

# UNIT - 5

## SYNTAX - DIRECTED TRANSLATION

### CONTENTS

- Syntax directed definitions
- Evaluation orders for SDD's
- Application of SDT
- SDT schemes.

### Syntax Directed Definitions:

A syntax directed definition in a context free grammars with attributes & rules. Attributes are associated with grammar symbols & rules with production. If 'x' is a symbol, 'a' is one of attributes then we write  $x.a$  to denote value of  $a$  at a particular part of tree. Note that attributes may be of many kinds: numbers, types, tables, references, strings, etc...

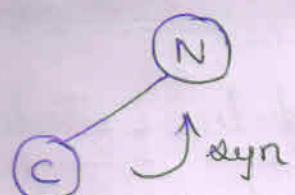
- \* 2 types of attributes:
  - ▷ Synthesized attr
  - ▷ Inherited attr

① Synthesized attr: A synth attr for a nonterminal A at a parse tree node N is defined by a semantic rule associated with the production at N. Note that the production must have A as its head. A synthesized attr at node N is defined only in terms of attribute values at the children of N or at N itself.

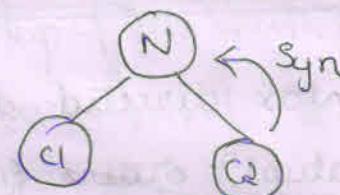
② Inherited attr An inherited attr for a nonterminal B at a parse tree node N is defined by a semantic rule associated with the production at the parent of N. Note that the production must have B as a symbol in its body. An inherited attr at node N is defined only

in terms of attr values at N's parent, N itself & N's sibling

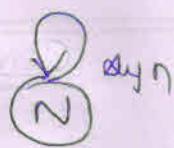
NOTE ① Synthesized  $\begin{cases} N \rightarrow \text{Node under consideration} \\ C \rightarrow \text{child} \end{cases}$



Case(i) Single child



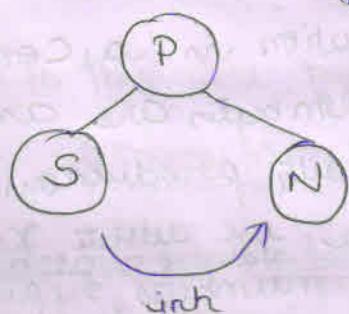
Case(ii) Rightmost child



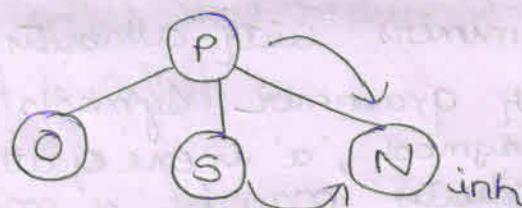
Case(iii) No child  
⇒ itself.

② Inherited

$\begin{cases} P \rightarrow \text{parent} \\ N \rightarrow \text{node under consideration} \\ S \rightarrow \text{Sibling} \\ O \rightarrow \text{operator} \end{cases}$



Case(i) sibling



Case(ii) : inherited from both parent of sibling.

③ Terminals can have Synthesized attributes but not inherited attributes.

\* Attr of terminals have lexical value that are supplied by lexical analysis.

\* Types of SAD : ① S Attributed SAD  
② L Attributed SAD

① S-Attributed SAD

- \* As SAD that involves only synthesized attr other it is called S-attributed SAD
- \* In s-attributed SAD, each rule consists

an attribute for the terminal at the head of a production from attribute taken from body of production.

\* S-attributed SOD can be implemented naturally in conjunction with an LR parser / bottom up parser.

### \* Annotated parse tree

A parse tree showing the value(s) of attribute(s) is called an annotated parse tree.

\* It is used in bottom-up parser.

\* Order of evaluation is postorder traversal.

\* Example of S-attributed SOD.

Production	Semantic rules
1) $\mathcal{E} \rightarrow E_n$	$E.\text{val} = E_n.\text{val}$
2) $E \rightarrow EI + P$	$E.\text{val} = EI.\text{val} + P.\text{val}$
3) $E \rightarrow T$	$E.\text{val} = T.\text{val}$
4) $T \rightarrow T, *F$	$T.\text{val} = T.\text{val} * F.\text{val}$
5) $T \rightarrow F$	$T.\text{val} = F.\text{val}$
6) $F \rightarrow (E)$	$F.\text{val} = E.\text{val}$
7) $F \rightarrow \text{digit}$	$F.\text{val} = \text{digit}. \text{lexval}$

### ② L-Attributed SOD

\* Example of mixed attributes / L-attributed SOD

Production	Semantic rules
1) $T \rightarrow FT'$	$T'.\text{inh} = F.\text{val}$
	$T.\text{val} = T'.\text{syn}$
2) $T' \rightarrow *FT''$	$T''.\text{inh} = T'.\text{inh} * F.\text{val}$
	$T'.\text{syn} = T''.\text{syn}$
3) $T'' \rightarrow E$	$T'.\text{syn} = T''.\text{inh}$
4) $F \rightarrow \text{digit}$	$F.\text{val} = \text{digit}. \text{lexval}$

- \* A SAD which has both synthesized & inherited attributes is called as L-attributed SAD.
- \* It is used in top down parsing.
- \* Order of evaluation is topological sorting.

### Evaluating Orders For SAD's

- \* A dependency graph is used to determine the order of computation of attributes.
- \* While an annotated parse trees shows the values of attributes, a dependency graph helps us to determine how those values can be computed.

### DEPENDENCY GRAPHS

A dependency graph predicts the flow of information among the attribute instances in a particular parse tree: An edge from one attribute instance to another means that the value of first is needed to compute the second.

- 1) For each parse tree node, say node  $X$ , the dependency graph has a node for each attribute associated with  $X$ .
- 2) If a semantic rule (defines) associated with a product ' $p$ ' defines the value of synthesized attribute  $A, b$  in terms of value of  $X, c$  then the dependency graph has an edge from  $X, c$  to  $A, b$ .
- 3) If a semantic rule associated with a product ' $p$ ' defines the value of inherited attribute  $B, c$  in terms of value of  $X, a$  then the

dependency graph has edge from X.c to B.c

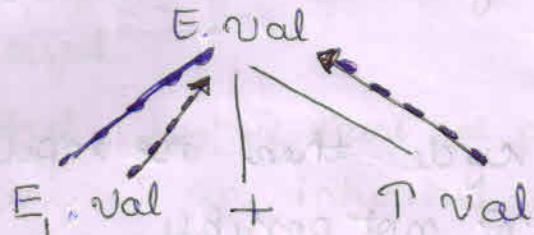
Eg1: production

$$E \rightarrow E_1 + T$$

Semantic Rule

$$E.\text{Val} = E_1.\text{Val} + T.\text{Val}$$

fig: E.Val is synthesized from E<sub>1</sub>.Val & T.Val



Eg2: Production

$$T \rightarrow FT'$$

$$T' \rightarrow *FT'$$

$$T' \rightarrow E$$

$$F \rightarrow \text{digit}$$

Semantic Rule.

$$T'.\text{inh} = F.\text{val}$$

$$T.\text{Val} = T'.\text{syn}$$

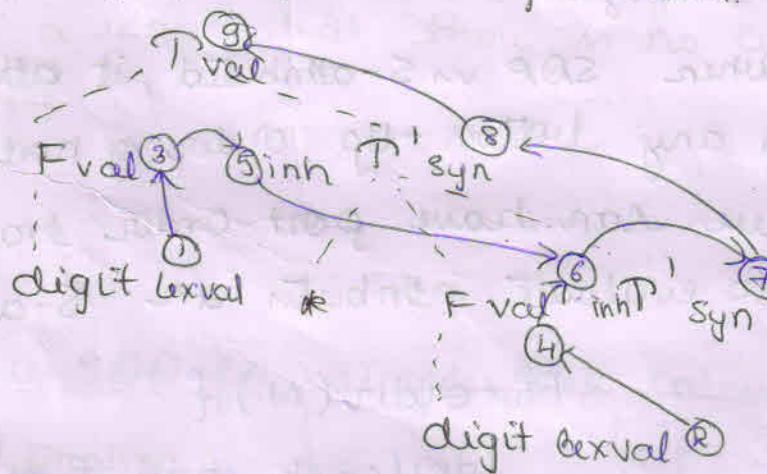
$$T'.\text{inh} = T'.\text{inh} * F.\text{val}$$

$$T'.\text{syn} = T'.\text{syn}$$

$$T'.\text{syn} = T'.\text{inh}$$

$$F.\text{val} = \text{digit.lexval}$$

fig: Dependency graph for above production.



Ordering the evaluation of attributes

- \* If the dependency graph has an edge from node M to N, then the attribute corresponding to M must be evaluated before attribute of N.

- \* Thus the only allowable orders of evaluation are those sequences of nodes  $N_1, N_2, \dots, N_k$  if there is an edge of the dependency graph from  $N_i$  to  $N_j$  then  $i < j$ .
- \* Such an ordering is called a topological sorting of a graph
- \* If there is any cycle then no topological sort, i.e., evaluation of SAD is not possible.
- \* Eg: For dependency graph for (Eq 2) in previous page  
topological sorting: 1, R, 3, 4, 5, 6, 7, 8, 9  
or  
1, 3, 5, R, 4, 6, 7, 8, 9

### S-Attributed Definitions

- An SAD is S-attributed if every attribute is synthesized.
- When SAD is S-attributed, it attributes evaluated in any bottom-up order of nodes of parse trees.
- we can have post-order traversal of parse tree to evaluate attributes in S-attributed definitn.

Postorder(N) {

for(each child C of N, from the left)  
postorder(c);

evaluate the attributes associated with  
node N;

}

- S-Attributed definitions can be implemented during bottom up parsing without the need to explicitly create parse trees.

## - Attributed Definitions

- \* A SOD is k-attributed if the edges in dependency graph goes from Left to Right but not from Right to Left.
- \* More precisely, each attribute must be either
  - synthesized.
  - Inherited , but if there is a production  $A \rightarrow X_1 X_2 \dots X_n$  & there is an inherited attribute  $X_i$ . a Competed by a rule associated with this product then the rule may only use:
    - Inherited attributes associated with the head A .
    - Either inherited or synthesized attr associated with the occurrence of symbols  $X_1 X_2 \dots X_{i-1}$  located to the left of  $X_i$ .
    - Inherited/Synthesized attr associated with this occurrence of  $X_i$  itself but only in such a way that there is no cycle in the graph.

## PROBLEMS

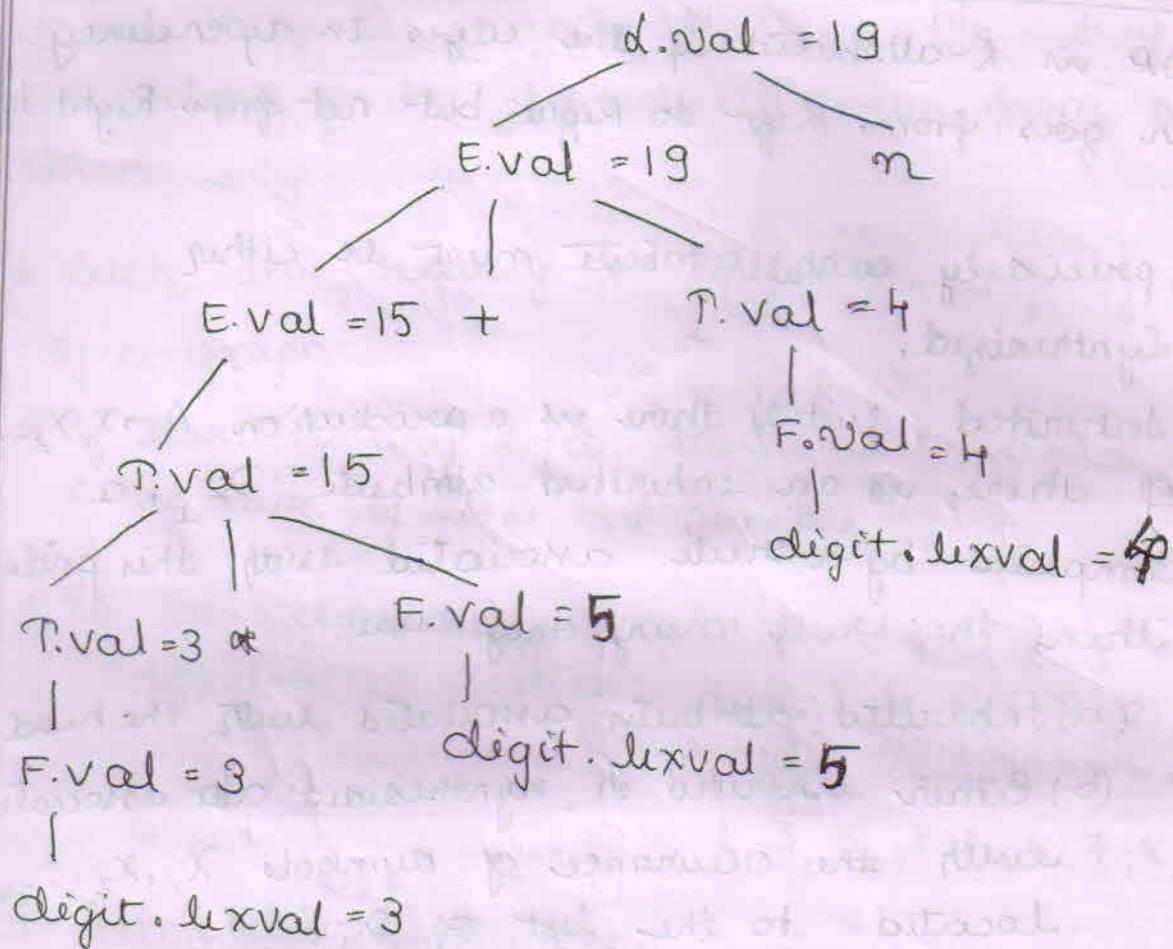
- 5.1. Write a SOD for simple desk calculator

Sol: SOD definition

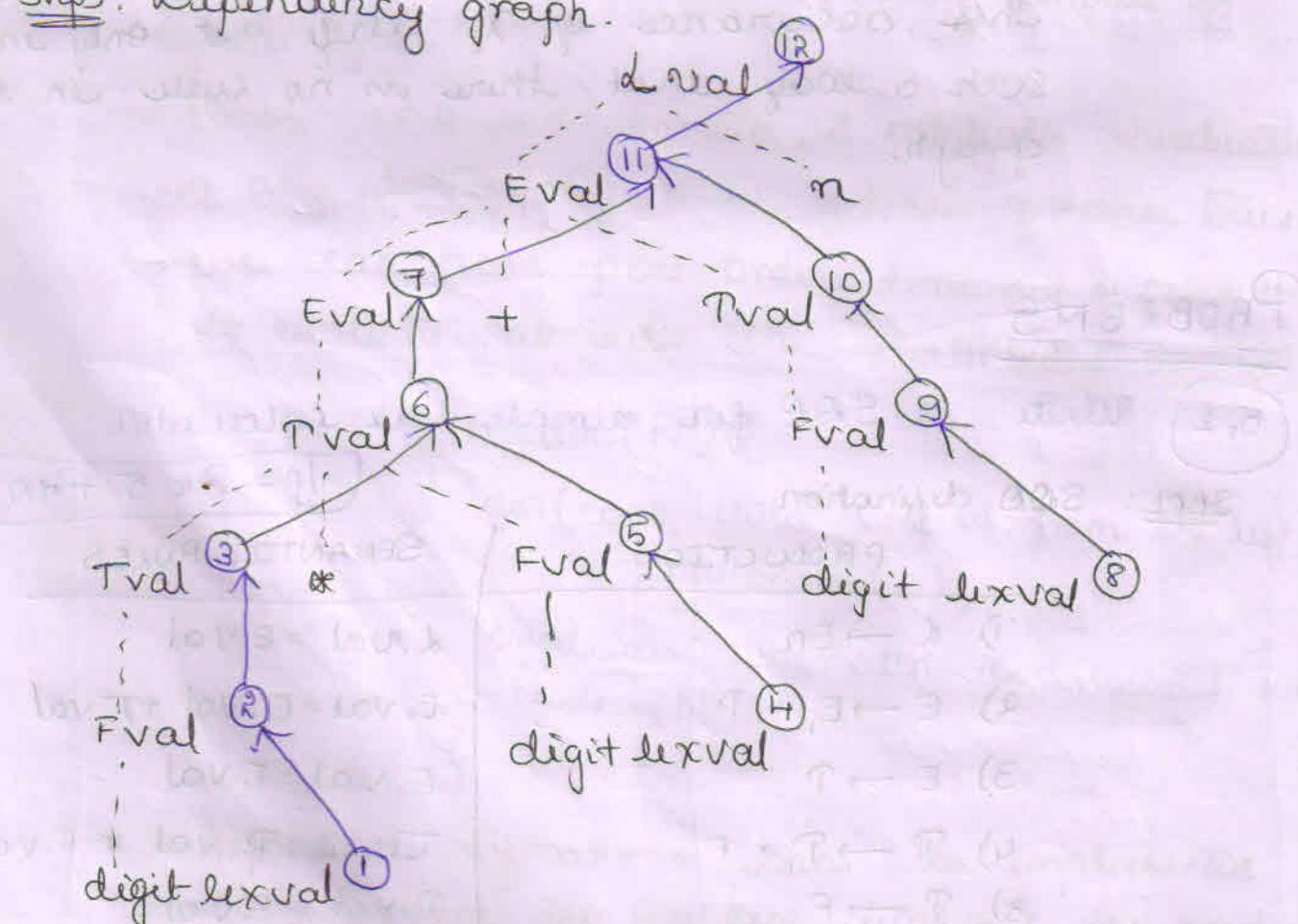
$$i/p = 3 * 5 + 4 n \}$$

PRODUCTION	SEMANTIC RULES
1) $K \rightarrow E_n$	$L.val = E.val$
2) $E \rightarrow E_1 + P$	$E.val = E_1.val + T.val$
3) $E \rightarrow P$	$E.val = T.val$
4) $T \rightarrow T_1 * F$	$T.val = T_1.val * F.val$
5) $T \rightarrow F$	$T.val = F.val$
6) $F \rightarrow ( E )$	$F.val = E.val$
7) $F \rightarrow \text{digit}$	$E.val = \text{digit.lexval}$

### Step 2: Annotated Parse tree



### Step 3: Dependency graph



Step 4: Topological Order: ① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨ ⑩ ⑪ ⑫

Write the SOD & construct annotated parse tree, dependency graph.

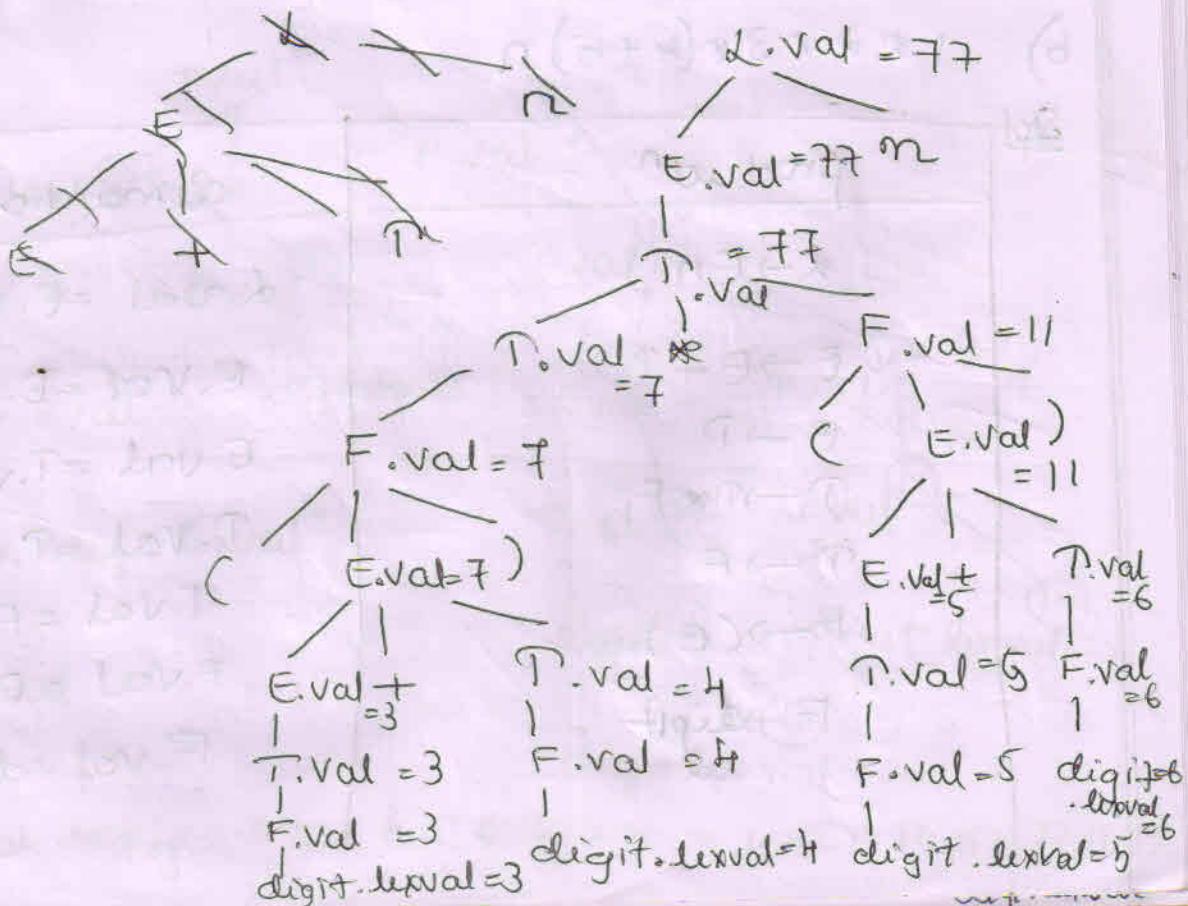
- a)  $(3+4)*(5+6)n$
- b)  $1*2*3*(4+5)n$
- c)  $(9+8*(7+6)+5)*4n$

d)  $(3+4)*(5+6)n$

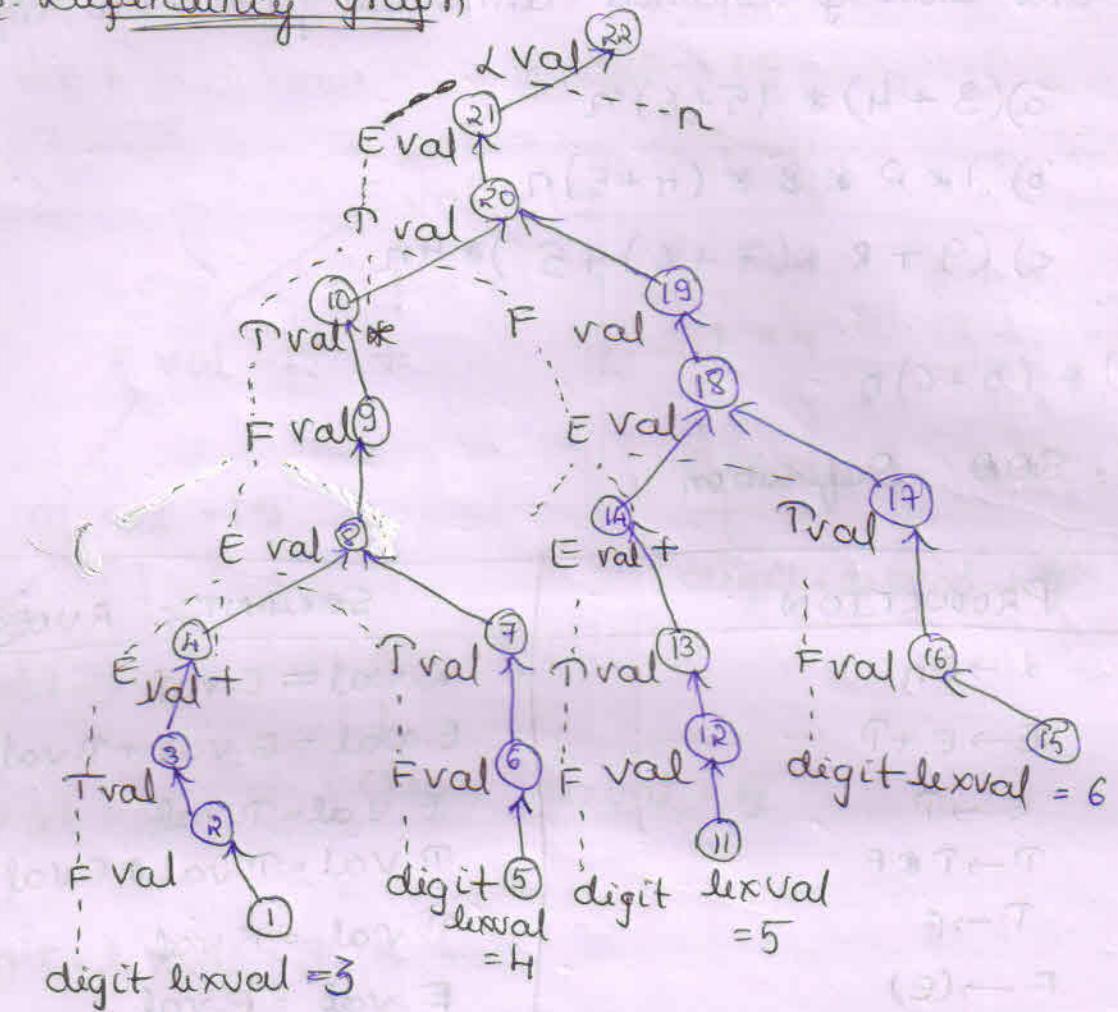
Step 1: SOD Definition

PRODUCTION	SEMANTIC RULES
$L \rightarrow E_n$	$L.val = E.val$
$E \rightarrow E + T$	$E.val = E.val + T.val$
$E \rightarrow T$	$E.val = T.val$
$T \rightarrow T * F$	$T.val = T.val * F.val$
$T \rightarrow F$	$T.val = F.val$
$F \rightarrow (E)$	$F.val = E.val$
$F \rightarrow \text{digit}$	$F.val = \text{digit.lexval}$

