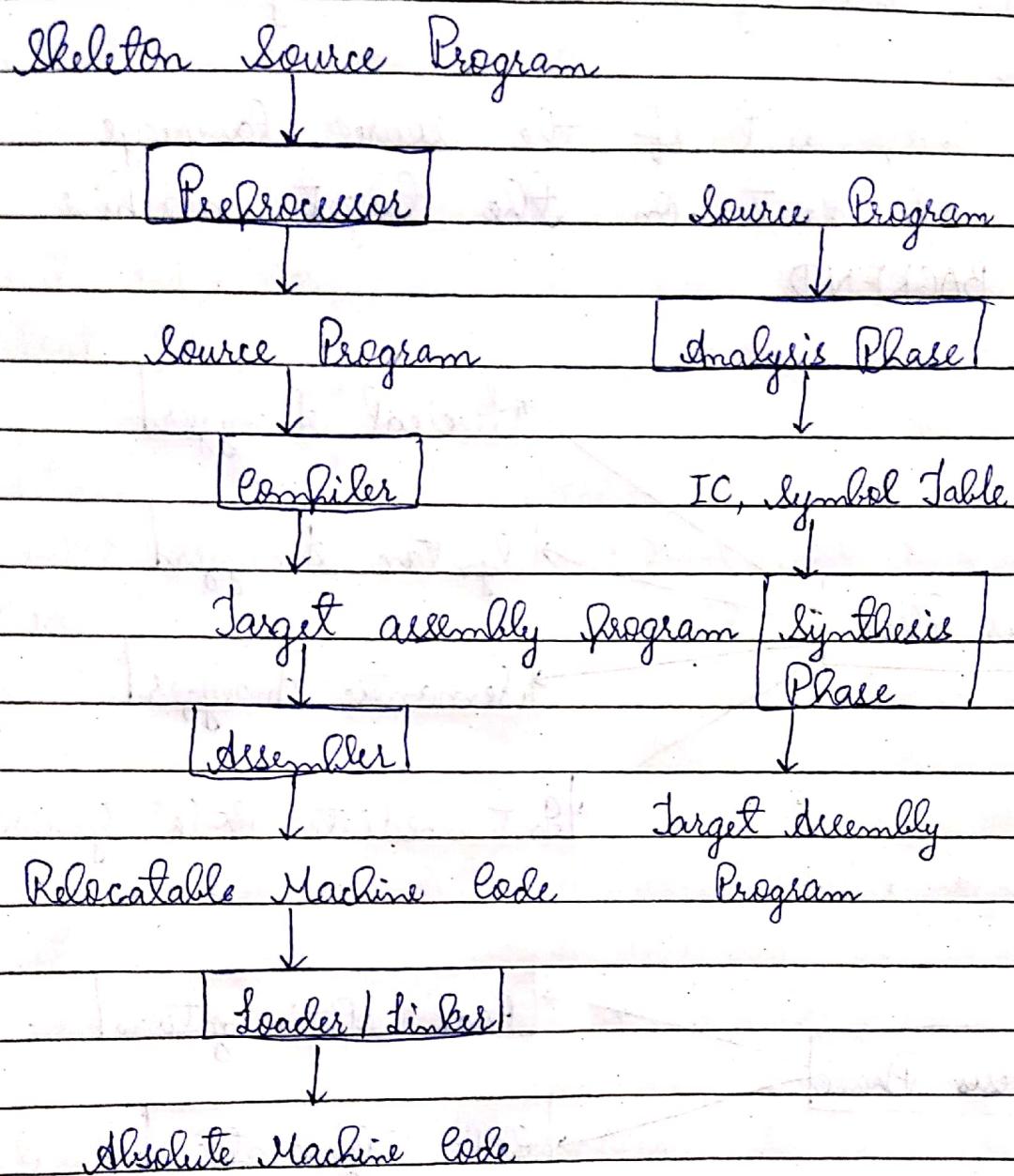


## \* Phases of Compiler

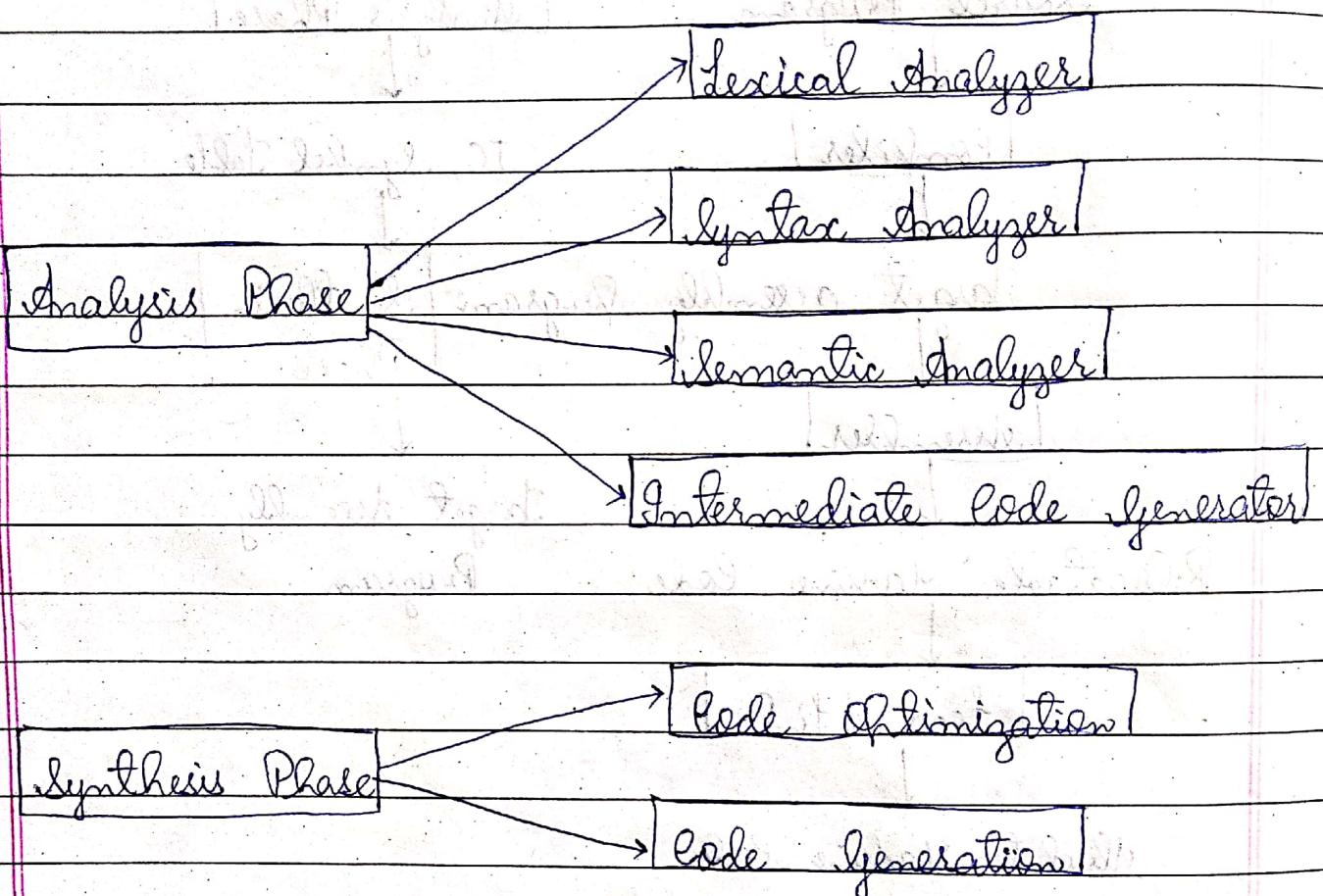


### Analysis Phase:

1. Perform analysis of input code
2. Generate errors and warnings
3. Largely dependent on the source language
4. Largely independent of the target machine
5. Creates FRONTEND

