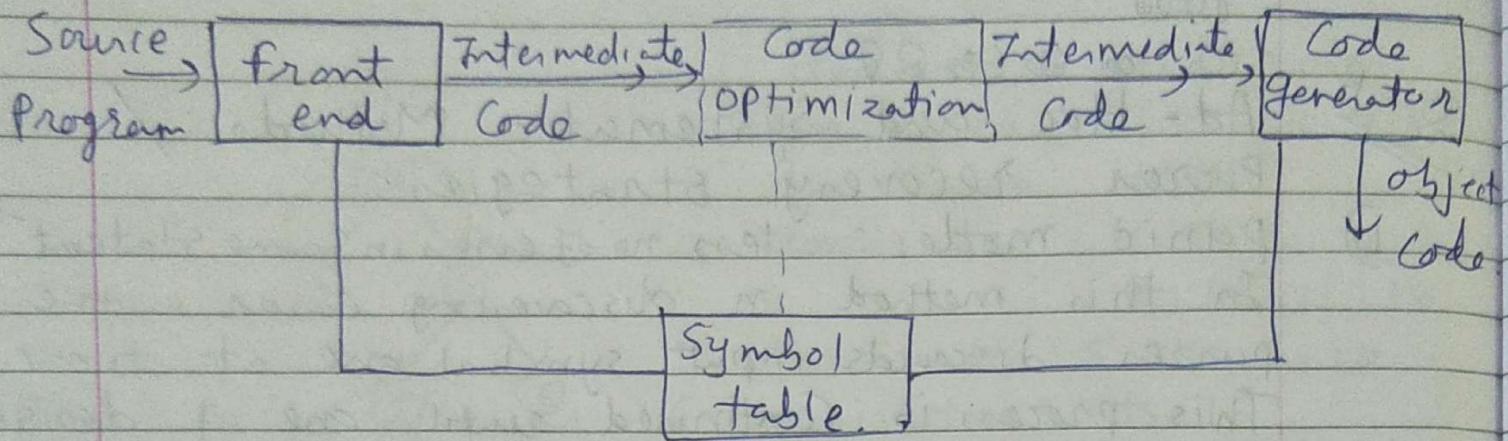


Code generation:



→ Code generation is the final phase in the process of compilation. It takes Intermediate Code as an input and generates target machine Code as output. The position of Code generator in compilation process is illustrated by following figure.

Properties of Code generation:

- Correctness
- High quality
- Efficient use of resources of the target machine
- Quick code generation.

Q-B issues in Code generation

- 1 Input to Common Code generator
- 2 Target programs.
- 3 Memory Management
- 4 Instruction Selection.



- 5 Register allocation.
 6 Choice of evaluation order.
 7 Approaches to code generation.

1) Input to the code generator:

- Code generation phase takes intermediate code as input.
- This Intermediate Code along with the symbol table information is used to determine the runtime addressing of the data objects.

↓
 denoted by names
 in TR
- The intermediate code may be in any form: three address code, quadruple, POSIX or syntax tree or DAG.
- Intermediate code generated by front end should be such that target machine can easily manipulate it.
- In the front end: necessary type checking and type conversion needs to be done. The detection of the semantic errors should be done before submitting the input to the code generator. The code generation phase requires complete error free intermediate code as input.

2) Target Program:

output of code generator: target code.

the target code comes in three forms such as:
absolute machine language.
relocatable machine language
assembly language.

advantage of producing target code in absolute machine code:

it can be placed directly at the fixed memory location and then can be executed immediately.

→ small programs can be quickly compiled.

ex: The WATFIV and PL/C are the compilers which produce the absolute code as output.

advantage of producing the relocatable machine code as output is that subroutines can be compiled separately → linked and loaded for execution by linking loader. If the target machine can not handle the relocation automatically then the compiler must provide explicit relocation information to the loader in order to link the separately compiled subroutine segments.

advantage of producing assembly code as output makes the code generation process easier.

The symbolic instructions and Macro facilities of assembler can be used to generate the code.

advantages to have assembly language as output for machines with small memory.

