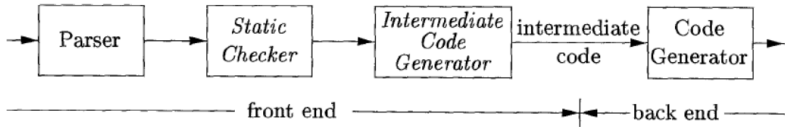


# Code generation & Optimization



# Code generation

## Input to the Code Generator

- (a) **Intermediate representation (IR)** of the source program produced by the front end,
- (b) **Symbol table** that is used to determine the **run-time** addresses of the data objects

## The Target Program

### Instruction-set architecture of the target machine

- Has a significant impact on the **difficulty of constructing a good code generator** that produces high-quality machine code.
- The most common target-machine architectures are RISC (reduced instruction set computer), CISC (complex instruction set computer)
- A RISC machine typically has many registers, three-address instructions, simple addressing modes, and a relatively simple instruction-set architecture.

## Instruction Selection

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

The **complexity** of performing this **mapping** is determined by a factors such as

- the level of the IR
  - the nature of the **instruction-set architecture**
  - the **desired quality** of the **generated code**.
- 
- If the **IR is high level**, the code generator may translate each IR statement into a sequence of machine instructions **using code templates**.
  - Such statement by statement code generation, however, often produces **poor code**

**Translation scheme:** If we do not care about the efficiency of the target program, instruction selection is straightforward.

For **each type of three-address statement**, we can design a **code skeleton** that defines the **target code** to be generated for that construct

# Instruction Selection – Example


example, every three-address statement of the form  $x = y + z$ , where  $x$ ,  $y$ , and  $z$  are statically allocated, can be translated into the code sequence

```
LD  R0, y      // R0 = y      (load y into register R0)
ADD R0, R0, z   // R0 = R0 + z (add z to R0)
ST  x, R0      // x = R0      (store R0 into x)
```

```
a = b + c
d = a + e
```

would be translated into

```
LD  R0, b      // R0 = b
ADD R0, R0, c   // R0 = R0 + c
ST  a, R0      // a = R0
LD  R0, a      // R0 = a
ADD R0, R0, e   // R0 = R0 + e
ST  d, R0      // d = R0
```

 Redundant

# Instruction Selection – Example

- On most machines, a **given IR program** can be implemented by **many different code sequences**,
  - Significant **cost differences** between the **different implementations**.
- A **naive translation** of the intermediate code may therefore **lead to correct but unacceptably inefficient target code**.
- For example, if the target machine has an **"increment" instruction (INC)**
- The three-address statement  **$a = a + 1$**  may be implemented more efficiently by the **single instruction INC a**,
  - Rather than by a more obvious sequence that loads **a** into a register, adds one to the register, and then stores the result back into **a**

```
LD  R0, a      // R0 = a
ADD R0, R0, #1  // R0 = R0 + 1
ST  a, R0      // a = R0
```

# Register Allocation

- A **key problem** in code generation is **deciding what values to hold in what registers**.
- Registers are the **fastest** computational unit on the target machine,
  - but we usually do not have **enough** of them to hold **all values**.
  - Values not held in registers need to reside **in memory**.
- **Instructions involving register** operands are invariably **shorter and faster** than those involving **operands in memory**,
  - **Efficient utilization of registers** is particularly important.

The **use of registers** is often subdivided into **two subproblems**:

1. **Register allocation**, during which we **select the set of variables** that will **reside in registers** at each point in the program.
2. **Register assignment**, during which we **pick the specific register** that a variable will reside in.

# Basic Blocks & Flow graphs

- Introduce a **graph representation of intermediate code** that is helpful for discussing code generation
  - Even if the graph is not constructed explicitly by a code-generation algorithm.
- Code generation **benefits from context**.
- We can do a **better job of register allocation** if we know how variables are **defined and used**.

# Basic Blocks & Flow graphs

The representation is constructed as follows:

1. Partition the intermediate code into *basic blocks*, which are maximal sequences of consecutive three-address instructions with the properties that
  - (a) The flow of control can only enter the basic block through the first instruction in the block. That is, there are no jumps into the middle of the block.
  - (b) Control will leave the block without halting or branching, except possibly at the last instruction in the block.
2. The basic blocks become the nodes of a *flow graph*, whose edges indicate which blocks can follow which other blocks.



# Basic Blocks

- We begin a **new basic block** with the **first instruction**
- Keep adding instructions
  - until we meet either a **jump, a conditional jump,**
  - or a **label** on the following instruction.
- In the **absence of jumps and labels**, control proceeds **sequentially** from one instruction to the next.
- **Task:** *Identify leaders*, that is, the **first instructions** in some **basic block**.

# Basic Blocks - Leaders

1. The first three-address instruction in the intermediate code is a leader.
2. Any instruction that is the target of a conditional or unconditional jump is a leader.
3. Any instruction that immediately follows a conditional or unconditional jump is a leader.

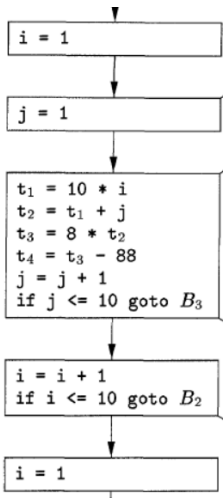
# Basic Blocks

```
for i from 1 to 10 do  
    for j from 1 to 10 do  
         $a[i, j] = 0.0;$   
for i from 1 to 10 do  
     $a[i, i] = 1.0;$ 
```

**leaders are instructions**  
**1, 2, 3, 10, 12, and 13**

```
1)   $i = 1$   
2)   $j = 1$   
3)   $t1 = 10 * i$   
4)   $t2 = t1 + j$   
5)   $t3 = 8 * t2$   
6)   $t4 = t3 - 88$   
7)   $a[t4] = 0.0$   
8)   $j = j + 1$   
9)  if  $j \leq 10$  goto (3)  
10)  $i = i + 1$   
11) if  $i \leq 10$  goto (2)  
12)  $i = 1$   
13)  $t5 = i - 1$   
14)  $t6 = 88 * t5$   
15)  $a[t6] = 1.0$   
16)  $i = i + 1$   
17) if  $i \leq 10$  goto (13)
```

# Basic Blocks



# Flow Graphs

- We represent the **flow of control** by a **flow graph**.
- The **nodes** of the flow graph are the **basic blocks**.
- There is an **edge from block B to block C** if and only if
  - it is possible for the **first instruction in block C** to immediately follow the **last instruction in block B**.

