



中国科学技术大学
University of Science and Technology of China

代码优化

《编译原理和技术》

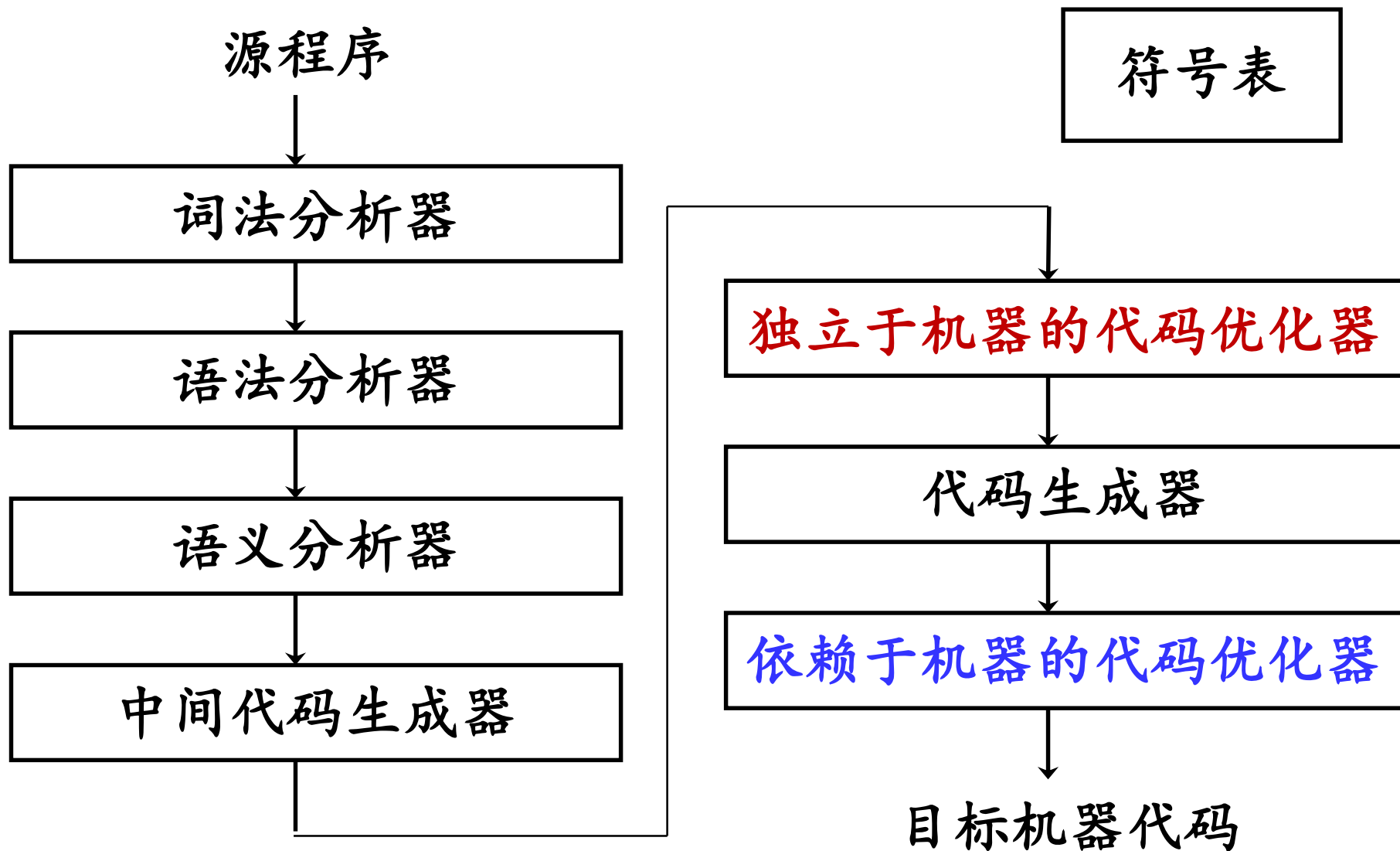
张昱

0551-63603804, yuzhang@ustc.edu.cn

中国科学技术大学
计算机科学与技术学院



本章内容





本章内容

□ 代码优化

- 通过**程序变换**（局部变换和全局变换）来改进程序
- 局部：基本块内；全局：跨基本块

□ 代码改进变换的标准

- 代码变换必须保程序的含义
- 采取安全稳妥的策略
- 通过变换减少程序的运行时间能平均达到某可度量的值
- 变换所作的努力是值得的

□ 本章介绍独立于机器的优化

- 重点：数据流分析及其一般框架



1. 优化的源头和种类

- 基本块内优化、全局优化
- 公共子表达式删除、复写传播、死代码删除
- 循环优化



优化的主要源头和主要种类

□ 主要源头:程序中存在程序员无法避免的冗余计算

如 $A[i][j]$ 、 $x.f1$ 等数据访问操作

- 编译后被展开成多步低级算术运算
- 对同一数据结构的多次访问会产生许多公共的低级运算

□ 主要种类

- 公共子表达式删除(common subexpression elimination)
- 复写传播(copy propagation)
- 死代码删除(dead code elimination)
- 代码外提(loop hoisting, code motion)
-



一个实例

$i = m - 1; j = n; v = a[n];$

while (1) {

do $i = i + 1;$

while($a[i] < v$);

do $j = j - 1;$

while ($a[j] > v$);

if ($i \geq j$) break;

$x = a[i]; a[i] = a[j]; a[j] = x;$

}

$x = a[i]; a[i] = a[n]; a[n] = x;$

(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) $j = j - 1$

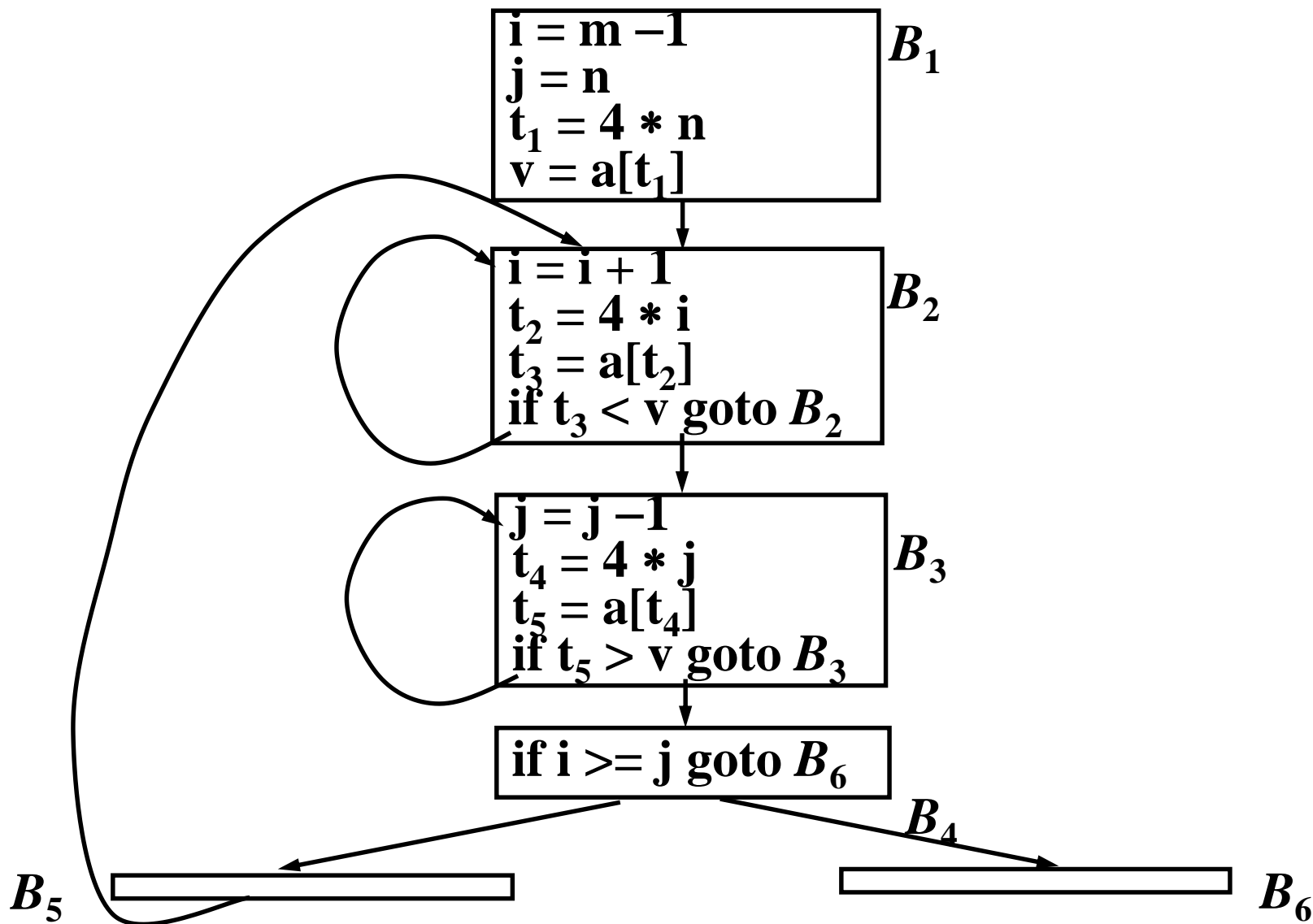
(10) $t4 = 4 * j$

(11) $t5 = a[t4]$

(12) ...



一个实例：程序流图



基本块内的优化

B_5 $x=a[i]; a[i]=a[j]; a[j]=x;$

公共子表达式删除(CSE)

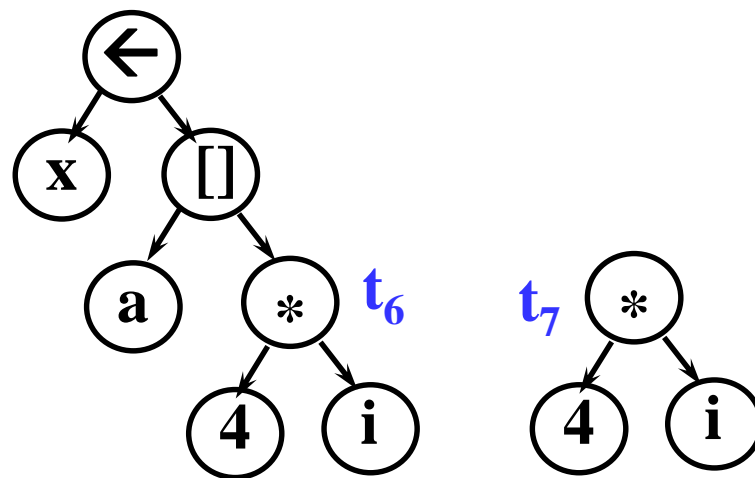
```

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
    
```

```

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

用DAG表示基本块或函数的一部分，
如LLVM中的SelectionDAG
用于简化代码，方便实现优化



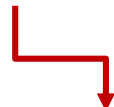
应用Value Numbering(值编码)决定两个计算是否等价,并去掉等价计算中的一个以进行保语义优化



基本块内的优化

B_5 $x=a[i]; a[i]=a[j]; a[j]=x;$

公共子表达式删除(CSE)



```
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
```

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

复写传播

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

复写传播变换的做法是在复写语句 $f=g$ 后, 尽可能用 g 代表 f



基本块内的优化

B_5 $x=a[i]; a[i]=a[j]; a[j]=x;$
公共子表达式删除(CSE)

复写传播本身不是优化, 但给其他
优化(常量合并、DCE等)带来机会

复写传播

死代码删除
(DCE)