## Synthesis Phase:

1. Takes IC, uses symbol table and generates target program
2. Largely independent of the source language
3. Largely dependent on the target machine
4. creates BACKEND



## 1. Lexical Analysis:

- Grouping into tokens (sequence of characters with meaning)
- Identifies valid words or lexemes in source program
- Removes white spaces, tabs, newlines
- Passes sequence of tokens to syntactic analysis phase
- Implemented as a finite automata
- Error generated: invalid identifier

→ Example :

(i)  $\text{position} = \text{initial} + \text{rate} * 60$

Output of lexical phase:

$\langle \text{id}, 1 \rangle \Rightarrow \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle \text{number}, 60 \rangle$

(ii)  $\text{int } n1, n2;$

$\text{float } f;$

$n1 = n1 + f^* 5 - n2;$

Output of lexical phase:

$\langle \text{int} \rangle \langle \text{id}, 1 \rangle \langle \text{id}, 2 \rangle \langle ; \rangle \langle \text{float} \rangle \langle \text{id}, 3 \rangle \langle ; \rangle$   
 $\langle \text{id}, 1 \rangle \Rightarrow \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle \text{number}, 5 \rangle$   
 $\langle - \rangle \langle \text{id}, 2 \rangle \langle ; \rangle$

Limitations:

- Can identify validity of lexemes (individual words) only
- Not powerful to analyze expression or statement  
(Can't match a pair of parenthesis)

## 2. Syntax Analysis (Hierarchical analysis, Parsing)

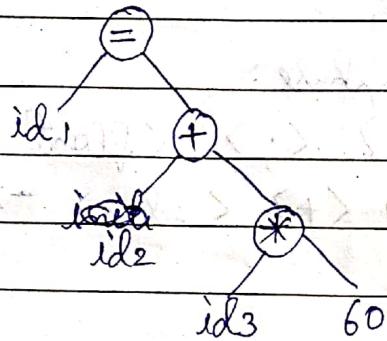
- Takes sequence of tokens from lexical analyzer
- Checks syntactic correctness of expression / blocks / statements / functions / programs
- Successful if it can identify grammar rules (CFG, CSG) for a sequence of tokens.
- Generates a parse tree from a sequence of tokens
- If parse tree is complete  $\Rightarrow$  construct is correct
- If parse tree is incomplete  $\Rightarrow$  syntax error in the construct

- uses recursive grammar rules
- Errors generated: missing parenthesis, missing semicolon, missing operand

Example:

- Position = initial + rate \* 60

Tokens: Position (id1), '='; initial (id2), '+', rate (id3)  
 '\*', 60



### 3. Semantic Analysis

- If program is syntactically correct, then check for semantic (meaning) correctness of the program
- Semantics depends on the programming language
- Checks for semantic errors
- Gathers type information
- Uses parse tree generated by the syntax phase
- Common checks are:

- Type of variables and type-casting
- Applicability of operators on operands - each operator has operands permitted by the source language  
 (String operators, only int index in array)
- Scope of variables and functions
- Determine definition of variables and functions  
 (function, overloading)

Example:

```
int a=5, b=2;  
float f;  
f = a/b;  
f = 5/2.0; //inbuilt-type casting
```

→ Errors:

- (i) float array index
- (ii) cannot convert 'int' to 'string'
- (iii) undeclared variable

#### 4. Intermediate Code Generator

→ After semantic analysis, Intermediate Code (IC) of source program will be created (low-level or machine-like IR)

→ Properties:

- (i) Easy to produce
- (ii) Easy to translate into target code

→ Variety of forms:

- (i) Three Address code (3A code)
- (ii) Prefix representation

(i) Three address code (3A code)

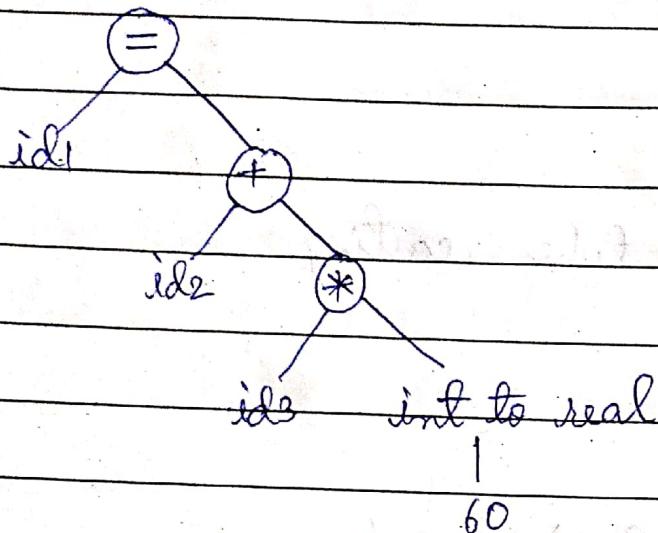
→ Each instruction has

- (a) Maximum three operands
- (b) Maximum one operator in addition to assignment

→ Compiler will decide order of evaluation

→ Compiler will generate a temp name to hold computed results

→ Example:  $\text{id}_1 = \text{id}_2 + \text{id}_3 * 60$



$t_1 = \text{int to real}(60)$

$t_2 = \text{id}_3 * t_1$

$t_3 = \text{id}_2 + t_2$

$\text{id}_1 = t_3$

## 5. Code Optimization

→ Improves the Intermediate code for

(i) Faster execution of machine code.