There are two ways that such an edge could be justified:

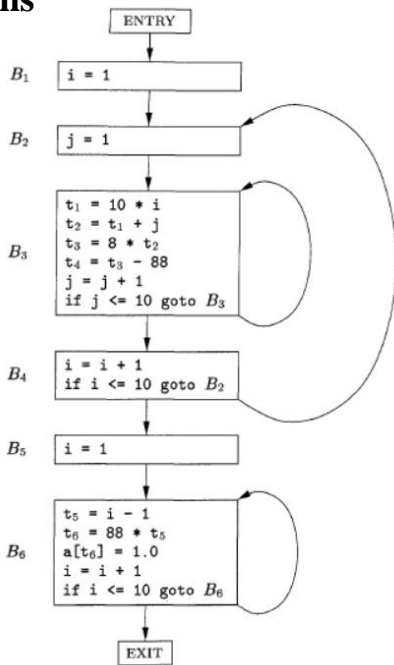
- There is a **conditional or unconditional jump** from the end of B to the beginning of C.
- **Block C immediately follows Block B** in the original order of the three-address instructions
  - B does not end in an unconditional jump
  - **Maybe due to labels**

We say that **B is a predecessor of C**, and **C is a successor of B**.

# Flow Graphs

- Often we add two nodes, called the **entry** and **exit**,
- There is an **edge from the entry** to the **first executable node** of the flow graph,
  - that is, to the **basic block** that comes from the **first instruction** of the intermediate code.
- There is an edge **to the exit** from **any basic block** that contains an instruction that could be the **last executed instruction** of the program.

# Flow Graphs



# Loops

Many code transformations depend upon the identification of "loops" in a flow graph. We say that a set of nodes  $L$  in a flow graph is a loop if

1. There is a node in  $L$  called the *loop entry* with the property that no other node in  $L$  has a predecessor outside  $L$ . That is, every path from the entry of the entire flow graph to any node in  $L$  goes through the loop entry.
2. Every node in  $L$  has a nonempty path, completely within  $L$ , to the entry of  $L$ .



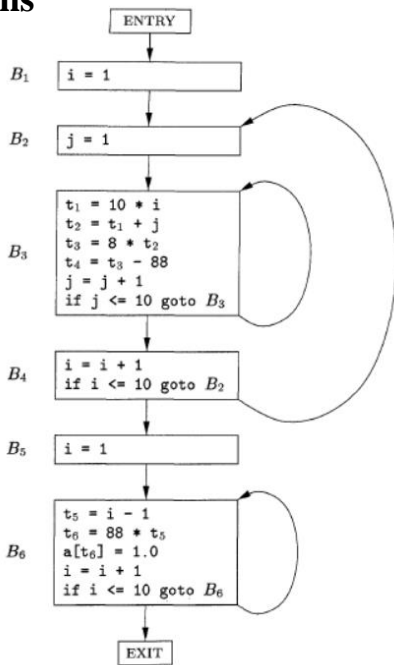
# Loops

1.  $B_3$  by itself.
2.  $B_6$  by itself.
3.  $\{B_2, B_3, B_4\}$ .

The first two are single nodes with an edge to the node itself. For instance,  $B_3$  forms a loop with  $B_3$  as its entry. Note that the second requirement for a loop is that there be a nonempty path from  $B_3$  to itself. Thus, a single node like  $B_2$ , which does not have an edge  $B_2 \rightarrow B_2$ , is not a loop, since there is no nonempty path from  $B_2$  to itself within  $\{B_2\}$ .

The third loop,  $L = \{B_2, B_3, B_4\}$ , has  $B_2$  as its loop entry. Note that among these three nodes, only  $B_2$  has a predecessor,  $B_1$ , that is not in  $L$ . Further, each of the three nodes has a nonempty path to  $B_2$  staying within  $L$ . For instance,  $B_4$  has the path  $B_4 \rightarrow B_3 \rightarrow B_2$ .

# Flow Graphs



# Next-Use Information

- Knowing when the value of a **variable will be used next** is essential for generating good code.
  - If the value of a variable that is currently in a register **will never be referenced subsequently**, then that register can be **re-assigned to another variable**.
- Suppose **three-address statement i assigns a value to x**.
- If **statement j has x as an operand**, and control can flow from statement i to j along a path that has no intervening assignments to x, then we say **statement j uses the value of x computed at statement i**.
  - We further say that **x is live at statement i**.
- We wish to determine for **each three-address statement  $x = y + z$  what the next uses of x, y, and z are**.
  - We store the information in the symbol table.
- Focus on the basic block containing this three-address statement.
- The algorithm to determine **liveness and next-use** information **makes a backward pass over each basic block**.

# Next-Use Information

**INPUT:** A basic block  $B$  of three-address statements. We assume that the symbol table initially shows all nontemporary variables in  $B$  as being live on exit.

**OUTPUT:** At each statement  $i: x = y + z$  in  $B$ , we attach to  $i$  the liveness and next-use information of  $x$ ,  $y$ , and  $z$ .

**METHOD:** We start at the last statement in  $B$  and scan backwards to the beginning of  $B$ . At each statement  $i: x = y + z$  in  $B$ , we do the following:

1. Attach to statement  $i$  the information currently found in the symbol table regarding the next use and liveness of  $x$ ,  $y$ , and  $y$ .
2. In the symbol table, set  $x$  to “not live” and “no next use.”
3. In the symbol table, set  $y$  and  $z$  to “live” and the next uses of  $y$  and  $z$  to  $i$ .

# Code Generator

- Algorithm that generates code for a **single basic block**
- It considers **each three-address instruction** in turn, and keeps track of **what values are in what registers** so it can avoid generating **unnecessary loads and stores**.
- Deciding how to **use registers to best advantage**
- In most machine architectures, some or all of the **operands** of an operation must be in **registers** in order to perform the operation.
- These are competing needs, since the number of registers available is **limited**.

# Code Generator

- We further assume that for each operator, there is exactly one machine instruction that takes the necessary operands in registers and performs that operation, leaving the result in a register. The machine instructions are of the form

*LD reg, mem*

*ST mem, reg*

*OP reg, reg, reg*

# Register and Address Descriptors

- Our 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.
- We require a **data structure** that tells us what program **variables** currently have **their value in a register, and in which register**
- We also need to know whether the **memory location for a given variable currently has the proper value for that variable**
  - Since a new value for the variable may have been computed in a register and not yet stored.

# Register Descriptors

- a **register descriptor** keeps track of the variable names whose current value is in that register.
- All register descriptors are empty. As the code generation progresses, each register will hold the value of zero or more names.

# Address Descriptors

For each program variable, an **address descriptor** keeps track of the **location or locations** where the current value of that variable can be found.

The **location** might be a **register, a memory address** etc.

The information can be stored in the **symbol-table entry** for that **variable name**.





# The Code-Generation Algorithm

- An essential part of the algorithm is a **function getReg(I)**,
  - which selects registers for each memory location associated with the three-address instruction I.
- Function **getReg** has access to the **register and address descriptors** for all the variables of the basic block
- While we do not know the total number of registers available for local data belonging to a basic block, we assume that there are **enough registers**

## Machine Instructions for Operations

For a three-address instruction such as  $x = y + z$ , do the following:

1. Use  $getReg(x = y + z)$  to select registers for  $x$ ,  $y$ , and  $z$ . Call these  $R_x$ ,  $R_y$ , and  $R_z$ .
2. If  $y$  is not in  $R_y$  (according to the register descriptor for  $R_y$ ), then issue an instruction  $LD\ R_y, y'$ , where  $y'$  is one of the memory locations for  $y$  (according to the address descriptor for  $y$ ).
3. Similarly, if  $z$  is not in  $R_z$ , issue an instruction  $LD\ R_z, z'$ , where  $z'$  is a location for  $z$ .
4. Issue the instruction  $ADD\ R_x, R_y, R_z$ .

