

# UCS2708 COMPILER DESIGN

## Unit 5

### Code Optimization

Principal sources of optimization - DAG - Optimization of basic blocks -

Global data flow analysis - Introduction to low level Virtual

Machine (LLVM) - Design of LLVM - Core libraries - Developing plugin in LLVM

#### \* Introduction to Code Optimization - Basic Blocks

→ A basic block is a sequence of consecutive TAC statements in which:

- (i) Control flow can only enter the block at the first statement
- (ii) Control flow can only exit the block from the last statement

For eg - The TAC below forms a basic block

More on code optimization

$t_1 := a * a$

$t_2 := a * b$

$t_3 := a * t_2$

$t_4 := t_1 + t_3$

$t_5 := b * b$

$t_6 := t_4 + t_5$

- can be machine dependent or independent
- can be applied at
  - source level
  - intermediate code level
  - target program level
- Principle Sources of Optimization
  1. Local Optimization
  2. Global Optimization
  3. Loop Optimization - local

## Partition Algorithm for Basic Blocks

Input: A sequence of TAC statements

Output: A list of basic blocks with each statement in exactly one block.

1. Determine the set of leaders

- ① The first statement is the leader
- ② Any statement that is the target of a goto is a leader
- ③ Any statement that immediately follows a goto is a leader

2. For each leader, its basic block consists of the leader and all statements up to but not including the next leader or the end of the program.

Example → Partition into basic blocks

- |                                  |          |                      |
|----------------------------------|----------|----------------------|
| (1) $\text{prod} = 0$            | → Leader | $B_1 = (1) \cup (2)$ |
| (2) $i = 1$                      |          |                      |
| (3) $t_1 = 4 * i$                | → Leader | $B_2 = (3) \cup (4)$ |
| (4) $t_2 = a[t_1]$               |          |                      |
| (5) $t_3 = 4 * i$                |          |                      |
| (6) $t_4 = b[t_3]$               |          |                      |
| (7) $t_5 = t_2 + t_4$            |          |                      |
| (8) $t_6 = \text{prod} + t_5$    |          |                      |
| (9) $\text{prod} = t_6$          |          |                      |
| (10) $t_7 = i + 1$               |          |                      |
| (11) $i = t_7$                   |          |                      |
| (12) IF ( $i \leq 20$ ) goto (3) |          |                      |

(3)

## Example Partition the TAC into basic blocks

1)  $i = 1 \rightarrow$  Leader

2)  $j = 1 \rightarrow$  Leader

3)  $t_1 = 10 * i \rightarrow$  Leader

4)  $t_2 = t_1 + j$

5)  $t_3 = 8 * t_2$

6)  $t_4 = t_3 - 88$

7)  $a[t_4] = 0.0$

8)  $j = j + 1$

9) If  $j <= 10$  goto (3)

10)  $i = i + 1 \rightarrow$  Leader

11) If  $i <= 10$  goto (2)

12)  $i = 1 \rightarrow$  Leader

13)  $t_5 = i - 1 \rightarrow$  Leader

14)  $t_6 = 88 * t_5$

15)  $a[t_6] = 1.0$

16)  $i = i + 1$

17) If  $i <= 10$  goto (13)