(ii) Space optimization

→ Optimizing compilers

→ Machine independent code optimization

→ Example: (i) 60.0 instead of int to real(60)

(ii) No need of  $t_3$

$t_1 = \text{id}_3 * 60.0$

$\text{id}_1 = \text{id}_2 + t_1$

→ Removes dead code like unreachable code

→ If anything is written after return statement in function

→ Loop optimizations - major source of optimization

→ Saving even a single line in loop improves execution time significantly

→ Example,

$b = 5;$

for ( $i=0; i < 10000; i++$ ) {

$a = b + 100 - i * 2;$

}

Optimized code:

$b = 5;$

$C = b + 100;$

for ( $i=0; i < 10000; i++$ ) {

$a = C - i * 2;$

}

## 6. Code Generation

- Takes an optimized Intermediate code as an input and converts it to the target code
- The target code can be:
  - (i) Relocatable Machine code
  - (ii) Assembly code
- Memory locations are selected for each variables of a program
- Intermediate code is translated to a sequence of machine instructions
- Part of backend of the compiler design
- Target code generation depends on:
  - (i) availability of machine instructions
  - (ii) addressing modes
  - (iii) number of registers (general and special purpose)

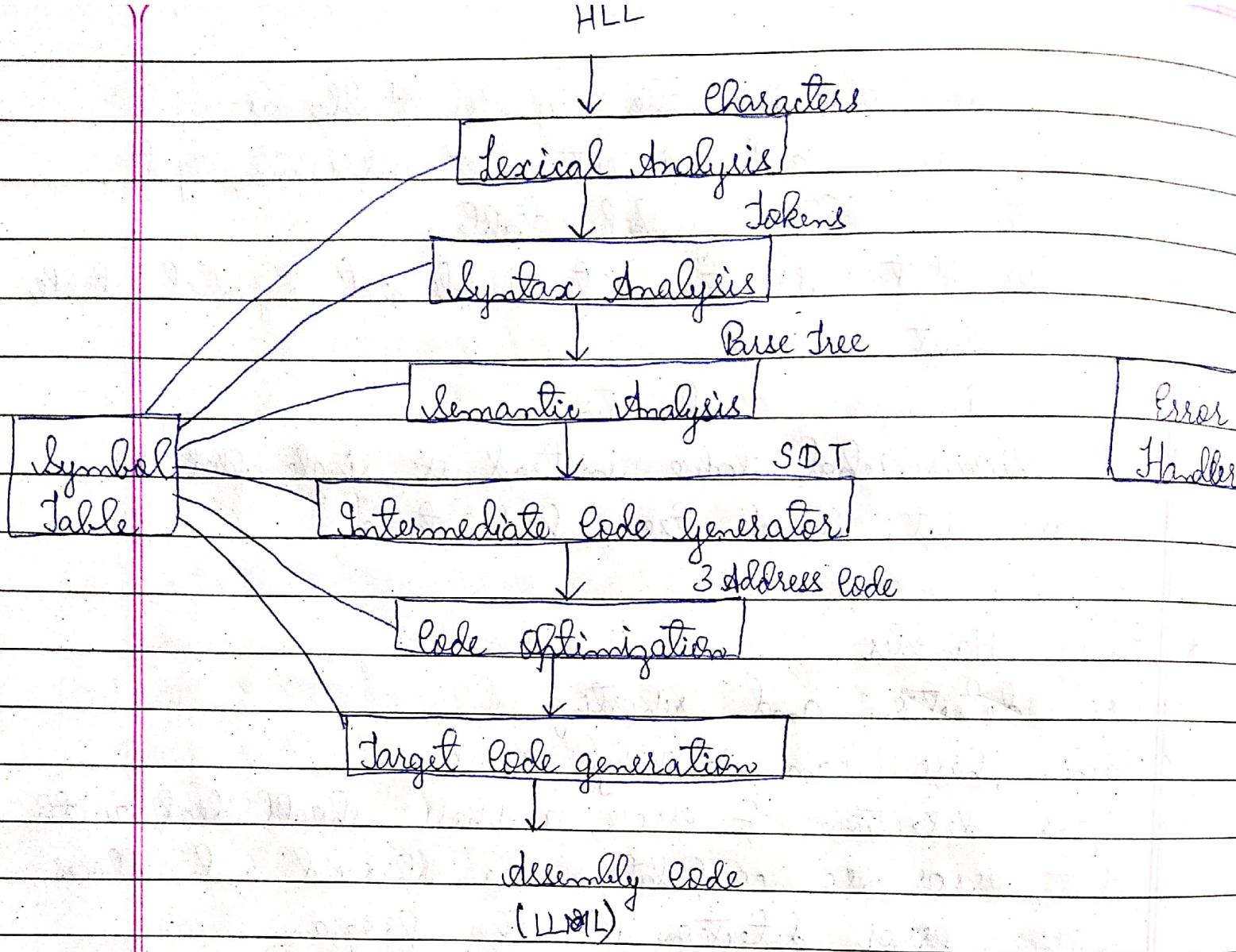
## \* Symbol Table Management:

- Interacts with all phases - used as a reference table by all the phases
- Compiler needs
  - (i) To record the identifiers (variables) used in source program
  - (ii) To collect information about various attributes of each identifier
  - (iii) Attributes for variable names - Name, type, class, size, relative offset within program, scope, value
  - (iv) Attributes for function names - Name, number of arguments, types of arguments, return type, method of passing each parameter
- It is a data structure containing a record for each identifier with fields for the attributes of the identifiers.
- Lexical analyzer: detects identifier and enters in the symbol table, but can't determine all its attributes.
- The remaining phases enter information about identifiers into symbol table and use the entered information in different ways.
- Semantic analyzer and Intermediate code generation needs to know type of variable.
- Code generator: enters and uses details about the storage assigned to identifier.
- During compilation, all phases often look up this table for definition of variable, example variable used is defined before that or not?

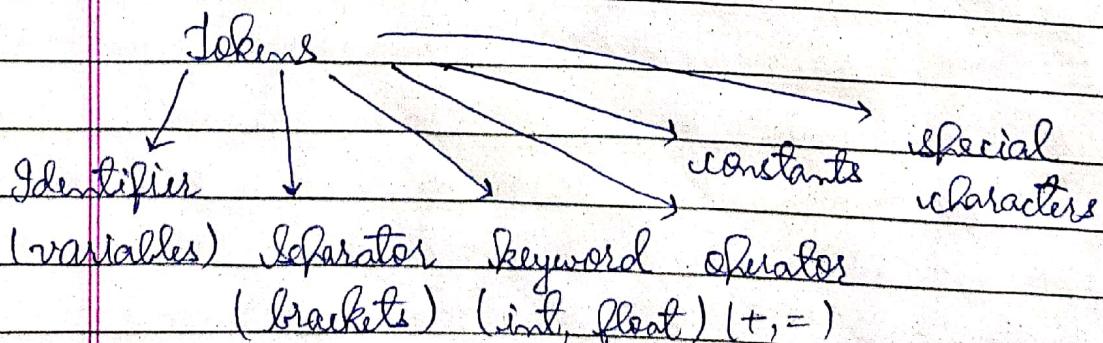
- Hence, operations in the symbol table are allocate, free, search, insert, set attribute, get attribute in the symbol table.
- Use the data structure to implement symbol table such that
  - (i) It minimizes search time
  - (ii) the hierarchical table instead of single flat one
  - (iii) linear list, search trees, hash table

### \* Error Handler

- Error detection and reporting
- Each phase can identify errors
- After detection of errors, a phase should deal with that error so compilation can proceed and allows further error detection in the program
- Semantic error - parser can progress
- Syntax error - parser reaches in error state
- If indicate regarding type of error
- Undo some processing already carried out by parser - error recovery
- Error detection and recovery - vital role in overall operations of compiler



