

# *CS293S Iterative Data-Flow Analysis*

Yufei Ding

# *Review: Computing Available Expressions*

## The Big Picture

1. Build a control-flow graph
2. Gather the initial data:  $DEEXPR(b)$  &  $EXPRKILL(b)$
3. Propagate information around the graph, evaluating the equation

$$\underline{AVAIL(b)} = \bigcap_{x \in \text{pred}(b)} \underline{(DEEXPR(x) \cup (AVAIL(x) \cap \overline{EXPRKILL(x)}))}$$

Entry point of block b

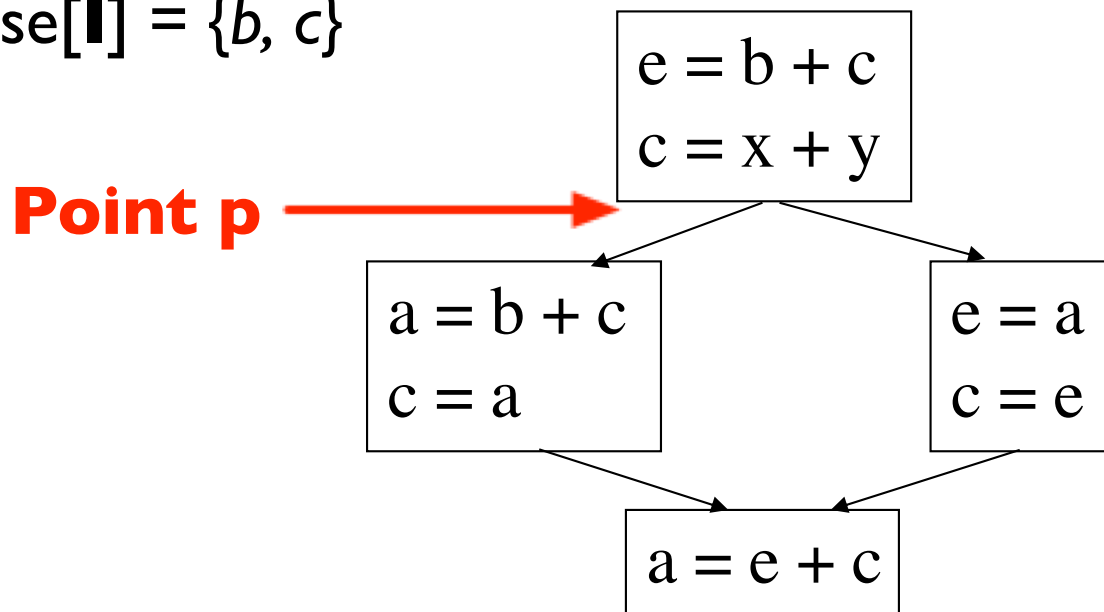
Exit point of block x

Works for loops through an iterative algorithm: finding the fixed-point.

All data-flow problems are solved, essentially, this way.

# Live Variables

- A variable  $v$  is live at a point  $p$  if there is a path from  $p$  to a **use of  $v$** , and that path does **not** contain a **redefinition of  $v$**
- Example:  $\blacksquare: a \leftarrow b + c$ 
  - A statement/instruction  $\blacksquare$  is a definition of a variable  $v$  if it may write to  $v$ .  $\text{def}[\blacksquare] = a$
  - A statement is a use of variable  $v$  if it may read from  $v$ .  $\text{use}[\blacksquare] = \{b, c\}$



# *Live Variables*

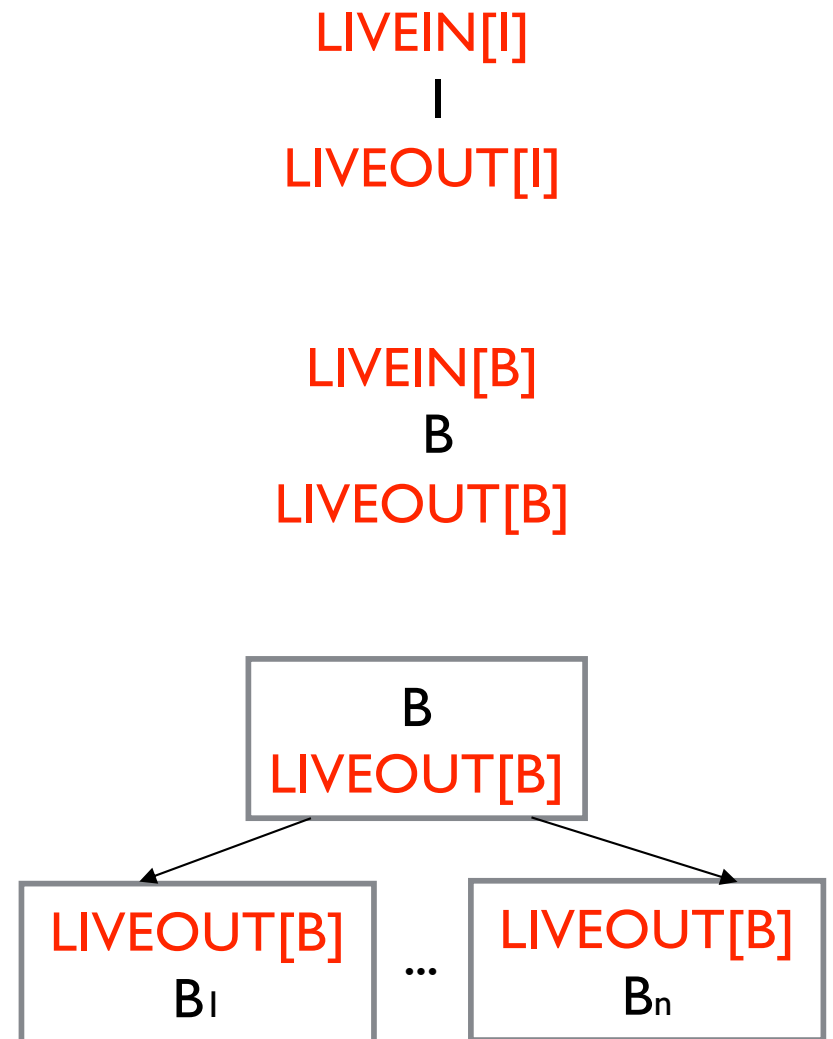
- A variable  $v$  is **live** at point  $p$  if and only if there is a path from  $p$  to a use of  $v$  along which  $v$  is not redefined.
- Usage
  - Global register allocation
  - Improve SSA construction
    - reduce # of f-functions
  - Detect references to uninitialized variables & defined but not used variables
  - Drive transformations
    - useless-store elimination

