

- HIGH LEVEL LANGUAGE → COMPILERS → LOW LEVEL LANGUAGE
(program) / MACHINE CODE
- BATCHWISE PROCESSING OF PROGRAMS

[ASSEMBLY LANGUAGE]

[PUNCH CARDS]

- JOHN BACKUS → COMMERCIAL COMPILERS → 1957 → 18 PERSON YRS.

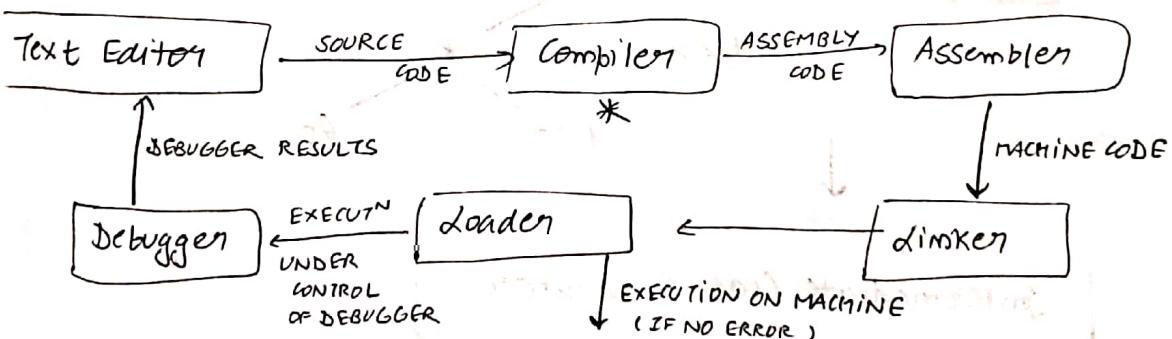
FORTRAN - I

- Assemblers are usually built into compilers

- compilers should have foll. criteria

- 1) convert one representation to another
- 2) Meaning should remain the same.

THE BIG PICTURE



* Compiler is a Translator

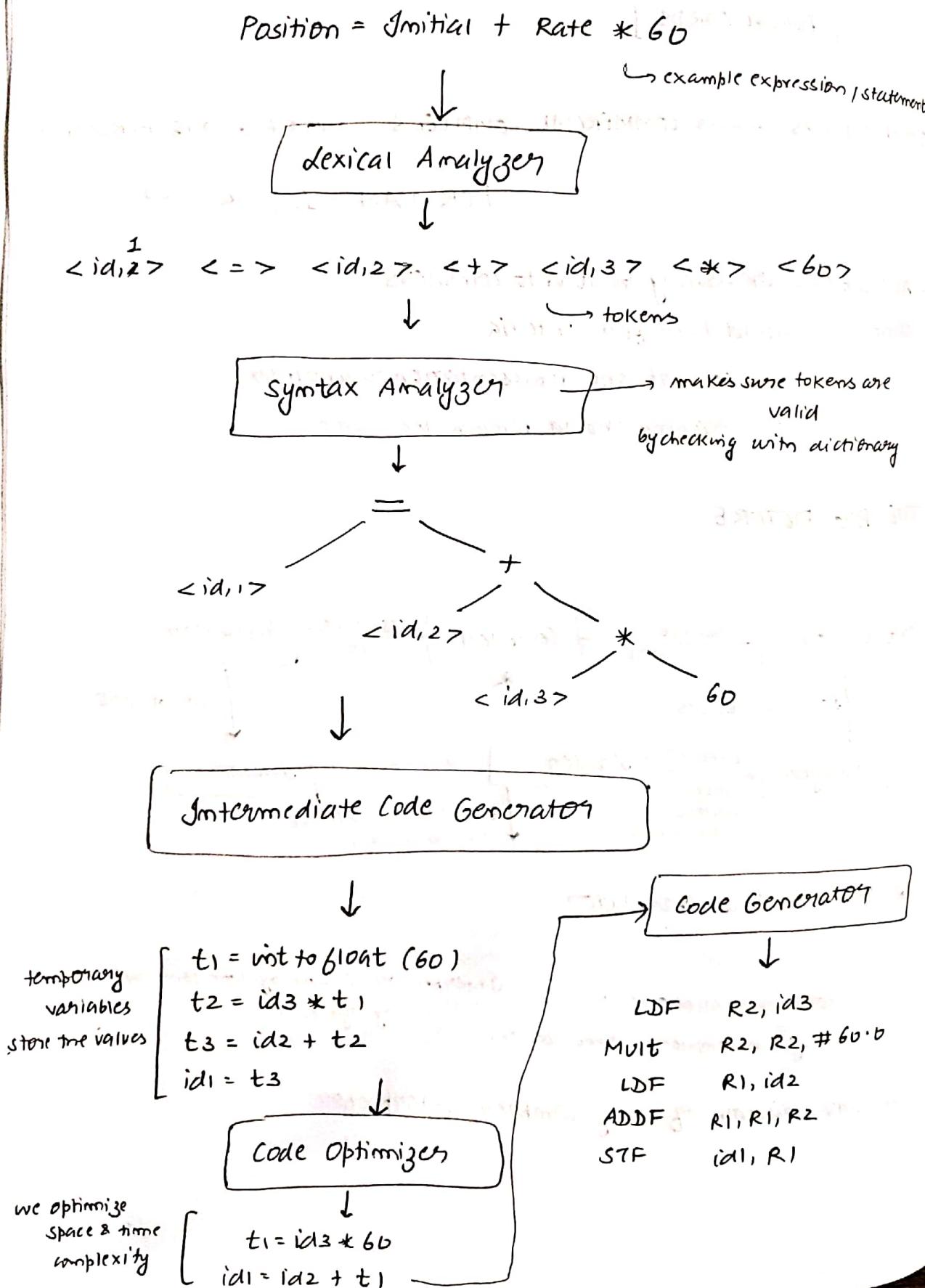
Compiler (all at once conversion)
eg: C# compiler, TurboC, DevC++

Interpreter (Line-by-line conversion)
eg: Python, R

HW adv., disadv. of using compiler, interpreter

- We must know about variables, expressions, constants. Impression of compiler, registers, memory etc. are not to be remembered while programming.

PHASES OF A COMPILER



- yacc, lex are used to build a compiler

COMPILERS : PRINCIPLES, TECHNIQUES & TOOL, ALFRED V ALIB, MONIEN,
ULLMAN

- each word in expression is mapped to symbol table by lexical analyzer during token generation.

ADVANTAGES OF COMPILER OVER INTERPRETOR

- Traditional compiler translates source code into directly executable machine instructions., once & for all.
- since compilation is one timer, compilers usually put more effort into optimizations & thus execute faster even after translation.
- compilers can perform global optimizations.
 - Code uses less CPU cycles,
 - runs faster & uses less power

DISADVANTAGES OF INTERPRETER

- An interpreter reads character by character & has to re-do all the work every time user wants to run a program through it.
- Interpreters translate only the part of program it is currently executing, thus they can't perform global optimizations.
- Many bugs (like type errors) are encountered by an interpreter program at run-time, ~~thus many bugs are delivered to customer~~

ADVANTAGES OF INTERPRETER OVER COMPILER

- source code is multi-platform, compiled code is processor specific
- It's usually easier to develop & debug in an interpreted environment
- Interpreters are easier to use, particularly for beginners, since errors are immediately displayed.

DISADVANTAGES OF COMPILER

- Hardware specific source code, :- multiple versions of source code must be maintained.
- Large applicatⁿs can take significant amt. of time to compile.
- difficult to debug code that is optimized.

LEXICAL ANALYSIS

#1

Generates a series of tokens

→ identifiers → variables, keywords

operators

whitespaces

Tabs

int a,b,c;

SYNTAX ANALYSIS

#2

Analysis of structure

We is going to school → syntactically correct, but not semantically

int a,b,c → analysis is done using parse tree

SEMANTIC ANALYSIS #3

Understanding the meaning

int a = b + c;
 ↑ ↑ ↑
 int float char

↑ ↑
typechecking

} Removes all the code ambiguity

eg Bytocode generation
in Java

INTERMEDIATE CODE GENERATION

(3 ADDRESS CODES)

#4

OPTIMIZATION (OPTIONAL) #5

TARGET MACHINE CODE GENERATION #6

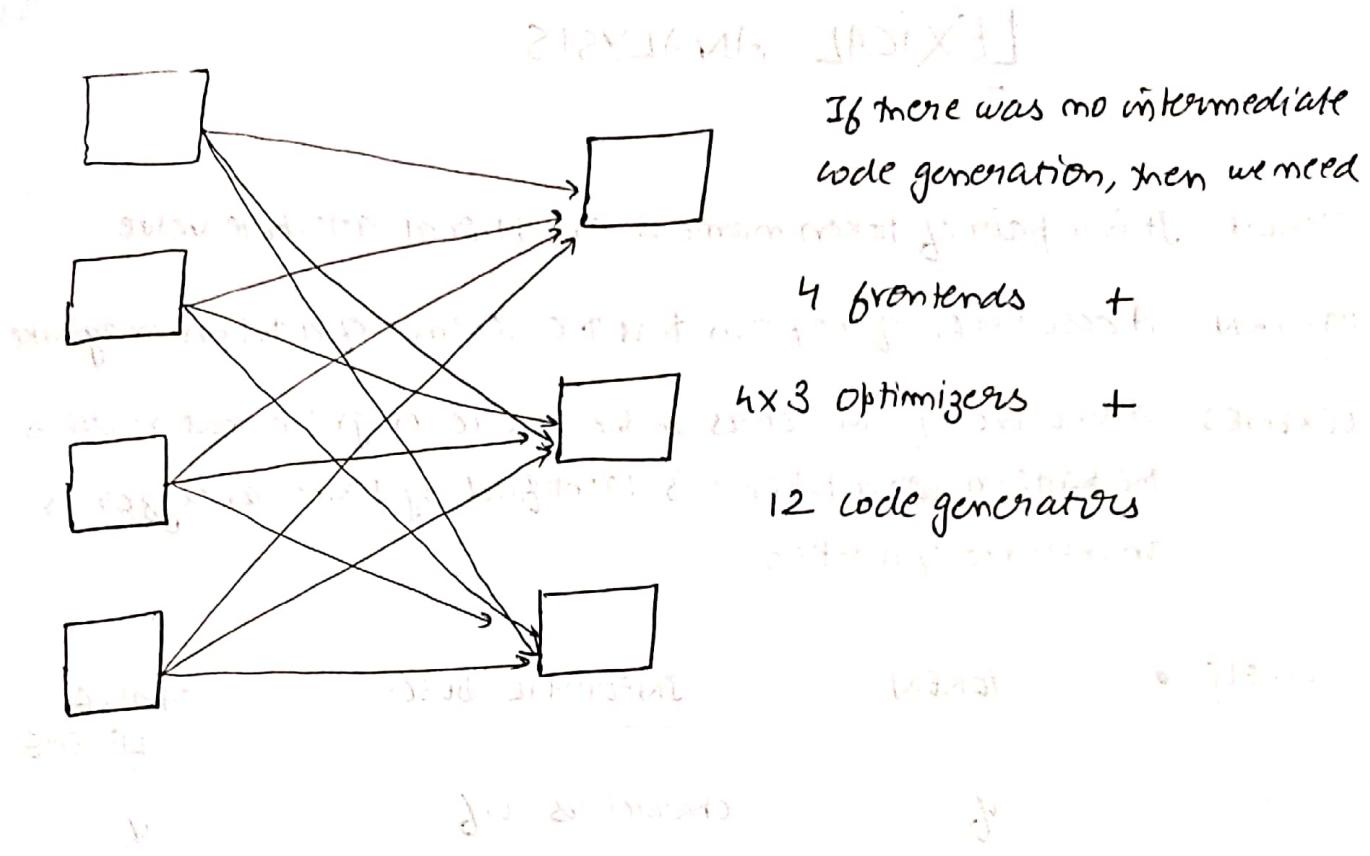
(Assembly lang. code)

#1, #2, #3 come under analysis phase (machine independent)
front end

#5, #6 comes under synthesis phase / back end (machine dependent)