- \* Lexical Analysis (Lexer, Tokenizer, Scanner)
    1. Tokenization
    2. Give error message
      - Exceeding length
      - Unmatched string
      - Illegal character
    3. Eliminate comments, white space (tab, blank space, new line)



Example,

1. int max()

{  
5 /\* Find max of a and b \*/ → comment

int a = 30, b = 20; 14.

if (a < b) 20

— return(b); 25

else

return(a); 31

{ (32)

2. printf("i = %.d, &i = %ox") ( i, &i );

3                  4        5      6      7      8      9      10

3. int main()

{ 5

x = y + z; 11

int x, y, z; 18

printf("sum %d %d", x);

{ (26)

4. main() { } 4

a = b ++ + -- - - + ++ = ;

printf("%d %d, a, b);

{ 15    16    17

\* Finding First() and Follow() from given grammar

First(A) contains all terminals present in first place of every string derived by A.

First(Terminal) = Terminal

First( $\epsilon$ ) =  $\epsilon$

E.g.  $S \rightarrow abc \mid defghi$

First(S) = {a, d, g}

E.g.  $S \rightarrow ABC \mid ghi \mid jkl$  ( $A, g, j$ )  $\Rightarrow (a, g, j) \cup (a, b, c, g, j)$

$A \rightarrow alb \mid c$  ( $a, b, c$ )

$B \rightarrow b$  ( $b$ )

$D \rightarrow \delta$  ( $\delta$ )

E.g.  $S \rightarrow ABC$

$A \rightarrow alb \mid c$

$B \rightarrow c \mid d \mid e$

$C \rightarrow e \mid f \mid e$

$F(C) = \{e, f, \epsilon\}$ ,  $F(B) = \{c, d, \epsilon\}$ ,  $F(A) = \{a, b, e\}$

$F(S) = F(A) = \{a, b, \epsilon\} \cup F(B) = \{c, d\} \cup F(C) = \{e, f, \epsilon\}$   
 $= \{a, b, c, d, e, f, \epsilon\}$

E.g.  $E \rightarrow TE'$

$E' \rightarrow *TE' | \epsilon$

$T \rightarrow FT'$

$T' \rightarrow \epsilon | +FT'$

$F \rightarrow id | ( E )$

$\text{first}(F) = \{id, ()\}$ ,  $\text{first}(T) = \text{first}(F) = \{id, ()\}$ ,

$\text{first}(T') = \{\epsilon, +\}$ ,  $\text{first}(E') = \{*, \epsilon\}$ ,

$\text{first}(E) = \text{first}(T) = \{id, ()\}$

$\text{Follow}(A)$  contains set of all terminals present immediate in right of A. It never contains  $\epsilon$ .  
 $\text{Follow}$  of start symbol is  $\$$ .

E.g.  $S \rightarrow ACD$

$C \rightarrow a1b$

$F_0(A) = \text{first}(C) = \{a, b\}$

$F_0(D) = F_0(S) = \{\$\}$

E.g.  $S \rightarrow a \overset{\curvearrowright}{S} b \overset{\curvearrowright}{S} | b \overset{\curvearrowright}{S} a \overset{\curvearrowright}{S} | \epsilon$

$F_0(S) = \{\$, b, a\}$

E.g.  $S \rightarrow A \overset{\curvearrowright}{a} A \overset{\curvearrowright}{b} | B \overset{\curvearrowright}{b} B \overset{\curvearrowright}{a}$

$A \rightarrow E$

$B \rightarrow E$

$F_0(A) = \{a, b\}$ ,  $F_0(B) = \{b, a\}$

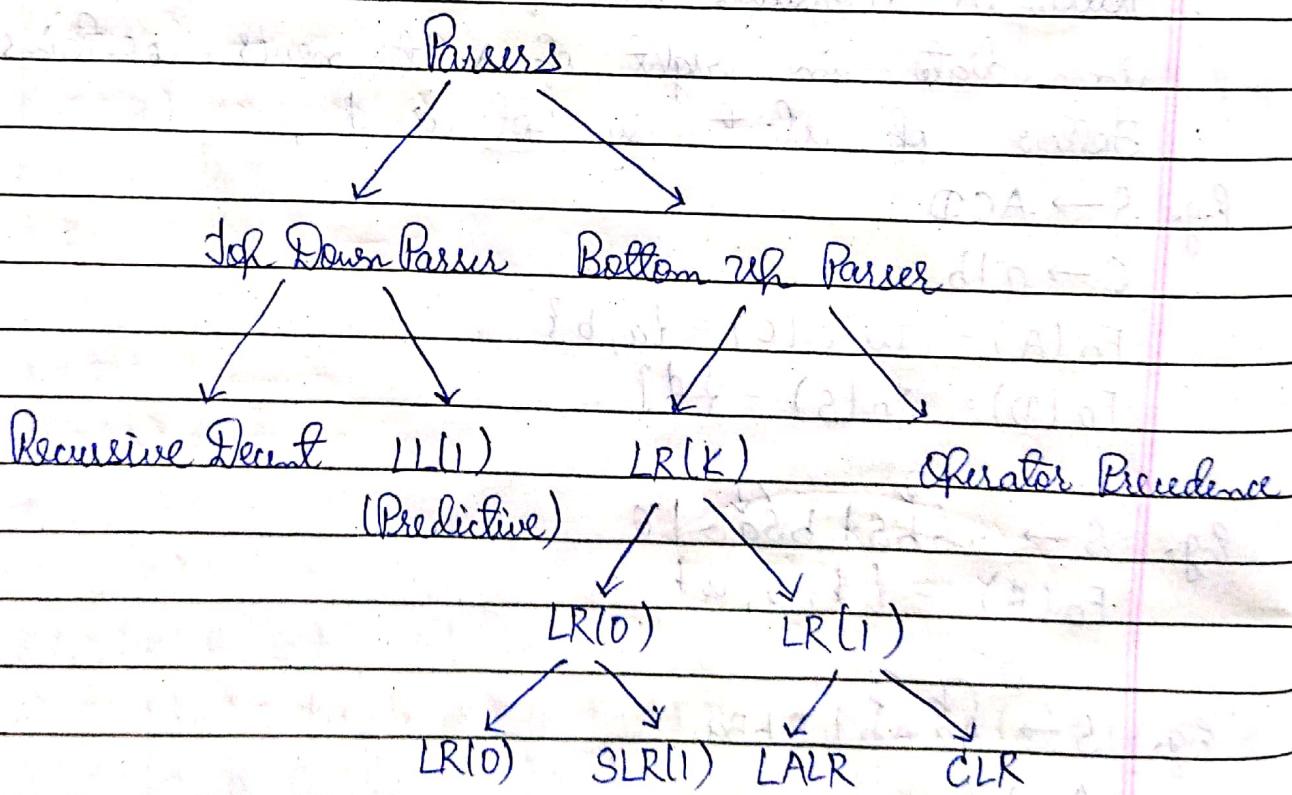
E.g.  $S \rightarrow ABC$   
 $A \rightarrow DEF$

$B, C \rightarrow \epsilon, D \rightarrow \epsilon, E \rightarrow \epsilon, F \rightarrow \epsilon$

$F_0(A) = \text{First}(B) = \text{First}(C) = \text{Follow}(S) = \{\$\}$

### \* Parsing

It is a process of deriving string from a given grammar.