# Machine Instructions for Copy Statements

**three-address copy statement of the form  $x = y$ .**

- We assume that **getReg** will always choose the **same register** for **both  $x$  and  $y$**
- If  $y$  is **not** already in that register  **$R_y$** ,
  - then generate the machine instruction **LD  $R_y, y$** .
  - If  $y$  was already in  $R_y$ , we **do nothing**.
- It is only necessary that we adjust the **register description** for  **$R_y$** 
  - So that it **includes  $x$**  as one of the values found there.

# Ending the Basic Block

- If the **variable is live on exit** from the block,
  - Or if we **don't know** which variables are live on exit,
  - then we assume that the **value of the variable** is **needed later**.
- In that case, for **each variable x** whose **address descriptor does not say** that its value is located in the **memory location for x**
- We must generate the instruction **ST x, R**, where **R is a register in which x value** exists at the end of the block.

# Managing Register and Address Descriptors

- As the code-generation algorithm issues load, store, and other machine instructions,
- It needs to **update the register and address descriptors**.
- The rules are as follows:

1. For the instruction LD  $R, x$



- (a) Change the register descriptor for register  $R$  so it holds only  $x$ .
- (b) Change the address descriptor for  $x$  by adding register  $R$  as an additional location.

2. For the instruction ST  $x, R$ , change the address descriptor for  $x$  to include its own memory location.

3. For an operation such as ADD  $R_x, R_y, R_z$  implementing a three-address instruction  $x = y + z$



- (a) Change the register descriptor for  $R_x$  so that it holds only  $x$ .
- (b) Change the address descriptor for  $x$  so that its only location is  $R_x$ . Note that the memory location for  $x$  is *not* now in the address descriptor for  $x$ .
- (c) Remove  $R_x$  from the address descriptor of any variable other than  $x$ .



# Managing Register and Address Descriptors

When we process a copy statement  $x = y$ , after generating the load for  $y$  into register  $R_y$ , if needed, and after managing descriptors as for all load statements (per rule 1):

- (a) Add  $x$  to the register descriptor for  $R_y$ .
- (b) Change the address descriptor for  $x$  so that its only location is  $R_y$ .



t = a - b

u = a - c

v = t + u

a = d

d = v + u

t = a - b

LD R1, a

LD R2, b

SUB R2, R1, R2

u = a - c

LD R3, c

SUB R1, R1, R3

v = t + u

ADD R3, R2, R1

a = d

LD R2, d

R1 R2 R3

--	--	--

a b c d t u v

a	b	c	d			
---	---	---	---	--	--	--

a	t	
---	---	--

a,R1	b	c	d	R2		
------	---	---	---	----	--	--

u	t	c
---	---	---

a	b	c,R3	d	R2	R1	
---	---	------	---	----	----	--

u	t	v
---	---	---

a	b	c	d	R2	R1	R3
---	---	---	---	----	----	----

u	a,d	v
---	-----	---

R2	b	c	d,R2		R1	R3
----	---	---	------	--	----	----

$t = a - b$

$u = a - c$

$v = t + u$

$a = d$

$d = v + u$

R1    R2    R3            a    b    c    d    t    u    v

a = d

LD R2, d

u	t	v
---	---	---

a	b	c	d	R2	R1	R3
---	---	---	---	----	----	----

u	a,d	v
---	-----	---

R2	b	c	d,R2		R1	R3
----	---	---	------	--	----	----

d = v + u

ADD R1, R3, R1

d	a	v
---	---	---

R2	b	c	R1			R3
----	---	---	----	--	--	----

exit

ST a, R2

ST d, R1

d	a	v
---	---	---

a,R2	b	c	d,R1			R3
------	---	---	------	--	--	----



# Design of the Function *getReg*

- There are many options,
- although there are also some **absolute prohibitions** against choices that lead to incorrect code due to the loss of the value of one or more live variables
- We use  $x = y + z$  as the generic example.
- First, we must **pick a register for y** and **a register for z**.
- The issues are the same, so we shall concentrate on picking register **Ry for y**.

1. If  $y$  is currently in a register, pick a register already containing  $y$  as  $R_y$ . Do not issue a machine instruction to load this register, as none is needed.
2. If  $y$  is not in a register, but there is a register that is currently empty, pick one such register as  $R_y$ .
3. The difficult case occurs when  $y$  is not in a register, and there is no register that is currently empty. We need to pick one of the allowable registers anyway, and we need to make it safe to reuse. Let  $R$  be a candidate register, and suppose  $v$  is one of the variables that the register descriptor for  $R$  says is in  $R$ . We need to make sure that  $v$ 's value either is not really needed, or that there is somewhere else we can go to get the value of  $R$ . The possibilities are:

- (a) If the address descriptor for  $v$  says that  $v$  is somewhere besides  $R$ , then we are OK.
- (b) If  $v$  is  $x$ , the value being computed by instruction  $I$ , and  $x$  is not also one of the other operands of instruction  $I$  ( $z$  in this example), then we are OK. The reason is that in this case, we know this value of  $x$  is never again going to be used, so we are free to ignore it.
- (c) Otherwise, if  $v$  is not used later (that is, after the instruction  $I$ , there are no further uses of  $v$ , and if  $v$  is live on exit from the block, then  $v$  is recomputed within the block), then we are OK.
- (d) If we are not OK by one of the first two cases, then we need to generate the store instruction  $ST\ v, R$  to place a copy of  $v$  in its own memory location. This operation is called a *spill*.

Since  $R$  may hold several variables at the moment, we repeat the above steps for each such variable  $v$ . At the end,  $R$ 's "score" is the number of store instructions we needed to generate. Pick one of the registers with the lowest score.

# Selection of the register $R_x$

Now, consider the selection of the register  $R_x$ . The issues and options are almost as for  $y$ , so we shall only mention the differences.

1. Since a new value of  $x$  is being computed, a register that holds only  $x$  is always an acceptable choice for  $R_x$ . This statement holds even if  $x$  is one of  $y$  and  $z$ , since our machine instructions allows two registers to be the same in one instruction.
2. If  $y$  is not used after instruction  $I$ , in the sense described for variable  $v$  in item (3c), and  $R_y$  holds only  $y$  after being loaded, if necessary, then  $R_y$  can also be used as  $R_x$ . A similar option holds regarding  $z$  and  $R_z$ .

The last matter to consider specially is the case when  $I$  is a copy instruction  $x = y$ . We pick the register  $R_y$  as above. Then, we always choose  $R_x = R_y$ .