Ans

$B_1 = (1)$

$B_2 = (2)$

$B_3 = (3) - (9)$

$B_4 = (10), (11)$

$B_5 = (12)$

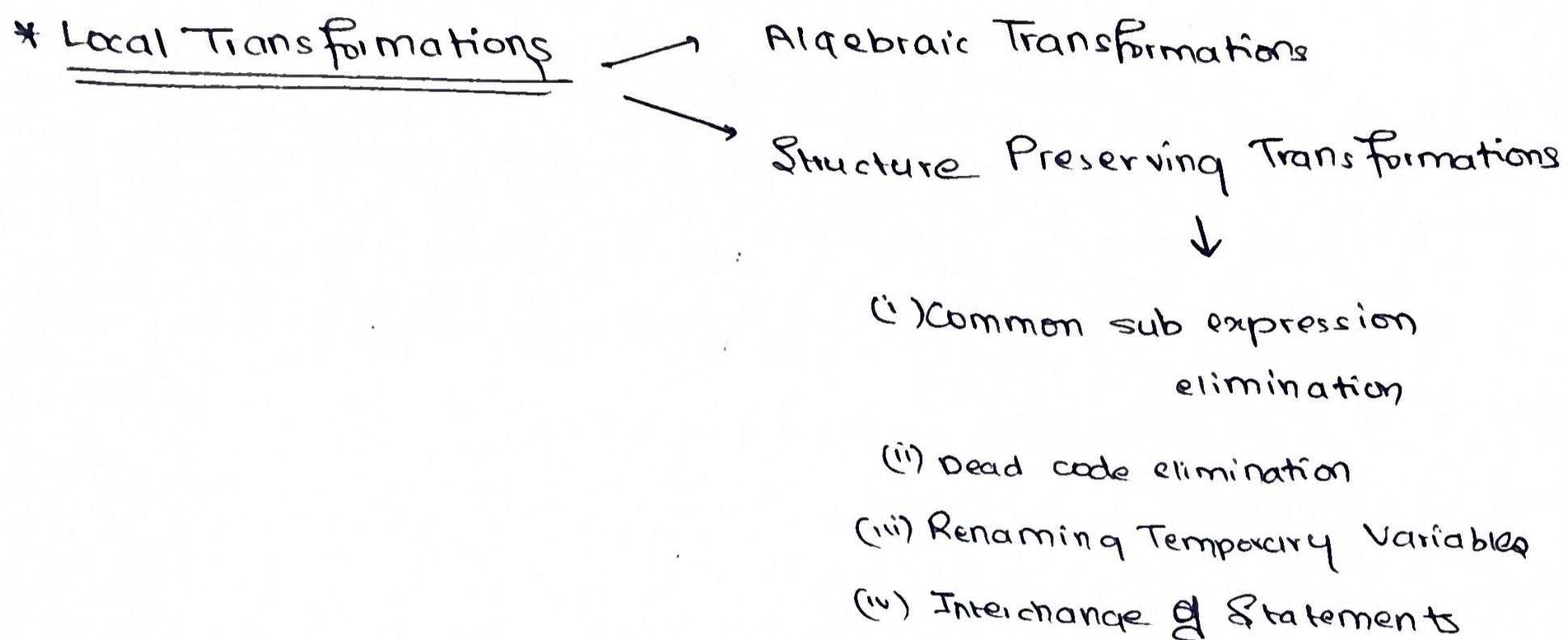
$B_6 = (13) - (16)$

## \*Transformations on Basic Blocks

→ A code-improving transformation is a code optimisation to improve speed or reduce code size.

→ Global transformations are performed across basic blocks

- Local Transformations are only performed on single basic blocks
- Transformations must be safe and preserve the meaning of the code. - A local transformation is safe if the transformed basic block is guaranteed to be equivalent to its original form.



## A. Structure Preserving Transformations

### ① Common Sub-Expression Elimination

eq1	$a = b + c$	$a = b + c$
	$b = a - d$	$\Rightarrow \quad b = a - d$
	$c = b + c$	$c = b + c$
	$d = a - d$	$d = b$

eq2	$t_1 = b * c$	$t_1 = b * c$
	$t_2 = a - t_1$	$\Rightarrow \quad t_2 = a - t_1$
	$t_3 = b * c$	$t_4 = t_2 + \cancel{t_1}$
	$t_4 = t_2 + t_3$	

## ② Dead Code Elimination

eq1 assume  $a$  is unused

$$\begin{array}{l} b = a + 1 \\ a = b + c \end{array} \Rightarrow b' = a + 1$$

Pq2 if true goto L2 } remove unreachable code  
 $b = x + y$

## ③ Renaming Temporary Variables

→ Temporary variables that are dead at the end of a block can be safely renamed.

eq1  $t_1 = b + c$

$$t_2 = a - t_1$$

$$t_1 = t_1 * d$$

$$d = t_2 + t_1$$

$$t_1 = b + c$$

$$t_2 = a - t_1$$

$$t_3 = t_1 * d$$

$$d = t_2 + t_3$$

change  $t_1$  to  $t_3$   
 so that each var  
 uniquely represents  
 one computation

## ④ Interchange of Statements

→ Independent statements can be reordered

$$t_1 = b + c$$

$$t_2 = a - t_1$$

$$t_3 = t_1 * d \Rightarrow$$

$$d = t_2 + t_3$$

$$t_1 = b + c$$

$$t_3 = t_1 * d$$

$$t_2 = a - t_1$$

$$d = t_2 + t_3$$

## B. Algebraic Transformations

→ includes constant folding

→ change arithmetic operations to transform blocks into algebraic equivalent forms

$$t_1 = a - a$$

$$t_1 = 0$$

$$t_2 = b + t_1$$

$$\Rightarrow t_2 = b$$

$$t_3 = 2 * t_2$$

$$t_3 = t_2 \ll 1$$

(Left shift by 1 means mul by 2)

→ Also eliminate:

$$x = x + 0 \Rightarrow x$$

$$x = x * 1 \Rightarrow x$$

→ Modify  $x = 4 * x^2 \Rightarrow x = 4 * 4$

## \* DAG and Loop Optimization

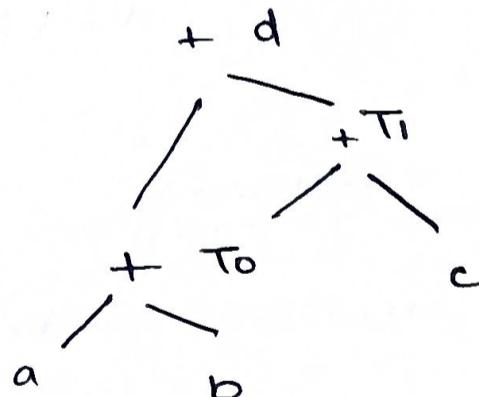
### \* DAG for Basic Blocks

Eq1  $T_0 = a + b$

$$T_1 = T_0 + c$$

$$d = T_0 + T_1$$

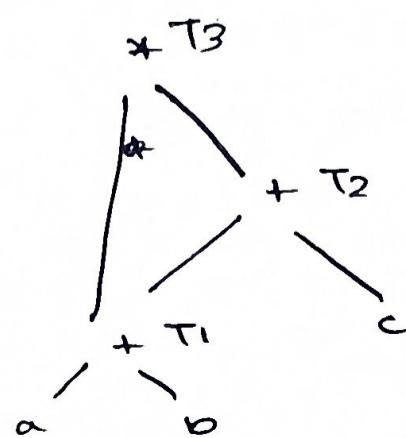
→



Eq2  $T_1 = a * b$

$$T_2 = T_1 + c$$

$$T_3 = T_1 * T_2$$



eq3 Simplify the given TAC using a DAG (7)