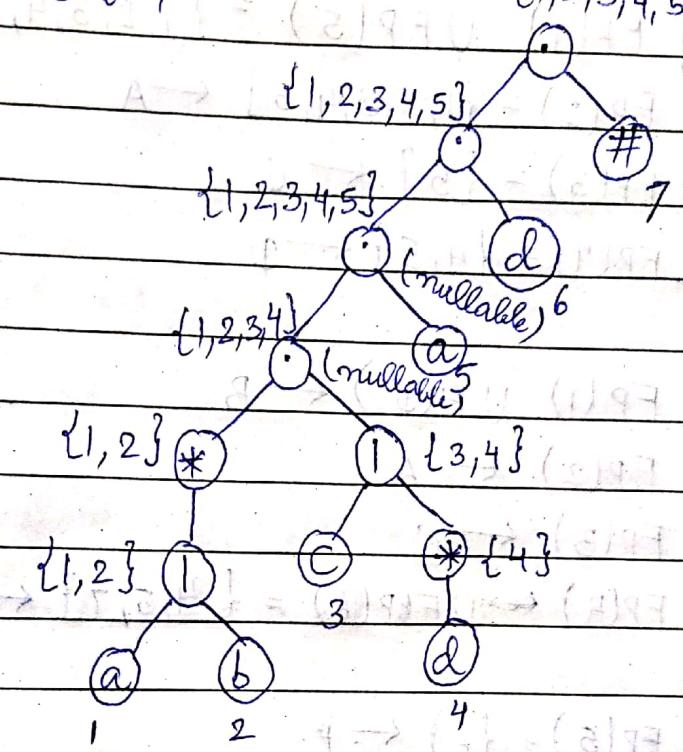


\* RE to DFA using syntax tree method

E.g.  $(ab)^* \cdot (cd)^* \cdot ad \#$

1 2 3 4 5 6 7

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



Union of

$\text{FIRSTPOS}(1) = \text{FIRSTPOS}(\text{child nodes})$

$\text{FIRSTPOS}(* \cdot *) = \text{FIRSTPOS}(\text{child node})$

$\text{FIRSTPOS}(\cdot) = \text{FIRSTPOS}(\text{left node})$

(OR)

Union of  $\text{FIRSTPOS}(\text{left node})$  and

$\text{FIRSTPOS}(\text{right node})$  if left node  
contains \* (that is nullable)

FOLLOWPOS - which symbols can follow the given symbol

$\text{FIRSTPOS}(\text{root}) = \{1, 2, 3, 4, 5\} \leftarrow A$

$\text{FOLLOWPOS}(1) = \{1, 2, 3, 4, 5\}$

$\text{FOLLOWPOS}(2) = \{1, 2, 3, 4, 5\}$

$\text{FOLLOWPOS}(3) = \{5\}$

$\text{FOLLOWPOS}(4) = \{4, 5\}$

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

$\text{FOLLOWPOS}(6) = \{7\}$

$\text{FOLLOWPOS}(7) = \#\emptyset$

same variable

$$A \left\{ \begin{array}{l} FP(1) \cup FP(5) = \{1, 2, 3, 4, 5, 6\} \leftarrow B \\ FP(2) = \{1, 2, 3, 4, 5\} \leftarrow A \\ FP(3) = \{5\} \leftarrow C \\ FP(4) = \{4, 5\} \leftarrow D \end{array} \right.$$

$$B \left\{ \begin{array}{l} FP(1) \cup FP(5) \leftarrow B \\ FP(2) \leftarrow A \\ FP(3) \leftarrow C \\ FP(4) \leftrightarrow FP(6) = \{4, 5, 7\} \leftarrow E \end{array} \right.$$

$$C \{ FP(5) = \{6\} \leftarrow F \}$$

$$D \{ FP(4) \leftarrow D \}$$

$$\{ FP(5) \leftarrow F \}$$

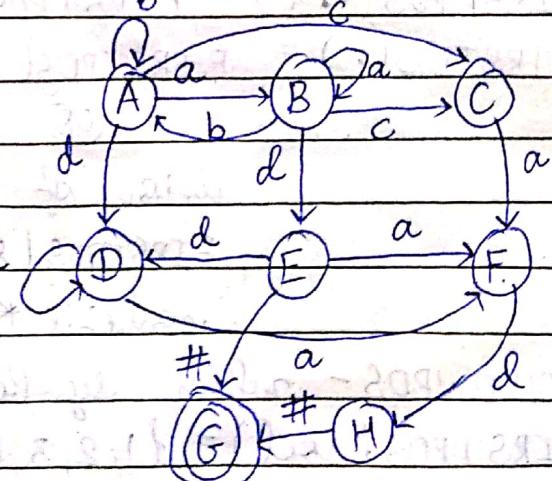
$$E \{ FP(4) \leftarrow D \}$$

$$\{ FP(5) \leftarrow F \}$$

$$\{ FP(7) \leftarrow \# \leftarrow G \}$$

$$F \{ FP(6) = \{7\} \leftarrow H \}$$

$$G \{ FP(7) \leftarrow \# \}$$



## \* LL(1) Parser (Top-Down)

(\*) How to check whether the string is LL(1) or not  
 e.g.  $S \rightarrow (L) | a$

$$L \rightarrow SL'$$

$$L' \rightarrow \epsilon, SL'$$

$$F(S) = \{(, a\} \quad F_0(S) = \{ \$, \geq, ), \}$$

$$F(L) = F(s) = \{(, a\} \quad F_0(L) = \{\}\}$$

$$F(L') = \{\epsilon, \geq, '\} \quad F_0(L') = \{\}\}$$

Make Parsing Table where rows are non-terminals and columns are terminals and  $\$$ .

	(	)	a	,	\$
S	$S \rightarrow (L)$	e	$S \rightarrow a$	e	$\epsilon$
L	$L \rightarrow SL'$	e	$L \rightarrow SL'$	e	$\epsilon$
L'	$L' \rightarrow \epsilon$	e	$L' \rightarrow SL'$	e	$\epsilon$
<u>FOLLOW</u>					

Make entry for every row of the table where terminal's first is non-terminal (column). The entry would be equation whose first is non-terminal. If first contains  $\epsilon$  then use follow of that non-terminal. All other entries will be  $e$ . If one cell maximum one entry then the grammar will be accepted by LL(1) Parser.

The above grammar is LL(1).