# Peephole Optimization

- Improve the quality of the target code by applying "**optimizing transformations**" to the target program
- Peephole optimization is done by **examining a sliding window** of target instructions (called the peephole) and
- **Replacing instruction** sequences within the peephole by a **shorter or faster sequence**,

## Eliminating Redundant Loads and Stores

```
LD a, R0  
ST R0, a
```

- Redundant loads and stores of this nature **would not be generated** by the **simple code generation** algorithm

However, a **naive code generation** algorithm would generate redundant sequences

- Note that **if the store instruction had a label**, we could not be sure that **the first instruction is always executed before the second**, so we could not remove the store instruction.

# Eliminating Unreachable Code

An **unlabeled instruction immediately following an unconditional jump** may be removed.

One obvious peephole optimization is to eliminate jumps over jumps

```
    if debug == 1 goto L1
    goto L2
L1: print debugging information
L2:
```

## After optimization

```
    if debug != 1 goto L2
    print debugging information
L2:
```

If `debug` is set to 0 at the beginning of the program, constant propagation would transform this sequence into

```
    if 0 != 1 goto L2
    print debugging information
L2:
```

# Flow-of-Control Optimizations

- Simple intermediate code-generation algorithms frequently produce
  - jumps to jumps, jumps to conditional jumps, or conditional jumps to jumps.

```
goto L1  
...
```

```
L1: goto L2
```

by the sequence

```
goto L2  
...
```

```
L1: goto L2
```



**Remove this, if no other jump at L1**


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

```
L1: goto L2
```

can be replaced by the sequence

```
    if a < b goto L2
    ...
```

```
L1: goto L2
```

 **Remove this, if no other jump at L1**

```
goto L1
```

```
...
```

```
L1: if a < b goto L2
```

```
L3:
```

may be replaced by the sequence

```
    if a < b goto L2
```

```
goto L3
```

```
...
```

```
L3:
```



- While the **number of instructions** in the two sequences is the **same**,
  - we **sometimes skip** the **unconditional jump** in the **second** sequence,
  - but never in the first sequence .



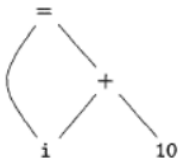
# Optimization of Basic Blocks

## DAG Representation of Basic Blocks

1. There is a node in the DAG for each of the initial values of the variables appearing in the basic block.
2. There is a node  $N$  associated with each statement  $s$  within the block. The children of  $N$  are those nodes corresponding to statements that are the last definitions, prior to  $s$ , of the operands used by  $s$ .
3. Node  $N$  is labeled by the operator applied at  $s$  and also attached to  $N$  is the list of variables for which it is the last definition within the block.
4. Certain nodes are designated *output nodes*. These are the nodes whose variables are *live on exit* from the block; that is, their values may be used later, in another block of the flow graph.

## DAG Representation of Basic Blocks: Value Number Method

$i = i + 10$



(a) DAG

1	id			to entry for i
2	num	10		
3	+	1	2	
4	=	1	3	
5		...		

(b) Array.

- The nodes of a syntax tree or DAG are stored in an array of records
- In this array, we refer to nodes by giving the integer index of the record for that node within the array.
- Each row of the array represents one record, and therefore one node.
- In each record, the first field is an operation code, indicating the label of the node.

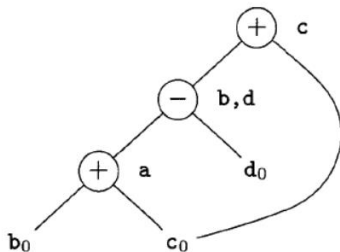
# Finding Local Common Subexpressions

$$a = b + c$$

$$b = a - d$$

$$c = b + c$$

$$d = a - d$$



**Common subexpressions** can be detected by noticing, as a new node ***M*** is about to be added, whether there is an **existing node *N*** with the same children, in the same order, and with the same operator.

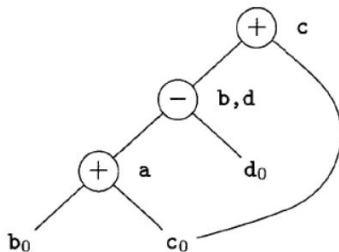
# Finding Local Common Subexpressions

$a = b + c$

$b = a - d$

$c = b + c$

$d = a - d$



Since there are **only three nonleaf nodes** in the DAG, the basic block in can be replaced by a block with only **three statements**.

$a = b + c$

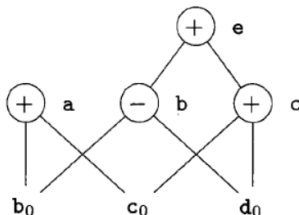
$d = a - d$

$c = d + c$

- If **b** is not live on exit from the block, then we do not need to compute **b** variable, and can use **d** to receive the value
- If both **b** and **d** are live on exit, then a fourth statement must be used to copy the value from one to the other

# Dead Code Elimination

```
a = b + c;  
b = b - d  
c = c + d  
e = b + c
```



If **a** and **b** are live but **c** and **e** are not, we can immediately **remove the root labeled e**.

Then, the node labeled **c** becomes a root and can be removed.

The roots labeled **a** and **b** remain, since they each have live variables attached

# The Use of Algebraic Identities

$$x + 0 = 0 + x = x$$

$$x \times 1 = 1 \times x = x$$

$$x - 0 = x$$

$$x/1 = x$$

*reduction in strength,*

EXPENSIVE

CHEAPER

$$x^2 = x \times x$$

$$2 \times x = x + x$$

$$x/2 = x \times 0.5$$

*constant folding.*

Thus the expression  $2 * 3.14$  would be replaced by 6.28.

$*$  is commutative

$$x * y = y * x.$$

# The Use of Algebraic Identities

**Associative laws** might also be applicable to **expose common subexpressions**

$$a = b + c$$

$$t = c + d$$

$$e = t + b$$

If  $t$  is not needed outside this block, we can change this sequence to

$$a = b + c$$

$$e = a + d$$

# Representation of Array References

$x = a[i]$

$a[j] = y$

$z = a[i]$

$z = x \quad ???$

An assignment from an array, like  $x = a[i]$ , is represented by creating a node with operator  $=[]$  and two children representing the initial value of the array,  $a_0$  in this case, and the index  $i$ . Variable  $x$  becomes a label of this new node.

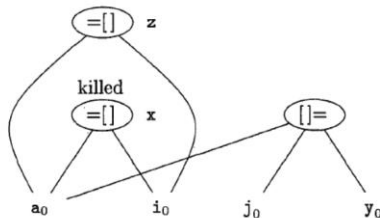
An assignment to an array, like  $a[j] = y$ , is represented by a new node with operator  $[]=$  and three children representing  $a_0$ ,  $j$  and  $y$ . There is no variable labeling this node.

this node *kills* all currently constructed nodes whose value depends on  $a_0$ . A node that has been killed cannot receive any more labels; that is, it cannot become a common subexpression.



# Representation of Array References

```
x = a[i]
a[j] = y
z = a[i]
```



An assignment from an array, like `x = a[i]`, is represented by creating a node with operator `=[]` and two children representing the initial value of the array, `a0` in this case, and the index `i`. Variable `x` becomes a label of this new node.

