

UNIT V

SYLLABUS:

Run-time Environment: Storage Organization, Stack Allocation strategies: Static, Stack, Heap allocation, Activation record.

Code optimization: Introduction, Principal sources of optimization, Flowgraphs, Techniques in global and local optimization.

Code Generation: Issues in code generation, DAG, Simple code generator.

RUN-TIME ENVIRONMENT

The final phase in the compiler model is the code generator. It takes as input an intermediate representation of the source program and produces as output an equivalent target program. The code generation techniques presented below can be used whether or not an optimizing phase occurs before code generation.

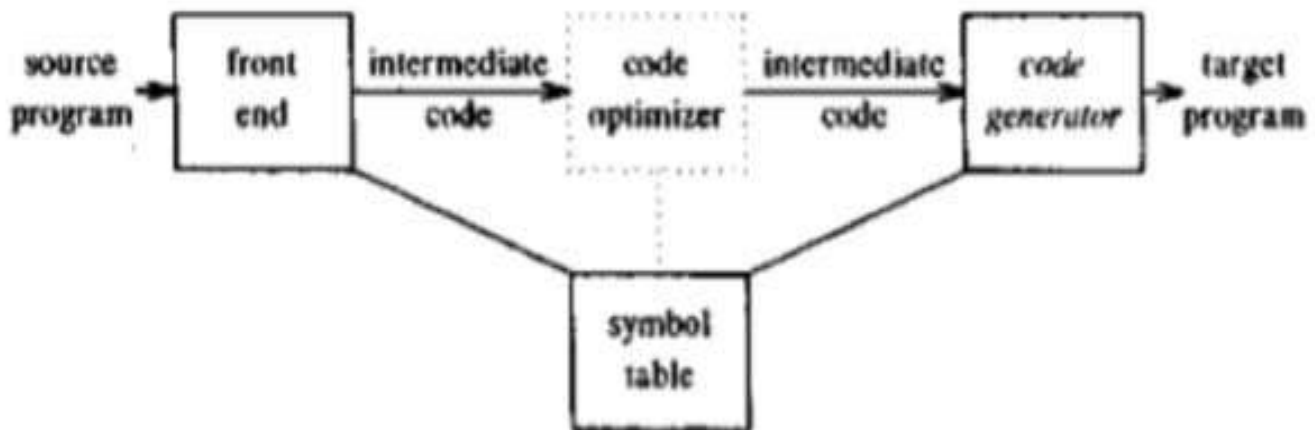


Figure 5.1 Position of code generator

STORAGE ORGANISATION

- The executing target program runs in its own logical address space in which each program value has allocation.
- The management and organization of this logical address space is shared between the compiler, operating system and target machine. The operating system maps the logical address into physical addresses, which are usually spread throughout memory.

Typical subdivision of run-time memory:

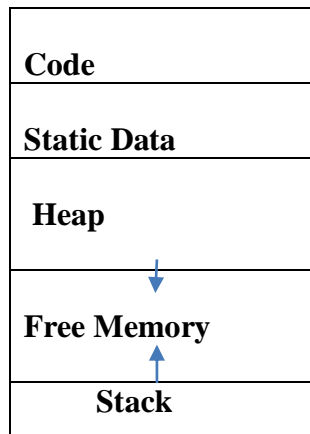


Figure 5.2: Typical subdivision of run-time memory:

- Run-time storage comes in blocks, where a byte is the smallest unit of addressable memory. Four bytes form a machine word. Multi-byte objects are stored in consecutive bytes and given the address of first byte.
 - This run-time storage might be subdivided to hold:
 1. The generated target code,
 2. Data objects, and
 3. A counterpart of the control stack to keep track of procedure activations.
- The storage layout for data objects is strongly influenced by the addressing constraints of the target machine.
- A character array of length 10 needs only enough bytes to hold 10 characters, a compiler may allocate 12 bytes to get alignment, leaving 2 bytes unused.
- This unused space due to alignment considerations is referred to as padding.
- The size of some program objects may be known at runtime and may be placed in an area called static.
- The dynamic areas used to maximize the utilization of space at runtime are stack and heap.

ACTIVATION RECORDS:

- Procedure calls and returns are usually managed by a runtime stack called the *control stack*.
- Each live activation has an activation record on the control stack, with the root of the activation tree at the bottom, the latter activation has its record at the top of the stack.
- The contents of the activation record vary with the language being implemented. The diagram below shows the contents of activation record.

Actual Parameters
Return Value
Control Link
Access Link
Machine Status
Local Data
Temporaries

Figure 5.3: General activation record.

- Temporary values such as those arising from the evaluation of expressions.
- Local data belonging to the procedure whose activation record this is.
- A saved machine status, with information about the state of the machine just before the call to procedures.
- An access link may be needed to locate data needed by the called procedure but found elsewhere.
- A control link pointing to the activation record of the caller.
- Space for the return value of the called functions, if any. Again, not all called procedures return a value, and if one does, we may prefer to place that value in a register for efficiency.
- The actual parameters used by the calling procedure. These are not placed in activation record but rather in registers, when possible, for greater efficiency.

STORAGE ALLOCATION STRATEGIES

The different storage allocation strategies are :

1. **Static allocation** – lays out storage for all data objects at compile time
2. **Stack allocation** – manages the run-time storage as a stack.
3. **Heap allocation** – allocates and deallocates storage as needed at run time from a data area known as heap

Static allocation

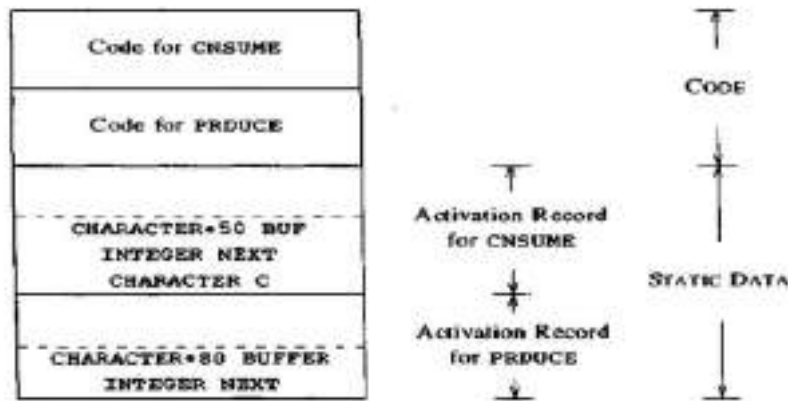
- In static allocation, names are bound to storage as the program is compiled, so there is no need for a run-time support package.
- Since the bindings do not change at run-time, every time a procedure is activated, its names are bound to the same storage locations.
- Therefore values of local names are *retained* across activations of a procedure. That is, when control returns to a procedure the values of the locals are the same as they were when control left the last time.
- From the type of a name, the compiler decides the amount of storage for the name and decides where the activation records go. At compile time, we can fill in the addresses at which the target code can find the data it operates on.

Some limitations of using static allocation:

1. The size of a data object and constraints on its position in memory must be known at compile time.
2. Recursive procedures are restricted, because all activations of a procedure use the same bindings for local names.
3. Data structures cannot be created dynamically, since there is no mechanism for storage allocation at run time.

FORTRAN uses static storage allocation

- **Consider the program**
- **Program CNSUME**
- **BUF : Array(1..80) of char**
- **NEXT : int**
- **C : char**
- **.....**
- **End**
- **Function PRDUCE**
- **BUFFER : Array(1..80) of char**
- **NEXT : int**
- **.....**
- **end**



Stack allocation

- All compilers for languages that use procedures, functions or methods as units of user-defined actions manage at least part of the run-time memory as a stack.
- Each time a procedure is called, space for its local variables is pushed onto a stack, and when the procedure terminates, that space is popped off the stack.

Calling sequences:

- Procedures called are implemented in what is called a calling sequence, which consists of code that allocates an activation record on the stack and enters information into its fields.
- A return sequence is similar to code to restore the state of machines so the calling procedure can continue its execution after the call.
- The code in a calling sequence is often divided between the calling procedure (caller) and the procedure it calls (callee).
- When designing calling sequences and the layout of activation records, the following principles are helpful:
 - Values communicated between caller and callee are generally placed at the beginning of the callee's activation record, so they are as close as possible to the caller's activation record.
 - Fixed-length items are generally placed in the middle. Such items typically include the control link, the access link, and the machine status fields.
 - Items whose size may not be known nearly enough are placed at the end of the activation record. The most common example is a dynamically sized array, where the value of one of the callee's parameters determines the length of the array.
 - We must locate the top-of-stack pointer judiciously. A common approach is to have it point to the end of fixed-length fields in the activation record. Fixed-length data can then be accessed by fixed offsets, known to the intermediate-code generator, relative to the top-of-stack point