```
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
```

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

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

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



□ 死代码

- 死代码指计算的结果绝不被引用的语句
- 一些优化变换可能会引起死代码

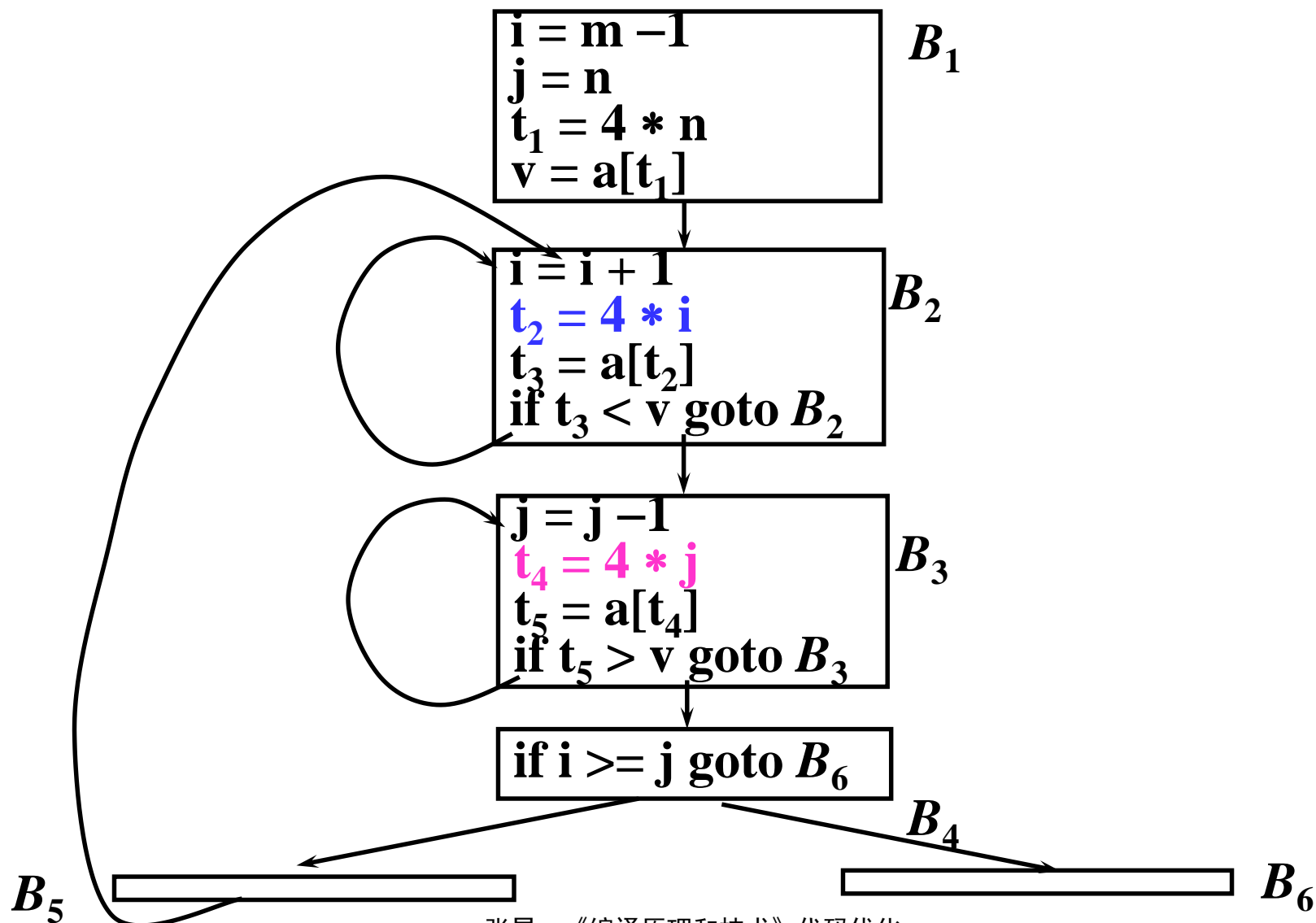
例： 为便于调试，可能在程序中加打印语句，测试后改成右边的形式

| | | |
|-----------------------------------|--|-----------------------------------|
| <code>debug = true;</code> | | <code>debug = false;</code> |
| <code>...</code> | | <code>...</code> |
| <code>if (debug) print ...</code> | | <code>if (debug) print ...</code> |

靠优化来保证目标代码中没有该条件语句部分



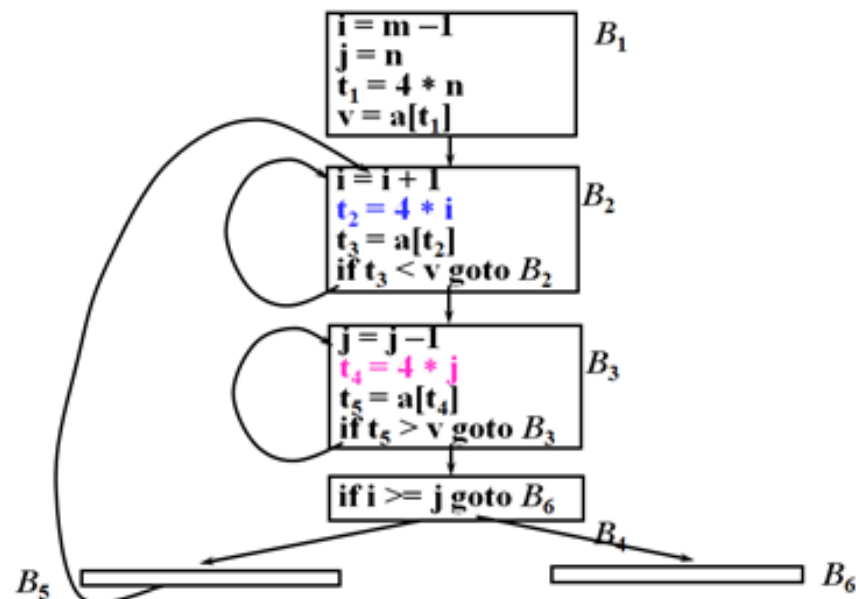
基本块间的优化



基本块间的优化

B_5 $x = a[i]; a[i] = a[j]; a[j] = x;$

全局CSE、复写传播、DCE



$t_6 = 4 * i$
 $x = a[t_6]$
 $t_7 = 4 * i$
 $t_8 = 4 * j$
 $t_9 = a[t_8]$
 $a[t_7] = t_9$
 $t_{10} = 4 * j$
 $a[t_{10}] = x$
 goto B_2

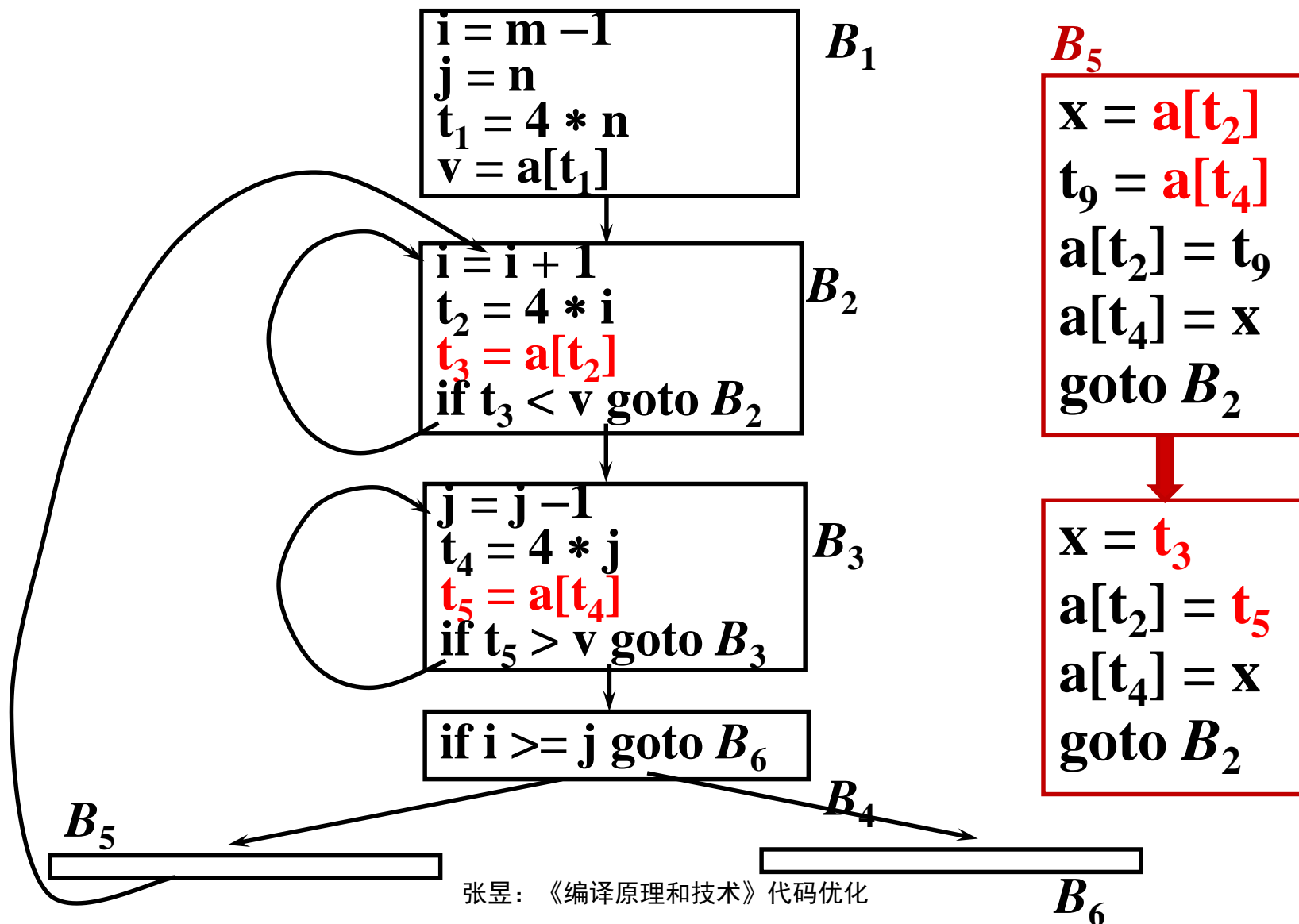
$t_6 = 4 * i$
 $x = a[t_6]$
 $t_8 = 4 * j$
 $t_9 = a[t_8]$
 $a[t_6] = t_9$
 $a[t_8] = x$
 goto B_2

$x = a[t_2]$
 $t_9 = a[t_4]$
 $a[t_2] = t_9$
 $a[t_4] = x$
 goto B_2

$x = t_3$
 ~~$t_9 = a[t_4]$~~
 $a[t_2] = t_5$
 $a[t_4] = x$
 goto B_2



基本块间的优化





实例：B6的优化

B_6 $x = a[i]; a[i] = a[n]; a[n] = x;$

B_1 : $t_1 = 4 * n$

B_2 : $t_2 = 4 * i, t_3 = a[t_2]$

$t_{11} = 4 * i$
 $x = a[t_{11}]$
 $t_{12} = 4 * i$
 $t_{13} = 4 * n$
 $t_{14} = a[t_{13}]$
 $a[t_{12}] = t_{14}$
 $t_{15} = 4 * n$
 $a[t_{15}] = x$

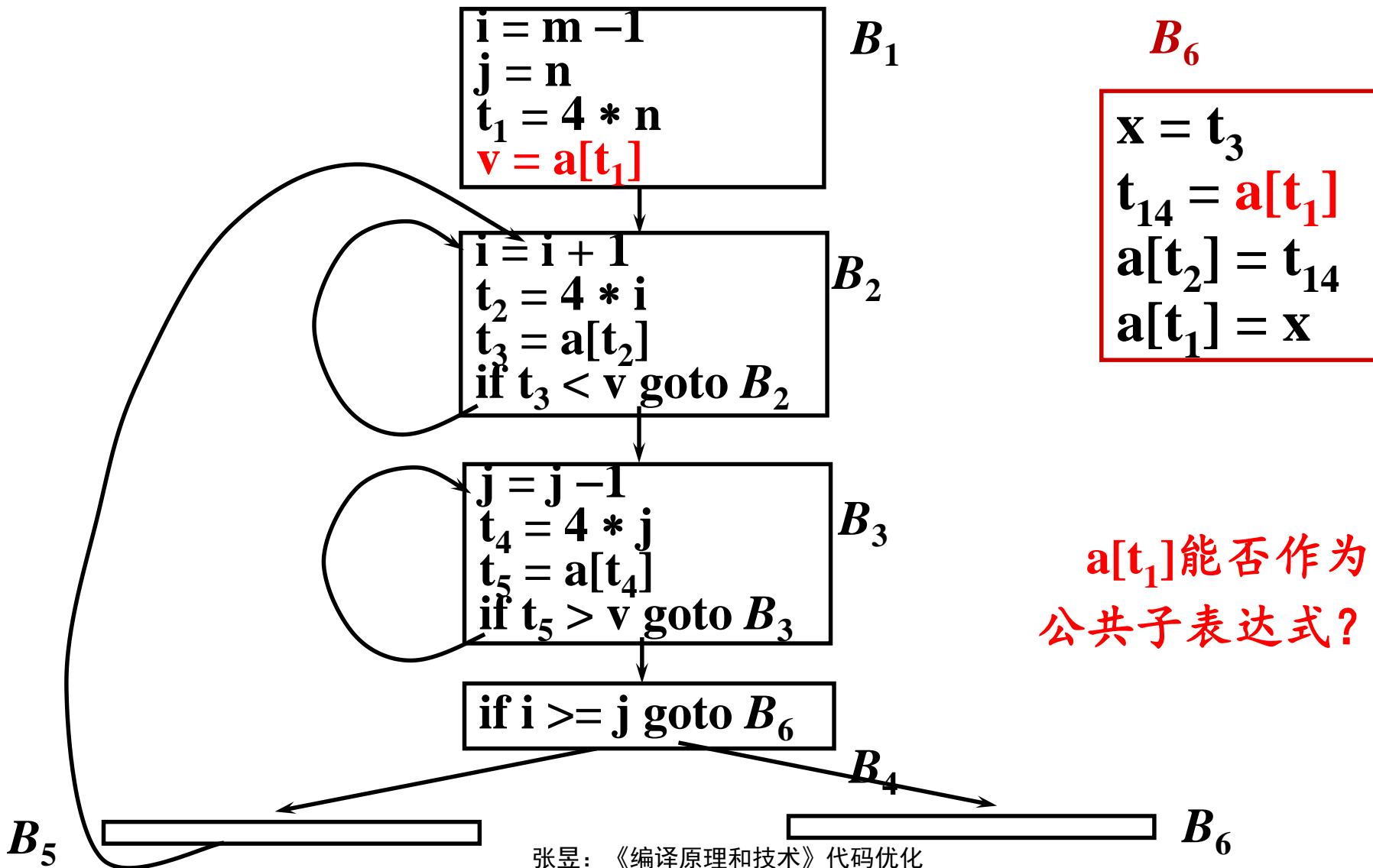
基本块内和
基本块间的优化



$x = t_3$
 $t_{14} = a[t_1]$
 $a[t_2] = t_{14}$
 $a[t_1] = x$



控制流对优化的影响

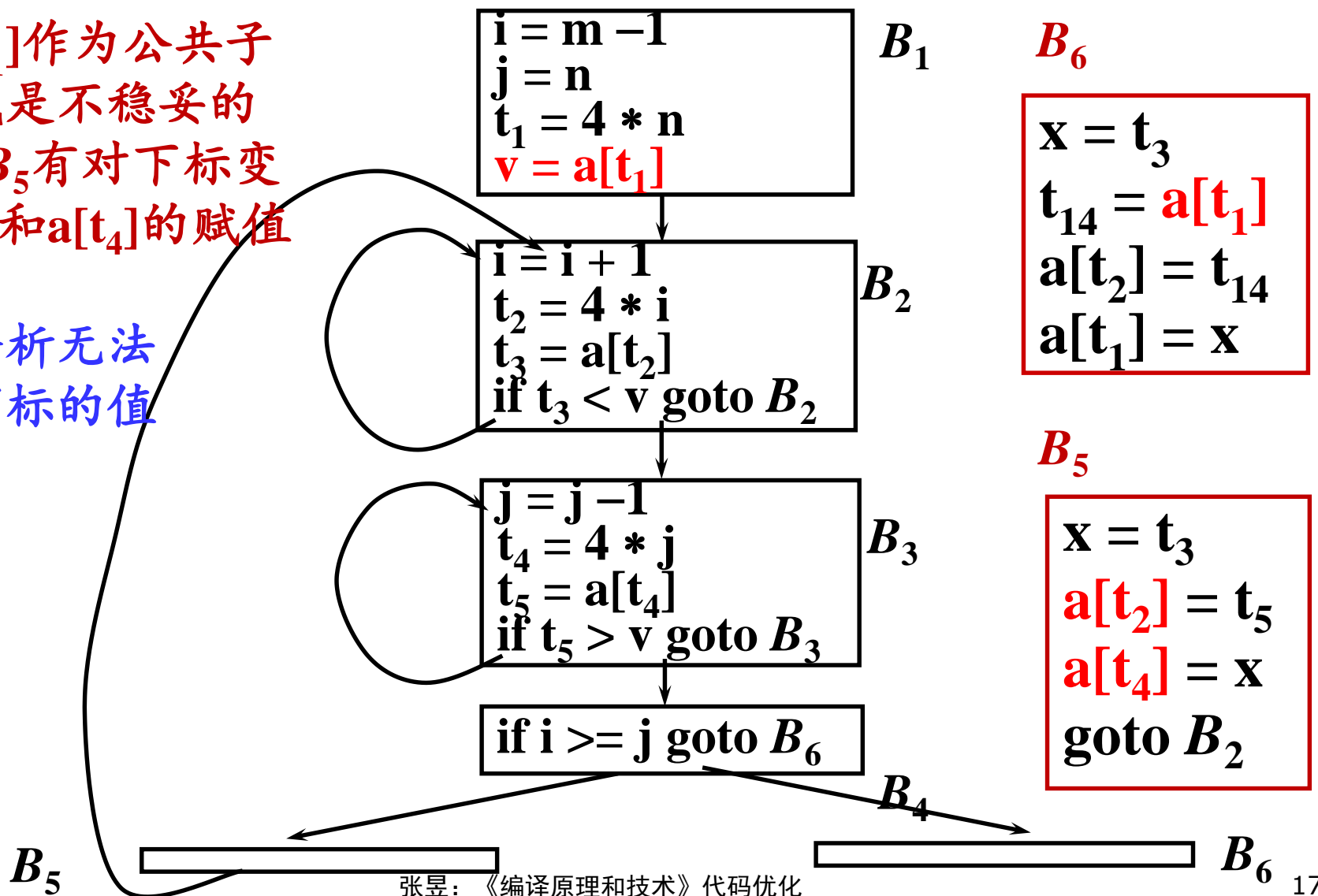




控制流对优化的影响

把 $a[t_1]$ 作为公共子
表达式是不稳妥的
因为 B_5 有对下标变
量 $a[t_2]$ 和 $a[t_4]$ 的赋值

静态分析无法
确定下标的值





□ 循环优化的主要技术

- 归纳变量删除(induction variable elimination)
- 强度削弱(strength reduction): 如将乘法变换成加法
- 代码外提(loop hoisting, code motion): 将循环不变的运算外提

□ 代码外提

例: `while (i <= limit - 2) ...` `limit - 2`是循环不变式

代码外提后变换成

`t = limit - 2;`

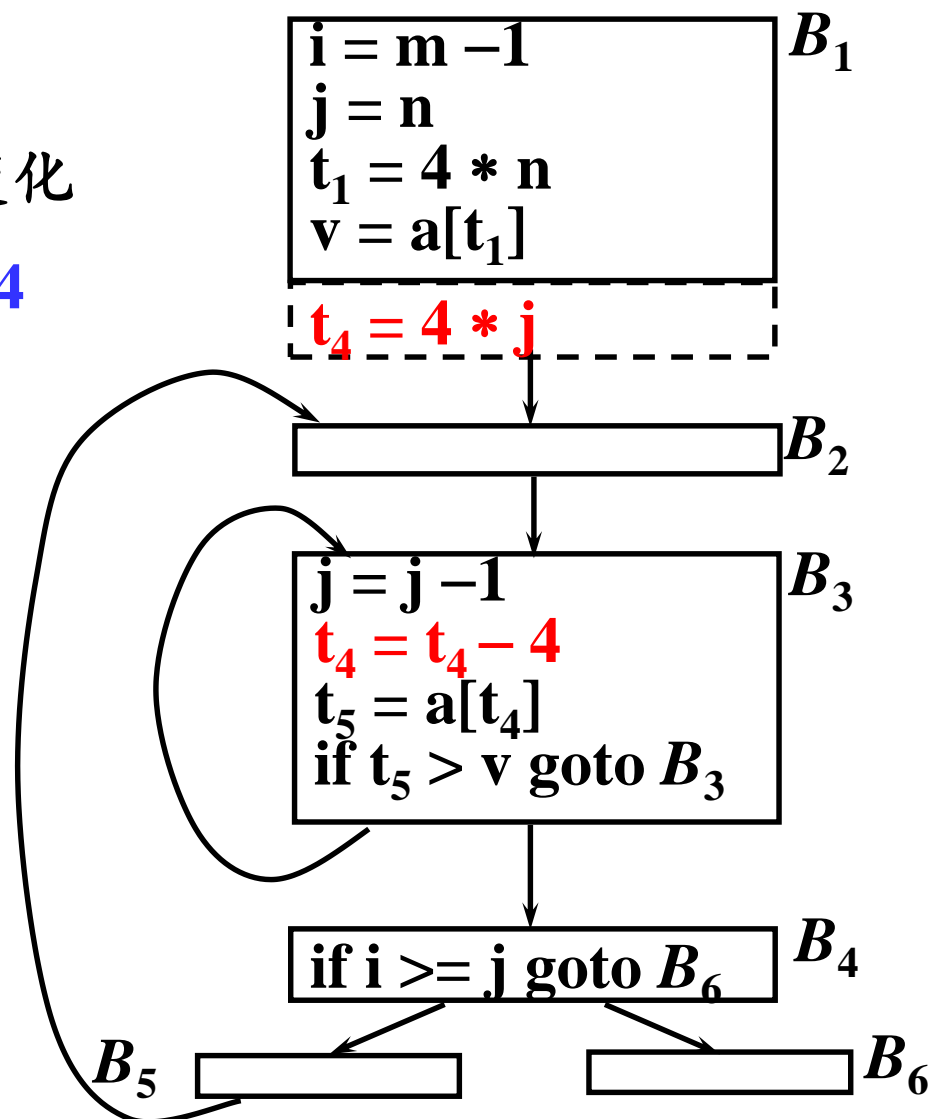
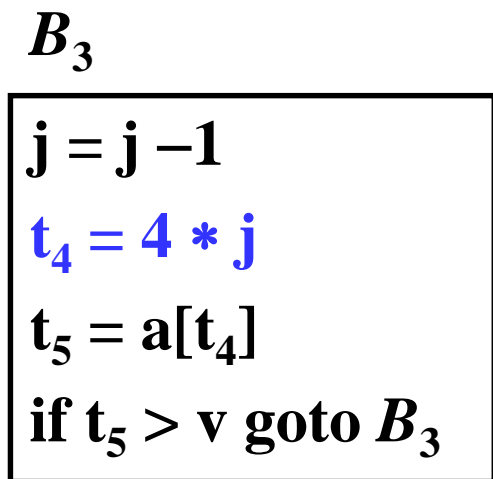
`while (i <= t) ...`



循环优化-强度削弱

□ 强度削弱

- j 和 t_4 的值步调一致地变化
- $t_4 = 4 * j$ 变换为 $t_4 = t_4 - 4$
- 在循环前增加 $t_4 = 4 * j$





循环优化-强度削弱

□ 强度削弱

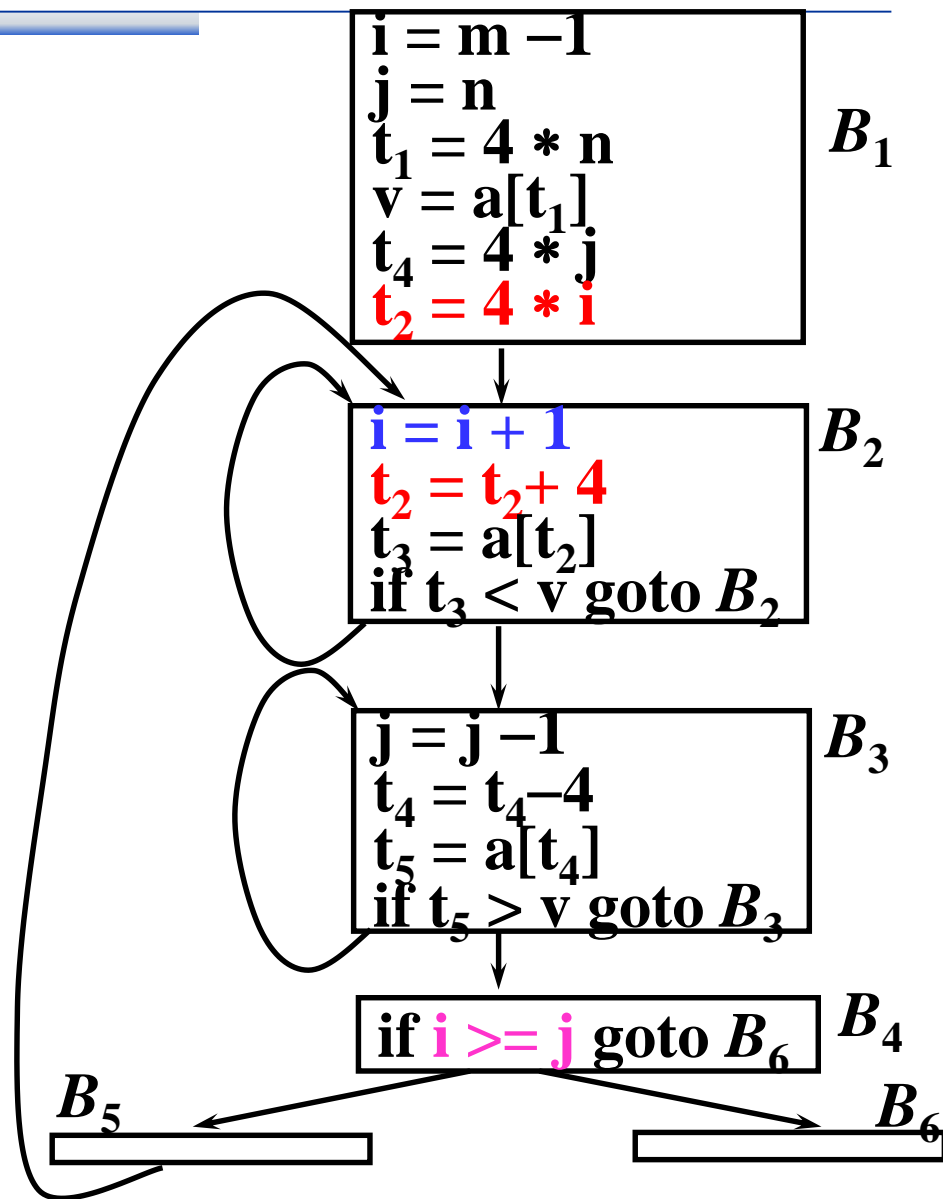
- B_2 也可以类似地变换

□ 归纳变量删除

- 循环控制条件 $i \geq j$
可以表示为 $t_2 \geq t_4$

B_2