An assignment to an array, like `a[j] = y`, is represented by a new node with operator `[]=` and three children representing `a0`, `j` and `y`. There is no variable labeling this node.

this node *kills* all currently constructed nodes whose value depends on `a0`. A node that has been killed cannot receive any more labels; that is, it cannot become a common subexpression.

# Optimal Code Generation for Expressions

- We can **choose registers** optimally when a **basic block** consists of a **single expression**
- We introduce a labelling scheme for the nodes of an expression tree

Two phases

- Labelling phase
- Code generation phase

# The Labelling Algorithm : **Ershov Numbers**

Labels each node of the tree with an integer:

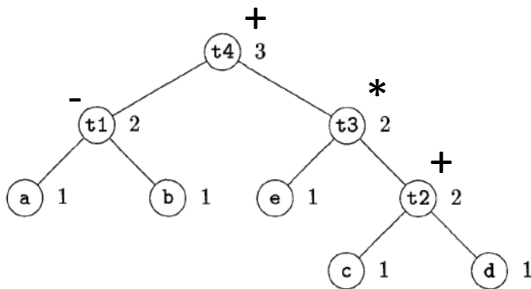
- No. of registers required to evaluate the tree with no intermediate stores to memory
- Consider binary trees

1. Label any leaf 1.
2. The label of an interior node with one child is the label of its child.
3. The label of an interior node with two children is
  - (a) The larger of the labels of its children, if those labels are different.
  - (b) One plus the label of its children if the labels are the same.

# The Labelling Algorithm : Ershov Numbers

tree for expression  $(a - b) + e \times (c + d)$

t1 = a - b  
t2 = c + d  
t3 = e \* t2  
t4 = t1 + t3



# Generating Code From Labeled Expression Trees

- **All operands must be in registers**, and registers can be used by both an **operand** and the **result** of an operation
- The **label** of a node is the **fewest registers** with which the expression can be evaluated using no stores of temporary results.

# Generating Code From Labeled Expression Trees

**METHOD:** The following is a recursive algorithm to generate the machine code. The steps below are applied, starting at the root of the tree. If the algorithm is applied to a node with label  $k$ , then only  $k$  registers will be used. However, there is a "base"  $b > 1$  for the registers used so that the actual registers used are  $R_b, R_{b+1}, \dots, R_{b+k-1}$ . The result always appears in  $R_{b+k-1}$ .

1. To generate machine code for an interior node with label  $k$  and two children with equal labels (which must be  $k-1$ ) do the following:
  - (a) Recursively generate code for the right child, using base  $b+1$ . The result of the right child appears in register  $R_{b+k}$ .
  - (b) Recursively generate code for the left child, using base  $b$ ; the result appears in  $R_{b+k-1}$ .
  - (c) Generate the instruction  $OP\ R_{b+k}, R_{b+k-1}, R_{b+k}$ , where  $OP$  is the appropriate operation for the interior node in question.

$K=4, b=1$

Right child ( $k-1$ )  $\rightarrow R_2, R_3, R_4$  (result)

Left child ( $k-1$ )  $\rightarrow R_1, R_2, R_3$  (result)

# Generating Code From Labeled Expression Trees

**METHOD:** The following is a recursive algorithm to generate the machine code. The steps below are applied, starting at the root of the tree. If the algorithm is applied to a node with label  $k$ , then only  $k$  registers will be used. However, there is a “base”  $b > 1$  for the registers used so that the actual registers used are  $R_b, R_{b+1}, \dots, R_{b+k-1}$ . The result always appears in  $R_{b+k-1}$ .

1. To generate machine code for an interior node with label  $k$  and two children with equal labels (which must be  $k - 1$ ) do the following:
  - (a) Recursively generate code for the right child, using base  $b + 1$ . The result of the right child appears in register  $R_{b+k}$ .
  - (b) Recursively generate code for the left child, using base  $b$ ; the result appears in  $R_{b+k-1}$ .
  - (c) Generate the instruction  $OP\ R_{b+k}, R_{b+k-1}, R_{b+k}$ , where  $OP$  is the appropriate operation for the interior node in question.

For a leaf representing operand  $x$ , if the base is  $b$  generate the instruction  $LD\ R_b, x$ .

# Generating Code From Labeled Expression Trees



K-1 Registers for right child  $R(b+1), R(b+2), \dots, R(b+k)$

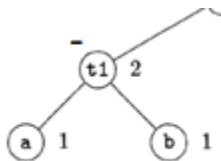
K-1 Registers for right child  $R(b), R(b+1), \dots, R(b+k-1)$



Result



# Generating Code From Labeled Expression Trees



Base b=1

```
LD R2, b
LD R1, a
SUB R2, R1, R2
```

# Generating Code From Labeled Expression Trees

2. Suppose we have an interior node with label  $k$  and children with unequal labels. Then one of the children, which we'll call the "big" child, has label  $k$ , and the other child, the "little" child, has some label  $m < k$ . Do the following to generate code for this interior node, using base  $b$ :
- (a) Recursively generate code for the big child, using base  $b$ ; the result appears in register  $R_{b+k-1}$ .
  - (b) Recursively generate code for the small child, using base  $b$ ; the result appears in register  $R_{b+m-1}$ . Note that since  $m < k$ , neither  $R_{b+k-1}$  nor any higher-numbered register is used.
  - (c) Generate the instruction  $OP\ R_{b+k-1}, R_{b+m-1}, R_{b+k-1}$  or the instruction  $OP\ R_{b+k-1}, R_{b+k-1}, R_{b+m-1}$ , depending on whether the big child is the right or left child, respectively.

$k=4, b=1$

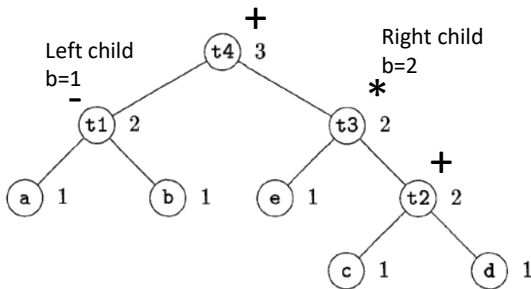
Big child ( $k$ )  $\rightarrow$  R1, R2, R3, R4 (result)

small child ( $m=3$ )  $\rightarrow$  R1, R2, R3 (result)

# Generating Code From Labeled Expression Trees

tree for expression  $(a - b) + e \times (c + d)$

t1 = a - b  
t2 = c + d  
t3 = e \* t2  
t4 = t1 + t3

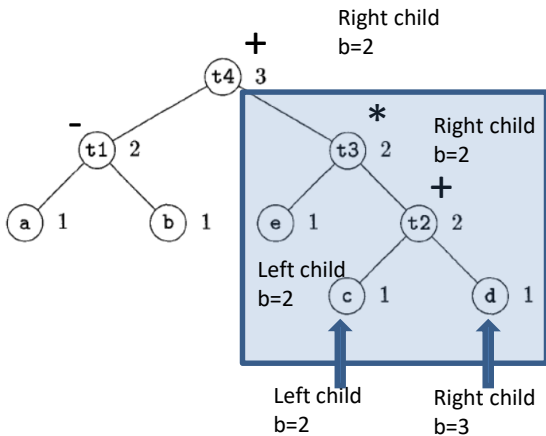


- Since the **label of the root is 3**, the **result will appear in R3**,
- Only R1, R2, and R3 will be used.
- The base for the root is b = 1.
- Since the root has children of equal labels, we generate code for the right child first, with base 2.