$$(1) t_1 = 4 * i$$

$$t_2 = a [t_1]$$

$$t_3 = 4 * i$$

$$t_4 = b [t_3]$$

$$t_5 = t_2 * t_4$$

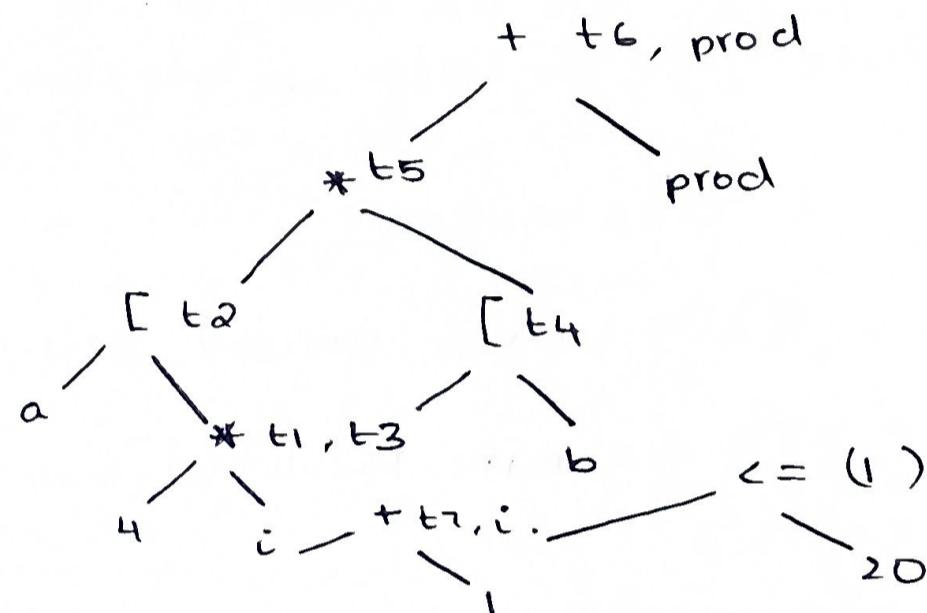
$$t_6 = \text{prod} + t_5$$

$$\text{prod} = t_6$$

$$t_7 = i + 1$$

$$i = t_7$$

If  $i \leq 20$  goto (1)



replacements to make

replace

$t_6$  with  $\text{prod}$

$t_3$  with  $t_1$

$t_7$  with  $i$

$$t_1 = 4 * i$$

$$t_2 = a [t_1]$$

eliminate

$t_6 = \text{prod} \rightarrow \{ \text{no need to}$

$i = t_7 \rightarrow \{ \text{double assignment, use } i \neq \text{prod}$

$t_3 = 4 * i \rightarrow \text{redundant}$

$$t_4 = b [t_1]$$

$$t_5 = t_2 * t_4$$

$$\text{prod} = \text{prod} + 5$$

$$i = i + 1$$

If  $i \leq 20$  goto (1)

## \* Loop Optimization

- 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 the loop.
- 3 techniques for loop optimization are:

- ① Code motion
- ② Induction variable elimination
- ③ Strength Reduction

### A. Code Motion

- This transformation takes an expression that yields the same result independent of the no. of times a loop is executed, and places the expression before the loop.

