

Subject: SPCC

Sem: II

Assignment no. 2

1. What are different phases of compiler?

Illustrate compiler's internal representation
 of source program for following statement
 after each phase Position = initial + rate * 60.

- ① Lexical Analysis: a) The first phase of a compiler is called lexical analysis or scanning.
 b) The lexical analyzer reads the stream of characters making up source program and groups the characters into meaningful sequences called lexemes.
- ② For each lexeme, the lexical analyzer produces as output a token of the form $\langle \text{token-name}, \text{attribute-value} \rangle$ that it passes on to the subsequent phase, syntax analysis. In the token,
 c) In the token, the first component token-name is an abstract symbol that is used during syntax analysis, and the second component attribute-value points to an entry in the symbol table for this token.
- ③ Information from the symbol-table entry is needed for semantic analysis and code generation.

Lexical analysis of Position = initial + rate * 60 is given by:

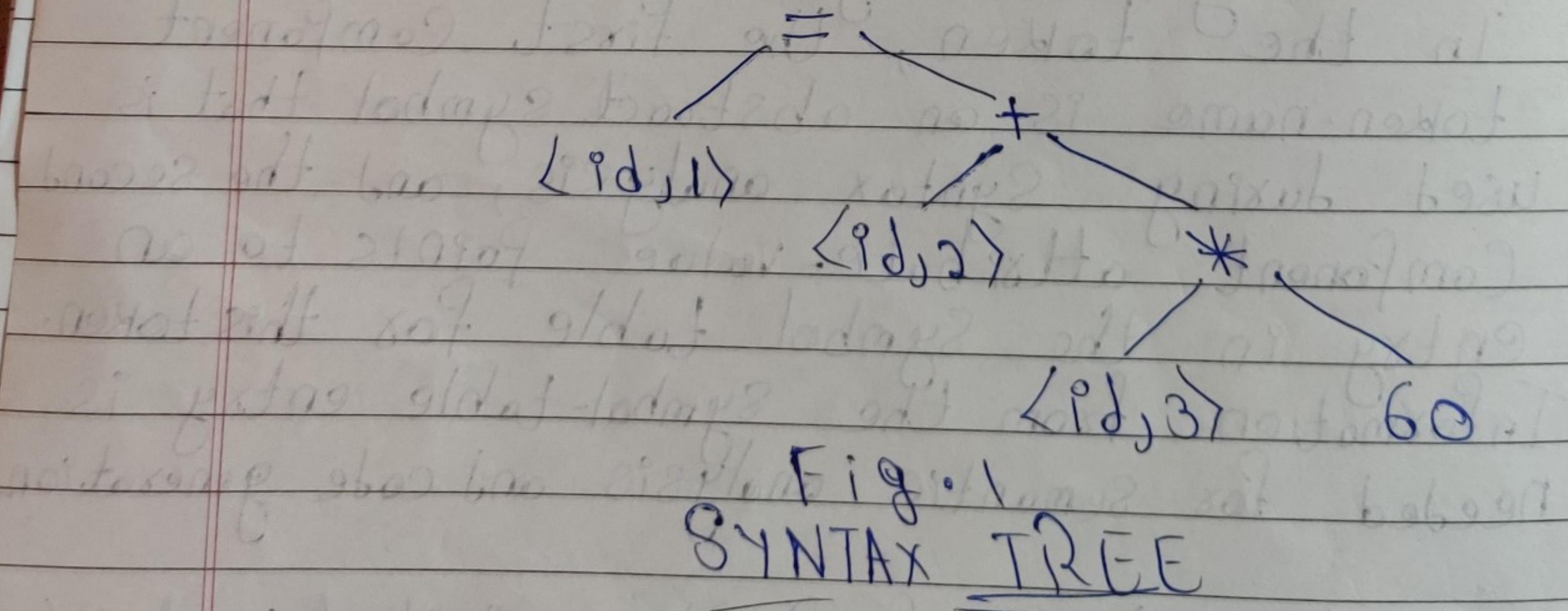
① - <id, 1> {=} <id, 2> {+} <id, 3> {*} <60>

For some tokens attributes values are implicit.

