# Unit - V

## Code Generation

→ The final phase in our compiler model is the code generator. It tak[es]

as i/p an intermediate representation of the source program & produce[s]
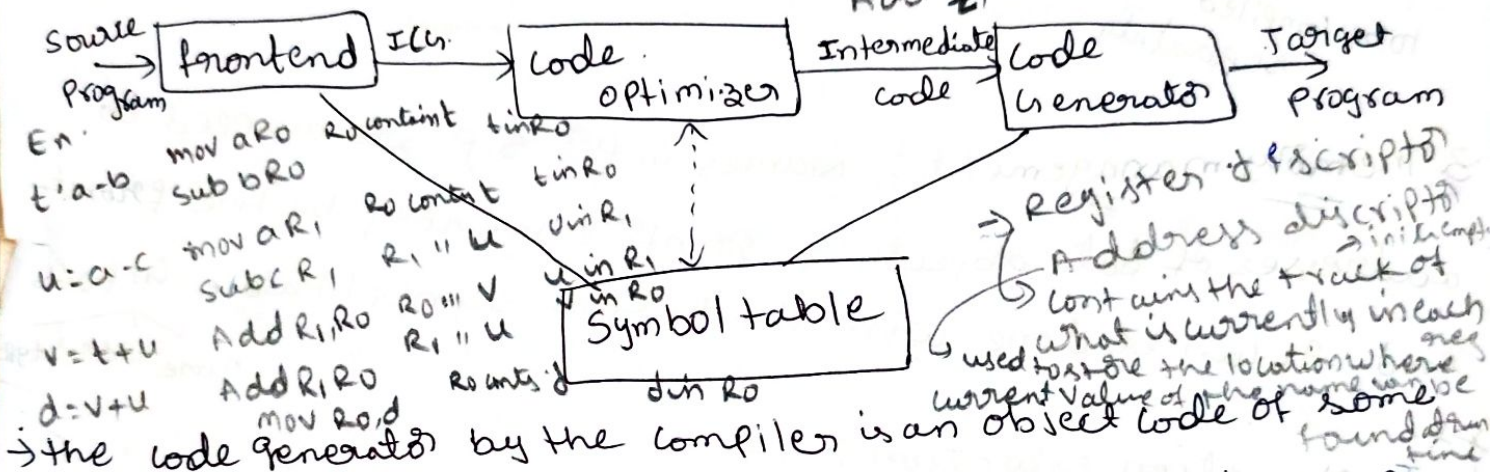
as o/p an equivalent target program (Assembly)

→ The requirements traditionally imposed on a code generator are

* The o/p code must be correct & of high quality.
  → it should worry the exact meaning of the s.c.

  * meaning that it should make effective use of the resources

  of the target machine. / It should worry the Exact meaning of the source code.
  → It should be efficient in terms of CPU usage & memory management.

  * the code generator itself should run efficiently. (using descri[p]tors)

  Ex:- $x := y + z$.   mov y, R0   mov R0, x
                       Add z, R0

Source → [ frontend ] ICu → [ Code optimizer ] → Intermediate code → [ Code Generator ] → Target Program

Ex.
t := a-b   mov a R0   R0 contains   t in R0
u := a-c   sub b R0   R0 contains   t in R0
           mov a R1   R1 " u       u in R1
v := t+u   subc R1                  u in R1
           Add R1,R0  R0 " v       v in R0
d := v+u   Add R1,R0  R1 " u       
           mov R0,d   R0 cnts d    d in R0

→ Register & script[or]
→ Address discrip[tor]
  → in Lcomp[?]
  contains the track of what is currently in each reg
used to store the location where current value of the name can be found dyn time

→ the code generated by the compiler is an object code of some

lower-level programming language, for Ex: Assembly language.

## Issues in the design of code generator :-

→ The important criterion is that it produce correct code. The design

goals are.   * correctness   * easily implemented   * tested & maintain[ed]

* In the code generation phase, various issues can arises.