# Generating Code From Labeled Expression Trees

tree for expression  $(a - b) + e \times (c + d)$

```
t1 = a - b  
t2 = c + d  
t3 = e * t2  
t4 = t1 + t3
```

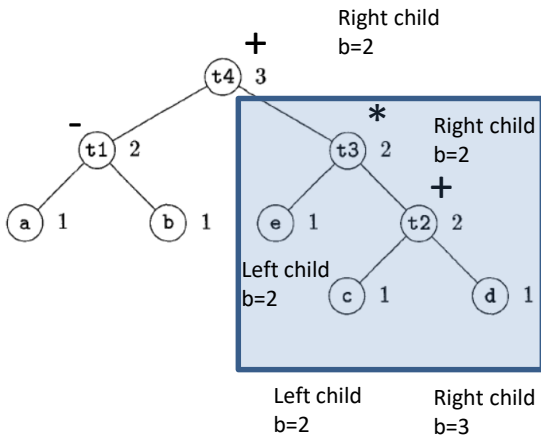


# Generating Code From Labeled Expression Trees

tree for expression  $(a - b) + e \times (c + d)$

```
t1 = a - b  
t2 = c + d  
t3 = e * t2  
t4 = t1 + t3
```

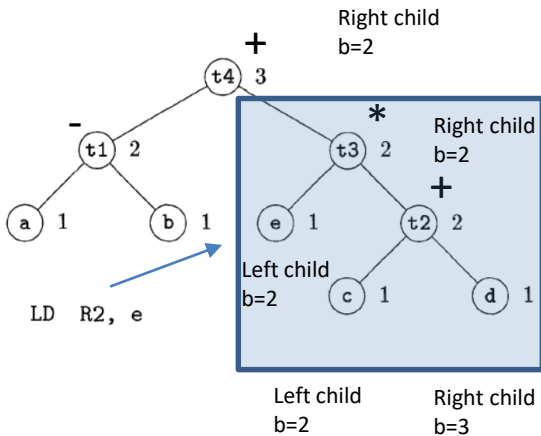
```
LD  R3, d  
LD  R2, c  
ADD R3, R2, R3
```



# Generating Code From Labeled Expression Trees

tree for expression  $(a - b) + e \times (c + d)$

```
t1 = a - b  
t2 = c + d  
t3 = e * t2  
t4 = t1 + t3
```

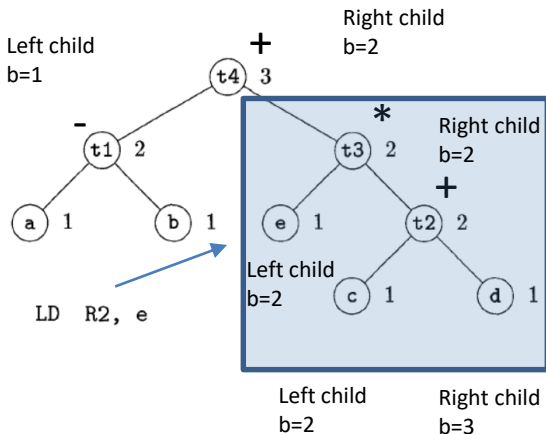


# Generating Code From Labeled Expression Trees

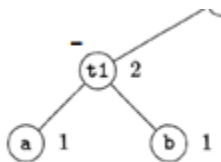
tree for expression  $(a - b) + e \times (c + d)$

```
t1 = a - b  
t2 = c + d  
t3 = e * t2  
t4 = t1 + t3
```

```
LD  R3, d  
LD  R2, c  
ADD R3, R2, R3  
LD  R2, e  
MUL R3, R2, R3
```



# Generating Code From Labeled Expression Trees



Base b=1

```
LD R2, b
LD R1, a
SUB R2, R1, R2
```

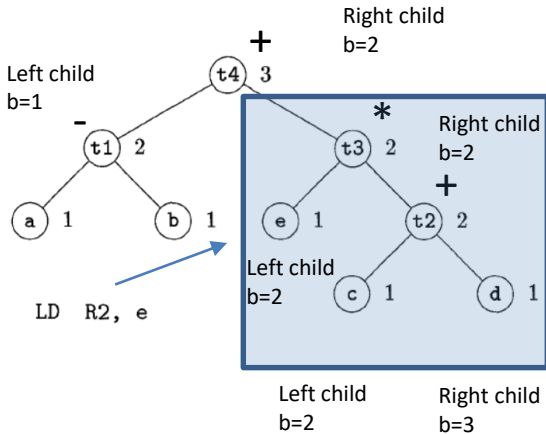


# Generating Code From Labeled Expression Trees

tree for expression  $(a - b) + e \times (c + d)$

```
t1 = a - b
t2 = c + d
t3 = e * t2
t4 = t1 + t3
```

```
LD  R3, d
LD  R2, c
ADD R3, R2, R3
LD  R2, e
MUL R3, R2, R3
LD  R2, b
LD  R1, a
SUB R2, R1, R2
ADD R3, R2, R3
```



# Insufficient Supply of Registers

**INPUT:** A labeled tree with each operand appearing once (i.e., no common subexpressions) and a number of registers  $r \geq 2$ .

**OUTPUT:** An optimal sequence of machine instructions to evaluate the root into a register, using no more than  $r$  registers, which we assume are  $R_1, R_2, \dots, R_r$ .

- Starting at the **root** of the tree, with **base b = 1**.
- For a node **N with label r or less**, the algorithm is **exactly the same**
- However, for interior nodes with a label **k > r**, we need to work on each side of the tree separately and **store the result**
- That result is **brought back into memory** just before node **N is evaluated**,
- And the final step will take place in **registers R(r-1) and R(r)**

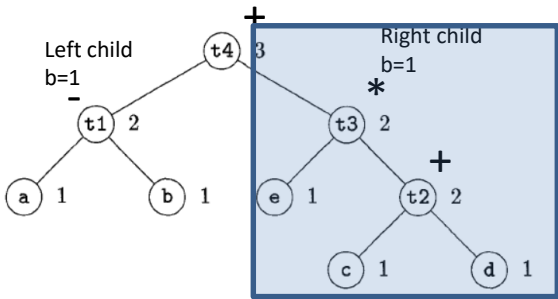
# Insufficient Supply of Registers

1. Node  $N$  has at least one child with label  $r$  or greater. Pick the larger child (or either if their labels are the same) to be the “big” child and let the other child be the “little” child.
2. Recursively generate code for the big child using base  $b = 1$ . The result of this evaluation will appear in register  $R_r$ .
3. Generate the machine instruction  $ST\ t_k, R_r$ , where  $t_k$  is a temporary variable used for temporary results used to help evaluate nodes with label  $k$ .
4. Generate code for the little child as follows. If the little child has label  $r$  or greater, pick base  $b = 1$ . If the label of the little child is  $j < r$ , then pick  $b = r - j$ . Then recursively apply this algorithm to the little child; the result appears in  $R_r$ .  **$j$  registers  $R(r-j), R(r-j+1), \dots, R(r-j+j)=R(r)$**
5. Generate the instruction  $LD\ R_{r-1}, t_k$ .
6. If the big child is the right child of  $N$ , then generate the instruction  $OP\ R_r, R_r, R_{r-1}$ . If the big child is the left child, generate  $OP\ R_r, R_{r-1}, R_r$ .