```
i = i + 1  
t2 = 4 * i  
t3 = a[t2]  
if t3 < v goto B2
```





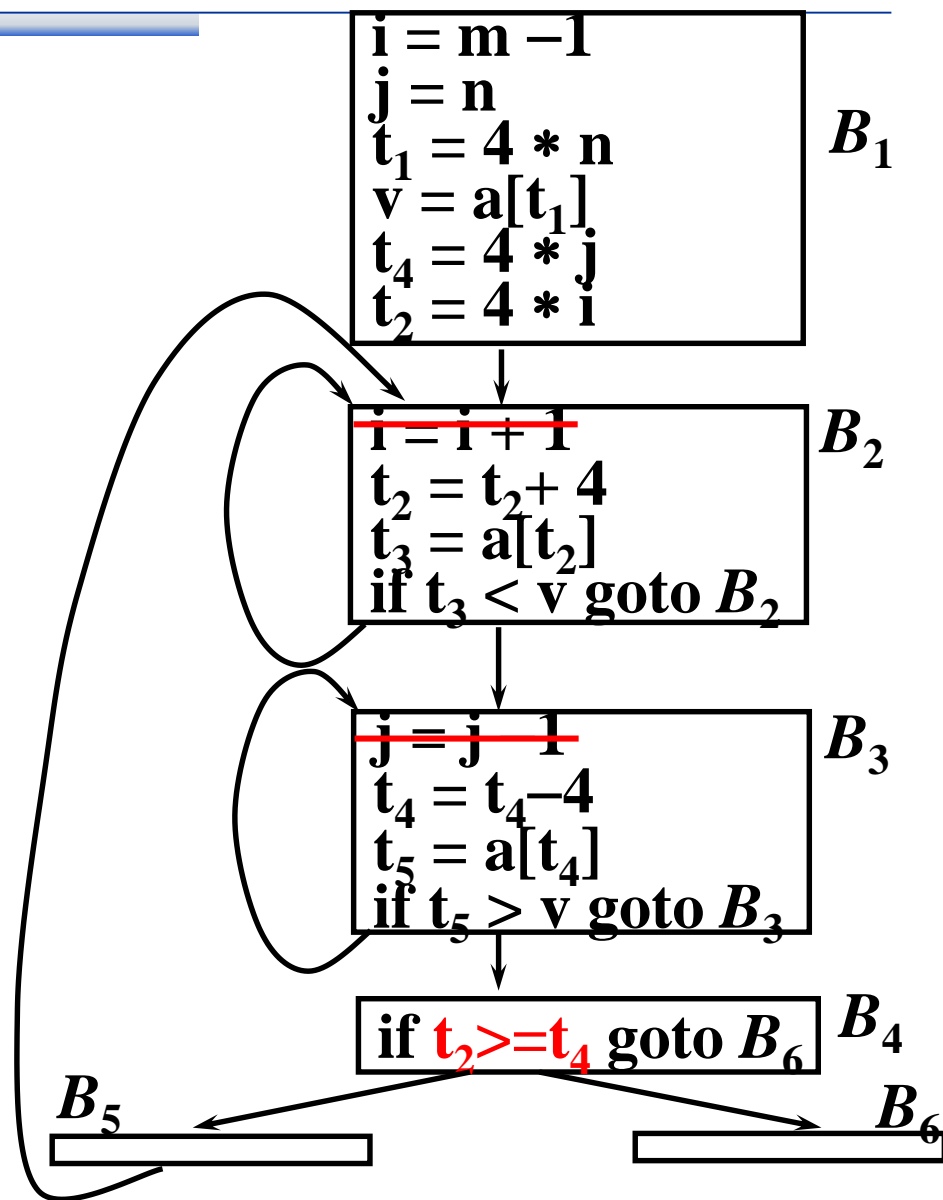
循环优化-强度削弱

□ 强度削弱

- B_2 也可以类似地变换

□ 归纳变量删除

- 循环控制条件 $i \geq j$
可以表示为 $t_2 \geq t_4$
- 归纳变量 i 和 j 可删除





代码优化总结

□ 局部优化

■ 基本块内的优化

- **窥孔(peephole)优化**：仅分析一个滑动窗口内的指令。
每次转换后可能还会暴露相邻窗口间的某些优化机会

- ✓ 冗余指令删除：如

`mov r0, a` `// r0 => a`

`mov a, r0` `// a => r0, 可删除`

- ✓ 删除死代码

`goto L1`

`goto L2` `// 语句前无标号, 死代码`



代码优化总结

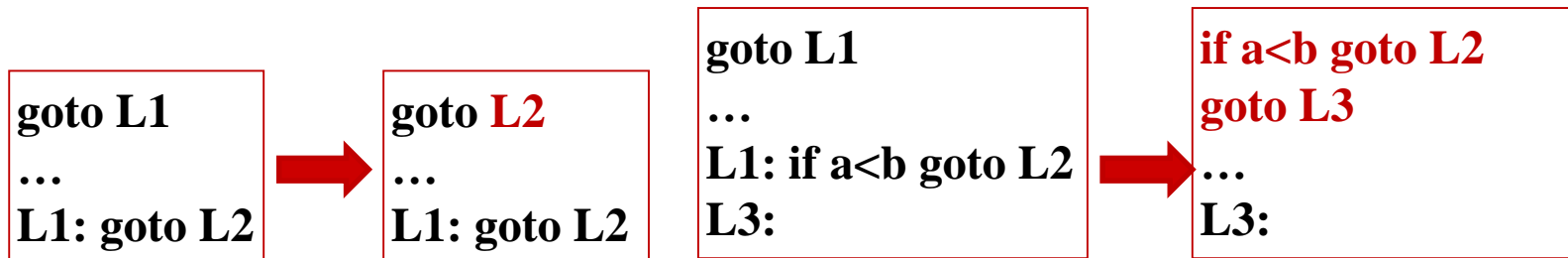
□ 局部优化

■ 基本块内的优化

- **窥孔(peephole)优化**: 仅分析一个滑动窗口内的指令。
每次转换后可能还会暴露相邻窗口间的某些优化机会

✓ 冗余指令删除、删除死代码

✓ 控制流优化





代码优化总结

□ 局部优化

■ 基本块内的优化

- **窥孔(peephole)优化**：仅分析一个滑动窗口内的指令。
每次转换后可能还会暴露相邻窗口间的某些优化机会

- ✓ 冗余指令删除、删除死代码、控制流优化
- ✓ 强度削弱、删除无用指令

`mul $8, r0 ==> shiftright $3, r0`

`add $0, r1 mul $1, r2 // 均可删除`

- ✓ 利用目标机指令特点

如 `inc`、`enter`(建立栈帧)、`leave`(清除栈帧)、
龙芯的乘加指令、向量扩展指令等等



代码优化总结

□ 局部优化

- 基本块内的优化
- **窥孔(peephole)优化**：仅分析一个滑动窗口内的指令。
每次转换后可能还会暴露相邻窗口间的某些优化机会

□ 全局优化

- 基本块间优化（过程内）
- 过程间优化：程序全局优化



2. 程序分析

- 控制流分析
- 数据流分析：到达-定值分析、数据流分析模式、活跃变量分析、可用表达式分析等



□ 控制流分析

- 发现每一个过程内的控制流层次结构
- 从构成过程的**基本块**开始，然后构造**流图**
- 使用**必经结点**来找出**循环**

□ 数据流分析

- 确定一个过程中与**数据处理**有关的全局信息
- 例如，常量传播分析是力求判定对一个特定变量的所有赋值在某个特定程序点是否总是给定相同的常数值。
如果是这样，则在那一点可用一个常数来替代该变量



数据流分析举例

程序点(5)的所有程序状态

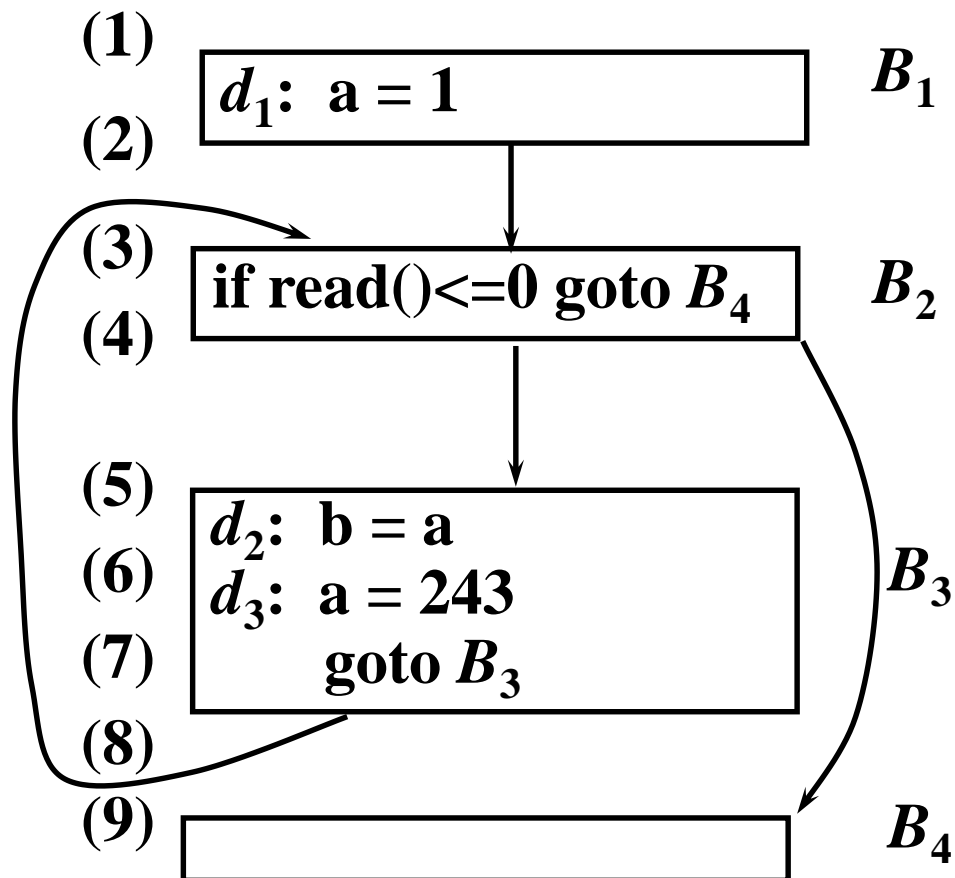
- $a \in \{1, 243\}$
- a 由 $\{d_1, d_3\}$ 定值

1. 到达一定值

$\{d_1, d_3\}$ 的定值到达点(5)

2. 常量合并

a 在点(5)不是常量





数据流分析

□ 数据流分析

- 分析程序行为时，必须在其流图上考虑**所有的执行路径**
(在调用或返回语句被执行时，还需要考虑执行路径在**多个流图之间的跳转**)
- 由于存在循环，从流图得到的程序**执行路径数无限**，且执行路径长度没有有限的上界
- 每个程序点的**不同状态数也可能无限**
程序状态：存储单元到值的映射

不可能清楚所有执行路径上的所有程序状态



数据流抽象

□ 数据流分析不打算

- 区分到达一个程序点的不同执行路径
- 掌握该程序点的每个完整的状态

□ 数据流分析要点

- 从程序状态中**抽取**特定数据流分析所需信息→**数据流值**
- 总结出用于该分析目的的一组**有限的事实**，且这组事实和到达这个程序点的路径无关，即从任何路径到达该程序点都有这样的事实
保守的：得到的信息不会误解程序的行为
- 分析的目的不同，从程序状态提炼的信息也不同



□ 数据流值：到达一个程序点的所有定值

- 每个定值语句有唯一编号，数据流值：定值编号的集合
可用来判断一个变量在某程序点是否为常量、是否无初值

□ 别名给到达-定值的计算带来困难

- 过程参数、数组访问、间接引用等都有可能引起别名
例如：若 $p == q$ ，则 $p \rightarrow next$ 和 $q \rightarrow next$ 互为别名
- 程序分析必须是稳妥的
- 本章其余部分仅考虑变量无别名的情况

□ 定值的注销(kill)

在一条执行路径上，对 x 的赋值会注销先前对 x 的所有赋值



到达-定值的计算

□ 基本块向下暴露的定值

■ $gen[B]$: 基本块 B 生成的定值

■ $kill[B]$: 基本块 B 注销的定值

$gen[B_1] = \{d_1, d_2, d_3\}$

$kill[B_1] = \{d_4, d_5, d_6, d_7\}$

$gen[B_2] = \{d_4, d_5\}$

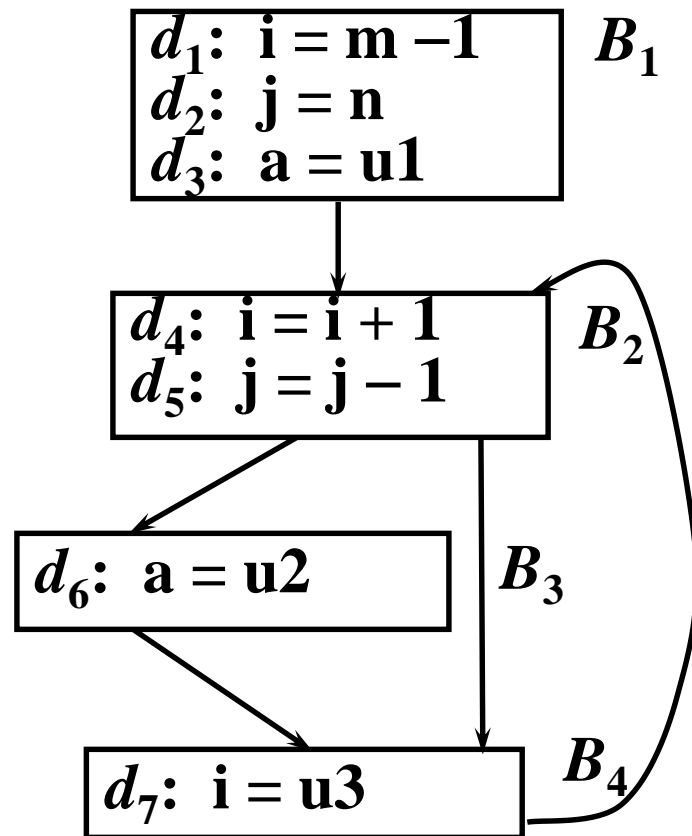
$kill[B_2] = \{d_1, d_2, d_7\}$

$gen[B_3] = \{d_6\}$

$kill[B_3] = \{d_3\}$

$gen[B_4] = \{d_7\}$

$kill[B_4] = \{d_1, d_4\}$





基本块的 *gen* 和 *kill* 的计算

- 对三地址指令 $d: u = v + w$ ，它的状态迁移函数是

$$f_d(x) = gen_d \cup (x - kill_d),$$

x 为到达指令入口点的定值

- 若 $f_1(x) = gen_1 \cup (x - kill_1)$, $f_2(x) = gen_2 \cup (x - kill_2)$

$$\begin{aligned} \text{则: } f_2(f_1(x)) &= gen_2 \cup (gen_1 \cup (x - kill_1) - kill_2) \\ &= (gen_2 \cup (gen_1 - kill_2)) \cup (x - (kill_1 \cup kill_2)) \end{aligned}$$

- 若基本块 B 有 n 条三地址指令

$$kill_B = kill_1 \cup kill_2 \cup \dots \cup kill_n$$

$$gen_B = gen_n \cup (gen_{n-1} - kill_n) \cup (gen_{n-2} - kill_{n-1} - kill_n) \cup \dots \cup (gen_1 - kill_2 - kill_3 - \dots - kill_n)$$



到达一定值的数据流方程

□ 数据流方程/等式(flow equation)

■ gen_B : B 中生成的且能到达 B 的结束点的定值

■ $kill_B$: B 注销的定值

■ $IN[B]$: 能到达 B 的开始点的定值集合

■ $OUT[B]$: 能到达 B 的结束点的定值集合

两组等式（根据 gen 和 $kill$ 定义 IN 和 OUT ）

■ $IN[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT[P]$

■ $OUT[B] = gen_B \cup (IN[B] - kill_B)$

■ $OUT[ENTRY] = \emptyset$

到达一定值方程组的迭代求解，最终到达不动点(MFP)



到达-定值的迭代计算算法

// 正向数据流分析

(1) $\text{OUT}[\text{ENTRY}] = \emptyset;$

(2) for (除了ENTRY以外的每个块 B) $\text{OUT}[B] = \emptyset;$

(3) while (任何一个OUT出现变化)

(4) for (除了ENTRY以外的每个块 B) {

(5) $\text{IN}[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P];$

(6) $\text{OUT}[B] = f_B(\text{IN}[B]);$

 // $f_B(\text{IN}[B]) = \text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$

(7) }

这里采用流不敏感
的算法，即没有按基本块的
执行次序依次处理
每个基本块



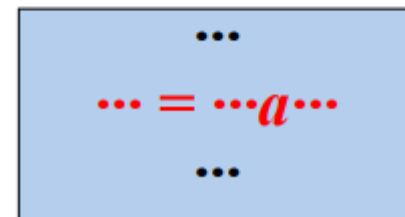
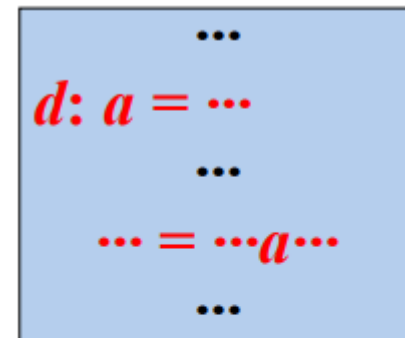
引用-定值(Use-Def)

□ 引用-定值链(简称ud链)

- 对于变量的每一个引用，记录到达该引用的所有定值

□ 通过到达-定值数据流分析计算并形成ud链

- 如果块B中变量 a 的引用之前有 a 的定值，那么只有 a 的最后一次定值会在该引用的ud链中
- 如果块B中变量 a 的引用之前无 a 的定值，那么该引用的ud链就是IN[B]中 a 的定值的集合





到达-定值计算示例

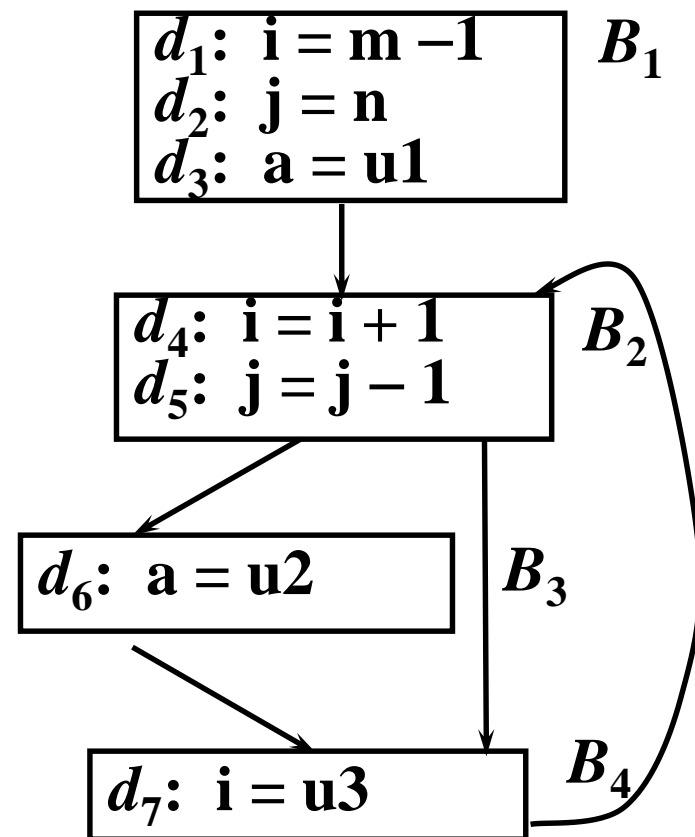
| | IN [B] | OUT [B] |
|-------|--------|----------|
| B_1 | | 000 0000 |
| B_2 | | 000 0000 |
| B_3 | | 000 0000 |
| B_4 | | 000 0000 |

$gen [B_1] = \{d_1, d_2, d_3\}$
 $kill [B_1] = \{d_4, d_5, d_6, d_7\}$

$gen [B_2] = \{d_4, d_5\}$
 $kill [B_2] = \{d_1, d_2, d_7\}$

$gen [B_3] = \{d_6\}$
 $kill [B_3] = \{d_3\}$

$gen [B_4] = \{d_7\}$
 $kill [B_4] = \{d_1, d_4\}$





到达-定值计算示例

$$\text{IN}[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P]$$
$$\text{OUT}[B] = \text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$$

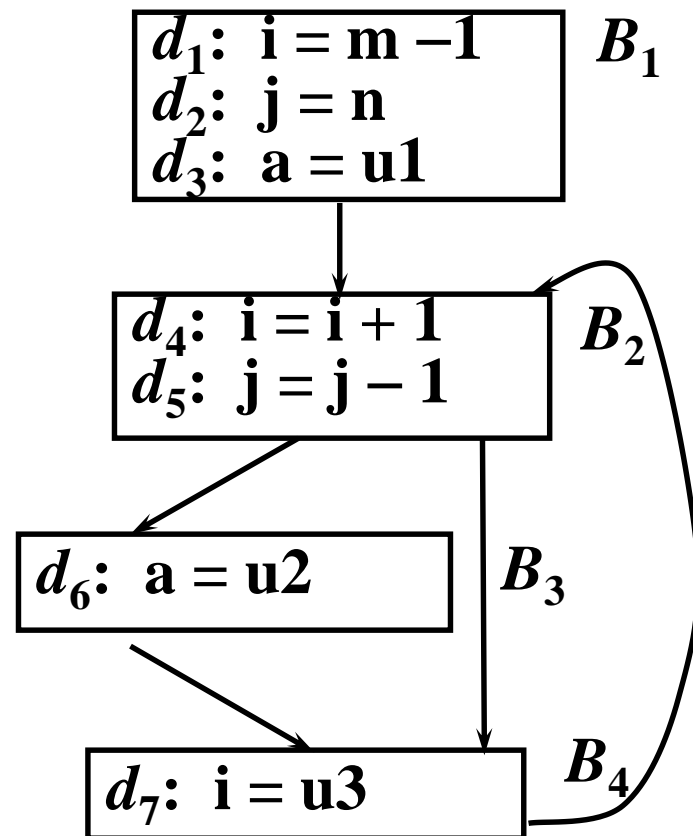
| | IN [B] | OUT [B] |
|-------|----------|----------|
| B_1 | 000 0000 | 000 0000 |
| B_2 | | 000 0000 |
| B_3 | | 000 0000 |
| B_4 | | 000 0000 |

$\text{gen}[B_1] = \{d_1, d_2, d_3\}$
 $\text{kill}[B_1] = \{d_4, d_5, d_6, d_7\}$

$\text{gen}[B_2] = \{d_4, d_5\}$
 $\text{kill}[B_2] = \{d_1, d_2, d_7\}$

$\text{gen}[B_3] = \{d_6\}$
 $\text{kill}[B_3] = \{d_3\}$

$\text{gen}[B_4] = \{d_7\}$
 $\text{kill}[B_4] = \{d_1, d_4\}$





到达-定值计算示例

$IN[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT[P]$

$OUT[B] = gen_B \cup (IN[B] - kill_B)$

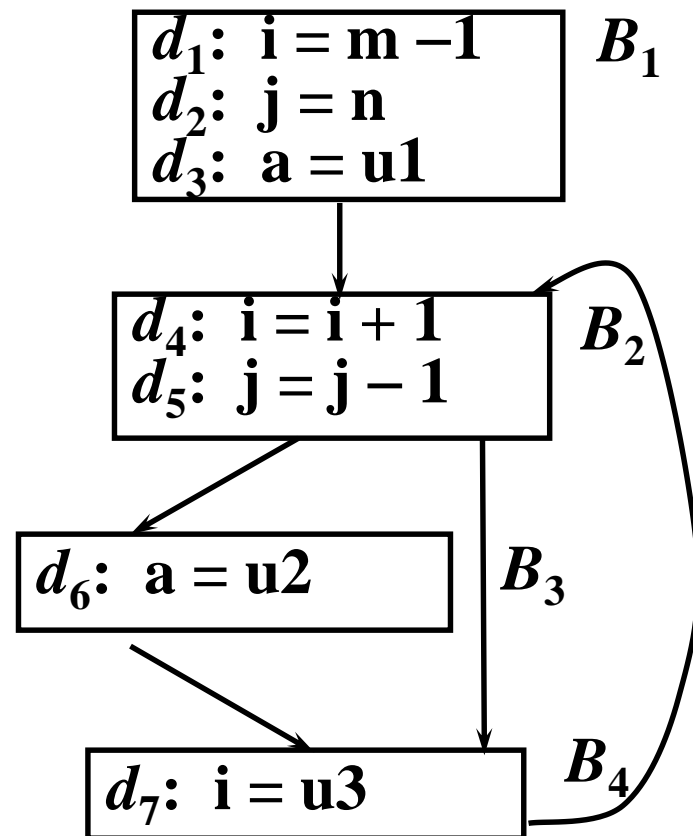
| | IN [B] | OUT [B] |
|-------|----------|----------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | | 000 0000 |
| B_3 | | 000 0000 |
| B_4 | | 000 0000 |

$gen [B_1] = \{d_1, d_2, d_3\}$
 $kill [B_1] = \{d_4, d_5, d_6, d_7\}$

$gen [B_2] = \{d_4, d_5\}$
 $kill [B_2] = \{d_1, d_2, d_7\}$

$gen [B_3] = \{d_6\}$
 $kill [B_3] = \{d_3\}$

$gen [B_4] = \{d_7\}$
 $kill [B_4] = \{d_1, d_4\}$





到达-定值计算示例

$$\text{IN}[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P]$$
$$\text{OUT}[B] = \text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$$

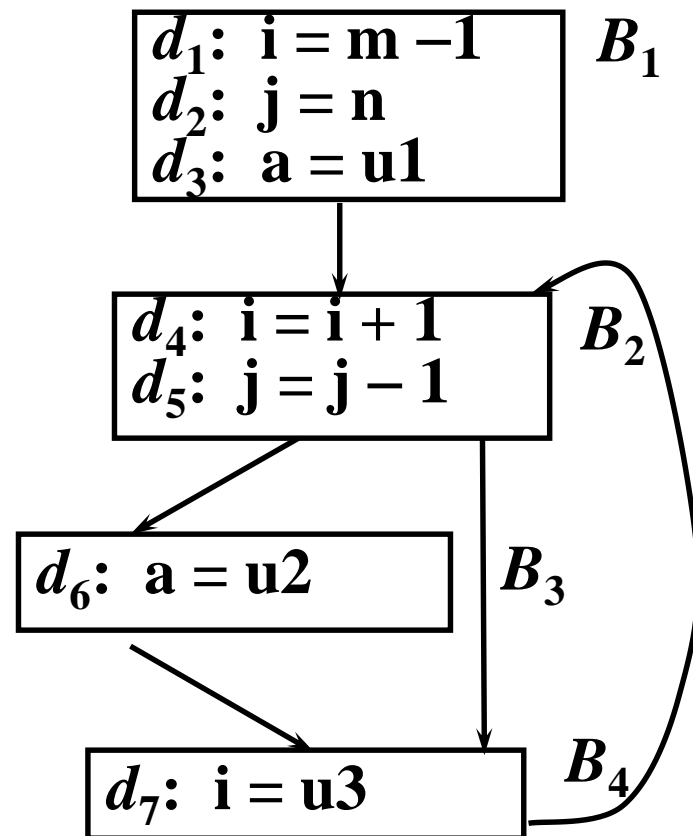
| | IN [B] | OUT [B] |
|-------|----------|----------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0000 | 000 0000 |
| B_3 | | 000 0000 |
| B_4 | | 000 0000 |

$\text{gen}[B_1] = \{d_1, d_2, d_3\}$
 $\text{kill}[B_1] = \{d_4, d_5, d_6, d_7\}$

$\text{gen}[B_2] = \{d_4, d_5\}$
 $\text{kill}[B_2] = \{d_1, d_2, d_7\}$

$\text{gen}[B_3] = \{d_6\}$
 $\text{kill}[B_3] = \{d_3\}$

$\text{gen}[B_4] = \{d_7\}$
 $\text{kill}[B_4] = \{d_1, d_4\}$





到达-定值计算示例

$IN[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT[P]$

$OUT[B] = gen_B \cup (IN[B] - kill_B)$

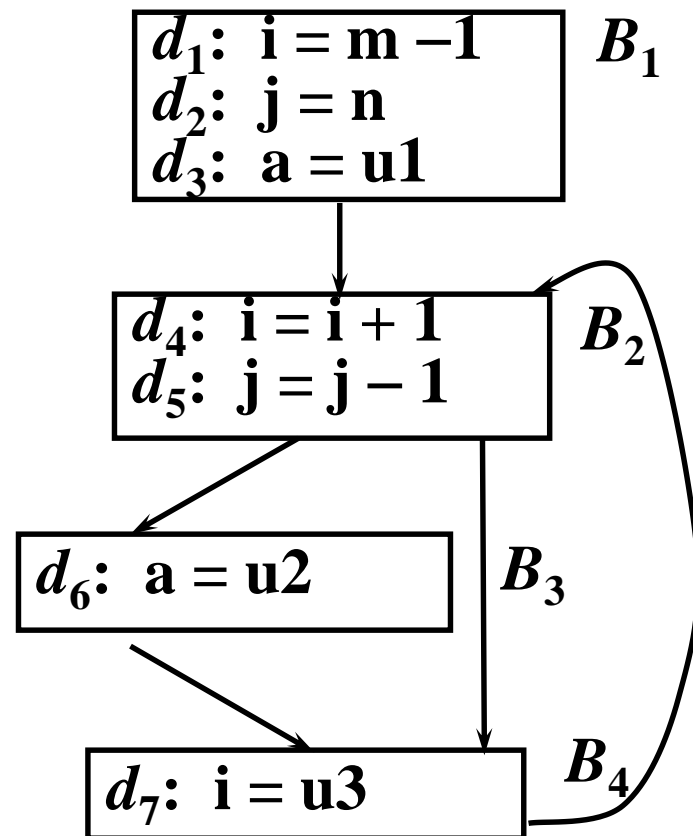
| | IN [B] | OUT [B] |
|-------|----------|----------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0000 | 001 1100 |
| B_3 | | 000 0000 |
| B_4 | | 000 0000 |

$gen [B_1] = \{d_1, d_2, d_3\}$
 $kill [B_1] = \{d_4, d_5, d_6, d_7\}$

$gen [B_2] = \{d_4, d_5\}$
 $kill [B_2] = \{d_1, d_2, d_7\}$

$gen [B_3] = \{d_6\}$
 $kill [B_3] = \{d_3\}$

$gen [B_4] = \{d_7\}$
 $kill [B_4] = \{d_1, d_4\}$





到达-定值计算示例

$$\text{IN}[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P]$$
$$\text{OUT}[B] = \text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$$