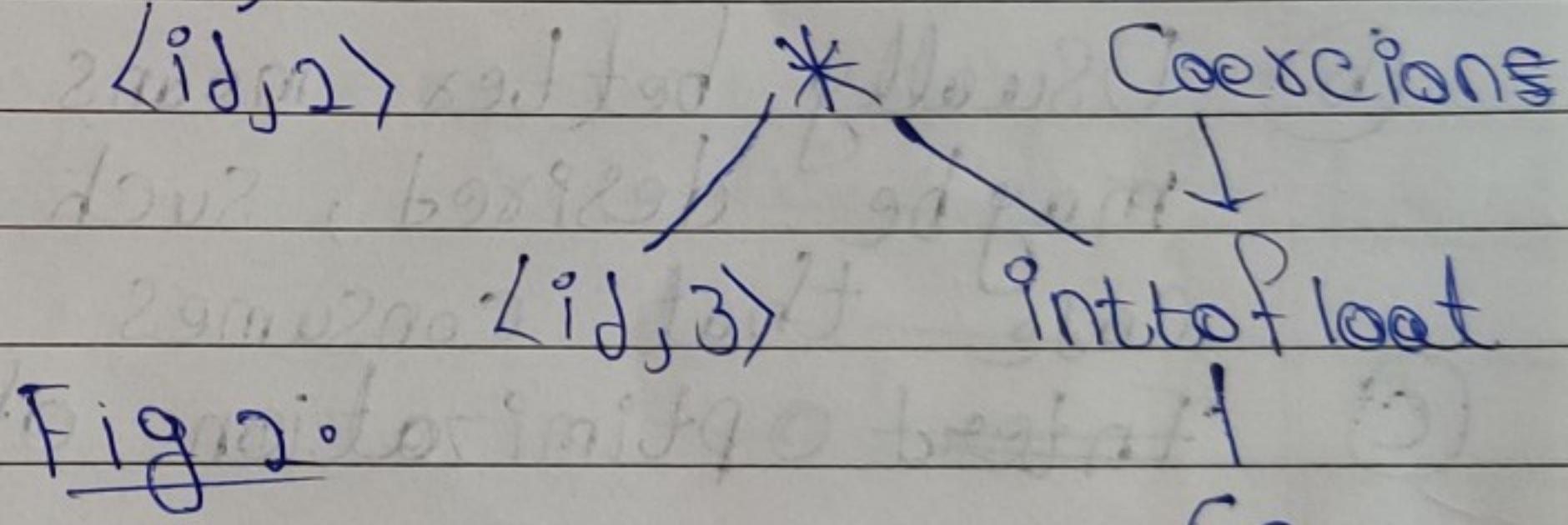
↓
fixed name

Pointed	1	Position initial state	Properties
	2		
	3		
			for 60 {number, 4} representation can also be used.

- ② Syntax Analysis : ① The second phase of compiler is Syntax analysis or Parsing.
- ③ The parser uses the first components of the tokens produced by the lexical analyzer to create a tree like intermediate representation that depicts the grammatical structure of the token stream.
- ④ The intermediate node represents operation & the children represent the arguments of the operation.
- ⑤ A syntax tree for the token stream ① is:



- ③ Semantic Analysis: ① The semantic analyzer uses the syntax tree and the information in the symbol table to check the source program for semantic consistency with the language definition.
- ② It also gathers & checks type information and saves it in either syntax tree or the symbol table, for subsequent use during intermediate-code generation.
- ③ The language specification may permit some type conversions called coercions.
- ④ Semantic Analysis of fig 1.:



Syntax tree after Semantic Analysis

- ⑤ Intermediate code generation: ① Many compilers generate an explicit low-level or machine-like intermediate representation, which we can think of as a program for an abstract machine.
- ② This intermediate representation should have two important properties: It should be easy to produce and it should be easy to translate into the target machine.

(c) We consider a intermediate form called three-address code, which consists of a sequence of assembly-like instructions with three operands per instruction. Each operand can act like a register.

(d) Intermediate code of syntax tree shown in Fig. 2:

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

$$t_2 = id_3 * t_1 \quad \text{--- (2)}$$

$$t_3 = id_2 + t_2$$

$$id_{1*} = t_3$$

(e) Code optimization: (a) The machine-independent code-optimization phase attempts to improve the intermediate code so that better target code will result.

(b) Usually better means faster, but other objectives may be desired, such as shorter code, or target code that consumes less power.

(c) Enter optimization of intermediate code in (d):

$$t_2 = id_3 * 60.0 \quad \text{--- (3)}$$

$$id_1 = id_2 + t_1$$

(f) Code Generation: The code generator takes as input an intermediate representation of the source program and maps it into the target language.

(b) If the target language is machine code, registers or memory locations are selected for each of the variables used by the program.

(c) Then, the intermediate instructions are translated into sequences of machine instructions that perform the same task.

(Q) Intermediate code

(D) Machine code for intermediate code in (Q)

LOF R₂, id₃ ← Immediate constant
 MULF R₂, R₂, #60.0
 LD F R₁, id₂
 ADDF R₁, R₁, R₂
 ST F id₁, R₁

2) Eliminate Left recursion in the following grammar (Remove Direct and Indirect recursion)

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid s d \mid \epsilon$$

Ans) Left recursion in A : $A \rightarrow Ac \Delta A \xrightarrow{*} Aa$

$$\therefore S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Aad \mid bd \mid \epsilon$$

$$\therefore A \rightarrow Ad_1 \mid Ad_2 \mid \dots \mid \beta_1 \mid \beta_2 \mid \dots$$

is equivalent to

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

$$A' \rightarrow d_1 A' \mid d_2 A' \mid d_3 A^3 \mid \dots \mid \epsilon$$

In our case, $d_1 = c$, $d_2 = ad$

$$\beta_1 = bd, \beta_2 = \epsilon$$

$$\therefore A \rightarrow Ad_1 \mid Ad_2 \mid \beta_1 \mid \beta_2$$

Can be re written as :

76-Afreen Sheikh

$$A \rightarrow B, A' \mid B, A'$$

$$AA' \rightarrow d, A' \mid d, A'$$

re substituting d, B & B'

$$S \rightarrow AaB$$

$$\therefore A \rightarrow bda' \mid A'$$

$$A' \rightarrow cA' \mid adA' \mid e$$

3. Consider the following grammar:

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

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

$F \rightarrow (E) \mid id$ to construct the Predictive Parser table.

Ans ① Eliminating Left Recursion

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

$$\therefore E \rightarrow TE'$$

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

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

$$\therefore T \rightarrow FT'$$

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

Non-recurs^r Left recursive Grammar:

$$E \rightarrow TE'$$

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

$$T \rightarrow FT' \mid F$$

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

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

② First and Follow:

$$\text{First}(E) = \{ \text{id}, C \}$$

$$\text{First}(E') = \{ +, e \}$$

$$\text{First}(T) = \{ C, \text{id} \}$$

$$\text{First}(T') = \{ *, e \}$$

$$\text{First}(F) = \{ C, \text{id} \}$$

$$\text{Follow}(E) = \{ (,) \}$$

$$\text{Follow}(E) = \{ (,), \$ \}$$

$$\text{Follow}(E') = \{) \}$$

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

$$\text{Follow}(T') = \{ +,), \$ \}$$

$$\text{Follow}(F) = \{ *, +,), \$ \}$$

Predictive Parser table:

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

$E \rightarrow TE' \rightarrow \text{First}(TE') = \text{First}(T) = \{c, id\}$ $E' \rightarrow E + TE' \rightarrow \text{First}(E + TE') = \{+\}$ $E' \rightarrow E \rightarrow \text{First}(E) = \{\} \Rightarrow \text{Follow}(E) = \{, \$\}$ $T' \rightarrow FT' \rightarrow \text{First}(F) = \{c, id\}$ $T' \rightarrow *FT' \rightarrow \text{First}(*) = \{* \}$ $F \rightarrow CE \rightarrow \text{First}(CE) = \{\} \Rightarrow \text{Follow}(T')$ $F \rightarrow \text{Follow}(T') = \{+, *\}$ $F \rightarrow CE \rightarrow \text{First}(C) = \{c\}$ $F \rightarrow id \rightarrow \text{First}(id) = \{id\}$

4. Differentiate Top-down and Bottom-up Parsing techniques. Explain Shift-reduce Parser with example.

Ans)

Top-down

Bottom-up

a) Top-down parsing can be viewed as the problem of constructing a parse tree for the input string, starting from the root and creating the nodes of the parse tree in pre-order.

A bottom-up parse corresponds to the construction of a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root.

b) Left most derivation is used.

Right most derivation is used.

c) It attempts to find the left most derivations for an input string.

It attempts to reduce the input string to start symbol of a grammar.

d) Its main decision is to select what production rule to use in order to construct the string.

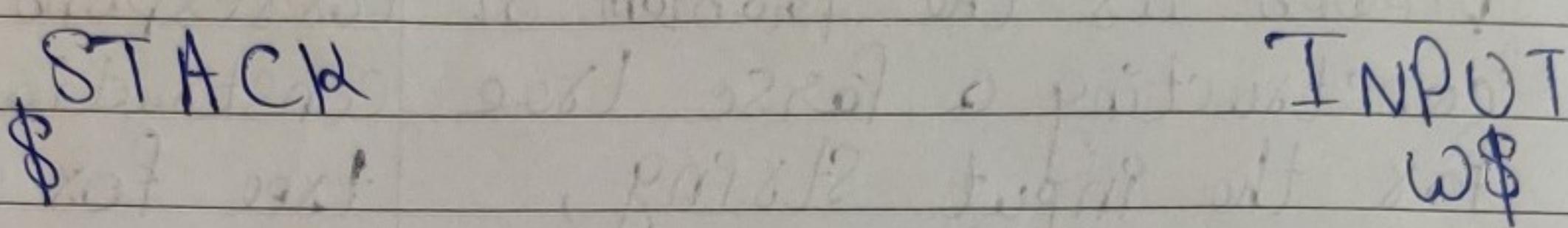
Its main decision is to select when to use a production rule to reduce the string to get the starting symbols.

Shift-reduce.

Shift-reduce Parser:

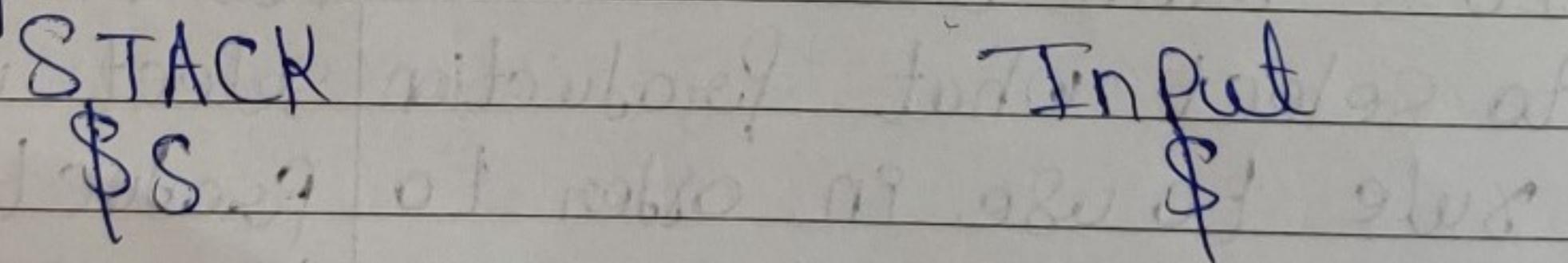
Shift-reducing Parsing is a form of bottom parsing in which a stack holds grammar symbols and an input buffer holds the rest of the string to be parsed.

We use \$ to mark the bottom of the stack and also the right end of the input. Conventionally, In bottom-up Parsing, we show the top of the stack on the right. Initially, the stack is empty, and the string ω is on the input, as follows:



During a left-to-right scan of the input string, the parser shifts zero or more input symbols onto the stack, until it is ready to reduce a string β of grammar symbols on top of the stack. It then reduces β to the head of the appropriate production.

The parser repeats this cycle until it has detected an error or until the stack contains the start symbol & the input is empty.