## *Live Variables at Special Points*

- For an instruction  $I$ 
  - $LIVEIN[I]$ : live variables at program point before  $I$
  - $LIVEOUT[I]$ : live variables at program point after  $I$
  
- For a basic block  $B$ 
  - $LIVEIN[B]$ : live variables at the entry point of  $B$
  - $LIVEOUT[B]$ : live variables at the exit point of  $B$
  
- If  $I =$  first instruction in  $B$ , then  $LIVEIN[B] = LIVEIN[I]$
- If  $I =$  last instruction in  $B$ , then  $LIVEOUT[B] = LIVEOUT[I]$

# How to Compute Liveness?

- Question 1: for each instruction  $I$ , what is the relation between  $LIVEIN[I]$  and  $LIVEOUT[I]$ ?
- Question 2: for each basic block  $B$ , what is the relation between  $LIVEIN[B]$  and  $LIVEOUT[B]$ ?
- Question 3: for each basic block  $B$  with successor blocks  $B_1, \dots, B_n$ , what is the relation between  $LIVEOUT[B]$  and  $LIVEOUT[B_1], \dots, LIVEOUT[B_n]$ ?

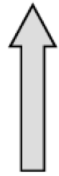


# Part 1: Analyze Instructions

- Question: what is the relation between the sets of live variables before and after an instruction I?

LIVEIN[I]  
|  
LIVEOUT[I]

Examples:



LIVEIN[I] = {y,z}	LIVEIN[I] = {y,z,t}	LIVEIN[I] = {x,t}
x = y+z;	x = y+z;	x = x+1;
LIVEOUT[I] = {z}	LIVEOUT[I] = {x,t}	LIVEOUT[I] = {x,t}

... is there a general rule?

# Analyze Instructions

## □ Two Rules:

- Each variable live after I is also live before I, unless I defines (writes) it.
- Each variable that I uses (reads) is also live before instruction I

## □ Mathematically:

$$\text{LIVEIN}[I] = ( \text{LIVEOUT}[I] - \text{def}[I] ) \cup \text{use}[I]$$

where:  $\text{def}[I]$  = variables defined (written) by instruction I

$\text{use}[I]$  = variables used (read) by instruction I

□ The information flows **backward!**

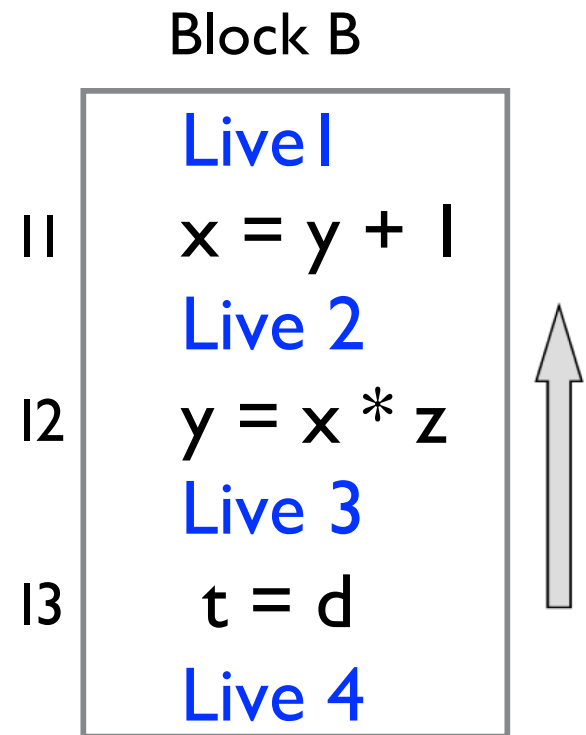
# Analyze block

- Example: block B with three instructions I1, I2, I3:

- $\text{Live1} = \text{LIVEIN}[B] = \text{LIVEIN}[I1]$
- $\text{Live2} = \text{LIVEOUT}[I1] = \text{LIVEIN}[I2]$
- $\text{Live3} = \text{LIVEOUT}[I2] = \text{LIVEIN}[I3]$
- $\text{Live4} = \text{LIVEOUT}[I3] = \text{LIVEOUT}[B]$

- Relation between Live sets:

- $\text{Live1} = (\text{Live2} - \{x\}) \cup \{y\}$
- $\text{Live2} = (\text{Live3} - \{y\}) \cup \{x, z\}$
- $\text{Live3} = (\text{Live4} - \{t\}) \cup \{d\}$



# Analyze Block

## □ Two Rules:

- Each variable live after B is also live before B, unless B defines (writes) it.
- Each variable  $v$  that B uses (reads) before any redefinition in B is also live before B