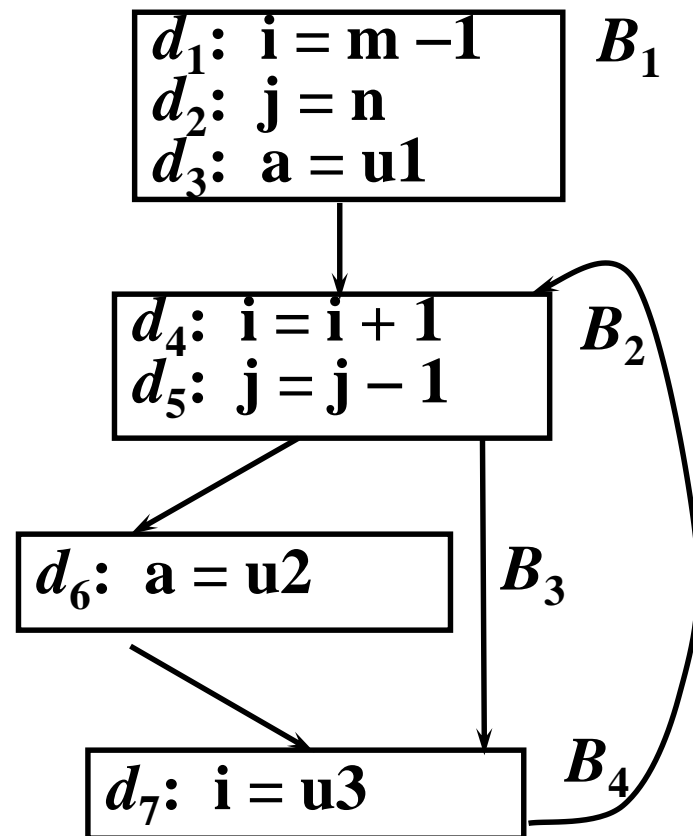
| | IN [B] | OUT [B] |
|-------|----------|----------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0000 | 001 1100 |
| B_3 | 001 1100 | 000 0000 |
| B_4 | | 000 0000 |

$\text{gen}[B_1] = \{d_1, d_2, d_3\}$
 $\text{kill}[B_1] = \{d_4, d_5, d_6, d_7\}$

$\text{gen}[B_2] = \{d_4, d_5\}$
 $\text{kill}[B_2] = \{d_1, d_2, d_7\}$

$\text{gen}[B_3] = \{d_6\}$
 $\text{kill}[B_3] = \{d_3\}$

$\text{gen}[B_4] = \{d_7\}$
 $\text{kill}[B_4] = \{d_1, d_4\}$





到达-定值计算示例

$IN[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT[P]$

$OUT[B] = gen_B \cup (IN[B] - kill_B)$

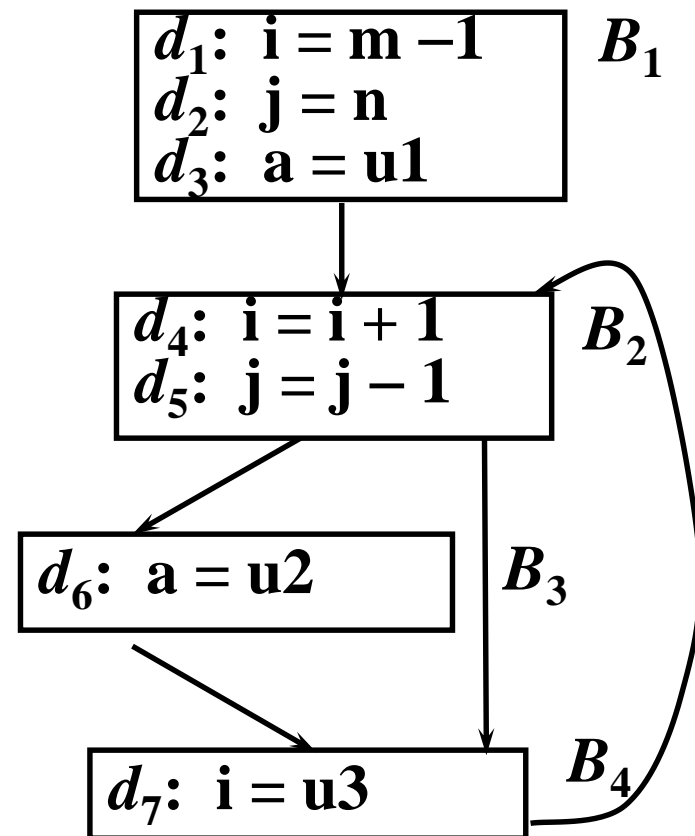
| | IN [B] | OUT [B] |
|-------|----------|-----------------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0000 | 001 1100 |
| B_3 | 001 1100 | 000 1110 |
| B_4 | | 000 0000 |

$gen [B_1] = \{d_1, d_2, d_3\}$
 $kill [B_1] = \{d_4, d_5, d_6, d_7\}$

$gen [B_2] = \{d_4, d_5\}$
 $kill [B_2] = \{d_1, d_2, d_7\}$

$gen [B_3] = \{d_6\}$
 $kill [B_3] = \{d_3\}$

$gen [B_4] = \{d_7\}$
 $kill [B_4] = \{d_1, d_4\}$





到达-定值计算示例

$$\text{IN}[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P]$$
$$\text{OUT}[B] = \text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$$

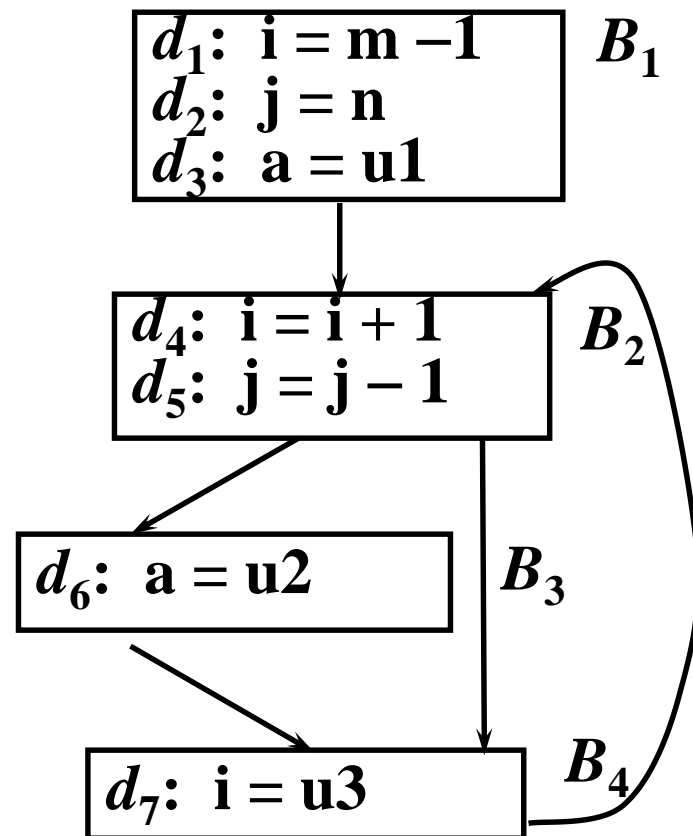
| | IN [B] | OUT [B] |
|-------|-----------------|----------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0000 | 001 1100 |
| B_3 | 001 1100 | 000 1110 |
| B_4 | 001 1110 | 000 0000 |

$\text{gen}[B_1] = \{d_1, d_2, d_3\}$
 $\text{kill}[B_1] = \{d_4, d_5, d_6, d_7\}$

$\text{gen}[B_2] = \{d_4, d_5\}$
 $\text{kill}[B_2] = \{d_1, d_2, d_7\}$

$\text{gen}[B_3] = \{d_6\}$
 $\text{kill}[B_3] = \{d_3\}$

$\text{gen}[B_4] = \{d_7\}$
 $\text{kill}[B_4] = \{d_1, d_4\}$





到达-定值计算示例

$IN[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT[P]$

$OUT[B] = gen_B \cup (IN[B] - kill_B)$

| | IN [B] | OUT [B] |
|-------|----------|-----------------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0000 | 001 1100 |
| B_3 | 001 1100 | 000 1110 |
| B_4 | 001 1110 | 001 0111 |

$gen [B_1] = \{d_1, d_2, d_3\}$

$kill [B_1] = \{d_4, d_5, d_6, d_7\}$

$gen [B_2] = \{d_4, d_5\}$

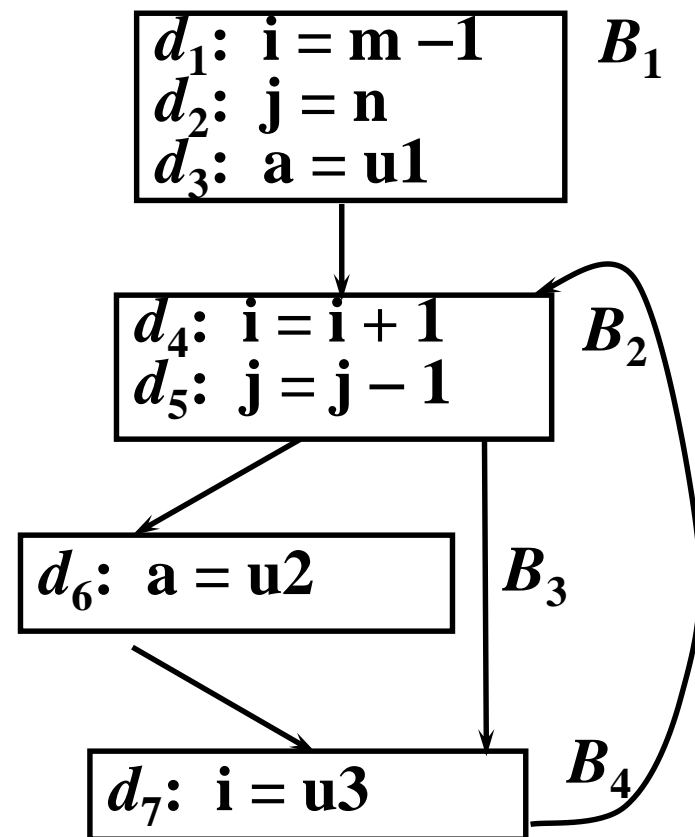
$kill [B_2] = \{d_1, d_2, d_7\}$

$gen [B_3] = \{d_6\}$

$kill [B_3] = \{d_3\}$

$gen [B_4] = \{d_7\}$

$kill [B_4] = \{d_1, d_4\}$





到达-定值计算示例

$$\text{IN}[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P]$$
$$\text{OUT}[B] = \text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$$

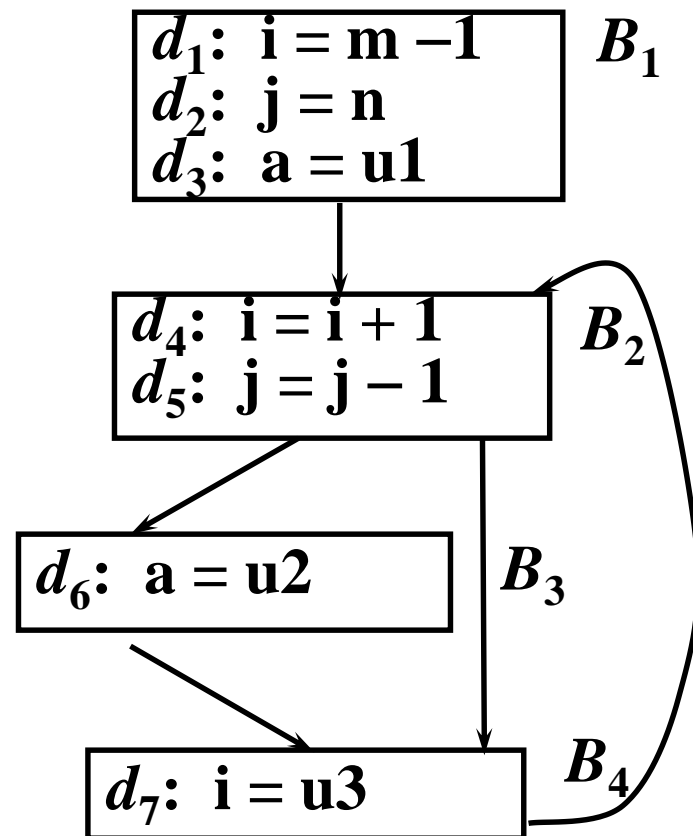
| | IN [B] | OUT [B] |
|-------|----------|----------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0111 | 001 1100 |
| B_3 | 001 1100 | 000 1110 |
| B_4 | 001 1110 | 000 0000 |

$\text{gen}[B_1] = \{d_1, d_2, d_3\}$
 $\text{kill}[B_1] = \{d_4, d_5, d_6, d_7\}$

$\text{gen}[B_2] = \{d_4, d_5\}$
 $\text{kill}[B_2] = \{d_1, d_2, d_7\}$

$\text{gen}[B_3] = \{d_6\}$
 $\text{kill}[B_3] = \{d_3\}$

$\text{gen}[B_4] = \{d_7\}$
 $\text{kill}[B_4] = \{d_1, d_4\}$





到达-定值计算示例

$IN[B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT[P]$

$OUT[B] = gen_B \cup (IN[B] - kill_B)$

| | IN [B] | OUT [B] |
|-------|----------|----------|
| B_1 | 000 0000 | 111 0000 |
| B_2 | 111 0111 | 001 1110 |
| B_3 | 001 1100 | 000 1110 |
| B_4 | 001 1110 | 001 0111 |

不再继续演示迭代计算

$gen [B_1] = \{d_1, d_2, d_3\}$

$kill [B_1] = \{d_4, d_5, d_6, d_7\}$

$gen [B_2] = \{d_4, d_5\}$

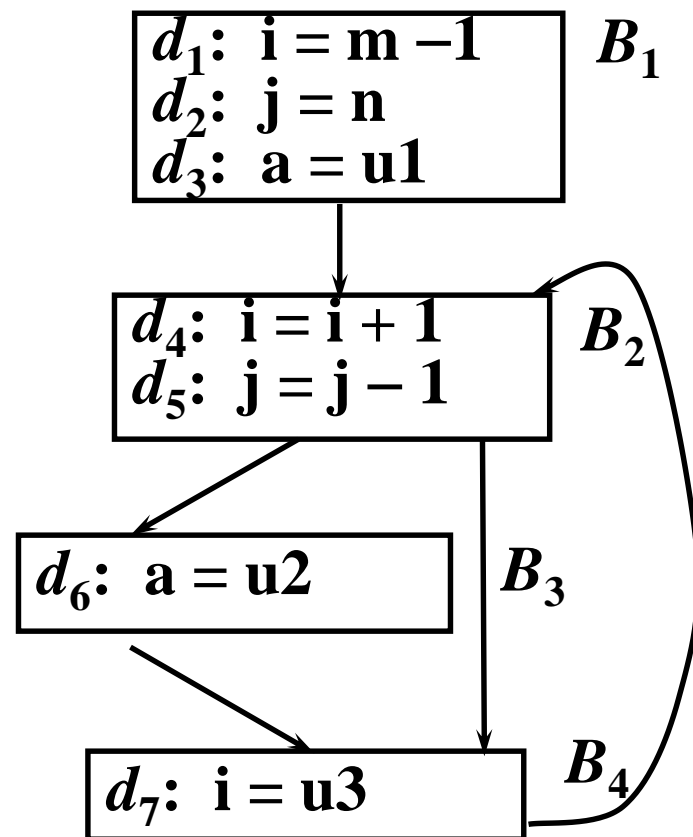
$kill [B_2] = \{d_1, d_2, d_7\}$

$gen [B_3] = \{d_6\}$

$kill [B_3] = \{d_3\}$

$gen [B_4] = \{d_7\}$

$kill [B_4] = \{d_1, d_4\}$





2. 程序分析

- 控制流分析
- 数据流分析：到达-定值分析、数据流分析模式、活跃变量分析、可用表达式分析等



到达-定值计算

□ 到达-定值数据流等式(方程)是**正向**的方程

■ $OUT [B] = gen [B] \cup (IN [B] - kill [B])$

■ $IN [B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT [P]$

某些数据流等式是反向的

□ 到达-定值数据流等式的**合流**运算是**求并集**

■ $IN [B] = \bigcup_{P \text{ 是 } B \text{ 的前驱}} OUT [P]$

某些数据流等式的合流运算是求交集

□ 对到达-定值数据流等式，迭代求它的**最小解**

某些数据流方程可能需要求最大解



数据流分析模式

□ 数据流值

- 数据流分析总把程序点和数据流值联系起来
- 数据流值代表在程序点能观测到的所有可能程序状态集合的一个抽象
- 语句 s 前后两点数据流值用 $IN[s]$ 和 $OUT[s]$ 来表示
- 数据流问题就是通过基于语句语义的约束（迁移函数）和基于控制流的约束来寻找所有语句 s 的 $IN[s]$ 和 $OUT[s]$ 的一个解



数据流分析模式

□ 迁移函数(transfer function) f

- 语句前后两点的数据流值受该语句的语义约束
- 若沿执行路径**正向**传播, 则 $\text{OUT}[s] = f_s(\text{IN}[s])$
- 若沿执行路径**逆向**传播, 则 $\text{IN}[s] = f_s(\text{OUT}[s])$

若基本块 B 由语句 s_1, s_2, \dots, s_n 依次组成, 则

- **正向**: $\text{IN}[s_{i+1}] = \text{OUT}[s_i], i = 1, 2, \dots, n-1$
逆向: $\text{OUT}[s_{i-1}] = \text{IN}[s_i], i = 2, 3, \dots, n$
- **正向**: $f_B = f_n \circ \dots \circ f_2 \circ f_1$
逆向: $f_B = f_1 \circ \dots \circ f_{n-1} \circ f_n$
- **正向**: $\text{OUT}[B] = f_B(\text{IN}[B])$
逆向: $\text{IN}[B] = f_B(\text{OUT}[B])$



□ 控制流约束

■ 正向传播

$$\text{IN}[B] = Y_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P]$$

■ 逆向传播

$$\text{OUT}[B] = Y_{S \text{ 是 } B \text{ 的后继}} \text{IN}[S]$$

□ 约束方程组的解通常不是唯一的

- 求解的目标是要找到满足这两组约束（控制流约束和迁移约束）的最“精确”解



2. 程序分析

- 控制流分析
- 数据流分析：到达-定值分析、数据流分析模式、活跃变量分析、可用表达式分析等



活跃变量(live-variable)

□ 定义

- 变量 x 的值在 p 点开始的某条执行路径上被引用，则说 x 在 p 点活跃，否则称 x 在 p 点已经死亡
- $IN[B]$: 块 B 开始点的活跃变量集合
- $OUT[B]$: 块 B 结束点的活跃变量集合
- use_B : 块 B 中有引用且在引用前无定值的变量集
- def_B : 块 B 中有定值的变量集

□ 应用

- 一种重要应用就是基本块的寄存器分配



□ 例

$use[B_2] = \{ i, j \}, def[B_2] = \{ i, j \}$

□ 活跃变量数据流等式

■ $IN[B] = use_B \cup (OUT[B] - def_B)$

■ $OUT[B] = \bigcup_{S \text{ 是 } B \text{ 的后继}} IN[S]$

■ $IN[EXIT] = \emptyset$

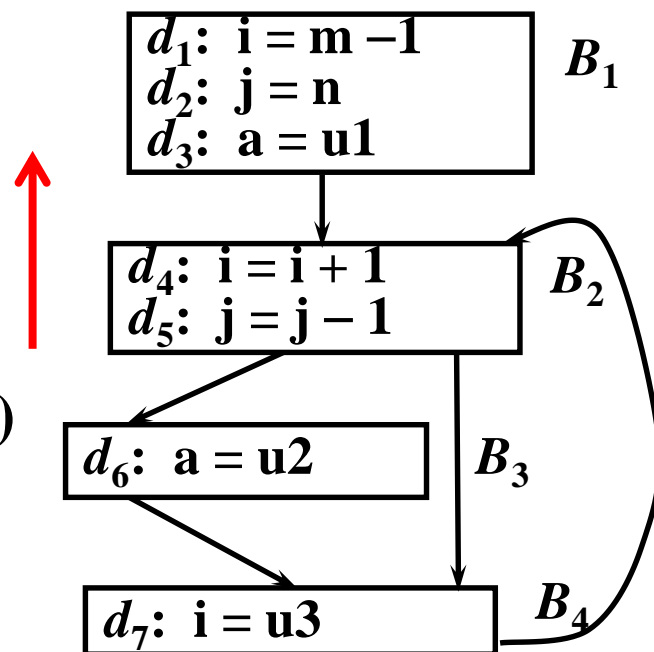
□ 和到达一定值等式之间的联系与区别

■ 汇合算符：都是集合并算符

■ 信息流动方向相反，IN和OUT的作用相互交换

■ use 和 def 分别取代 gen 和 $kill$

■ 仍然需要最小解





定值-引用(Def-Use)

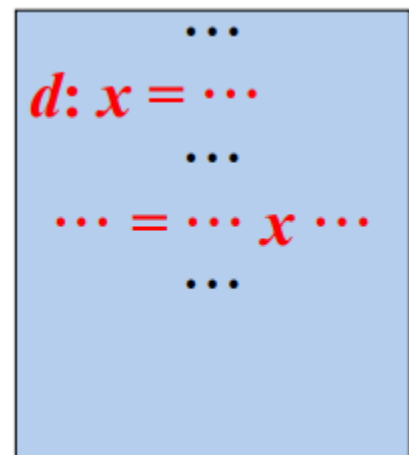
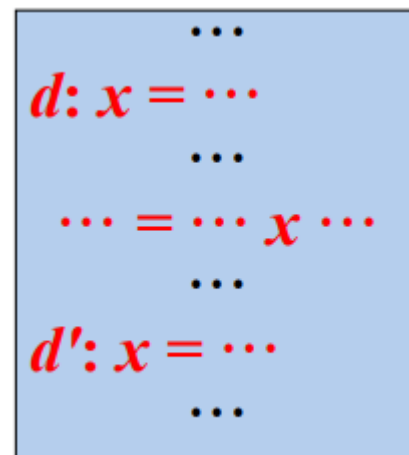
□ 定值-引用链(du链)

- 设变量 x 有某定值 d , 该定值能到达的所有引用 u 的集合称为 x 在 d 处的定值-引用链, 简称du链

□ 通过活跃变量分析来计算du链

OUT[B]: 从B的末尾处能到达的引用集合

- 如果B中 x 的定值 d 之后有 x 的第一个定值 d' , 则 d 和 d' 之间 x 的所有引用构成 d 的du链
- 如果B中 x 的定值 d 之后无 x 的新的定值, 则B中 d 之后 x 的所有引用以及OUT[B]中 x 的所有引用构成 d 的du链





LLVM IR的UD和DU链

□ Value类

- 代表一个带类型的数据

Constant, Argument, Instruction, Function 等

- 维护一个该数据使用者的列表

□ 该表描述def-use信息

- 主要方法

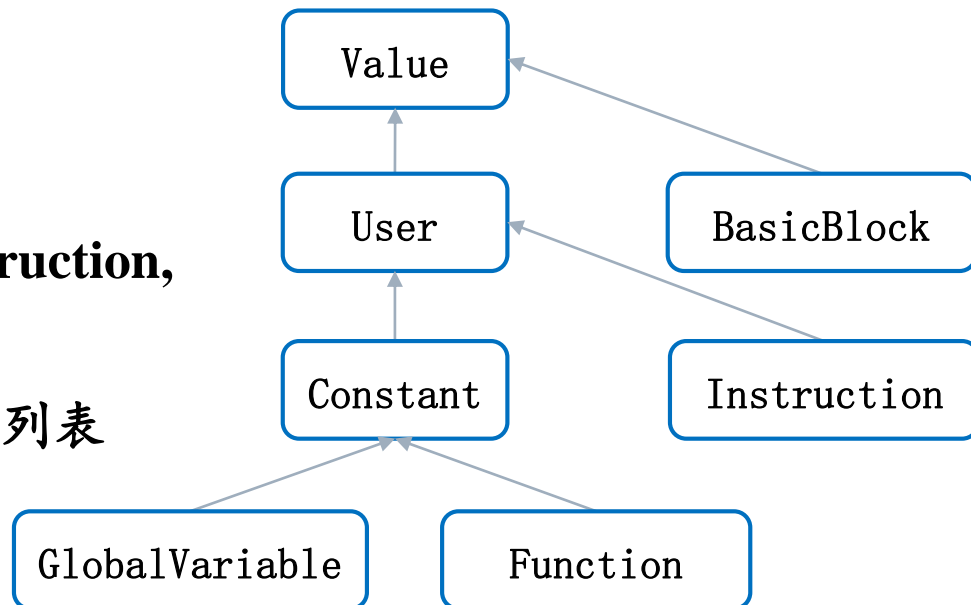
□ `Type* getType()`: 获取 Value 的类型

□ use 迭代器

- `Value::use_iterator/Value::const_use_iterator`

- `use_size(), use_empty(), use_begin(), use_end()`

□ use 替换: `void replaceAllUsesWith(Value *V)`



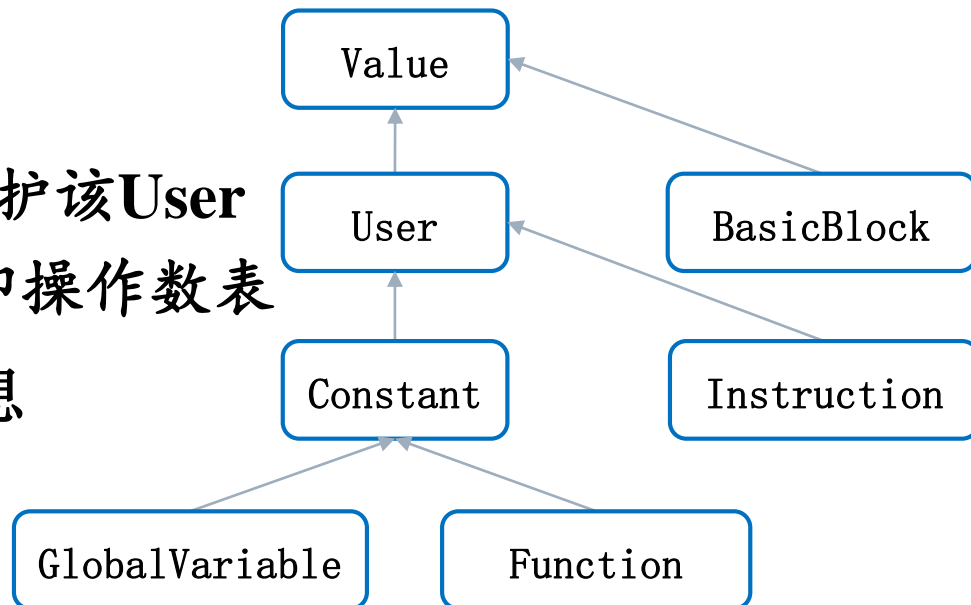


LLVM IR的UD和DU链

□ User类

- 代表一个**使用者**，维护该User使用的所有Value，即操作数表

□ 提供了use-def信息



■ 主要方法

□ 获取操作数

- Value *getOperand(unsigned i), 获取第i个操作数
- unsigned getNumOperands(), 获取操作数的个数

□ 操作数迭代器, User::op_iterator

- op_begin(), op_end()



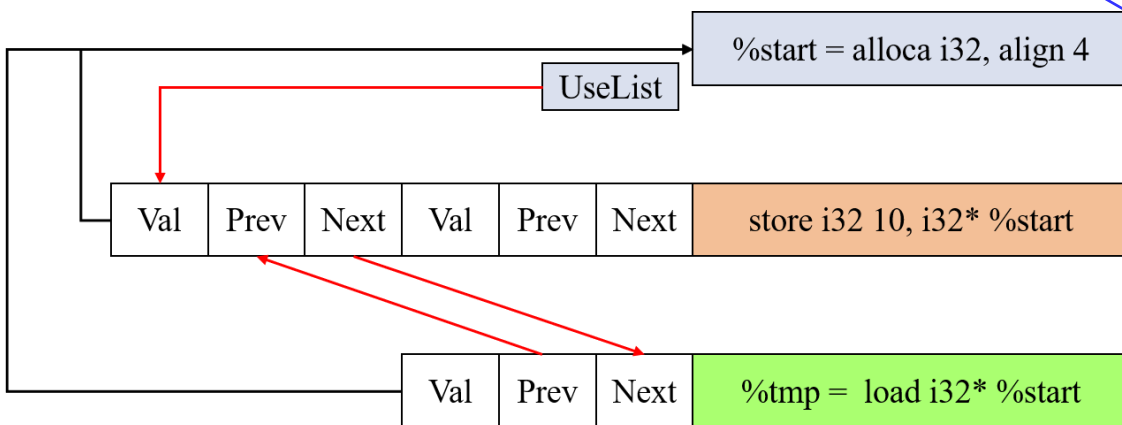
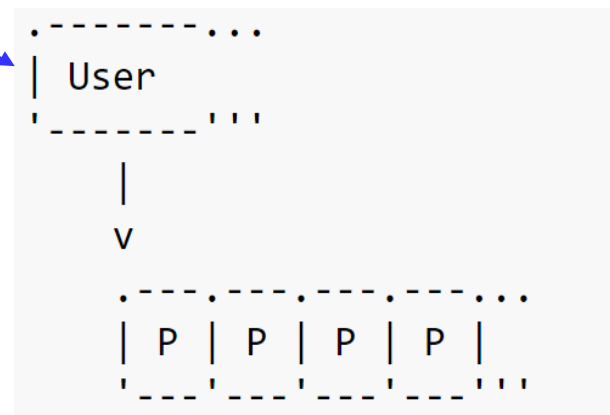
LLVM IR的UD和DU链

■ Value、Use、User组成的双向链保存DU、UD关系

- 每个Value对象有一个Use ** UseList双向链表，维护所有的使用者
- 每个Use对象代表一条从User→Value的边
- 每个User对象维护一个Value列表

■ 固定数目操作数的User

■ 可变数目操作数的User





LLVM IR的UD和DU链

■ Value、Use、User组成的双向链保存DU、UD关系

□ 每个Value对象有一个Use ** UseList双向链表，维护所有的使用者

□ 每个Use对象代表一条从User→

□ 每个User对象维护一个Value列

■ 固定数目操作数的User

■ 可变数目操作数的User

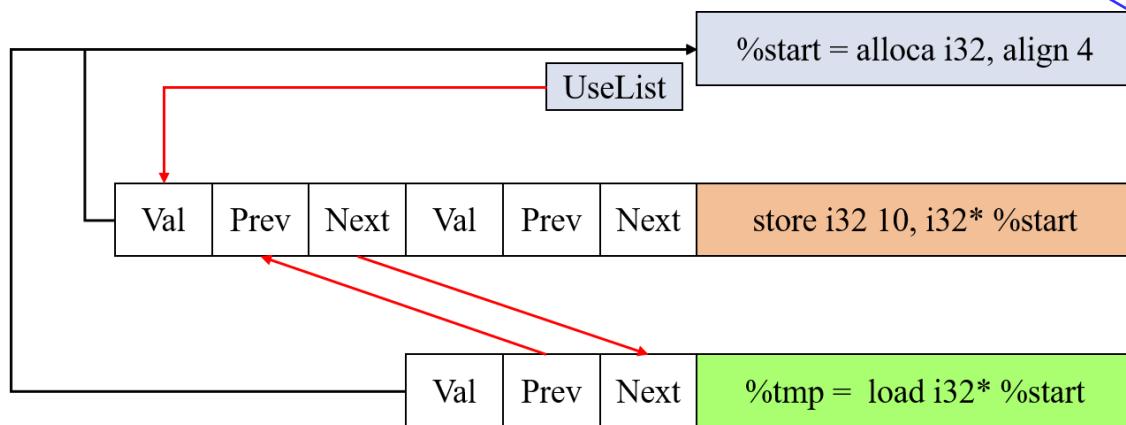
遍历使用者