# Insufficient Supply of Registers

assume that  $r = 2$ ; that is, only registers R1 and R2 are available

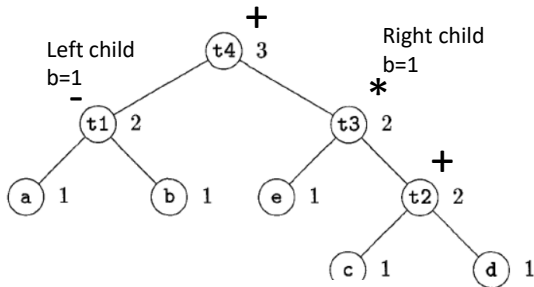
```
LD  R2, d
LD  R1, c
ADD R2, R1, R2
LD  R1, e
MUL R2, R1, R2
```



# Insufficient Supply of Registers

assume that  $r = 2$ ; that is, only registers R1 and R2 are available

```
LD  R2, d
LD  R1, c
ADD R2, R1, R2
LD  R1, e
MUL R2, R1, R2
```



Generate the machine instruction `ST  $t_k$ ,  $R_r$` , where  $t_k$  is a temporary variable used for temporary results used to help evaluate nodes with label  $k$ .

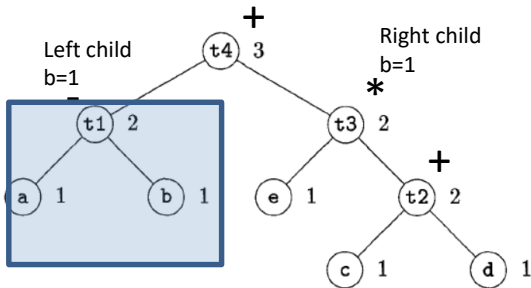
```
ST  t3, R2
```

**Spill**

# Insufficient Supply of Registers

assume that  $r = 2$ ; that is, only registers R1 and R2 are available

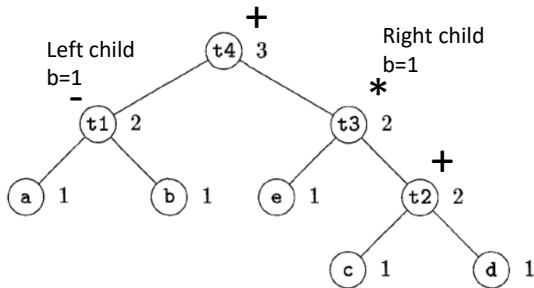
```
LD  R2, b
LD  R1, a
SUB R2, R1, R2
```



# Insufficient Supply of Registers

assume that  $r = 2$ ; that is, only registers R1 and R2 are available

```
LD R2, d
LD R1, c
ADD R2, R1, R2
LD R1, e
MUL R2, R1, R2
ST t3, R2
LD R2, b
LD R1, a
SUB R2, R1, R2
LD R1, t3
ADD R2, R2, R1
```



# Optimization – Beyond basic blocks

```
void quicksort(int m, int n)
/* recursively sorts a[m] through a[n] */
{
    int i, j;
    int v, x;
    if (n <= m) return;
    /* fragment begins here */
    i = m-1; j = n; v = a[n];
    while (1) {
        do i = i+1; while (a[i] < v);
        do j = j-1; while (a[j] > v);
        if (i >= j) break;
        x = a[i]; a[i] = a[j]; a[j] = x; /* swap a[i], a[j] */
    }
    x = a[i]; a[i] = a[n]; a[n] = x; /* swap a[i], a[n] */
    /* fragment ends here */
    quicksort(m,j); quicksort(i+1,n);
}
```

(1)	i = m-1	(16)	t7 = 4*i
(2)	j = n	(17)	t8 = 4*j
(3)	t1 = 4*m	(18)	t9 = a[t8]
(4)	v = a[t1]	(19)	a[t7] = t9
(5)	i = i+1	(20)	t10 = 4*j
(6)	t2 = 4*i	(21)	a[t10] = x
(7)	t3 = a[t2]	(22)	goto (5)
(8)	if t3<v goto (5)	(23)	t11 = 4*i
(9)	j = j-1	(24)	x = a[t11]
(10)	t4 = 4*j	(25)	t12 = 4*i
(11)	t5 = a[t4]	(26)	t13 = 4*n
(12)	if t5>v goto (9)	(27)	t14 = a[t13]
(13)	if i>=j goto (23)	(28)	a[t12] = t14
(14)	t6 = 4*i	(29)	t15 = 4*n
(15)	x = a[t6]	(30)	a[t15] = x



```

(1)  i = m-1
(2)  j = n
(3)  t1 = 4*n
(4)  v = a[t1]
(5)  i = i+1
(6)  t2 = 4*i
(7)  t3 = a[t2]
(8)  if t3<v goto (5)
(9)  j = j-1
(10) t4 = 4*j
(11) t5 = a[t4]
(12) if t5>v goto (9)
(13) if i>=j goto (23)
(14) t6 = 4*i
(15) x = a[t6]

```



$x = a[i]$  is translated as

```

t6 = 4*i
x = a[t6]

```

```

(16) t7 = 4*i
(17) t8 = 4*j
(18) t9 = a[t8]
(19) a[t7] = t9
(20) t10 = 4*j
(21) a[t10] = x
(22) goto (5)
(23) t11 = 4*i
(24) x = a[t11]
(25) t12 = 4*i
(26) t13 = 4*n
(27) t14 = a[t13]
(28) a[t12] = t14
(29) t15 = 4*n
(30) a[t15] = x

```



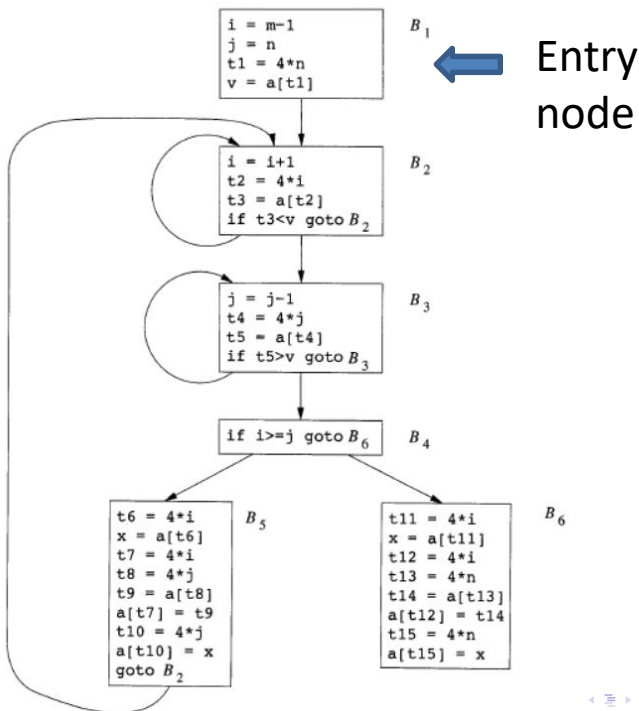
$a[j] = x$  becomes

```

t10 = 4*j
a[t10] = x


```

- Notice that every array access in the original program translates into a pair of steps, consisting of a multiplication and an array subscription operation.
- Short program fragment translates into a rather long sequence of three-address operations.