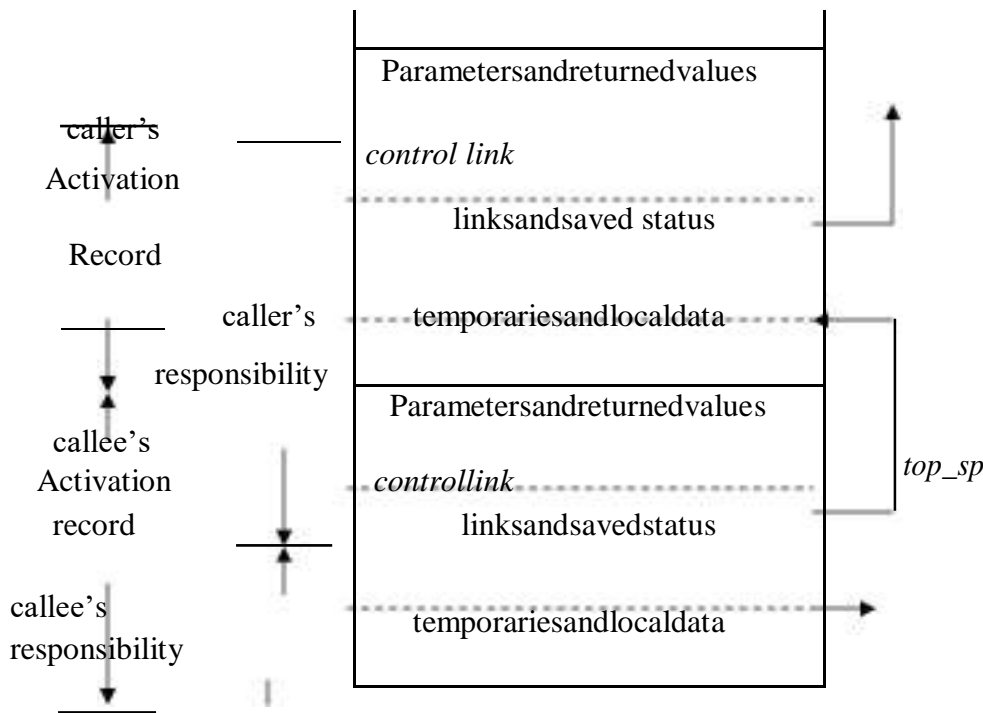


Figure 5.5: Division of tasks between caller and callee



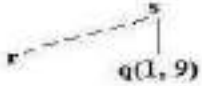
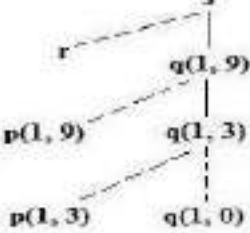
- The calling sequence and its division between caller and callee are as follows.
 - ☐ The callee evaluates the actual parameters.
 - ☐ The caller stores a return address and the old value of *top_sp* into the callee's activation record. The caller then increments the *top_sp* to the respective positions.
 - ☐ The callee saves the register values and other status information.
 - ☐ The callee initializes its local data and begins execution.
- As suitable, the corresponding return sequence is:
 - ☐ The callee places the return value next to the parameters.
 - ☐ Using the information in the machine-status field, the callee restores *top_sp* and other registers, and then branches to the return address that the caller placed in the status field.
 - ☐ Although *top_sp* has been decremented, the caller knows where the return value is, relative to the current value of *top_sp*; the caller therefore may use that value.

Variable length data on stack:

- The run-time memory management system must deal frequently with the allocation of space for objects, the sizes of which are not known at the compile time, but which are local to a procedure and thus may be allocated on the stack.
- The reason to prefer placing objects on the stack is that we avoid the expense of garbage collecting

their space.

- The same scheme works for objects of any type if they are local to the procedure called and have a size that depends on the parameters of the call.

Position in Activation Tree	Activation Record on the Stack	Remarks
	<div>s</div> <hr/> <div>a : array</div>	Frame for s
	<div>s</div> <hr/> <div>a : array</div> <hr/> <div>r</div> <hr/> <div>i : integer</div>	r is activated
	<div>s</div> <hr/> <div>a : array</div> <hr/> <div>q(1, 9)</div> <hr/> <div>i : integer</div>	Frame for r has been popped and q(1, 9) pushed
	<div>s</div> <hr/> <div>a : array</div> <hr/> <div>q(1, 9)</div> <hr/> <div>i : integer</div> <hr/> <div>q(1, 3)</div> <hr/> <div>i : integer</div>	control has just returned to q(1, 3)

Heap allocation

Stack allocation strategy cannot be used if either of the following is possible:

- The values of local names must be retained when an activation ends.
 - A called activation outlives the caller.
- Heap allocation parcels out pieces of contiguous storage, as needed for activation records or other objects.
 - Pieces may be deallocated in any order, so over the time the heap will consist of alternate areas that are free and in use

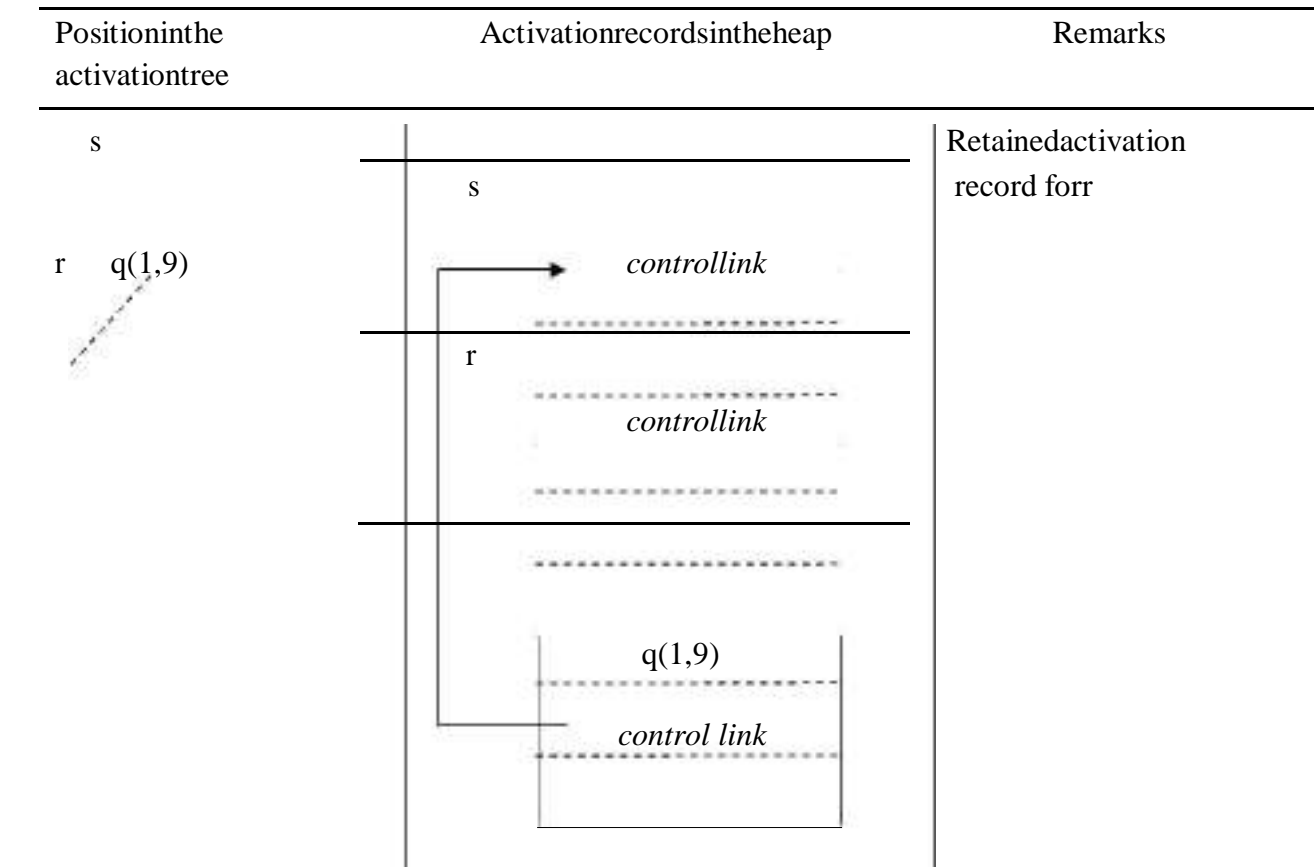


Figure 5.6: Heap allocation

- Therecordforanactivationofprocedurerisretainedwhentheactivationends.
- Therefore,therecordforthenewactivationq(1,9)cannotfollowthatfor sphysically.
- Iftheretainedactivationrecordforrisdeallocated,therewillbefreespaceintheheapbetweentheactivationrecords for s andq.
- ☐ Forlargeblocksofstorageusetheheapmanager. Thisapproachresultsinfastallocation anddeallocationofsmallamountsofstorage,sincetakingandreturningablockfromlinkedlistare efficientoperations.

CODE GENERATION

ISSUESINTHEDESIGNNOFACODEGENERATOR

Thefollowingissuesariseduringthecodegenerationphase:

1. Input tocode generator
2. Targetprogram
3. Memorymanagement

4. Instruction selection
5. Register allocation
6. Evaluation order

1. Input to code generator:

- The input to the code generation consists of the intermediate representation of the source program produced by front end, together with information in the symbol table to determine run-time addresses of the data objects denoted by the names in the intermediate representation.
- Intermediate representation can be:
 - a. Linear representations such as postfix notation
 - b. Three address representations such as quadruples
 - c. Virtual machine representations such as stack machine code
 - d. Graphical representations such as syntax trees and DAGs.
- Prior to code generation, the front end must be scanned, parsed and translated into intermediate representation along with necessary type checking. Therefore, input to code generation is assumed to be error-free.

2. Target program:

- The output of the code generator is the target program. The output may be:
 - a. Absolute machine language
 - It can be placed in a fixed memory location and can be executed immediately.
 - b. Relocatable machine language
 - It allows subprograms to be compiled separately.
 - c. Assembly language
 - Code generation is made easier.

3. Memory management:

- Names in the source program are mapped to addresses of data objects in run-time memory by the front end and code generator.
- It makes use of symbol table, that is, a name in a three-address statement refers to a symbol-table entry for the name.
- Labels in three-address statements have to be converted to addresses of instructions. For example,
 - j : goto i generates jump instruction as follows:
 - if $i < j$, a backward jump instruction with target address equal to location of code for quadruple i is generated.
 - if $i > j$, the jump is forward. We must store on a list for quadruple i the location of the first machine instruction generated for quadruple j . When i is processed, the machine locations for all instructions that forward jump to i are filled.

4. Instruction selection:

- The instructions of target machines should be complete and uniform. is considered.
- The quality of the generated code is determined by its speed and size.
- The former statement can be translated into the latter statement as shown below:

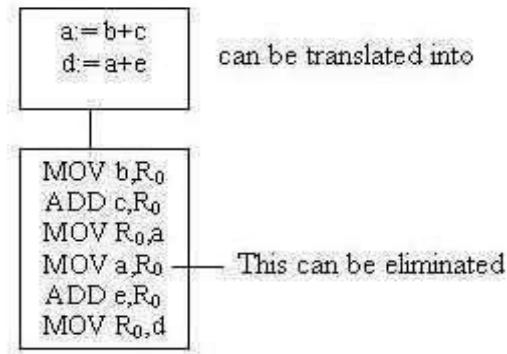


Figure 5.3: Code Translation

5. Register allocation

- Instructions involving register operands are shorter and faster than those involving operands in memory.
- The use of registers is subdivided into two subproblems:

Register allocation—the set of variables that will reside in registers at a point in the program is selected.