## □ Mathematically:

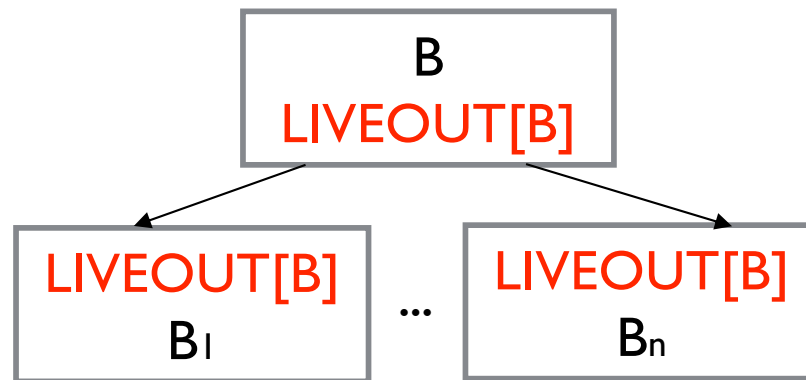
$$\text{LIVEIN}[B] = (\text{LIVEOUT}[B] - \text{VarKill}(B)) \cup \text{UEVar}(B)$$

where:

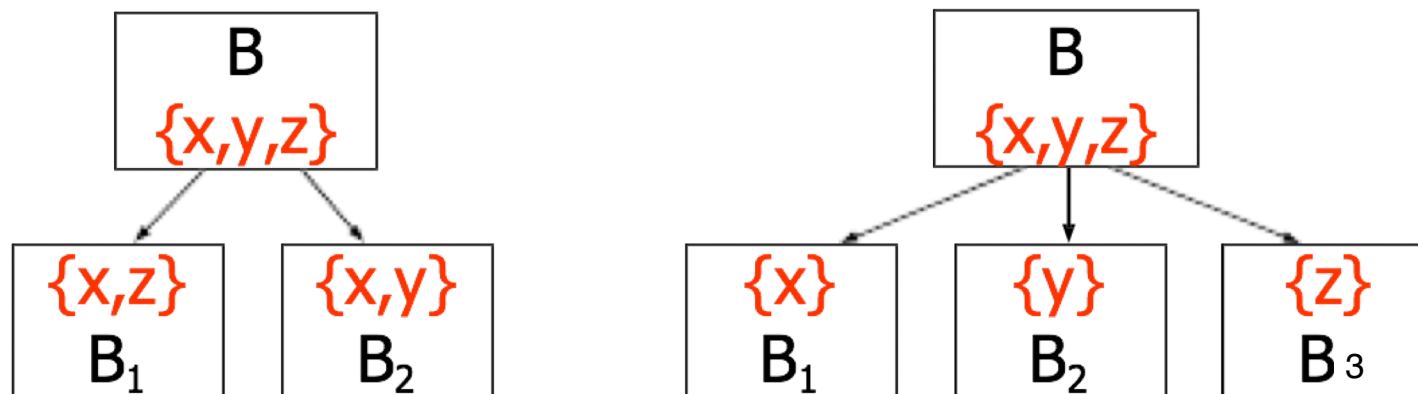
- $\text{VARKILL}(B)$  = variables that are defined in B
- $\text{UEVAR}(B)$  variables that are used in B before any redefinition in B, i.e., upward-exposed variables

# Analyze CFG

- Question: for each basic block  $B$  with successor blocks  $B_1, \dots, B_n$ , what is the relation between  $\text{LIVEOUT}[B]$  and  $\text{LIVEIN}[B_1], \dots, \text{LIVEIN}[B_n]$ ?



- Example:



- General rule?

# Analyze CFG

- **Rule:** A variables is live at end of block B if it is live at the beginning of **one (or more)** successor blocks
- **Mathematically:**

$$\begin{aligned} LIVEOUT[B] &= \bigcup_{B' \in succ(B)} LIVEIN[B'] \\ &= \bigcup_{B' \in succ(B)} ((LIVEOUT[B'] - VARKILL(B')) \cup UEVAR(B')) \end{aligned}$$

- Again, information flows **backward:** from successors B' of B to basic block

## *Equations for Live Variables*

- **LIVEOUT(B)** contains the name of every variable that is live on exit from n (a basic block)
- **UEVAR(B)** contains the upward-exposed variables in n, i.e. those that are used in n before any redefinition in n
- **VAR KILL(B)** contains all the variables that are defined in n
- Equation ( $n_f$  is the exit node of the CFG)

$$LIVEOUT[B] = \bigcup_{B' \in succ(B)} ((LIVEOUT[B'] - VAR KILL(B')) \cup UEVAR(B'))$$

Note:  $A - B = A \cup \overline{B}$

## *Three Steps in Data-Flow Analysis*

- Build a CFG
- Gather the initial information for each block (i.e., (UEVAR and VARKILL))
- Use an iterative fixed-point algorithm to propagate information around the CFG

# Algorithm

## // Get initial sets

```
for each block b
  UEVAR(b) = ∅
  VARKILL(b) = ∅
  for i=1 to number of instr in b
    (assuming inst I is “x= y op z”)
    if y ∉ VARKILL(b) then
      UEVAR(b) = UEVAR(b) ∪ {y}
    if z ∉ VARKILL(b) then
      UEVAR(b) = UEVAR(b) ∪ {z}
      VARKILL(b) = VARKILL(b) ∪ {x}
```

## // update LiveOut version 1

```
set LIVEOUT(bi) to ∅ for all blocks
Worklist ← {all blocks}
while (Worklist ≠ ∅)
  remove a block b from Worklist
  recompute LIVEOUT(b)
  if LIVEOUT(b) changed then
    Worklist ← Worklist ∪ pred(b)
```

$$LIVEOUT[B] = \bigcup_{B' \in succ(B)} ((LIVEOUT[B'] - VARKILL(B')) \cup UEVAR(B'))$$