Step 2: Annotated Parse Tree



### Step 3: Dependency Graph



### Step 4: Topological sorting

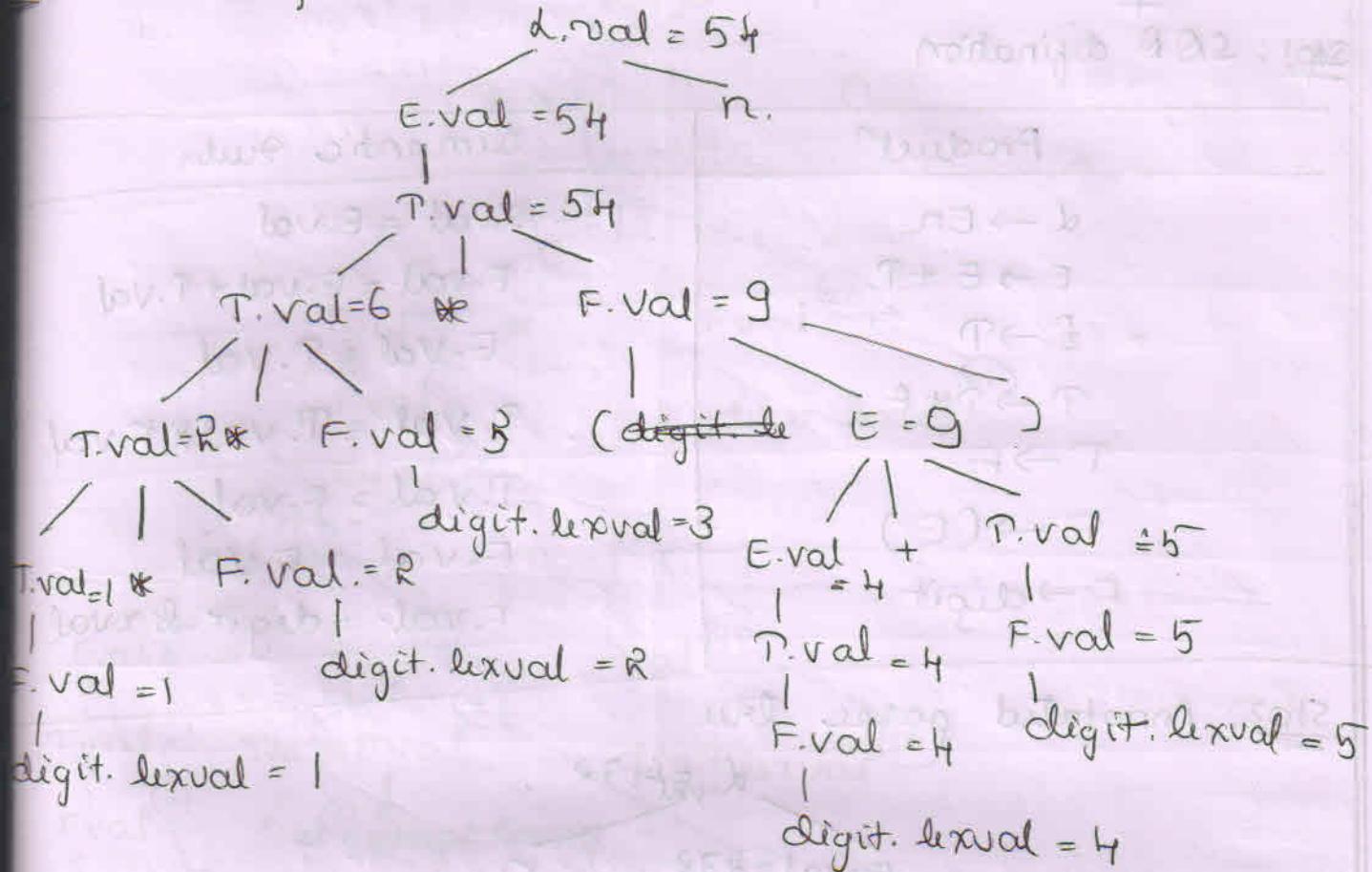
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15  
 16 17 18 19 20 21 22.

b)  $1 * 2 * 3 * (4 + 5) n$

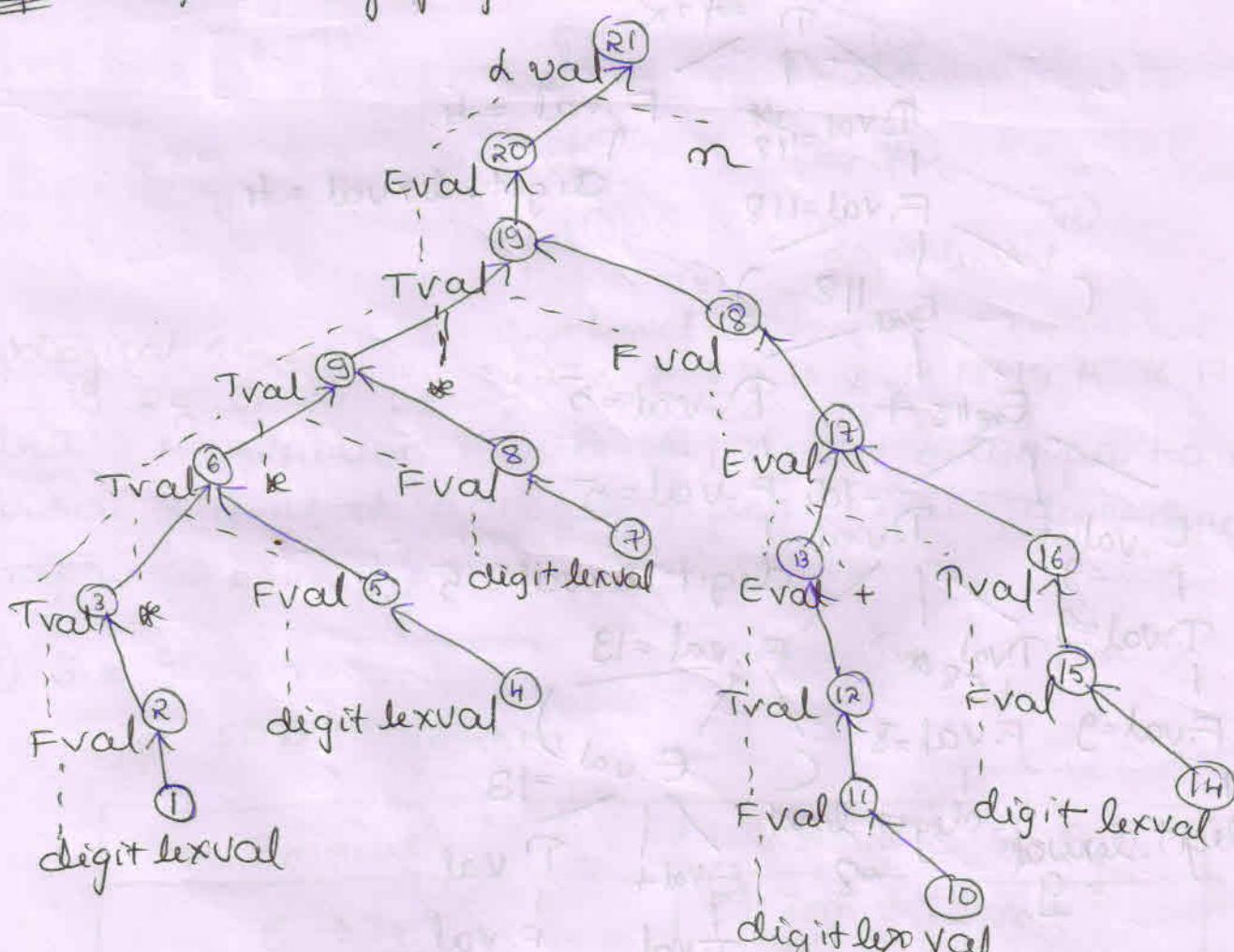
### Step 5:

product^n	Semantic rules
$A \rightarrow E^n$	$L\text{-val} = E\text{-val}$
$E \rightarrow E + T$	$E\text{-val} = E\text{-val} + T\text{-val}$
$E \rightarrow T$	$E\text{-val} = T\text{-val}$
$T \rightarrow T * F$	$T\text{-Val} = T\text{-val} * F\text{-val}$
$T \rightarrow F$	$T\text{-Val} = F\text{-val}$
$F \rightarrow (E)$	$F\text{-val} = E\text{-val}$
$P \rightarrow \text{digit}$	$P\text{-val} = \text{digit lexval}$

## 2. Annotated parse tree



### Step 3: Dependency graph



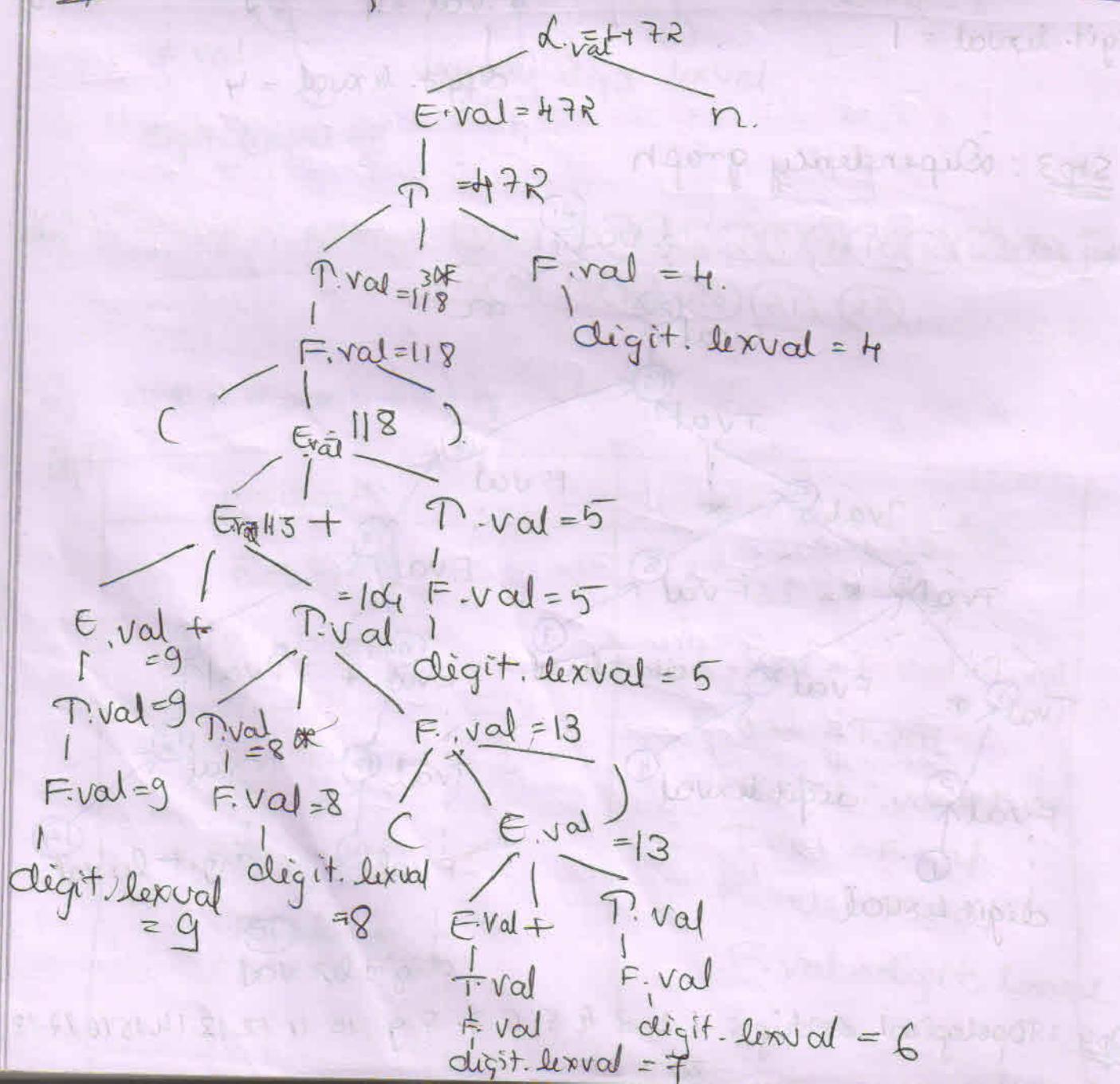
Sig: Topological sorting: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19  
20 21

$$c) (9 + 8 * (7 + 6) + 5) \text{ & L.H.S}$$

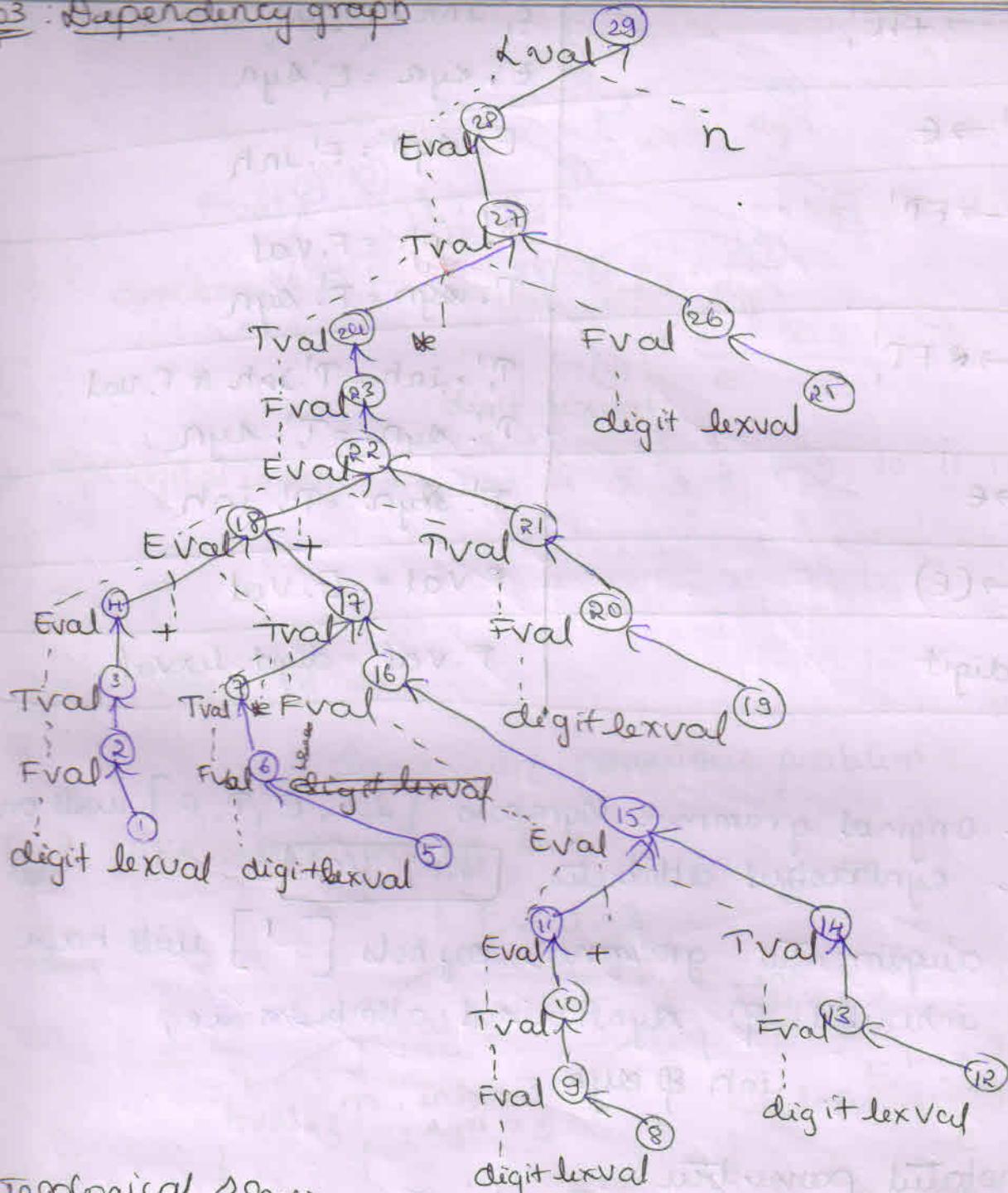
Step 1: SOD definition

Product <sup>n</sup>	Semantic rules
$L \rightarrow E_n$	$L.val = E.val$
$E \rightarrow E + T$	$E.val = E.val + T.val$
$E \rightarrow T$	$E.val = T.val$
$T \rightarrow T * F$	$T.val = T.val * F.val$
$T \rightarrow F$	$T.val = F.val$
$F \rightarrow (E)$	$F.val = E.val$
$F \rightarrow \text{digit}$	$F.val = \text{digit.lexval}$

Step 2: Annotated parse tree



### Step 3: Dependency graph



Topological sorting: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

Write L-attributed SAD (8) Write a SAD for top down parser & construct annotated parse tree, dependency graph for given I/p.