➤ **Register assignment**—the specific register that a variable will reside in is picked.

- Certain machines require even-odd register pairs for some operands and results. For example, consider the division instruction of the form:

Dx, y

where, x —dividend even register y —even/odd register pair—
divisor
even register holds the remainder
odd register holds the quotient

6. Evaluation order

- The order in which the computations are performed can affect the efficiency of the target code. Some computation orders require fewer registers to hold intermediate results than others.

ASIMPLECODEGENERATOR

- A code generator generates target code for a sequence of three-address statements and effectively uses registers to store operands of the statements.

For example: consider the three-address statement $a := b + c$

It can have the following sequence of codes:

ADDR _j , R _i	Cost=1	//if R _i contains b and R _j contains c (or)
ADDc, R _i	Cost=2	//if c is in a memory location
(or)		
MOVc, R _j	Cost=3	//move c from memory to R _j and add
ADDR _j , R _i		

Register and Address Descriptors:

- A register descriptor is used to keep track of what is currently in each register. The register descriptors show what initially all the registers are empty.
- An address descriptor stores the location where the current value of the name can be found at runtime.

A code-generation algorithm:

The algorithm takes as input a sequence of three-address statements constituting a basic block. For each three-address statement of the form $x := yopz$, perform the following actions:

1. Invoke a function *getreg* to determine the location L where the result of the computation $yopz$ should be stored.
2. Consult the address descriptor for y to determine y' , the current location of y . Prefer the register for y' if the value of y is currently both in memory and a register. If the value of y is not already in L , generate the instruction **MOV** y', L to place a copy of y in L .
3. Generate the instruction **OP** z', L where z' is a current location of z . Prefer a register to a memory location if z is in both. Update the address descriptor of x to indicate that x is in location L . If x is in L , update its descriptor and remove x from all other descriptors.
4. If the current values of y or z have no next uses, are not live on exit from the block, and are in registers, alter the register descriptor to indicate that, after execution of $x := yopz$, those registers will no longer contain y or z .

The algorithmic sequence of *getreg* function can be,

1. if x value is in register that register is returned.
2. If (1) fails, new register is returned.

3. If(2) fails, and the operation needs a special register, that register value is temporarily moved to the memory and the register is returned.
4. If(3) fails, finally memory location is returned.

Generating Code for Assignment Statements:

- The assignment $d := (a-b) + (a-c) + (a-c)$ might be translated into the following three-address code sequence:

```

    t := a - b      u :
    = a - c        v := t
    + u d := v + u

```

with `dlive` at the end.

Code sequence for the example is:

Statements	Code Generated	Register descriptor	Address descriptor
		Register empty	
t:=a-b	MOV a,R0 SUB b,R0	R0 contains t	t in R0
u:=a-c	MOV a,R1 SUB c,R1	R0 contains t R1 contains u	t in R0 u in R1
v:=t+u	ADD R1,R0	R0 contains v R1 contains u	u in R1 v in R0
d:=v+u	ADD R1,R0 MOV R0,d	R0 contains d	d in R0 d in R0 and memory

Generating Code for Indexed Assignments

The table shows the code sequences generated for the indexed assignment statements **a:=b[i]** and **a[i]:=b**

Statements	Code Generated
a:=b[i]	MOV b(Ri),R
a[i]:= b	MOV b,a(Ri)

Generating Code for Pointer Assignments

The table shows the code sequences generated for the pointer assignments **a := *p** and ***p := a**

Statements	Code Generated
a:=*p	MOV *Rp,a
*p:=a	MOV a,*Rp

Generating Code for Conditional Statements:

Statement	Code
if x < y goto z	CMP x, y CJ<z /* jump to z if condition code is negative */
x := y +z if x <0 goto z MOV y, R0	ADD z, R0 MOV R0,x CJ<Z

DAG:

THE DAG REPRESENTATION FOR BASIC BLOCKS

- A DAG for a basic block is a directed acyclic graph with the following labels on nodes:
 1. Leaves are labeled by unique identifiers, either variable names or constants.
 2. Interior nodes are labeled by an operator symbol.
 3. Nodes are also optionally given a sequence of identifiers for labels to store the computed values.
- DAGs are useful data structures for implementing transformations on basic blocks.
- It gives a picture of how the value computed by a statement is used in subsequent statements.
- It provides a good way of determining common sub - expressions.

Algorithm for construction of DAG

Input: A basic block

Output: A DAG for the basic block containing the following information:

1. A label for each node. For leaves, the label is an identifier. For interior nodes, an operator symbol.
2. For each node a list of attached identifiers to hold the computed values. Case (i) $x := y \text{ OP } z$
Case (ii) $x := \text{OP } y$ Case (iii) $x := y$

Method:

Step1: If y is undefined then create node(y).

If z is undefined, create node(z) for case (i).

Step2: For the case (i), create a node(OP) whose left child is node(y) and right child is node(z). (

Checking for common sub expression). Let n be this node.

For case (ii), determine whether there is node(OP) with one child node(y). If not create such a node.

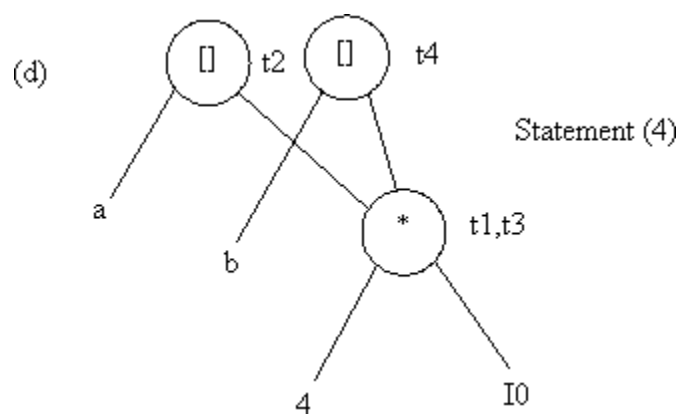
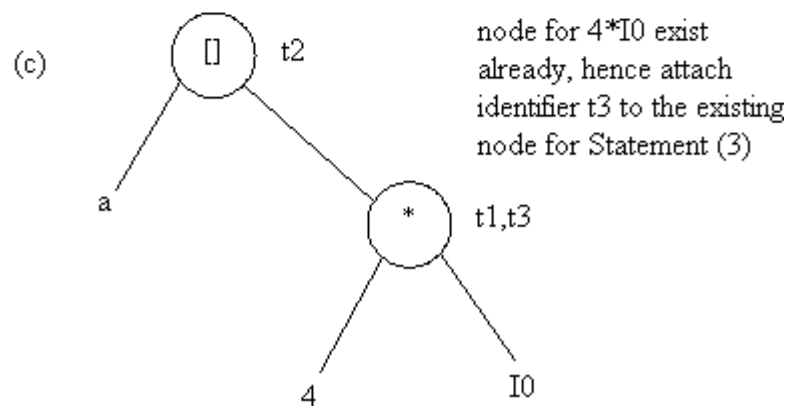
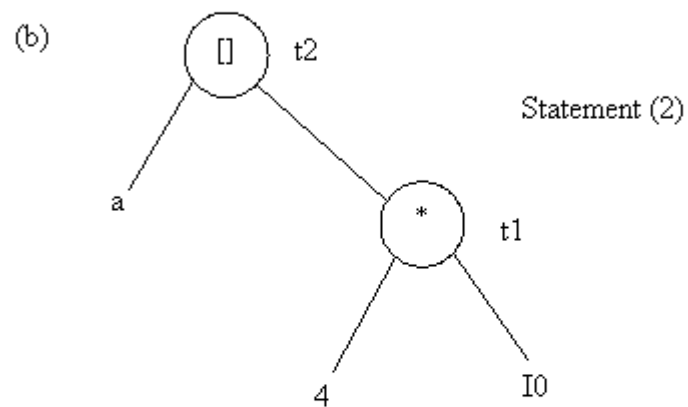
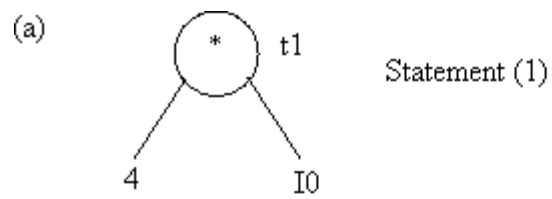
For case (iii), node n will be node(y).

Step3: Delete x from the list of identifiers for node(x). Append to the list of attached identifiers for the node n found in step 2 and set node(x) to n .