For eg.

```
while (i <= n - 2)  
{  
}
```

Evaluation of  $n - 2$  is loop independent. Place it outside

$$t = n - 2$$

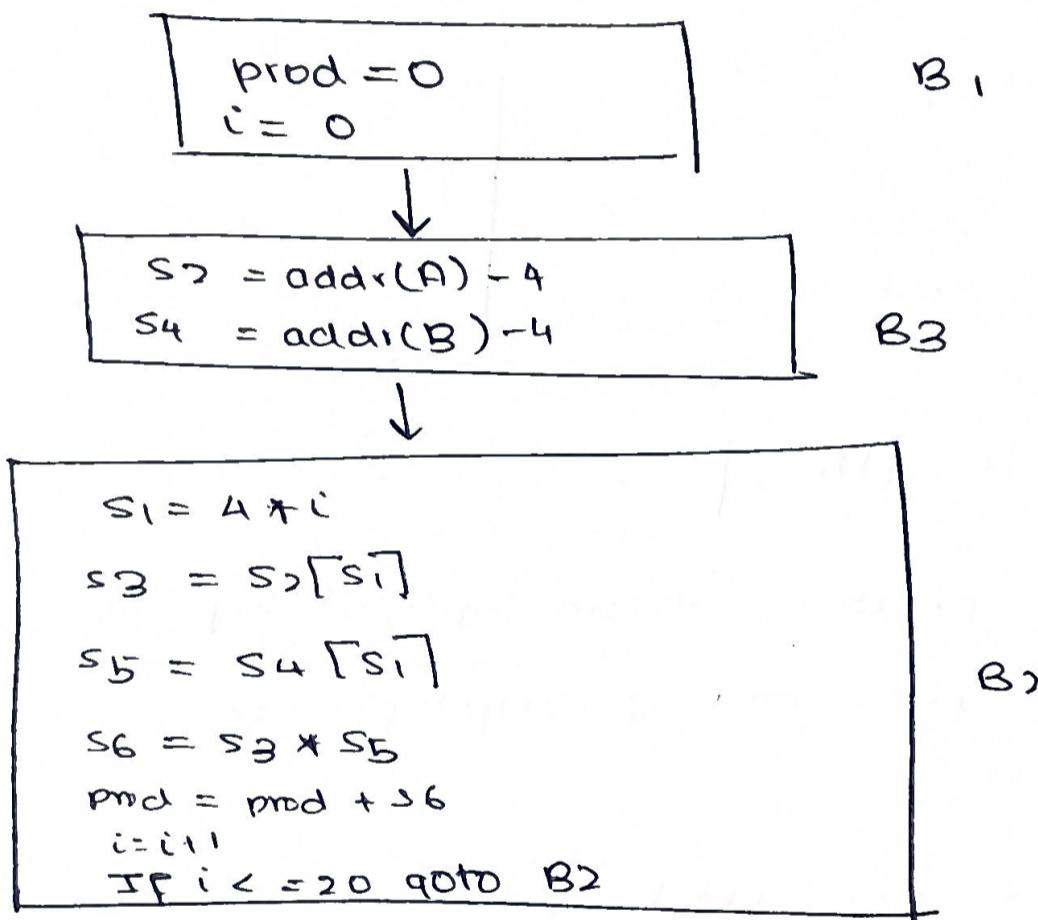
```
while (i <= t) {
```

}

## B. | Induction Variable Elimination

- An induction variable is one that changes in a regular pattern (eg. increments / decrements in each iteration)
- Observe patterns where the induction var is used and eliminate it entirely, by replacing its use with an alternate computation

Consider the following example:



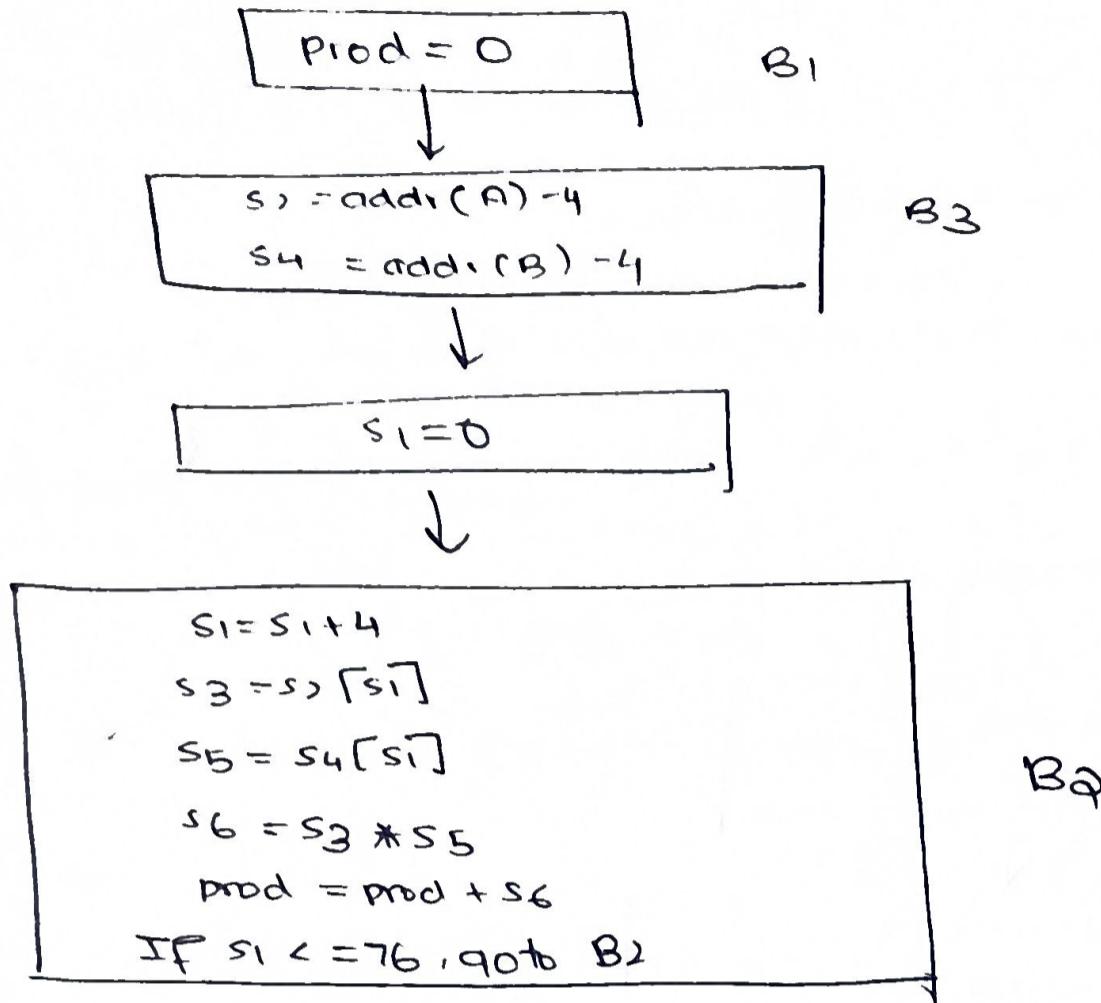
In this code,  $i$  takes the value 1, 2, ..., 26

⇒  $S1$  would take the values 4, 8, ..., 80

⇒ Replace  $A * i$  and  $i = i + 1$  with  $S1 = S1 + 4$

In the if block,  $i <= 20$  can be replaced with  $S1 <= 76$   
after calculating value

∴ The new code block would be:



### c. Reduction in Strength

→ Higher strength operators can be replaced by lower strength operators - e.g. replacement of multiplication by repeated addition

For e.g.