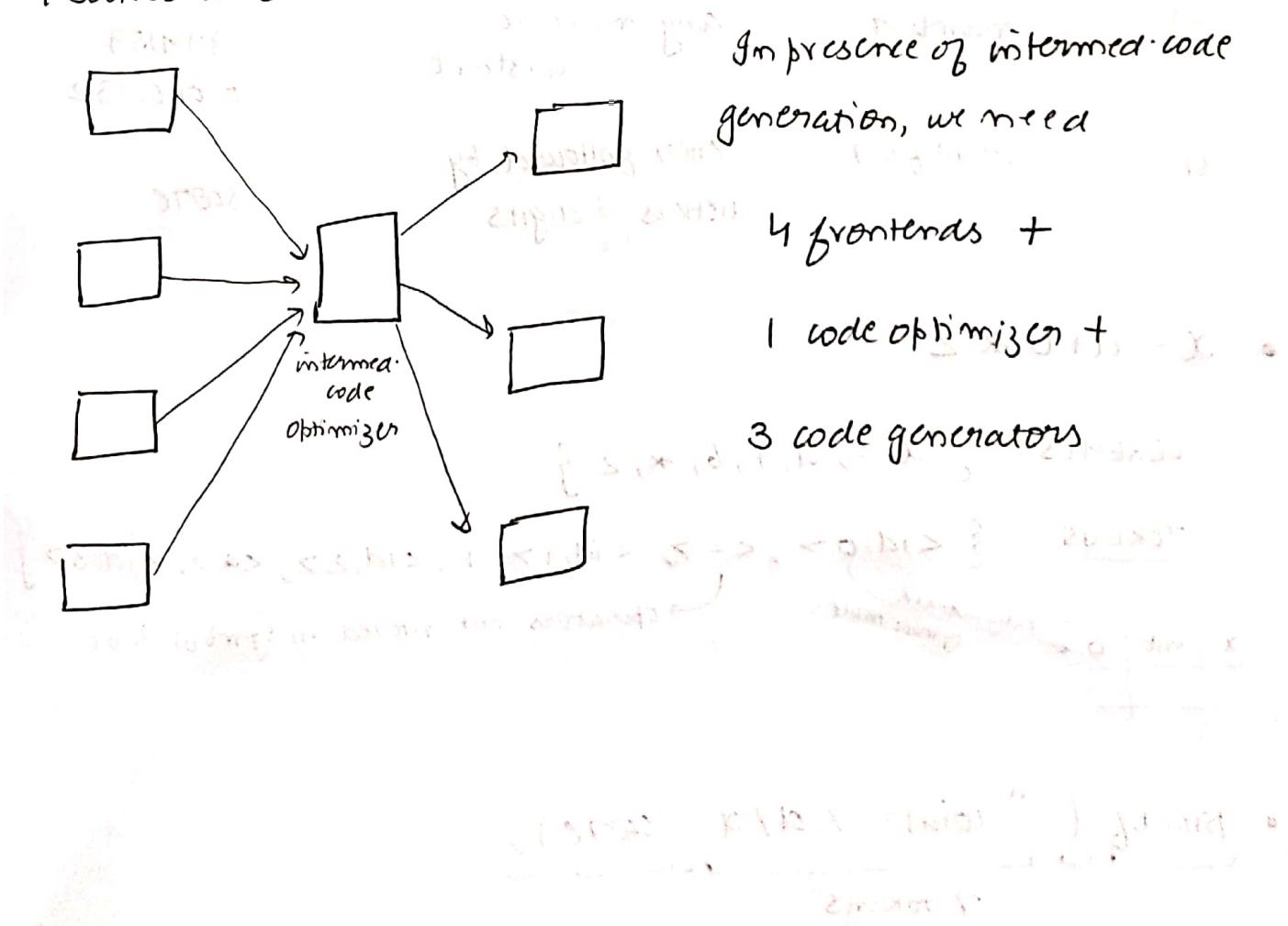
4 SOURCE LANG.

3 TARGET MACHINES



4 SOURCE LANG.

3 TARGET M/C



LEXICAL ANALYSIS

TOKEN: It is a pair of token name & an optional attribute value

PATTERN: A description of the form that the lexemes of a token may take.

LEXEMES: A sequence of characters in the source program that matches the pattern for a token & is identified by lexical analyzer as an instance of a token.

<u>EXAMPLE</u>	<u>TOKEN</u>	<u>INFORMAL DESC.</u>	<u>SAMPLE LEXEME</u>
1)	if	characters i, f	i
2)	number	any numeric constant	3.14159 0.0000132
3)	identifier	letter followed by letters & digits	score
•	$x = a + b * 2$		

LEXEMES: { $x, =, a, +, b, *, 2$ }

TOKENS: { $\langle id, 0 \rangle, \langle = \rangle, \langle id, 1 \rangle, +, \langle id, 2 \rangle, \langle * \rangle, \langle id, 3 \rangle$ }

x	int	0
:	:	
:	:	
:	:	

entry index in symbol table

operators not entered in symbol table

• printf ("Total = $y \cdot a / x$ ", score);
7 tokens

Lexical Analyzer performs the following tasks in order

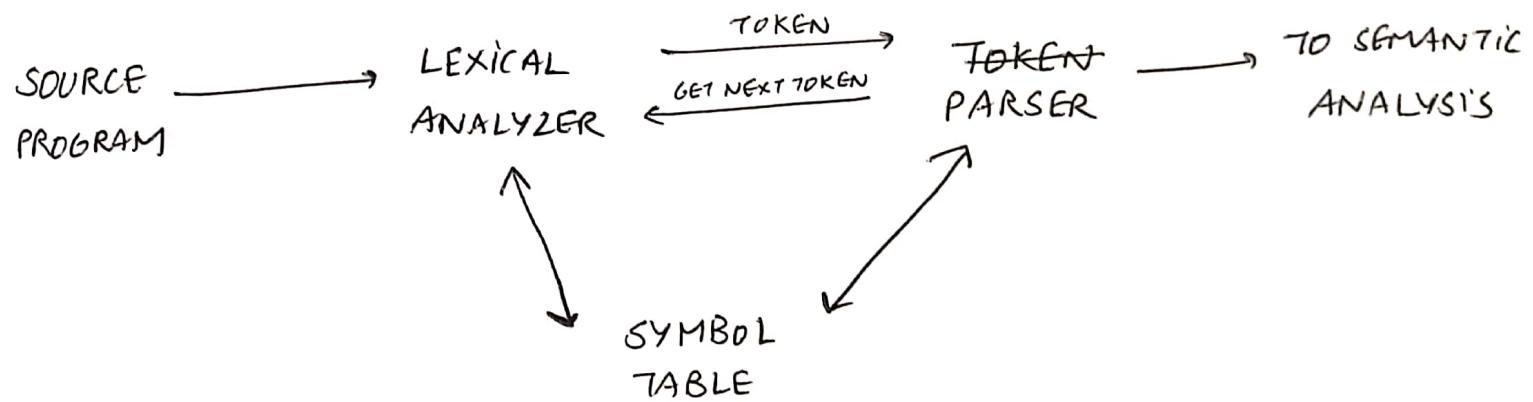
1) Scanning

2) Lexical Analysis

LEXICAL ERROR

$f_i(a == x)$

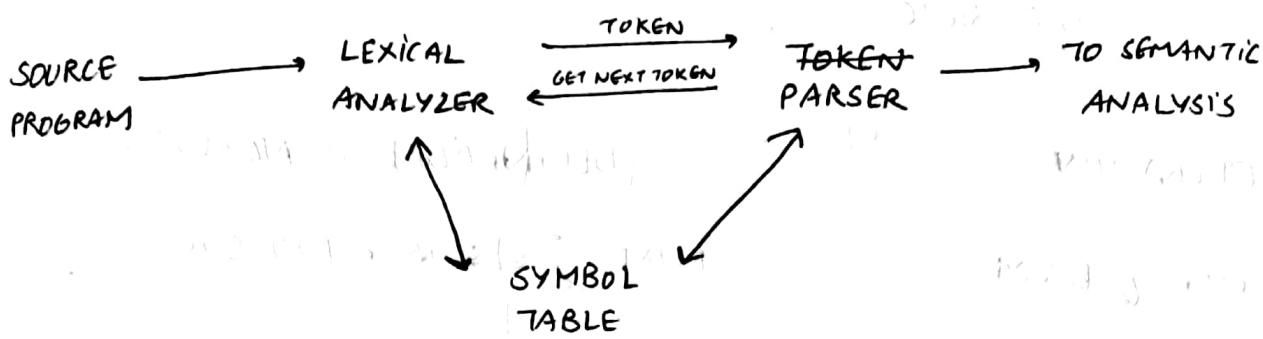
considered as an identifier but later on will display errors as this is a function name. Lexical analyzer won't display this as an error



LEXICAL ERROR

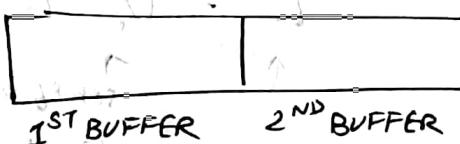
$f_i(a == x)$ → stored in secondary memory

considered as an identifier but later on will display error as this is a function name. Lexical analyzer won't display this as an error

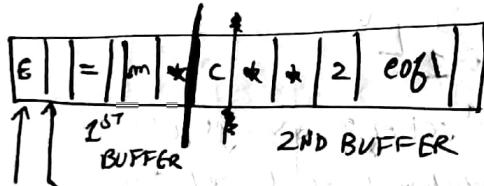


TOKEN: { token name, optional attribute }

INPUT BUFFERING → TWO Buffer Scheme



Eg: $E = m * c * * 2$



forward ptr moves to word character by character while lexeme begins. pointer remains static.

- All Keywords are stored in symbol table.

LEXICAL ANALYSIS → SPECIFICATION → Regular expressions

LANGUAGE → alphabet → set of symbols
finite
eg {0,1} ASCII, UNICODE

Keyword, variable, digit
recognition patterns

↪ finite sequence of symbols drawn from the alphabet

any countable set of strings over that alphabets
some fixed

OPERATIONS OF LANGUAGES

- 1) Union
- 2) Concatenation
- 3) Closure

$\epsilon \rightarrow$ Empty string

OPERATION

Union of L & M

DEFINITION & NOTATION

$$L \cup M = \{ s \mid s \text{ is in } L \text{ or } s \text{ is in } M \}$$

Concatenation of L & M

$$L \cdot M = \{ st \mid s \text{ is in } L \text{ & } t \text{ is in } M \}$$

Kleene closure of L

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

set of all strings obtained by concatenating L, zero or more times

Positive closure of L

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

set of all strings obtained by concatenating L ≥ 1 times

A notation called Regular Expressions is commonly used for describing all the languages that can be built from these operators applied to the symbols of some alphabet.

TWO RULES THAT FORM THE BASIS OF DEF^N OF REG. EXPR. (RE)
OVER SOME ALPHABET Σ & THE LANGUAGES THAT THOSE
EXPRESSIONS DENOTE :

- 1) ϵ is a RE & $L(\epsilon)$ is $\{\epsilon\}$, i.e. language whose sole member is the empty string. → string whose length = 0, it's not null

2) If ' a ' is a symbol in Σ , then ' a ' is a RE & $L(a) = \{a\}$, i.e. the language with one string of length one.

Suppose M & S are RE denoting the languages $L(M)$ & $L(S)$

respectively then,

- $(M) | (S)$ is a RE denoting the language $L(M) \cup L(S)$
- $(M)(S)$ is a RE denoting the language $L(M) L(S)$
- $(M)^*$ is a RE denoting $(L(M))^*$
- (M) is a RE denoting $L(M)$

ORDER OF PRECEDENCE : 1) * }
 2) Concatenation }
 3) Union }

EG: Let $\Sigma = \{a, b\}$

- The RE $a|b$ denotes the language $\{a, b\}$
- $(a|b)(a|b)$ denotes $\{aa, bb, ab, ba\}$ the language

of all strings of length two over the alphabet Σ

Another RG for the same language $\rightarrow aa|bb|ab|ba$

• a^* denotes the language consisting of all strings of ≥ 0 a 's, i.e. $\{\epsilon, a, aa, aaa, \dots\}$

• $(a|b)^*$ denotes the set of all strings consisting of ≥ 0 instances of a or b $\{\epsilon, a, b, aa, bb, aab, aaab, \dots\}$

• $a|a^*b$ denotes the language $\dots \{a, ab, aaab, aab, b\}$
 ' b ' is also a valid string \rightarrow zero or more a 's followed by a single ' b '

SOME ALGEBRAIC RULES THAT HOLD FOR ARBITRARY RE

LAW

$$r_1 s = s r_1$$

DESCRIPTION

r_1 is commutative

$$r_1(s t) = (r_1 s)t \quad |' \text{ is associative}$$

$$r_1(s t) = (r_1 s)t$$

concatenation is associative

$$r_1(s t) = r_1 s |^r_1 t$$

concatenation distributes over

$$(s t)r = s r | t r$$

ϵ is the identity for concatenation

$$r_1^* = (r_1 \epsilon)^*$$

ϵ is guaranteed in a closure

$$r_1^{**} = r_1^*$$

* is idempotent

$$1) r_1 = \emptyset, L(r_1) = \{\emptyset\}, \emptyset \text{ is the empty set}$$

$$2) r_1 = \epsilon, L(r_1) = \{\epsilon\}$$

$$3) r_1 = a, L(r_1) = \{a\}$$

$$4) r_1 = a+b, L(r_1) = \{a, b\}$$

$$5) r_1 = ab, L(r_1) = \{ab\}$$

$$6) r_1 = a+b+c, L(r_1) = \{a, b, c\}$$

$$7) r_1 = (ab+a) \cdot b, L(r_1) = \{abb, ab\}$$

$$8) r_1 = a^+, L(r_1) = \{a, aa, aaa, aaaa, \dots\}$$

$$a) r_1 = a^k, L(r_1) = \{a \text{ k times } y\}$$

- 10) $\pi = (a+ba)(b+a)$, $L(\pi) = \{ab, aa, bab, baa\}$

11) $\pi = (a+\epsilon)(b+\phi)$, $L(\pi) = \{ab, b\}$

12) $\pi = (a+b)(a+b)$, $L(\pi) = \{aa, ab, ba, bb\}$

13) $\pi = (a+b)^*$, $L(\pi) = \{\epsilon, a, b, aa, bb, ab, ba, \dots\}$

14) $\pi = (a+b)^*(a+b)$, $L(\pi) = \{a, b, aa, bb, ab, ba, \dots\}$

15) $\pi = a^* \cdot a^*$, $L(\pi) = \{\epsilon, a, aa, aaa, \dots\}$

16) $\pi = (ab)^*$, $L(\pi) = \{ab, \epsilon, abab, abab, \dots\}$

? : ZERO OR ONE INSTANCE eg: $L(\{1\}) \cup \{\epsilon\}$

CHARACTER CLASSES : A RE $a_1 | a_2 | \dots | a_m \rightarrow [a_1, a_2, \dots, a_m]$

$[abc]$ is shorthand for $a|b|c$

$$[a_3] = \cdots \quad \dots ab|c\cdots|_3$$

Eg : C identifiers are strings of letters, digits & underscores.