3) Memory Management :

Both the front end and code generator performs the task of mapping the names in the source program to addresses to the data objects in run time memory.

→ names used in three address code.

↳ refers to entry in the symbol table.

The type in declaration statement → determines the amount of storage needed to store the declared identifier.

→ if three address code contains the labels then those labels can be converted into equivalent memory addresses.

for instance if a reference to goto j is encountered in three address code then appropriate jump instruction can be generated by computing the memory address for label j.

→ 4) Instruction Selection :

The uniformity and completeness of instruction set is an important factor for code generator.

→ two important factors in selection of instructions:

 
 Speed of instruction machine idioms.

If we do not consider the efficiency of target code then the instruction selection becomes a straightforward task.

Ex. $x := a + b$

Mov a, R₀
Add b, R₀
Mov R₀, x

In the above example the code is generated line by line.

→ Such a line by line code generation process generates the poor code because redundancy can be achieved by subsequent lines ~~can~~

for ex :

$x := y + z$
 $a := x + t$

Code for above statements :

Mov y, R₀
Add z, R₀ → y + z
Mov R₀, x
Mov x, R₀
Add t, R₀ → y + z + t
Mov R₀, a.



efficient code :

Mov y, R₀
Add z, R₀
Add t, R₀ → y + z + t
Mov R₀, a.

5) Register allocation :

- If the instruction contains register operands then such a use becomes shorter and faster than that of using operands in the memory.
- There are two important activities done while using registers.

1. Register allocation : During register allocation select appropriate set of variables that will reside in registers.

2. Register assignment : During register assignment, pick up the specific register in which corresponding variable will reside.

Ex: Consider the three address code:

$$t_1 := a + b$$

$$t_1 := t_1 * c$$

$$t_1 := t_1 / d$$

Machine Code :

MOV a, R0

ADD b, R0

MUL c, R0

DIV d, R0

MOV R0, t1

6) choice of evaluation order :

The evaluation order is an important factor.

Some order requires less number of registers to hold the intermediate results than others. Picking up the best order is one of the difficulties in code generation.

~~few~~ Avoid this problem by referring the order in which three address code is generated by semantic rules.

Approaches to code generation:

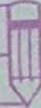
→ The most important factor for a code generation is that it should produce the correct code. With this approach of code generation various algorithms for generating code are designed.

Q-B → DAG Representation of basic blocks:

Directed acyclic graphs are useful data structures for implementing transformations on basic blocks.

A DAG for basic block is directed acyclic graph with the following labels on nodes:

- Leaves are labeled by unique identifiers, either variable names or constants. From the operator applied to a name we determine whether the l-value or r-value of a name is needed; most leaves represent r-values. The leaves represent initial values of names and we



Subscript them with o to avoid confusion with labels denoting current values of names as in (3) below.

- 2 Interior nodes are labeled by an operator symbol.
- 3 Nodes are also optionally given a sequence of identifiers for labels. The intention is that interior nodes represent computed values and the identifiers labeling a node are deemed to have that value.

Three address code for Block B₂.