```
Function *F = ...;  
for (User *U : F->users())  
  if (Instruction *Inst =  
      dyn_cast<Instruction>(U)) {...}
```



遍历操作数

```
Instruction *pi = ...;  
for (Use &U : pi->operands()) {  
  Value *v = U.get();  
  // ...  
}
```





2. 程序分析

- 控制流分析
- 数据流分析：到达-定值分析、数据流分析模式、活跃变量分析、可用表达式分析等



□ 可用表达式(available expressions)

$x = y + z$

...

...

...

p

没有对
y和z的
定值

$y + z$ 在 p 点
可用

$x = y + z$

...

$y = \dots$

...

p

$y + z$ 在 p 点
不可用

$x = y + z$

...

$z = \dots$

...

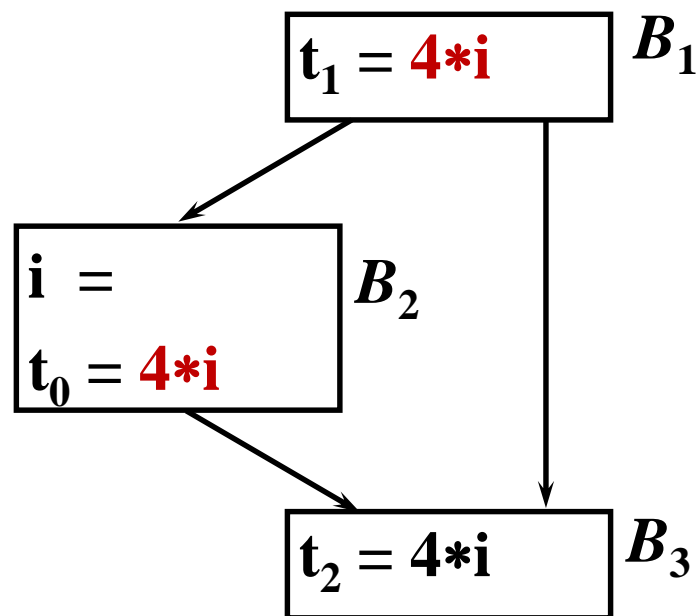
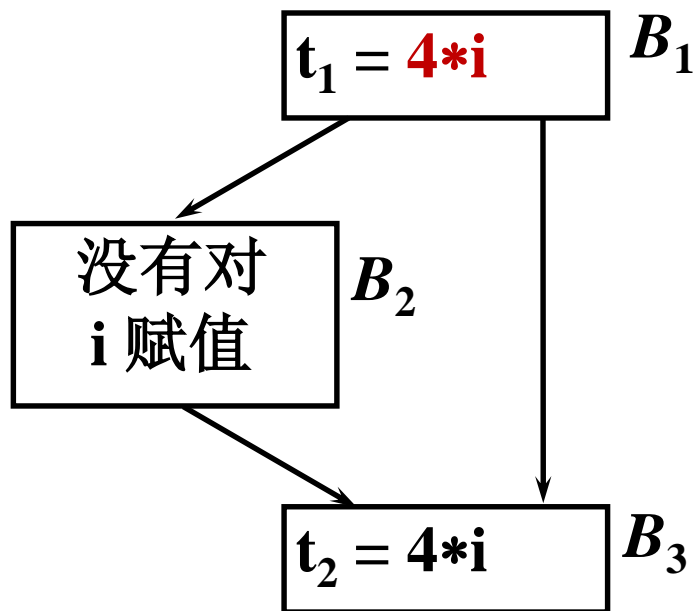
p

$y + z$ 在 p 点
不可用



可用表达式

下面两种情况下， $4*i$ 在 B_3 的入口都可用





可用表达式

□ 定义

- 若到点 p 的每条执行路径都计算 $x + y$ ，并且计算后没有对 x 或 y 赋值，那么称 $x + y$ 在点 p 可用
- e_gen_B : 块 B 产生的可用表达式集合
- e_kill_B : 块 B 注销的可用表达式集合
- $IN[B]$: 块 B 入口的可用表达式集合
- $OUT[B]$: 块 B 出口的可用表达式集合

□ 应用

- 公共子表达式删除



□ 数据流等式

- $OUT [B] = e_gen_B \cup (IN [B] - e_kill_B)$
- $IN [B] = \mathbf{Y}_{P \text{ 是 } B \text{ 的前驱}} OUT [P]$
- $IN [ENTRY] = \emptyset$

□ 同先前的主要区别

- 汇合算符：使用 \cap 而不是 \cup
- 求最大解而不是最小解



值编码 Value Numbering

□ Value Numbering

- Balke 1968 or Ershov 1954
- 早期面向低级的线性IR（如三地址码）
- 现在也存在面向树或图的等价方法

□ 优化的范围

- 基本块：Local value numbering
- 扩展的基本块：Superlocal value numbering (SVN)
- 支配关系：基于支配的值编码 (DVN)



Local Value Numbering

□ 每个值不同的表达式有一个值编号

Original Code

$a \leftarrow x + y$

$z \leftarrow y$

$d \leftarrow 17$

$c \leftarrow x + z$

With VNs

$a^3 \leftarrow x^1 + y^2$

$z^2 \leftarrow y^2$

$d^4 \leftarrow 17$

$c^3 \leftarrow x^1 + z^2$

Hash Table for VN

$\{ \langle x, 1 \rangle, \langle y, 2 \rangle, \langle \langle +, 1, 2 \rangle, 3 \rangle, \langle a, 3 \rangle \}$

$\{ \dots, \langle z, 2 \rangle \}$

$\{ \dots, \langle z, 2 \rangle, \langle 17, 4 \rangle, \langle d, 4 \rangle \}$

$\{ \dots, \langle z, 2 \rangle, \langle 17, 4 \rangle, \langle d, 4 \rangle, \langle c, 3 \rangle \}$

□ 冗余消除 (Redundancy Elimination)

Rewritten

$a \leftarrow x + y$

$z \leftarrow y$

$d \leftarrow 17$

$c \leftarrow a$

值编号作为key,
value为存放值的变
量或者常量

Hash Table for Rewritten

$\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle \}$

$\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle \}$

$\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle, \langle 4, 17 \rangle \}$

$\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle, \langle 4, 17 \rangle \}$



LVN: 解决重写问题

Original Code

$a \leftarrow x + y$
 $z \leftarrow y$
* $b \leftarrow x + y$
 $a \leftarrow 17$
* $c \leftarrow x + y$

Two redundancies
marked by *

With VNs

$a^3 \leftarrow x^1 + y^2$
 $z^2 \leftarrow y^2$
* $b^3 \leftarrow x^1 + y^2$
 $a^4 \leftarrow 17$
* $c^3 \leftarrow x^1 + y^2$

Rewritten

$a \leftarrow x + y$
 $z \leftarrow y$
* $b \leftarrow a$
 ~~$a \leftarrow 17$~~
* ~~$c \leftarrow a$~~

Hash Table for Rewritten

$\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle \}$
 $\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle \}$
 $\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle, \langle 4, 17 \rangle \}$
 $\{ \langle 1, x \rangle, \langle 2, y \rangle, \langle 3, a \rangle, \langle 4, 17 \rangle \}$

Optional Solutions:

- Use $c^3 \leftarrow b^3$
- Save a^3 in t^3
- Rename around it (best)

根据引用的值编号来查找值
并计算，记录当前定值所关
联的值编号到值的映射
对a重写时，并不去查询
value为a的pair



LVN: 值编码的重命名

□ 重命名：给每个值唯一的名字 => SSA

Original Code

$a_0 \leftarrow x_0 + y_0$
 $z_0 \leftarrow y_0$
* $b_0 \leftarrow x_0 + y_0$
 $a_1 \leftarrow 17$
* $c_0 \leftarrow x_0 + y_0$

With VNs

$a_0^3 \leftarrow x_0^1 + y_0^2$
 $z_0^2 \leftarrow y_0^2$
* $b_0^3 \leftarrow x_0^1 + y_0^2$
 $a_1^4 \leftarrow 17$
* $c_0^3 \leftarrow x_0^1 + y_0^2$

Rewritten

$a_0 \leftarrow x_0 + y_0$
 $z_0 \leftarrow y_0$
* $b_0 \leftarrow a_0$
 $a_1 \leftarrow 17$
* $c_0 \leftarrow a_0$

Hash Table for Rewritten

$\{ \langle 1, x_0 \rangle, \langle 2, y_0 \rangle, \langle 3, a_0 \rangle \}$
 $\{ \langle 1, x_0 \rangle, \langle 2, y_0 \rangle, \langle 3, a_0 \rangle \}$
 $\{ \langle 1, x_0 \rangle, \langle 2, y_0 \rangle, \langle 3, a_0 \rangle, \langle 4, 17 \rangle \}$
 $\{ \langle 1, x_0 \rangle, \langle 2, y_0 \rangle, \langle 3, a_0 \rangle, \langle 4, 17 \rangle \}$

Result:

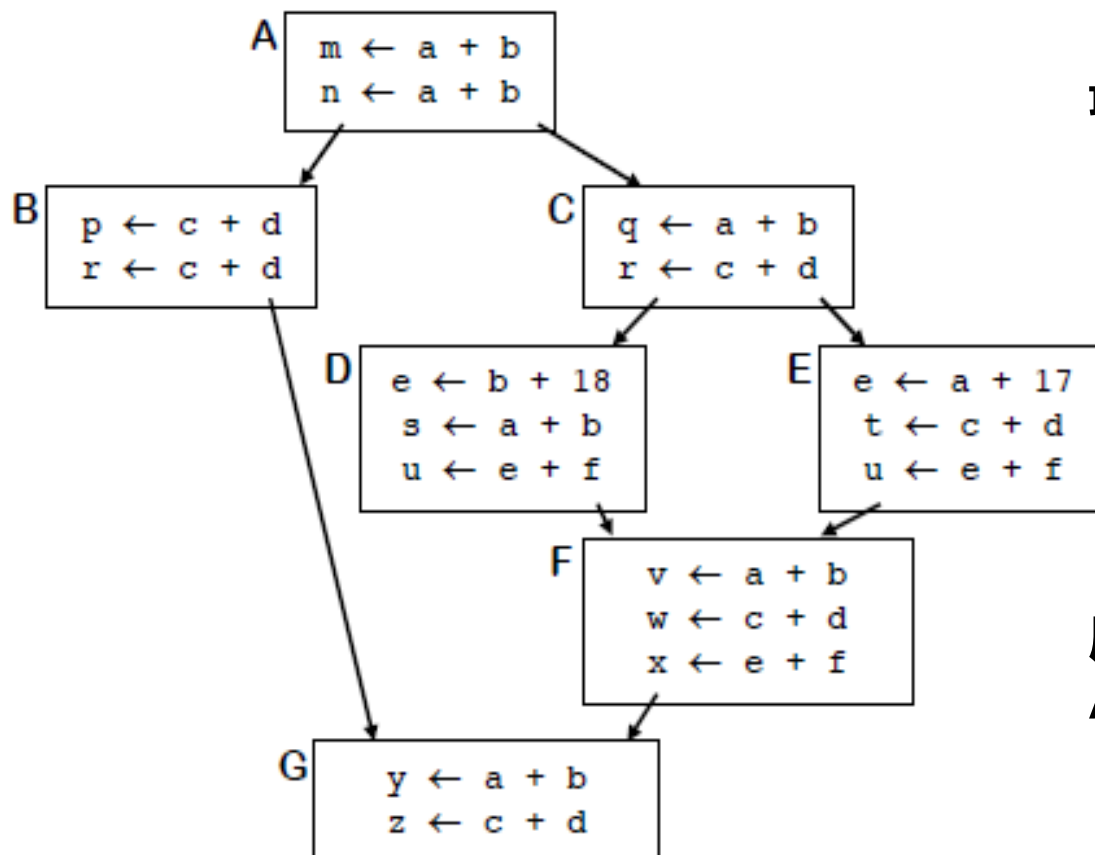
- a_0 is available
- Rewriting just works



Superlocal Value Numbering

□ 扩展的基本块EBB

- B_1, B_2, \dots, B_n , $2 \leq i \leq n$, B_i 有唯一的前驱



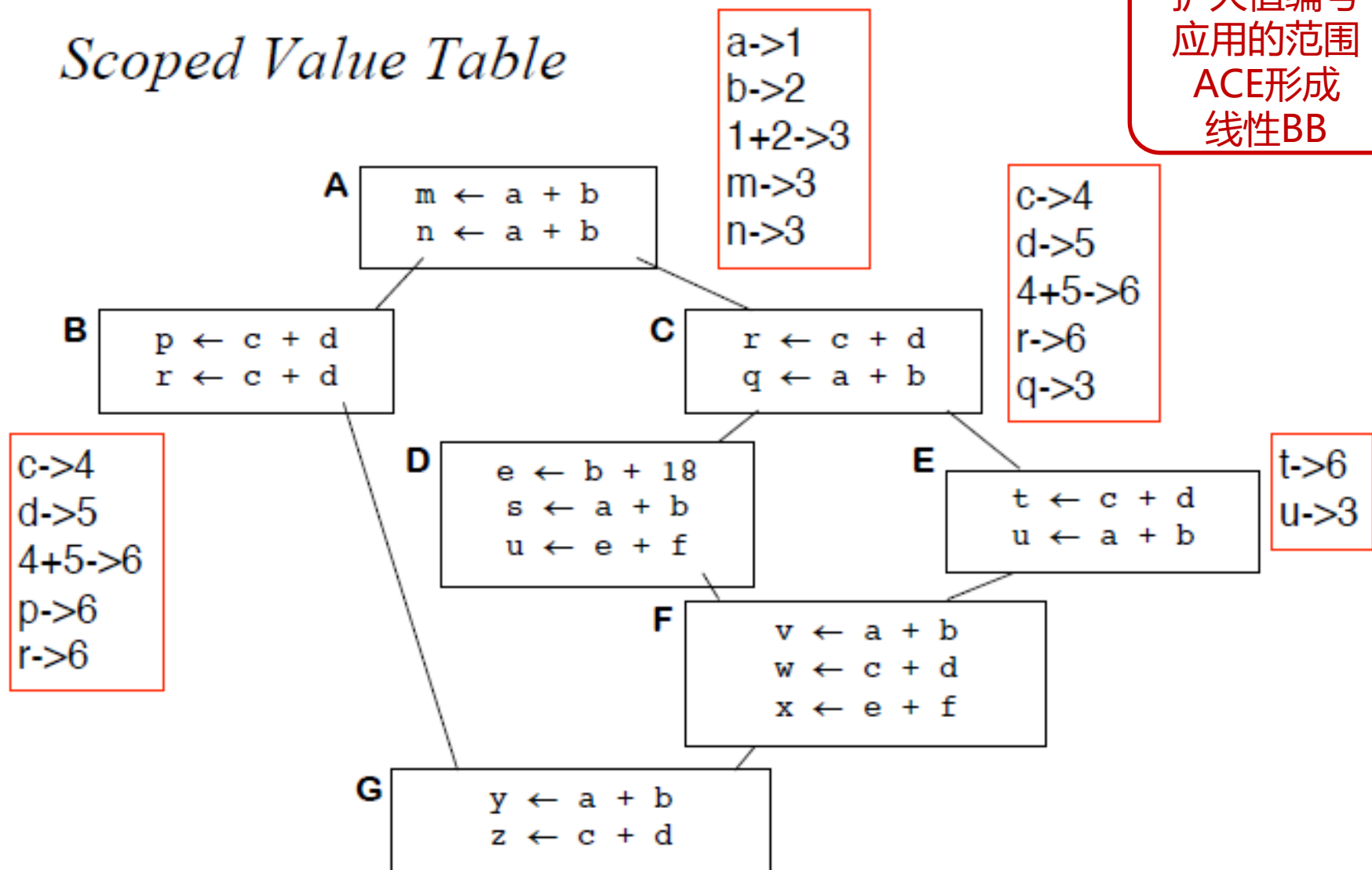
最大的EBB: ABCDE、F、G

应用局部化到EBB的路径:
AB、ACD、ACE、F、G



Superlocal Value Numbering

Scoped Value Table





数据流问题小结

见书上表9.2

□ 数据流问题的组成

- 数据流值的论域、数据流的方向、迁移函数、边界条件、汇合算符、数据流等式

| | 数据流值的定义域 | 方向 | 应用 |
|--------|---------------------------|----|----------------------------|
| 到达-定值 | 定值点集合 | 正向 | 用于DCE的def-use链 |
| 活跃变量 | 变量集合 | 反向 | 寄存器分配、未初始化检测、SSA构造、无用存储消除等 |
| 可用表达式 | 表达式集合 | 正向 | CSE |
| 非常量表达式 | 表达式集合 | 反向 | 循环外提 |
| 常量传播 | $\langle v, c \rangle$ 集合 | 正向 | 常量折叠 |



流敏感(flow-sensitivity)

□ 流不敏感分析 (flow-insensitive analysis)

- 不考虑程序中语句执行的顺序。
- 把程序中语句随意交换位置 (即: 改变控制流), 如果分析结果始终不变, 则该分析为流不敏感分析。

前面的数据流分析算法中(4)没有规定对基本块的操作次序

□ 流敏感分析 (flow-sensitive analysis)

考虑程序中过程内的控制流情况 (顺序、分支、循环)



上下文敏感、域敏感、路径敏感

- 上下文不敏感分析Context-insensitive analysis
 - 在过程调用时候忽略调用的上下文
- 上下文敏感分析Context-sensitive analysis
 - 考虑过程调用的上下文(调用点的**实参、返回值的**使用)
- 域(不)敏感分析field-sensitive analysis
 - 分析中是否考虑结构体中的不同域、数组中的不同下标元素
- 路径(不)敏感分析path-sensitive analysis
 - 是否依据分支语句的不同谓词来计算不同的分析信息



中国科学技术大学
University of Science and Technology of China

3. 符号执行



□ 能分析程序中所有可能的运行

- 有许多有趣的想法和工具
- 很多只是在论文上展示有好的效果
- 学术界推出的被企业认可的商用工具寥寥无几

Dawson Engler: [coverity](#) [CACM2010]



但是，开发者使用起来....

- 不容易，论文中的结果描述的是静态分析专家所用的
- 有生命力的商用工具：要能处理误报(false positives)、开发者的困惑、错误管理、...

[CACM2010] A few billion lines of code later: using static analysis to find bugs in the real world, Communication of the ACM, 53(2):66-75, 2010.

张昱：《编译原理和技术》代码优化



符号执行的引入

□ 抽象的作用

- 让静态分析对所有可能的运行进行建模，但引入保守性
- *-敏感性方法试图改进之，但是远远不够

□ 静态分析的抽象 \neq 开发者的抽象

□ 测试

- 每个测试只能考察一种可能的执行
- 希望测试用例更具有有一般性，但是却没有保障

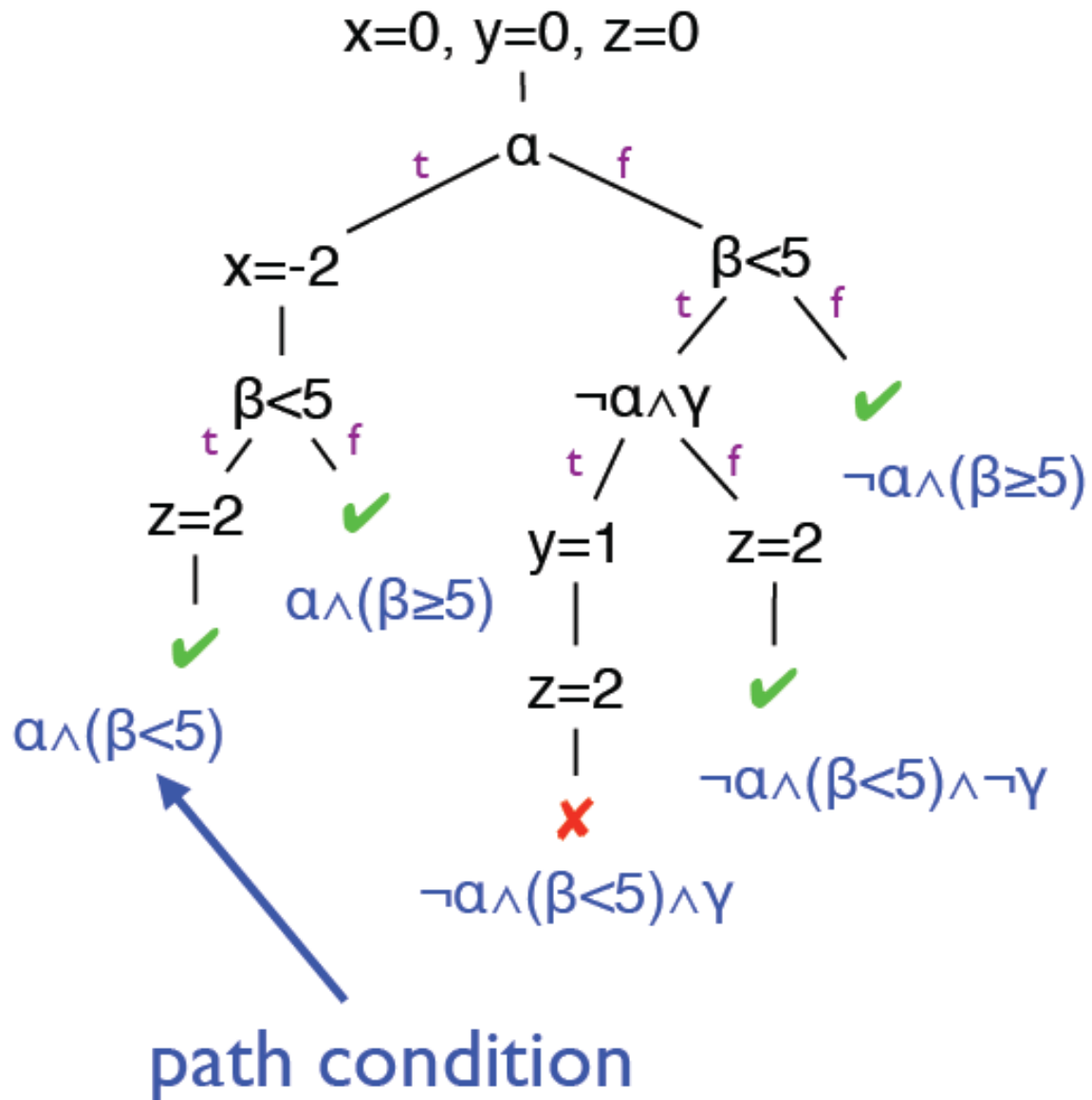
□ 符号执行：对测试的泛化，引入符号值来计算

[计算机学报2015] 软件安全漏洞检测技术, 38(4):717-732, 2015.



符号执行举例