Regular definition ?

letter → A|B|...|z|a|b|...|z| → letter → [A-zA-Z-]
 digit → 0|1|...|9 → digit → [0-9]
 id → letter (letter|digit)* → id → letter (letter|digit)*

Eg: Unsigned numbers (integers or floating pt) : are strings

such as 5280, 0.01234, 6.336 E4, or 1.89 E-4

Regular definition ?

digit \rightarrow 0111...19

digits → digit⁺

Optional Fraction → . digits | E

OptionalExponent → (E (+1-1e) digits | e)

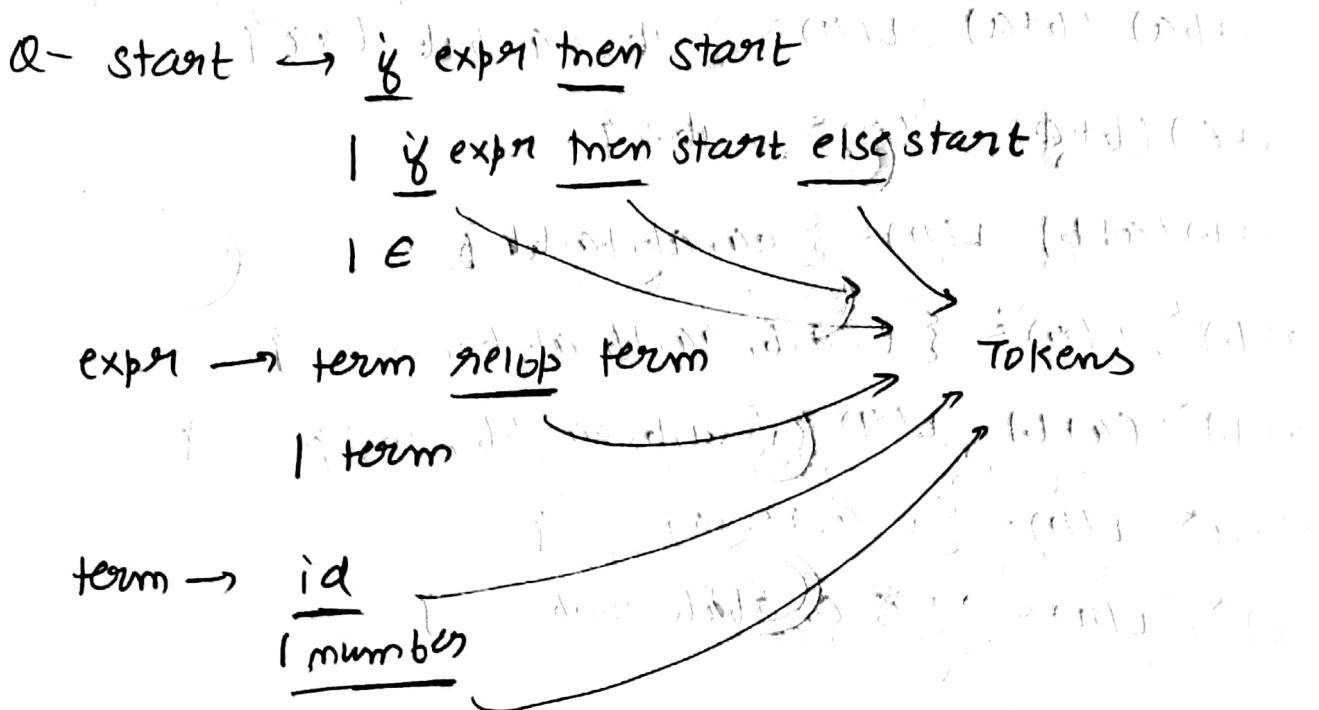
number → digits Optional Fraction optional Exponent

digit → [0-9]

digits → digit⁺

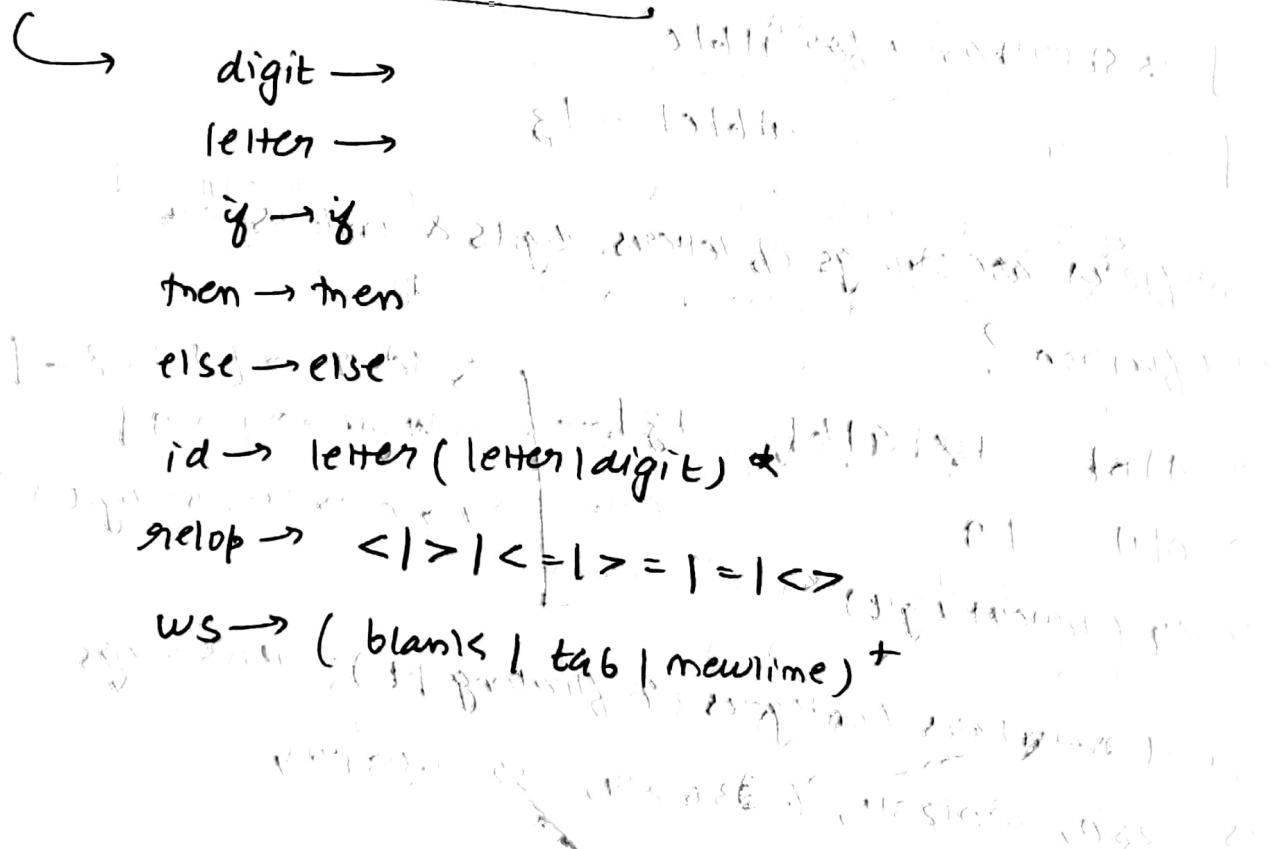
number → digits (· digits)?

(\in [+ -] ? digits)?



Eg: Annex A grammar for branching statement

Patterns for branching statement



Blank → space
 Tab → horizontal tab
 Newline → line feed

GOAL FOR LEXICAL ANALYZER

(3) Lexical Analyzer

can do it in 1986

<u>LEXEMES</u>	<u>TOKEN NAME</u>	<u>ATTRIBUTE VALUE</u>
Any ws	blank	-
if	if	-
then	then	-
else	else	-
Any id	id	pointer to table entry
Any number	number	" "
<	relOp	LT
<=	relOp	LE
=	relOp	EQ
>	relOp	NE
>=	relOp	GT
>=	relOp	GE

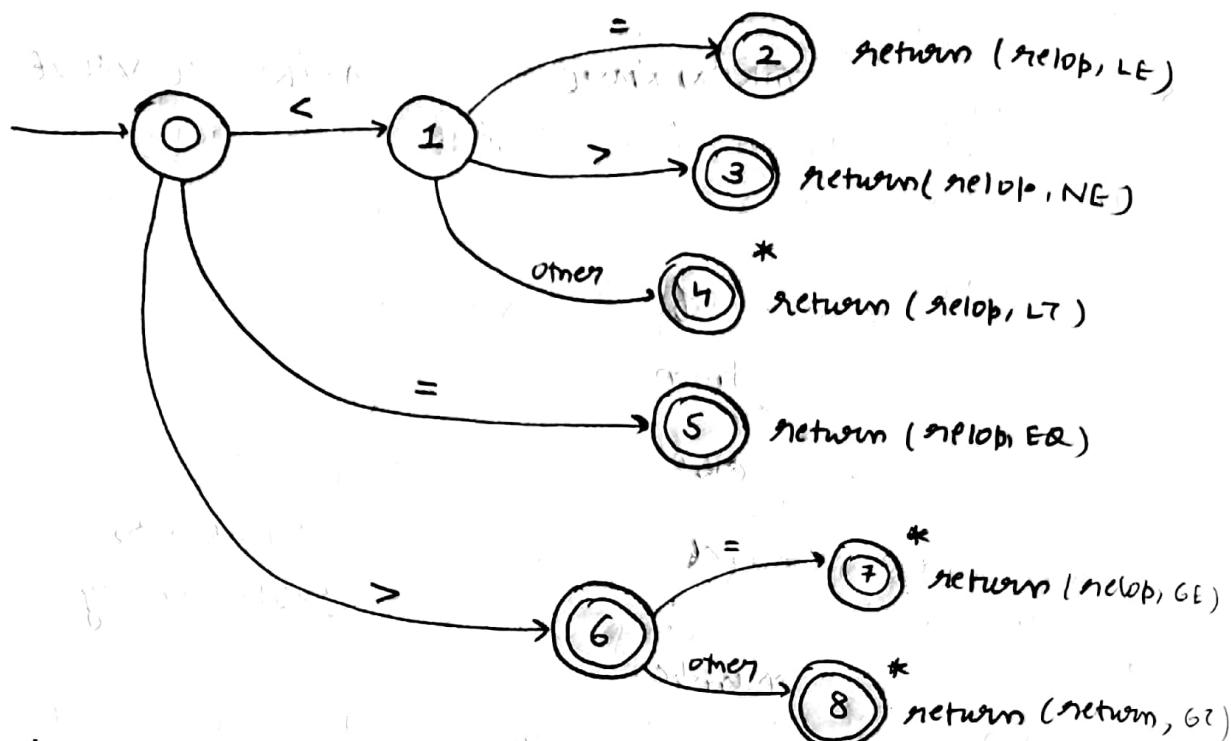
Transition Diagram: diff. states

like a flowchart

States & edges: actions to be taken given a particular IIP



TRANSITION DIAGRAM EG



* → retract back one step (like forward ptr & lexeme begins in 2/1 buffering
take 1 step back, till 1 back step, we'll have a token)

→ returns object of type TOKEN

TOKEN getRelop()

```

{ TOKEN getToken = new( RELOP )           ↴ symbolic code for RELOP
  while(1)
    /* repeat character processing until a failure or return occurs */
    switch(state)
      { case 0: C = nextCharacter();
        if (C == '<') state = 1;
        else if (C == '=') state = 5;
        else if (C == '>') state = 6;
        else fail(); /* lexeme not a relop */
        break;
      case 1:   /* implementation of state 1 */
      case 5:   /* implementation of state 5 */
      case 6:   /* implementation of state 6 */
      case 7:   /* implementation of state 7 */
      case 8:   /* implementation of state 8 */
        }
      }
    
```