① 3 \* 5

Step 4: SAD Definition

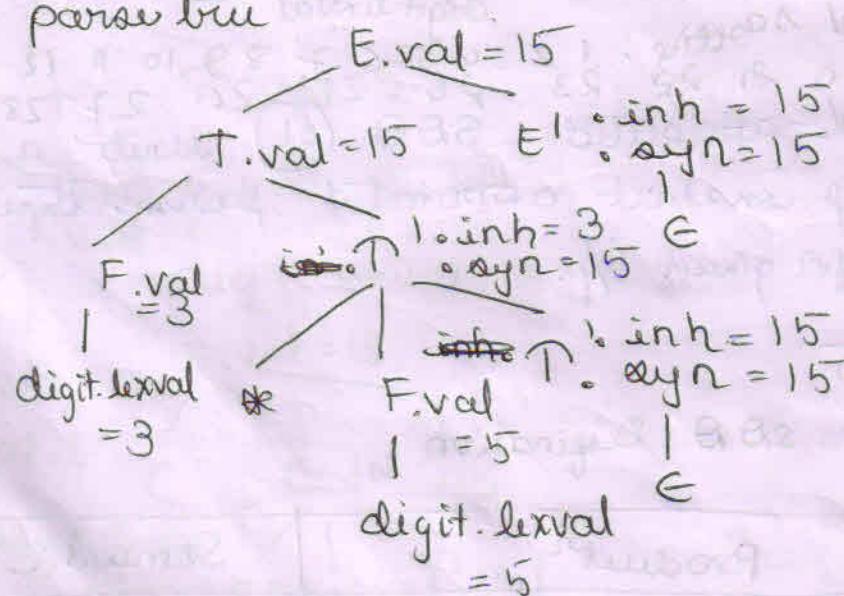
Product	Semantic rules
$E \rightarrow TE^1$	$E'.inh = T'.Val$ $E'.Val = E'.syn$

$E' \rightarrow +TE'$  $E' \rightarrow E$  $T \rightarrow FT'$  $T' \rightarrow *FT'$  $T' \rightarrow E$  $F \rightarrow (E)$  $F \rightarrow \text{digit}$  $E'_i.inh = E.iinh + T.inh$  $E'_i.syn = E'_i.syn$  $E'_i.syn = E'_i.inh$  $T'_i.inh = F.val$  $T'_i.syn = T'_i.syn$  $T'_i.inh = T.inh * F.val$  $T'_i.syn = T'_i.syn$  $T'_i.syn = T'_i.inh$  $F.val = E.val$  $F.val = \text{digit.lexval}$ 

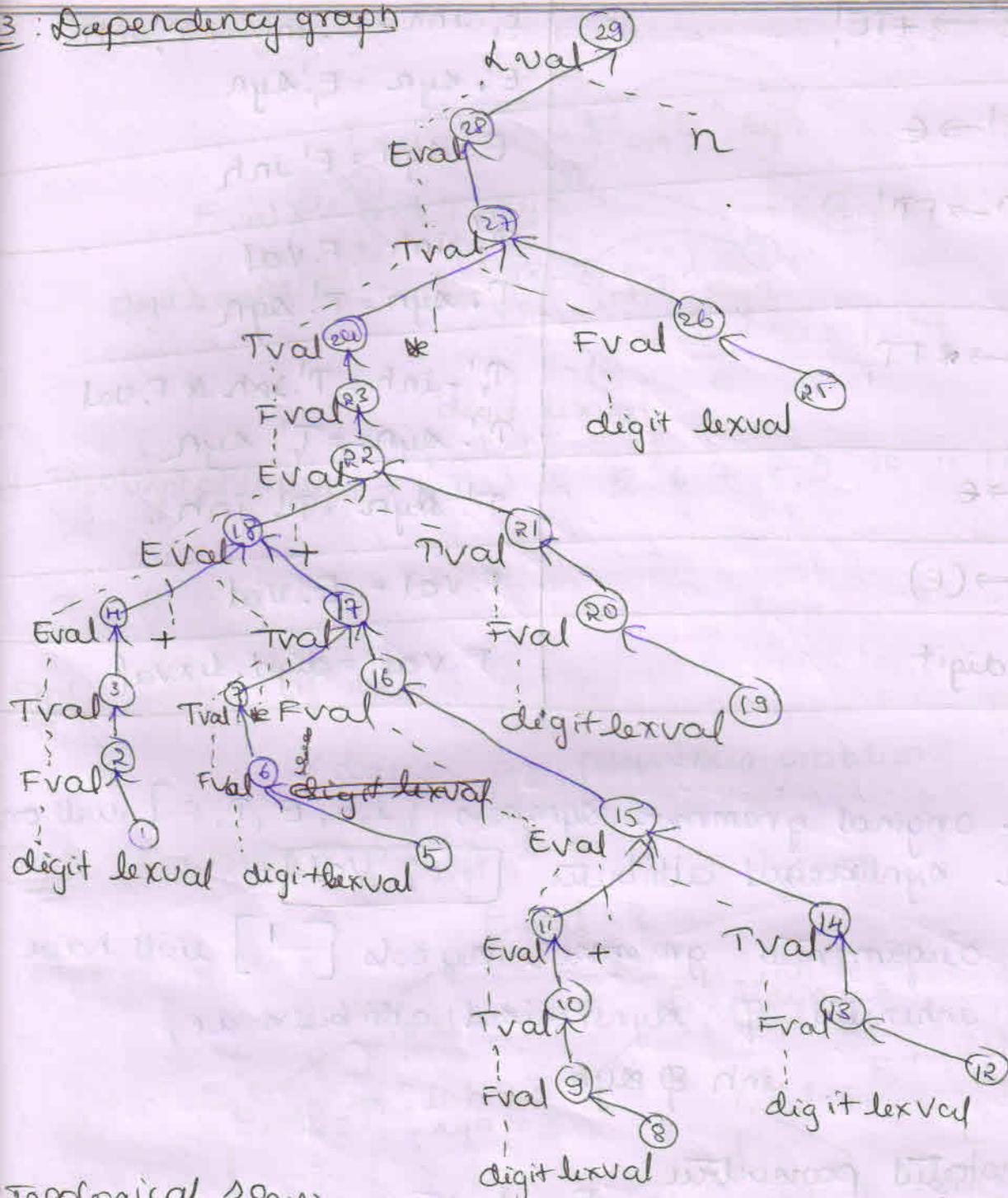
### NOTE

- The original grammar symbols [ $i.e., E, T, F$ ] will only have synthesized attributes  $\boxed{i.e., val.}$
- The augmented grammar symbols [ $'$ ] will have both inherited & synthesized attributes i.e.,  $inh \& syn$

### Step 2: Annotated parse tree



### Q3 : Dependency graph



Topological sorting: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18  
 19 20 21 22 23 24 25 26 27 28 29

Write L-attributed SAD (a) write a SAD for top down parser & construct annotated parse tree, dependency graph for given i/p.

① 3 \* 5

Step : SAD Definition

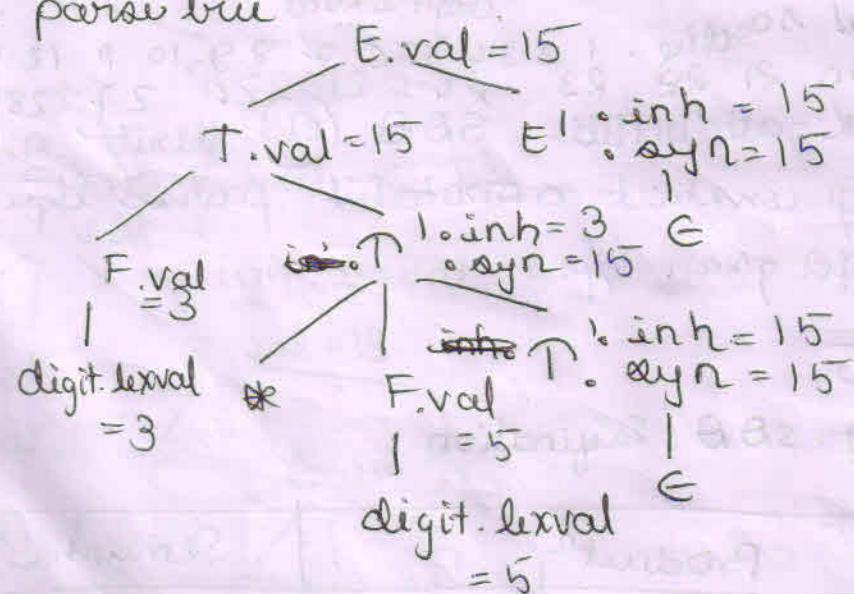
Product	Semantic rules.
$E \rightarrow TE'$	$E'.inh = T.Val$ $E.syn = E'.syn$

$E' \rightarrow +TE'$  $E'_i.inh = E.iinh + T.inh$  $E' \rightarrow E$  $E'_i.syn = E'_i.syn$  $T \rightarrow FT'$  $E'_i.syn = E'_i.inh$  $T' \rightarrow *FT'$  $T'_i.inh = F.Val$  $T' \rightarrow E$  $T'_i.syn = T'_i.syn$  $F \rightarrow (E)$  $T'_i.inh = T'_i.inh * F.Val$  $F \rightarrow \text{digit}$  $T'_i.syn = T'_i.syn$  $T'_i.syn = T'_i.inh$  $F.Val = E.Val$  $F.Val = \text{digit.lexval}$ 

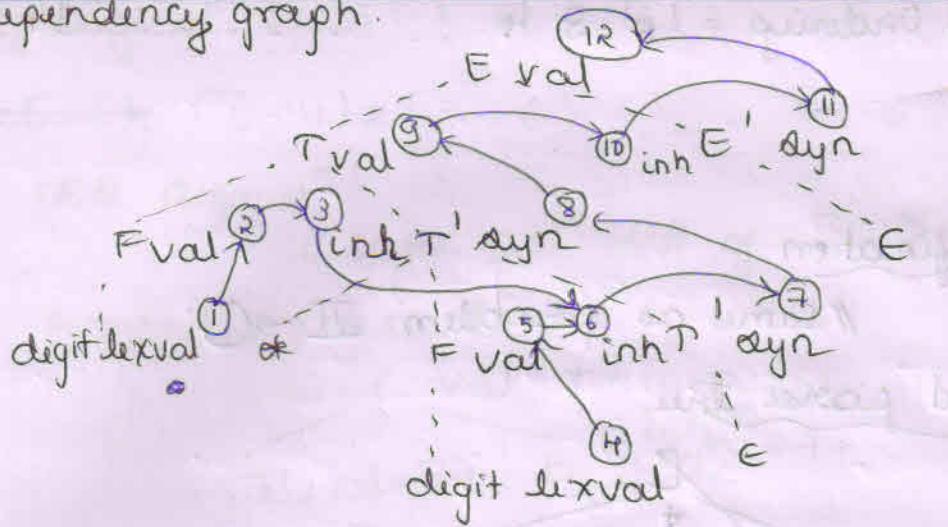
### NOTE

- The original grammar symbols [i.e.,  $E, T, F$ ] will only have synthesized attributes i.e.,  $Val$ .
- The augmented grammar symbols [ $'$ ] will have both inherited & synthesized attributes i.e.,  $inh \& syn$ .

### Step 2: Annotated parse tree



Op3: Dependency graph.



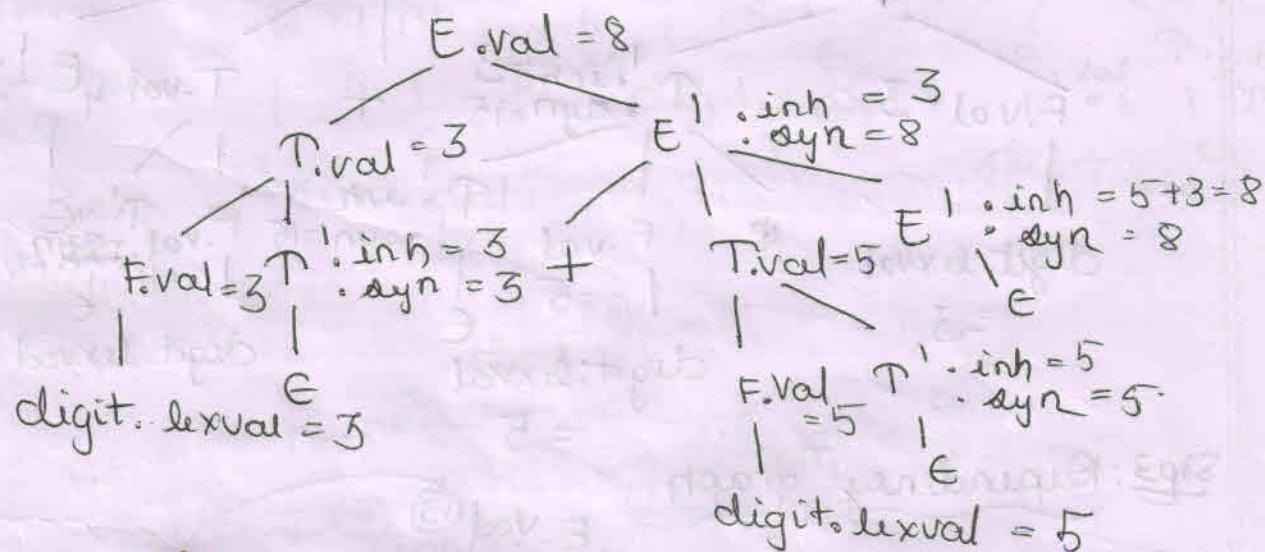
Op4: Topological Order  $\rightarrow 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11 \ 12$ .

Op5: 3 + 5

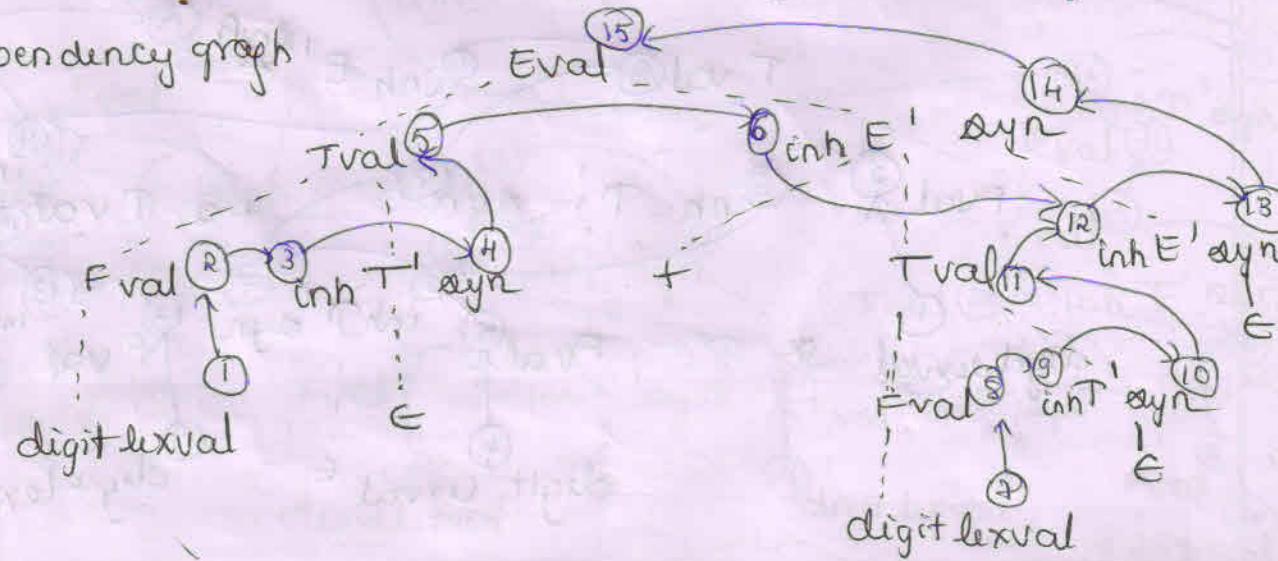
Step1: SDA definition

// same as previous problem.

Step2: Annotated parse tree



Op3: Dependency graph



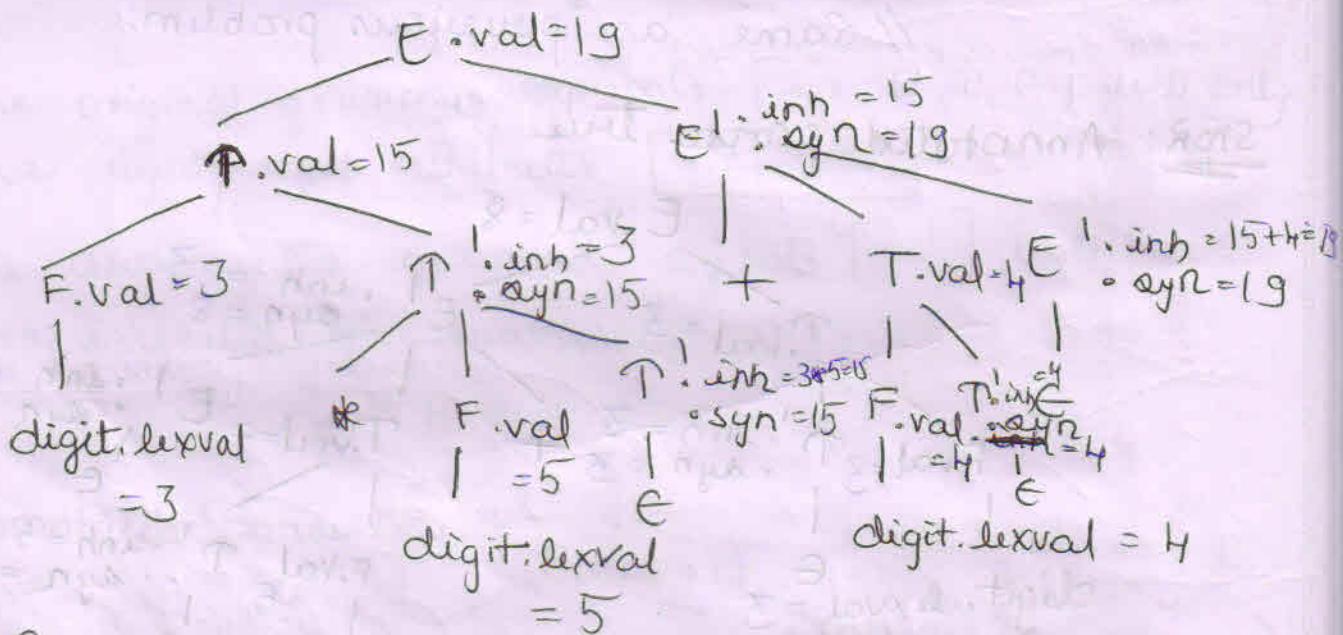
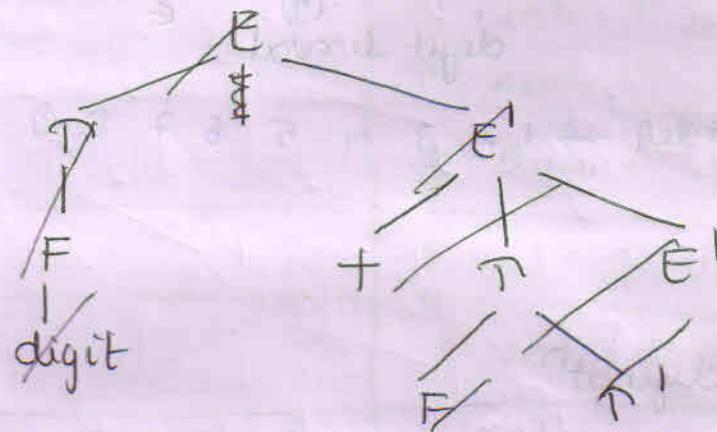
Step 1: Topological Ordering = 1 2 3 4 ... 14 15