```
1. int a =  $\alpha$ , b =  $\beta$ , c =  $\gamma$ ;  
2. // symbolic  
3. int x = 0, y = 0, z = 0;  
4. if (a) {  
5.     x = -2;  
6. }  
7. if (b < 5) {  
8.     if (!a && c) { y = 1; }  
9.     z = 2;  
10. }  
11. assert(x+y+z!=3)
```





符号执行

□ 每个符号执行路径代表多个实际的程序运行

- 实际运行的具体值满足符号执行的路径条件

相比测试，符号执行可以覆盖更多的程序执行空间

□ 符号执行的发展

- 早期的工作：James C. King. [Symbolic execution and program testing](#). CACM, 19(7):385–394, 1976. (高引)
- 问题：1980's 当时的机器内存小且慢
- 符号执行代价极高：大量可能的执行路径，需要求解器判断哪些路径是可行的，程序状态有许多位



符号执行现状

- 计算机：速度快、内存便宜
- 有非常强大的SMT/SAT求解器
 - SMT= Satisfiability Modulo Theories =SAT++
 - 快速解决非常多的问题：检查断言、削减可行路径
 - **求解器**：[Z3](#)(已集成到LLVM中)、[STP](#)、[Yices](#)等
- 近10年来的有代表性论文
 - [KLEE](#) (OSDI2008, [C. Cadar](#)、[D. Engler](#)等,用于LLVM)
 - [[CACM2013](#)] Symbolic execution for software testing(30y+)
 - [[CSUR2018](#)] A Survey of Symbolic Execution Techniques



□ SAT求解器是SMT求解器的核心

- SMT(可满足性模理论): 接受不同格式的等式系统
- SAT: 必须是合取CNF范式的布尔等式
- 理论上, 所有SMT查询可以约减到SAT查询
- 实践上, SMT和高级优化是关键
- 简单的等式 $x+0=x$
- 数组理论: $\text{read}(x, \text{write}(42, x, A)) = 42$
 - x 是下标, A 是数组, $\text{write}(42, x, A)$ 表示 $A[x]=42$ 后的数组
- 缓存: 记住求解器的查询; 删除无用变量; ...



- 两个关键的系统
 - DART (Godefroid and Sen, PLDI 2005)
 - EXE (Cadar, Ganesh, Pawlowski, Dill, and Engler, CCS2006)
- SAGE: Microsoft的concolic执行器(动态符号执行)
- Mayhem(CMU), Angr(UCSB), Triton
- Java Symbolic PathFinder
- KLEE



简单命令式语言IMP的符号执行