Fig Sketch of
implementation
of relop
transition diagram

→ case 1: `C = nextChar();`

if (`C == '='`) state = 2;
else if (`C == '>'`) state = 3;
else state = 4;
break;

case 2: `netToken.setAttribute = LE;`
`return (netToken);`

case 3: `netToken.setAttribute = NE;`
`return (netToken);`

case 4: `netact();`
`netToken.setAttribute = LT;`
`return (netToken);`

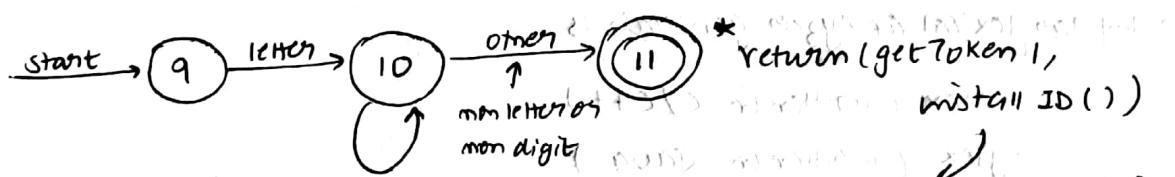
case 5:

case 6:

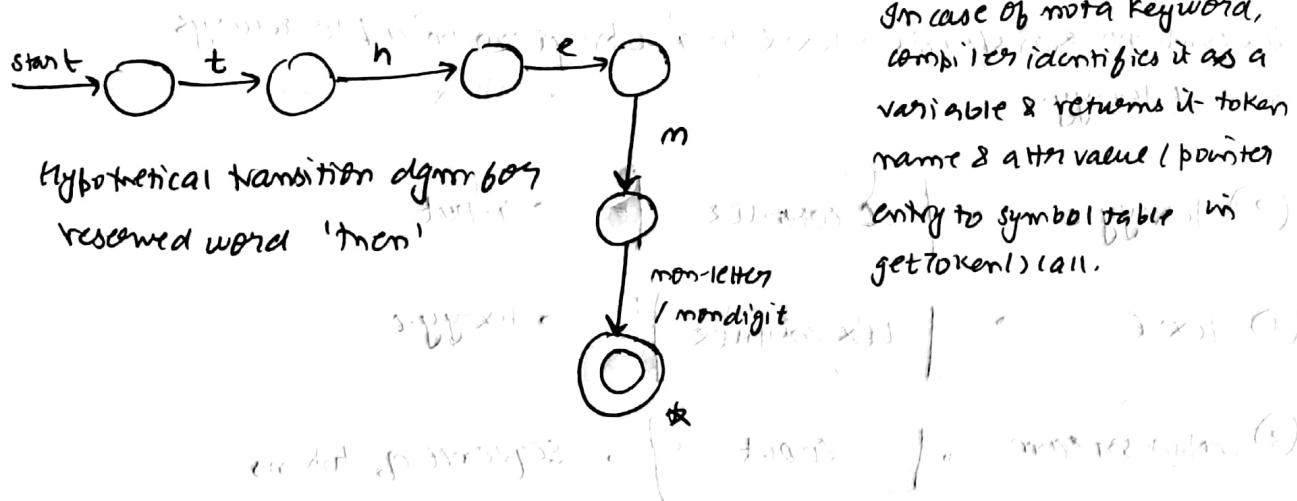
case 7:

RECOGNITION OF RESERVED WORDS & IDENTIFIERS

↳ 1) SYMBOL TABLE → 32 keywords in C language



2) TRANSITION DIAGRAM (hypothetical eg)



CONFLICT RESOLUTION IN LEX

- 1) Always prefer a longer prefix to a shorter prefix
- 2) If longest prefix matches 2 or more patterns, prefix the pattern listed first in lex program

e.g. for ① → if 3 abc

] IIP lexemes

$$if \rightarrow if$$

identifier → letter (letter | digit)*

→ here if 3 can match keyword if but acc. to longer prefix strategy, it will resolve in favour of identifier i.e. longer prefix

for ② y] IIP lexeme

→ matches both if & identifier but will resolve in favour of pattern listed first hence matches the keyword & conflict resolves

LEXICAL ANALYZER GENERATORS

Process of generating a lexical analyzer can be automated

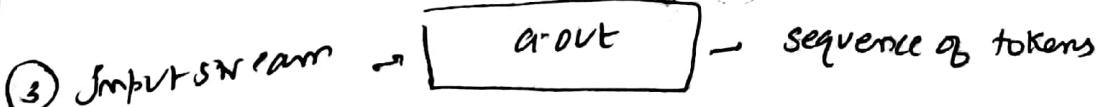
2 popular lexical analyzer generators:

-Flex (written in C/C++)

-Jflex (written in Java)

lex.yylex → a file that simulates transition diagram

lex compiler transforms lex.yylex to a C program in a file always named lex.yylex



STRUCTURE OF LEX PROGRAMS

STRUCTURE OF LEX PROGRAMS

contains
regular
expressions → declarations (contains pattern defn)

%%

translation rules (contains (pattern, {actions}, etc))

%%

Auxiliary functions (rest of necessary definitions)

yyval & other variables

yytext

- The program contains a % - delimiter to mark beginning of rules & one rule

PROGRAM TO DELETE FROM I/P ALL BLANKS OR TABS AT ENDS OF LINES

%.

[\t] + \$

lex programs recognise only regex

You writes ~~parser~~ that accept a large set of CFG.

Absolute main - lex program → % . (nodeln, n rules)

If operator characters are to be used as text characters, we must use escape '\ ' characters before them.

lex default action → copying I/P to O/P.

A pattern which is not matched for a string is copied as it is to O/P.

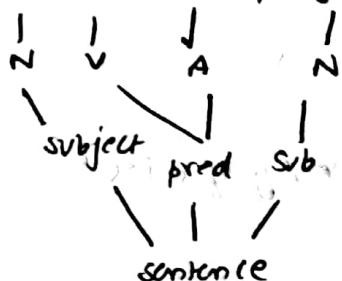
lex - turns rules to a program

SYNTAX ANALYSIS

PARSING

source code \rightarrow scanner \rightarrow ...
 \downarrow ...

He wrote the program



ROLE OF PARSER

- 1) Must distinguish b/w valid & invalid tokens
- 2) ~~Type checking~~

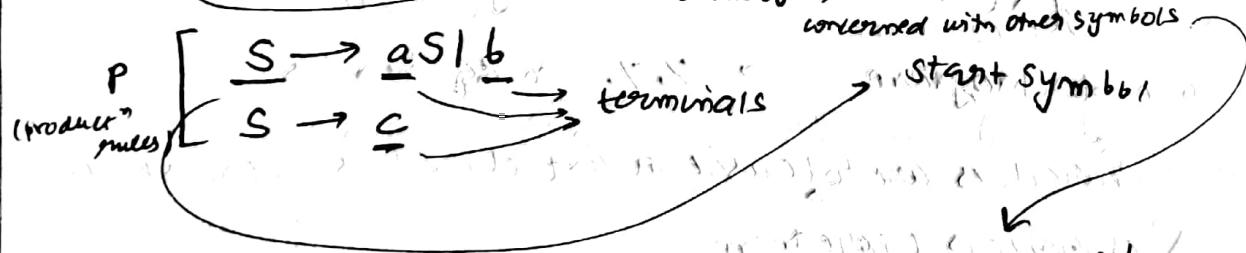
Type checking & initializatⁿ of variables is handled by semantic analysis

Regular languages can't handle all cases

e.g.: strings of balanced parentheses $(((-))_-)$

\therefore we require push down automata / Table driven parsers.

A grammar G contains (variables (V), terminals (T), Prodⁿ Rules (P), Context Free Gram. CFG \rightarrow why? only a single symbol is involved Start (S)).



Derivation of a Variable

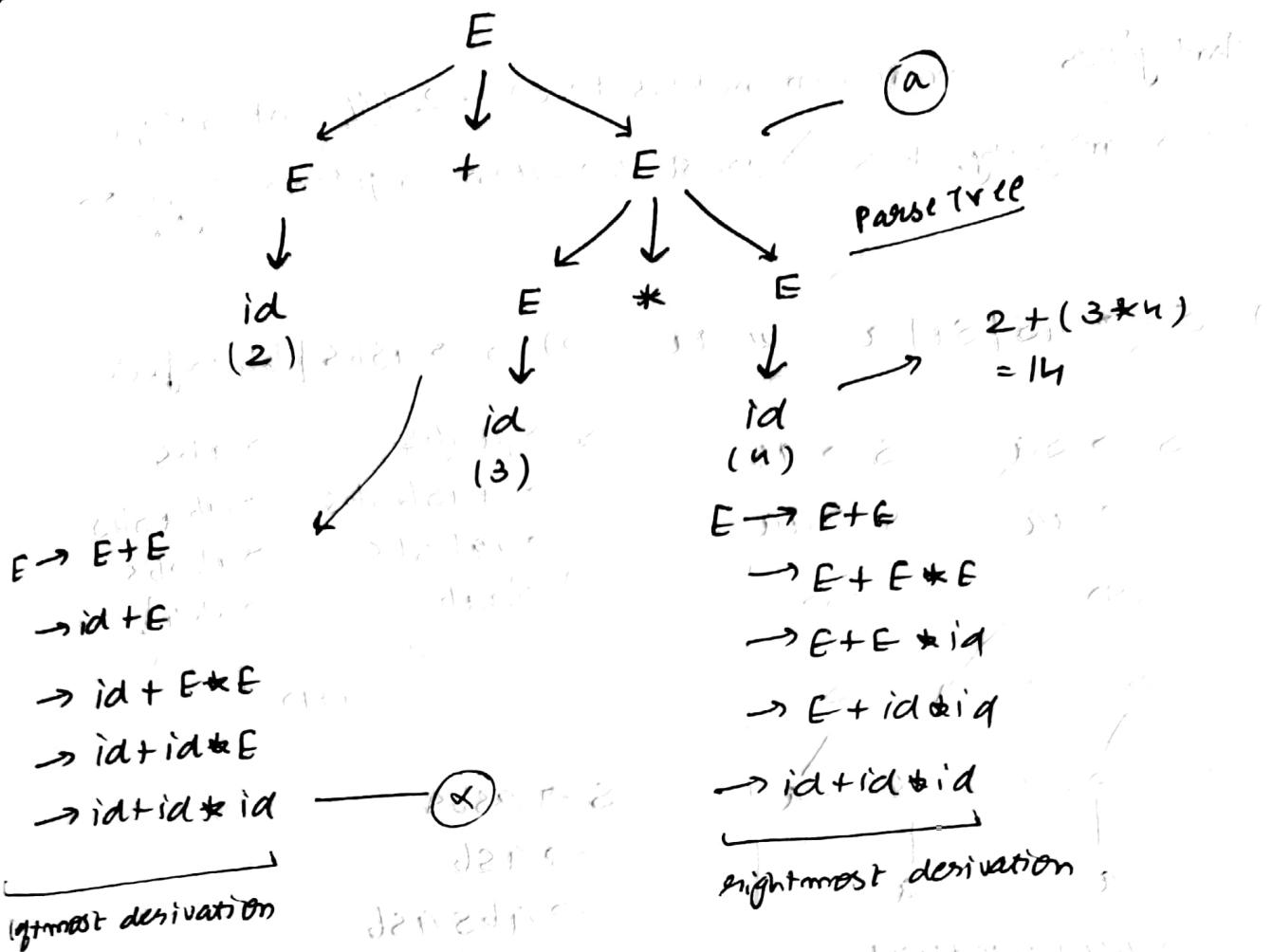
- \rightarrow leftmost derivation $\xrightarrow{1}$
 \rightarrow rightmost derivation $\xrightarrow{2}$

① we start expanding from leftmost

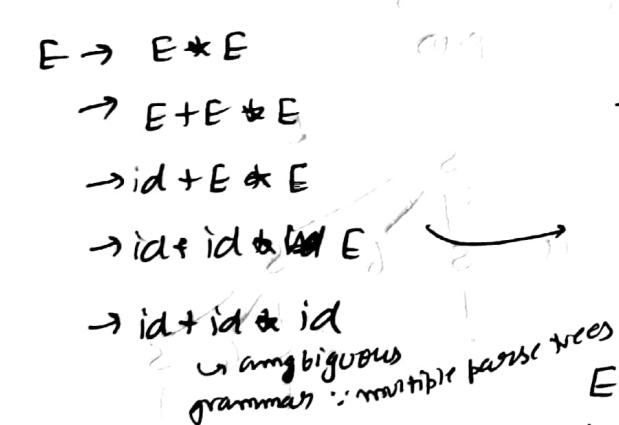
e.g. G : E \rightarrow E+E / E*E / id

string : id + id * id

check whether this string $\in G$ or not?



Another method for deriving (a) in leftmost derivation form



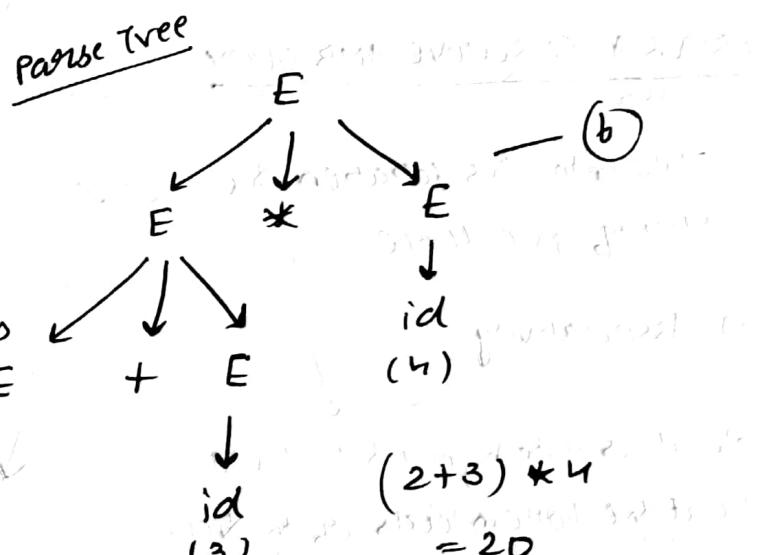
How does parser know

which tree to generate

from (a) & (b)

→ Ambiguity to be removed

while programming.