1. Input to the code Generator :-

It consists of the intermediate representation of the source

program produced by front end, together with info in the S.T

to determine run-time addresses of the data objects denoted

by the names in the intermediate representation.
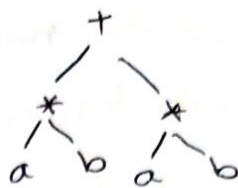
→ Intermediate representation can be

1. Postfix $abt$ — triple
2. 3-address code — quadruple, Indirect triple
3. syntax tree / dags.

* Frontend produces low-level intermediate representation directly. manipulated by the machine instructions.

$a*b+a*b$



* The code generator phase needs complete error free Intermediate code as an i/p.

2. Target program :- The O/P of the CG is the target program

The O/P may be directly placed in a fixed memory. It allows subprograms to be compiled separately.

* Absolute machine code Language → .exe
* Relocatable machine Language , ob's (Linker)
* Assembly Language , Nemonics
  · asm (loader)

3. memory management :- Names in the S.P are mapped to addresses of data objects in runtime memory by the front-end. It can be done with the help of symbol table info's

name address type

** 
4. Instruction selection :-

→ The instruction set of the target machine should be complete & uniform.

→ When you consider the efficiency of target machine then the instruction speed & machine idioms are important factors

→ The quality of the generated code can be determined by its speed & size.

Eg₂: $x:=x+1$

MOV x, Ro
INC x { ADD #1, Ro
MOV Ro, x.

Eg₃:- $x:=y+z, w=x+s$

MOV y, Ro
ADD z, Ro
MOV Ro, x
MOV x, Ro — this can be eliminated
ADD s, Ro
MOV Ro, w

Eg₁:- $a:=b+c$

MOV b, Ro    Ro→b
ADD c, Ro
MOV Ro, a    a→Ro

**5. Register allocation:-** instructions involving register operands are shorter & faster than those involving operands in memory.

→ The following subproblems arise when we use registers.

* **Register allocation:** set of var's that reside in registers/. (allocation of the memory to the variables)

* **Register Assignment:-** The specific register that a variable will reside is selected / we have to select the specific register for a particular variable. (Assigns the values to the variable)

Ex:-
$$t = x + y$$
$$t = t * z$$
$$t = t / w$$

⇒
| | |
|---|---|
| MOV x, R0 | R0 → x |
| ADD y, R0 | R0 + y |
| MUL z, R0 | R0 * z |
| DIV w, R0 | R0 / w |
| mov R0, t | t → R0 |

⇒ Here we didn't allocate Reg. for temp. variable & R0 is a single reg is used to perform all operations.

**6. Evaluation order:-** The code generator decides the order in which the instruction will be Executed.. The order in which computations are performed can effect the efficiency of the target code.

Ex:- $a + b - (c + d) * e$

| Three-address code | Code | Reordered three-address code | Code | |
|---|---|---|---|---|
| $t_1 := a + b$ | 1. MOV a, R0 | $t_2 := c + d$ | 1. MOV c, R0 | The needed |
| $t_2 := c + d$ | 2. ADD b, R0 | $t_3 := t_2 * e$ | 2. ADD d, R0 | Inst's reduce |
| $t_3 := t_2 * e$ | 3. MOV R0, t_1  R0free | $t_1 := a + b$ | 3. MOV e, R1 | The no. of |
| $t_4 := t_1 - t_3$ | 4. MOV c, R1 | $t_4 := t_1 - t_3$ | 4. MUL R0, R1 | Final code by |
| | 5. ADD d, R1 | | 5. MOV a, R0 | 2. Thus |
| | 6. MOV e, R0 | | 6. ADD b, R0 | saved in |
| | 7. MUL R1, R0  mul:R0 | | 7. SUB R1, R0 | cost. |
| | 8. MOV t₁, R₁ | | 8. MOV R0, t_4 | |
| | 9. SUB R0, R1 | | | |
| | 10. MOV R4, t4 | | | |

# Generic Code Generation Algorithm :-

It generates target code for a sequence of instructions.

→ Code generation algorithm considers each three-address instruction in turn and decides what loads are necessary to get the needed operands into registers.

→ After generating the loads, it generates the operation itself, then, if there is a need to store the result into a memory location, it also generates that store.

→ The code generator has to track both the registers & addresses while generating the code, for both <u>availability</u> & <u>location of</u> values the following two descriptors are used.

