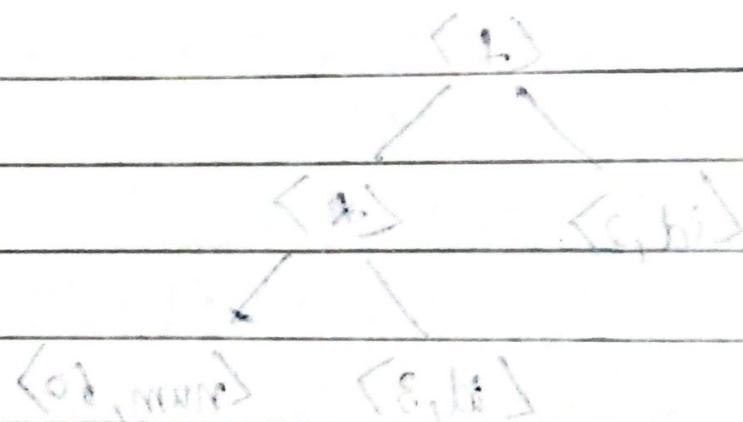
MULF R<sub>2</sub>, R<sub>1</sub>, #60.0

~~BBDE~~

LDF R<sub>3</sub>, id 2

ADDF R<sub>4</sub>, R<sub>3</sub>, R<sub>2</sub>

STF id1, R<sub>4</sub>

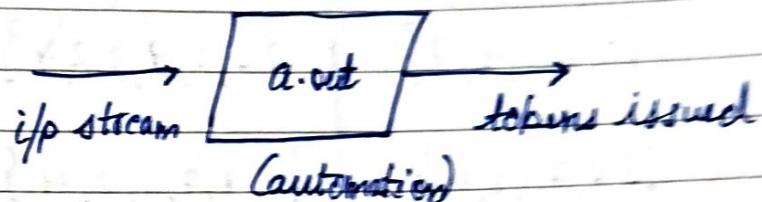
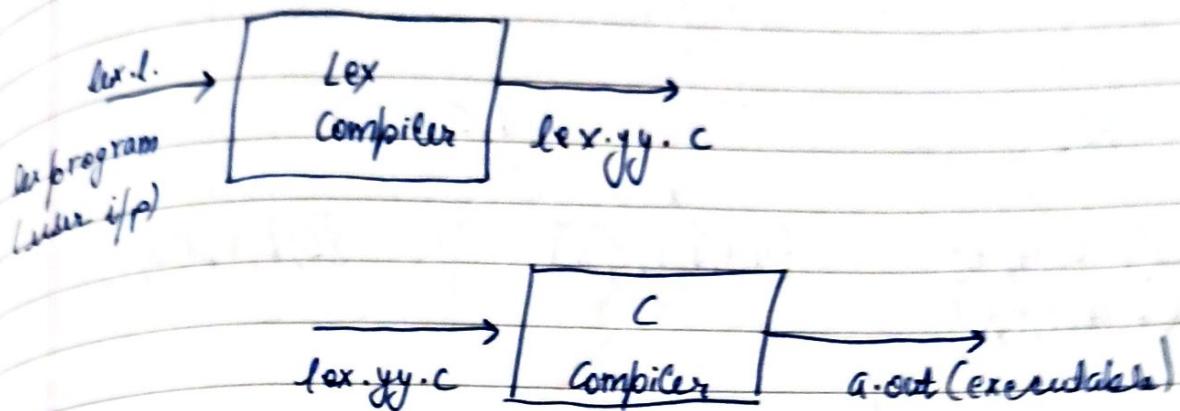


total bits 500

- \* 1. Lexical categories - NFA
- \* 2. I/P buffering
- \* 3. Thompson algorithm & subset construction
- \* 4. Lex