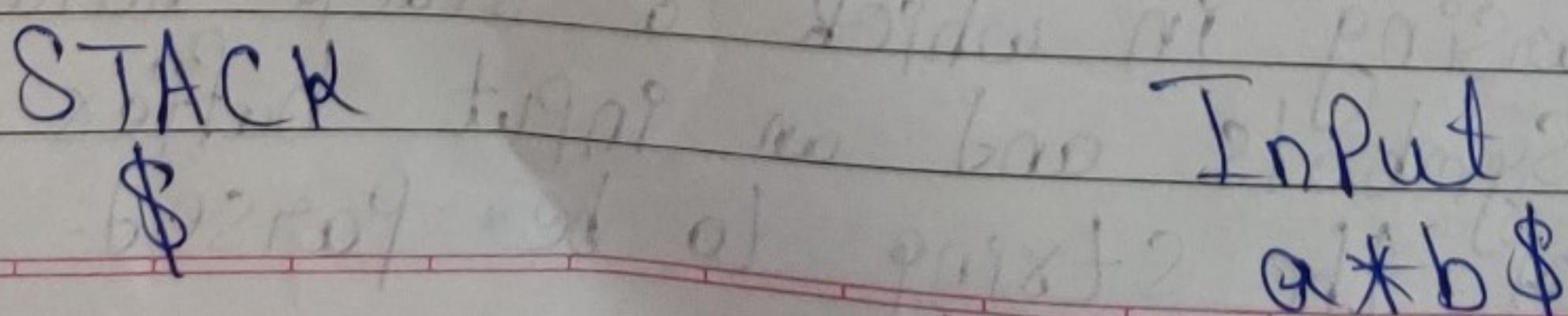


E.g. Consider, Grammar:

$$E \rightarrow ET \mid T, \quad T \rightarrow T * F \mid F, \quad F \rightarrow (E) \mid id$$

and input : $a * b$

Initial Configuration



Scan the Postfix string from the left,

STACK	Input	Action
a	a * b \$	Shift
F	* b \$	Reduce
T	* b \$	Reduce
T *	* b \$	Shift
T * b	* b \$	Shift
T * F	* b \$	Reduce