③  $3 * 5 + 4$

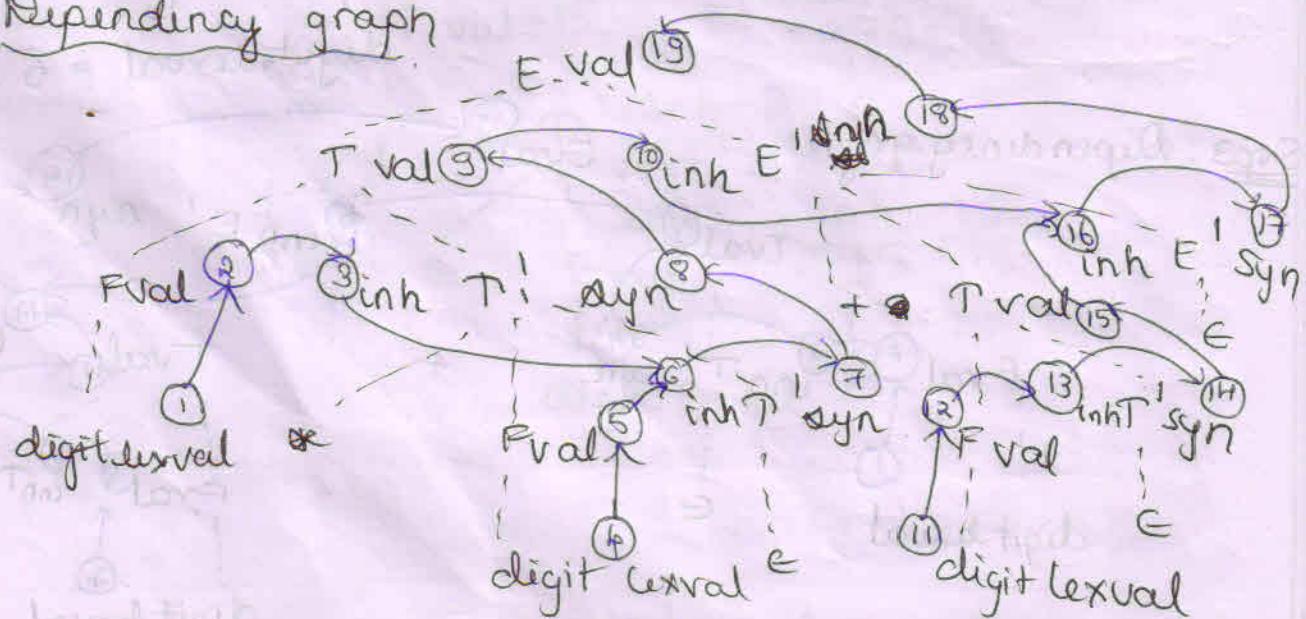
Step 1: SOD definition

// same as problem II  $\rightarrow$  ①

Step 2: Annotated parse tree



Step 3: Dependency graph



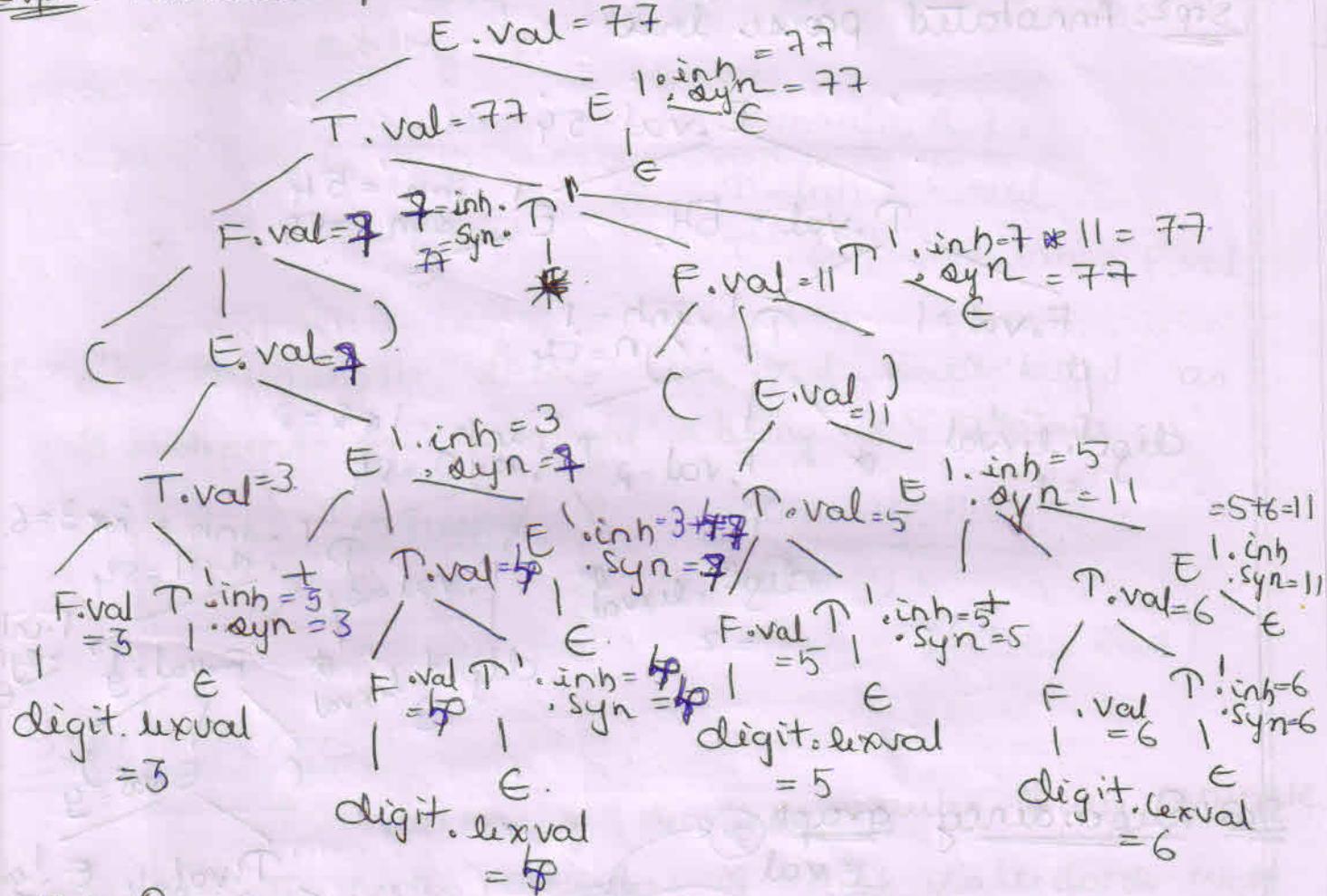
Topological Order: 1 2 3 4

19

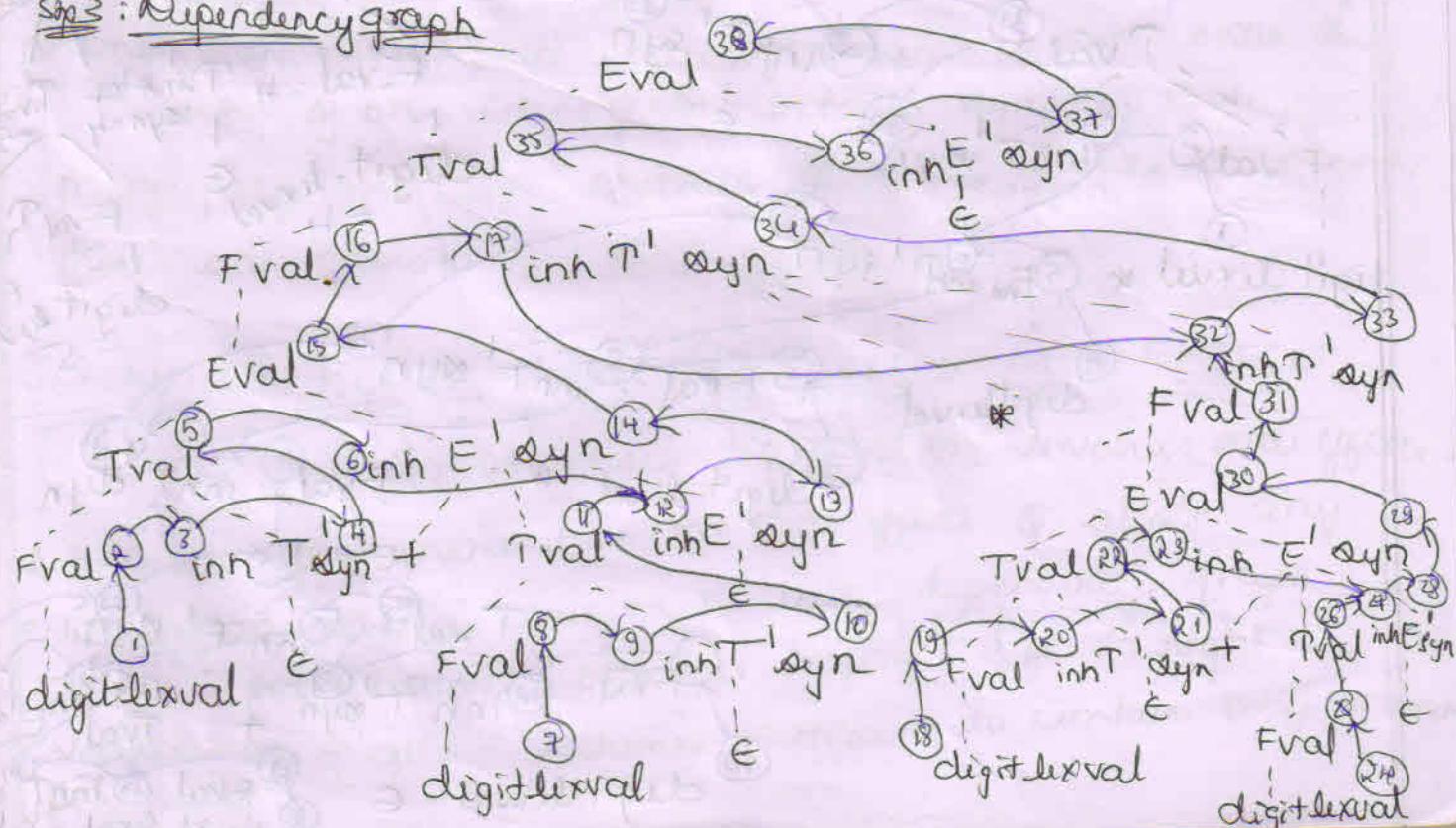
$$\cancel{3+5} \leftarrow (3+4) * (5+6)$$

Step: SAD definition  
// same as SAD of  $\Pi \rightarrow \mathbb{I}$

sfp2: Annotated parse tree



### Sip 3 : Dependency graph



# Topological Ordering

① ② ③ ④

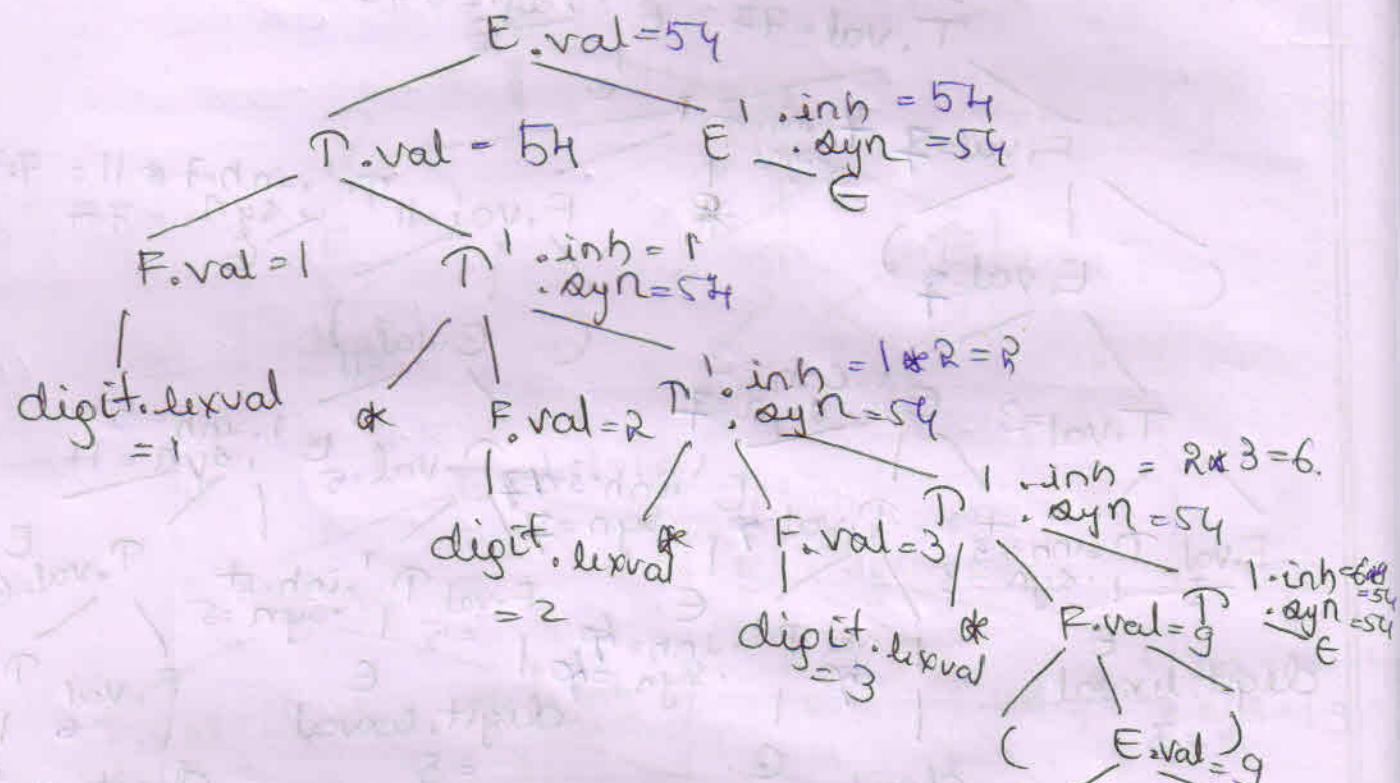
37 38

⑤ 1\*2+3\*(4+5)

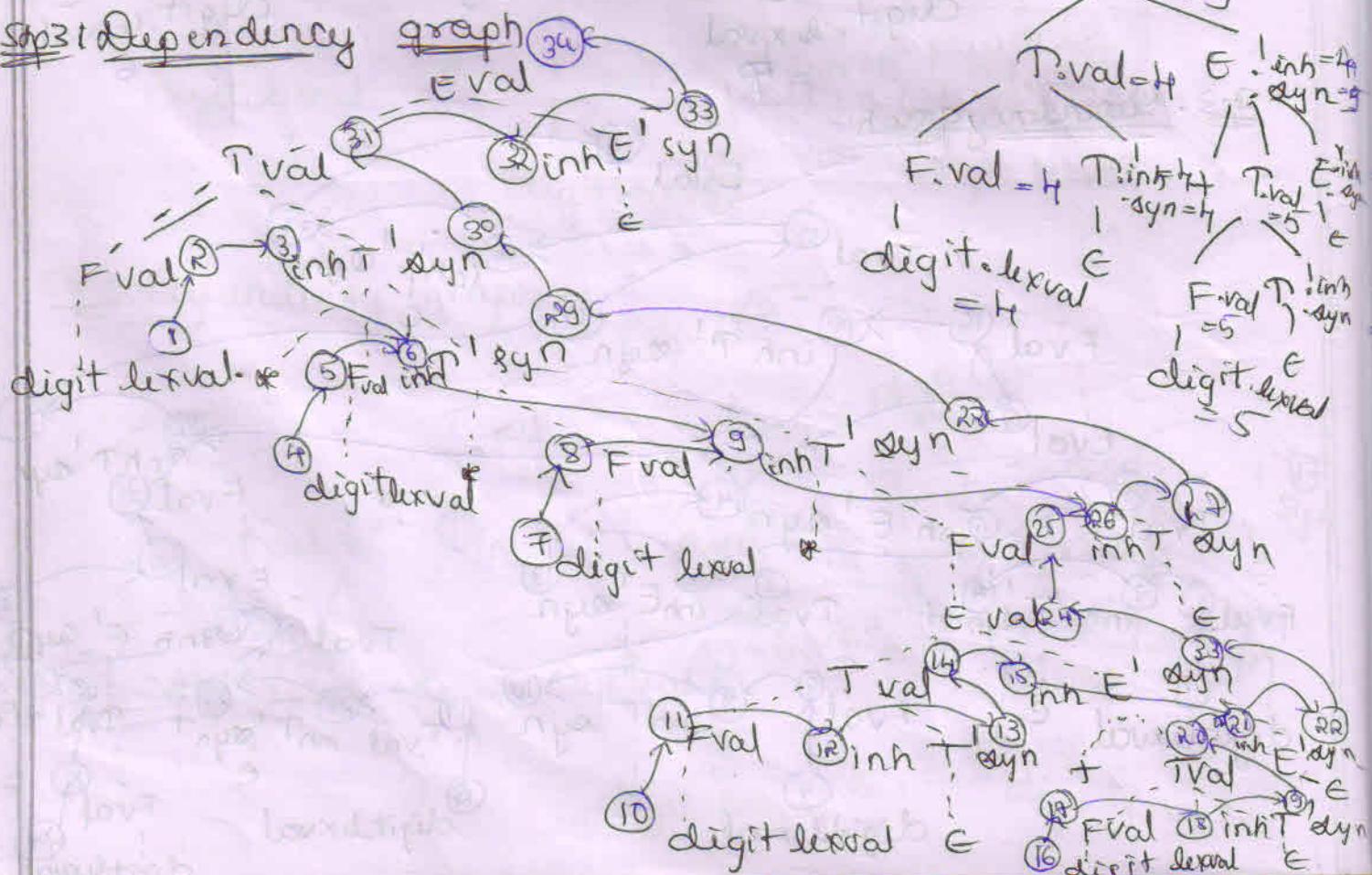
Step 1: SDD definition

// same as II → I

Step 2: Annotated parse tree



Step 3: Dependency graph



Eg1: The following definition is ~~not~~ attributed. Here the inherited attribute of  $T'$  gets its value from its left sibling  $F$ . Similarly  $T_1'$  gets its value from its parent  $T'$ 's left sibling  $F$

### Production

$$T \rightarrow FT'$$

$$T' \rightarrow FT_1'$$

### Semantic Rules

$$T'.inh = F.val$$

$$T_1'.inh = T'.inh * F.val$$

Eg2: The definitions below are not ~~not~~ attributed as  $B.i$  depends on its right sibling  $C$ 's attribute.

### Production

$$A \rightarrow BC$$

### Semantic Rules

$$A.s = B.b$$

$$B.i = f(C.c, A.s)$$

## SIDE EFFECTS

Evaluation of semantic rules may generate intermediate codes, ~~or~~ Eg: A disk calculator might print a result, a code generator might enter the type of an identifier into a symbol table, may perform type checking & may issue error msg. These are known as side effects.

## SEMANTIC RULES WITH CONTROLLED SIDE EFFECTS

In practise translation involves side effects.

Attribute grammars has no side effects & allow any evaluation order consistent with dependency graph whereas translation schemes impose left to right evaluation & allow scheme actions to contain any program fragment.

## Ways to Control Side Effects

1. permit incidental side effects that do not constrain attribute evaluation.

In other words, permit side effects when attr evaluta<sup>?</sup> based on any topological sort of the dependency graph produces a correct translation.

2. Impose constraints on allowable evaluation orders so that the same translation is produced for any allowable order.

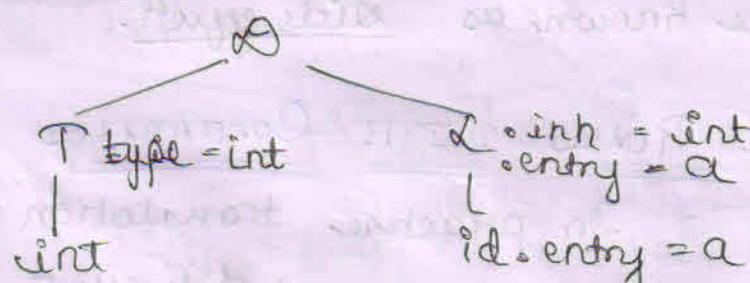
Write an SAD for simple type declaration

a) ifp : int

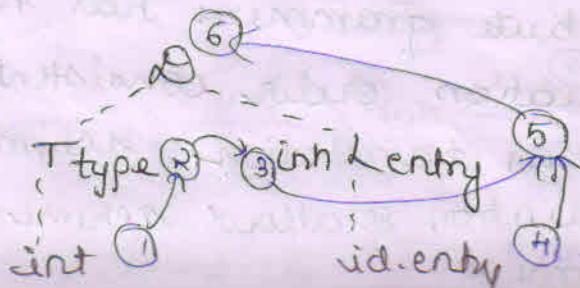
Step1:

Production	Semantic rules
$\Theta \rightarrow T \alpha$	$\alpha \cdot \text{inh} = T \cdot \text{type}$
$T \rightarrow \text{int}$	$T \cdot \text{type} = \text{int}$
$T \rightarrow \text{float}$	$T \cdot \text{type} = \text{float}$
$\alpha \rightarrow \alpha, \text{id}$	$\alpha \cdot \text{inh} = \text{id} \cdot \text{entry}$ $\text{addtype}(\text{id} \cdot \text{entry}, \alpha \cdot \text{inh})$
$\alpha \rightarrow \text{id}$	$\text{addtype}(\text{id} \cdot \text{entry}, \alpha \cdot \text{inh})$

Step2: Annotated parse tree



Step3: Dependency graph



### Explanation:

Non terminal  $\Delta$  represents a declaration, which from product 1, consists of a type  $T$  followed by a list  $L$  of identifiers.  $T$  has one attribute:  $T.type$ , which is the type in the declaration  $\Delta$ . Nonterminal  $\kappa$  has one attribute, which call  $\kappa.inh$  to emphasize that it is an inherited attribute. The purpose of  $\kappa.inh$  is to pass the declared type down the list of identifiers, so that it can be the appropriate symbol table entries.

Product  $\textcircled{1}$  &  $\textcircled{3}$  each evaluate the synthesized attribute  $T.type$  giving it the appropriate value, integer or float. This type is passed to the attribute  $\kappa.inh$  in the rule of product 1. Product  $\textcircled{4}$  passes  $\kappa.inh$  down the parser tree i.e., the value of  $\kappa.inh$  is compared at a parse tree node by copying the value of  $\kappa.inh$  from the parent of that node, the parent corresponds to the head of product. Product  $\textcircled{4}$  &  $\textcircled{5}$  also have a rule in which a function  $\text{addtype}$  is called with  $R$  arguments:

- 1)  $\text{id.entry}$  is a lexical value that points to a symbol table object.
- 2)  $\kappa.inh$ , the type being assigned to every identifier on the list.

The function  $\text{addType}$  properly installs the type  $\kappa.inh$  as the type of the represented identifier. Note that the side effect, adding the type info to the table, does not affect the evaluation order.

b) int a, b

Step 1

Product<sup>n</sup>

D → T &

& T → int

T → float

L → L, id

L → id

semantic rules.

L.inh = T.type

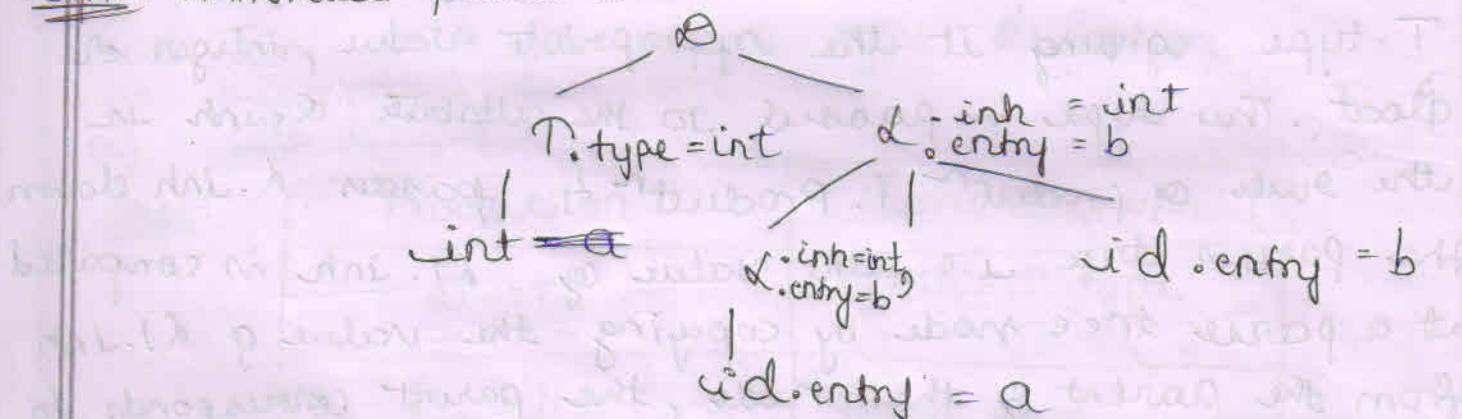
T.type = int

T.type = float

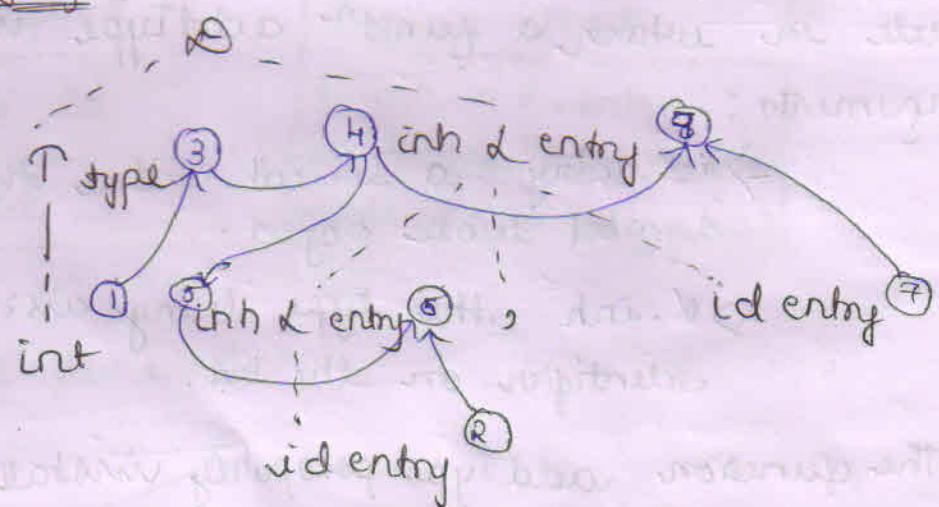
addtype(id.entry, L.inh)

addtype(id.entry, L.inh)

Step 2: Annotated parse tree



Step 3: Dependency graph



Step 4: Topological order ① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨

## UNIT-5

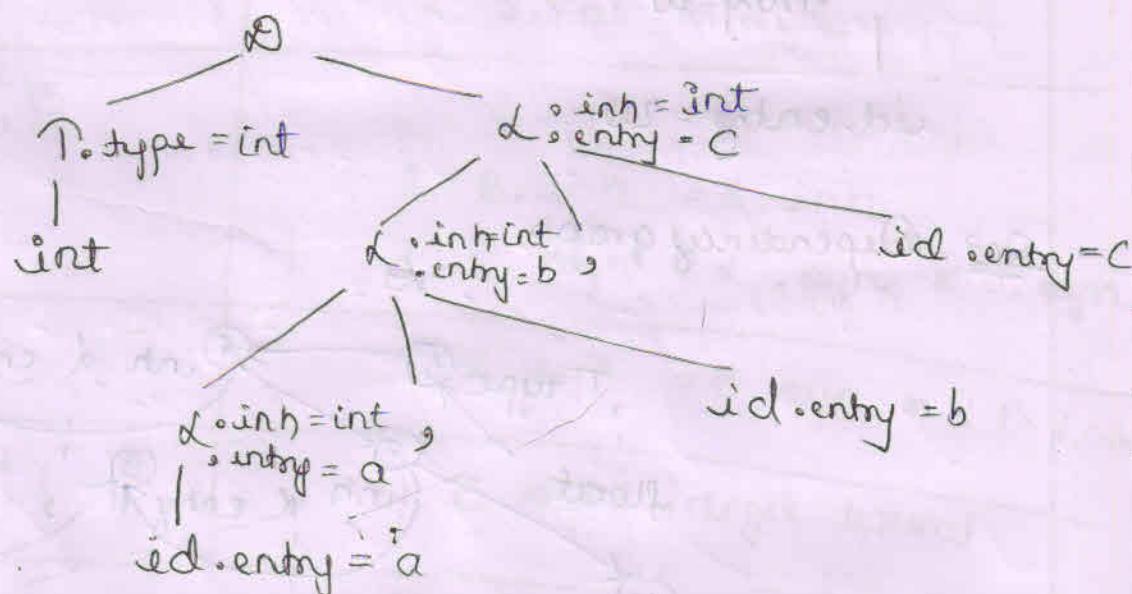
[continued . . . .]