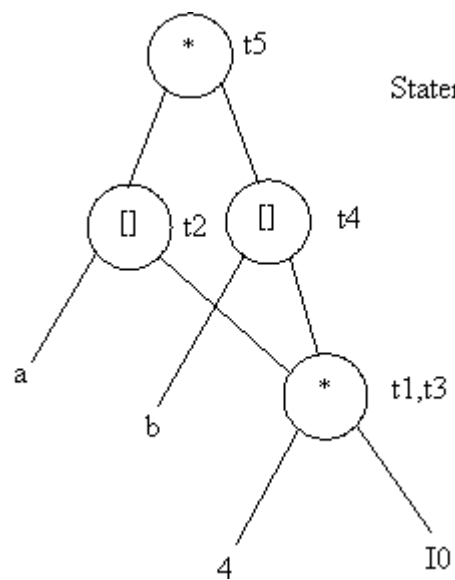
Example: Consider the block of three- address statements:

1. $t_1 := 4 * i$
2. $t_2 := a[t_1]$
3. $t_3 := 4 * i$
4. $t_4 := b[t_3]$
5. $t_5 := t_2 * t_4$
6. $t_6 := \text{prod} + t_5$
7. $\text{prod} := t_6$
8. $t_7 := i + 1$
9. $i := t_7$
10. if $i \leq 20$ goto(1)

Stages in DAG Construction

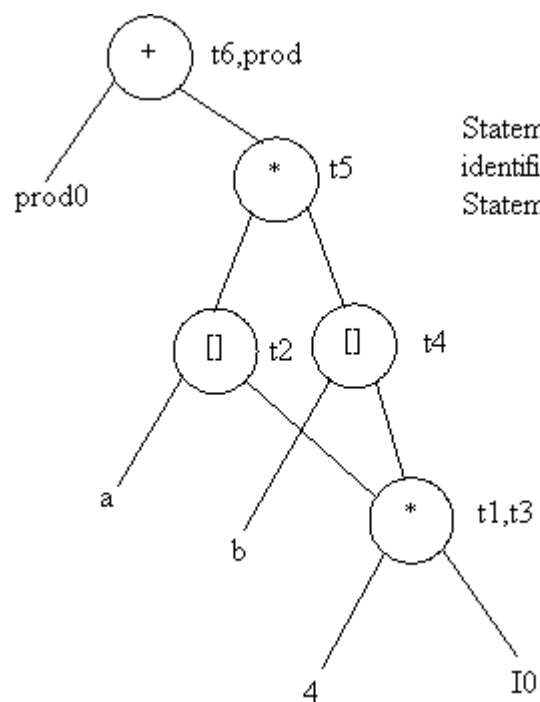


(e)



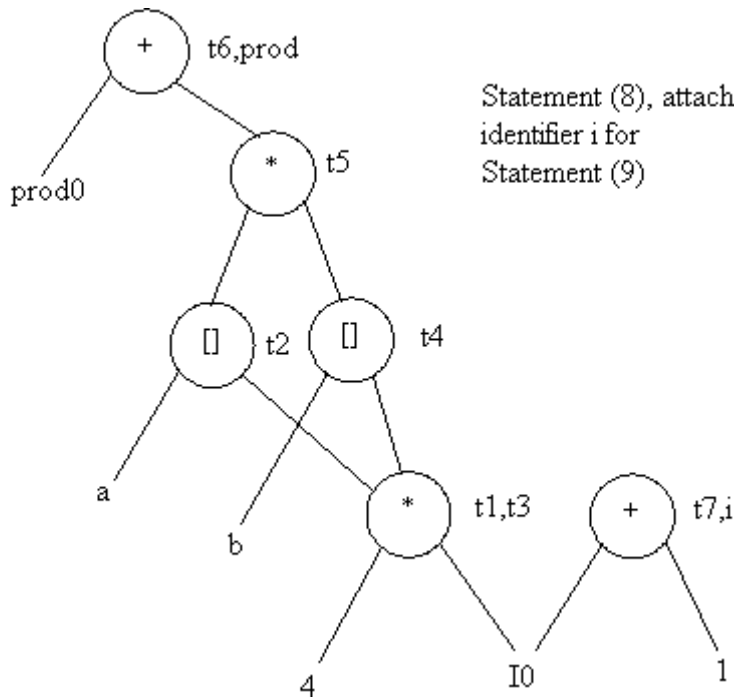
Statement (5)

(f)

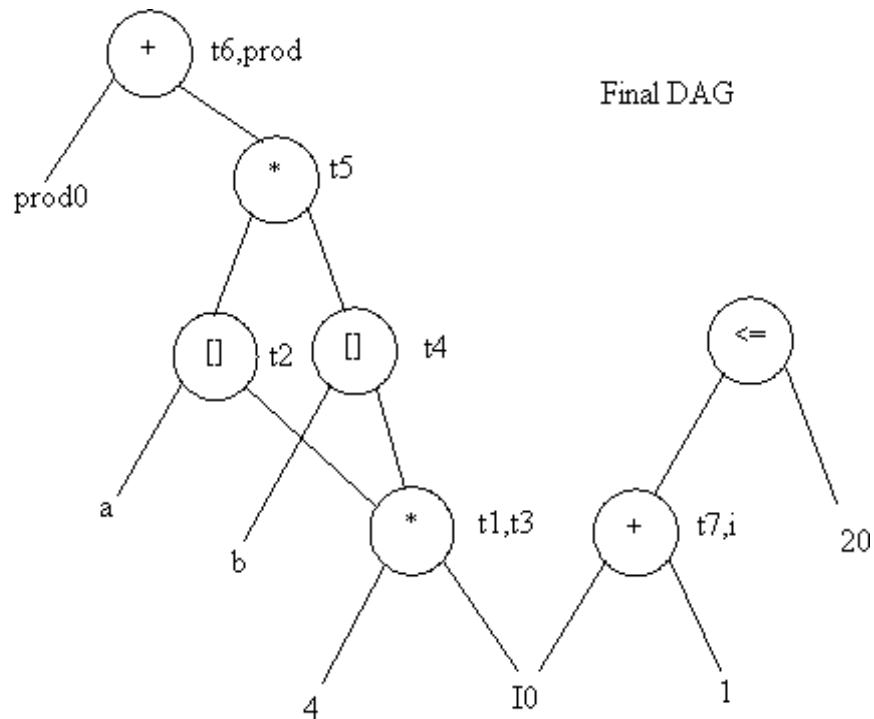


Statement (6), attach
identifier prod for
Statement (7)

(g)



(h)



Application of DAGs:

1. We can automatically detect common sub expressions.
2. We can determine which identifiers have their values used in the block.
3. We can determine which statements compute values that could be used outside the block.

GENERATING CODE FROM DAGs

The advantage of generating code for a basic block from its dag representation is that, from a dag we can easily see how to rearrange the order of the final computation sequence than we can starting from a linear sequence of three-address statements or quadruples.

Rearranging the order

The order in which computations are done can affect the cost of resulting object code.

For example, consider the following basic block: $t_1 := a + b$

$t_2 := c + d$

$t_3 := e - t_2$

$t_4 := t_1 - t_3$

Generated code sequence for basic block:

MOV a , R₀

ADD b , R₀

MOV c , R₁

ADD d , R₁

MOV R₀ , t₁

MOV e , R₀

SUB R₁ , R₀

MOV t₁ , R₁

SUB R₀ , R₁

MOV R₁ , t₄

Rearranged basic block:

Now t₁ occurs immediately before t₄.

$t_2 := c + d$

$t_3 := e - t_2$

$t_1 := a + b$

$t_4 := t_1 - t_3$

Revised code sequence:

MOV c , R₀

ADD d , R₀

MOV a , R₀

SUB R₀ , R₁

MOV a , R₀

ADD b , R₀

SUB R₁ , R₀

MOV R₀ , t₄

In this order, two instructions MOV R₀ , t₁ and MOV t₁ , R₁ have been saved.

A Heuristic ordering for Dags:

The heuristic ordering algorithm attempts to make the evaluation of a node immediately follow the evaluation of its leftmost argument.

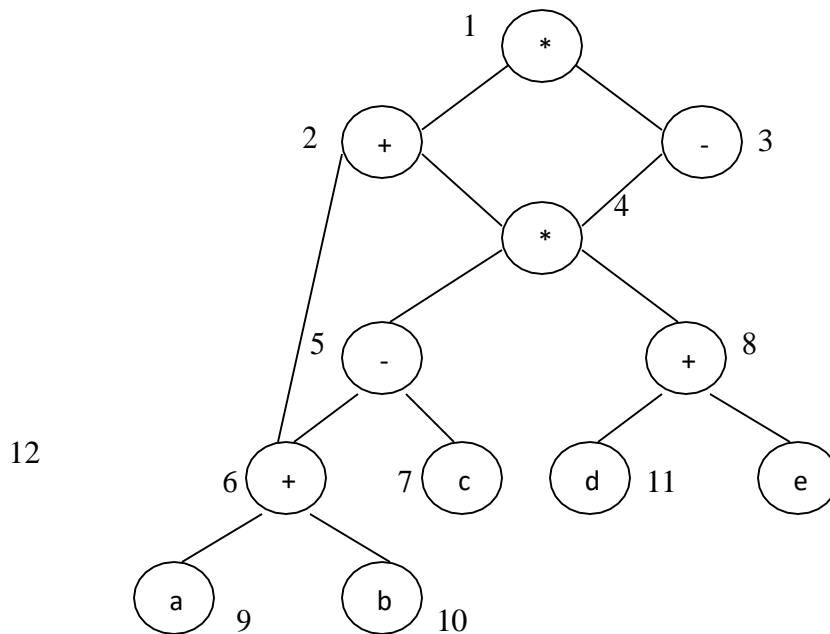
The algorithm shown below produces the ordering in reverse.

Algorithm:

- 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 leftmost child m of n has no unlisted parents and is not a leaf do begin
- 5) list m ;
- 6) $n := m$

endend

Example: Consider the DAG shown below:



Initially, the only node with no unlisted parents is 1 so set $n=1$ at line (2) and list 1 at line (3).

Now, the left argument of 1, which is 2, has its parents listed, so we list 2 and set $n=2$ at line (6).

Now, at line (4) we find the leftmost child of 2, which is 6, has an unlisted parent 5. Thus we select a new n at line (2), and node 3 is the only candidate. We list 3 and proceed down its left chain, listing 4, 5 and 6. This leaves only 8 among the interior nodes so we list that.

The resulting list is 1234568 and the order of evaluation is 8654321.

Code sequence:

$t_8 := d + e$

$t_6 := a + b$

$t_5 := t_6 - c$

$t_4 := t_5 * t_8$

$t_3 := t_4 - e$

$t_2 := t_6 + t_4$

$t_1 := t_2 * t_3$

This will yield an optimal code for the DAG on machine whatever be the number of registers.

CODE OPTIMIZATION

INTRODUCTION: The code produced by the straight forward compiling algorithms can often be made to run faster or take less space, or both. This improvement is achieved by program transformations that are traditionally called optimizations. Compilers that apply code-improving transformations are called optimizing compilers.

Optimizations are classified into two categories. They are

- Machine independent optimizations:
- Machine dependent optimizations:

Machine independent optimizations:

- Machine independent optimizations are program transformations that improve the target code without taking into consideration any properties of the target machine.