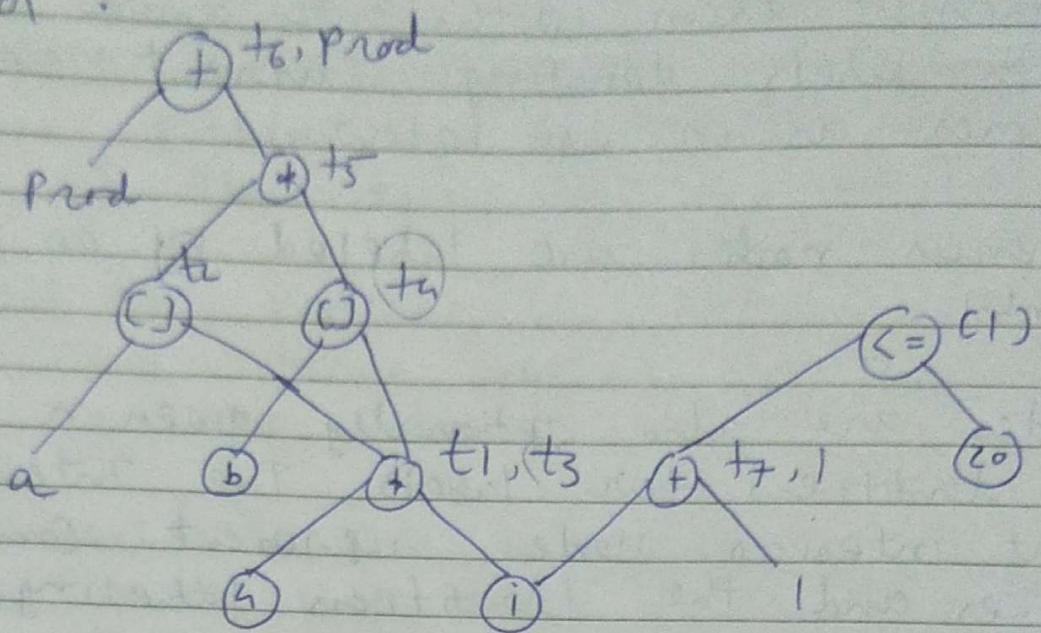
<sup>Basic
block</sup>
2 $t_1 := 4 * i$
 $t_2 := a[t_1]$
 $t_3 := 4 * i$
 $t_4 := b[t_3]$
 $t_5 := t_2 + t_4$
 $t_6 := \text{prod} + t_5$
 $\text{prod} := t_6$
 $t_7 := i + 1$
 $i = t_7$
 $\text{if } i \leq 20 \text{ goto C1}$.

<sup>Basic
block</sup>
① prod = 0
 $i = 1$

Prepared by kinjal patel

internal node : is any node of tree that has child node and is thus not leaf node
 Intermediate node between root and leaf : Internal node
 Date: _____
 Page No. _____

DATA :



→ A dag for a basic block has following labels on the nodes:

leaves are labeled by unique identifiers, either Variable names or Constants.

→ Interior nodes are labeled by an operator symbol.

→ Nodes are also optionally given a sequence of identifiers for labels.

Algorithm for Construction of DAOT:

Case (i) $x := y \text{ op } z$

Case (ii) $x := \text{op } y$

Case (iii) $x := y$

With the help of following steps DAOT can be constructed :

Step: 1 If y is undefined then create node(y). Similarly if z is undefined create a node(z)

Step: 2 For the case(i) create a node(op) whose left child is node(y) and node(z) will be the right child. Also check for any Common Subexpression. For the case(ii) determine whether is a node labeled op , such node will have a child node(y). In case(iii) node n will be node(y).

Step: 3 Delete x from list of identifiers for node(n). Append x to list of attached identifiers for node n found in 2.

Application of DAOT:

DAOT is used in:

- 1 Determining the Common Sub expression
- 2 Determining which names are used inside the block and computed outside the block

- 3 Determining which statements of the block could have their computed value outside the block.
- 4 Simplifying the list of quadruples by eliminating the common sub expressions and not performing the assignment of the form $x := y$ unless and until it is a must.

→ Peephole Optimization:

→ Definition:

Peephole optimization is a simple and effective technique for locally improving target code. This technique is applied to improve the performance of the target program by examining the short sequence of target instructions and replacing these instructions by shorter or faster sequence.

Peephole: a small window is moved over the target code and transformations can be made. And hence is the name.

Characteristics of peephole optimization can be applied on the target code using following characteristics:

- 1 Redundant instruction elimination
- 2 flow of control optimization
- 3 Algebraic Simplification
- 4 Use of machine idioms
- 5 Strength reduction.

1 Redundant instruction elimination:
 Especially the redundant loads and stores can be eliminated in this types of transformations:

for example :

Mov R₀, X
 Mov X, R₀

→ We can eliminate the second instruction since X is already in R₀. But if (Mov X, R₀) is a label statement then we cannot remove it.