For (i=1, i<=50, i++)

{

C = i \* 7

}

replace as

For  $t = 7$   
(i=1, i<=50, i++)

{

C = t

$t = t + 7$

}

## \* Peephole Optimization

- Peephole optimization is a technique for improving code by examining small, localized sections of target code and replacing inefficient or redundant instructions with more efficient ones.
- A small moving 'peephole' is used to analyse and optimise sequences of instructions
- Limited to short code sequences rather than the entire program

Techniques include

- Eliminate redundant loads & stores
- Eliminate unreachable code
- Flow of control optimization
- Simplify algebraic expressions
- Strength Reduction

### A. | Eliminate redundant loads & stores

consider:

(1) `MOV R0, a`  
 (2) `MOV a, R0`

(2) can be removed if it does not have a label

### B. | Eliminate unreachable code

Consider the following

```
#define debug 0
if (debug) {
    print()
}
goto L1
print()
```

The print statements are inaccessible

### c. Flow of Control Optimization

Consider:

Goto L1

:

L1 : Goto L2

If under all cases, L2 follows L1, remove one goto

Make it goto L2

### d. Simplify algebraic expressions

$$x := x + 0 \Rightarrow x$$

$$x := x \times 1 \Rightarrow x$$

### e. Strength Reduction

$$(i) x * 2 \rightarrow x \ll 1$$

$$(ii) x / 2 \rightarrow x \gg 1$$

$$(iii) ADD #1, R \rightarrow INCR$$

$$(iv) x ^ 2 \rightarrow x * x$$

Example - Identify the code optimization techniques and optimize

the code given below.

$$1. A = 2 * (22.0 / 7.6) * r$$

$$2. c = a * b$$

→ constant folding

$$x = a$$

$$d = a * b + 4$$

$$A = 6.28 * r$$

→ Redundant code elimination

$$c = a * b$$

~~$$x = a$$~~

$$d = a * b + 4$$

3.  $a = 200$   
 while ( $a > 0$ ) {  
      $b = x + 4$ ;  
     if ( $a \cdot b = 0$ ) {  
         printf ("%d", a);  
     }  
 }

→ Dead code elimination Loop optimization - code motion  
~~remove~~ → move  $b = x + 4$  outside loop

4.  $c = a * b$   
 $x = a$   
 $d = a * b + 4$

→ common subexpression elimination

$$d = c + 4$$

5.  $i = 1$   
 while ( $i < 100$ ) {  
      $y = i * 4$   
 }

→ Induction variable elimination

$y = 4$   
 while ( $y < 400$ ) {  
      $y = y + 4$   
 }

## \* Global Flow Data Analysis

- A technique to optimize code by analyse the flow of data across the entire program.
  - Specifically it:
    - (i) Tracks where variables are defined and used throughout the program
    - (ii) Identifies relationships between definitions and the corresponding uses (called Use-Definition (Ud) Chains)
    - (iii) Helps eliminate redundant or unnecessary computations
  - For Ud-chaining Reaching definitions should be found
  - Then data-flow equations have to be determined iteratively
- Reaching Definition → A definition of a variable is said to 'reach' a point in the program if there exists a path from the definition to that point, such that the value assigned to the variable is still valid, and not redefined along the way.
- Algorithm
1. Assign a number to each equation
  2. Compute each of the following for each block 'B':
    - A.  $\text{GEN}[B]$  - set of generated definitions within basic block B that reach the end of the block.

B. KILL(B) - set of definitions outside of B that have a definition within B

C. Use bit vectors for representation

3. Compute  $\boxed{IN[B]}$  for all blocks B - all definitions reaching the point just before the first statement of block B

4. Compute  $\boxed{OUT[B]}$  for all blocks B - the set of definitions reaching the point just after the last statement of B