# Global Common Subexpressions

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



```
t6 = 4*i
x = a[t6]
t7 = 4*i
t8 = 4*j
t9 = a[t8]
a[t7] = t9
t10 = 4*j
a[t10] = x
goto B2
```

(a) Before,

$B_5$

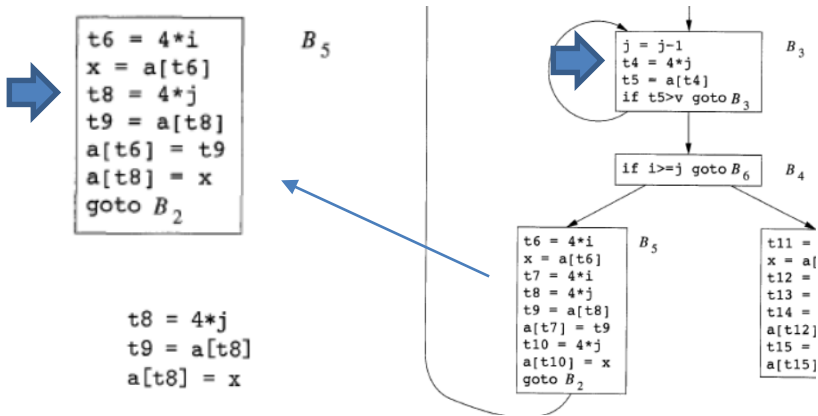
```
t6 = 4*i
x = a[t6]
t8 = 4*j
t9 = a[t8]
a[t6] = t9
a[t8] = x
goto B2
```

$B_5$

(b) After,

Compute the common subexpressions  $4 * i$  and  $4 * j$ , respectively

# Global Common Subexpressions

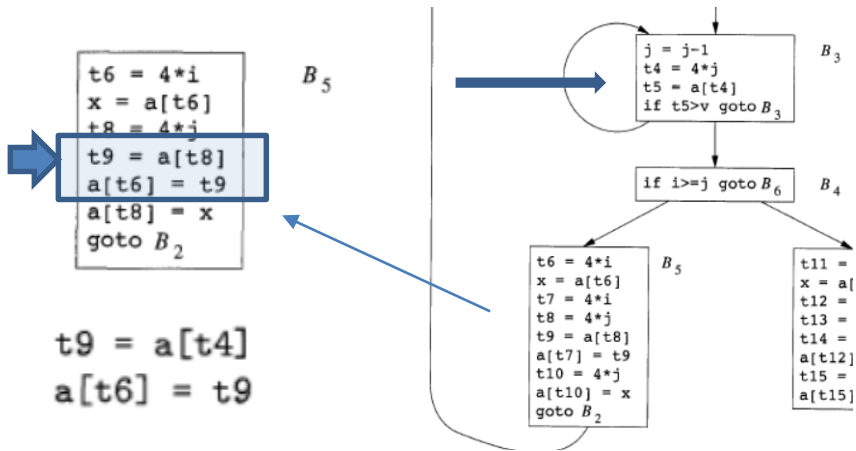


in  $B_5$  can be replaced by

```
t9 = a[t4]
a[t4] = x
```

- Control passes from the evaluation of  $4 * j$  in  $B_3$  to  $B_5$ ,
- No change to  $j$  and no change to  $t4$ , so  $t4$  can be used if  $4 * j$  is needed.

# Global Common Subexpressions

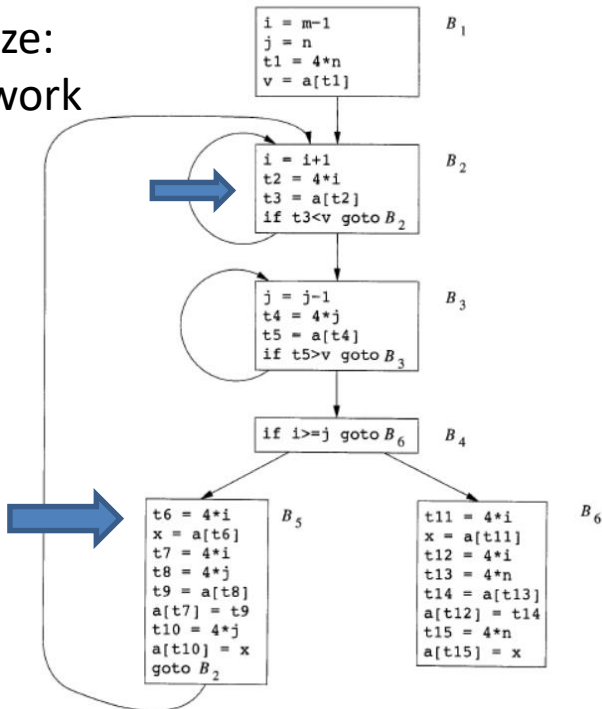


in  $B_5$  therefore can be replaced by

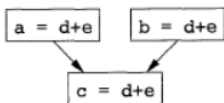
`a[t6] = t5`

`a[t4]`, a value computed into a temporary `t5`, retains its value as control leaves  $B_3$  and then enters  $B_5$

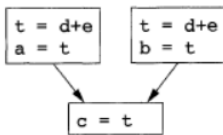
# Optimize: Homework



# Copy Propagation



(a)





(b)


We must use a new variable  $t$  to hold the value of  $d + e$   
After the copy statement  $u = v$ , **use  $v$**  for  $u$

```
x = t3 ←  
a[t2] = t5  
a[t4] = t3 ←  
goto B2
```

# Dead-Code Elimination

<code>x = t3</code>			
<code>a[t2] = t5</code>			<code>a[t2] = t5</code>
<code>a[t4] = t3</code>			<code>a[t4] = t3</code>
<code>goto B<sub>2</sub></code>			<code>goto B<sub>2</sub></code>

`if (debug) print ...`

`debug = FALSE`  Copy propagation

- Drop this code segment
- Constant folding



# Code Motion

- Loops are a very important place for optimizations, especially the inner loops where programs tend 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.
- Code Motion takes an expression
  - That yields the **same result** independent of the number of times a loop is executed (**a loop-invariant computation**) and
  - Evaluates the expression before the loop

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

Code motion will result in the equivalent code

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