# Algorithm

## // Get initial sets

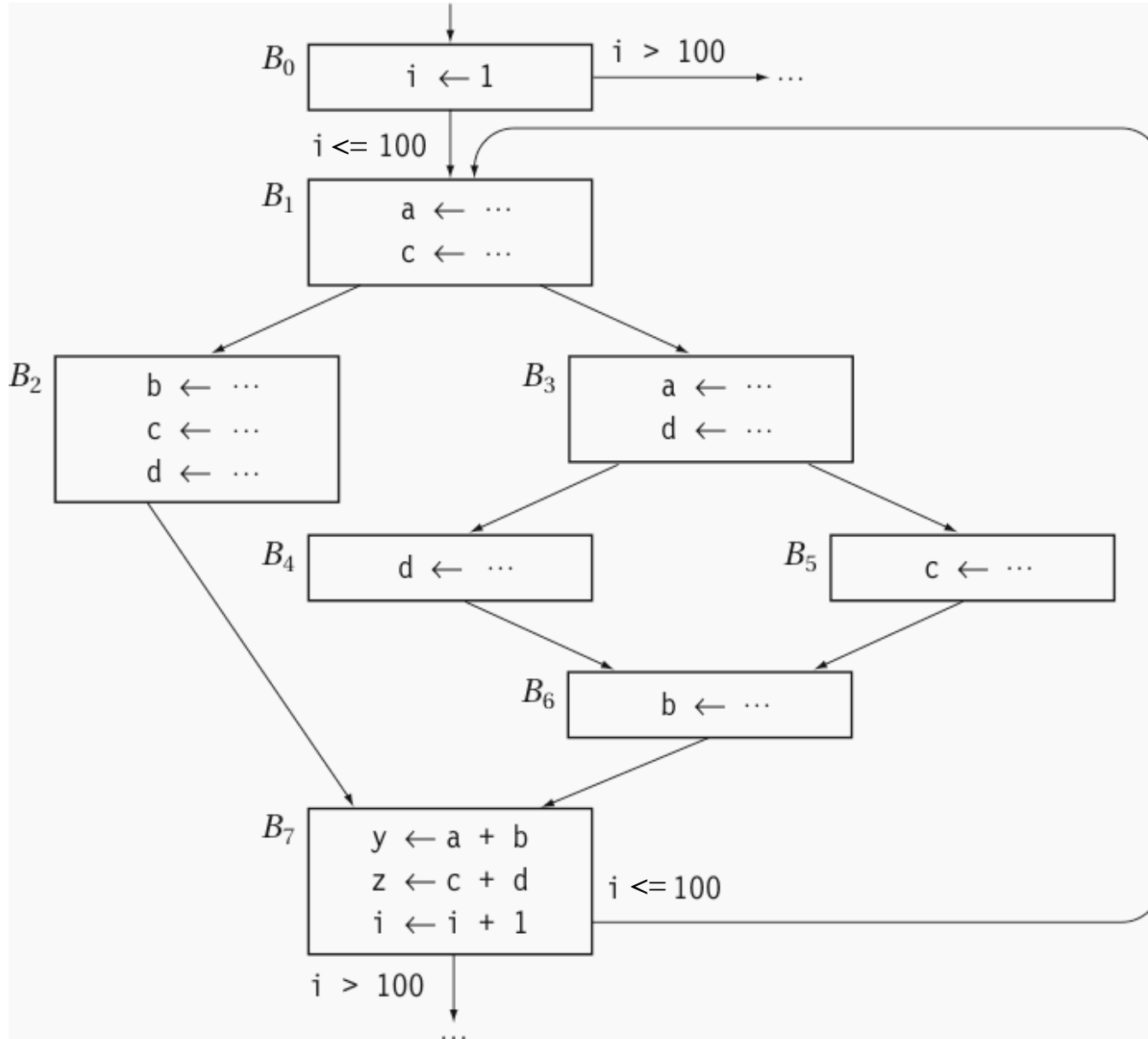
```
for each block b
  UEVAR(b) =  $\emptyset$ 
  VARKILL(b) =  $\emptyset$ 
  for i=1 to number of instr in b
    (assuming inst I is “x= y op z”)
    if  $y \notin \text{VARKILL}(b)$  then
      UEVAR(b) = UEVAR(b)  $\cup$  {y}
    if  $z \notin \text{VARKILL}(b)$  then
      UEVAR(b) = UEVAR(b)  $\cup$  {z}
  VARKILL(b) = VARKILL(b)  $\cup$  {x}
```

## // update LiveOut version2

```
set LIVEOUT(bi) to  $\emptyset$  for all blocks
changed = true
while (changed)
  changed = false
  for i = 1 to N (number of blocks)
    recompute LIVEOUT(i)
    if LIVEOUT(i) changed then
      changed = true
```

$$\text{LIVEOUT}[B] = \bigcup_{B' \in \text{succ}(B)} ((\text{LIVEOUT}[B'] - \text{VARKILL}(B')) \cup \text{UEVAR}(B'))$$

# Example



## *Example (cont.)*

	B0	B1	B2	B3	B4	B5	B6	B7
UEVar	∅	∅	∅	∅	∅	∅	∅	a,b,c,d,i
VarKill	i	a, c	b, c, d	a, d	d	c	b	y, z, i

## Example (cont.)

Can the algorithm converge in fewer iterations?