E.g.  $S \rightarrow aSbS \mid bSaS \mid \epsilon$

$$F(S) = \{a, b, \epsilon\}$$

$$FO(S) = \{b, a, \$\}$$

a	b	\$
$S \rightarrow aSbS$	$S \rightarrow bSaS$	$S \rightarrow \epsilon$

Here, there are two entries twice for one cell so the above grammar is not LL(1).

\* How LL(1) parser works?

$$S \rightarrow AA$$

$$A \rightarrow aA$$

$$A \rightarrow b$$

a	b	\$
$S \rightarrow AA$	$S \rightarrow AA$	$\epsilon$
$A \rightarrow aA$	$A \rightarrow b$	$\epsilon$

String: ababb and stack is \$

Initially, there is first non-terminal on the top of stack followed by \$ and string is also followed by \$. So at every point; non-terminal at top of stack is replaced with the function present in row of non-terminal of stack and column of terminal of string where pointer is. If top of stack and string pointer has same terminal

Then stack is popped and string pointer is incremented. If there are no rules present for a given non-terminal in stack according to table then grammar is not LL(1). If top of stack becomes \$ and string pointer also becomes \$ then string is LL(1).

Note: While replacing non-terminal with rule, all the contents are added in reverse order.

	Stack	ab ab \$	Rule
①	\$ S	$\uparrow \uparrow \uparrow \uparrow$	$S \rightarrow AA - A$
	\$ AA		$A \rightarrow aA$
	\$ AAA	(pop & increment)	
②	\$ AA	$\uparrow \uparrow \uparrow \uparrow$	$A \rightarrow b - A$
	\$ Ab	(pop & increment)	
	\$ A	$\uparrow \uparrow \leftarrow T$	$A \rightarrow aA - T$
	\$ Aa	(pop & increment)	
	\$ A	$\uparrow \uparrow \leftarrow T$	$A \rightarrow b - A$
	\$ b	(pop & increment)	
	\$	(Accepted)	

- ① Check table for rules in row S and column a.
- ② Check table for rule in row A and column b.

## \* Conditions for Top-Down Parser

The grammar should not contain:

- Left Recursion
- Ambiguity (Left Factoring) Non-Determinism

### 1. Eliminating Left Recursion

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

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

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

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid B_1 \mid B_2 \mid \dots \mid B_n$$

$$A \rightarrow B_1 A' \mid B_2 A' \mid B_3 A' \mid \dots \mid B_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid e$$

$$E \rightarrow TAE'$$

$$T \rightarrow FT'$$

$$E' \rightarrow +TA'E' \mid -TA'E' \mid \epsilon$$

$$S \rightarrow abA \mid abc \mid Bba \mid Sdd \mid bs$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$S \rightarrow abAS' \mid Bbas' \mid bss'$$

$$S' \rightarrow aes' \mid dds' \mid e$$

$$A \rightarrow a$$

$$B \rightarrow b$$

2. Eliminating left - factoring

$$A \rightarrow \alpha B_1 | \alpha B_2 | \dots | \alpha B_n$$



$$A \rightarrow \alpha A'$$

$$A' \rightarrow B_1 | B_2 | \dots | B_n$$

e.g.  $S \rightarrow Aa | Ab | AS$

$$S \rightarrow AS'$$

$$S' \rightarrow ab | b$$

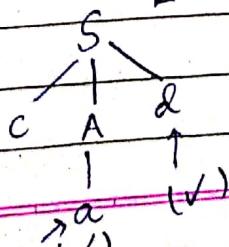
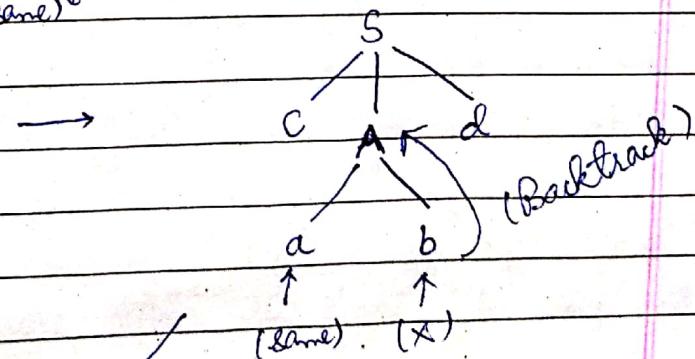
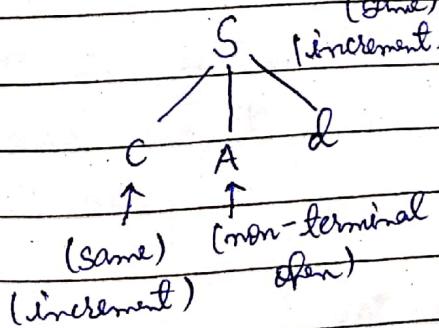
\* Recursive Descent Parser (Top - Down)

A Recursive Descent Parser is a top-down parser built from a set of mutually recursive procedures where each such procedure implements one of the non-terminals of the grammar.