1. Register descriptor :- for each available register, a register descriptor keeps track of the variable names whose current value is in that register. initially all registers are empty. These registers are used for local use within a basic block.

2. Address descriptor :- for each program variable, an address descriptor track of the location where the current value of the name can be found. It may be a <u>register</u> /stack location/ memory address.

Ex:- $d = (a-b) + (a-c) + (a-c)$

$t_1 = a - b$
$t_2 = a - c$
$t_3 = t_1 + t_2$
$d = t_3 + t_2$

$t_u = 4 * j$

$t_u = u * j$

$t_u = t_u + 4$

$\underset{y_j}{\underbrace{}}$

$\phi = 2 * a$

$\lambda = a + a$

$j = j$
$j = j - 1$
$j = 1$
$j = 0$
$t_u = 0$

| Statements | Code Generator | Reg. Descriptor | Address descriptor |
|---|---|---|---|
| $t_1 = a - b$ | MOV a, R0 <br> Sub b, R0 | Reg's are empty <br> R0 contains $t_1$ | $t_1$ in R0 |
| $t_2 = a - c$ | MOV a, R1 <br> Sub c, R1 | R0 contains $t_1$ <br> R1 contains $t_2$ | $t_1$ in R0 <br> $t_2$ in R1 |
| $t_3 = t_1 + t_2$ | ADD R1, R0 | R0 contains $t_3$ <br> R1 contains $t_2$ | $t_2$ in R1 <br> $t_3$ in R0 |
| $d = t_3 + t_2$ | ADD R1, R0 <br> MOV R0, d | R0 contains d | d in R0 <br> d in R0 & memory |

## Generating code from DAG:
## Code Generation using DAG :-

→ Directed Acyclic Graph (DAG) is a tool that depicts the structure of basic blocks, helps to see the flow of values flowing among the basic blocks, and offers optimization too.

* leaf nodes represent identifiers, names or constants.

* Interior nodes represent operators.

* Interior nodes also represent the results of expressions

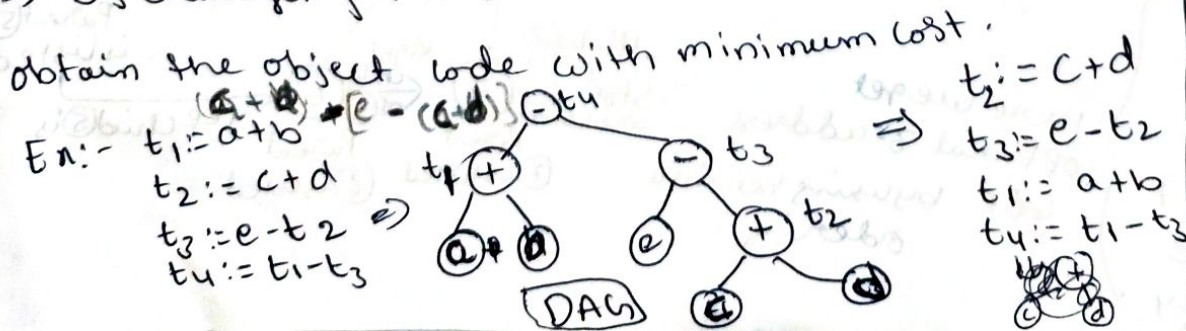→ Generating code from DAG is much simpler than the linear seq of three addr code

→ methods generating code from DAG's are :

→ using DAG we can re arrange some sequence of instructions & generate an efficient code

1. Rearranging order.   2. Heuristic ordering   3. Labeling algorithm

1. Rearranging order: The order of 3-address code effects the cost of the object code being generated.

→ By changing the order in which computations are done we can obtain the object code with minimum cost.