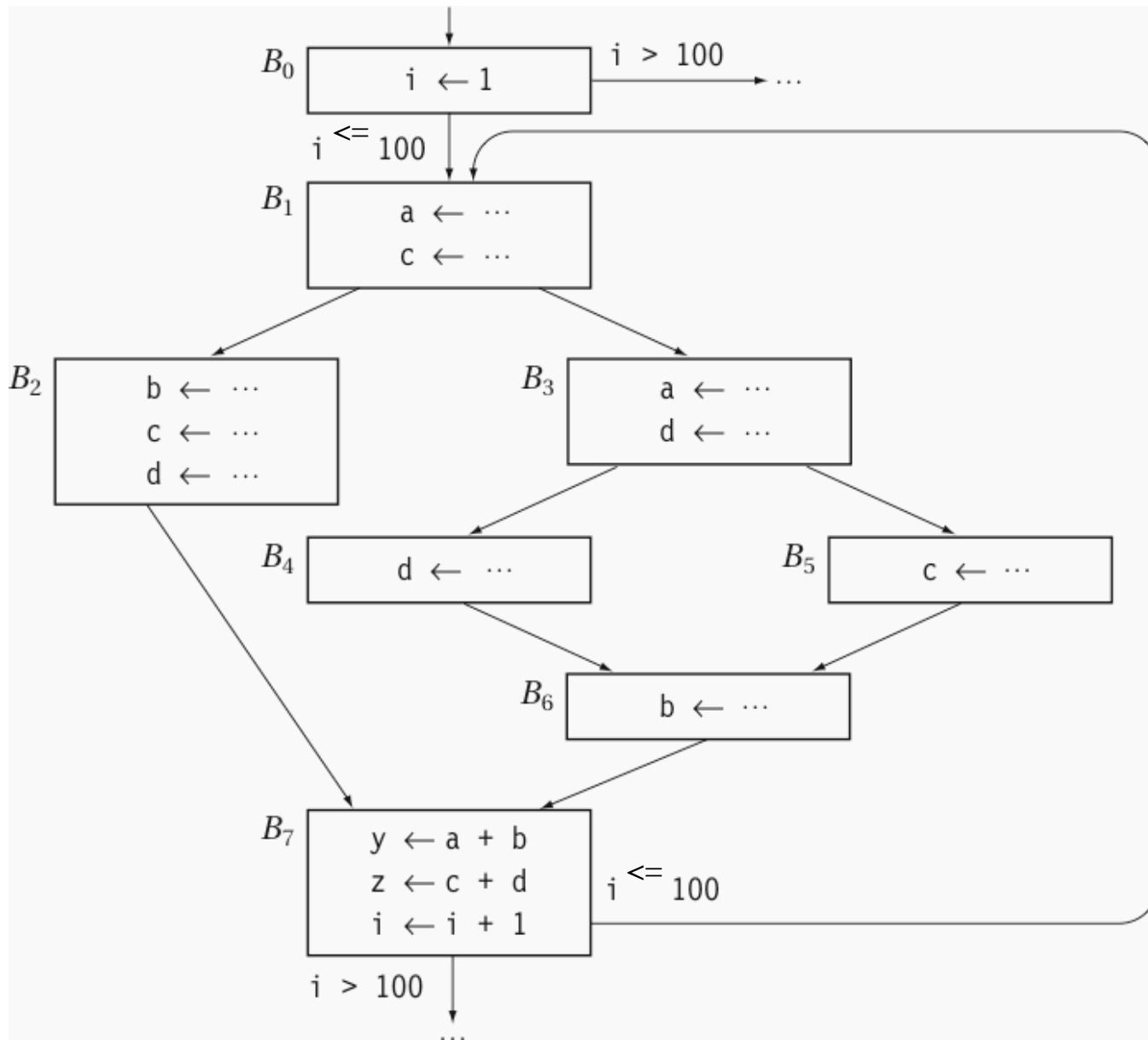
### LiveOut (b)

iteration	B0	B1	B2	B3	B4	B5	B6	B7
0	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
1	$\emptyset$	$\emptyset$	a,b,c,d,i	$\emptyset$	$\emptyset$	$\emptyset$	a,b,c,d,i	$\emptyset$
2	$\emptyset$	a,i	a,b,c,d,i	$\emptyset$	a,c,d,i	a,c,d,i	a,b,c,d,i	i
3	i	a,i	a,b,c,d,i	a,c,d,i	a,c,d,i	a,c,d,i	a,b,c,d,i	i
4	i	a,c,i	a,b,c,d,i	a,c,d,i	a,c,d,i	a,c,d,i	a,b,c,d,i	i
5	i	a,c,i	a,b,c,d,i	a,c,d,i	a,c,d,i	a,c,d,i	a,b,c,d,i	i

$$LIVEOUT[B] = \bigcup_{B' \in succ(B)} ((LIVEOUT[B'] - VARKILL(B')) \cup EVAR(B'))$$

$$LIVEOUT[B] = \bigcup_{B' \in succ(B)} ((LIVEOUT[B'] - VARKILL(B')) \cup UEVAR(B'))$$