e.g.  $S \rightarrow cAd$

$$A \rightarrow ab | a$$

Input string: cad



E.g.  $S \rightarrow aSbS \mid bSas \mid e$

String: abab

$S \rightarrow aSbS$

$\Rightarrow a \cdot aSbS \ bS \quad \times \text{(Backtracking)}$

$\Rightarrow a \ bSas \ bS$

$\Rightarrow ab \ aSbS \ aSbS \quad \times \text{(Backtracking)}$

$\Rightarrow ab \ bSas \ aSbS \quad \times \text{(Backtracking)}$

$\Rightarrow ab \ e \ aSbS$

$\Rightarrow aba \ aSbS \ bS \quad \times \text{(Backtracking)}$

$\Rightarrow aba \ bSas \ bS \quad \times \text{(Backtracking)}$

$\Rightarrow aba \ e \ bS$

$\Rightarrow abab \ S$

$\Rightarrow abab \ aSbS \quad \times \text{(Backtracking)}$

$\Rightarrow abab \ bSas \quad \times \text{(Backtracking)}$

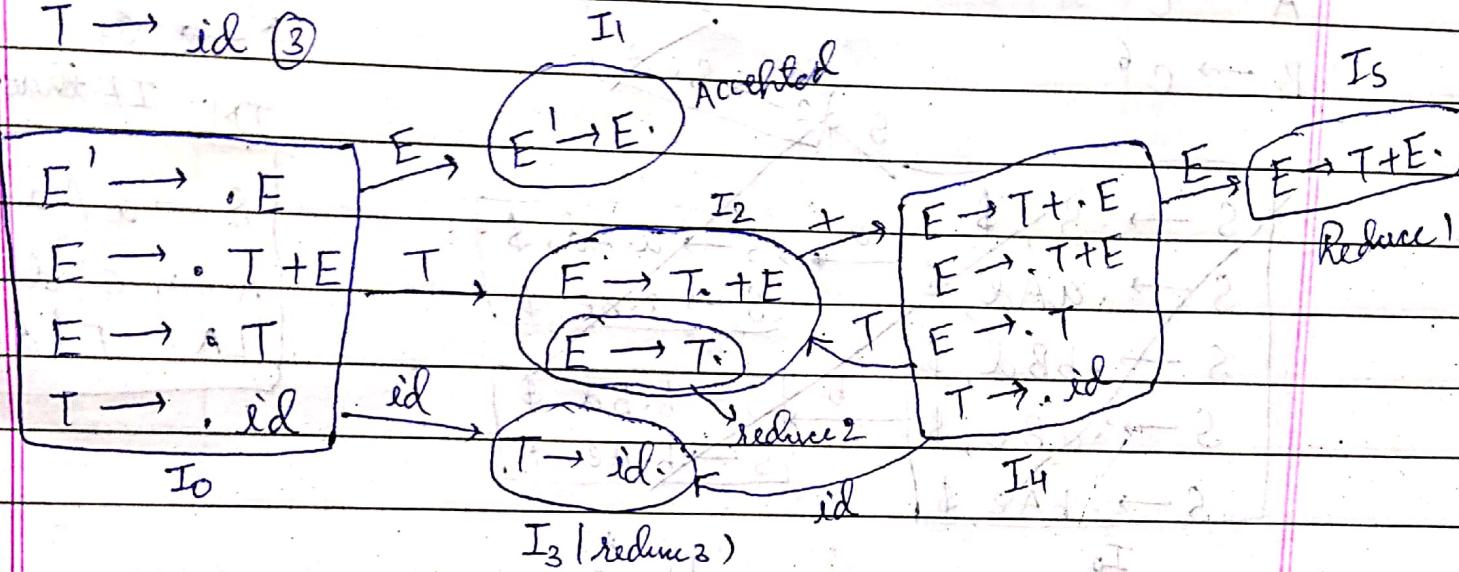
$\Rightarrow abab \ e$

$\Rightarrow abab$

\* SLR(1) (Bottom-up)

E.g.  $E \rightarrow T + E \mid T$

$T \rightarrow id \quad (3)$



state	id + \$	$E \rightarrow T$
0	S <sub>3</sub>	
1		Accepted
2	S <sub>4</sub>	$\text{FO}(E)$
3	S <sub>4</sub>	S <sub>3</sub>
4	S <sub>3</sub>	5 $\leftarrow$ 2
5		21

Find follow of LHS of reduce equation

$E \rightarrow T$ .

$\text{FO}(E) = \{ \$ \}$

$\text{FO}(T) = \{ \$, + \}$

Write reduce in row having same state number and shift in the row from where equation is going.

SLR(1) ✓ (Not more than one entries in any cell)

\* CLR (Bottom - Up)

IMP

$S \rightarrow a^1 Ad | b^2 Bd | a^3 Be | b^4 Ae$

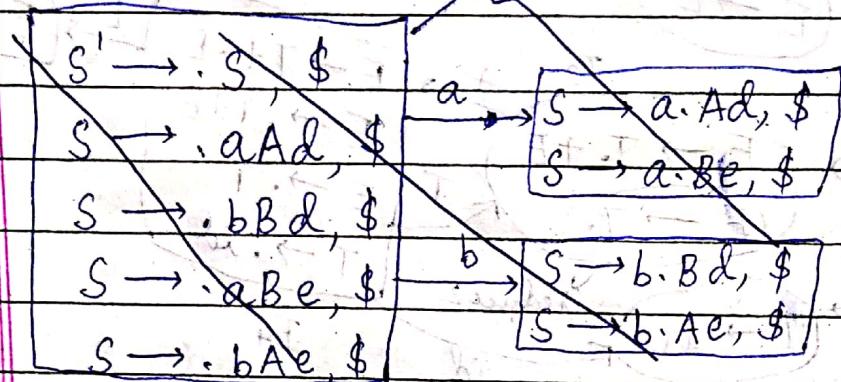
$A \rightarrow c^5$

$B \rightarrow c^6$

[By default,  
lookahead is \$]

~~\$~~  $\rightarrow$   $\$$

IMP IF there NT after \$



$S \rightarrow a.Ad, \$$

$A \rightarrow .c, d$

$$LA = F(d, \$) = d$$

$I_0$

$I_1 *$

$S' \rightarrow .S, \$$

$S \rightarrow .aAd, \$$

$S \rightarrow .bBd, \$$

$S \rightarrow .aBe, \$$

$S \rightarrow .bAe, \$$

$S \rightarrow .S, \$$

$S \rightarrow a.Ad, \$$

$A \rightarrow .c, d$

$S \rightarrow a.Be, \$$

$B \rightarrow .c, e$

$I_4$

$I_6 *$

$A \rightarrow S \rightarrow aAd, \$$

$B \rightarrow S \rightarrow ab.e, \$$

$A \rightarrow C., d$

$B \rightarrow C., e$

$I_0$

$b$

$I_3$

$I_9$

$S \rightarrow b.Bd, \$$

$B \rightarrow .c, d$

$S \rightarrow b.Ae, \$$

$A \rightarrow .c, e$

$S \rightarrow bB.d, \$$

$A \rightarrow S \rightarrow bA.e, \$$

$C \rightarrow B \rightarrow .c, d$

$A \rightarrow C., e$

If state is accepted then write reduce in column of its lookahead.

State	a	b	c	d	e	\$	A	B
0	$s_2$	$s_3$						
1							$s_1$	
2								
3				$s_6$	$s_7$	$s_8$		
4				$s_9$				
5						$s_{10}$		
6							$s_{11}$	
7							$s_5$	$s_6$
8							$s_{12}$	
9							$s_{13}$	
10								$s_{10}$
11								$s_{11}$
12								$s_{12}$
13								$s_{10}, s_4$

### \* LALR (Bottom-up)

LALR is similar to CLR but in LALR similar states are merged.