6.2.25

Lex : Automatic lexical analyzer generator



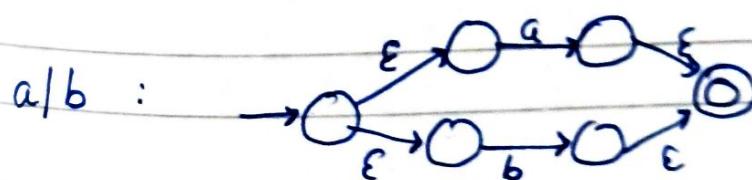
Thompson algo & subset construction

g.r.e.  $\rightarrow$  nfa

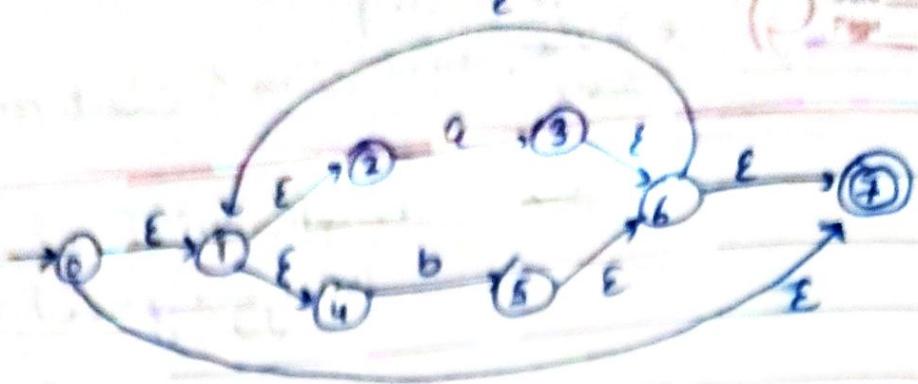
nfa  $\rightarrow$  dfa

- i) Construct nfa from g.r.e. using Thompson
- ii) convert that nfa to dfa using subset construction

$$\text{g.r.e.} : (a/b)^*$$



$(a/b)^*$



P/S.

$$\{0, 1, 2, 4, 7\}$$

$$\{4, 6, 7, 1, 2, 3\}$$

a

$$\{4, 6, 7, 1, 2, 3\}$$

N.S.

b

$$\{5, 6, 7, 1, 2, 3\}$$

P.S.