E	\$	Reduce
		Halt, Success- ful Parsing completion

Q6) Construct SLR Parsing table for following grammar.

$$S \rightarrow (S)S$$

$$S \rightarrow \epsilon$$

Solⁿ: Augmented:

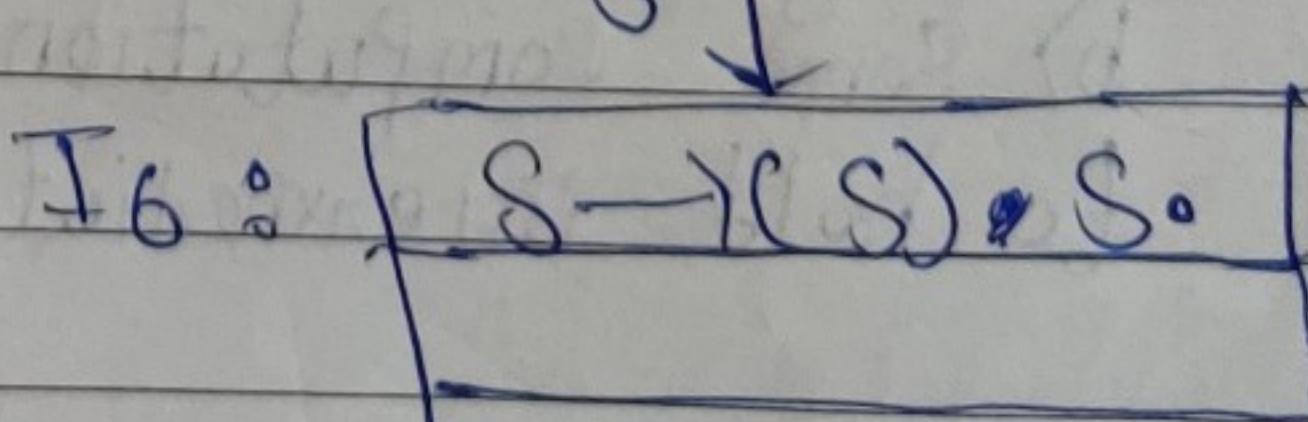
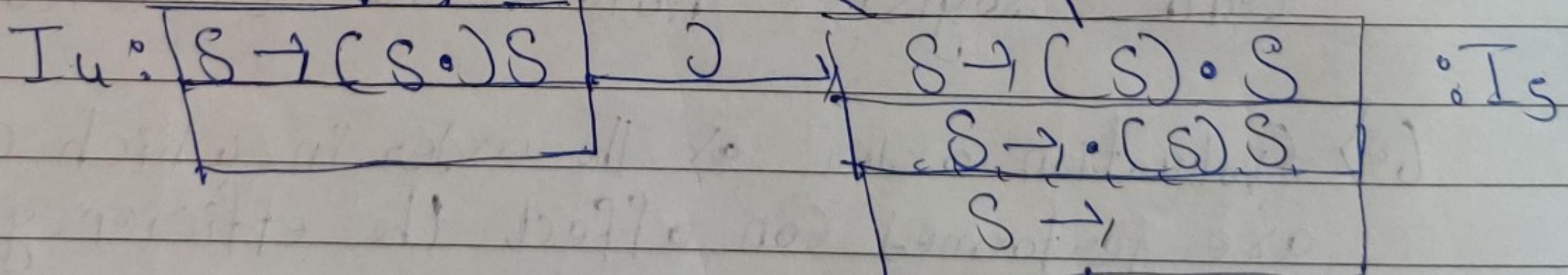
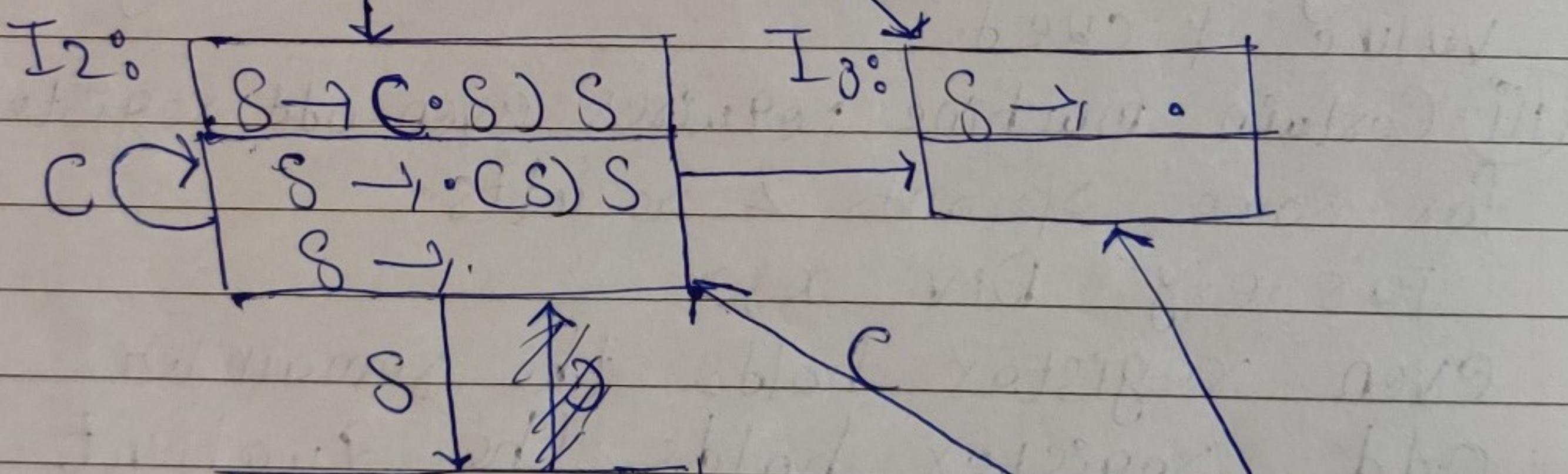
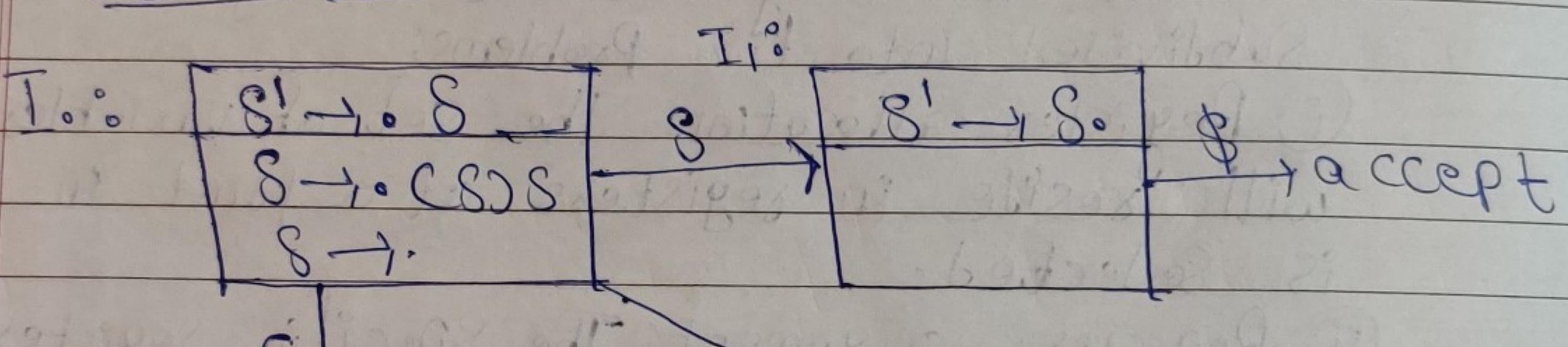
$$S' \rightarrow S \quad (1)$$