```
a ::= n | X | -a | a0+a1 | a0-a1 | a0×a1 | a0/a1  
b ::= bv | ¬b | b0^b1 | b0∨b1 | a0=a1 | a0<a1 | a0>a1  
c ::= skip | input(s) | X:=a | if b then c else c  
    | c0; c1 | while b do c | assert b
```

- n 是整数、 X 是变量、 bv 是布尔值
- c 是命令、 a 和 b 分别是整型和布尔型表达式
- Sym-while的OCaml实现: <https://github.com/Isweet/sym-while>
 - [syswhile.ml](#) 包含总控程序, 即`main`, 具体执行或符号执行
 - [lexer.mll](#)(词法描述)、[parser.mly](#)(语法描述)、[ast.ml](#)(AST)
 - [concrete.ml](#): 解释执行, 状态是变量到整数的映射, [Imp.run s](#)
 - [symbol.ml](#), [symbolic.ml](#): 符号执行, 状态时变量到符号表达式和路径条件的映射



□ 符号表达式可能包含变量

[ast.ml](#)

```
type arith =  
  | AVar of string  
  | AEnum of int  
  | AENegate of arith  
  | AEPlus of arith * arith  
  | AEMinus of arith * arith  
  | AEMult of arith * arith  
  | AEDiv of arith * arith
```

```
type boolean =  
  | BTrue  
  | BFalse  
  | BNot of boolean  
  | BAnd of boolean * boolean  
  | BOr of boolean * boolean  
  | BELT of arith * arith  
  | BEGT of arith * arith  
  | BEEq of arith * arith
```



符号状态

□ 具体状态：变量到整数

```
type conc_state = int StringMap.t
```

[concrete.ml](#)

□ 符号状态：变量到符号表达式和路径条件 [symbolic.ml](#)

```
type sym_state = (Symbol.int_t StringMap.t) * Symbol.t
```

■ **Symbol.t**:布尔型符号表达式类型

■ **Symbol.int_t**:整型符号表达式类型



基于fork的符号执行

□ 如何判断哪个分支可行(feasible)?

```
let rec eval (s : stmt) (s_st : sym_state) : answer =  
  let (env, pc) = s_st in
```

```
| Sif (b, s1, s2) ->  
  let l = t_of_boolean b env in          (* branch cond *)  
  let cond_true = LAnd (l, pc) in        (* ... and path cond *)  
  let cond_false = LAnd ((LNot l), pc) in  
  let sat_true = check (z3_of_t cond_true) in  
  let sat_false = check (z3_of_t cond_false) in  
  (match sat_true with                    (* might do both branches *)  
   | Some _ -> ...(eval s1 (env, cond_true))  
   | None -> ...);  
  (match sat_false with  
   | Some _ -> ...(eval s2 (env, cond_false))  
   | None -> ...);
```



符号执行策略

□ 顶层策略

- 初始化状态：pc=0, 路径条件为true, 状态为empty
- 对每条语句符号求值
- 一旦执行分叉，则对两个分支都分别求值（DepthFS）
- 执行完后，返回多个符号状态

□ 路径爆炸(path explosion): 分支、循环

□ 搜索策略

DFS（容易在某部分stuck）、BFS；树→图

优先权搜索、随机搜索、覆盖引导的启发式搜索、分代...

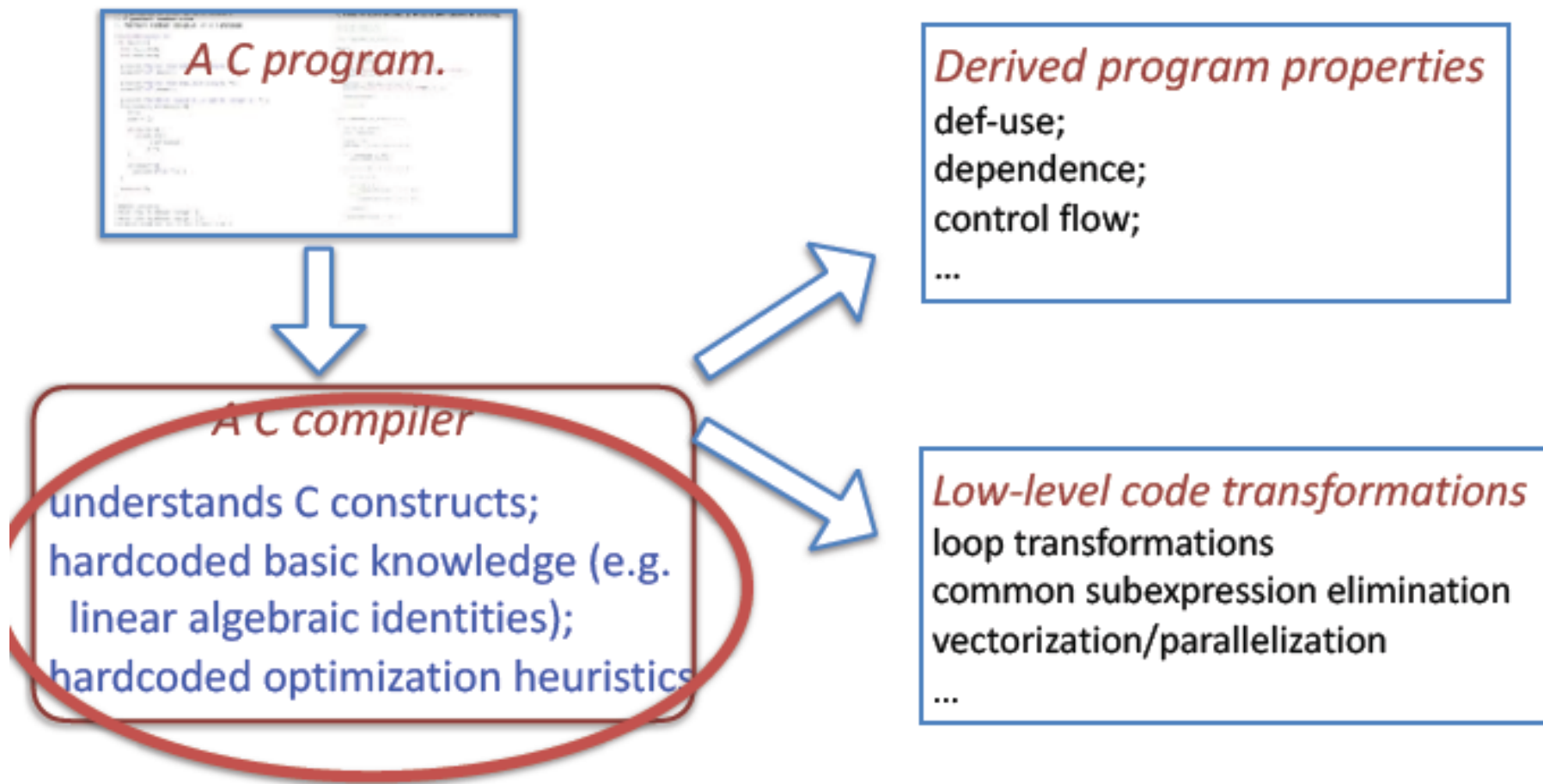


中国科学技术大学
University of Science and Technology of China

4. 面向AI的优化



传统的优化编译器





高级语义驱动的编译

□ 高级语义

- DSL: 如矩阵运算 $b = a + c$

- 泛化的编译技术

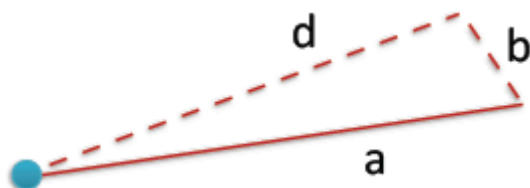
 - 强度削弱

 - 循环冗余消除

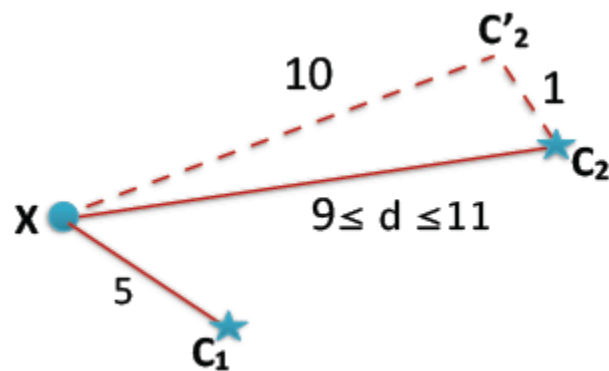


■ 强度削弱

传统: $b/2 \rightarrow b \gg 1$



$$a - b \leq d \leq a + b$$





```
for (i=0; i<N; i++) {  
    a[i] = b/c;  
}
```

Elusive to existing techniques.