(c) float a, b, c. or int a, b, c <Exercise 5.2.2>

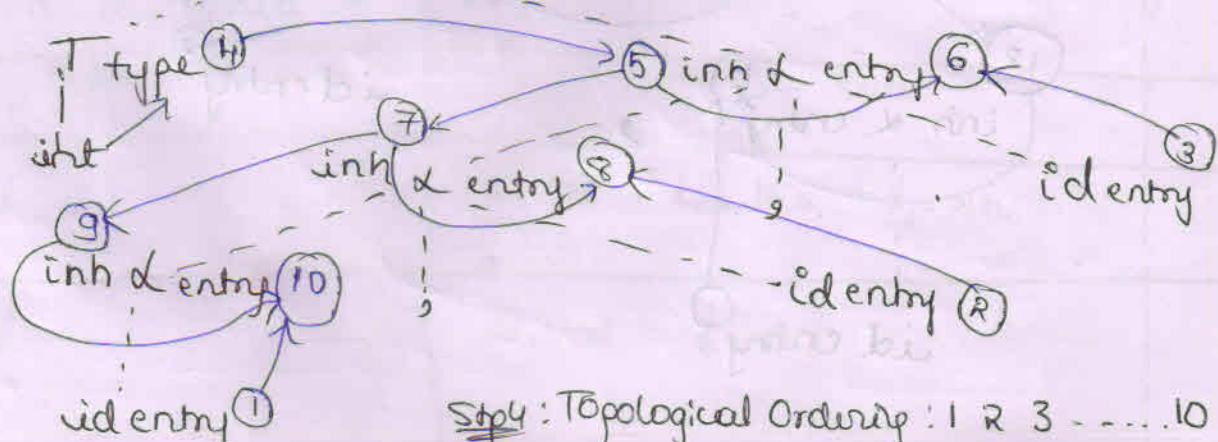
Step 1: SAD definition

Production	SAD
$S \rightarrow T \&$	$\text{d.inh} = T.\text{type}$
$T \rightarrow \text{int}$	$T.\text{type} = \text{int}$
$T \rightarrow \text{float}$	$T.\text{type} = \text{float}$
$d \rightarrow d, \text{id}$	$d.\text{inh} = d.\text{type}$ $d.\text{type}(\text{id.entry}, d.\text{inh})$
$\& \rightarrow \text{id}$	$\&.\text{type}(\text{id.entry}, d.\text{inh})$

Step 2: Annotated parse tree



Step 3: Dependency graph



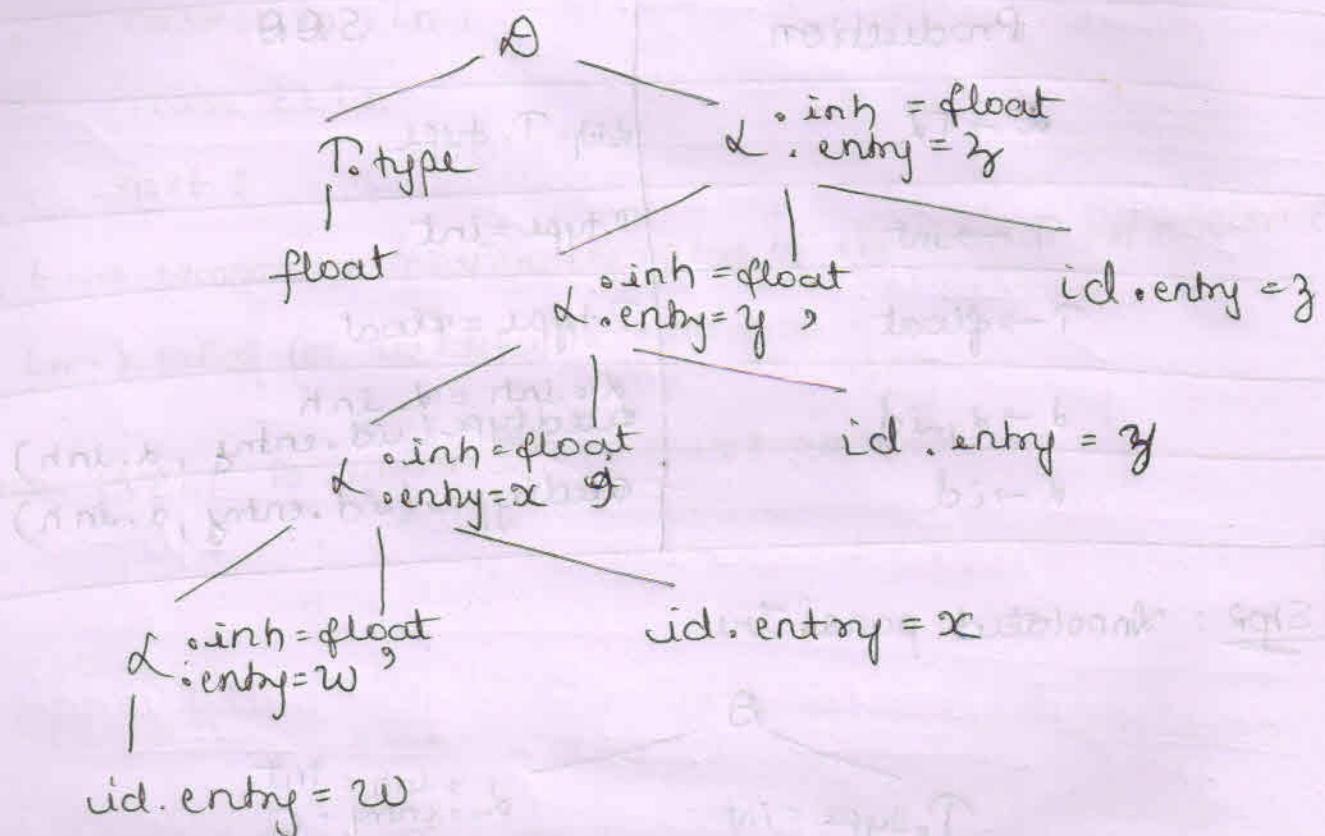
Step 4: Topological Ordering: 1 2 3 ... 10

(d) float w, x, y, z

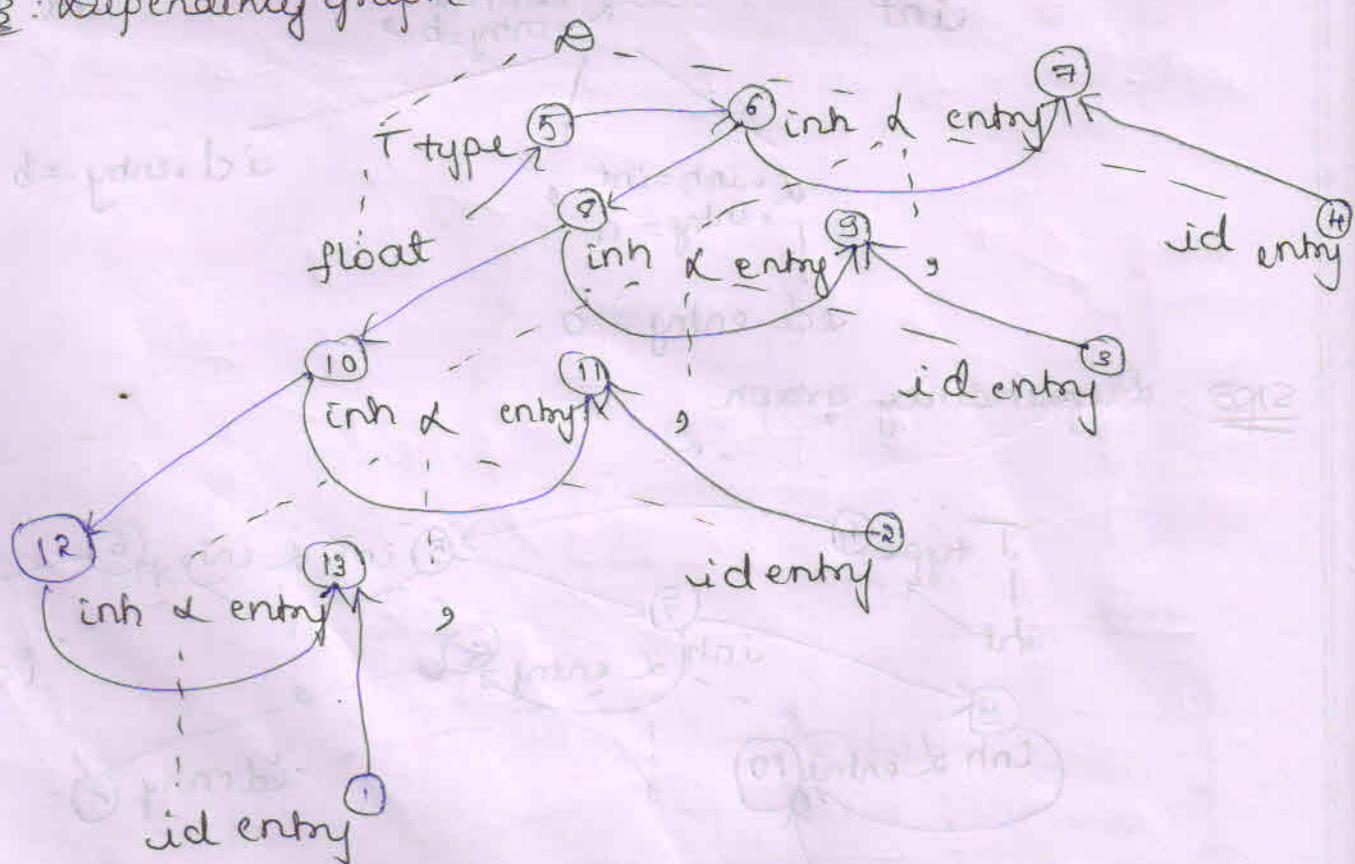
Step 1 SAD definition

// Same as previous problem.

Step 2: Annotated parse tree



Step 3: Dependency graph



Exercise 5.2.4 This grammar generates binary no with a decimal point

$$S \rightarrow d \cdot d \mid \epsilon$$

$$d \rightarrow 0 \mid 1$$

$$B \rightarrow 0 \mid 1$$

Design an L-attributed SAD to compute S.val, the decimal no value of input string. For eg., translating string 01.101 should be the decimal no 5.625.

Sols

### Production

$$1) S \rightarrow d \cdot d,$$

$$2) S \rightarrow \epsilon$$

$$3) d \rightarrow 0, 1$$

$$4) \epsilon \rightarrow B$$

$$5) B \rightarrow 0 \mid 1$$

### SAD

$$1. d.inh = 0$$

$$2. d_1.inh = -1$$

$$3. S.val = d_1.syn + d_2.syn$$

$$1. d.inh = 0$$

$$2. S.val = d.syn$$

$$1. d_1.inh = d.inh + 1$$

$$2. B.inh = d.inh$$

$$3. d_1.syn = d_1.syn * B.syn$$

$$1. \epsilon.syn = B.syn * 2^{\text{d.inh}}$$

$$1. B.syn = \text{digit.} \Delta \text{xval}$$

Exercise 5.2.5 Design an S-attributed SAD for grammar of translation described in 5.2.4

### Production

### Semantic Rule

$$1) S \rightarrow d \cdot d,$$

$$S.val = d.lhs + d_1.rhs$$

$$2) S \rightarrow d$$

$$S.val = d.lhs$$

3)  $L \rightarrow L, B$

1)  $L.lhs = L.lhs + (2^{\text{L.lhs-exponent}} * B.val)$

2)  $L.rhs = L.rhs + (2^{\text{L.rhs-exponent}} * B.val)$

3)  $L.lhs\_exponent = L.lhs\_exponent + 1$

4)  $L.rhs\_exponent = L.rhs\_exponent + 1$

4)  $L \rightarrow B$

1)  $L.lhs = 2^{\text{L.lhs-exponent}} * B.val$

2)  $L.rhs = 2^{\text{L.rhs-exponent}} * B.val$

3)  $L.lhs\_exponent = 0$

4)  $L.rhs\_exponent = -1$

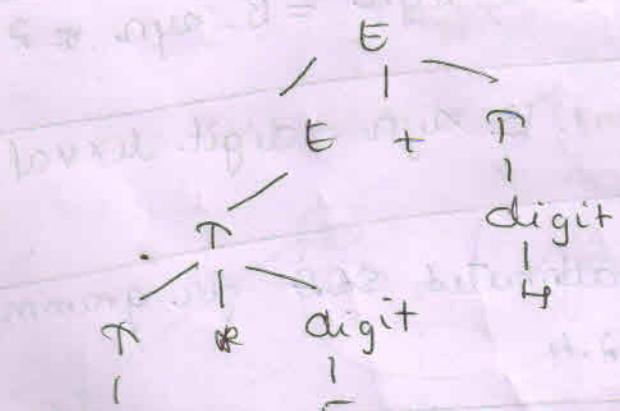
5)  $B \rightarrow 0|1$

$B.val = \text{digit.lexval}$

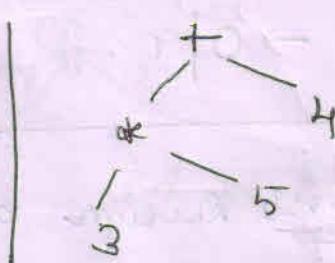
## Application Of Syntax-Directed Translation

### 1. Construction Of Syntax Tree

SDD's are useful for construction of Syntax tree.  
A syntax tree is condensed form of parse tree.



Parse tree



Syntax tree

\* Syntax trees are useful for expressing programming language constructs like expressions & statements.

- \* They help computer design by decoupling parsing from translation.
- \* Each node of a syntax tree represent a construct; the children of the node represent the meaningful components of a construct.
 

Eg: A syntax tree node representing an expression  $E_1 + E_2$  has label + & 2 children representing the sub expression  $E_1 \oplus E_2$
- \* Each node is implemented by objects with suitable no of fields; each object will have an opfield that is the label of node with additional fields as follow:
  - If the node is a leaf, an addition field holds the lexical value for the leaf. This is created by function Leaf(op, val).
  - If the node is an interior node, there are as many fields as the node has children in Syntax tree. This is created by function Node(op, c<sub>1</sub>, c<sub>2</sub>, ..., c<sub>k</sub>)

Example: The S-attributed definit<sup>n</sup> in fig below constructs Syntax trees for a simple expr grammar involving only binary operators + & -. As usual these operators are at the same precedence level & are jointly left associative. All nonterminals have one synthesized attr node, which represents a node of the syntax tree.

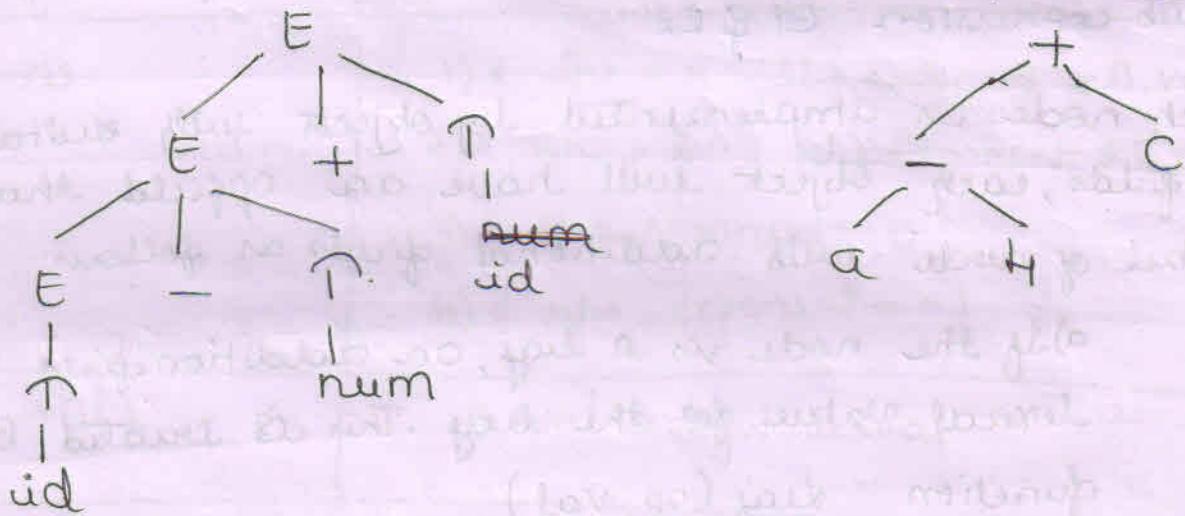
Q-H+C

PRODUCTION	SEMANTIC RULES
1. $E \rightarrow E_1 + T$	$E.\text{node} = \text{new Node}('+' , E_1.\text{node}, T.\text{node})$
2. $E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node}('-', E_1.\text{node}, T.\text{node})$
3. $E \rightarrow T$	$E.\text{node} = T.\text{node}$

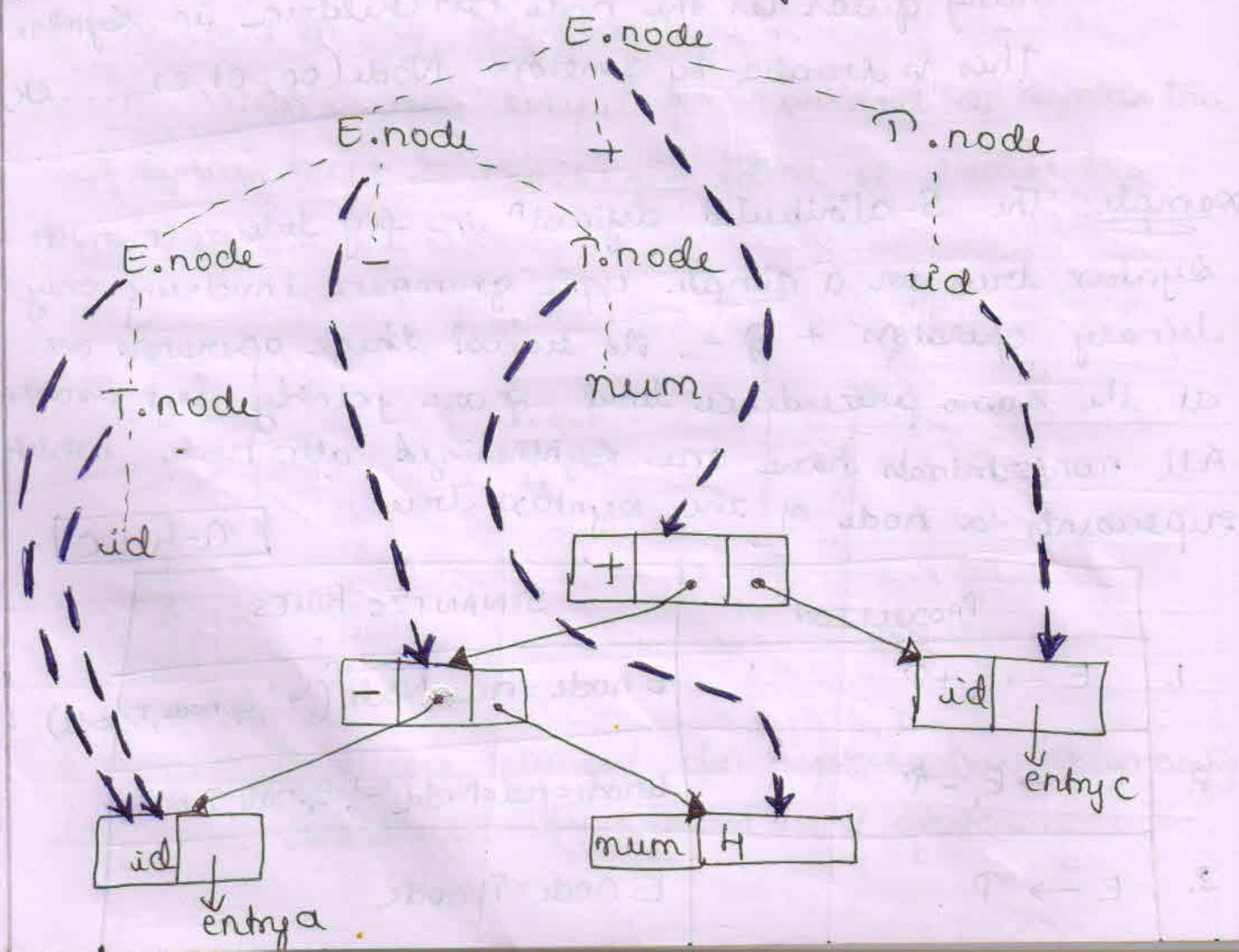
$T \rightarrow (E)$	$T.\text{node} = E.\text{node}$
$T \rightarrow \text{id}$	$T.\text{node} = \text{newLeaf}(\text{id}, \text{id}.entry)$
$T \rightarrow \text{num}$	$T.\text{node} = \text{newLeaf}(\text{num}, \text{num}.val)$

Step 2 : Parse tree

Syntax tree



Step 3 : Syntax tree for  $a - b + c$  using above SNO :



Steps in construction of the syntax tree for  $a - b + c$   
 If the rules are evaluated during a post order traversal of the parse tree, or with reduct<sup>n</sup> during a bottom up parse, then the sequence of steps shown below ends with  $p_5$  pointing to the root of the constructed syntax tree.