$$\{0, 1, 2, 4, 7\}$$

a

$$\{3, 6, 7, 1, 2, 4\}$$

b

$$\{5, 6, 7, 1, 2, 4\}$$

$$\{3, 6, 7, 1, 2, 4\}$$

a

$$\{3, 6, 7, 1, 2, 4\}$$

b

$$\{5, 6, 7, 1, 2, 4\}$$

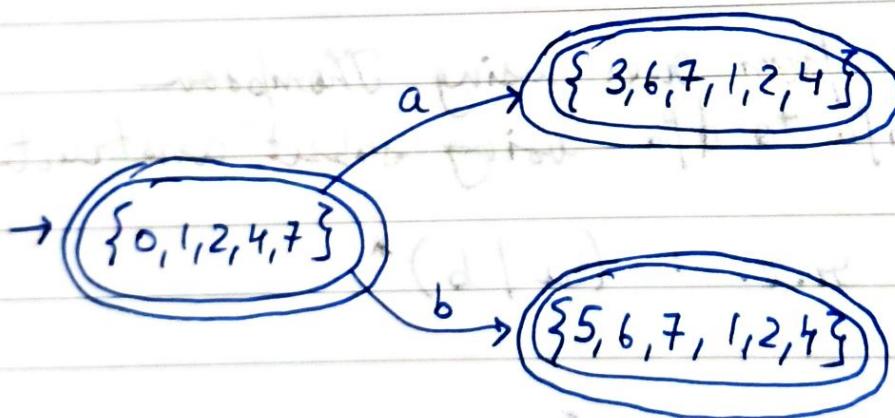
$$\{5, 6, 7, 1, 2, 4\}$$

a

$$\{3, 6, 7, 1, 2, 4\}$$

b

$$\{5, 6, 7, 1, 2, 4\}$$



\* Must draw rules before answering

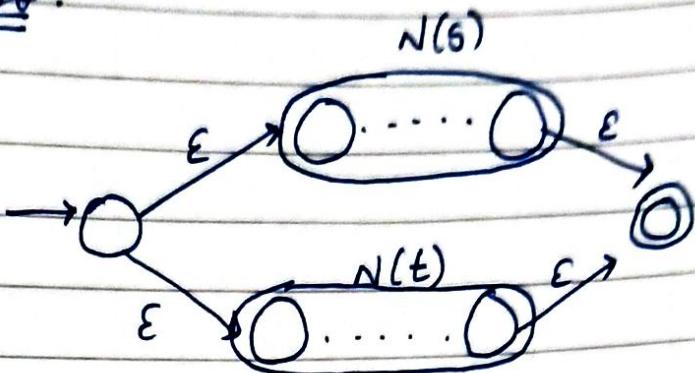
classmate

Date \_\_\_\_\_

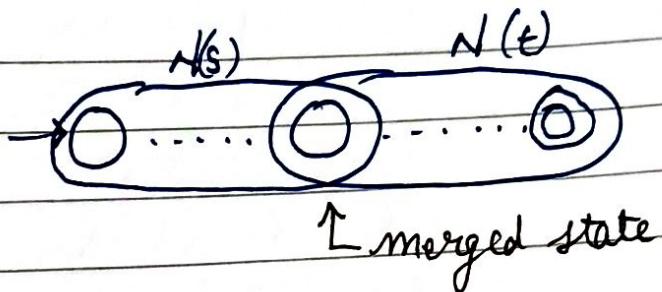
Page \_\_\_\_\_

Rules for Thompson algo :-

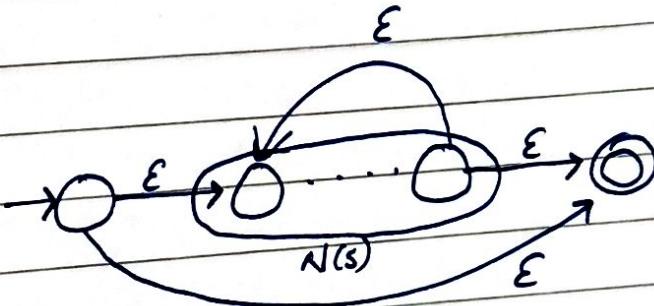
(1)  $s/t$



(2)  $s \cdot t$



(3)  $s^*$



→ Transition table of ~~DFA~~: (i) <sup>no</sup> blanks  
(ii) one entry per cell

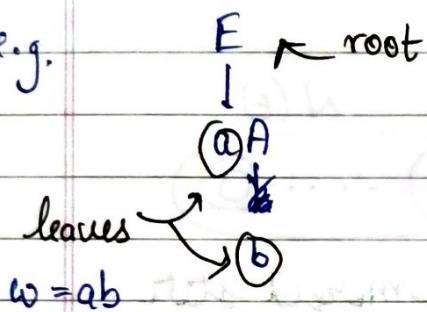
## Syntax Analysis (arrangement of tokens)

→ Parser: program which says yes/no for a stream.

Q Design a parser given grammar G & string w

~~e.g.~~ ↳ if parser gives green flag, we can make tree.

e.g.



Parser

Top down

e.g. LL(1)

Bottom-up

e.g. LR(0), SLR(1)

→ Designing LL(1) parser :

↳ given G & w

e.g.

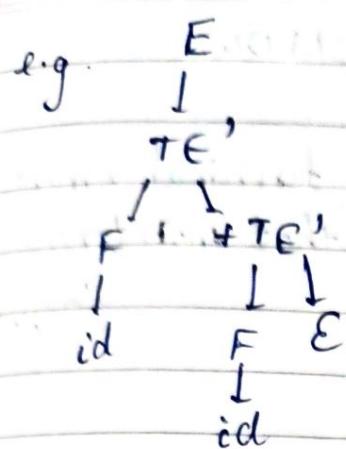
$$\begin{aligned} G : \quad E &\rightarrow TE' \\ F' &\rightarrow +TE' / \epsilon \\ T &\rightarrow F \\ F &\rightarrow (E) / \text{id} \end{aligned}$$