We can eliminate the unreachable instructions for example, following is a piece of C code.

Sum = 0;
 if (Sum == 1)
 {Printf("1/d", sum);}

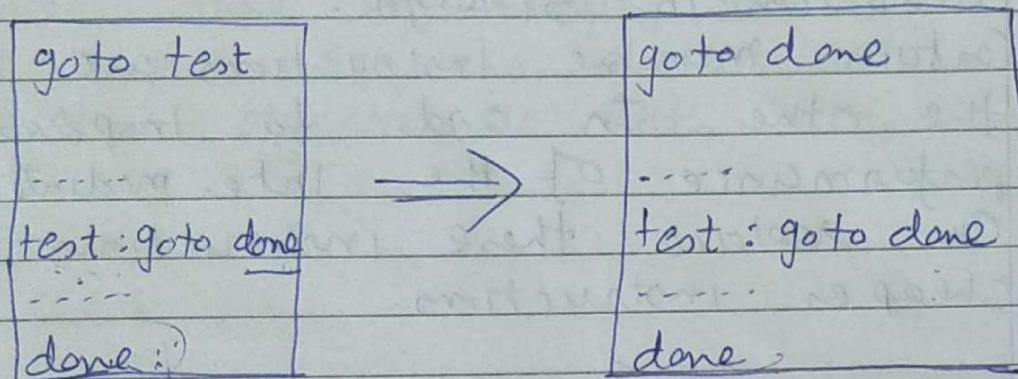
Now this if statement will never get executed hence we can eliminate such a unreachable code.

```
int func(int a, int b)
{
    c = a + b;
    return c;
    printf("%d", c);
}
```

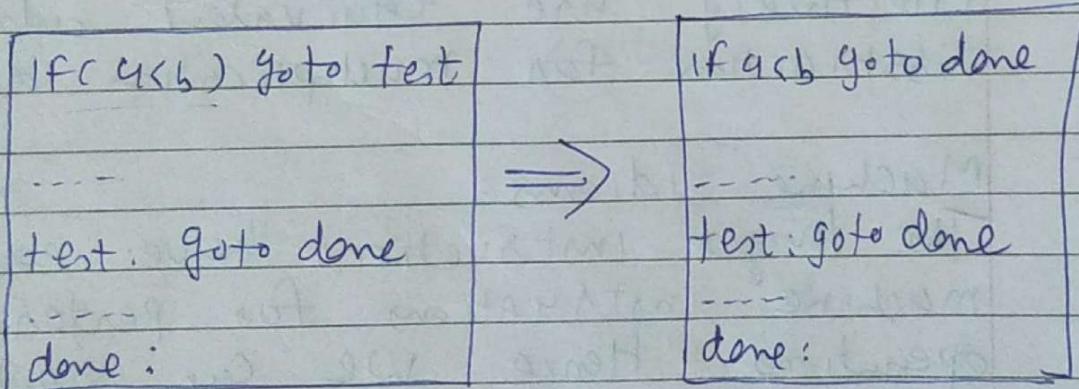
2 flow of control optimization:

Using peephole optimization unnecessary jumps can be eliminated.

for ex:



→ Thus we reduce one jump by this transformation:



3

Algebraic Simplification:
Peephole optimization is an effective technique for algebraic simplification

The statement such as:

$$x := x + 0$$

or

$$x := x * 1$$

Can be eliminated by peephole optimization

4

Reduction in strength:

Certain machine instructions are cheaper than the other. In order to improve the performance of the intermediate code we can replace these instructions by equivalent cheaper instructions.

Ex: x^2 is cheaper than $x + x$.