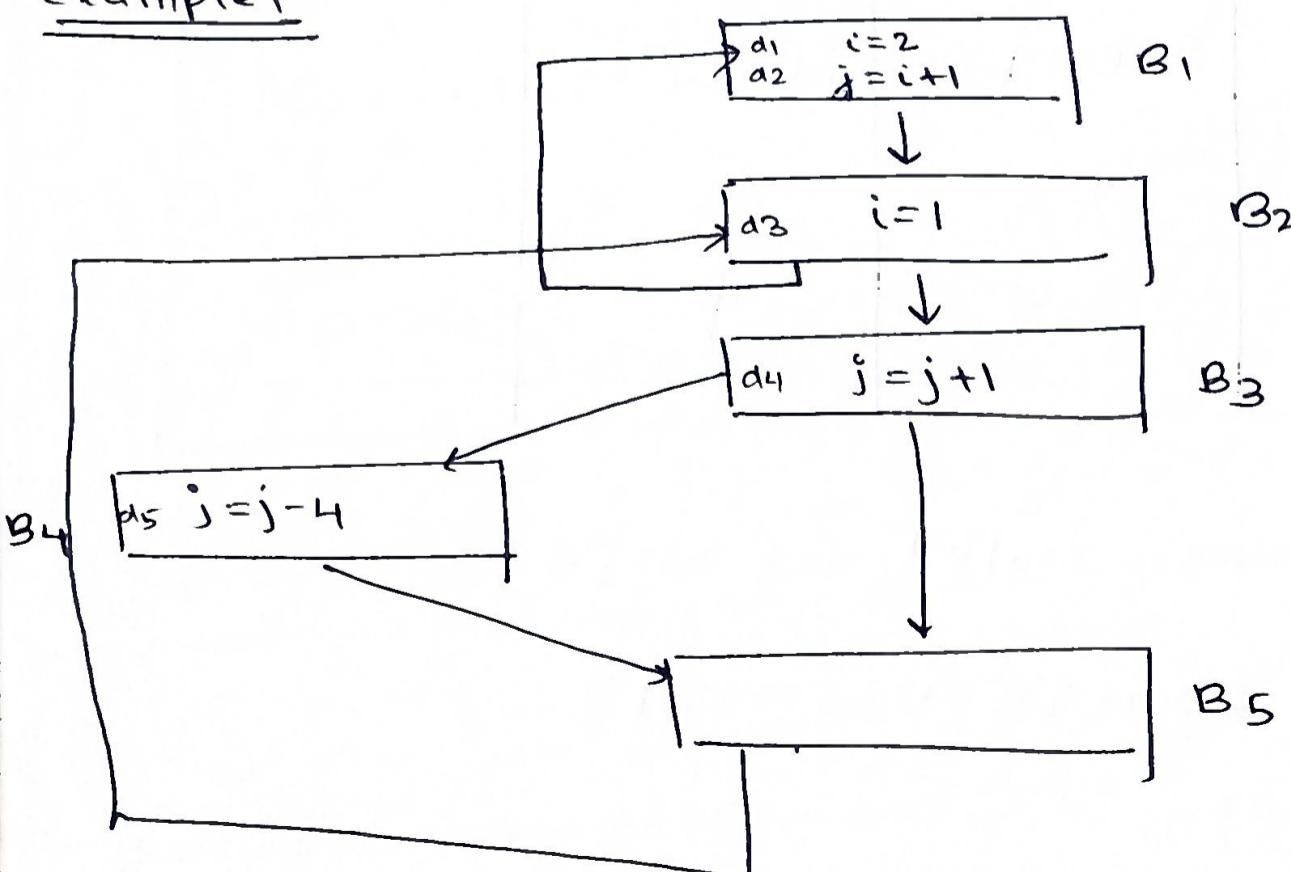
5. Repeat 3,4,B iteratively until convergence.

The data flow equations to compute IN and OUT are:

$$OUT[B] = IN[B] - (KILL[B] \cup GEN[B])$$

$$IN[B] = \bigcup OUT[P], P \text{ is a predecessor.}$$

Example



Ans Step 1 : Assign numbers to the definitions

$$d_1 \rightarrow i = 2$$

$$d_2 \rightarrow j = i + 1$$

$$d_3 \rightarrow i = 1$$

$$d_4 \rightarrow j = j + 1$$

$$d_5 \rightarrow j = j - 4$$

write predecessor or

$$P(B_1) = B_2$$

$$P(B_2) = B_1, B_5$$

$$P(B_3) = B_2$$

$$P(B_4) = B_3$$

$$P(B_5) = B_3, B_4$$

Step 2 Compute GEN and KILL for each block

	GEN[B]	Bit Vector	KILL[B]	Bit Vector
B <sub>1</sub>	{d <sub>1</sub> , d <sub>2</sub> }	11 000	{d <sub>3</sub> , d <sub>4</sub> , d <sub>5</sub> }	00 111
B <sub>2</sub>	{d <sub>3</sub> }	00 100	{d <sub>1</sub> }	100 00
B <sub>3</sub>	{d <sub>4</sub> }	000 10	{d <sub>2</sub> , d <sub>5</sub> }	01001
B <sub>4</sub>	{d <sub>5</sub> }	00001	{d <sub>2</sub> , d <sub>4</sub> }	01010
B <sub>5</sub>	φ	00000	φ	00000

Step 3 : Iteratively compute IN[B] and OUT[B]

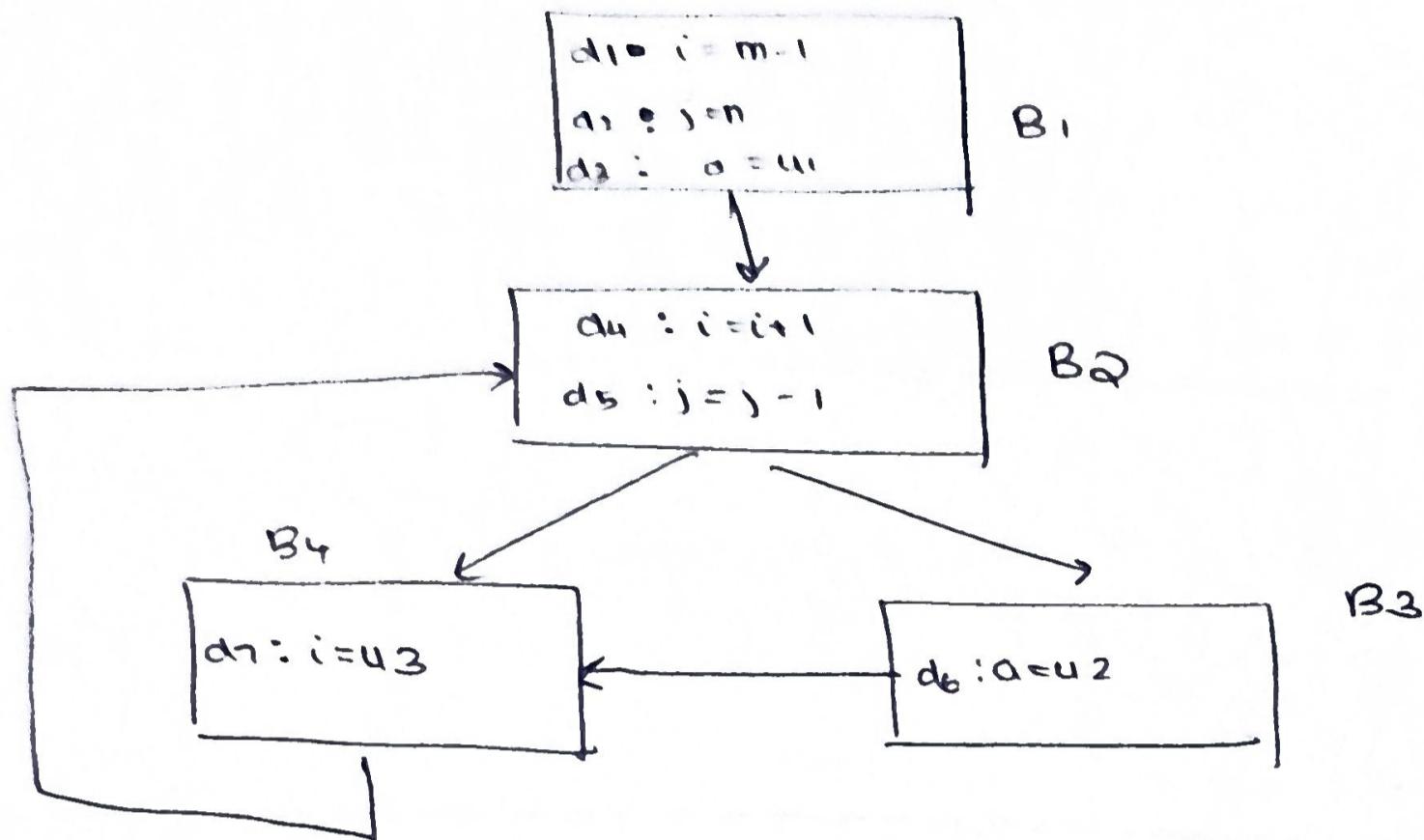
$$OUT[B] = IN[B] - \cancel{KILL[B]} \cup \cancel{GEN[B]}$$