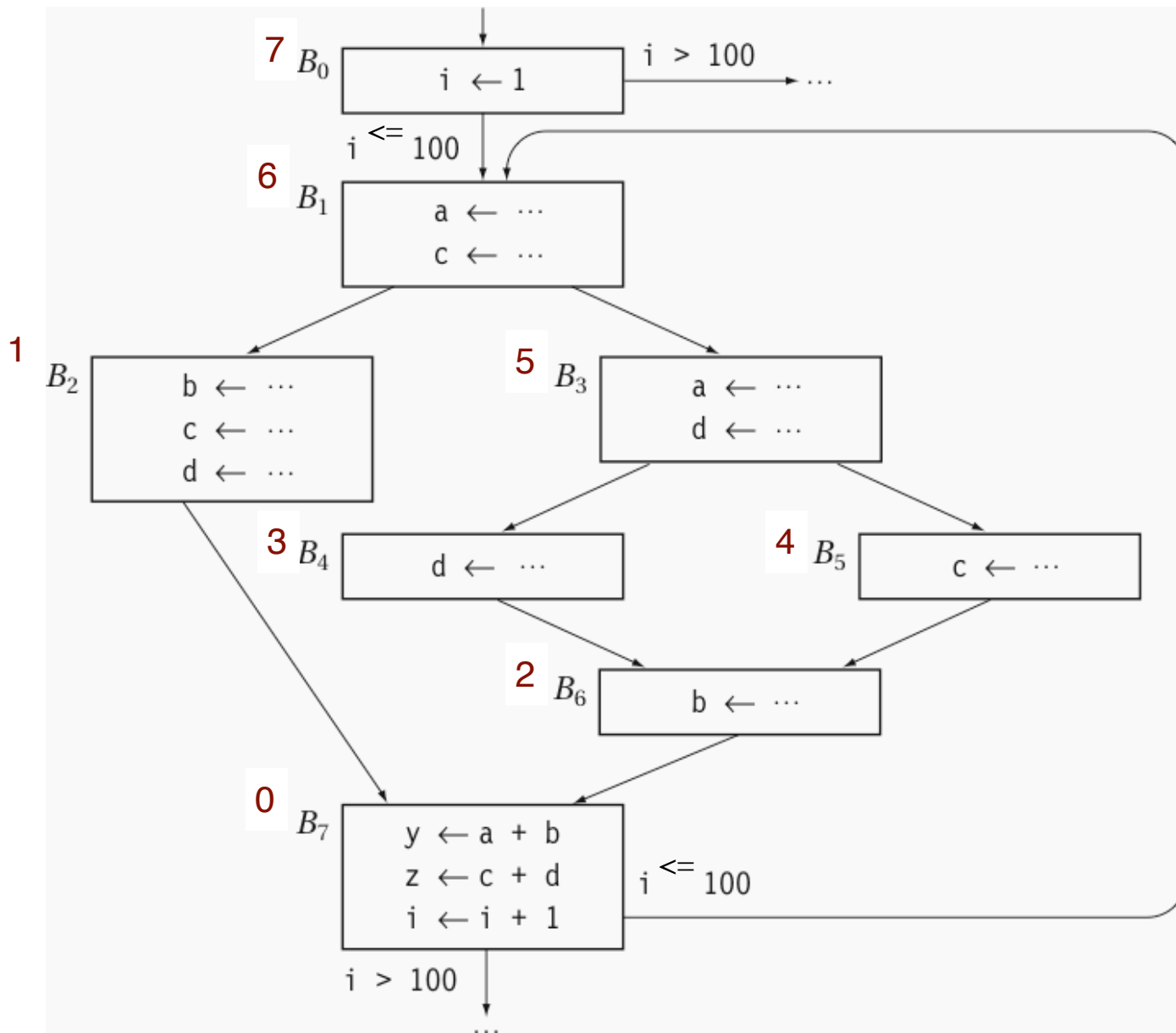
**Preorder:  
parents  
first.**  
w/o  
considering  
backedges.



$$LIVEOUT[B] = \bigcup_{B' \in succ(B)} ((LIVEOUT[B'] - VARKILL(B')) \cup UEVAR(B'))$$

**Postorder:  
children  
first.**

w/o  
considering  
backedges.



# Algorithm

// Get initial sets

```
for each block b
  UEVAR(b) = ∅
  VARKILL(b) = ∅
  for i=1 to number of instr in b
    (assuming inst I is “x= y op z”)
    if y ∉ VARKILL(b) then
      UEVAR(b) = UEVAR(b) ∪ {y}
    if z ∉ VARKILL(b) then
      UEVAR(b) = UEVAR(b) ∪ {z}
      VARKILL(b) = VARKILL(b) ∪ {x}
```

// update LiveOut version2

```
set LIVEOUT(bi) to ∅ for all blocks
changed = true
while (changed)
  changed = false
  for i = 0 to N
    // different orders could be used
    recompute LIVEOUT(i)
    if LIVEOUT(i) changed then
      changed = true
```

$$LIVEOUT[B] = \bigcup_{B' \in succ(B)} ((LIVEOUT[B'] - VARKILL(B')) \cup UEVAR(B'))$$

# *Postorder (5 iterations becomes 3)*

iteration	B0	B1	B2	B3	B4	B5	B6	B7
0	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
1	i	a,c,i	a,b,c,d,i	a,c,d,i	a,c,d,i	a,c,d,i	a,b,c,d,i	$\emptyset$
2	i	a,c,i	a,b,c,d,i	a,c,d,i	a,c,d,i	a,c,d,i	a,b,c,d,i	i
3	i	a,c,i	a,b,c,d,i	a,c,d,i	a,c,d,i	a,c,d,i	a,b,c,d,i	i

# Order

*Parent relation does not consider backedges.*

- **Preorder**: visit parents before children.
  - also called reverse postorder
- **Postorder**: visit children before parents.
  
- Forward problem (e.g., AVAIL):
  - A node needs the info of its predecessors.
  - Preorder on CFG.
- Backward problem (e.g., LIVEOUT):
  - A node needs the info of its successors.
  - Postorder on CFG.