Machine dependent optimizations:

- Machine dependent optimizations are based on register allocation and utilization of special machine-instruction sequences.

The criteria for code improvement transformations:

- ✓ Simply stated, the best program transformations are those that yield the most benefit for the least effort.
- ✓ The transformation must preserve the meaning of programs. That is, the optimization must not change the output produced by a program for a given input, or cause an error such as division by zero, that was not present in the original source program. At all times we take the “safe” approach of missing an opportunity to apply a transformation rather than risk changing what the program does.
- ✓ A transformation must, on the average, speed up programs by a measurable amount. We are also interested in reducing the size of the compiled code although the size of the code has less importance than it once had. Not every transformation succeeds in improving every program, occasionally an “optimization” may slow down a program slightly.
- ✓ The transformation must be worth the effort. It does not make sense for a compiler writer to expend the intellectual effort to implement a code improving transformation and to have the compiler expend the additional time compiling source programs if this effort is not repaid when the target programs are executed. “Peephole” transformations of this kind are simple enough and beneficial enough to be included in any compiler.

Organization for an Optimizing Compiler:

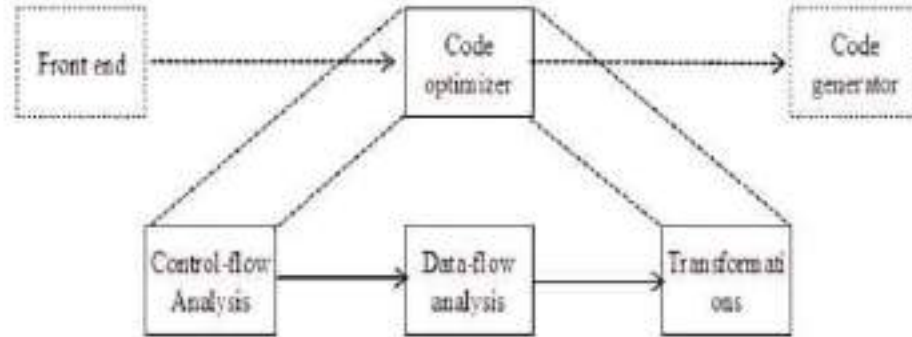


Figure 5.7: Organization for an Optimizing Compiler

- Flow analysis is a fundamental prerequisite for many important types of code improvement.
- Generally control flow analysis precedes data flow analysis.
- Control flow analysis (CFA) represents flow of control usually in form of graphs, CFA constructs such as:
 - Control flow graph
 - Call graph
- Data flow analysis (DFA) is the process of asserting and collecting information prior to program execution about the possible modification, preservation, and use of certain entities (such as values or attributes of variables) in a computer program.

PRINCIPAL SOURCES OF OPTIMISATION

- A transformation of a program is called local if it can be performed by looking only at the statements in a basic block; otherwise, it is called global.
- Many transformations can be performed at both the local and global levels. Local transformations are usually performed first.

Function-Preserving Transformations

- There are a number of ways in which a compiler can improve a program without changing the function it computes.
- The transformations:
 - ✓ Common subexpression elimination,
 - ✓ Copy propagation,
 - ✓ Dead-code elimination, and
 - ✓ Constant folding

Are common examples of such function-preserving transformations. The other transformations come up primarily when global optimizations are performed.

- Frequently, a program will include several calculations of the same value, such as an offset in an array. Some of the duplicate calculations cannot be avoided by the programmer because they lie below the level of detail accessible within the source language.

➤ CommonSubexpressionselimination:

- An occurrence of an expression E is called a common sub-expression if E was previously computed, and the values of variables in E have not changed since the previous computation. We can avoid recomputing the expression if we can use the previously computed value.
- For example

```
t1: = 4*i
t2:=a[t1]
t3: = 4*j
t4: = 4*i
t5: = n
t6:=b[t4]+t5
```

The above code can be optimized using the common sub-expression elimination as $t_1: = 4*i$

```
t2:=a[t1]
t3: = 4*j
t5: = n
t6:=b[t1]+t5
```

The common sub expression $t_4: = 4*i$ is eliminated as its computation is already in t_1 . And value of i is not been changed from definition to use.

➤ CopyPropagation:

- Assignments of the form $f := g$ called copy statements, or copies for short. The idea behind the copy-propagation transformation is to use g for f , whenever possible after the copy statement $f := g$. Copy propagation means use of one variable instead of another. This may not appear to be an improvement, but as we shall see it gives us an opportunity to eliminate x .
- For example:

```
x=Pi;
A=x*r*r;
```

The optimization using copy propagation can be done as follows: $A=Pi*r*r$;
Here the variable x is eliminated

➤ Dead-Code Eliminations:

- A variable is live at a point in a program if its value can be used subsequently; otherwise, it is dead at that point. A related idea is dead or useless code, statements that compute

Values that never get used. While the programmer is unlikely to introduce any dead code intentionally, it may appear as the result of previous transformations. An optimization can be done by eliminating dead code.

Example:

```
i=0;
if(i=1)
{
a=b+5;
}
```

Here, 'if' statement is dead code because this condition will never get satisfied.

➤ Constant folding:

- We can eliminate both the test and printing from the object code. More generally, deducing at compile time that the value of an expression is a constant and using the constant instead is known as constant folding.
- One advantage of copy propagation is that it often turns the copy statement into dead code.
- For example,

$a = 3.14157/2$ can be replaced by

$a = 1.570$ thereby eliminating a division operation.

➤ Loop Optimizations:

- We now give a brief introduction to a very important place for optimizations, namely loops, especially the inner loops where program tends to spend the bulk of their time. The running time of a program may be improved if we decrease the number of instructions in an inner loop, even if we increase the amount of code outside that loop.
- Three techniques are important for loop optimization:
- code motion, which moves code outside a loop;
- Induction-variable elimination, which we apply to replace variables from inner loop.
- Reduction in strength, which replaces an expensive operation by a cheaper one, such as a multiplication by an addition.

➤ Code Motion:

- An important modification that decreases the amount of code in a loop is code motion. This transformation takes an expression that yields the same result independent of the number of times a loop is executed (a loop-invariant computation) and places the expression before the loop. Note that the notion “before the loop” assumes the existence of an entry for the loop. For example, evaluation of limit-2 is a loop-invariant computation in the following while-statement:

```
while (i<= limit-2)          /*statement does not change limit*/
```

Code motion will result in the equivalent of

```
t=limit-2;
while(i<=t)          /*statementdoesnotchangelimitort*/
```

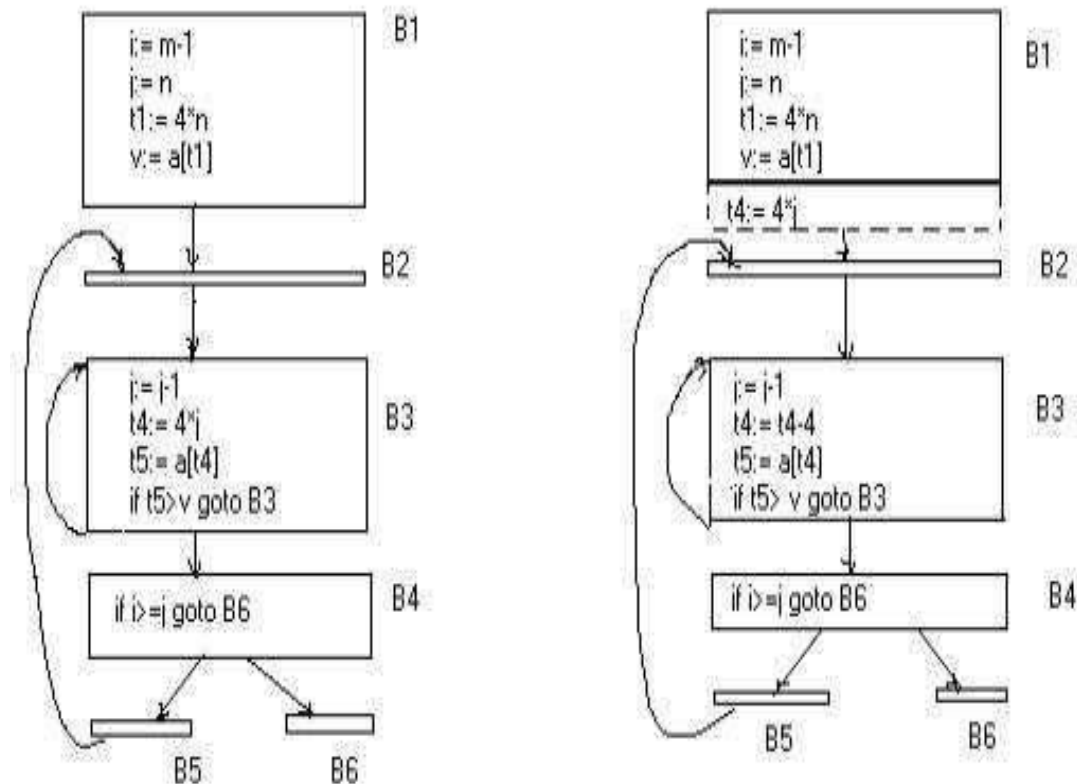
➤ InductionVariables:

- Loopsareusuallyprocessedinsideout.ForexampleconsiderthelooparoundB3.
- Note that the values of j and t_4 remain in lock-step; everytime the value of j decreases by 1, that of t_4 decreases by 4 because $4*j$ is assigned to t_4 . Such identifiers are called induction variables.
- Whentherearetwoor moreinduction variables inaloop,itmaybe possibletogetridof all but one, by the process of induction-variable elimination. For the inner loop around B3 in Fig. we cannot get rid of either j or t_4 completely; t_4 is used in B3 and j in B4. However, we can illustrate reduction in strength and illustrate a part of the process of induction-variableelimination.EventuallyjwillbeeliminatedwhentheouterloopofB2

-B5isconsidered.