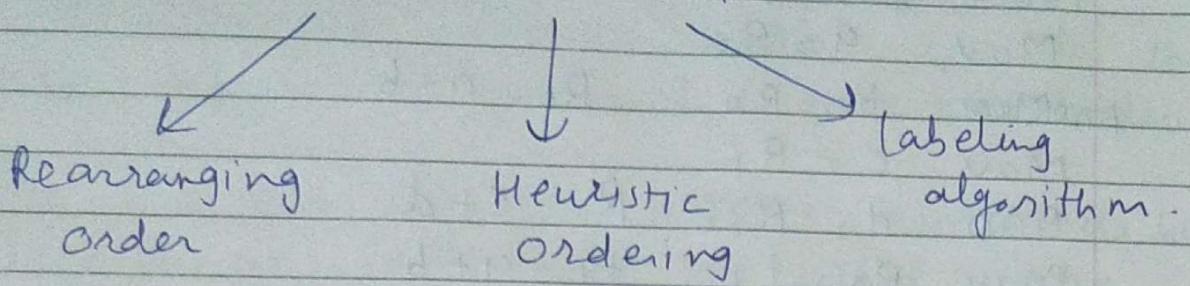
Addition and Subtraction is cheaper than multiplication and division. So, we can effectively use equivalent addition and subtraction for multiplication division

5

Machine Idioms:

The target instructions have equivalent machine instructions for performing some operations. Hence we can replace these instructions by equivalent machine instruction in order to improve the efficiency. for ex:

generating Code from DAt:



I Rearranging order :

The order of three address code affects the cost of the object code being generated. In the sense that by changing the order in which computations are done we can obtain the object code with minimum cost.

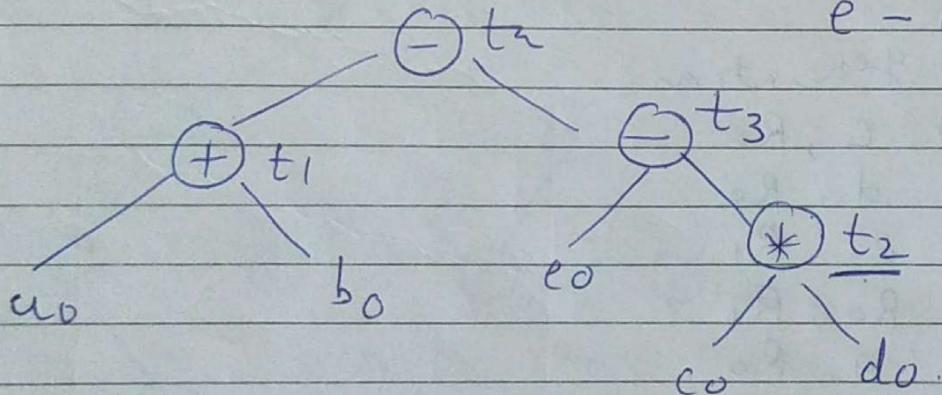
for example :

$$t_1 := a + b$$

$$t_2 := c + d$$

$$t_3 := e \bar{+} t_2$$

$t_4 := t_1 \bar{*} t_3$ for the expression $(a+b) - (e - (c + d))$



Code generation :

Mov q, R₀
 AND R₀, b, R₀ R₀: q + b
 Mov c, R₁
 ADD d, R₁ R₁: c + d
 Mov R₀, t₁ t₁: q + b
 Mov e, R₀ R₀: e
 SUB R₁, R₀ R₀: ~~e~~ - c + d
 Mov t₁, R₁ R₁: q + b
 SUB R₀, R₁ R₁: q + b - ~~e~~ - (c + d)
 Mov R₁, t₂ t₂: q + b - ~~e~~ - (c + d)