Ex :- 
$$t_1 := a + b$$
$$t_2 := c + d$$
$$t_3 := e - t_2$$
$$t_4 := t_1 - t_3$$

⇒
$$t_2 := c + d$$
$$t_3 := e - t_2$$
$$t_1 := a + b$$
$$t_4 := t_1 - t_3$$


DAG

| | |
|---|---|
| MOV a, R0 | MOV c, R0 |
| ADD b, R0 | ADD d, R0 |
| MOV c, R1 | ADD e, R1 |
| ADD d, R1 | ⇒ SUB R0, R |
| MOV R0, t1 | MOV a, R0 |
| MOV e, R0 | ADD b, R0 |
| SUB R1, R0 | SUB R1, R |
| MOV t1, R1 | MOV R0, t |
| SUB R0, R1 | Ass. code |
| MOV R1, | |

Simple , n . .
+ euristic ordering :-

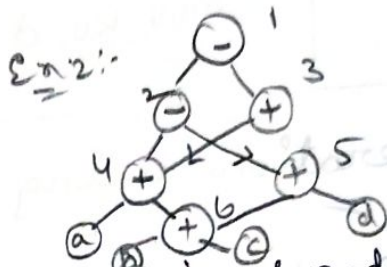obtain all the interior nodes. consider these interior nodes as unlisted nodes.

Algorithm :-
Algm H-DAG()
While (unlisted interior nodes remain) do
{
  Pickup an unlisted node n, whose parents have been listed.
                                                                    n root.
  List n;
  while ( the left most child m of 'n' has no unlisted parent
                                                AND is not leaf)
  {
    list m;
    m:= m;
  }
}

order = reverse of the order of listing of nodes;

→ Given the DAG first numbered (from top to bottom
  & from left to right then consider the unlisted
                                          interior nodes;

Interior nodes = 1, 2, 3, 4, 5, 6, 7
unlisted nodes = 1, 2, 3, 4, 5, 6, 7.

(1)
Pickup an unlisted node, whose parents have been listed

$\sum n \geq 2:-$

$t_6 = b + c$
$t_5 = t_6 + b$
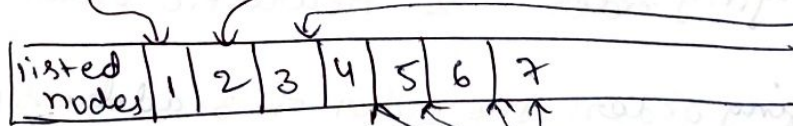$t_4 = a + t_6$
$t_3 = t_4 + t_5$
$t_2 = t_4 - t_5$
$t_1 = t_2 - t_3$



Ea:-

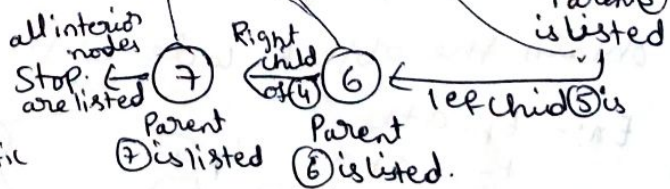Reverse order of listing nodes: 7, 6, 5, 4, 3, 2, 1

$t_7 = d + e$
$t_6 = a + b$
$t_5 = t_6 - c$
$t_4 = t_5 * t_7$
$t_3 = t_4 - e$
$t_2 = t_6 + t_4$

Here we get
optimal 3-address
code by using heuristic
ordering.

(1)  leftmost   (2)        leftmost  (6)      move      (3)
root   →       Parent(2)is(1)  →      Parent    right    Parent(3)
noParent(1)is(2)  listed ao   (1)is(6)  (5)isnot  child(1)  is listed  R
list(1)                               listed
                                      not listed(6)           left
                                                             ↓ child
                                                             (4)
                                                         Parent(3)u
listed nodes | 1 | 2 | 3 | 4 | 5 | 6 | 7                  listed.

                                                             left
                                                             ↓ child
                                                             (4)

                                                             (5)
                                                         Parent(5)
all interior                                             is listed
Stop:  ← (7)      Right                                   ↓
are listed        child
              Parent  of(4) (6) ←    left child(5)is
              (7)is listed  Parent
                           (6)is listed.

3. **Labeling Algorithm:-** The labeling algorithm generates the optimal code for given expression in which minimum registers are required.