Example:

As the relationship $t_4 := 4*j$ surelyholds after such an assignment to t_4 in Fig. and t_4 is not changed elsewhere in the inner loop around B3, it follows that just after the statement $j := j-1$ the relationship $t_4 := 4*j-4$ musthold. Wemaythereforereplacetheassignment $t_4 := 4*j$ by $t_4 := t_4 - 4$. Theonlyproblemis that t_4 doesnothavea valuewhenweenterblock B3 for the firsttime. Since wemustmaintaintherelationship $t_4 = 4*j$ onentrytotheblockB3, we place an initialization of t_4 at the end of the block where j itself is



before

after

Initialized, shown by the dashed addition to block B1 in second Fig.

- Thereplacementofamultiplicationbyasubtractionwillspeeduptheobjectcodeif multiplication takes more time than addition or subtraction, as is the case on many machines.

➤ ReductionInStrength:

- Reduction in strength replaces expensive operations by equivalent cheaper ones on the target machine. Certain machine instructions are considerably cheaper than others and can often be used as special cases of more expensive operators.
- For example, x^2 is invariably cheaper to implement as $x*x$ than as a call to an exponentiation routine. Fixed-point multiplication or division by a power of two is cheapertoimplementas ashift.Floating-pointdivisionbyaconstantcanbeimplemented as multiplication by a constant, which may be cheaper.

OPTIMIZATIONOFBASICBLOCKS

Therearetwotypesofbasicblockoptimizations.Theyare:

- Structure-PreservingTransformations
- AlgebraicTransformations

Structure-PreservingTransformations:

TheprimaryStructure-PreservingTransformationonbasicblocksare:

- Commonsub-expressionelimination
- Deadcodeelimination
- Renamingoftemporaryvariables
- Interchangeoftwoindependentadjacentstatements.

➤ Commonsub-expressionelimination:

Common sub expressions need not be computed over and over again. Instead they can be computedonceand kept in store from where it's referencedwhenencounteredagain – of course providing the variable values in the expression still remain constant.

Example:

```
a:=b+c
b:=a-d
c:=b+c
d:=a-d
```

The2ndand4thstatementscomputethesameexpression: $b+c$ anda- d Basic block can be transformed to

```
a:=b+c
b:=a-d
c:=a
d:=b
```

➤ Deadcodeelimination:

It's possible that a large amount of dead (useless) code may exist in the program. This might be especially caused when introducing variables and procedures as part of construction or error-correction of a program – once declared and defined, one forgets to remove them in case they serve no purpose. Eliminating these will definitely optimize the code.

➤ Renamingoftemporaryvariables:

- A statement $t := b + c$ where t is a temporary name can be changed to $u := b + c$ where u is another temporary name, and change all uses of t to u .
- In this we can transform a basic block to its equivalent block called normal-form block.

➤ Interchangeoftwoindependentadjacentstatements:

- Two statements

$t_1 := b + c$

$t_2 := x + y$

Can be interchanged or reordered in its computation in the basic block when value of t_1 does not affect the value of t_2 .

Algebraic Transformations:

- Algebraic identities represent another important class of optimizations on basic blocks. This includes simplifying expressions or replacing expensive operation by cheaper ones i.e. reduction in strength.

- Another class of related optimizations is constant folding. Here we evaluate constant expressions at compile time and replace the constant expressions by their values. Thus the expression $2 * 3.14$ would be replaced by 6.28 .
- The relational operators $<=$, $>=$, $<$, $>$, $+$ and $=$ sometimes generate unexpected common sub expressions.
- Associative laws may also be applied to expose common sub expressions. For example, if the source code has the assignments

$a := b + c$

$e := c + d + b$

the following intermediate code may be generated:

$a := b + c$

$t := c + d$

$e := t + b$

Example:

$x := x + 0$ can be removed

$x := y * 2$ can be replaced by a cheaper statement $x := y * y$

- The compiler writer should examine the language carefully to determine what rearrangements of computations are permitted, since computer arithmetic does not always obey the algebraic identities of mathematics. Thus, a compiler may evaluate $x*y-x*z$ as $x*(y-z)$ but it may not evaluate $a+(b-c)$ as $(a+b)-c$.

FLOWGRAPH

A compiler first converts the source code of any programming language into an intermediate code. It is then converted into basic blocks. After dividing an intermediate code in basic blocks, the flow of control among basic blocks is represented by a flow graph.

Properties of Flow Graphs

1. The control flow graph is process-oriented.
2. A control flow graph shows how program control is passed among the blocks.
3. The control flow graph depicts all of the paths that can be traversed during the execution of a program.
4. It can be used in software optimization to find unwanted loops.

Representation of Flow Graphs

Flow graphs are directed graphs. The nodes/blocks of the control flow graph are the basic blocks of the program. There are two designated blocks in Control Flow Graph:

1. Entry Block: The entry block allows the control to enter in the control flow graph.
2. Exit Block: Control flow leaves through the exit block.

An edge can flow from one block A to another block B if:

1. the first instruction of the B's block immediately follows the last instruction of the A's block.
2. there is a conditional/unconditional jump from A's end to the starting of B.
3. B follows X in the original order of the three-address code, and A does not end in an unconditional jump.

Let's see an example,

Consider the source code for converting a 10 x 10 matrix to an identity matrix.

```
for r from 1 to 10 do
    for c from 1 to 10 do
        a [ r, c ] = 0.0;
```

```
for r from 1 to 10 do
    a [ r, c ] = 1.0;
```

The following are the three address codes for the above source code:

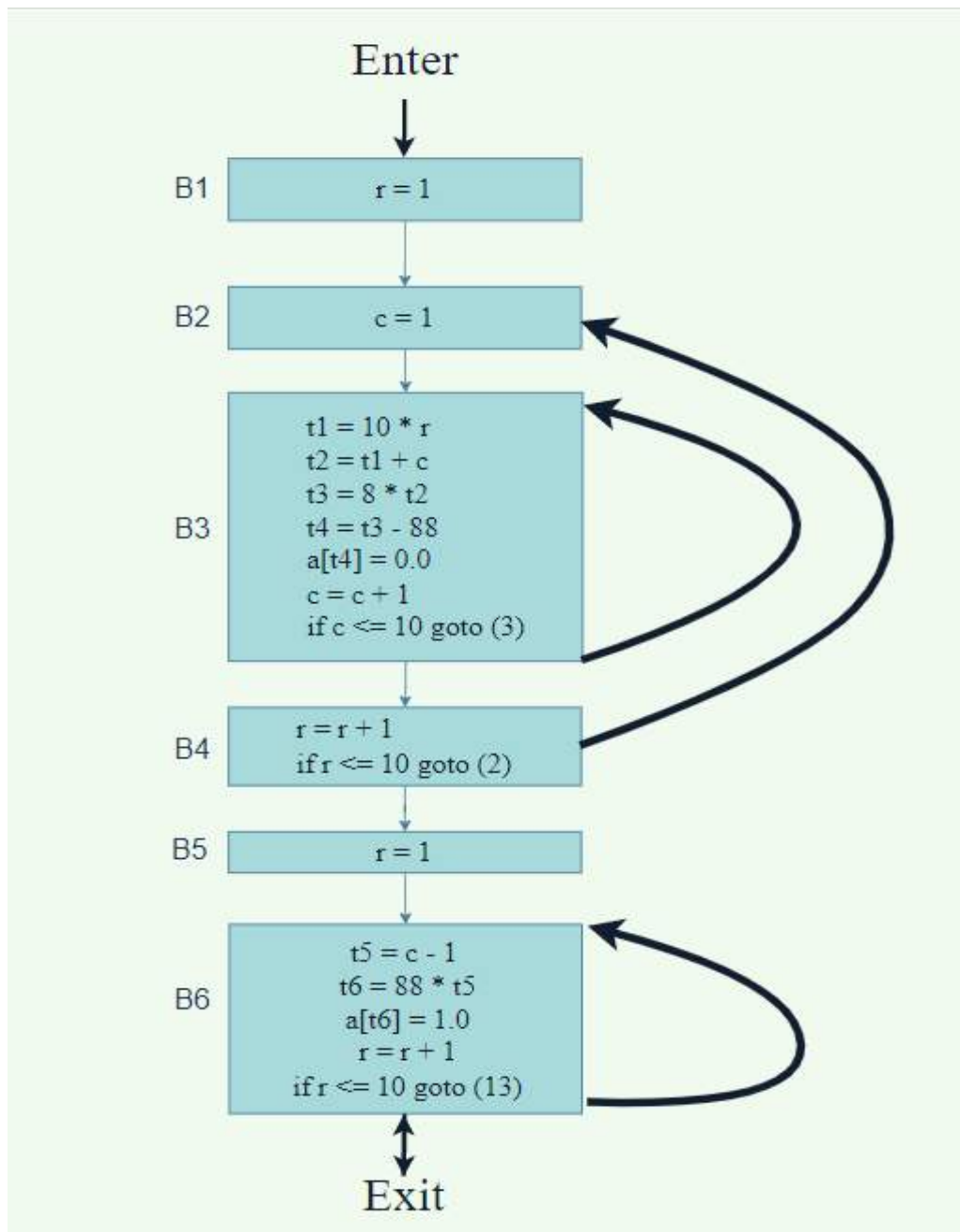
- 1) $r = 1$
- 2) $c = 1$
- 3) $t1 = 10 * r$
- 4) $t2 = t1 + c$
- 5) $t3 = 8 * t2$
- 6) $t4 = t3 - 88$
- 7) $a[t4] = 0.0$
- 8) $c = c + 1$
- 9) if $c \leq 10$ goto (3)
- 10) $r = r + 1$
- 11) if $r \leq 10$ goto (2)

```
12) r = 1
13) t5 = c - 1
14) t6 = 88 * t5
15) a[t6] = 1.0
16) r = r + 1
17) if r <= 10 goto (13)
```

There are six basic blocks for the above-given code, which are:

- B1 for statement 1
- B2 for statement 2
- B3 for statements 3-9
- B4 for statements 10-11
- B5 for statement 12
- B6 for statements 13-17.