$$IN[B] = \cup OUT[P]$$

remove definitions from in that are in kill

Block	IN	OUT initially	IN	OUT	IN	OUT
B <sub>1</sub>	∅	some ac GEN $\{d_1, d_2\}$	$U_{OUT}[B_2]$ BUT $= \{d_3\}$	$IN[B_1] -$ $(IN[B] \cup GEN[B])$ $\{d_3\} \setminus \{d_3, d_4, d_5\}$ $+ \{d_1, d_2\} \setminus$ $\boxed{\{d_1, d_2\}}$	$\{d_2, d_3\}$	$\{d_2, d_3\} -$ $\{d_3, d_4, d_5\}$ $+ \{d_1, d_2\}$ $\boxed{\{d_1, d_2\}}$
B <sub>2</sub>	∅	$\{d_3\}$	$U_{OUT}[B_1]$ $= \{d_1, d_2\}$	$\{d_1, d_2\} - \{d_1\}$ $+ \{d_3\}$ $= \boxed{\{d_2, d_3\}}$	$\{d_1, d_2, d_3, d_4,$ $d_5\}$	$\{d_1, d_2, d_3, d_4, d_5\}$ $- \{d_1\} + \{d_3\}$ $\boxed{\{d_2, d_3, d_4, d_5\}}$
B <sub>3</sub>	∅	$\{d_4\}$	$U_{OUT}[B_2]$ $= \{d_2, d_3\}$	$\{d_2, d_3\} -$ $\{d_2, d_5\} + \{d_4\}$ $\{d_3\} + \{d_4\} = \boxed{\{d_3, d_4\}}$	$\{d_2, d_3\} *$	$\{d_2, d_3\} -$ $\{d_2, d_5\} + \{d_4\}$ $\boxed{\{d_3, d_4\}}$
B <sub>4</sub>	∅	$\{d_5\}$	$U_{OUT}[B_3]$ $= \{d_3, d_4\}$	$\{d_3, d_4\} -$ $\{d_2, d_4\} + \{d_5\}$ $= \boxed{\{d_3, d_5\}}$	$\{d_3, d_4\}$	$\{d_3, d_4\} -$ $\{d_2, d_4\} + \{d_5\}$ $\boxed{\{d_3, d_5\}}$
B <sub>5</sub>	∅	$\{d_4\}$	$\{d_3, d_5\}$ $d_4$	$\{d_3, d_5\} -$ $\{d\} + \{d\}$ $\boxed{\{d_3, d_4, d_5\}}$	$\{d_3, d_4, d_5\}$	→ If you calc again, you'll see it doesn't change ③

Example : Analyze the data flow for the following example.



Ans Step 1 Compute GEN & KILL

	GEN[B]	Bit Vector	Kill[B]	Bit Vecto
B1	{d1, d2, d3}	1110000	{d4, d5, d6, d7}	00001111
B2	{d4, d5}	0001100	{d1, d2, d7}	1100001
B3	{d6}	0000010	{d3}	0010000
B4	{d7}	0000001	{d1, d4}	1001000

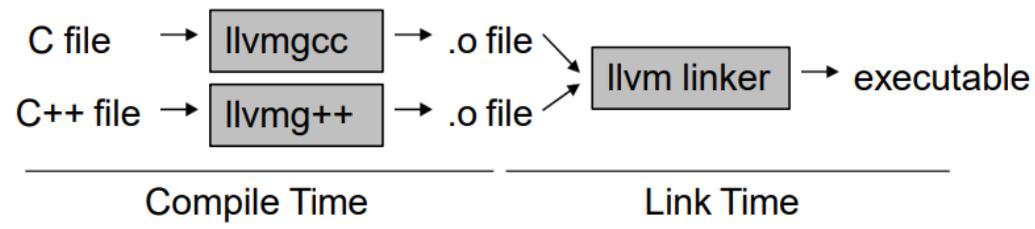
## 1. What is LLVM?