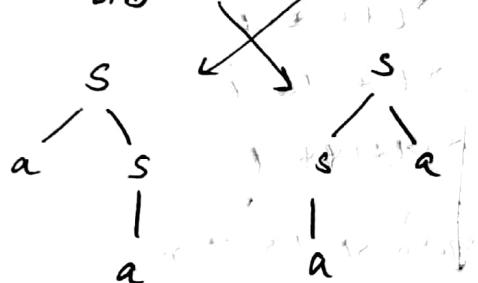


Ambiguous — more than one parse tree \rightarrow 2 different interpretations

To check ambiguity, check if leftmost derivation = rightmost derivation

$$1) S \rightarrow as \mid sa \mid a \quad w=aa$$

$$\begin{array}{l} S \rightarrow sa \\ \quad \quad \quad \rightarrow aa \\ \hline LMD \end{array} \quad \begin{array}{l} S \rightarrow as \\ \quad \quad \quad \rightarrow aa \\ \hline RMD \end{array}$$



2 diff. parse trees

\therefore ambiguous

$$2) S \rightarrow asbs \mid bsas \mid \dots \in w=abab$$

$$\begin{array}{l} S \rightarrow asbs \\ \quad \quad \quad \rightarrow abs \\ \quad \quad \quad \rightarrow aabsbs \\ \quad \quad \quad \rightarrow aaabsbs \\ \quad \quad \quad \rightarrow aab \\ \hline LMD \end{array}$$

$$\begin{array}{l} S \rightarrow bsas \\ \quad \quad \quad \rightarrow ab \\ \quad \quad \quad \rightarrow absasb \\ \quad \quad \quad \rightarrow absab \\ \quad \quad \quad \rightarrow abab \\ \hline RMD \end{array}$$

CRITERIA TO REMOVE AMBIGUITY

OR RULES

1) Take into consideration, the order of precedence

2) Associativity

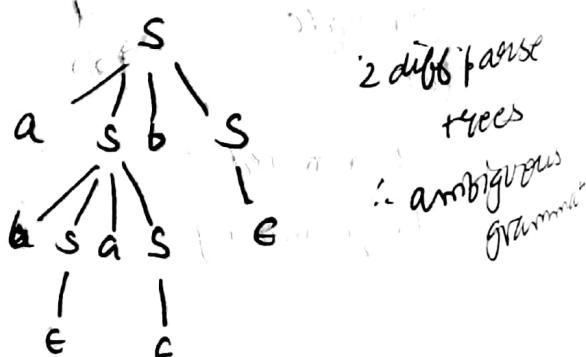
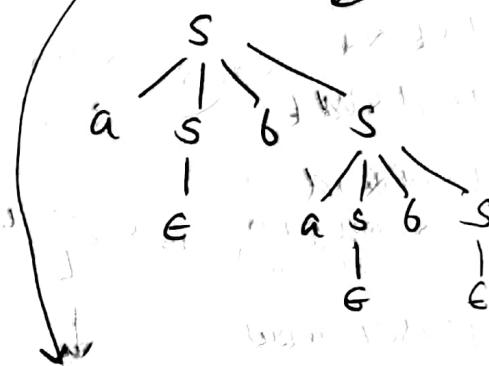
operators with higher precedence

lie at the bottom levels of the tree.

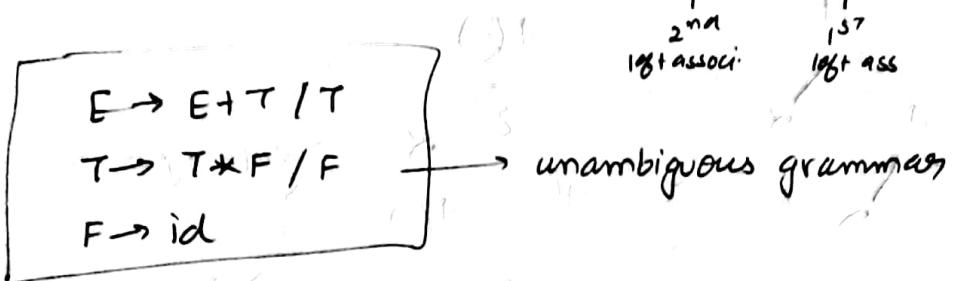
when operator is:

left associative \Rightarrow left recursive

right associative \Rightarrow right recursive

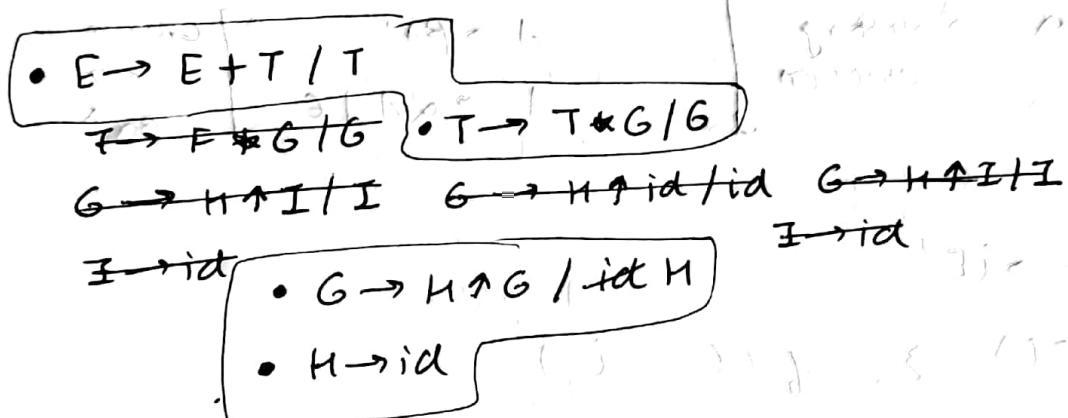


Q- Remove ambiguity from $E \rightarrow E+E \mid E*E \mid id$



Q- + * \uparrow (power operator)
 3RD 2ND IST precedence
 left associativity right associativity

Generate a grammar which follows the above rules



PARSING TECHNIQUES / ALGOS

e.g. of Top Down Approach → Recursive Descent Parsing

can't handle left recursion

Elimination of Left Recursion Methods

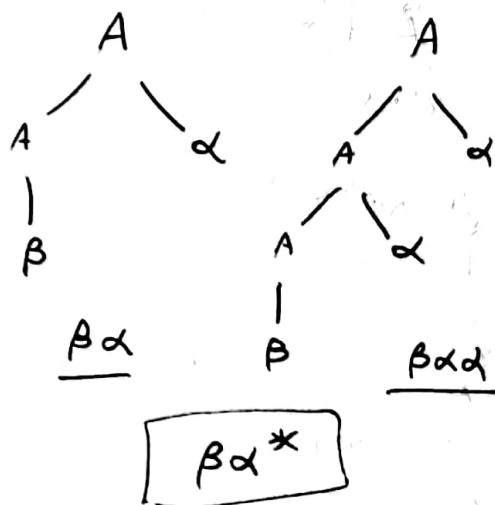
$$A \rightarrow A \propto | \beta |$$

 left recursive

α & β ane terminals

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

↑ right
regressive



A()

Σ

$\alpha \rightarrow$ stopping criterion

A()

g

$\alpha^* \beta$

NO LOOP

A()

{

no stopping
criterion
here

A()

3

∞ LOOP

(α) → stopping
criterion

AIM : to remove left recursion
while keeping grammar same

$A \rightarrow \beta A'$

$A' \rightarrow \alpha A' | E$

Grammar

609

$\beta \alpha^*$

$E \rightarrow i E'$

$E() \Sigma$ if ($i == 'i'$), then i ,

$\Sigma \text{ match}('i') ;$

$E'() ;$

$\text{expr} \rightarrow \text{expr} + \text{terminal} \mid \text{terminal}$

$A \quad \quad A \quad \quad \alpha' \quad \quad \beta$

left recursion
eliminated

$\text{expr} \rightarrow \text{terminal expr}'$

$\text{expr}' \rightarrow + \text{terminal expr}' \mid \epsilon$

$$① S \rightarrow SOSIS \mid 01$$

$$\begin{matrix} S \\ \alpha \\ \beta \end{matrix}$$

first is A root is 2

$$\begin{matrix} S \rightarrow \beta S' \\ S' \rightarrow \alpha S' \mid \epsilon \end{matrix}$$

$$② L \xrightarrow{\alpha} L, \underbrace{B}_{A \propto} \xrightarrow{\beta} B$$

$$\begin{matrix} L \rightarrow \alpha L' \\ L' \rightarrow \beta L' \mid \epsilon \end{matrix}$$

GRAMMAR

Unambiguous

Ambiguous

left recursive

Right recursive

Deterministic

Non-deterministic

left factoring: converting

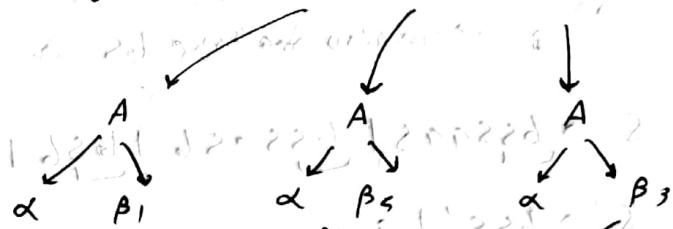
non-deterministic to
deterministic

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$$

Deterministic

from $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3$



backtracking problem

∴ we should try to keep grammar
deterministic

- $S \rightarrow \underline{assbs} \mid \underline{asasb} \mid \underline{abb} \mid b$ (non-deterministic)
- thus non-deterministic

to make it deterministic:

$$S \rightarrow as'1b$$

$$S' \rightarrow \underline{ssbs} \mid \underline{sasb} \mid bb$$

$$S' \rightarrow ss''1bb$$

$$S'' \rightarrow sbs \mid asb$$

- $S \rightarrow \underline{bssaas} \mid \underline{bssasb} \mid \underline{bsb} \mid a$
- non-deterministic

to make it deterministic

$$S \rightarrow bs'1a$$

$$S' \rightarrow \underline{ssaas} \mid \underline{ssasb} \mid \underline{sb}$$

$$S \rightarrow ss''$$

$$S'' \rightarrow \underline{sas} \mid \underline{sasb} \mid b$$

$$S'' \rightarrow ss'''1b$$

$$S''' \rightarrow \underline{aas} \mid \underline{asb}$$

$$S''' \rightarrow as'''$$

$$S''' \rightarrow as \mid sb$$

we can also take bs common

$$S \rightarrow \underline{bssaas} \mid \underline{bssasb} \mid \underline{bsb} \mid a$$

$$S \rightarrow bss'1a$$

$$S' \rightarrow \underline{sas} \mid \underline{sasb} \mid b$$

$$S' \rightarrow sas''1b$$

$$S'' \rightarrow as \mid sb$$

TOP-DOWN APPROACH

BOTTOM-UP APPROACH



Eg: $S \rightarrow aABC$

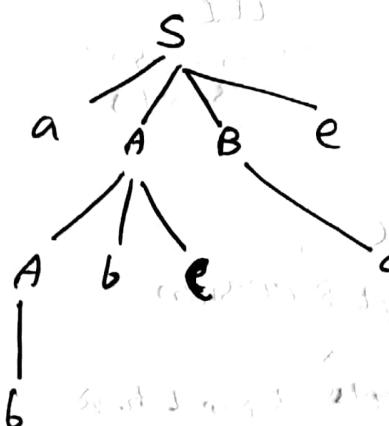
$A \rightarrow A\epsilon | b$

$B \rightarrow d$

given $w \rightarrow abbede$

we can generate parse tree using these approaches

TOP-DOWN PARSE TREE :



Q - what production rule

BOTTOM-UP PARSE TREE



Q - When to reduce?