# *Comparison with AVAIL*

- Common

  - Three steps

  - Fixed-point algorithm finds solution

- Differences

  - AVAIL: domain is a set of expressions

**Domain**

  - LIVEOUT: domain is a set of variables

  - AVAIL: forward problem

**Direction**

  - LIVEOUT: backward problem

  - AVAIL: intersection of all paths (**all path** problem)

**May/Must**

    - Also called Must Problem

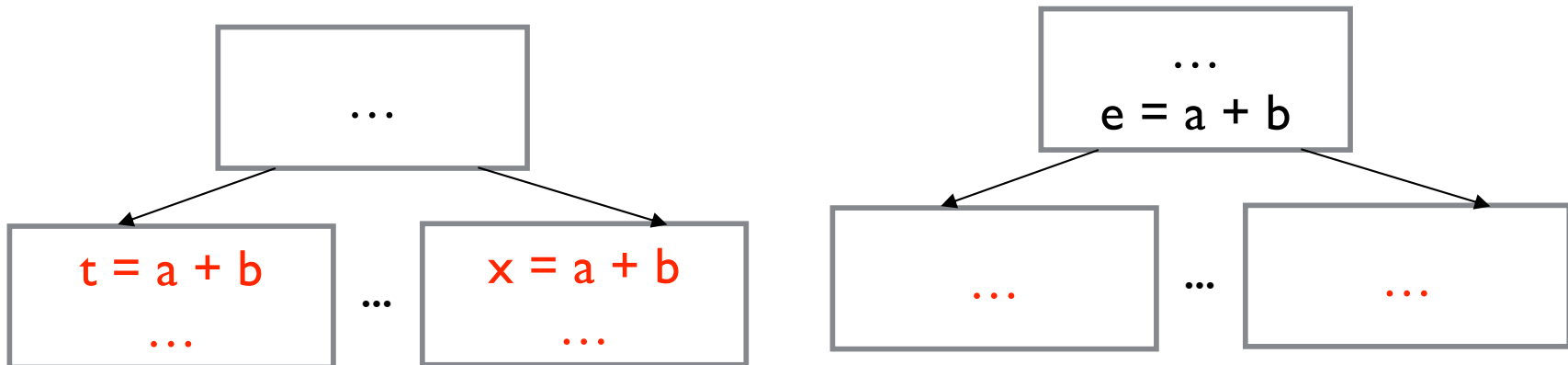
  - LIVEOUT: union of all paths (**any path** problem)

    - Also called May Problem

## *Other Data Flow Analysis*

# Very Busy Expressions

- Def:  $e$  is a very busy expression at the exit of block  $b$  if
  - $e$  is evaluated and used along every path that leaves  $b$ , and
  - evaluating  $e$  at the end of  $b$  produces the same result
- useful for code hoisting
- saves code space



# *Very Busy Expressions*

- VERYBUSY(b) contains expressions that are very busy at end of b
- UEEXPR(b): up exposed expressions (i.e. expressions defined in b and not subsequently killed in b)
- EXPRKILL(b): killed expressions

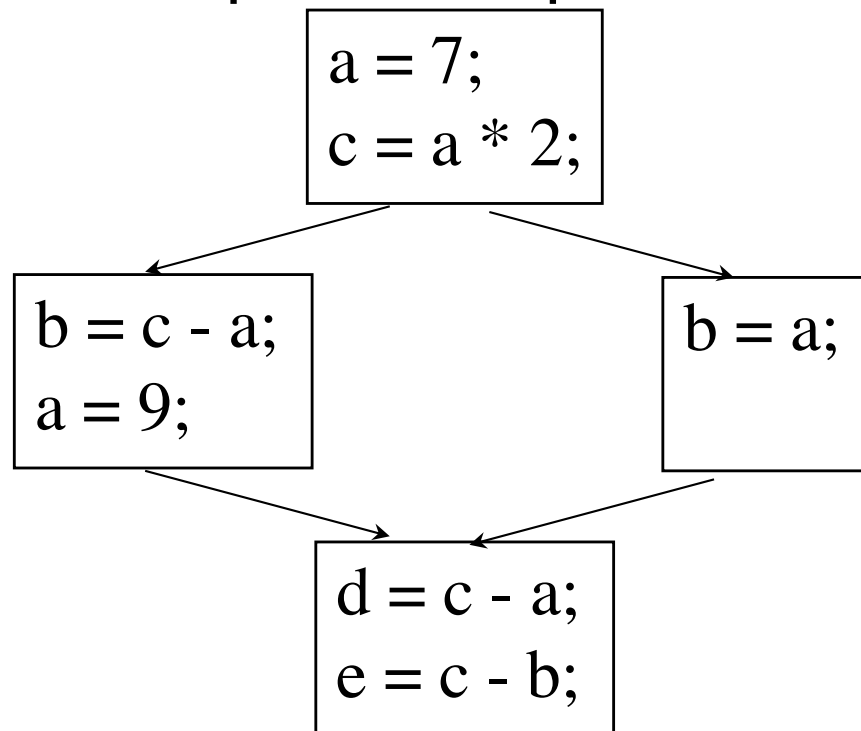
A **backward** flow problem, domain is the set of **expressions**

$$\text{VERYBUSY}(b) = \bigcap_{s \in \text{succ}(b)} \text{UEEXPR}(s) \cup (\text{VERYBUSY}(s) \cap \overline{\text{EXPRKILL}(s)})$$

$$\text{VERYBUSY}(n_f) = \emptyset$$

# Constant Propagation

- Def of a constant variable  $v$  at point  $p$ :
  - Along every path to  $p$ ,  $v$  has same known value
- Specialize computation at  $p$  based on  $v$ 's value



# *Constant Propagation: Another Data Flow Problem*

Domain is the set of pairs  $\langle v_i, c_i \rangle$  where  $v_i$  is a variable and  $c_i \in C$

$$\text{CONSTANTS}(b) = \bigwedge_{p \in \text{preds}(b)} f_p(\text{CONSTANTS}(p))$$

- $\bigwedge$  performs a pairwise meet on two sets of pairs
- $f_p(x)$  is a block specific function that models the effects of block  $p$  on the  $\langle v_i, c_i \rangle$  pairs in  $x$

A **forward** flow problem, domain is the set of **pairs**  $\langle v, c \rangle$ .

$C$ : constants or  $\perp$ .