- LLVM stands for **Low-Level Virtual Machine**. It is a **compiler infrastructure** that provides reusable libraries to build compilers.
- Key Features:
  - Implemented in **C++** for flexibility and speed.
  - Supports **multiple front-ends** (e.g., C, C++, ADA) and **back-ends** (e.g., ARM, x86, MIPS).
  - Produces an **intermediate representation (IR)** that is language and target independent.
  - Open-source, first released in 2003, and widely adopted for various compiler-related tasks.

### Illustration:

Think of LLVM as a toolkit to:

- Transform programming languages into an optimized intermediate form (LLVM IR).
- Convert this optimized form into executable machine code for specific architectures.



## 2. LLVM Compiler System

The system is divided into:

- **LLVM Compiler Infrastructure:**
  - Reusable components for **building compilers**.
  - Reduces time/cost by avoiding building everything from scratch.
  - Builds static compilers, Just-In-Time (JIT) compilers, and optimizers.
- **LLVM Compiler Framework:**
  - Creates end-to-end compilers (e.g., compilers for C and C++).
  - Supports both native (machine code) and intermediate code generation.

---

## 3. Three Primary Components

## 1. LLVM Virtual Instruction Set:

- Centralized **Intermediate Representation (IR)** that is independent of programming languages and hardware.
- Represents code in a **persistent, target-independent format**.

## 2. Well-Integrated Libraries:

- Include libraries for **code generation**, optimization, debugging, profiling, garbage collection, etc.

## 3. Tools Built from Libraries:

- Assemblers, debuggers, linkers, code generators, etc., simplify compiler development.

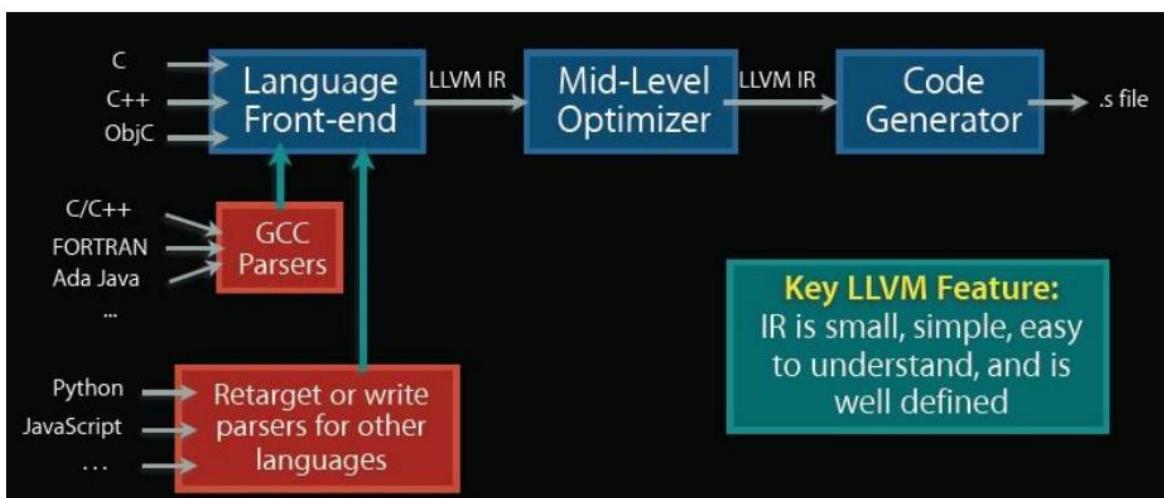
---

## 4. LLVM C/C++ Compiler

- LLVM provides **clang** (C/C++ front-end compiler) that:
  - Uses LLVM IR as an intermediate step.
  - Translates source code into **.o files** containing IR or machine code.
  - Supports GCC makefiles and uses GCC parsers (e.g., gcc 3.4).

### Compile Time Workflow:

1. C/C++ source code → LLVM front-end (e.g., llvmgcc).
2. IR passes through multiple **optimizers**.
3. Resulting optimized IR is linked by llvm-linker into:
  - Executable machine code (native) or bytecode (for JIT execution).



---

## 5. Compiler Organization with LLVM IR

- The **compiler pipeline**:
    - Language-specific **front-end** converts source code to LLVM IR.
    - **Mid-level optimizers** improve LLVM IR using optimizations like inlining, dead code elimination, etc.
    - **Code generators** create target-specific machine code.
- 

## 6. Compile-Time Events

- Key Steps:
  - **Front-End Parsing**:
    - Converts source code into Abstract Syntax Tree (AST) or LLVM IR.
    - Example: C files use cc1, while C++ files use cc1plus.
  - **Optimization**:
    - Runs **analysis and transformation passes** to optimize code.
    - Example passes: Dead Argument Elimination, Loop Invariant Code Motion (LICM), and Memory Promotion.
  - **IR Verification**:
    - Verifies correctness and validity of LLVM IR before code generation.

---

## 7. Link-Time Events

- During linking:
  - The **LLVM linker** combines .o files and applies final optimizations.
  - **Native code generators** (e.g., llc) convert LLVM IR to machine code.
  - Optionally, the code can remain in LLVM IR or bytecode format for further JIT compilation.