⟶ In this labelling algorithm is used to find the min. no. of reg used for Executing the code

⟶ Using labeling algorithm the labeling can be done to tree by visiting nodes in bottom up order. ( leaf node ⟶ interior ⟶ root node)
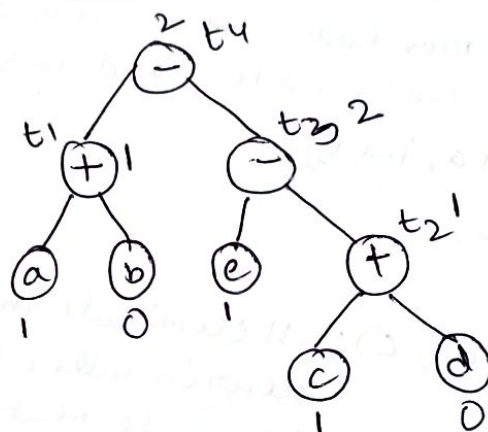
⟶ for computing the label at node 'n' with the label $L_1$ to left child then the value is '1'. and label $L_2$ to right child then the value is '0'.

$$Label (n) = max (L_1, L_2) \text{ if } L_1 \text{ not equal to } L_2$$

$$L_1 \neq L_2$$

$$Label (n) = L_1 + 1 \text{ if } L_1 = L_2$$

3-address code

Eg:- $t_1 := a + b$
$t_2 := c + d$
$t_3 := e - t_2$
$t_4 := t_1 - t_3$

(L-R-Root)
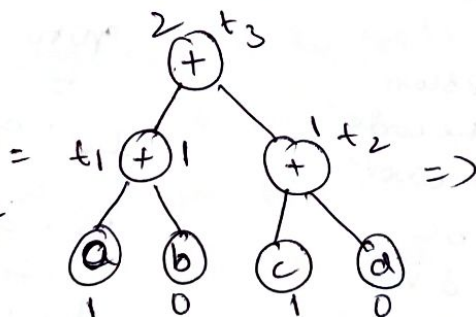Post order traversal : $a\ b\ t_1\ e\ c\ d\ t_2\ t_3\ t_4$

```
MOV   a, R0
ADD   b, R0
MOV   R0, t1
MOV   C, R1
ADD   d, R1
MOV   R1, t2
MOV   e, R0
SUB   R0, R1
ADD   t1, R0
SUB   R0, R1
MOV   R1, t4
```
two registers R0, R1

Eg2:- $t_1 = a + b$
$t_2 = c + d$
$t_3 = t_1 + t_2$

```
MOV  a, R0
ADD  b, R0
MOV  C, R1
ADD  d, R1
ADD  R0, R1
MOV  R1, t3
```

**Algorithm:-** for computing the label at node n

1. If n is a leaf then ⑥.
2. If n is the leftmost child of its parent then label(n) := 1.
3. else lable (n) := 0. else begin /* n is an interior node */
   So label (n₁) ≥ label (n₂) ≥ ... ≥ label (nₖ);
4. label(n) := max (label(nᵢ) + i - 1) where 1 ≤ i ≤ k.
   5. end

# Peephole optimization technique:-

→ This technique works locally on source code to transform
it into an optimized code.

→ It is a short seq of target instruction that can be replaced
by shorter or faster seq. instructions in a window (Peephole).

→ This tech is applied to improve the performance of prgm by Examing
a short seq of inst in a peep hole.

characteristics of Peephole optimization:-

1. Redundant instruction elimination.
2. Removal of unreachable code
3. Flow of control optimization
4. Algebraic simplifications
5. Machine idioms. (Ex: $a = a + 1 \to Inc\ a$, $a = a - 1 \to DEC\ a$)

→ Eg:- mov Ro, a     &  mov Ro, a     Here a's already in Register Ro
       mov a, Ro        so remove    2nd instruction & consider one.

→ Eg:- int add (int a, int b)
       {  int c = a + b;
          return c;
          Printf("%d", c);    // eliminate this inst. becoz once return
       }                       occur in called fn control moves to calling fn
                               & Execute next stmt in calling fn. so it's
                               not a optimized code.

→ Eg3:- Here using peephole optimization unnessary Jumps can be eliminated

• goto L1                   multiple Jumps can          goto L3

L1: goto L2                 → make the code           L1: goto L3
                            inefficient above
L2: goto L3                 code can be               L2: goto L3
                            replaced by
L3: mov a, Ro                                          L3: mov a, Ro

Ex4:- $a = a + 0$ (or) $a = a \times 1$ either, $a = a^2$ replaced with $a = a \times a$
the above stmts can be eliminated becoz by Executing those stmts
the result 'a' won't changes.