- 1)  $P_1 = \text{new Leaf}(\text{id}, \text{entry}-a)$
- 2)  $P_2 = \text{new Leaf}(\text{num}, +)$
- 3)  $P_3 = \text{new Node}(' - ', P_1, P_2)$
- 4)  $P_4 = \text{new Leaf}(\text{id}, \text{entry}-c)$
- 5)  $P_5 = \text{new Node}(' + ', P_3, P_4)$

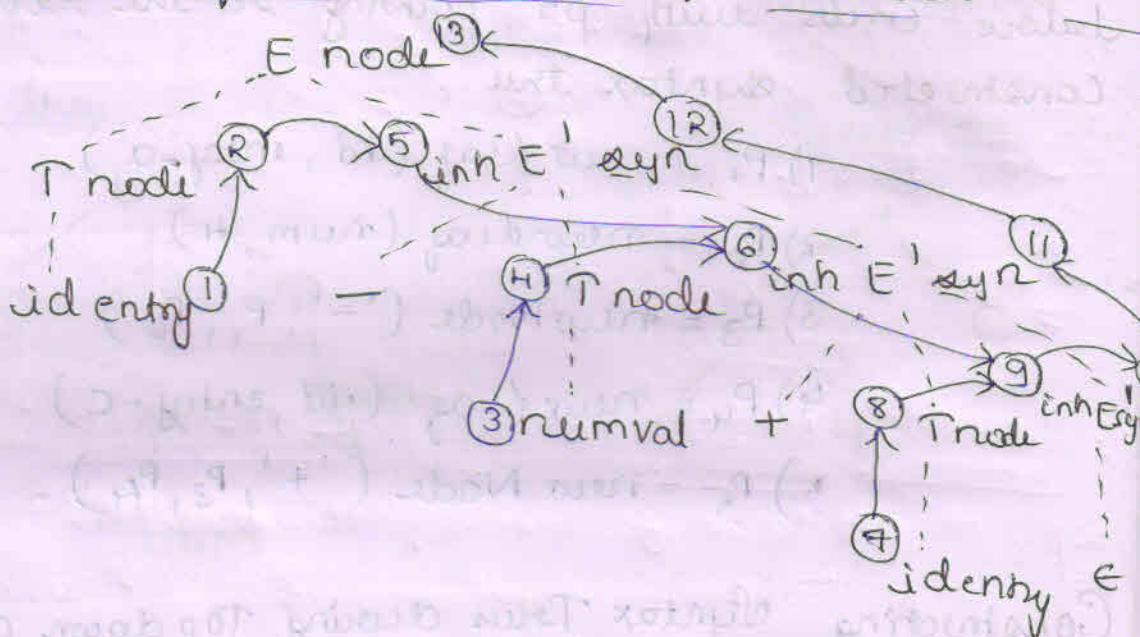
### Constructing Syntax Trees during Top down parsing

With a grammar designed for top-down parsing, the syntax trees are constructed, using the same sequence of steps, even though the structure of the parse trees differs significantly from that of syntax trees. The L-attributed definition below performs the same translation as the S-attributed definition shown before.

1) $E \rightarrow TE'$	$E.\text{node} = E^1.\text{syn}$ $E^1.\text{inh} = T.\text{node}$
2) $E^1 \rightarrow +TE^1$	$E^1.\text{inh} = \cancel{E^1.\text{inh} + P_1}$ $\text{new Node}( '+ ', E^1.\text{inh}, T.\text{node})$ $E^1.\text{syn} = E^1.\text{syn}$
3) $E^1 \rightarrow -TE^1$	$E^1.\text{inh} = \text{new Node}( ' - ', E^1.\text{inh}, T.\text{node})$ $E^1.\text{syn} = E^1.\text{syn}$
4) $E^1 \rightarrow E$	$E^1.\text{syn} = E^1.\text{inh}$
5) $T \rightarrow (E)$	$T.\text{node} = E.\text{node}$

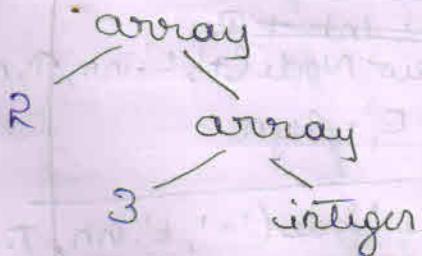
6) $T \rightarrow id$	$T.\text{node} = \text{new leaf}(id, id.\text{entry})$
7) $T \rightarrow num$	$T.\text{node} = \text{new leaf}(num, num.\text{val})$

Dependency Graph for a - H + C with Inherited SAD



### STRUCTURE OF A TYPE

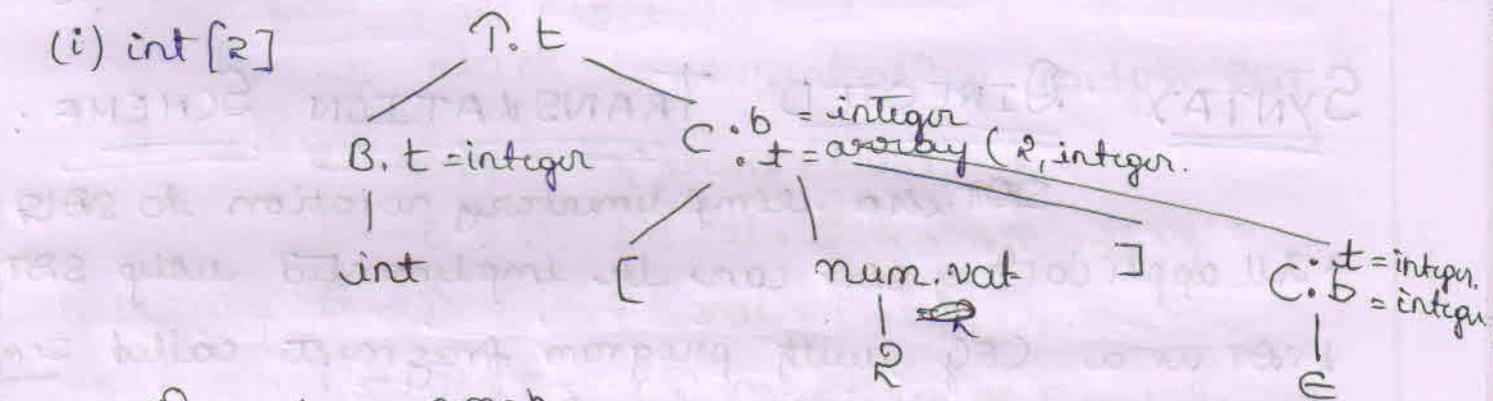
This is an example of how inherited attributes can be used to carry info from one part of the parse tree to another. In C, the type  $\text{int}[2][3]$  can be read as "array of 2 arrays of 3 integers". The corresponding type expression  $\text{array}(2, \text{array}(3, \text{integer}))$  is represented by the tree as shown below.



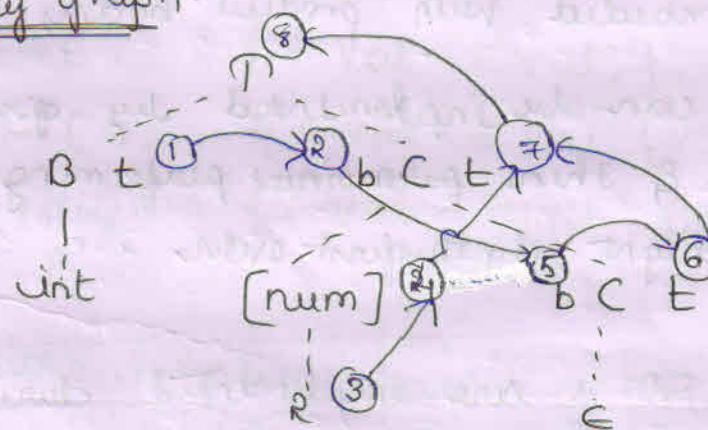
product	semantic rules
1) $T \rightarrow BC$	$T.t = C.t$ $C.b = B.t$
2) $B \rightarrow \text{int}$	$B.t = \text{integer}$
3) $B \rightarrow \text{float}$	$B.t = \text{float}$
4) $C \rightarrow [\text{num}]C$ ,	$C.t = \text{array}(\text{num}.val, C.b)$ $C.b = C.b$
5) $C \rightarrow \epsilon$	$C.t = C.b$

- The non terminals  $B \otimes T$  have a synthesized attribute  $t$  representing a type
- The non terminal  $C$  has 2 attributes: an inherited attr ( $b$ )  
of a synthesized attr ( $t$ ).
- The inherited attribute,  $b$  pass a basic type down the tree
- They synthesized attribute,  $t$  accumulate the result.
- An annotated parse tree for i/p:  $\text{int}[2]$

(i)  $\text{int}[2]$

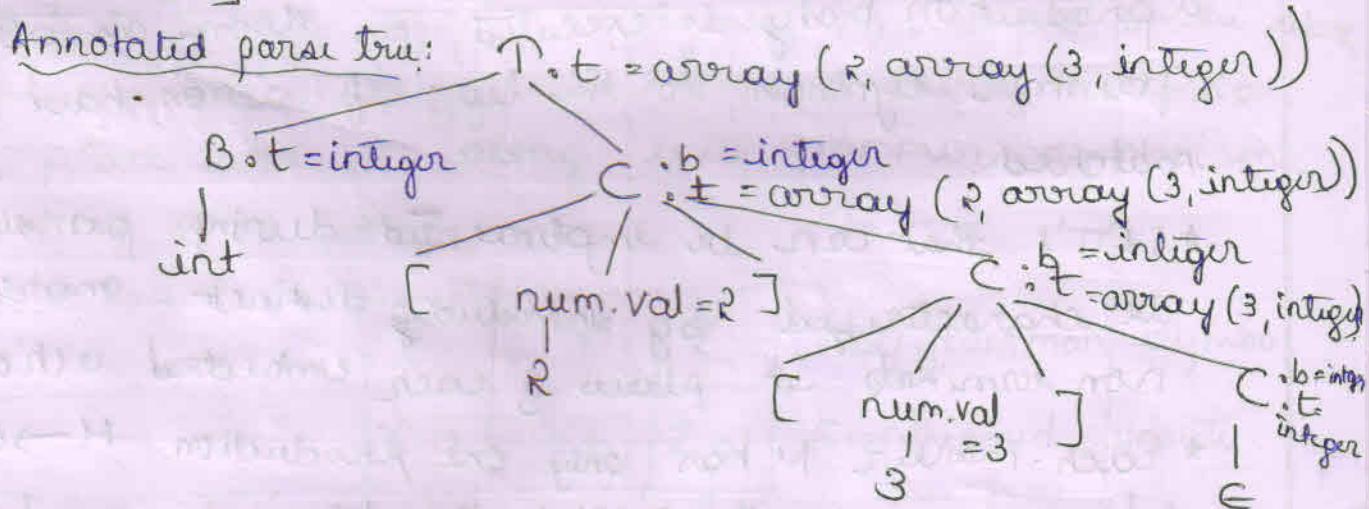


Dependency graph

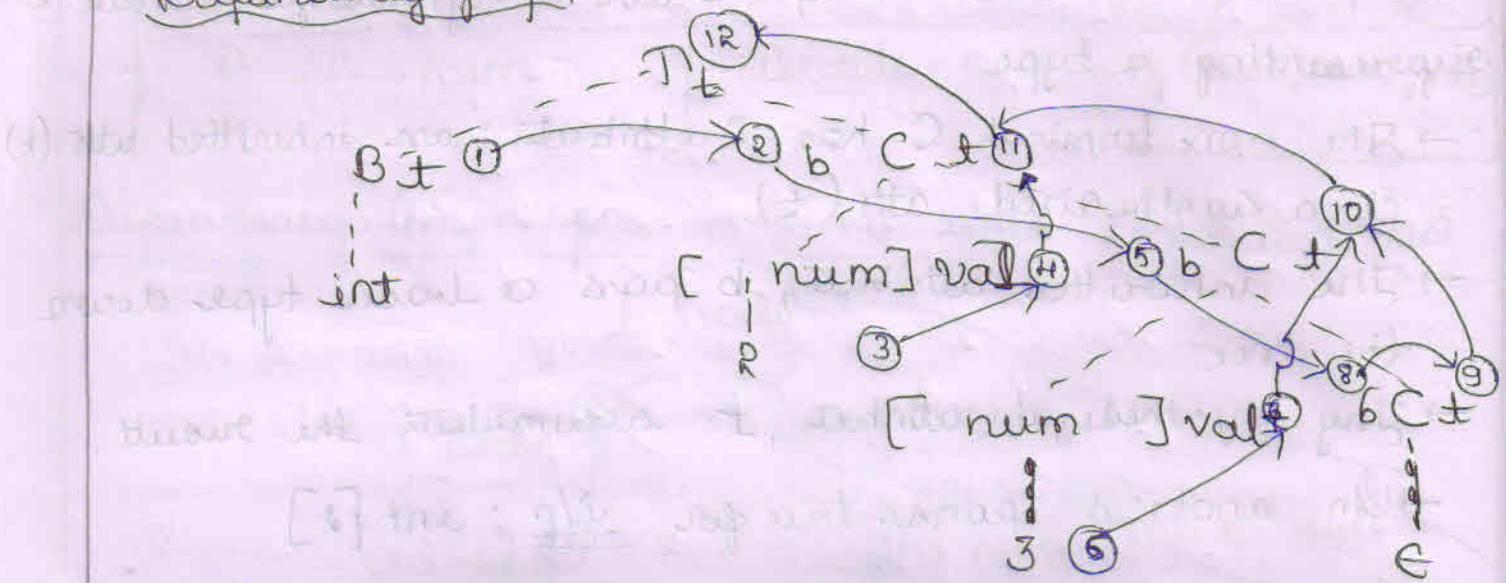


(ii)  $\text{int}[2][3]$

Annotated parse tree:



## Dependency graph



## SYNTAX DIRECTED TRANSLATION SCHEMES :

SDT is a complementary notation to SDD.

- \* All applicat<sup>n</sup> of SDD can be implemented using SDT.
- \* SDT is a CFG with program fragments called semantics actions embedded with production bodies.
- \* Any SDT can be implemented by first building a parse tree & then performing performing the actions in a left to right, depth first order i.e., during preorder traversal.
- \* Typically SDT's are implemented during parsing without building parse tree. During parsing, an action in a production body is executed as soon as all the grammar symbols to the left of action have been matched.
- \* SDT's that can be implemented during parsing can be characterized by introducing distinct marker non terminals in place of each embedded action.
- \* Each marker M has only one production  $M \rightarrow G$ .
- \* If grammar with marker non terminals can be parsed by a given method, then SDT can be implemented

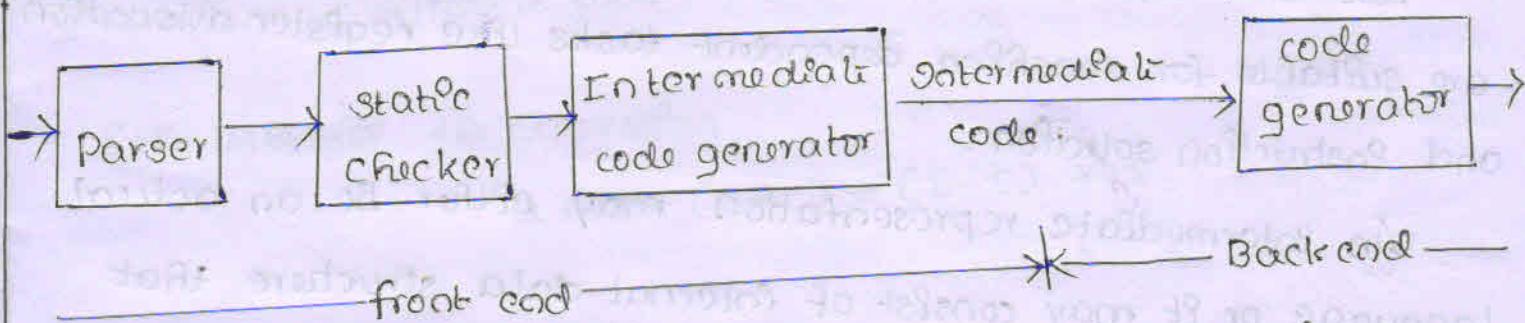
## UNIT-6

# INTERMEDIATE CODE

## GENERATION

### Intermediate Code generation.

In the analysis-synthesis model of a compiler, the front end analyzes a source program and creates an intermediate representation, from which the back end generates target code.

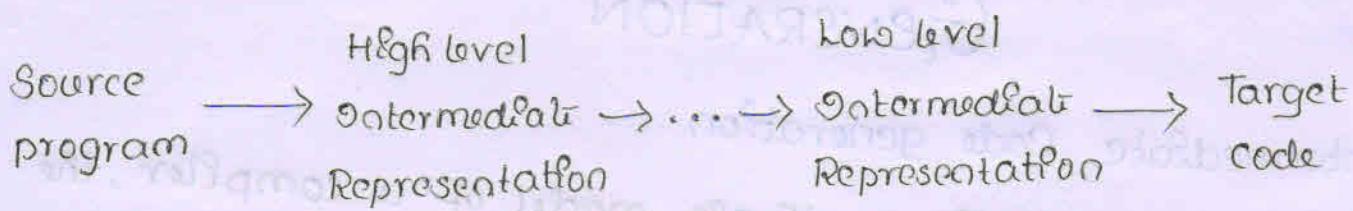


### Logical structure of a compiler front end

Parsing, static checking and intermediate code generation are done sequentially; sometimes they can be combined and folded into parsing.

Static checking includes type checking, which ensures that operators are applied to compatible operands. It also includes any syntactic checks that remain after parsing. Ex: It ensures that a break-statement in C is enclosed within a while-, for- or switch-statement; an error is reported if such an enclosing statement does not exist.

In the process of translating a program in a given source language into code for a given target machine, a compiler may construct a sequence of intermediate representation as



High level representations are close to the source language and are well suited to tasks like static type checking.

Ex: Syntax tree

Low level representations are close to the target machine & are suitable for machine dependent tasks like register allocation and instruction selection.

An intermediate representation may either be an actual language or it may consist of internal data structures that are shared by phases of the compiler.

### Variants of Syntax Trees

Nodes in a syntax tree represent constructs in the source program; the children of a node represent the meaningful components of a construct.

A directed acyclic graph (DAG) for an expression identifies the common subexpressions of the expression.

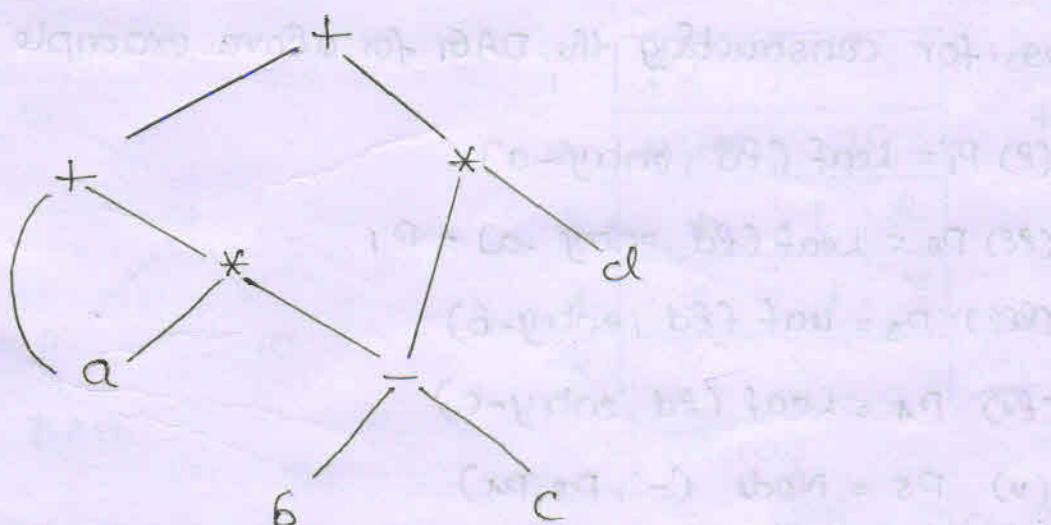
# ① Directed Acyclic Graphs for Expressions.

In DAG leaves represents the atomic operands and interior nodes represents the operators. as in the syntax tree A node  $N$  in a DAG has more than one parent if  $N$  represents a common subexpression ; But in the syntax tree , the tree for the common subexpression would be duplicated as many times as the subexpression appears. In the original expression .

DAG gives the compiler important clues regarding the generation of efficient code to evaluate the expressions .

Ex: DAG for the expression

$$a + a * (b - c) + (b - c) * d$$



↳ The leaf for 'a' has 2 parents, because 'a' appears twice in the expression

↳ The 2 occurrence of the common subexpression  $b - c$  are represented by one node , the node labeled '-'

## SDD to produce DAG.

PRODUCTION	SEMANTIC RULES
(2) $E \rightarrow E_1 + T$	$E.\text{node} = \text{new Node}('+', E_1.\text{node}, T.\text{node})$
(2E) $E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node}('-', E_1.\text{node}, T.\text{node})$
(2EE) $E \rightarrow T$	$E.\text{node} = T.\text{node}$
(2V) $T \rightarrow (E)$	$T.\text{node} = E.\text{node}$
(V) $T \rightarrow Ed$	$T.\text{node} = \text{new Leaf}(Ed, Ed.\text{entry})$
(V2) $T \rightarrow \text{num}$	$T.\text{node} = \text{new Leaf}(\text{num}, \text{num}.val)$

It will construct a DAG, before creating a new node, these functions first check whether an identical node already exists.

If a previously created identical  $\text{ex}$  node exists, the existing node is returned.

Steps for constructing the DAG for above example .

$$(2) P_1 = \text{Leaf}(Ed, \text{entry}-a)$$

$$(2E) P_2 = \text{Leaf}(Ed, \text{entry}-a) = P_1$$

$$(2EE) P_3 = \text{Leaf}(Ed, \text{entry}-b)$$

$$(2V) P_4 = \text{Leaf}(Ed, \text{entry}-c)$$

$$(V) P_5 = \text{Node}('-', P_3, P_4)$$

$$(V2) P_6 = \text{Node}('* ', P_1, P_5)$$

$$(V2E) P_7 = \text{Node}('+ ', P_1, P_6)$$

$$(2) P_8 = \text{Leaf}(Ed, \text{entry}-b) = P_3$$

$$(2X) P_9 = \text{Leaf}(Ed, \text{entry}-c) = P_4$$

$$(X) P_{10} = \text{Node}('-', P_3, P_4) = P_5$$

$$(X2) P_{11} = \text{Leaf}(Ed, \text{entry}-d)$$

(Ex8)  $P_{12} = \text{Node} ('*', P_5, P_{11})$

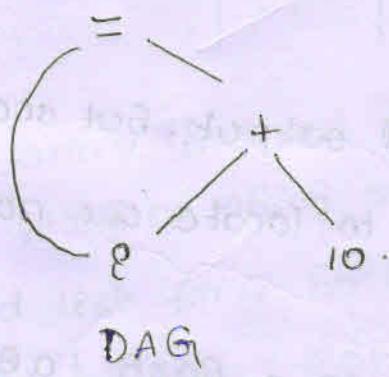
(Ex8)  $P_{13} = \text{Node} ('+', P_7, P_{12})$

When the call to Leaf (8th entry-a) is repeated at step 8, the node created by the previous call is returned, so  $P_2 = P_1$ .

### The Value-Number Method for Constructing DAG's

- \* The nodes of a DAG are stored in an array of records
- \* Each row of array represents one record & therefore one node.
- \* In each record, the first field is an operation code, indicating the label of the node. Leaves have one additional field which holds the lexical value and interior nodes have 2 additional fields indicating the left and right children.