The control flow graph of the above-given basic blocks is:



Explanation:

- B1 is the start point for the control flow graph because B1 contains the starting instructions.
- Because B1 does not end with unconditional jumps, and the B2 block's leader immediately follows B1's leader, B2 is the only successor of B1.
- There are two successors to the B3 block. The conditional jump in the last instruction of block B3 is targeted at the first instruction of the B3 block; therefore, one is block B3 itself. Another is block B4 due to conditional jump at the end of the B3 block.
- The last block, B6, is the exit point of the control flow graph.

LOOPSINFLOWGRAPH

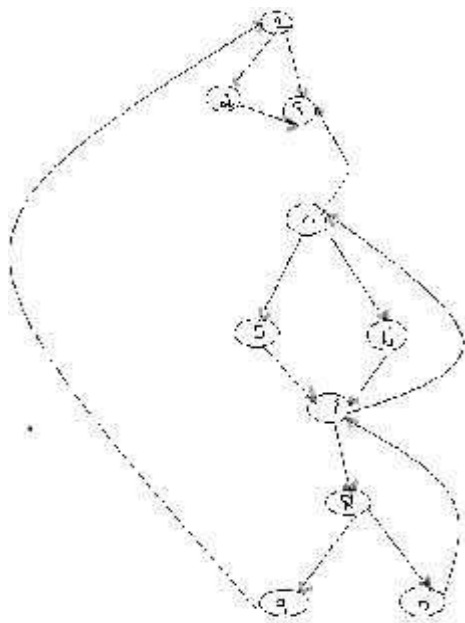
A graph representation of three-address statements, called a flow graph, is useful for understanding code-generation algorithms, even if the graph is not explicitly constructed by a code-generation algorithm. Nodes in the flow graph represent computations, and the edges represent the flow of control.

Dominators:

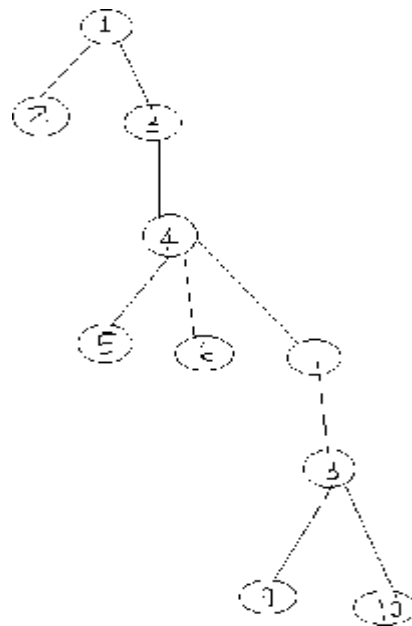
In a flow graph, a node d dominates node n , if every path from initial node of the flow graph to n goes through d . This will be denoted by $d \text{ dom } n$. Every initial node dominates all the remaining nodes in the flow graph and the entry of a loop dominates all nodes in the loop. Similarly every node dominates itself.

Example:

- *In the flow graph below,
- *Initial node, node 1 dominates every node.
- *node 2 dominates itself
- *node 3 dominates all but 1 and 2.
- *node 4 dominates all but 1, 2 and 3.
- *node 5 and 6 dominates only themselves, since flow of control can skip around either by going through the other.
- *node 7 dominates 7, 8, 9 and 10.
- *node 8 dominates 8, 9 and 10.
- *node 9 and 10 dominates only themselves.



- The way of presenting dominator information is in a tree, called the dominator tree in which the initial node is the root.
- The parent of each other node is its immediate dominator.
- Each node dominates only its descendants in the tree.
- The existence of dominator tree follows from a property of dominators; each node has a unique immediate dominator in that is the last dominator of n on any path from the initial node to n .
- In terms of the dom relation, the immediate dominator m has the property is $d \neq n$ and $d \text{ dom } n$, then $d \text{ dom } m$.



$$D(1) = \{1\}$$

$$D(2) = \{1, 2\}$$

$$D(3) = \{1, 3\}$$

$$D(4) = \{1, 3, 4\}$$

$$D(5) = \{1, 3, 4, 5\}$$

$$D(6) = \{1, 3, 4, 6\}$$

$$D(7) = \{1, 3, 4, 7\}$$

$$D(8) = \{1, 3, 4, 7, 8\}$$

$$D(9) = \{1, 3, 4, 7, 8, 9\}$$

$$D(10) = \{1, 3, 4, 7, 8, 10\}$$

NaturalLoop:

- One application of dominator information is in determining the loops of a flow graph suitable for improvement.
- The properties of loops are
 - A loop must have a single entry point, called the header. This entry point dominates all nodes in the loop, or it would not be the sole entry to the loop.
 - There must be at least one way to iterate the loop (i.e.) at least one path back to the header.
- One way to find all the loops in a flow graph is to search for edges in the flow graph whose heads dominate their tails. If $a \rightarrow b$ is an edge, b is the head and a is the tail. These types of edges are called as back edges.

Example:

In the above graph,

$7 \rightarrow 4$ $4 \text{ DOM } 7$

$10 \rightarrow 7$ $7 \text{ DOM } 10$

$4 \rightarrow 3$

$8 \rightarrow 3$

$2 \rightarrow 1$

- The above edges will form a loop in flow graph.
- Given a back edge $n \rightarrow d$, we define the natural loop of the edge to be the set of nodes that can reach n without going through d . Node d is the header of the loop.

❖ Algorithm: Constructing the natural loop of a back edge.

Input: A flow graph G and a back edge $n \rightarrow d$.

Output: The set loop consisting of all nodes in the natural loop $n \rightarrow d$.

Method: Beginning with node n , we consider each node $m \neq d$ that we know is in loop, to make sure that m 's predecessors are also placed in loop. Each node in loop, except for d , is placed once on stack, so its predecessors will be examined. Note that because d is put in the loop initially, we never examine its predecessors, and thus find only those nodes that reach n without going through d .

```
Procedure insert( $m$ );  
if  $m$  is not in  $\text{loop}$  then begin     $\text{loop}$   
   $:= \text{loop} \cup \{m\}$ ; push  $m$  onto  
   $\text{stack}$   
end;
```

$\text{stack} := \text{empty};$

```

loop:={d};
insert(n);
while stack is not empty do begin
  pop m, the first element of stack, off stack;
  for each predecessor p of m do insert(p)
end

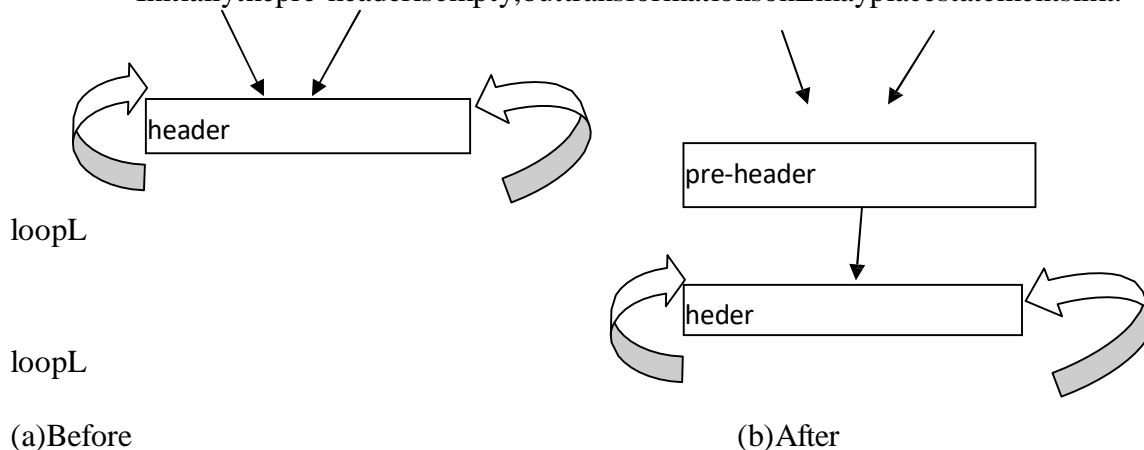
```

Innerloop:

- If we use the natural loops as “the loops”, then we have the useful property that unless two loops have the same header, they are either disjoint or one is entirely contained in the other. Thus, neglecting loops with the same header for the moment, we have a natural notion of inner loop: one that contains no other loop.
- When two natural loops have the same header, but neither is nested within the other, they are combined and treated as a single loop.

Pre-Headers:

- Several transformations require us to move statements “before the header”. Therefore begin treatment of a loop L by creating a new block, called the preheader.
- The pre-header has only the header as successor, and all edges which formerly entered the header of L from outside L instead enter the pre-header.
- Edges from inside loop L to the header are not changed.
- Initially the pre-header is empty, but transformations on L may place statements in it.



Reducible flow graphs:

- Reducible flow graphs are special flow graphs, for which several code optimization transformations are especially easy to perform, loops are unambiguously defined, dominators can be easily calculated, data flow analysis problems can also be solved efficiently.
- Exclusive use of structured flow-of-control statements such as if-then-else, while-do, continue, and break statements produces programs whose flow graphs are always reducible.

- The most important properties of reducible flow graphs are that there are no jumps into the middle of loops from outside; the only entry to a loop is through its header.
- *Definition:*

A flow graph G is reducible if and only if we can partition the edges into two disjoint groups, *forward* edges and *back* edges, with the following properties.

- The forward edges form an acyclic graph in which every node can be reached from the initial node of G .
- The back edges consist only of edges where the head dominates the tail.
- Example: The above flow graph is reducible.
- If we know the relation DOM for a flow graph, we can find and remove all the back edges.
- The remaining edges are forward edges.
- If the forward edges form an acyclic graph, then we can say the flow graph is reducible.
- In the above example, remove the five back edges $4 \rightarrow 3$, $7 \rightarrow 4$, $8 \rightarrow 3$, $9 \rightarrow 1$ and $10 \rightarrow 7$ whose heads dominate their tails; the remaining graph is acyclic.
- The key property of reducible flow graphs for loop analysis is that in such flow graphs every set of nodes that we would informally regard as a loop must contain a back edge.