Ex5:- It is the process of using powerful features of CPU inst

    Ex:- $a = a + 1 \to Inc\ a$
         $a = a - 1 \to Dec\ a$

# Simple code Generator (or) Generic code generation algorithm:-

It generates target code for a sequence of instructions.

→ It uses a function getReg(s) to assign registers to Variables.

→ for each operator in a three address code there is a corresponding target language operator (+, -, * etc).

→ It uses 2 data structures. 1. Register descriptor
2. Address descriptor

Register descriptor:- (RD) used to keep track of which Variable is stored in a register. Initially all registers are empty. for Example if we are using Ro & R₁ for storing the values, then register descriptor tells what is currently in Ro, R₁.

→ RD maintain the value stored in register & inform the availability of registers to code generator.

Address descriptor:- used to keep track of location where the variable is stored. location may be register, memory address, etc.

→ It maintain the memory location which are used in program.

Algorithm:- [used to generate code for the single basic block]
For Each three instr address code of the form $x = y\ op\ z$, & Assumes that 'L' is the location where the o/p of $y\ op\ z$ is to be stored.

1. call function getReg(s) to gets the location of L.

2. Determine the present location of 'y' by consulting address descriptor of y. If 'y' is not present in location 'L' then generate the instruction mov y, L to copy value of y to L.

3. The present location of z is determined using step 2 and the instr is generated as OP z, L. update address descriptor of 'x' to indicate that it is stored in L; if y op z i.e assigned to 'x'. So,

4. Now 'L' contains the value of y op z i.e assigned to 'x'. So, y & z have no next uses and not live on exit, update the descriptors to remove y & z.

→ Ex:- $d = (a-b) + (a-c) + (a-c)$ → simple code Generation of code gene ating

three address code : $t_1 = a-b$

$\qquad t_2 = a-c$

$\qquad t_3 = t_1 + t_2$

$\qquad d = t_3 + t_2$

| Statement | Code Generation | Register descriptor | Address descripto |
|-----------|-----------------|---------------------|-------------------|
| $t_1 = a-b$ | mov a, Ro <br> sub b, Ro | Ro contains $t_1$ | $t_1$ in Ro |
| $t_2 = a-c$ | mov a, R₁ <br> Sub c, R₁ | Ro contains $t_1$ <br> R₁ contains $t_2$ | $t_1$ in Ro <br> $t_2$ in R₁ |
| $t_3 = t_1 + t_2$ | add R₁, Ro | Ro contains $t_3$ <br> R₁ contains $t_2$ | $t_3$ in Ro <br> $t_2$ in R₁ |
| $d = t_3 + t_2$ | add R₁, Ro <br> ~~mov Ro,d~~ | Ro contains d <br> R₁ contains $t_2$ | d in Ro <br> ↳ final result |

In this getReg get Reg is used to select register for each memory for each memory location associated with the three address stmts. This function decides the status of register & location of value name.

→ Basic block is a collection of 3 address stmts, keeps tack of register allocation & assignment.

* Register allocation :- is a process of deciding what value a register must hold.

* Register Assignment is a process of picking up a registr to store a value of variable.

Operations of get Reg():- There are 3 types of operations

1. Load operations:- LD r, x loads value in location x in register r.

2. Store operations :- ST x, r stores value in register r in location x.

3. computation operations:- Add, Sub, mul, div, Inc etc.

# Gobal register allocation:-

1. The code generation algorithm used registers to hold values for the duration of a singleblock. however, all the live variables were stored at the end of each block.

2. To save some of these stores and corresponding loads, we might arrange to assign registers to frequently used variables and keep these registers ~~registers~~ consistent across block boundries (globally).

3. Since programs spend most of their time inner loops, a natural approach to global register assignment is try to keep a frequent used value in a fined register throuout a loop.

4. one strategy for global register allocation is to assign some fined no. of registers to hold the most active values in each inner loop.

usage count → count a savings of one for each use of $x$ in loop $L$.