(6 productions)

$$\begin{aligned} \therefore \text{terminals} : \quad &\{ \text{id}, (, ) , + \} \\ \text{non-terminal} : \quad &\{ E, T, E', F \} \\ \text{start} : \quad &E \end{aligned}$$

↳ info about identifier is stored in symbol table

classmate



$$w = id + id \notin$$

↑ EOF  
(mandatory)

Step-1 : parsing table (from grammar)  
Step-2 : table of sequence of moves (from table & i/p string)

## FIRST

↳ first terminal you see from L to R  
(ε allowed)

## FOLLOW

↳ which terminal immediately follows that non-terminal.  
(ε not allowed)

① Let  $X \rightarrow Y_1, Y_2 \dots Y_K$

↳ then,  $\text{FIRST}(X) = Y_1$ , if  $Y_1$  is terminal  
=  $\text{FIRST}(Y_1)$  otherwise

↳ if  $\text{FIRST}(Y_i) = \epsilon$ , then,  $\text{FIRST}(X) = \text{FIRST}(Y_2)$

...  
if all are ε, then  $\text{FIRST}(X) = \{\epsilon\}$

② FOLLOW has 3 rules

(i)  $\text{FOLLOW}(E) = \{\$\dots\}$

(ii) for production  $X \rightarrow \alpha B \beta$ , then

$\text{FOLLOW}(B) = \text{FIRST}(\beta)$

(iii) if  $\text{FIRST}(\beta) = \{\epsilon\}$ ,  
 $(= X \rightarrow \alpha B)$

$\text{FOLLOW}(B)$  will inherit from  $\text{FOLLOW}(X)$

## ① LL(1) Parsing table :-

	$\text{id}$	$($	$)$	$+$	$\$$
$E$	$E \rightarrow TE'$	$E \rightarrow TE'$	-	-	-
$T$	$T \rightarrow F$	$T \rightarrow F$	-	-	-
$E'$	-	-	$E' \rightarrow E$	$E' \rightarrow +TE'$	$E' \rightarrow \epsilon$
$F$	$F \rightarrow \text{id}$	$F \rightarrow (E)$	-	-	-

process  
for  
designing  
LL(1) for each production  $A \rightarrow \alpha$

1. enter  $A \rightarrow \alpha$  in row ( $A$ ) & column : FIRST ( $\alpha$ )
2. if  $\text{FIRST}(\alpha) = E$ , then enter  $A \rightarrow \alpha$  in FOLLOW ( $E$ ) column & row ( $\alpha$ )

e.g.  $E \rightarrow TE' \therefore \text{row} = E, \text{column} = \text{FIRST}(TE') = \text{FIRST}(T)$   
 $= \text{FIRST}(F) = \{C, \text{id}\}$

e.g.  $E' \rightarrow \epsilon, \text{FOLLOW}(E') = \text{FOLLOW}(E)$   
 $= \{\$, )\}$

Sequence of moves ( $\omega = \text{id} + \text{id} \$$ )

Stack	Input	Comments
$E \$$	$\text{id} \rightarrow \text{id} \$$	$(E, \text{id}) \text{ all}$
$\uparrow$	$\uparrow$	
$T E' \$$	$\text{id} + \text{id} \$$	$(T, \text{id}) \text{ cell}$
$\uparrow$	$\uparrow$	
$FE' \$$	$\text{id} + \text{id} \$$	$(F, \text{id}) \text{ cell}$
$\text{id} E' \$$	$\text{id} + \text{id} \$$	MATCH!
$E' \$$	$+ \text{id} \$$	$(E', +) \text{ cell}$
$+ TE' \$$	$+ \text{id} \$$	<del>MATCH!</del>
$TE' \$$	$\text{id} \$$	$(T, \text{id}) \text{ cell}$
$FE' \$$	$\text{id} \$$	$(F, \text{id}) \text{ cell}$
$\text{id} E' \$$	$\text{id} \$$	MATCH!
$E' \$$	$\$$	$(E', \$) \text{ cell}$
$\$$	$\$$	YES

Note: in the above procedure, you will get syntax error if

- (i) you don't get  $\$ | \$$  YES
- (ii) You reach empty cell
- (iii) mismatch of terminals.