PEEPHOLE OPTIMIZATION

- A statement-by-statement code-generation strategy often produces target code that contains redundant instructions and suboptimal constructs. The quality of such target code can be improved by applying “optimizing” transformations to the target program.
- A simple but effective technique for improving the target code is peephole optimization, a method for trying to improve the performance of the target program by examining a short sequence of target instructions (called the peephole) and replacing these instructions by a shorter or faster sequence, whenever possible.
- The peephole is a small, moving window on the target program. The code in the peephole need not be contiguous, although some implementations do require this. It is characteristic of peephole optimization that each improvement may spawn opportunities for additional improvements.
- We shall give the following examples of program transformations that are characteristic of peephole optimizations:
 - ✓ Redundant-instructions elimination
 - ✓ Flow-of-control optimizations
 - ✓ Algebraic simplifications
 - ✓ Use of machine idioms
 - ✓ Unreachable Code

Redundant Loads and Stores:

If we see the instructions sequence

(1) MOV R0, a

(2) MOV a,R0

we can delete instructions (2) because whenever (2) is executed. (1) will ensure that the value of a is already in register R0. If (2) had a label we could not be sure that (1) was always executed immediately before (2) and so we could not remove (2).

Unreachable Code:

- Another opportunity for peephole optimizations is the removal of unreachable instructions. An unlabeled instruction immediately following an unconditional jump may be removed. This operation can be repeated to eliminate a sequence of instructions. For example, for debugging purposes, a large program may have within it certain segments that are executed only if a variable debug is 1. In C, the source code might look like:

```
#define debug 0
```

```
....
```

```
If ( debug ) {  
Print debugging information  
}
```

- In the intermediate representations the if-statement may be translated as:

```
If debug =1 goto L2 goto L2
```

```
L1: print debugging information
```

```
L2: ..... (a)
```

- One obvious peephole optimization is to eliminate jumps over jumps. Thus no matter what the value of debug; (a) can be replaced by:

```
If debug ≠1 goto L2
```

```
Print debugging information
```

```
L2: ..... (b)
```

- As the argument of the statement of (b) evaluates to a constant true it can be replaced by

```
36
```

```
If debug ≠0 goto L2
```

```
Print debugging information
```

```
L2: ..... (c)
```

- As the argument of the first statement of (c) evaluates to a constant true, it can be replaced by goto L2. Then all the statement that print debugging aids are manifestly unreachable and can be eliminated one at a time.

Flows-Of-Control Optimizations:

- The unnecessary jumps can be eliminated in either the intermediate code or the target code by the following types of peephole optimizations. We can replace the jump sequence

```
goto L1
```

```
....
```

```
L1: gotoL2 by the
```

```
sequence
```

```
goto L2
```

```
....
```

```
L1: goto L2
```

- If there are now no jumps to L1, then it may be possible to eliminate the statement L1:goto L2 provided it is preceded by an unconditional jump. Similarly, the sequence

```
if a < b goto L1
```

```
....
```

```
L1: goto L2 can be
```

replaced by
If a < b goto L2

....

L1: goto L2

- Finally, suppose there is only one jump to L1 and L1 is preceded by an unconditional goto. Then the sequence goto L1

.....

37

L1: if a < b goto L2

L3: (1)

- May be replaced by
If a < b goto L2 goto L3

.....

L3: (2)

- While the number of instructions in (1) and (2) is the same, we sometimes skip the unconditional jump in (2), but never in (1). Thus (2) is superior to (1) in execution time

Algebraic Simplification:

- There is no end to the amount of algebraic simplification that can be attempted through peephole optimization. Only a few algebraic identities occur frequently enough that it is worth considering implementing them. For example, statements such as

x := x+0 Or

x := x * 1

- Are often produced by straightforward intermediate code-generation algorithms, and they can be eliminated easily through peephole optimization.

Reduction in Strength:

- Reduction in strength replaces expensive operations by equivalent cheaper ones on the target machine. Certain machine instructions are considerably cheaper than others and can often be used as special cases of more expensive operators.
- For example, x^2 is invariably cheaper to implement as $x*x$ than as a call to an exponentiation routine. Fixed-point multiplication or division by a power of two is cheaper to implement as a shift. Floating-point division by a constant can be implemented as multiplication by a constant, which may be cheaper.

$X_2 \rightarrow X*X$

Use of Machine Idioms:

- The target machine may have hardware instructions to implement certain specific operations efficiently. For example, some machines have auto-increment and auto-decrement addressing modes. These add or subtract one from an operand before or after using its value.
- The use of these modes greatly improves the quality of code when pushing or popping a stack, as in parameter passing. These modes can also be used in code for statements like
i := i+1.

i:=i+1 \rightarrow i++

i:=i-1 \rightarrow i--

Techniques in local loops and global optimization:

Loop Optimization is the process of increasing execution speed and reducing the overheads associated with loops. It plays an important role in improving cache performance and making effective use of parallel processing capabilities. Most execution time of a scientific program is spent on loops.

Loop Optimization is a machine independent optimization.

Decreasing the number of instructions in an inner loop improves the running time of a program even if the amount of code outside that loop is increased.

Loop Optimization Techniques

In the compiler, we have various loop optimization techniques, which are as follows:

1. Code Motion (Frequency Reduction)

In frequency reduction, the amount of code in the loop is decreased. A statement or expression, which can be moved outside the loop body without affecting the semantics of the program, is moved outside the loop.

Example:

Before optimization:

```
while(i<100)
{
  a = Sin(x)/Cos(x) + i;
  i++;
}
```

After optimization:

```
t = Sin(x)/Cos(x);
while(i<100)
{
  a = t + i;
  i++;
}
```

2. Induction Variable Elimination

If the value of any variable in any loop gets changed every time, then such a variable is known as an induction variable. With each iteration, its value either gets incremented or decremented by some constant value.

Example:

Before optimization:

```
B1
i:= i+1
x:= 3*i
y:= a[x]
if y< 15, goto B2
```

In the above example, i and x are linked, if i is incremented by 1 then x is incremented by 3. So, i and x are induction variables.

After optimization:

```
B1
i:= i+1
x:= x+4
y:= a[x]
```


if $y < 15$, goto B2

3. Strength Reduction

Strength reduction deals with replacing expensive operations with cheaper ones like multiplication is costlier than addition, so multiplication can be replaced by addition in the loop.

Example:

Before optimization:

```
while (x<10)
{
y := 3 * x+1;
a[y] := a[y]-2;
x := x+2;
}
```

After optimization:

```
t= 3 * x+1;
while (x<10)
{
y=t;
a[y]= a[y]-2;
x=x+2;
t=t+6;
}
```

4. Loop Invariant Method

In the loop invariant method, the expression with computation is avoided inside the loop. That computation is performed outside the loop as computing the same expression each time was overhead to the system, and this reduces computation overhead and hence optimizes the code.

Example:

Before optimization:

```
for (inti=0; i<10;i++)
t= i+(x/y);
...
end;
```

After optimization:

```
s = x/y;
for (inti=0; i<10;i++)
t= i+ s;
...
end;
```

5. Loop Unrolling

Loop unrolling is a loop transformation technique that helps to optimize the execution time of a program. We basically remove or reduce iterations. Loop unrolling increases the program's speed by eliminating loop control instruction and loop test instructions.

Example:

Before optimization:

```
for (inti=0; i<5; i++)  
printf("Pankaj\n");
```

After optimization:

```
printf("Pankaj\n");  
printf("Pankaj\n");  
printf("Pankaj\n");  
printf("Pankaj\n");  
printf("Pankaj\n");
```

6. Loop Jamming

Loop jamming is combining two or more loops in a single loop. It reduces the time taken to compile the many loops.

Example:

Before optimization:

```
for(inti=0; i<5; i++)  
    a = i + 5;  
for(inti=0; i<5; i++)  
    b = i + 10;
```

After optimization:

```
for(inti=0; i<5; i++)  
{  
    a = i + 5;  
    b = i + 10;  
}
```

7. Loop Fission

Loop fission improves the locality of reference, in loop fission a single loop is divided into multiple loops over the same index range, and each divided loop contains a particular part of the original loop

Example:

Before optimization:

```
for(x=0;x<10;x++)  
{  
    a[x]=...  
    b[x]=...  
}
```

After optimization:

```
for(x=0;x<10;x++)  
    a[x]=...  
for(x=0;x<10;x++)  
    b[x]=...
```

8. Loop Interchange

In loop interchange, inner loops are exchanged with outer loops. This optimization technique also improves the locality of reference.

Example:

Before optimization:

```
for(x=0;x<10;x++)  
for(y=0;y<10;y++)  
a[y][x]=...
```

After optimization:

```
for(y=0;y<10;y++)  
for(x=0;x<10;x++)  
a[y][x]=...
```

9. Loop Reversal

Loop reversal reverses the order of values that are assigned to index variables. This help in removing dependencies.

Example:

Before optimization:

```
for(x=0;x<10;x++)  
a[9-x]=...
```

After optimization:

```
for(x=9;x>=0;x--)  
a[x]=...
```

10. Loop Splitting

Loop Splitting simplifies a loop by dividing it into numerous loops, and all the loops have some bodies but they will iterate over different index ranges. Loop splitting helps in reducing dependencies and hence making code more optimized.

Example:

Before optimization

```
for(x=0;x<10;x++)  
if(x<5)  
a[x]=...  
else  
b[x]=...
```

After optimization:

```
for(x=0;x<5;x++)  
a[x]=...  
for(;x<10;x++)  
b[x]=...
```

11. Loop Peeling

Loop peeling is a special case of loop splitting, in which a loop with problematic iteration is resolved separately before entering the loop.

Before optimization:

```
for(x=0;x<10;x++)  
if(x==0)  
a[x]=...  
else  
b[x]=...
```

After optimization:

```
a[0]=...  
for(x=1;x<100;x++)  
b[x]=...
```

12. Unswitching

Unswitching moves a condition out from inside the loop, this is done by duplicating loop and placing each of its versions inside each conditional clause.

Before optimization:

```
for(x=0;x<10;x++)  
if(s>t)  
a[x]=...  
else  
b[x]=...
```

After optimization:

```
if(s>t)  
for(x=0;x<10;x++)  
a[x]=...  
else  
for(x=0;x<10;x++)  
b[x]=...
```