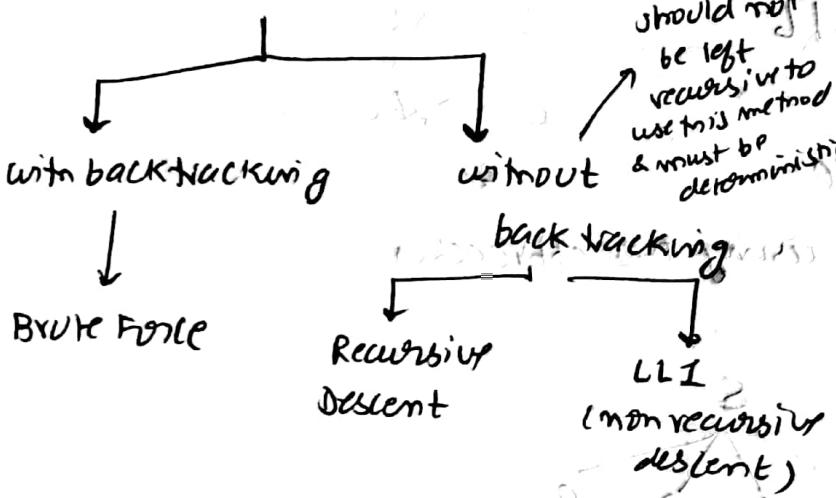
↓
Word Tree

↓
Parse Tree

← START here
go upwards

PARSERS

Top down approach (parsers)



Bottom up approach (parsers)

operator precedence parser

LR parser

$100\% head = 1$

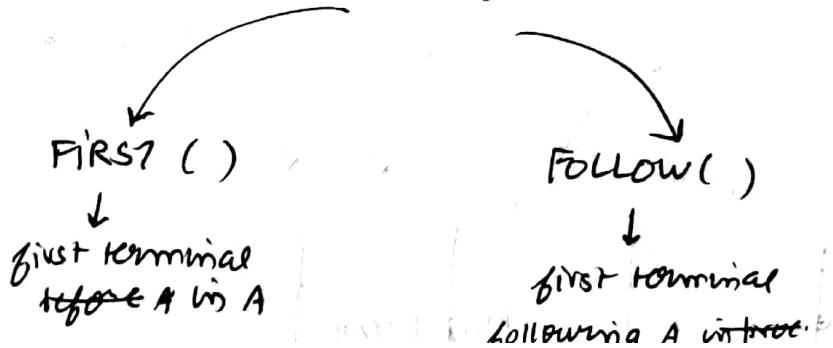
LL(1) : Deterministic
Eliminate left Recursion
LMD
Scanning of symbols from L to R

\$ always present in stack & buffer to keep check on when to stop

structures used (components)

- stack
- I/P buffer
- LL(1) parsing algo
- LL(1) parsing table

Two parsing methodologies under LL(1)



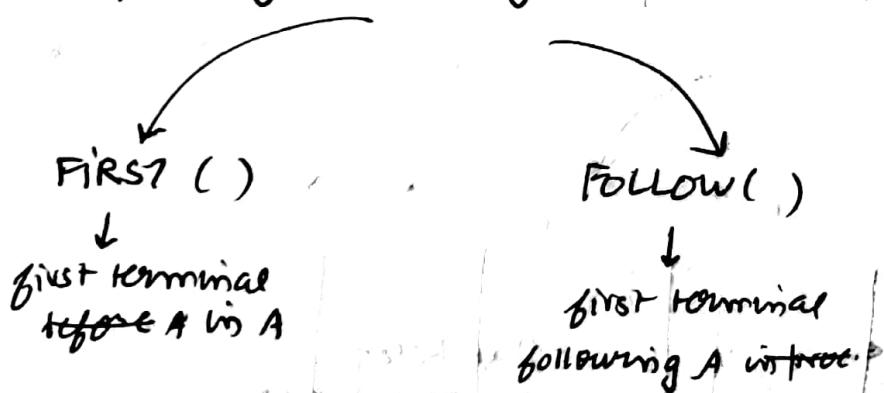
data structures which helps in implementing parsing algorithm

LL(1) : Deterministic
 Eliminate left Recursion
 LMD
 Scanning of symbols from L to R

structures used (components)
 - stack 
 - G/P buffer 
 - LL(1) parsing algo
 - LL(1) parsing table

\$' always present in stack & buffer to keep check on when do we arrive at the end of the algo. / when to stop

Two parsing methodologies under LL(1)



$$S \rightarrow aABC D$$

$$A \rightarrow b$$

$$B \rightarrow c$$

$$C \rightarrow d$$

$$D \rightarrow e$$

$$\text{FIRST}(S) \rightarrow \{a\}$$

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

$$\text{FIRST}(B) \rightarrow \{c\}$$

$$\text{FIRST}(C) \rightarrow \{d\}$$

$$\text{FIRST}(D) \rightarrow \{e\}$$

$$S \rightarrow ABCD$$

$$A \rightarrow b | G$$

$$B \rightarrow c$$

$$C \rightarrow d$$

$$D \rightarrow e$$

$$\text{FIRST}(S) \rightarrow \{b, c, d\}$$

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

$$\text{FIRST}(B) \rightarrow \{c\}$$

$$\text{FIRST}(C) \rightarrow \{d\}$$

$$\text{FIRST}(D) \rightarrow \{e\}$$

data structure
 which helps in implementing
 parsing algorithm

$S \rightarrow ABCD$ $\text{follow}(S) \rightarrow \$$ $A \rightarrow b | \epsilon$ $\text{follow}(A) \rightarrow \{\epsilon, c\}$ $B \rightarrow c$ $\text{follow}(B) \rightarrow \{\epsilon, d\}$ $C \rightarrow d$ $\text{follow}(C) \rightarrow \{\epsilon, e\}$ $D \rightarrow \epsilon | e$ $\text{follow}(D) \rightarrow \text{follow}(S) \rightarrow \{\epsilon, \$\}$

$\text{follow}(\cdot) \rightarrow$ what is the terminal which can follow a variable in process of derivation

$\epsilon \rightarrow$ should not be present, $\text{follow}(S) \rightarrow \$$

EXAMPLE $S \rightarrow ABCDE$ $A \rightarrow a | \epsilon$ $B \rightarrow b | \epsilon$ $C \rightarrow c$ $D \rightarrow d | \epsilon$ $E \rightarrow e | \epsilon$

<u>FIRST</u>	<u>FOLLOW</u>
$\{\epsilon, a, b, c\}$	$\{\$, \epsilon\}$
$\{\epsilon, a, e\}$	$\{b, c\}$
$\{b, e\}$	$\{c\}$
$\{c\}$	$\{d, e, \$\}$
$\{d, e\}$	$\{e, \$\}$
$\{e, \epsilon\}$	$\{\$\}$

Q- $E \rightarrow E + T \mid T$ $T \rightarrow T * F \mid F$ $F \rightarrow (E) \mid \text{id}$

$$\boxed{A = \beta A' \\ A' = \alpha A' | E}$$

 $E \rightarrow E + T \mid T \quad \begin{matrix} \beta \\ \alpha \end{matrix}$ $T \rightarrow T * F \mid F \quad \begin{matrix} \beta \\ \alpha \end{matrix}$ $E \rightarrow T E'$ $T \rightarrow F T' \quad \begin{matrix} \beta \\ \alpha \end{matrix}$ $E' \rightarrow + T E' \mid \epsilon$ $T' \rightarrow * F T' \mid \epsilon$ $E \rightarrow T E'$ $E : \{\epsilon, \text{id}, (\}\}$ $\{\$, \epsilon\}$ $E' \rightarrow + T E' \mid \epsilon$ $E' : \{+, \epsilon\}$ $\{\text{id}, (\}\}$ $T \rightarrow F T' \mid \epsilon$ $\{\epsilon, (\}\}$ $\{\ast, \epsilon\}$ $T' \rightarrow * F T' \mid \epsilon$ $\{\epsilon, \epsilon\}$ $\{\$, \epsilon\}$ $F \rightarrow (E) \mid \text{id}$ $\{\text{id}, (\}\}$ $\{\ast, \epsilon\}$ $\text{follow}(T) = \text{first}(E')$ $\text{follow}(T') = \text{follow}(T)$ $\text{follow}(E') = \text{follow}(E)$ $\text{follow}(F) = \text{first}(T'), \text{follow}(T')$

LL(1) PARSING

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

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

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



$$E \rightarrow TE'$$

$$C \rightarrow +TE' \mid E$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid E$$

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

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

$$A \rightarrow \beta A'$$

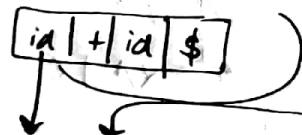
$$A' \rightarrow \alpha A' \mid E$$

PREDICTIVE PARSING

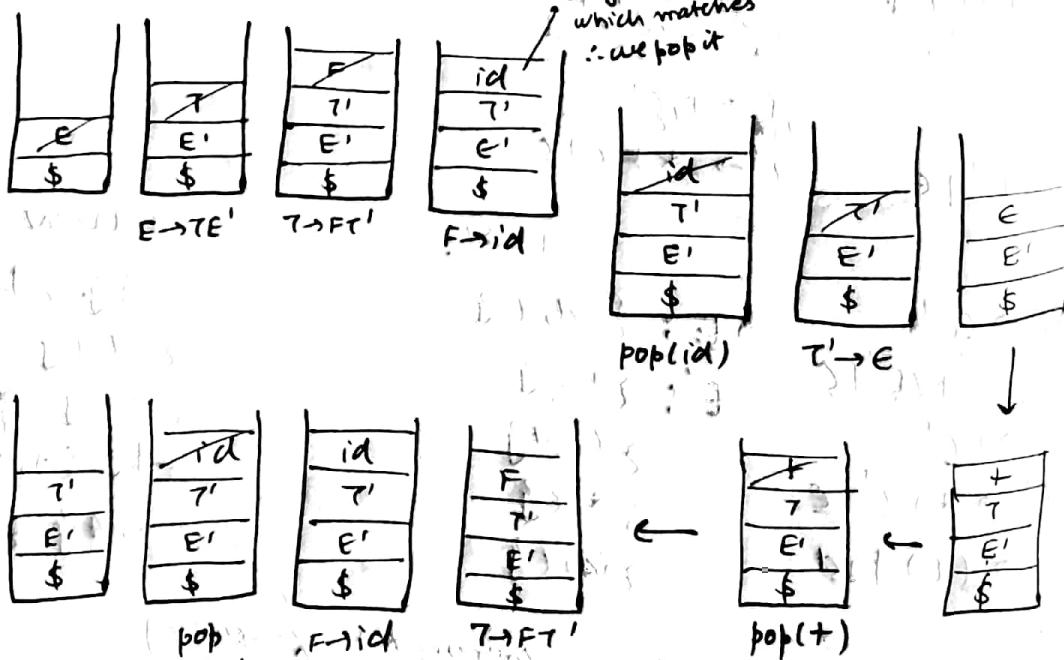
GENERATE PARSING TABLE

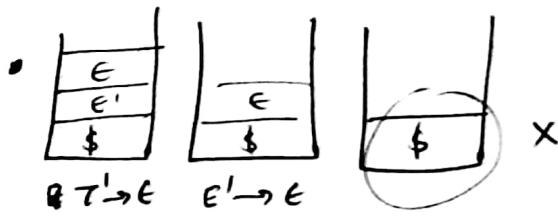
	id	+	*	()	\$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$			$E' \rightarrow E$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow FT'$		$T' \rightarrow E$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$				$F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

consider



we get id here
which matches
∴ we pop it





STACK	INPUT	ACTION
$E \$$	$id + id \$$	$E \rightarrow T E'$
$TE' \$$		$T \rightarrow F T'$
$FT'e' \$$		$F \rightarrow id$
$idT'E' \$$	$id + id \$$	$pop(id)$
$T'E' \$$	$id + id \$$	$T' \rightarrow E$
$E' \$$		$E' \rightarrow + T E'$
$+ T E' \$$		$pop(+)$
$TE' \$$	$id + id \$$	$T \rightarrow F T'$
$FT'E' \$$	$id + id \$$	$F \rightarrow id$
$idT'E' \$$		$pop(id)$
$T'E' \$$	$id + id \$$	$T' \rightarrow E$
$E' \$$	$id + id \$$	$E' \rightarrow E$
$\$$	$\$$	success

RECURSIVE DESCENT PARSING

$E \rightarrow iE'$
 $E' \rightarrow +iE' \mid \epsilon$

function $E()$
 {
 if ($\lambda = 'i'$)
 { match('i'); E'(); }
 else
 {
 G
 +
 i
 }
 }

$E'()$

{ if ($\lambda == '+'$)
 { match('+');
 match('i');
 E'(); }
 }
}
else

return // for E

match(char t)

{ if ($\lambda == t$)
 $\lambda = \text{getchar}();$
else
 printf("error");
}
}

main()

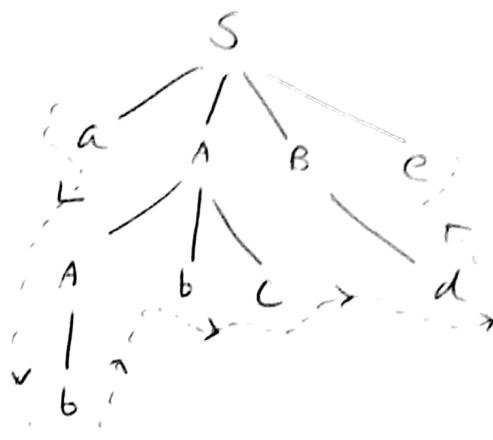
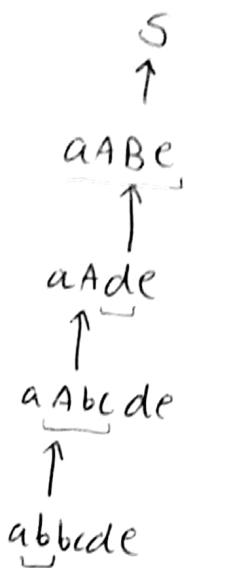
{
 E();
 if ($\lambda == \$$)
 printf("Parsing success");
 }
 $\lambda = \text{getchar}();$ // assuming to be global

[t] + \$ { }

BOTTOM-UP PARSING (we derive the start symbol)

Eg: $S \rightarrow aABe$
 $A \rightarrow Abc | b$
 $B \rightarrow d$

Given string, w - abbcde



RIGHTMOST DERIVATⁿ IN REVERSE

Handle: the substring that is being substituted

Handle Pruning: A handle is substring that matches body of product & whose reduction represents one step along the reverse of rightmost derivation

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

$T \rightarrow T * F \mid F$ if $w = id_1 * id_2$

$F \rightarrow (E) \mid id$

Table →

RIGHT SENTENTIAL FORMHANDLEREDUCING PRO.id₁ * id₂id₁

F → id

F * id₂

F

T → F

T * id₂id₂

F → id

T * F

T * F

T → T * F

T

T

E → T

E —

— Parsing successful as start symbol
achieved

SHIFT REDUCE PARSING: It's a form of bottom up parsing where a stack holds grammar symbols & an input buffer holds the rest of string to be parsed. The handle always appears at the top of stack just before its identity as the handle.