Ex: DAG for  $E = E + 10$  allocated in an Array



DAG

	Op		
1	Op		
2	Num	10	
3	+	1	2
4	=	1	3
5			

Array

to entry  
for E

- \* In this array, we refer to nodes by giving the integer index of the record for that node within the array.
- \* This integer historically has been called the "value number".
- \* This integer refers to the expression represented by the node for the node or for the expression represented by the node.
- \* For above example node labeled + has value number 3 & its left & right children have value numbers 1 & 2 respectively.

Suppose that nodes are stored in an array & each node is referred to by its value numbers. Let the signature of an interior node be the triple  $\langle \text{op}, l, r \rangle$  where op is the label, l is its left child's value number & r its right child's value number. A unary operator may be assumed to have  $r=0$ .

ALGORITHM: To construct the nodes of a DAG using value number method.

INPUT: Label op, node l and node r

OUTPUT: The value number of a node in the array with signature  $\langle \text{op}, l, r \rangle$

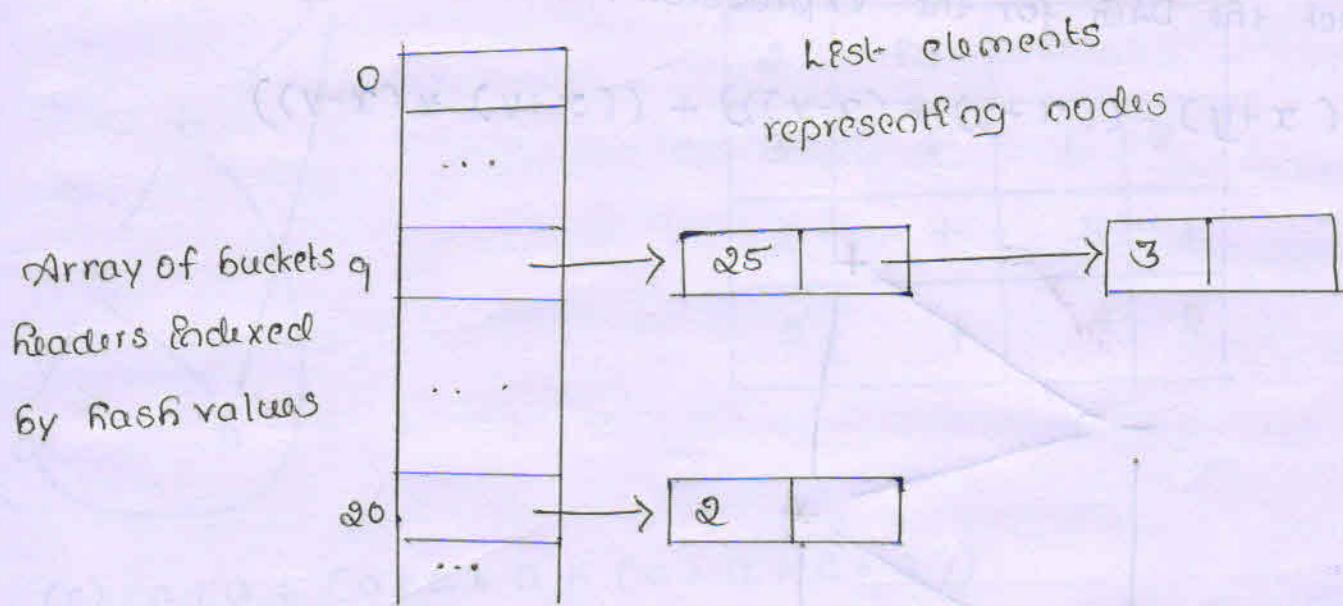
METHOD: Search the array for a node M with label op, left child l and right child r. If there is such a node, return the value number of M. If not, create in the array a new node M with label op, left child l and right child r & return its value number.

Above algorithm yields the desired output, but searching the entire array every time we are asked to locate one node is expensive.

A more efficient approach is to use a hash table, in which the nodes are put into "buckets" each of which typically will have only a few nodes. It supports dictionaries which is an abstract data type that allows us to insert & delete elements of a set & to determine whether a given element is currently in the set.

To construct a hash table for the nodes of a DAG, we need a hash function  $R$  that computes the index of the bucket for a signature  $\langle op, l, r \rangle$ .

The bucket index  $R(\langle op, l, r \rangle)$  is computed deterministically from  $op, l$  &  $r$  so that we may repeat the calculation & always get to the same bucket index for node  $\langle op, l, r \rangle$ .  
The buckets can be implemented as linked list as,



An array indexed by hash value, holds the bucket readers, each of which points to the first cell of a list. Within the linked list for a bucket, each cell holds the value number of one of the nodes that hash to that bucket. That is, node  $\langle op, l, r \rangle$  can be found on the list whose Reader is at index  $R(\langle op, l, r \rangle)$  of the array.

Thus, given the output input node  $op, l$  &  $r$  we compute the bucket index  $R(\langle op, l, r \rangle)$  & search the list of cells in this bucket for the given input node.

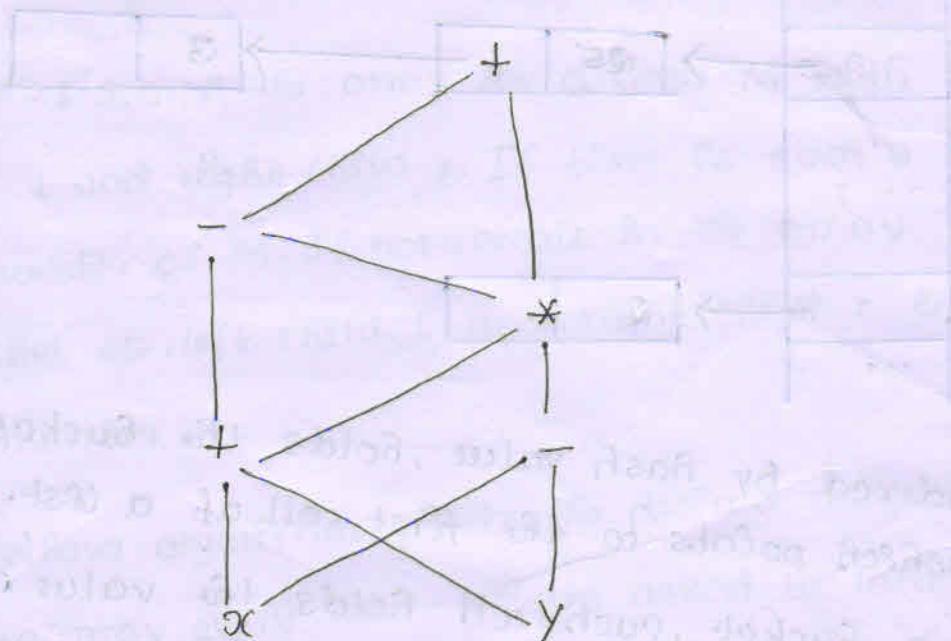
For each value number 'v' found in a cell, we must-

check whether the signature  $\langle op, l, r \rangle$  of the input node matches the node with value number  $v$  in the list of the cells. If we find a match, we return  $v$ . If we find no match, we know no such node can exist in any other bucket, so we create a new cell, add it to the list of cells for bucket- $i$ , index  $R\langle op, l, r \rangle$  & return the value number in that new cell.

Problems.

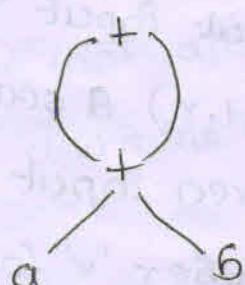
Construct the DAG for the expression.

$$((x+y) - ((x+y) * (x-y))) + ((x+y) * (x-y))$$



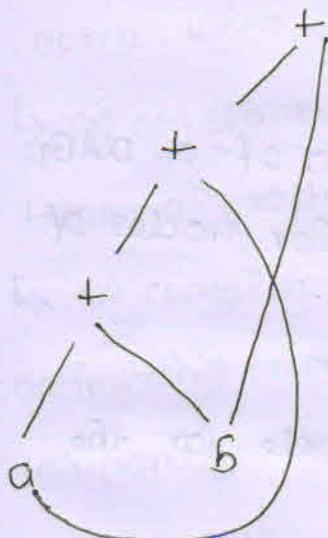
Construct the DAG & identify the value number for the sub expressions of the following expressions, assuming + associates from the left.

(a)  $a+b+(a+b)$



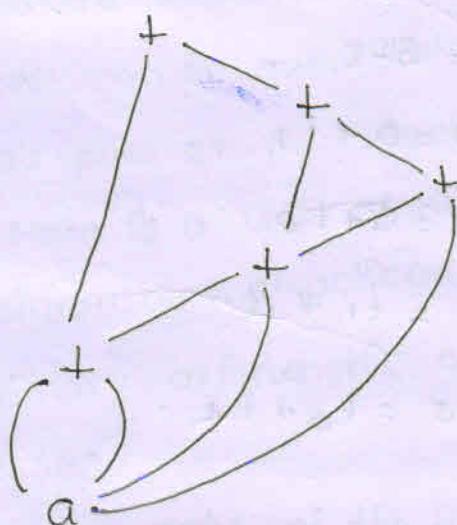
1	$\text{Ed}$	a
2	$\text{Ed}$	b
3	+	1   2
4	+	3   3

(b)  $a + b + a + b$



1	$\text{Ed}$	a	$\text{Ed}$
2	$\text{Ed}$	b	
3	+	1   2	
4	+	3   1	
5	+	4   2	

(c)  $a + a + (a + a + a + (a + a + a + a))$



1	$\text{Ed}$	a	
2	+	1   1	
3	+	2   1	
4	+	3   1	
5	+	3   4	
6	+	2   5	

## Three - Address Code

In 3-address code, there is at most one operator on the right & side of an instruction, e.g., no built up arithmetic expression are permitted.

Thus a source-language expression like  $x+y * z$  might be translated into the sequence of 3 address instructions

$$t_1 = y * z$$

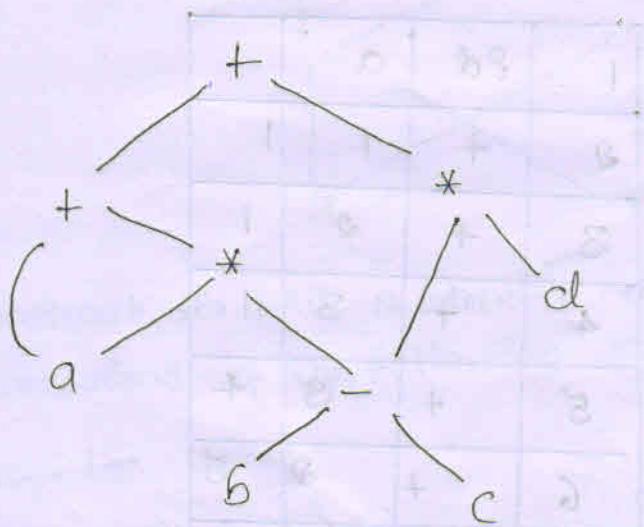
$$t_2 = x + t_1$$

3 ADDRESS (3)

where  $t_1$  &  $t_2$  are compiler generated names.

3 address code is a linearized representation of a DAG in which explicit names correspond to the interior nodes of the graph.

Ex: Write DAG & its corresponding 3 address code for the expression  $a + a * (b - c) + (b - c) * d$ .



DAG

$$t_1 = b - c$$

$$t_2 = a * t_1$$

$$t_3 = a + t_2$$

$$t_4 = t_1 * d$$

$$t_5 = t_3 + t_4$$

3 -address code

## Addresses and Instructions

3-address code is built from 2 concepts : address & instructions.

An address can be one of the following.

- ↳ A name: For convenience, we allow source program names to appear as addresses in 3-address code. In an implementation, a source name is replaced by a pointer to its symbol table entry, where all information about the name is kept.
- ↳ A constant: A compiler must deal with many different types of constants and variables.
- ↳ A compiler generated temporary: It is useful, especially in optimizing compilers, to create a distinct name each time a temporary is needed.

Symbolic labels will be used by instructions that alter the flow of control. A symbolic label represents the index of a 3-address instruction in the sequence of instructions. Actual indexes can be substituted for the labels, either by making a separate pass or by "backpatching".

- Here is a list of the common 3-address instruction forms
- (E) Assignment instructions of the form  $x = y \text{ op } z$  where  $op$  is a binary arithmetic or logical operation &  $x, y$  &  $z$  are addresses.
  - (UE) Assignments of the form  $x = \text{op } y$ , where  $\text{op}$  is a unary operation
  - (PE) Copy instructions of the form  $x = y$ , where  $x$  is assigned the value of  $y$ .

(cv) An unconditional jump goto L. The 3-address instruction with label L is the next to be executed.

(cv) Conditional jumps of the form If  $x$  goto L and If False  $x$  goto L. These instructions execute the instruction with label L next if  $x$  is true and false respectively.

(cvf) Conditional jumps such as If  $x$  relop  $y$  goto L, which apply a relational operator ( $<$ ,  $=$ ,  $\geq$  etc) to  $x$  &  $y$  & execute the instruction with label L next if  $x$  stands in relation relop to  $y$ . If not, the 3 address instruction following If  $x$  relop  $y$  goto L is executed next, in sequence.

(cvff) Procedures calls & returns are implemented using the following instructions: param  $x$  for parameters; call  $p,n$  &  $y=call p,n$  for procedure & function calls respectively & return  $y$ , where  $y$  representing a returned value, is optional.

param  $x_1$

param  $x_2$

...

param  $x_n$

call  $p,n$

The integer  $n$  indicating the no. of actual parameters in call  $p,n$  is not redundant because calls can be nested.

(cvfff) Indexed copy instructions of the form  $x=y[?]$  and  $x[?] = y$ . The instruction  $x=y[?]$  sets  $x$  to the value in the location  $y[?]$ . The instruction  $x[?] = y$  sets the contents of the locations  $y$  units beyond  $x$  to the value of  $y$ .

(Ex) Address & pointer assignments of the form  $x = \&y$ ,  $x = *y$  and  $*x = y$ .

Ex % Consider the statement

do  $\ell = \ell + 1$ ; while ( $a[\ell] < v$ );

2 possible translations of this statement are

L :  $t_1 = \ell + 1$

$\ell = t_1$

$t_2 = \ell * 8$

$t_3 = a[t_2]$

if  $t_3 < v$  goto L.

100 :  $t_1 = \ell + 1$

101 :  $\ell = t_1$

102 :  $t_2 = \ell * 8$

103 :  $t_3 = a[t_2]$

104 : if  $t_3 < v$  goto 100

(b) position number.

(a) symbolic labels

The translation in (a) uses a symbolic label L, attached to the first instruction. The translation in (b) shows position number for the instructions, starting arbitrarily at position 100. In both translations, the last instruction is a conditional jump to the first instruction. The multiplication  $\ell * 8$  is appropriate for an array of elements that each take 8 units of space.

### Quadruples

A quadruple has 4 fields, which we call op, arg<sub>1</sub>, arg<sub>2</sub>, & result. The op field contains an internal code for the operator.

Ex %  $x = y + z$  is represented by placing + in op, y in arg<sub>1</sub>, z in arg<sub>2</sub> & x in result.

The following are some exceptions to this rule.

(e) Instructions with unary operators like  $x = \text{minus } y$  or  $x = y$  do not use arg2. Note that for a copy statement like  $x = y, op \leftarrow$ , while for most other operations, the assignment operator is employed.

(ee) Operators like param use neither arg2 nor result.

(eee) Conditional & unconditional jumps put the target label in result.

(eee) Conditional & unconditional jumps put the target label in result.

Ex: Write quadruples for  $a = b * -c + b * -c$ .

$t_1 = \text{minus } c$

$t_2 = b * t_1$

$t_3 = \text{minus } c$

$t_4 = b * t_3$

$t_5 = t_2 + t_4$

$a = t_5$

(a) 3-address code

	op	arg1	arg2	result
0	minus	c		$t_1$
1	*	6	$t_1$	$t_2$
2	minus	c		$t_3$
3	*	6	$t_3$	$t_4$
4	+	$t_2$	$t_4$	$t_5$
5	=	$t_5$		a

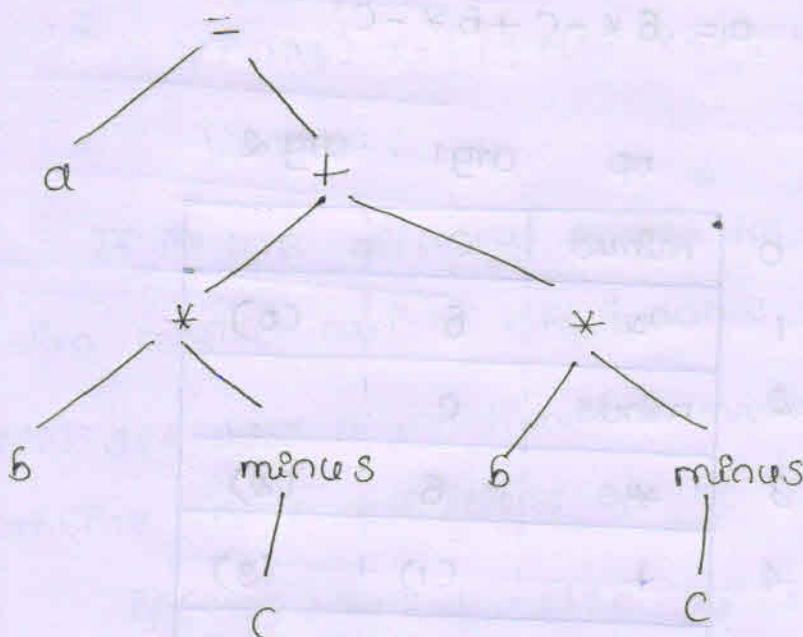
(b) Quadruples.

The special operator minus is used to distinguish the unary minus operator.

### Triples

A triples has only 3 fields, op, arg1 & arg2. Note that the result field in quadruples is used primarily for temporary names. Using triples, we refer to the result of the operation  $x op y$  by its position, rather than by an explicit temporary names.

Ex: write triples for  $a = b * c + b * c$



	op	arg1	arg2
0	minus	'c'	' '
1	*	'b'	(0)
2	minus	'c'	' '
3	*	'b'	(2)
4	+	(1)	(3)
5	=	'a'	(4)

(6) triples.

(a) Syntax tree

A ternary operator like  $x[e] = y$  requires 2 entries in the triple structure, for example, we can put  $x$  &  $e$  in one triple &  $y$  in the next.

The benefits of quadruples over triples can be seen in an optimizing compiler, where instructions are often moved around. With quadruples, if we move an instruction that computes a temporary  $t$ , then the instructions that use  $t$  require no change. With triples, the result of an operation is referred to by its position, so moving an instruction may require us to change all references to that result.

### Indirect triples.

Indirect triples consist of a listing of pointers to triples, rather than a listing of triples themselves.

With indirect triples, an optimizing compiler can move an

instruction by reordering the instruction list without affecting the triples themselves.

Ex: write indirect triples for  $a = 6 * -c + 6 * -c$

Instruction	Section
35 (0)	
36 (1)	
37 (2)	
38 (3)	
39 (4)	
40 (5)	
...	

	op	arg1	arg2
0	minus	c	
1	*	6	(0)
2	minus	c	
3	*	6	(2)
4	+	(1)	(3)
5	=	a	(4)
			...

### Static Single Assignment Form

Static single assignment form (SSA) is an intermediate representation that facilitates certain code optimization. 2 distinctive aspects of SSA that distinguish SSA from 3-address code

- (e) All assignments in SSA are to variables with distinct names.

$$\text{Ex: } p = a + b$$

$$p_1 = a + b$$

$$q_1 = p - c$$

$$q_1 = p_1 - c$$

$$p = q_1 * d$$

$$p_2 = q_1 * d$$

$$p = e - p$$

$$p_3 = e - p_2$$

$$q_2 = p + q_1$$

$$q_2 = p_3 + q_1$$

3-address code

static single assignment form

The same variable may be defined in different control flow paths in a program. For example, the source program

```
if(flag) x=-1; else x=1;
```

```
y = x * a;
```

If we use different names for  $x$  in the true part & false then conflict arises which name should use in  $y = x * a$ .  
(P2) SSA uses a notational convention called  $\phi$ -function to combine the definitions of  $x$ .

```
if(flag) x1=-1; else x2=1;
```

```
x3 =  $\phi(x_1, x_2);$ 
```

Here  $\phi(x_1, x_2)$  has the value  $x_1$  if the control flow passes through the true part of the conditional & the value  $x_2$  if the control flow passes through the false part.

Translate the arithmetic expression  $a - (b + c)$  into

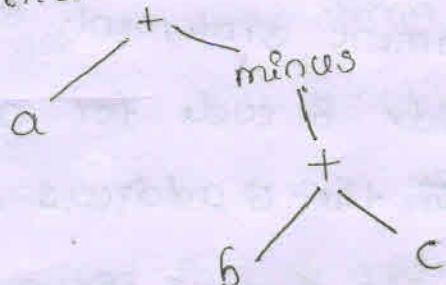
(a) A syntax tree

(b) Quadruples

(c) Triples

(d) Indirect triples.

(a) Syntax tree



$$t_1 = b + c$$

$$t_2 = \text{minus } t_1$$

$$t_3 = a + t_2$$

### (6) Quadruples

	op	arg 1	arg 2	result
1 0	+	b	c	t <sub>1</sub>
2 1	minus	t <sub>1</sub>		t <sub>2</sub>
3 2	=	a	t <sub>2</sub>	t <sub>3</sub>

### (c) triples

	op	arg 1	arg 2
0	+	b	c
1	minus	(0)	
2	=	a	(1)

### (d) Indirect triples.

Instructions

	op	arg 1	arg 2
35	(0)	+	b
36	(1)	minus	(0)
37	(2)	=	a

### Translation of Expressions.

An expression with more than one operator, like  $a+b*c$ , will translate into instructions with at most one operator per instruction. An array reference  $A[i][j]$  will expand into a sequence of 3-address instructions that calculate an address for the reference.

### Operations within Expressions

The following syntax-directed definition builds up the 3-address code for an assignment statement & using attributes code for s & attributes addr & code for an expression. Attributes s.code & e.code denotes the 3 address code for s & e respectively. Attribute e.addr denotes the address

## QUESTIONS

1. Define quadruples, triples and static single assignment form.

A quadruples has 4 fields, op, arg1, arg2 & result. The op field contains an incremental code for the operation.