$$S \rightarrow (S)S \quad (2)$$

$$S \rightarrow \epsilon \quad (3)$$

$$\text{Follow}(S') = \{\$\} ; \text{Follow}(S) = \{\$,)\}$$

Items Sets:



State	Actions			GOTO
	C	S	\$	
0	S2	S3		8
1			acc	1
2	S2	S3		4
3		x3	x3	
4		S5		
5	S2	S2	S3	6
6		g12	g12	

Q5) For the given grammar below construct operator precedence relations matrix, assuming *, + are binary operators and id as terminal symbol & E as non-terminal symbol.

$$E \rightarrow E+E ; E \rightarrow E * E ; E \rightarrow id.$$

Sol: ① Construct operator precedence relation matrix

$$\Sigma = \{ id, +, *, \$ \}$$

	id	+	*	\$
id	< . >	• >	• >	• >
+	< . >	• >	< . >	• >
*	< . >	• >	• >	• >
\$	< . >	• >	< . >	

Apply operator precedence parsing algorithm to obtain skeletal syntax tree for the statement: id + id * id

Input string \$ id + id * id \$

\$ \Rightarrow \$ < . id . > + < . id . > * < . id . > \$

\$ \Rightarrow \$ E + E * E \$

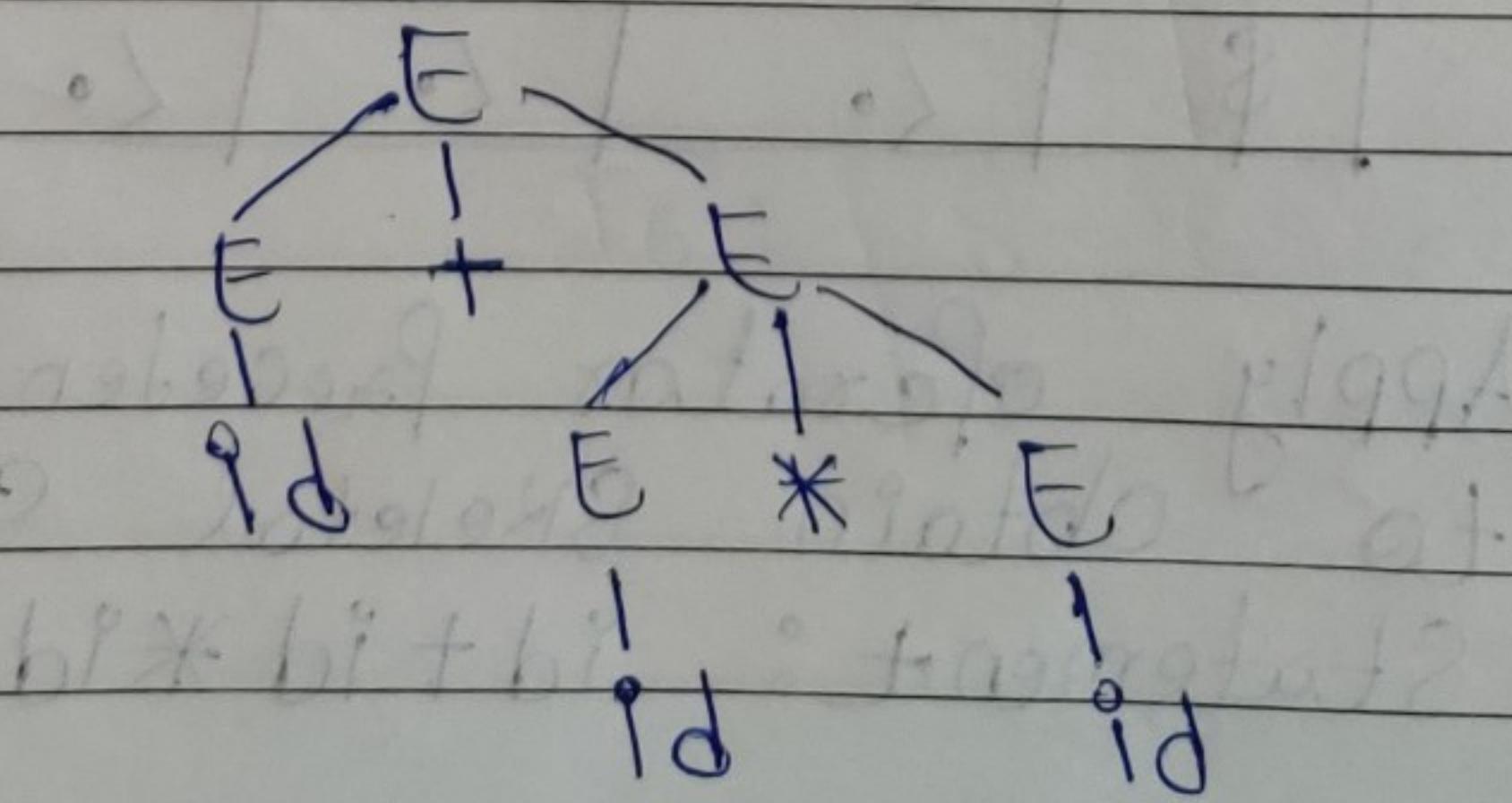
\$ \Rightarrow \$ + * \$ \Rightarrow \$ < . + < . * . > \$

\$ \Rightarrow \$ + \$ < . + \$ \Rightarrow \$ < . + . > \$ \Rightarrow \$ \$ Accepted

using stack: {a < . b ; a ÷ b} \Rightarrow Shift ; a \Rightarrow b \Rightarrow Reduce

Stack	Input String	a op b	Action
id	id + id * id \$	\$ < . id	Shift
E	+ id * id \$	id . > +	Reduce
E+	+ id * id \$	\$ < . +	Shift
E+id	* id \$	id . > *	Reduce
E+E	* id \$	+ < . *	Shift
E+E*	id \$	* < . id	Shift
E+E*id	\$	id . > \$	Reduce
E+E*E	\$	* . > \$	Reduce
E+E	\$	+ . > \$	Reduce
E	\$	\$ + \$	Accept

Tree:



Q.7) Write a short note on Syntax-Directed definitions.

- Ans) ① A Syntax-directed definition (SDD) is a Context-free grammar together with attributes and rules.
- ② Attributes are associated with grammar symbols & rules are associated with the productions. If X is a symbol and a is one of its attributes, then we write $X.a$ to denote the value of ' a ' at a particular parse-tree node labeled X .
- ③ Attributes may be of any kind: numbers, types, table, references, or strings, for instance.
- ④ There are two kinds of nonterm attributes for nonterminals:
- ⑤ A synthesized attribute for a nonterminal A at a parse-tree Node N is defined by semantic rule associated with the production at N .

- (b) An inherited attribute for a non-terminal B at a parse-tree Node N is defined by a semantic rule associated with the production at the parent of N.
- (c) An SDD that involves only synthesized attributes is called S-attributed.
- (d) In an S-attributed SDD, value of head is computed with the production body i.e. parent value is computed using children value.
- (e) The "L-attributed translations" (Left-to-Right) which encompasses virtually all translations that can be performed during parsing is mostly applied.
- (f) e.g. of infix-to-postfix translation:

Production	Semantic Rule
$E \rightarrow E_1 + T$	$E\text{-code} = E_1\text{-code} T\text{-code} +$

Code is an attribute & || used for concatenation. E & E_1 are same non-terminal but subscript is used for distinguish between instances.

- (g) Rules can be written as follow:

$$E \rightarrow E_1 + T \{ \text{print}'+' \}$$

The code fragment appears between '{ }' which is called as semantic action

- (h) The position of action in a production body determines the order in which the action is executed.