Initially

STACK

\$

INPUT

w \$

Final state

\$ S

\$

CONFIGURATION OF A SHIFT-REDUCE PARSER ON IIP id₁ * id₂

E → E + T | T

T → T * F | F

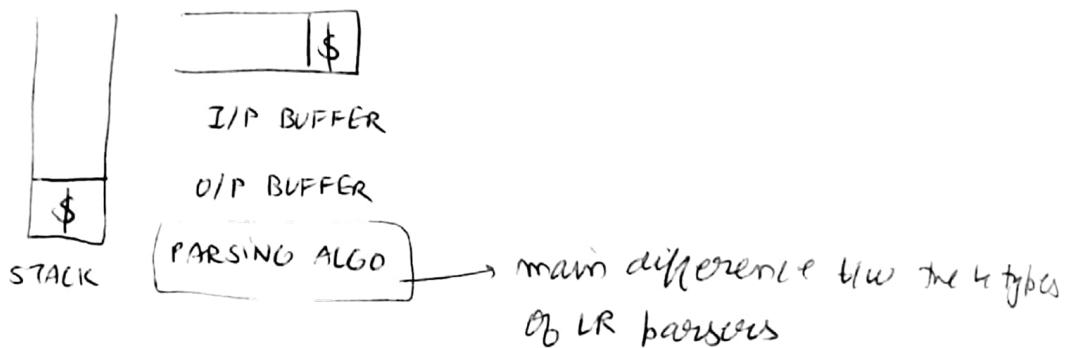
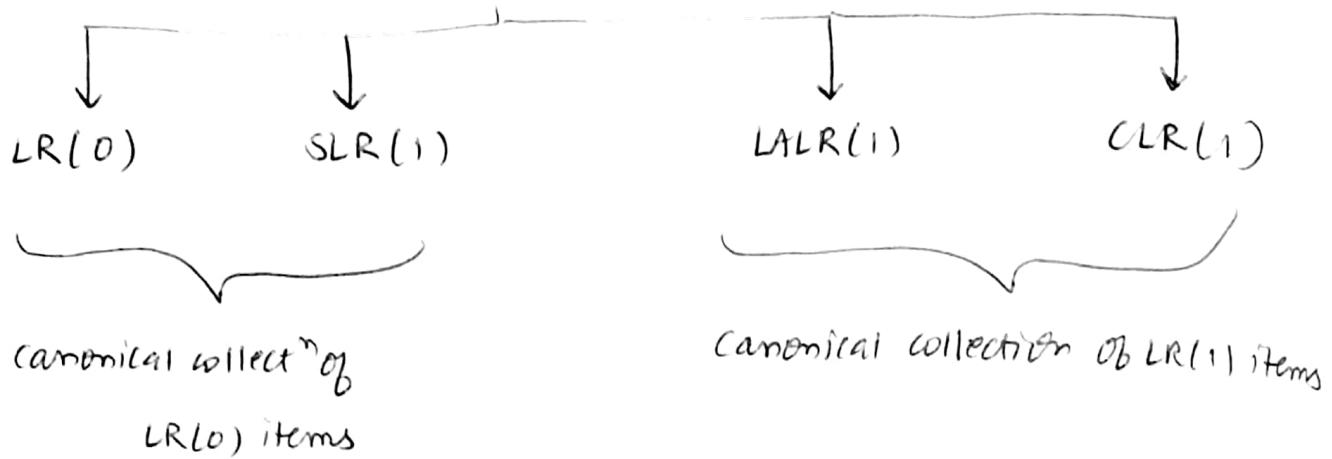
F → (E) | id

w = id₁ * id₂

<u>STACK</u>	<u>INPUT</u>	<u>ACTION</u>
\$	id ₁ * id ₂ \$	shift
\$ id ₁	* id ₂ \$	reduce by F → id
\$ F	* id ₂ \$	reduce by T → F
\$ T	* id ₂ \$	shift
\$ T *	id ₂ \$	shift
\$ T * id ₂	\$	reduce by F → id
\$ T * F	\$	reduce by T → T * F
\$ T	\$	reduce by E → T
\$ E	\$	accept

4 ACTIONS : shift, reduce, accept or error

LR PARSERS



LR(0) PARSER

Eg: $S \rightarrow AA$
 $A \rightarrow aA \mid b$

We include an additional production rule $S' \rightarrow S$

Augmented Grammar

$\left\{ \begin{array}{l} S' \rightarrow S \\ S \rightarrow AA \\ A \rightarrow aA \mid b \end{array} \right.$

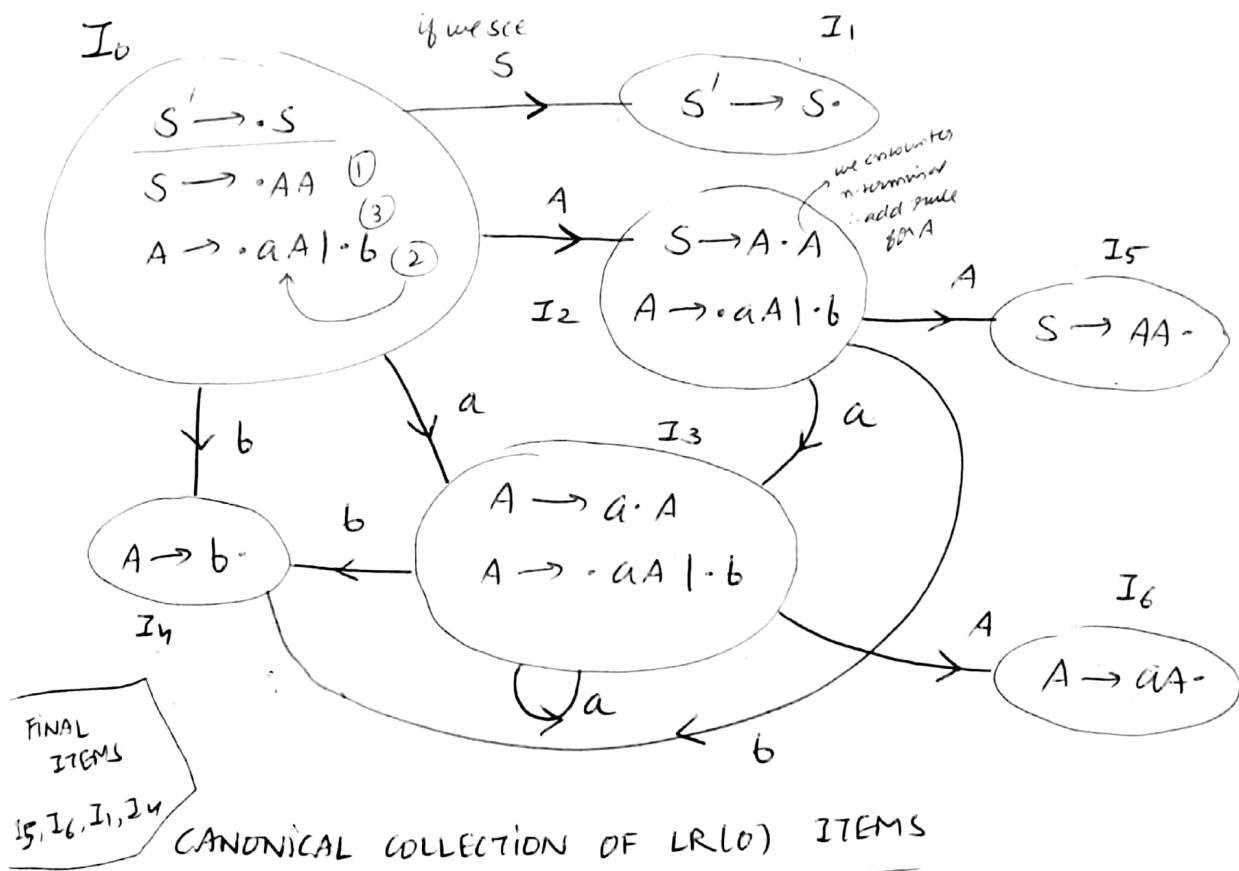
LR(0) items

$S \rightarrow \underline{\cdot} AA$ item
 $S \rightarrow A \cdot A$
 $S \rightarrow AA \cdot$

now ready for reduction

$S' \rightarrow \cdot S$
 $S \rightarrow \cdot AA$
 $A \rightarrow \cdot aA \mid \cdot b$

adding a new production rule for
 encountering a non-terminal
 is called closure
 terminals \downarrow no further rules



	ACTION (SHIFT, REDUCE)			GO-TO	
	a	b	\$	A	S
0	S_3	S_4		2	1
1				accept	
2	S_3	S_4		5	
3	S_3	S_4		6	
4	η_3	η_3	η_3		
5	η_1	η_1			
6	η_2	η_2	η_2		

LR(0) PARSING :

$$S' \rightarrow S$$

$$S \rightarrow \cdot A A \quad ①$$

$$A \rightarrow \cdot a A \quad ② \quad | \quad b \quad ③$$

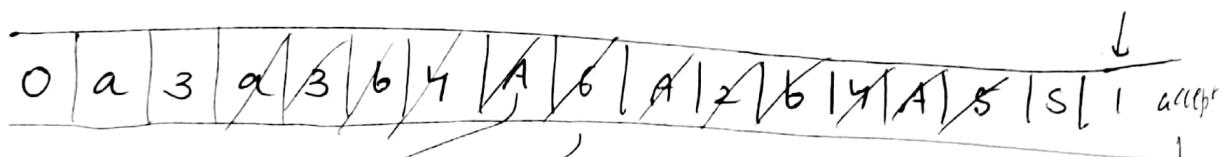
$$w \rightarrow aabb \quad \$$$

$$\text{Follow}(A) \rightarrow \{ q, b, \$ \}$$

	a	b	\$	A	S
0	s_3	s_4		2	1
1			accept		
2	s_3	s_4		5	
3	s_3	s_4		6	
4	q_3	q_3	q_3		
5	q_1	q_1	q_2		
6	q_2	q_2	q_2		

$\begin{array}{c} A \\ / \quad \backslash \\ A \quad A \\ / \quad \backslash \\ A \quad A \\ / \quad \backslash \\ s \end{array}$

for $w \rightarrow aabb \quad \$$



check q, b where it is q_3 i.e. reduce
we get b through prod. no. ③ :- length of
production = 1