Suppose we rearranged the order of the statements so that the computation of t₁ occurs immediately before that of t₂ is:

t₂ := c + d
 t₃ := e - t₂
 t₁ = q + b
 t₄ := t₁ - t₃

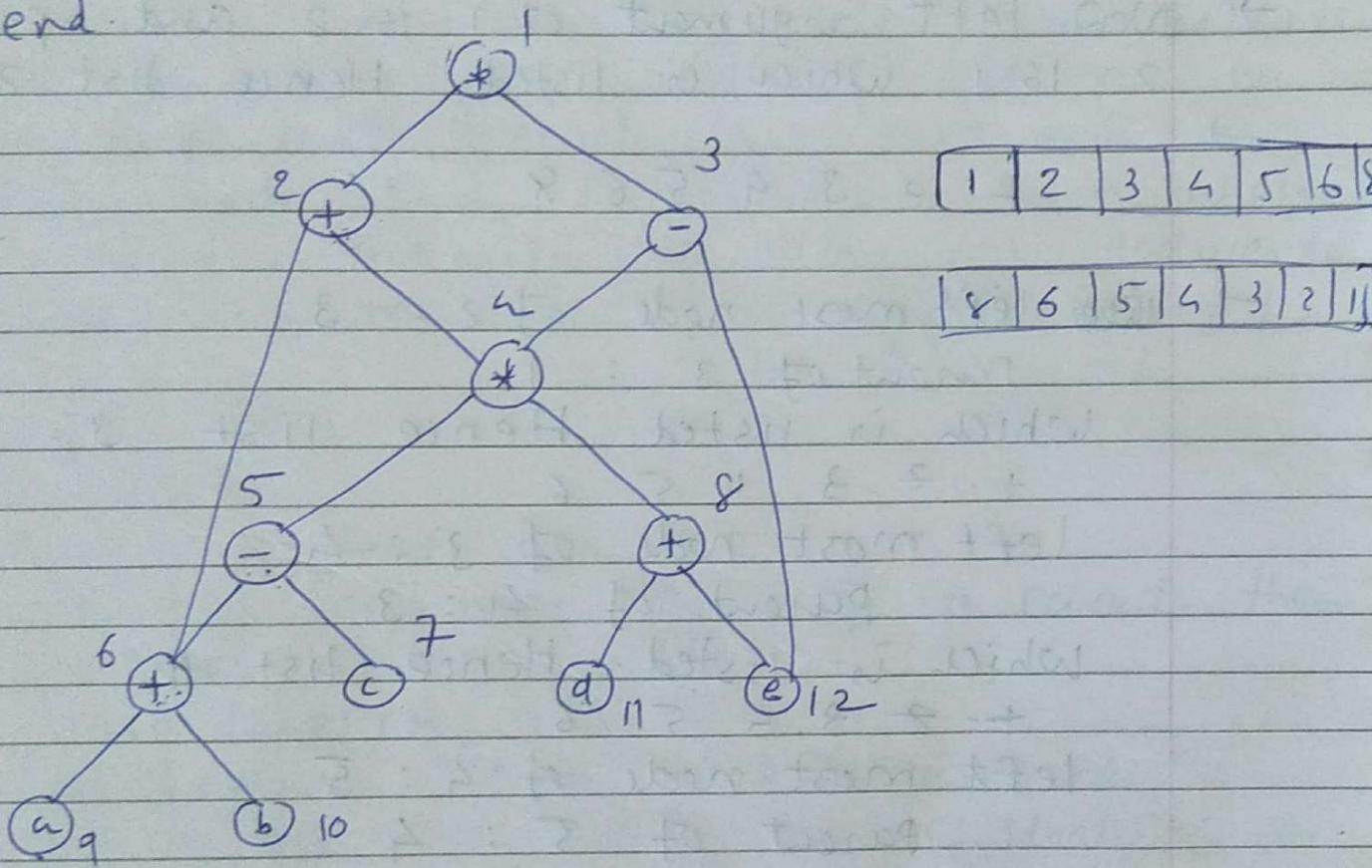
Code generation:

Mov c, R₀
 ADD d, R₀
 Mov e, R₁
 SUB R₀, R₁
 Mov q, R₀
 ADD b, R₀
 SUB R₁, R₀
 Mov R₀, t₂

Heuristic Ordering:

Obtain all the interior nodes. Consider those interior nodes as unlisted interior nodes.