8) Explain various form of intermediate code used by compiler

(Ans) The various intermediate code forms are:

- a) Polish notation
- b) Syntax tree
- c) Three address codes: Triples and Quadtuples.
- d) Quadtuples

a) Polish notation: The postfix notation for an expression E can be defined inductively as follows:

1. If E is a variable or constant, then the postfix notation for E is E itself.
2. If E is an expression of the form $E_1 \text{ op } E_2$, where op is any binary operator, then the postfix notation for E is $E_1 E_2 \text{ op}$, where E_1 and E_2 are the postfix notation for E_1 and E_2 .
3. If E is a parenthesized expression of the form (E_1) then the postfix notation for E is the same as the postfix notation for E_1 .

e.g.: $(9-5) + 2$ Postfix: $E = (E_1) + E_2$

$$(E_1) = E_1 = 9 - 5 = 95 -$$

$$E_2 = 2$$

$$\therefore E = (E_1) + E_2 = (E_1) E_2 + = 95 - 2 +$$

To evaluate postfix expression we repeatedly scan the postfix string from the left,

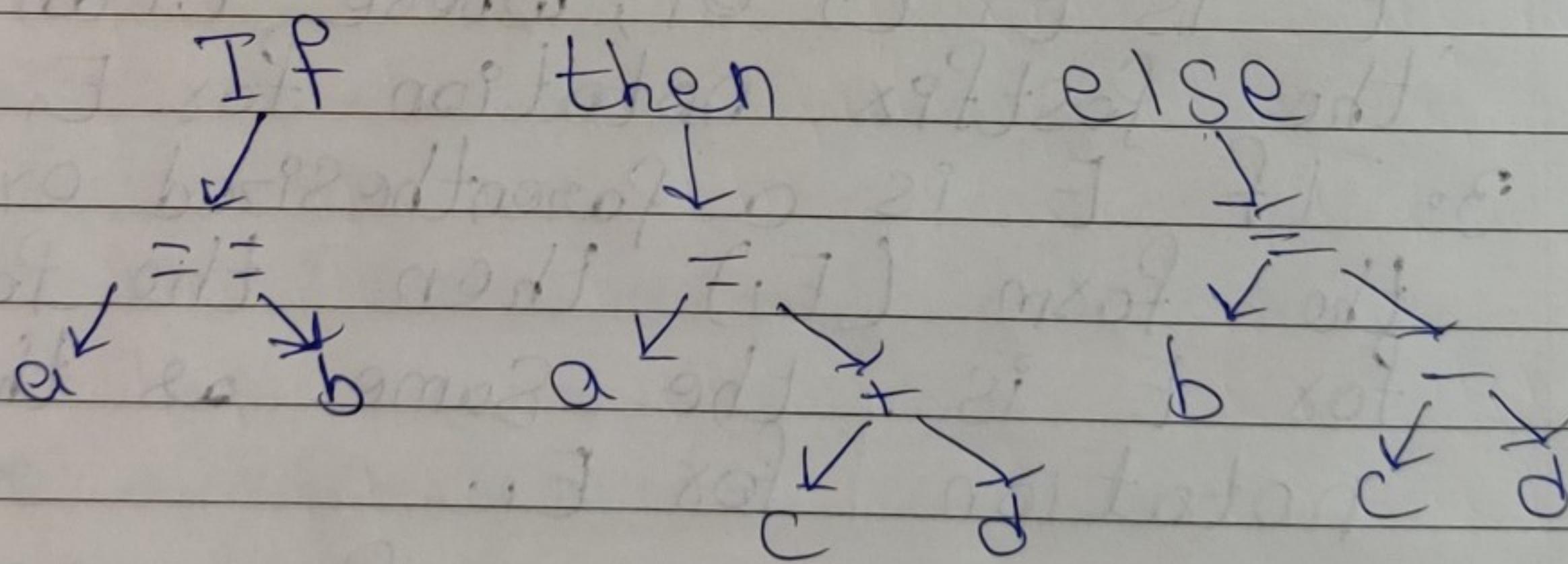
until we find an operator. Then, look to the left we find operands. Then, look to the left for the proper number of operands, & replace them by the result. Then repeat the process, continuing to the right & searching for another operator.

e.g.: $95 - 2 + \Rightarrow 9 - 5 \Rightarrow 4 2 + \Rightarrow 4 + 2 \Rightarrow 6$

b) Syntax tree: The parse tree itself is a useful intermediate language representation for a source program, especially in optimizing compilers where the intermediate code needs to extensively restructure.

We use an abstract syntax tree a tree in which each leaf represents an operand & each interior node an operator.

e.g.: If $a = b$ then $a = c + d$ else $b = c - d$



c) Three-Address Code:

- ① In three address code, there is atmost one operator on the right side of an instruction. e.g.: $a = b + c$.

- ② No more than 2 variables appear on right side & 1 on left of assignment operators.

- ③ It eliminates the ambiguity by evaluating only one operation at a time.

i) Quadruples: A quadruple has four fields: OP, arg1, arg2 & result. The evaluation is result of arg1 & arg2 if OP is stored in result.

ii) Triples: A triple has three fields: OP, arg1 & arg2. Triples store result of arg1 & arg2 in its position rather than explicit temporary name.

e.g. $a = b * c + b * c$

$$\begin{array}{c}
 a \\
 + \\
 \diagup \quad \diagdown \\
 b \quad c \\
 * \\
 \diagup \quad \diagdown \\
 b \quad c
 \end{array} \Rightarrow a = b * c + b * c$$

Triple address	Instruction	OP	arg1	arg2
35	(0)	minus	c	(0)
36	(1)	*	b	(0)
37	(2)	minus	c	(0)
38	(3)	*	b	(0)
39	(4)	+	(1)	(0)
40	(5)	=	a	(0)

Quadruple address	OP	arg1	arg2	result	OP	arg1	arg2	result
0	minus	c		t1	4	+	t2	t4
1	*	b	t1	t2	5	=	t5	a
2	min	c		t3				
3	*	b	t3	t4				

Q.) Explain the principal sources of optimization.