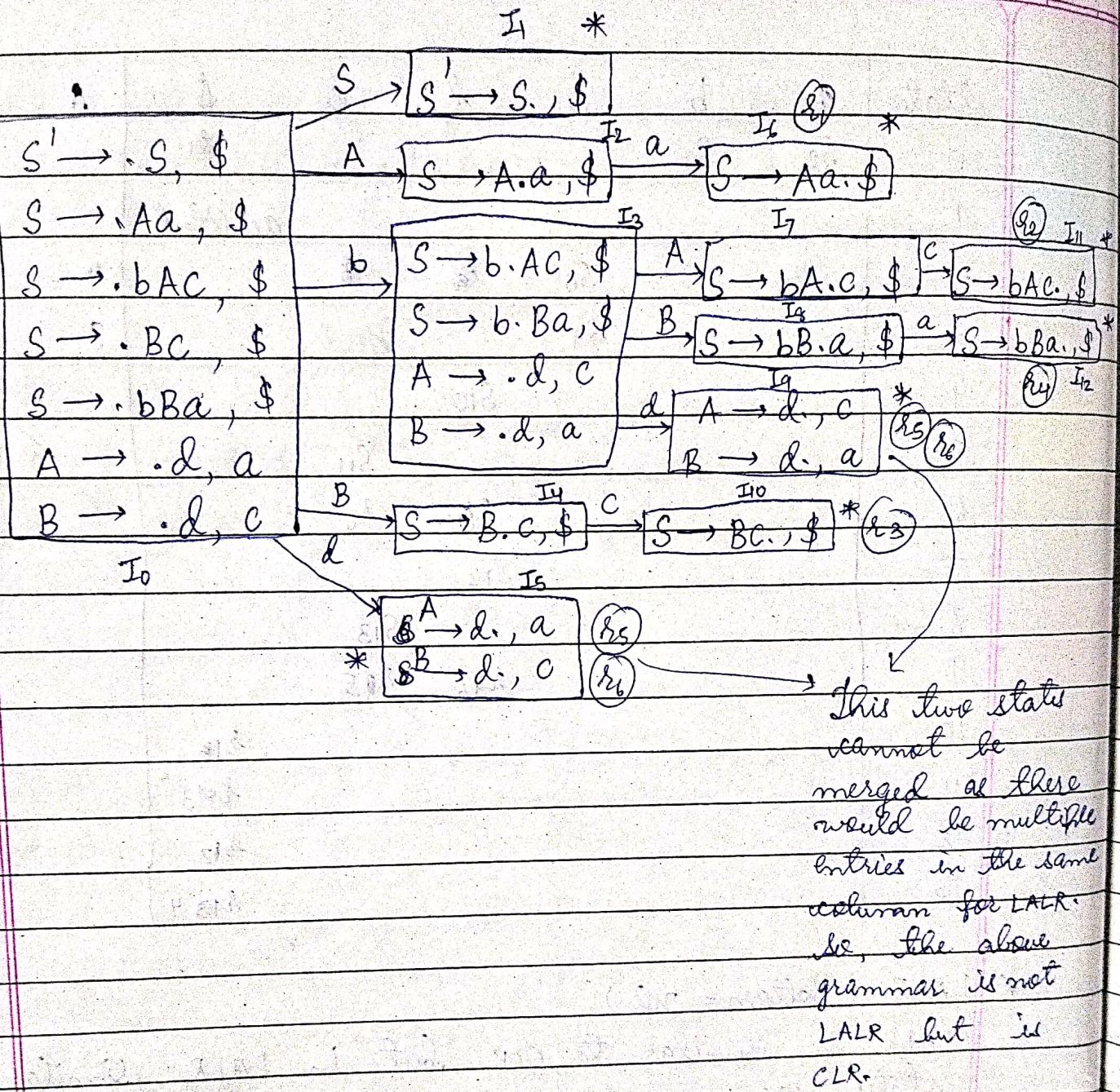
$$S \rightarrow Aa \mid bAc \mid Bc \mid bBa$$

$$A \rightarrow d^5$$

$$B \rightarrow d^6$$

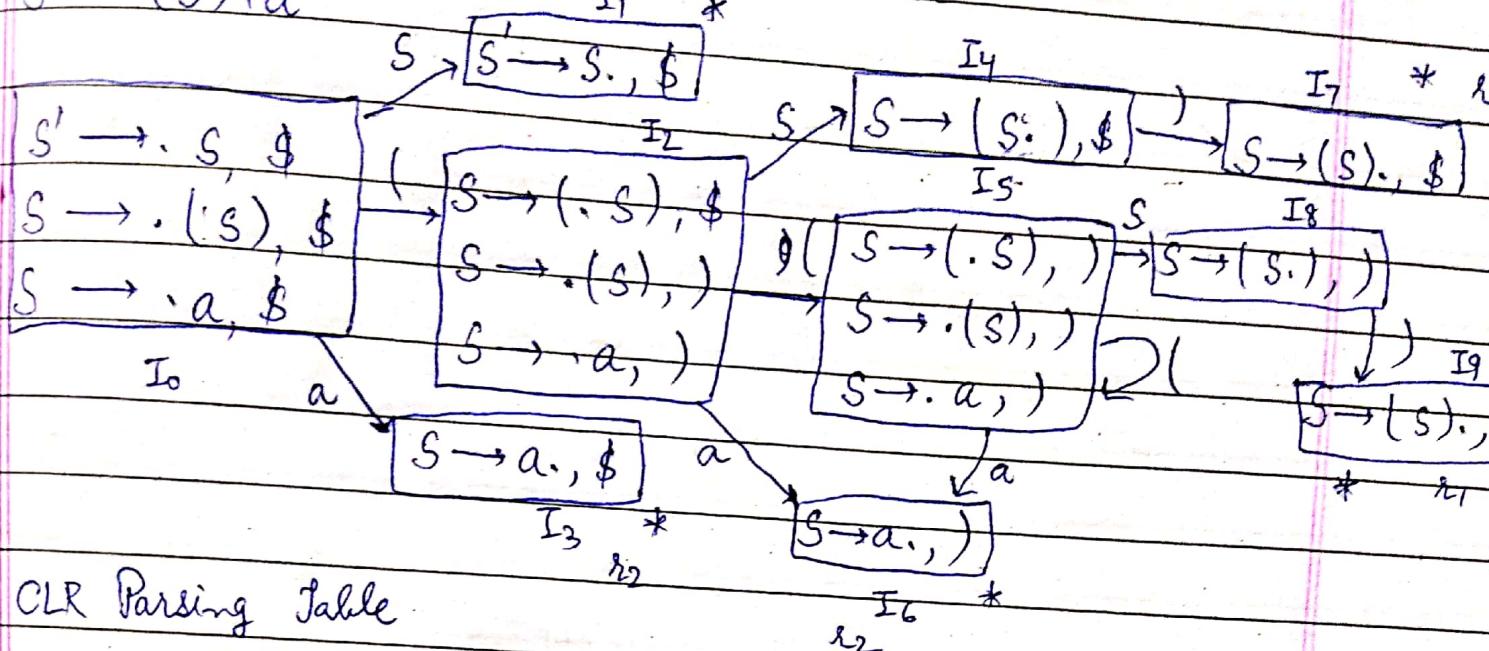
In LALR, two reduce with different look ahead can be merged.

IMP Sometimes, a grammar can be CLR but is not LALR.



\* Difference between CLR and LALR.

$$S \rightarrow (S) | a^2$$



CLR Parsing Table

	(	)	a	\$	S
0	S <sub>2</sub>		S <sub>3</sub>		1
1				accepted	
2	S <sub>5</sub>	65	S <sub>6</sub>		4
3				r <sub>2</sub>	
4		S <sub>7</sub>			
5	S <sub>5</sub>		S <sub>6</sub>		8
6		r <sub>2</sub>			
7				r <sub>1</sub>	
8		S <sub>9</sub>			
9		r <sub>1</sub>			

LALR Parsing Table

	(	)	a	\$	S
0	S <sub>2</sub>		S <sub>3</sub>		1
1			accepted		
2	S <sub>5</sub>		S <sub>6</sub>		48
3		r <sub>2</sub>		r <sub>2</sub>	
4		S <sub>7</sub>			
5		r <sub>1</sub>		r <sub>1</sub>	