```
w = w0;  
while (d > 0.01) {  
    d = 0;  
    for (i = 0; i < M; i++) {  
        d += a[i] + b[i] * w;  
    }  
    w = w - 0.001 * d;  
}
```

reduction over loop *i*

Stay the same across *while* loop, but vary across *for* loop.

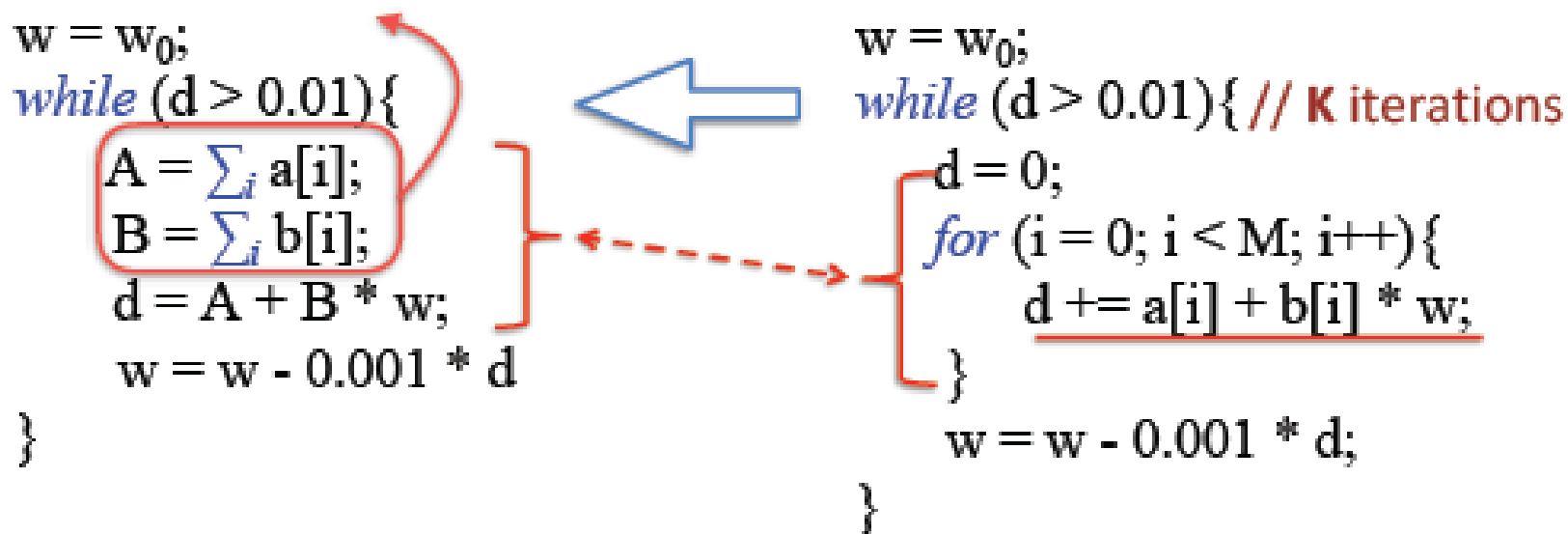
Stay the same across *for* loop, but change across *while* loop



循环冗余

- An equivalent form of the example.

$$O(M * K) \longrightarrow O(M + K)$$





中国科学技术大学
University of Science and Technology of China

5. 数据流分析的基础



数据流分析框架

□ 数据流分析框架 (D, V, \wedge, F) 包括

- 数据流分析的 **方向** D ，它可以是正向或逆向
- **数据流值的论域**：半格 V 、**汇合算子** \wedge
- V 到 V 的 **迁移函数族** F ，包括适用于边界条件（ENTRY 和 EXIT 结点）的常函数

□ 半格 (V, \wedge)

- 是一个集合 V 和一个二元交运算（汇合运算） \wedge ，满足：
- 幂等性：对所有的 x ， $x \wedge x = x$
- 交换性：对所有的 x 和 y ， $x \wedge y = y \wedge x$
- 结合性：对所有的 x, y 和 z ， $x \wedge (y \wedge z) = (x \wedge y) \wedge z$



半格(semilattices)

□ 半格有顶元 \top (最大元素), 可以还有底元 \perp (最小元素)

■ 对 V 中的所有 x , $\top \wedge x = x$

■ 对 V 中的所有 x , $\perp \wedge x = \perp$

□ 偏序关系: 集合 V 上的关系 \leq

■ 自反性: 对所有的 x , $x \leq x$

■ 反对称性: 对所有的 x 和 y , 如果 $x \leq y$ 且 $y \leq x$, 那么 $x = y$

■ 传递性: 对所有的 x, y 和 z , 如果 $x \leq y$ 且 $y \leq z$, 那么 $x \leq z$

此外, 关系 $<$ 的定义

$x < y$ 当且仅当 $(x \leq y)$ 并且 $(x \neq y)$



□ 半格和偏序关系之间的联系

- 半格 (V, \wedge) 的汇合运算 \wedge 确定了半格值集 V 上一种偏序 \leq ：
对 V 中所有的 x 和 y ， $x \leq y$ 当且仅当 $x \wedge y = x$
- 若 $x \wedge y$ 等于 g ，则 g 就是 x 和 y 的最大下界

例 半格的论域 V 是先前全域 U 的幂集

- 汇合运算为集合并： \emptyset 是顶元， U 是底元，偏序关系是 \supseteq
- 汇合运算为集合交： U 是顶元， \emptyset 是底元，偏序关系是 \subseteq
- 按偏序 \leq 意义上的最大解是最精确的
 - (1) 到达-定值：最精确的解含最少定值
 - (2) 可用表达式：最精确的解含最多表达式

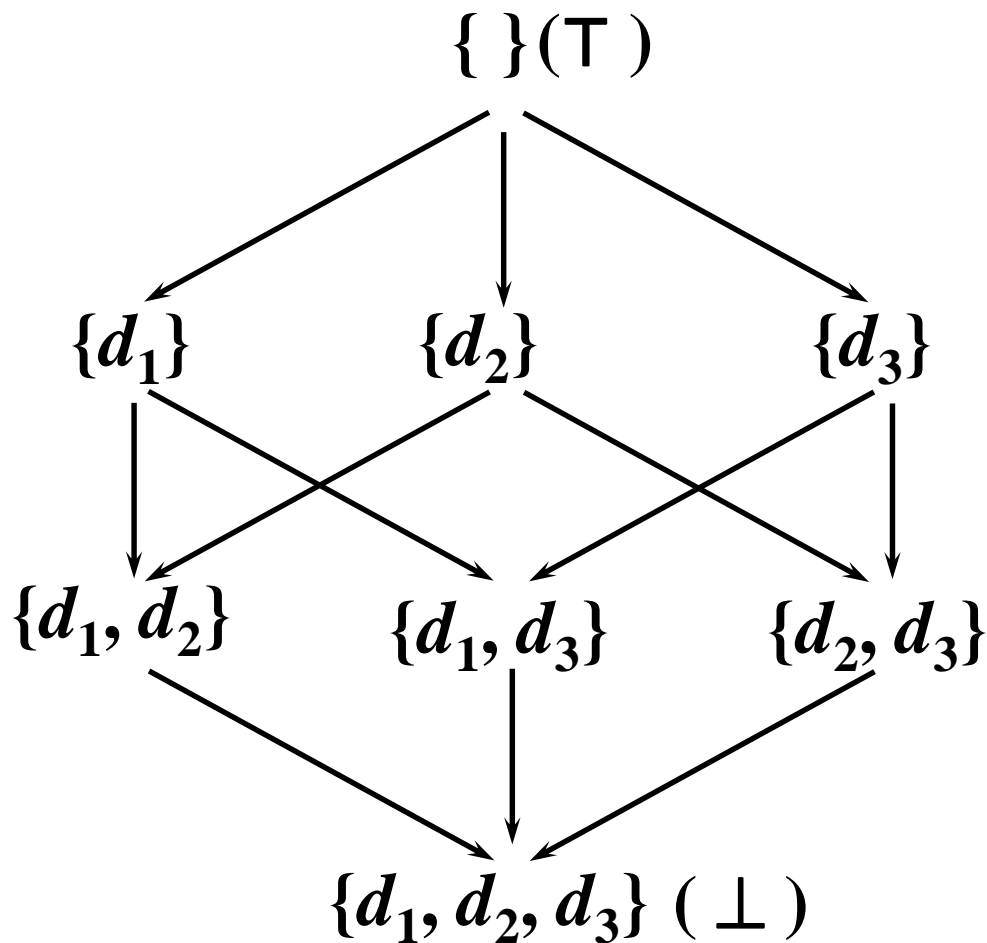


□ 格图

- 结点是 V 中的元素
- 如果 $y \leq x$, 则有从 x 朝下到 y 的有向边

右图是定值子集之间形成的格:

- 到达-定值的 \leq 是 \supseteq
- $x \wedge y$ 的最大下界是
 $x \cup y$





如何降低格图的规模?

□ 到达-定值格图存在的问题

- 数据流值的集合是定值集合的幂集

=>格图结点数随变量数呈指数级增长

- 每个变量的定值可达性独立于其他变量的定值可达性

定值半格表示为从每个变量的简单定值半格构造出的积半格

□ 积半格(假定 (A, \wedge_A) 和 (B, \wedge_B) 是半格)

- 论域是 $A \times B$

- 汇合运算 \wedge : $(a, b) \wedge (a', b') = (a \wedge_A a', b \wedge_B b')$



□ 半格的高度

- 偏序集合 (V, \leq) 中的一个上升链是序列 $x_1 < x_2 < \dots < x_n$

- 半格的高度就是其中最长上升链中 $<$ 的个数

例，在一个有 n 个定值的程序中，到达-定值的高度是 n

□ 数据流分析算法的收敛

- 半格的高度有限 \Rightarrow 数据流分析迭代算法收敛

- 半格的值论域有限 \Rightarrow 半格的高度有限

- 半格的值论域无限 \Rightarrow 半格的高度可能有限

如，常量传播算法中使用的半格



□ 迁移函数族 $F : V \rightarrow V$ 有下列性质

- F 包括恒等函数 I ，即对 V 中所有的 x ，有 $I(x) = x$
- F 封闭于复合，即对 F 中任意两个函数 f 和 g ， $g \circ f \in F$
- 若 F 中所有函数 f 都有单调性，即

$x \leq y$ 蕴涵 $f(x) \leq f(y)$ ，或 $f(x \wedge y) \leq f(x) \wedge f(y)$

则称框架 (D, V, \wedge, F) 是单调的

- 框架 (D, V, \wedge, F) 的分配性

对 F 中所有的 f ， $f(x \wedge y) = f(x) \wedge f(y)$

框架单调 \Rightarrow 所求得的解是数据流方程组的最大不动点



□ 例 到达-定值分析

若 $f_1(x) = G_1 \cup (x - K_1)$, $f_2(x) = G_2 \cup (x - K_2)$

■ 若 G 和 K 是空集, 则 f 是恒等函数

■
$$f_2(f_1(x)) = G_2 \cup ((G_1 \cup (x - K_1)) - K_2)$$
$$= (G_2 \cup (G_1 - K_2)) \cup (x - (K_1 \cup K_2))$$

因此 f_1 和 f_2 的复合 f 为 $f = G \cup (x - K)$ 的形式

■ 分配性可以由检查下面的条件得到

$$G \cup ((y \cup z) - K) = (G \cup (y - K)) \cup (G \cup (z - K))$$

分配性: $f(y \wedge z) = f(y) \wedge f(z)$



一般框架的迭代算法

□ 以正向数据流分析为例

- (1) $\text{OUT}[\text{ENTRY}] = v_{\text{ENTRY}};$
- (2) for (除了ENTRY以外的每个块 B) $\text{OUT}[B] = \text{T};$
- (3) while (任何一个OUT出现变化)
- (4) for (除了ENTRY以外的每个块 B) {
- (5) $\text{IN}[B] = \bigwedge_{P \text{ 是 } B \text{ 的前驱}} \text{OUT}[P];$
- (6) $\text{OUT}[B] = f_B(\text{IN}[B]);$
- (7) }



数据流解的含义

- 结论：算法所得解是理想解的稳妥近似
- 理想解所考虑的路径
 - 执行路径集：流图上每一条路径都属于该集合
若流图有环，则执行路径数是无限的
 - 程序可能的执行路径集：程序执行所走的路径属于该集合——这是理想解所考虑的路径集
 - 可能的执行路径集 \subseteq 执行路径集
 - 寻找所有可能执行路径是不可判定的
- 以下讨论以正向数据流分析为例



□ 理想解

若路径 $P = \text{ENTRY} \rightarrow B_1 \rightarrow B_2 \rightarrow \dots \rightarrow B_k$, 定义

- $f_P = f_{k-1} \circ \dots \circ f_2 \circ f_1$

- $\text{IDEAL}[B] = \bigwedge_{P \text{ 是从 ENTRY 到 } B \text{ 的一条可能路径}} f_P(v_{\text{ENTRY}})$

□ 有关理解解的结论

- 任何大于理想解IDEAL的回答一定是不对的
- 任何小于或等于IDEAL的值是稳妥的
- 在稳妥的值中, 越接近IDEAL的值越精确



□ MFP最大不动点解

maximal fixed point

- 访问每个基本块（不一定按照程序执行时的次序）
- 在每个汇合点，把汇合运算作用到当前得到的数据流值，所用的一些初值是人工引入的

□ MOP执行路径上的解

meet over paths

- $\text{MOP}[B] = \bigwedge_{P \text{ 是从 ENTRY 到 } B \text{ 的一条路径}} f_P(v_{\text{ENTRY}})$
- MOP解汇集了所有可能路径的数据流值，包括那些不可能被执行路径的数据流值
- 对所有的块 B ， $\text{MOP}[B] \leq \text{IDEAL}[B]$



MFP与MOP的联系

□ MFP与MOP的联系

- MFP访问基本块未必遵循执行次序

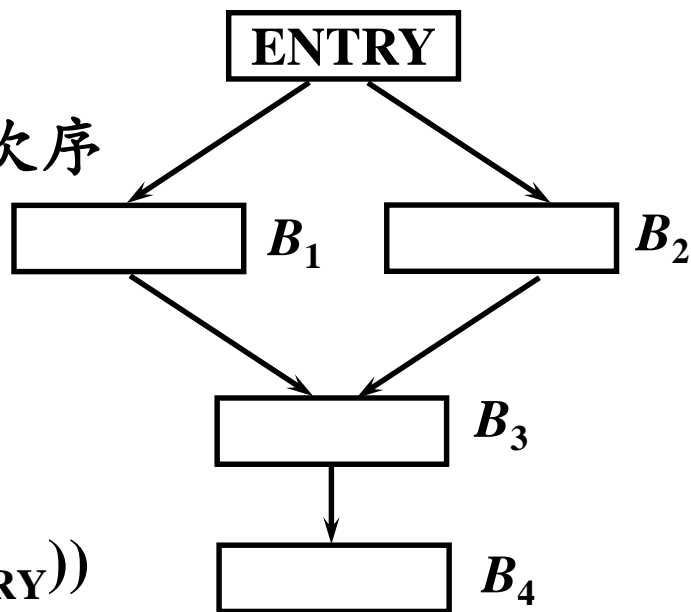
由各块的初值和迁移函数的
单调性保证结果一致

- MFP较早地使用汇合运算

$$\text{IN}[B_4] = f_3(f_1(v_{\text{ENTRY}}) \wedge f_2(v_{\text{ENTRY}}))$$

$$\text{而MOP}[B_4] = (f_3 \circ f_1)(v_{\text{ENTRY}}) \wedge (f_3 \circ f_2)(v_{\text{ENTRY}})$$

在数据流分析框架具有分配性时, 二者的结果是一样的



□ $\text{MFP} \preceq \text{MOP} \preceq \text{IDEAL}$



例题 1

一个C语言程序如下，右边是优化后的目标代码

| | | |
|-------------------------|-----------------------|-----------------------|
| main() | pushl %ebp | |
| { | movl %esp,%ebp | |
| int i,j,k; | movl \$1,%eax | -- j=1 |
| i=5; | movl \$6,%edx | -- k=6 |
| j=1; | L4: | |
| while(j<100){ | addl %edx,%eax | -- j=j+6 |
| k=i+1; | cmpl \$99,%eax | |
| j=j+k; | jle .L4 | -- while(j≤99) |
| } | | |
| } | | |

完成了哪些优化？



例题 1

一个C语言程序如下，右边是优化后的目标代码

| | | |
|-------------------------|-----------------------|-----------------------|
| main() | pushl %ebp | |
| { | movl %esp,%ebp | |
| int i,j,k; | movl \$1,%eax | -- j=1 |
| i=5; | movl \$6,%edx | -- k=6 |
| j=1; | L4: | |
| while(j<100){ | addl %edx,%eax | -- j=j+6 |
| k=i+1; | cmpl \$99,%eax | |
| j=j+k; | jle .L4 | -- while(j≤99) |
| } | | |
| } | | |

复写传播、常量合并、代码外提、删除无用赋值
对**i**，**j**和**k**分配内存单元也成为多余，从而被取消



例题 2

一个C语言程序

main()

{

 long i,j;

 while (i) {

 if (j) { i = j; }

 }

} 生成的汇编码见右边

为什么会有连续跳转？

pushl %ebp

movl %esp,%ebp

subl \$8,%esp

.L2:

 cmpl \$0,-4(%ebp)

 jne .L4

 jmp .L3

.L4:

 cmpl \$0,-8(%ebp)

je .L5

 movl -8(%ebp),%eax

 movl %eax,-4(%ebp)

.L5:

jmp .L2

.L3:

while E1 do S1

L2: E1的代码

 真转 L4

 无条件转 L3

L4: S1的代码

 JMP L2

L3:

if E2 then S2

 E2的代码

 假转 L5

 S2的代码

L5:



例题 2

一个C语言程序

main()

{

long i,j;

while (i) {

if (j) { i = j; }

}

} 生成的汇编码见右边

为什么会有连续跳转?

```
pushl %ebp
movl %esp,%ebp
subl $8,%esp
```

.L2:

```
cmpl $0,-4(%ebp)
jne .L4
jmp .L3
```

.L4:

```
cmpl $0,-8(%ebp)
je .L5
movl -8(%ebp),%eax
movl %eax,-4(%ebp)
```

.L5:

```
jmp .L2
```

.L3:

嵌套时代码结构变成

L2: E1的代码
真转 L4
无条件转 L3

L4: S1的代码
E2的代码
假转 L5
S2的代码

L5: JMP L2

L3:

if E2 then S2

E2的代码
假转 L5
S2的代码

L5:



例题 2

一个C语言程序

```
main()
```

```
{
```

```
    long i,j;
```

```
    while (i) {
```

```
        if (j) { i = j; }
```

```
    }
```

```
}
```

优化编译的汇编码见右边

```
    pushl %ebp
```

```
    movl %esp,%ebp
```

```
.L7:
```

```
    testl %eax,%eax
```

```
    je .L3
```

```
    testl %edx,%edx
```

```
    je .L7
```

```
    movl %edx,%eax
```

```
    jmp .L7
```

```
.L3:
```




例题 3 尾递归

求最大公约数的函数

```
long gcd(p,q)
```

```
long p,q;
```

```
{
```

```
    if (p%q == 0)
```

```
        return q;
```

```
    else
```

```
        return gcd(q, p%q);
```

```
}
```

- 其中的递归调用称为尾递归
- 对于尾递归，编译器应怎样产生代码，使得这种递归调用所需的时空开销大大减少？
- 计算实在参数 q 和 $p\%q$ ，存放在不同的寄存器中
- 将上述寄存器中实在参数的值存入当前活动记录中形式参数 p 和 q 的存储单元
- 转到本函数第一条语句的起始地址继续执行



例题 3 尾递归

求最大公约数的函数

long gcd(p,q)

long p,q;

{

if (p%q == 0)

return q;

else

return gcd(q, p%q);

}

movl 8(%ebp),%esi p

movl 12(%ebp),%ebx q

.L4:

movl %esi,%eax

cld

扩展为64位

idivl %ebx

movl %edx,%ecx

p%q

testl %ecx,%ecx

p%q

je .L2

movl %ebx,%esi

q⇒p

movl %ecx,%ebx

p%q⇒q

jmp .L4

.L2:



例题 4

Program \rightarrow **Stmt**

Stmt \rightarrow **id** **:=** **Exp** | **read** (**id**) | **write** (**Exp**) |
Stmt ; **Stmt** |
while (**Exp**) **do begin Stmt end** |
if (**Exp**) **then begin Stmt end**
else begin Stmt end

Exp \rightarrow **id** | **lit** | **Exp OP Exp**

定义**Stmt**的两个属性

- *MayDef*表示它可能定值的变量集合
- *MayUse*表示它可能引用的变量集合
- 写一个语法制导定义或翻译方案，它计算**Stmt**的上述*MayDef*和*MayUse*属性



例题 4

Stmt \rightarrow **id** := **Exp**

{ Stmt.MayDef = {id.name} ;

Stmt.MayUse = Exp.MayUse }

Stmt \rightarrow **read (id)**

{ Stmt.MayUse = \emptyset ; Stmt.MayDef = {id.name} }

Stmt \rightarrow **write (Exp)**

{ Stmt.MayDef = \emptyset ; Stmt.MayUse = Exp.MayUse }

Stmt \rightarrow **Stmt₁ ; Stmt₂**

{ Stmt.MayUse = Stmt₁.MayUse \cup Stmt₂.MayUse ;

Stmt.MayDef = Stmt₁.MayDef \cup Stmt₂.MayDef }



Stmt \rightarrow **if** (**Exp**) **then** **begin** Stmt₁ **end**
else **begin** Stmt₂ **end**
{ Stmt.*MayUse* = Stmt₁.*MayUse* \cup
Stmt₂.*MayUse* \cup Exp.*MayUse*;
Stmt.*MayDef* = Stmt₁.*MayDef* \cup Stmt₂.*MayDef* **}**

Stmt \rightarrow **while** (**Exp**) **do** **begin** Stmt₁ **end**
{ Stmt.*MayUse* = Stmt₁.*MayUse* \cup Exp.*MayUse*;
Stmt.*MayDef* = Stmt₁.*MayDef* **}**

Exp \rightarrow **id** **{** Exp.*MayUse* = {id.name} **}**

Exp \rightarrow **lit** **{** Exp.*MayUse* = \emptyset **}**

Exp \rightarrow Exp₁ OP Exp₂
{ Exp.*MayUse* = Exp₁.*MayUse* \cup Exp₂.*MayUse* **}**



例题 4

基于 $MayDef$ 和 $MayUse$ 属性, 说明 $Stmt_1; Stmt_2$ 和 $Stmt_2; Stmt_1$ 在什么情况下有同样的语义

$Stmt_1.MayDef \cap Stmt_2.MayUse = \emptyset$ and
 $Stmt_2.MayDef \cap Stmt_1.MayUse = \emptyset$ and
 $Stmt_1.MayDef \cap Stmt_2.MayDef = \emptyset$