Ans) ① There are number of ways in which a compiler can improve a program without changing the function it computes.

Function preserving transformations example:

a) Common Sub expressions elimination.

b) Copy propagation;

c) Dead-code Elimination

d) Constant folding

The other transformations come up primarily when global optimizations are performed.

② Common Sub expressions elimination:

a) An occurrence of an expression E is called a common sub-expression if E was previously computed, and the values in E haven't changed since the previous computation.

We can avoid recomputing the expression if we can use the previously computed value.

b) For example:

$$t_1 := 4 * i$$

$$t_2 := a[t_1]$$

$$t_3 := 4 * i \leftarrow \text{common}$$

$$t_4 := n$$

$$t_5 := b[t_3] + t_4$$

Optimized

$$t_1 := 4 * i$$

$$t_2 := a[t_1]$$

$$t_3 := n$$

$$t_4 := b[t_1] + t_3$$

③ Copy Propagation: a) Copy Propagation means use of one variable instead of another. This may not appear to be an improvement, but it gives an opportunity to eliminate x.

For e.g.: $x = p_i$; $A = x * x * x$;

optimized : $A := \text{Pi} * x * x$
 ~~x~~ eliminated.

- ④ Dead-code elimination: ② A variable is live at a point in a program if its value can be used subsequently; otherwise, it is dead at that point.

⑤ Example:

```
int compute_dead(int c){  
    int a := c * 2;  
    int b := a * c; → dead code  
    return a;  
}
```

in above code b is assigned but never used it can eliminated from there.

- ⑤ Constant folding: ② Deducing at compile time that the value of an expression is a constant and using the constant instead is known as Constant folding.

⑥ e.g.: $a := 3.14157 / 2$ optimized, $a := 1.570$

- ⑦ Loop optimizations: ② In loops, especially in the inner loops, programs tend to spend the bulk of their time. The running time of a program may be improved if the number of expressions in an inner loop is decreased, even if we increase the amount of code outside the loop.

- ⑧ These techniques are important for loop optimization.

i) Code motion

ii) Induction-variable elimination.

iii) Reduction in Strength.

- ⑦ Code Motion: ① An important modification that decreases the amount of code in a loop is code motion.
 ② This transformation takes an expression that yields the same result independent of the number of times a loop is executed (loop invariant) and places the expression in the outer loop if exists or "before the loop" if entry point exists.

e.g.:

~~while ($i < limit - 2$) {
 # not changing → Stmt;
 # limit } }~~

~~$t = limit - 2$~~

~~while ($i < t$) {
 Stmt;
 }~~

$limit - 2$ will remain same throughout the loop instead of calculating it each time we move it out & calculated once.

- ⑧ Induction Variables: ① Loops are usually processed inside out. For example: $j = k;$

② for e.g.: $j = k$ if $V = \text{some value}$

~~do {~~

~~$j = j - 1$~~

~~$t_1 = 4 * j$~~

~~$t_2 = a[t_1]$~~

~~} while $t_2 > V$~~

In above example values of j & t_1 remain in lock-step; every time $j = j - 1 \Rightarrow t_1$ decreases by 4. i.e. $t_1 = 4 * j$, $t_1 = 4(j - 1) = 4j - 4$. j , only implies previous value replaced with new one. Such identifiers are called induction.

① When there are two or more induction variables in a loop, it may be possible to get rid of all but one, by the process of induction Variable elimination.

For e.g. $j = k; t_1 = 4 * j; v = \text{some value}$

$\text{do } \{$

$j = j - 1; t_1 = t_1 + 4$

$t_2 = a[t_1]$

while $t_2 > v$

Here, we get ride of multiplication by substituting subtraction which is much cheaper & initialize t_1 at entry point of loop because it had no origin before assigning it value.

② Reduction in Strength: Reduces expensive operations by equivalent cheaper ones on the target machine.

b) For e.g.: $x * x$ is cheaper to implement than x^2 which required a call to exponentiation routine.

x^2 \Rightarrow Optimized $\Rightarrow x * x$

Q) Explain the various issues in the design of code generation.

Ans) ① Input to code generator: @ The input to the code generation consists of, the intermediate representation of the source produced by front end, together with the information in the symbol table to determine run-time address of the data objects denoted by the names in the intermediate representation.

② Intermediate representation can be:

- i) Postfix notation
- ii) Three address code
- iii) Virtual Machine (stack machine code)
- iv) Syntax tree.

③ Target program: @ The output of the code generator is the target program. The output may be:

- i) Absolute machine language: It can be placed in a fixed memory location and can be executed immediately.
- ii) Relocatable machine language: It allows subprograms to be compiled separately.
- iii) Assembly language: Code generation is made easier.

④ Memory management: @ Names in the source program are mapped to addresses of data objects in run-time memory by front-end & code generator.

b) It makes use of symbol table, that is, a name in a three-address statement refers to a symbol-table entry for the name.

c) Labels in three-address statements have to be converted to addresses for instructions.

- 4.) Instruction Selection: a) The instructions of target machine should be complete & uniform.
 b) Instruction speeds & machine idioms are important factors when efficiency of target program is considered.
 c) The quality of the generated code is determined by its speed and size.

- 5.) Register allocation: a) Instructions involving register operands are shorter & faster than those involving operands in memory. The use of registers ~~are~~ is subdivided into two problems:

- i) Register allocation: The set of variables that will reside in registers at a point in the program is selected.
- ii) Register assignment: The specific register that a value picked.

iii) Certain machine requires even-odd register pairs for some operands & results.

For e.g.: $\text{DIV } x, y$

Even register holds the remainder.

Odd register holds the quotient.

- 6.) Evaluation order: a) The order in which computations are performed can affect the efficiency of the target code.

- b) Some computation orders require fewer registers to hold intermediate results than others.