→ if $x$ is allocated a register, then count a savings of two for each block $L$.

→ formula for the benefit to be realised from allocating a register of $x$ allocated to a reg. to '$x$' within loop is:
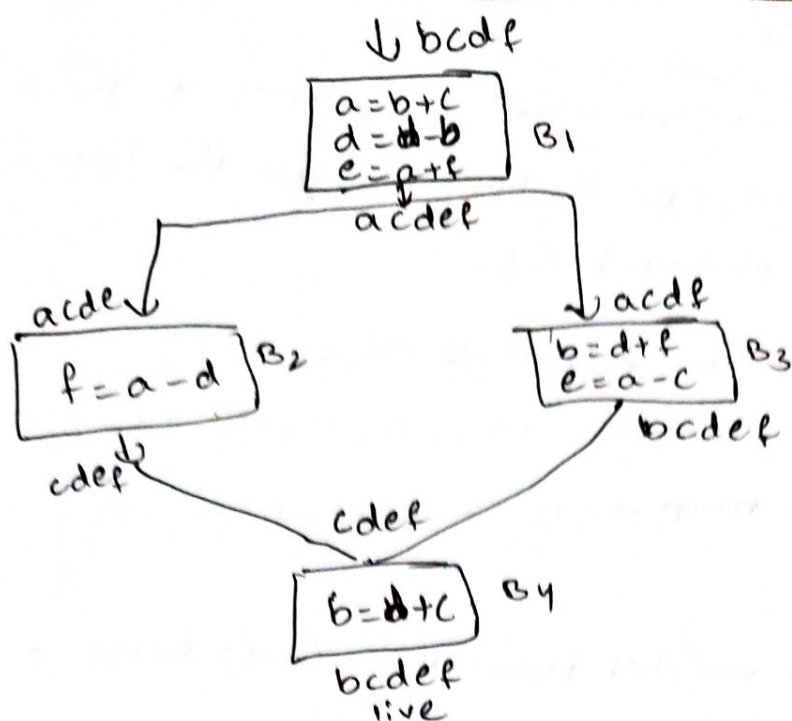
$$\sum_{blocks\ B\ in\ L} [use(x,B) + 2 * live(x,B)]$$

and is not preceded by an assignment to $x$ in the same block.

where $use(x,B)$ is the no. of times $x$ is used in $B$ prior to any definition of $x$.

$live(x,B) = 1$, if $x$ is live on Exit from $B$ and assigned a value in $B$. $live(x,B) = 0$, otherwise.

Eg:-

↓ bcdf

$$a = b + c$$
$$d = d - b$$
$$e = a + f$$
B1

acdef

acde ↓                    ↓ acdf

$$f = a - d$$  B2          $$b = d + f$$  B3
                          $$e = a - c$$

cdef ↓                    bcdef

            cdef

        $$b = d + c$$  B4

        bcdef
        live

|  | B1 | B2 | B3 | B4 | save units of cost |

$$a = (0 + 2*1 = 2) + (1 + 2*0) + (1 + 2*0) + (0 + 0) = 4$$
                        $$= 1$$           $$= 1$$

$$b = (2 + 2*0 = 2) + (0 + 2*0 = 0) + (0 + 2*1 = 2) + (0 + 2*1 = 2) = 6$$

$$c = (1 + 2*0 = 1) + (0 + 2*0 = 0) + (1 + 2*0 = 1) + (1 + 2*0 = 1) = 3$$

$$d = (0 + 2*1 = 2) + (1 + 2*0 = 1) + (1 + 2*0) + (1 + 2*0 = 1) = 5$$
                                    $$= 1$$

$$e = (0 + 2*1 = 2) + (0 + 2*0 = 0) + (0 + 2*1 = 2) + (0 + 2*0 = 0) = 4$$

$$f = (1 + 2*0 = 1) + (0 + 2*1 = 2) + (1 + 2*0 = 1) + (0 + 2*0 = 0) = 4$$

```
        d    e
a  R0   R1   R2        a, b, d higher units to save so
   ⇕    ↑    ↑
   @    b    d         here fix R0, R1, & R2 Registers
 c or f
```

→ c, e, f are used to use any registers R0, R1, R2