$\perp$ : non-constant or unknown value

$$CONSTANTS(b) = \bigwedge_{p \in \text{preds}(b)} f_p(CONSTANTS(p))$$

Meet operation  $\langle v, c_1 \rangle \wedge \langle v, c_2 \rangle$

□  $\langle v, c_1 \rangle$  if  $c_1 = c_2$ , else  $\langle v, \perp \rangle$

$\perp$ : non-constant or unknown value

What about  $f_p$  ?

- if  $p$  has only one statement, update the constant set with
  - the results if operands are all constants
  - $\perp$  if the result is unknown or non-constant
- If  $p$  has  $n$  statements then
  - $f_p(CONSTANTS(p)) = f_n(f_{n-1}(f_{n-2}(\dots f_2(f_1(CONSTANTS(p))))))$ ,
    - where  $f_i$  is the function generated by the  $i$ th statement in  $p$

$$CONSTANTS(b) = \bigwedge_{p \in \text{preds}(b)} f_p(CONSTANTS(p))$$

Meet operation  $\langle v, c_1 \rangle \wedge \langle v, c_2 \rangle$

□  $\langle v, c_1 \rangle$  if  $c_1 = c_2$ , else  $\langle v, \perp \rangle$

$\perp$ : non-constant or unknown value

Formal definition of  $f_p$ :

□ If  $p$  has one statement then

□  $x \leftarrow y$  with  $CONSTANTS(p) = \{\dots \langle x, l_1 \rangle, \dots \langle y, l_2 \rangle \dots\}$

then  $f_p(CONSTANTS(p)) = \{CONSTANTS(p) - \langle x, l_1 \rangle\} \cup \langle x, l_2 \rangle$

□  $x \leftarrow y \text{ op } z$  with  $CONSTANTS(p) = \{\dots \langle x, l_1 \rangle, \dots \langle y, l_2 \rangle \dots, \dots \langle z, l_3 \rangle \dots\}$

then  $f_p(CONSTANTS(p)) = \{CONSTANTS(p) - \langle x, l_1 \rangle\} \cup \langle x, l_2 \text{ op } l_3 \rangle$

□ If  $p$  has  $n$  statements then

$$f_p(CONSTANTS(p)) = f_n(f_{n-1}(f_{n-2}(\dots f_2(f_1(CONSTANTS(p))) \dots)))$$

where  $f_i$  is the function generated by the  $i^{\text{th}}$  statement in  $p$

# *Data-Flow Analysis Frameworks*

- Generalizes and unifies data flow problems.
- Important components:
  - ◆ Direction D: forward or backward.
  - ◆ A *Semilattice*: a domain  $V$  and a *meet* operator  $\wedge$  that captures the effect of path confluence.
  - ◆ A transfer function  $F(m)$ : compute the effect of passing through a basic block and include function value at boundary conditions.

A *semilattice* is an algebra  $\mathcal{S} = (S, *)$  satisfying, for all  $x, y, z \in S$ ,

- (1)  $x * x = x$ ,
- (2)  $x * y = y * x$ ,
- (3)  $x * (y * z) = (x * y) * z$ .

# Examples

□ (D, V, F, ^)

□ LIVE

◆ D: backward

◆ V: all variables

◆  $F_m$ :  $UEVAR(m) \cup (LIVEOUT(m) \cap \overline{VARKILL(m)})$  ;  $LIVEOUT(n_f) = \phi$

◆  $\wedge$  :  $\cup$

□ AVAIL

◆ D: forward, V: all expressions

◆  $F_m$ :  $DEEXPR(m) \cup (AVAIL(m) \cap \overline{EXPRKILL(m)})$  ;  $AVAIL(n_o) = \phi$

◆  $\wedge$  :  $\cap$

## Summary

	Domain	Direction	Uses
AVAIL	Expressions	Forward	GCSE
LIVEOUT	Variables	Backward	Register alloc. Detect uninit. Construct SSA Useless-store Elim.
VERYBUSY	Expressions	Backward	Hoisting
CONSTANT	Pairs $\langle v, c \rangle$	Forward	Constant folding

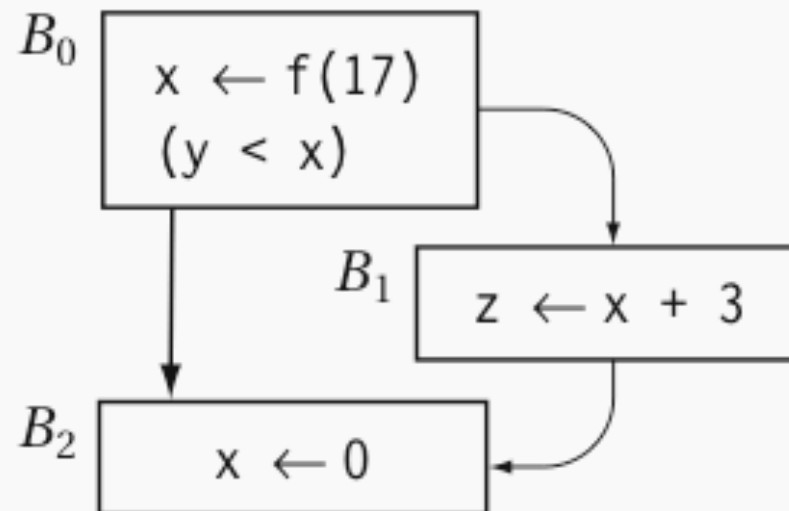
# *Why to Study Data Flow Analysis*

- Data-flow analysis
  - A collection of techniques for compile-time reasoning about the run-time flow of values.
- Backbone of scalar optimizing compilers

## *Limitation of Data-Flow Analysis*

- Imprecision from pointers, and procedure calls
- Assume all paths will be taken

```
x ← f(17)
if (y < x) then
  z ← x + 3
x ← 0
```



If  $y$  is always no less than  $x$ ,  $x$  is not live before  $B_2$ .  
But data-flow analysis may not figure that out.