- (1) While unlisted interior nodes remain do begin
 - (2) Select an unlisted node n , all of whose parents have been listed;
 - (3) list n ;
 - (4) while the left most child of m of n has no unlisted parents and is not a leaf do begin
 - (5) list m ;
 - (6) $n := m$
- end
end.



$$\begin{aligned}
 t_8 &:= d + e \\
 t_6 &:= a + b \\
 t_5 &:= t_6 - c \\
 t_3 &:= t_5 * t_8 \\
 t_2 &:= t_3 - e \\
 t_1 &:= t_2 * t_3
 \end{aligned}$$

→ Consider the unlisted interior nodes 1 2 3 4 5 6 8

→ Initially the only node with unlisted parent is 1.

→ Now left argument of 1 is 2 and parent of 2 is 1 which is listed. Hence list 2.

→

+ 2 3 4 5 6 8

→ left most node of 2 : 3

parent of 3 : 1

which is listed. Hence list 3.

+ 2 3 4 5 6

left most node of 3 : 4

parent of 4 : 3

which is listed. Hence list 4.

+ 2 3 4 5 6

left most node of 4 : 5

parent of 5 : 4

which is listed. Hence list 5

only node 8 is remaining from the unlisted interior nodes we will list it reversing the list : 8 6 5 4 3 2 1.



Labeling algorithm:

We use the term "left leaf" to mean a node that is leaf and the left most descendent of its parent. All other leaves are referred to as right leaves.

labeling can be done by visiting node in a bottom up order so that a node is not visited until its children all labeled. The order in which the first three nodes are created is suitable if the parse tree is used as intermediate code. So in this case the labels can be computed as syntax directed translation. The algorithm below is for computing the label at node n . In the important special case that n is a binary node and its children have labels l_1 and l_2 , formula of line (6) reduces to

$$\text{label}(n) = \max(l_1, l_2) \text{ IF } l_1 \neq l_2$$

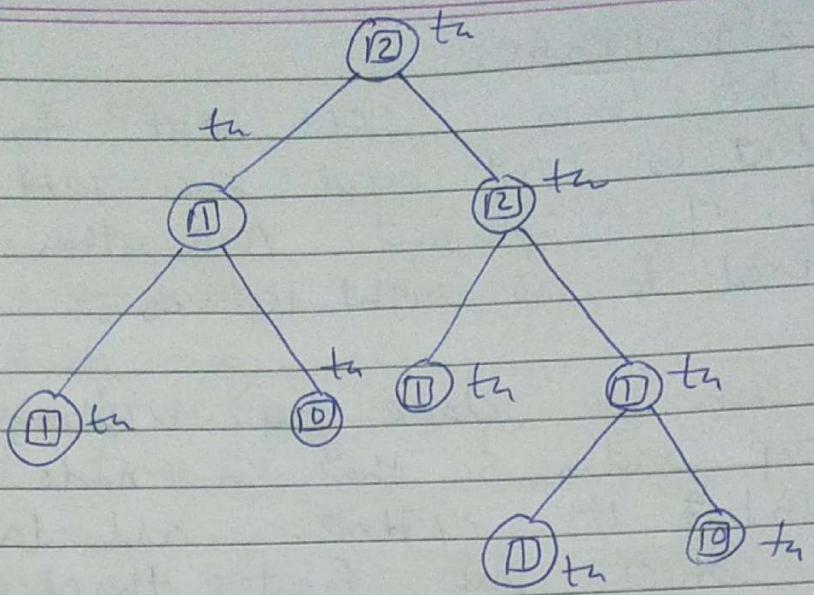
$$l_1 + 1 \text{ IF } l_1 = l_2$$

Algorithm :

```

1 if n is a leaf then
2   if n is left most child of its parents then
3     label(n) = 1
4   else label(n) = 0
5 else begin /* n is interior node */
6   let  $n_1, n_2, \dots, n_k$  be the children of n
   ordered by label.
7   so  $\text{label}(n_1) \geq \text{label}(n_2) \geq \dots \geq \text{label}(n_k)$ 
8    $\text{label}(n) = \max(\text{label}(n_i) + i - 1)$ 
9 end.

```



Prepared by kinjal patel