Ex: the quadruples for  $a = (b * -c) + (b * -c)$

$t_1 = \text{minus } c$

$t_2 = b * t_1$

$t_3 = \text{minus } c$

$t_4 = b * t_3$

$t_5 = t_2 + t_4$

$a = t_5$

	op	arg1	arg2	result
0	minus	c		$t_1$
1	*	b	$t_1$	$t_2$
2	minus	c		$t_3$
3	*	b	$t_3$	$t_4$
4	+	$t_2$	$t_4$	$t_5$
5	=	$t_5$		a

A triples has only 3 fields op, arg1 & arg2

Ex: the triples for  $a = 6 * -c + 6 * -c$

	op	arg1	arg2
0	minus	c	
1	*	b	(0)
2	minus	c	
3	*	b	(2)
4	+	(1)	(3)
5	=	a	(4)

Static single statement assignment form is an intermediate representation that facilitates certain code optimizations

Ex:

$$p = a + b$$

$$q = p - c$$

$$p = q * d$$

$$p = e - p$$

$$q = p + q$$

(a) 3-address code

$$p_1 = a + b$$

$$q_1 = p_1 - c$$

$$p_2 = q_1 * d$$

$$p_3 = e - p_2$$

$$q_2 = p_3 + q_1$$

(b) Static single assignment form.

- Q. Develop SDD to produce directed acyclic graph for an expression show the steps for constructing the DAG for the expression  $a + a * (b - c) + (b - c) * d$ .

Syntax directed definition is,

PRODUCTION	SEMANTIC RULES
(E) $E \rightarrow E_1 + T$	$E.\text{node} = \text{new Node} ('+', E_1.\text{node}, T.\text{node})$
(EE) $E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node} ('-', E_1.\text{node}, T.\text{node})$
(ET) $E \rightarrow T$	$E.\text{node} = T.\text{node}$
(EV) $T \rightarrow (E)$	$T.\text{node} = E.\text{node}$
(V) $T \rightarrow Pd$	$T.\text{node} = \text{new Leaf} (pd, pd.\text{entry})$
(VI) $T \rightarrow \text{num}$	$T.\text{node} = \text{new Leaf} (\text{num}, \text{num}.\text{val})$

Steps for constructing the DAG

(8)  $P_1 = \text{Leaf} (\&d, \text{entry}-a)$

(88)  $P_2 = \text{Leaf} (\&d, \text{entry}-a) = P_1$

(888)  $P_3 = \text{Leaf} (\&d, \text{entry}-b)$

(8v)  $P_4 = \text{Leaf} (\&d, \text{entry}-c)$

(v)  $P_5 = \text{Node} (^-, P_3, P_4)$

(v8)  $P_6 = \text{Node} (^*, P_1, P_5)$

(v88)  $P_7 = \text{Node} (^+, P_1, P_6)$

(8x) (v888)  $P_8 = \text{Leaf} (\&d, \text{entry}-b) = P_3$

(8x)  $P_9 = \text{Leaf} (\&d, \text{entry}-c) = P_4$

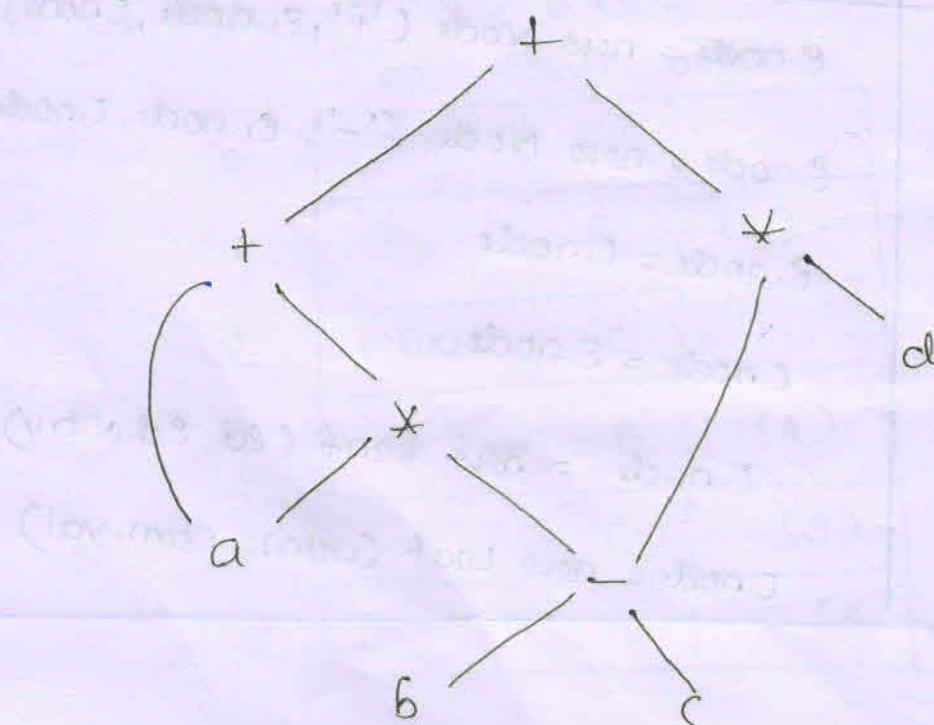
(x)  $P_{10} = \text{Node} (^-, P_3, P_4) = P_5$

(x8)  $P_{11} = \text{Leaf} (\&d, \text{entry}-d)$

(x88)  $P_{12} = \text{Node} (^*, P_5, P_{11})$

(x888)  $P_{13} = \text{Node} (^+, P_7, P_{12})$

DAG

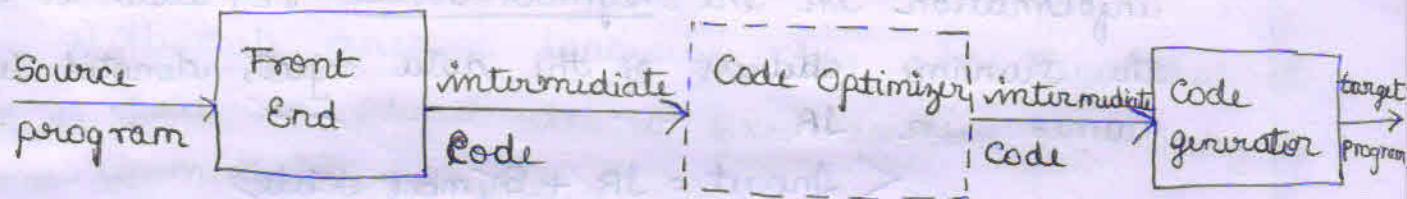


# UNIT - 8

## CODE GENERATION

### INTRODUCTION :

- \* Code generation is the final phase in the compiler design.
- \* The code optimizer accepts intermediate code representation which is generated from the front end of the compiler & produces another intermediate code representation which is optimized.
- \* Code generator takes intermediate representation produced by code optimizer along with supplementary information in symbol table of the source program & produce as output an equivalent target program.



- \* Code generator has 3 main tasks:

- 1) Instruction selection
- 2) Register allocation & assignment
- 3) Instruction Ordering

### 1) INSTRUCTION SELECTION :

Choose appropriate target machine codes instructions to implement the IR [intermediate representation statements]

### 2) REGISTER ALLOCATION & ASSIGNMENT :

Decide what values to keep in which registers

### 3) INSTRUCTION ORDERING

Decide in what order to schedule the execution of instructions.

8.1

### ISSUES IN THE DESIGN OF CODE GENERATOR:

1) Input to the code generator

2) The target program

3) Instruction selection

4) Register Allocation

5) Evaluation Order

1) Input to the code generator

\* Input to the code generator is the intermediate representation of the source program produced by the front end along with information in the symbol table i.e., used to determine the runtime address of the data objects denoted by the names in IR.

< Input = IR + Symbol table >

\* IR has several choices

(a) 3-address representation : quadruples, triples, indirect triples

(b) Virtual machine representation : byte codes of stack machine codes

(c) Linear representation such as postfix notation

(d) Graphical representation such as syntax trees or DAG's

\* Assumptions made are

(i) Front end produces low-level IR, i.e., values of names in it can be directly manipulated by the machine instruction.

(ii) Syntactic & semantic errors have been already detected

## 2) The Target Program:

- \* The output of code generator is target program.
- \* The instruction set architecture of the target machine has a significant impact on the design of code generator.
- \* Most common architectures are:
  - (a) CISC: It has few registers, has maximum of 2 operands & variety of addressing mode, variable length instructions & instruction with side effects.
  - (b) RISC: It has many registers, has maximum of 3 operands with simple addressing modes, & relatively simple instruction set architecture.
- \* Output may take variety of forms.
  - a) Absolute machine language [Executable code]
  - b) Relocatable machine language [object files for linker]
  - c) Assembly language [facilitates debugging]
- a) Absolute machine language has advantage that it can be placed in a fixed location in memory & immediately executed.
- b) Relocatable machine language program allows subprograms to be compiled separately.
- c) Producing Assembly language program as output makes the process of code generation somewhat easier.

## 3) Instruction Selection

The code generator must map the IR program into a code sequence that can be executed by the target machine.

\* The complexity of performing this mapping is determined by the factors such as:

- (i) the level of the IR
- (ii) the nature of the instruction set architectures.
- (iii) the desired quality of the generated code.

(i) the levels of the IR:

- > If the [IR is high level], use code templates to translate each IR statements into a sequence of machine instruction.
- > produces poor code, needs further optimizat<sup>n</sup>.
- > If the [IR is low level], use ~~code, this~~ <sup>low level</sup> information to generate more efficient code sequence.

(ii) the nature of the instruction set architectures has strong effect on difficulty of instruct<sup>r</sup> select<sup>r</sup>.

- > Uniformity & completeness of the instruct<sup>r</sup> set are imp factors.
- > If we do not care about the efficiency of the target program, instruct<sup>r</sup> select<sup>r</sup> is straightforward.

> For eg:

$x = y + z \Rightarrow$	LD R <sub>0</sub> , y
	ADD R <sub>0</sub> , R <sub>0</sub> , z
	ST x, R <sub>0</sub>

∴ produces redundant LD & store

eg2:

$a = b + c \Rightarrow$	LD R <sub>0</sub> , b
$d = a + b$	ADD R <sub>0</sub> , R <sub>0</sub> , c

ST a, R <sub>0</sub>
LD R <sub>0</sub> , a

ADD R <sub>0</sub> , R <sub>0</sub> , c
ST d, R <sub>0</sub>

REDUNDANT

(iii) the quality of the generated code is determined by its speed & size.

> For eg:

$$\begin{array}{l} a = a + 1 \rightarrow LD R_0, a \\ \quad \quad \quad ADD R_0, R_0, \#1 \\ \quad \quad \quad ST a, R_0 \end{array} \quad \begin{array}{l} \text{replaced} \\ \text{by} \end{array} \quad \begin{array}{l} INC a \end{array}$$

#### 4) Register Allocation:

\* Instruc<sup>n</sup> involving register operands are usually shorter & faster than those involving operands in memory.

\* 2 subproblems:

(i) Register allocation: Select the set of variables that will reside in registers at each point in the program.

(ii) Register assignment: Select specific register that a variable will reside in.

\* Complications imposed by the hardware architecture  
Eg: Register pairs for multiplication & division.

\* Multiplication instr<sup>n</sup> is of the form

$$M \boxed{x, y}$$

where  $x$  → Multiplicand, is the odd register of an even/odd register pair.  
 $y$  → Multiplier, is ~~the~~ a single register.

⇒ Product → occupies the entire even/odd register pair.

\* Division instr<sup>n</sup> is of the form

$$D \boxed{x, y}$$

where  $x$  → dividend, occupies even register  
 $y$  → divisor, occupies odd/even register  
⇒ Quotient → stored in even register  
remainder → stored in even register

Eg: two 3-address code sequences

$$t = a + b$$

$$t = t * c$$

$$t = t / d$$

$$t = a + b$$

$$t = t + c$$

$$t = t / d$$

Optimal machine-Code sequences

L R1, a

A R1, b

M R0, c

D R0, d

ST R1, t

L R0, a

A R0, b

A R0, c

SRDA R0, 3R

D R0, d

ST R1, t

### 5) Evaluation Order:

- \* The order in which computations are performed can effect the efficiency of the target code.
- \* When instructions are independent their evaluation order can be changed.
- \* Some computation orders require fewer registers to hold intermediate results than others.
- \* However picking a best order in the general case is a difficult NP-complete problem.

ADDITIONAL INFORMATION: Eg

$$t_1 = a + b$$
$$a + b - (c + d) * e \Rightarrow t_2 = c + d$$

$$t_3 = e * t_2$$

$$t_4 = t_1 - t_3$$

Reorder↓

$$t_2 = c + d$$

$$t_3 = e * t_2$$

$$t_1 = a + b$$

$$t_4 = t_1 - t_3$$

MOV R0, a

ADD R0, b

MOV R1, R0

MOV R1, c

ADD R1, d

MOV R0, e

MUL R0, R1

MOV R1, t1

SUB R1, R0

MOV B4, t4, R1

MOV R0, c

ADD R0, d

MOV R1, e

MUL R1, R0

MOV R0, a

ADD R0, b

SUB R0, R1

MOV t4, R0

## THE

### 8.2 THE TARGET LANGUAGE:

For designing a good code generator, we need to have familiarity with target machine & its instruction set. Instead of generating code on a specific target machine, a general machine consisting of many registers are considered.

#### A SIMPLE TARGET MACHINE MODEL:

The characteristics of target machine mode with instruction format & instruction set are shown below:

\* Our hypothetical machine:

(i) It is a 3-address machine with the following format

OP destination, Source1, Source2

#### NOTE:

A 3 address instruction can have 2 operands or 1 operand also but it can have max of 3 operands

(ii) The target machine is byte addressable i.e., it can access 8 bit of info from specific address

(iii) It has n no of Registers denoted by

$R_0, R_1, R_2, \dots, R_{n-1}$

\* Various types of instruction that are used by target m/c

(i) Load Instruction

(ii) Store Instruction

(iii) Computational Instruction

(iv) Unconditional Instruction

(v) Conditional Instruction

(i) Load Instruction: Used to copy the data into destination operand which must be a register.

SYNTAX: LD dest, addr

where addr operand  $\rightarrow$  register or memory locat

(ii) Store instruction: Used to copy the data into mem locat<sup>n</sup> specified in the destinat<sup>n</sup> operand.

SYNTAX: ST dst, or

where dst  $\rightarrow$  destination & it is a mem location

or  $\rightarrow$  register.

Computational operation.

(iii) Arithmetic instruction: They are performed using these instruction.

SYNTAX: OP dst, Src1, Src2.

where 1<sup>st</sup> operand, dst  $\rightarrow$  destination

2<sup>nd</sup> & 3<sup>rd</sup> operand  $\rightarrow$  Operands where R values fetched for Operat<sup>n</sup> to be per

Eg1: ADD R0, R1, RR //  $R0 = R1 + RR$

Eg2: SUB R0, R0, R1 //  $R0 = R0 - R1$

Eg3: MUL RR, R0, R1 //  $RR = R0 * R1$

(iv) Unconditional Jumps: The branch instruct<sup>n</sup> without any condit<sup>n</sup> are called unconditional jumps.

SYNTAX: BR label

where BR  $\rightarrow$  BRanch instruct<sup>n</sup>

(v) Conditional Jumps: Based on the value stored in a register i.e., whether it is true or zero or -ve, if branching takes place, then the branch inst<sup>n</sup> are called conditional jumps.

SYNTAX: Bcond or, label

where B stands from Branch,

Cond can be LT, GT, LTX, GTEX

Less than or equal  
Less than  
Greater than  
Greater than or equal

or → register, contains value such as 0, +ve or -ve.

Eg1: BL RO, T1

// Branch to T1, if RO contains +ve value

Eg2: BLTR RI, TR

// Branch to TR, if RI contains either 0 or -ve value

- \* Different addressing modes supported by generalized target machine:

1) Direct addressing mode

2) Indexed —————

3) Integer Indexed —————

4) Indirect —————

5) Immediate —————

#### (i) Direct A/M:

Address of the data to be accessed is directly present in the instructn, i.e., location is identified by a variable name x.

Eg: LD P LD RI, x

// Load value stored in memory locat<sup>n</sup> x into RI

(ii) Indexed A/M: The data can be accessed from a memory locat<sup>n</sup> using index. This addressing mode is useful for accessing arrays, where a is the base address of the array & register holds the index value

Eg: LD RI, a(RR)

// Access the data stored in  
RI = contents(a + contents(RR))

(iii) Indexed A/M where memory locat<sup>n</sup> is integer

It is same as previous one except that a memory locat<sup>n</sup> is identified as integer.

Eg: LD RI, 100(RR)

// RI = contents(100 + contents(RR))

(iv) Indirect A/M: Contents of the data can be accessed by differencing using \* operators as shown below:

LD R1, \*(R2)

// R2 contains memory loc<sup>n</sup>  
the data stored in that  
memory loc<sup>n</sup> is copied in  
register R1

LD R1, \*100(R2)

// R1 = contents(contents(100+  
contents(R

(v) Immediate A/M: The data to be manipulated is directly present in the instruction & preceded by

LD R1, #100

// R1  $\leftarrow 100$

### EXERCISE:

code for

1. Generate 3 address statement for  $x = y - z$

LD R1, y // R1 = y

LD RR, z // RR = z

ADD R1, R1, RR // R1 = R1 + RR

ST x, R1 // ~~R1~~ x = R1

code for

2. Generate 3 address statement  $x = *p$

LD R1, p // R1  $\leftarrow p$

LD RR, 0(R1) // RR = contents(0 + contents(R1))

ST x, RR // x = RR

(iv) Indirect A/M: Contents of the data can be accessed by differencing using \* operators as shown below:

LD R1, \*(R2)

// R2 contains memory loc<sup>n</sup>  
the data stored in that  
memory loc<sup>n</sup> is copied in  
register R1

LD R1, \*100(R2)

// R1 = contents(contents(100+  
contents(R

(v) Immediate A/M: The data to be manipulated is directly present in the instruction & preceded by

LD R1, #100

// R1  $\leftarrow 100$

### EXERCISE:

code for

1. Generate 3 address statement for  $x = y - z$

LD R1, y // R1 = y

LD RR, z // RR = z

ADD R1, R1, RR // R1 = R1 + RR

ST x, R1 // ~~R1~~ x = R1

code for

2. Generate 3 address statement  $x = *p$

LD R1, p // R1  $\leftarrow p$

LD RR, 0(R1) // RR = contents(0 + contents(R1))

ST x, RR // x = RR

3. Generate code for 3 address statement  $*p = y$

LD R1, p //  $R1 = p$   
LD RR, y //  $RR = y$   
ST O(R1), RR // contents(O + contents(R1)) = RR

4. Generate m/c code for 3 address statement  $b = a[i]$

LD R1, i //  $R1 = i$   
MUL R1, R1, 8 //  $R1 = R1 * 8$   
LD RR, a[R1] //  $RR = \text{contents}(a + \text{contents}(R1))$   
ST b, RR //  $b = RR$

5. Generate m/c code for 3 address statement  $a[j] = c$

LD R1, j //  $R1 = j$   
LD RR, c. //  $RR = c$   
MUL R1, R1, 8 //  $R1 = R1 * c$   
ST a[R1], RR //  $\text{contents}(a + \text{contents}(R1)) = RR$

6. Generate m/c code for 3 address statement  
if  $x < y$  goto L

LD R1, x //  $R1 = x$   
LD RR, y //  $RR = y$   
SUB R1, R1, RR //  $R1 = R1 - RR$   
BUTZ R1, M // if  $R1 < 0$  jump to M

## Program & Instruction Cost

\* For simplicity we take the cost of an instruction to be one plus the costs associated with the addressing modes of the operands.

- \* A/M involves registers have zero additional cost.
- \* A/M involving memory locat<sup>n</sup> or constant have additional cost of 1.

\* For example:

- a) LD R0, RI  $\rightarrow$  cost = 1
- b) LD R0, M  $\rightarrow$  cost = 2.
- c) LD RI, \*100(RR)  $\rightarrow$  cost = 3

\* Cost of Addressing mode:

Mode A/M	Form	Address	Added Cost
① Absolute direct A/M	M	M	1
② Register direct A/M	R	R	0
③ Indexed A/M	C(R)	C + contents(R)	1
④ Indirect register A/M	*R	contents(R)	0
⑤ Indirect indexed A/M	*C(R)	contents(C + contents(R))	1
⑥ Immediate A/M	#C	N/A	1

NOTE: Cost of each statement = 1 + cost (Addressing mode)

## EXERCISES (8.2)

1. Determine the costs of the following instruction sequence

LD R0, Y	Cost = 1 + cost(A.M)
LD RI, Y	Cost = 1 + 1 = 2
ADD R0, R0, RI	Cost = 1 + 1 = 2
ST x, R0	Cost = 1 + 0 = 1
	Cost = 1 + 1 = 2
	<u>Total Cost = 7</u>

2. LD R0, i

MUL R0, R0, 8

LD RI, a(R0)

ST b, RI

LD R0, i	Cost = Cost (A.M) + 1.
MUL R0, R0, 8	Cost = 1 + 1 = 2
LD RI, a(R0)	Cost = 1 + 1 = 2
ST b, RI	Cost = 1 + 1 = 2
	<u>Total cost = 8</u>

3. LD R0, C

LD RI, i

MUL RI, RI, 8

ST a(RI), R0

LD R0, C	Cost = cost (A.M) + 1
LD RI, i	1 + 1 = 2
MUL RI, RI, 8	1 + 1 = 2
ST a(RI), R0	1 + 1 = 2
	<u>Total cost = 8</u>

4. LD R0, P  
 LD RI, O(R0)  
 ST x, RI

$$\text{Cost} = \text{Cost(A.M)} + 1$$

~~S = 1~~  
~~R = 1~~  
~~L = 0~~  
~~C = 1~~

LD R0, P  
 LD RI, O(R0)  
 ST x, RI

$$1 + 1 = 2$$

$$1 + 1 = 2$$

$$1 + 1 = 2$$

$$\text{Total Cost} = 6$$

5. LD R0, P  
 LD RI, x  
 ST O(R0), RI

$$\text{Cost} = \text{Cost(A.M)} + 1$$

~~S = 1 + 1 = 2~~  
~~R = 1 + 1 = 2~~  
~~L = 1 + 1 = 2~~  
~~C = 1 + 1 = 2~~

LD R0, P  
 LD RI, x  
 ST O(RI), RI

$$1 + 1 = 2$$

$$1 + 1 = 2$$

$$\text{Total Cost} = 6$$

6. LD R0, x  
 LD RI, y  
 SUB R0, R0, RI  
 BLTR \*R3, R0

$$\text{Cost} = 1 + \text{Cost(A.M)}$$

~~S = 1 + 1 = 2~~  
~~R = 1 + 1 = 2~~  
~~L = 1 + 1 = 2~~  
~~C = 1 + 1 = 2~~

LD R0, x  
 LD RI, y  
 SUB R0, R0, RI  
 BLTR \*R3, R0

$$1 + 1 = 2$$

$$1 + 1 = 2$$

$$1 + 0 = 1$$

$$1 + 1 = 2$$

$$\text{Total Cost} = 7$$

<:: indirect A