↓
↓
parsing success

pointer is at
aab b
↑

∴ we pop out 2 × length of production
i.e. we pop q & b

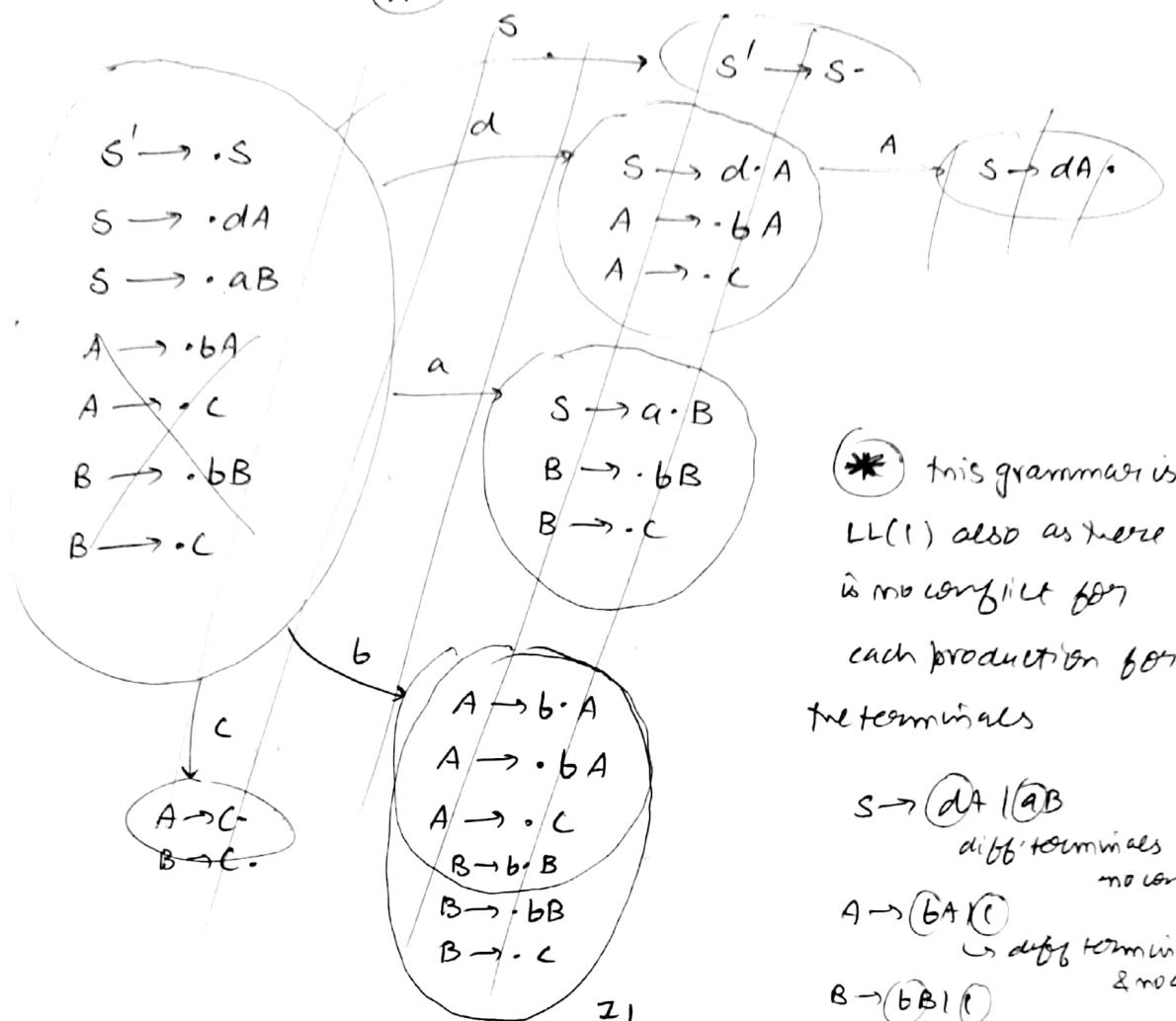
∴ we must
 q_2 ← length

if prod. ② = 2 ∴ we pop out 2 symbols i.e. $q, 3, A, b$

$$\begin{array}{l} \text{Q- } S \rightarrow dA \mid aB \\ \quad A \rightarrow bA \mid c \\ \quad B \rightarrow bB \mid c \end{array}$$

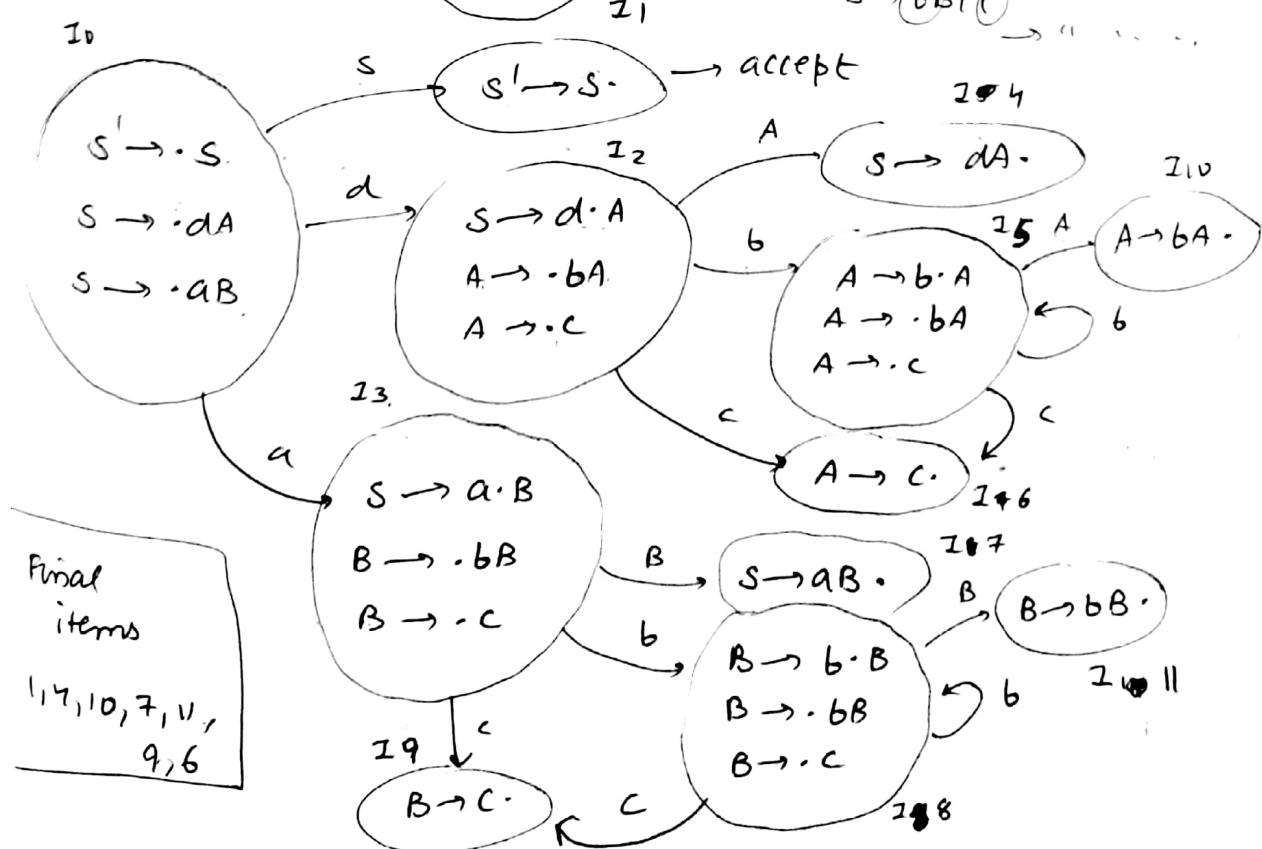
1) LR(0) items

2) table for LR(0) & SLR(1)



* this grammar is LL(1) also as there is no conflict for each production for the terminals

$$\begin{array}{l} S \rightarrow dA \mid aB \\ \text{diff terminals!} \\ \text{no conflict} \\ A \rightarrow bA \mid c \\ \text{diff terminals} \\ \text{& no conflict} \\ B \rightarrow bB \mid c \end{array}$$



	a	b	c	d	\$	ACTION	GO-TO
0	s_3				s_2		A
1						accept	B
2		s_5	s_6				4
3		s_1	s_9				7
4	η_1	η_{11}	η_1	η_1			
5					s_5	s_6	10
6	η_2	η_2	η_2	η_2			
7		s_8	s_9				11
8	η_6	η_6	η_6	η_6			
9	η_3	η_3	η_3	η_3			
10	η_5	η_5	η_5	η_5			
11							

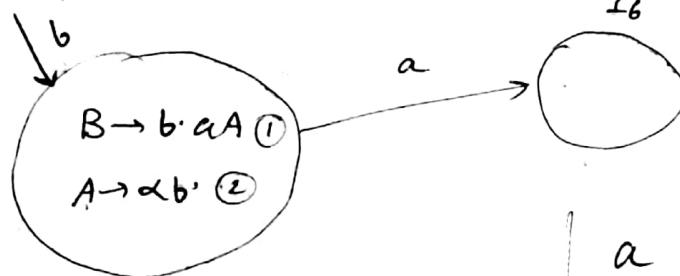
→ this grammar is LR(0) & SLR(1)

in SLR(1), η_1, \dots, η_6 is placed at the follow of LHS of production

Not every grammar is LR(0) ~~and every grammar is SLR(1)~~

POSSIBLE CASES

1) I₅

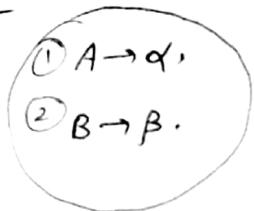


NOT LR(0)

	a	b
5	η_2 / s_6	η_2

shift-reduce conflict

2)

I₅

5	a	b	
	α_1/α_2	β_1/β_2	<u>NOT LR(0)</u>

reduce-reduce conflict

Not every grammar is SLR(1)Possible Cases1) I₅

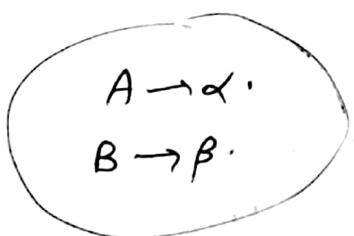
\$	a	b	
	s_6/s_1		

$$\text{Follow}(A) = \{a\}$$

↑
shift-reduce conflict

shift-reduce conflict
shift occurs if some other

production also conflicts in I₅

2) I₅

when $\text{follow}(A) \cap \text{follow}(B) \neq \emptyset$

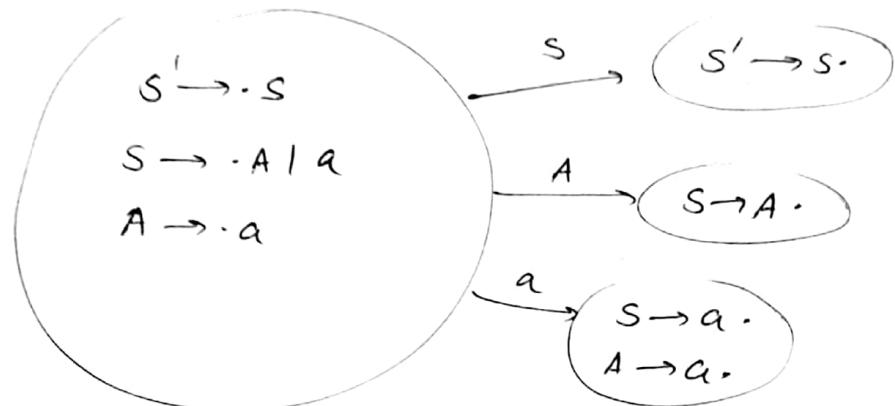
5	a	b	
	α_1/α_2	β_1/β_2	

↑
reduce-reduce conflict

Q- $S \rightarrow A \mid a$

$A \rightarrow a$

This is not LL(1) as $S \rightarrow \cdot A$ & $S \rightarrow \cdot a$ point to same terminal



case 2) of LR(0)

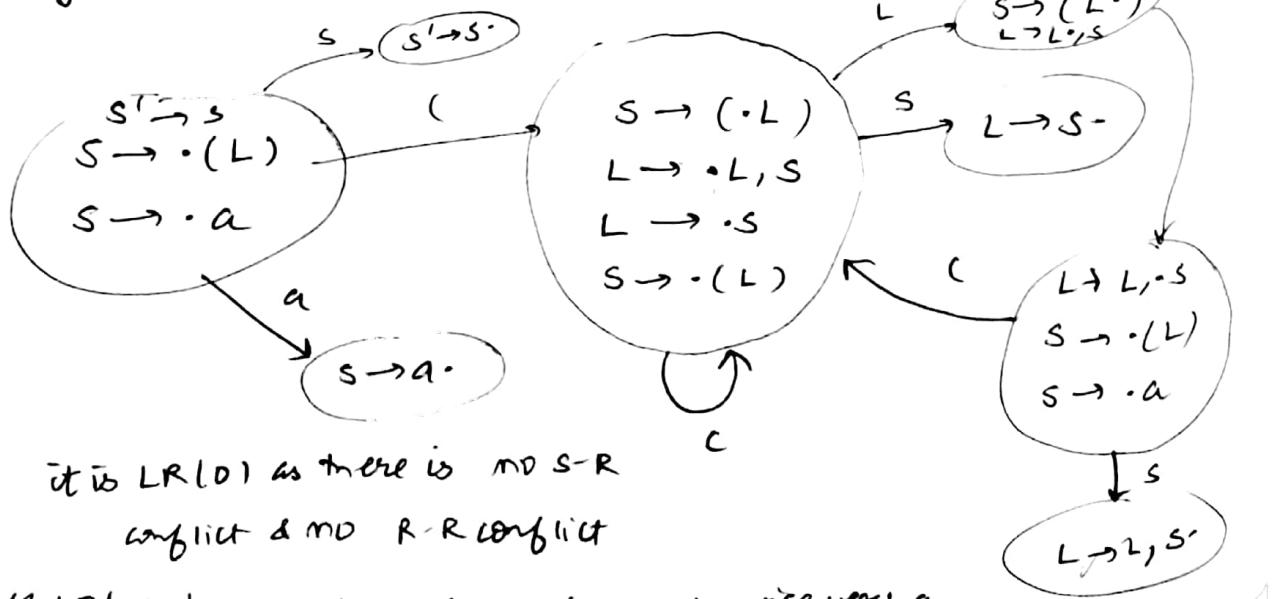
∴ not LR(0)

$\text{Follow}(S) = \text{Follow}(A) = \{\$\}$ ∴ not SLR(1)

Q- $S \rightarrow (L) \mid a$, $L \rightarrow S$
 $L \rightarrow L, S \mid S$

This is not LL(1) as $L \rightarrow S$ & $S \rightarrow (L) \mid a$ have same terminal
⇒ also $L \rightarrow L, S$ is left recursive ∴ not LL(1)

Always check for left recursion first



it is LR(0) as there is no S-R conflict & no R-R conflict

If LR(0) then SLR(1) always but not vice versa