



南方科技大学  
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

# Chapter 8: Introduction to Data-Flow Analysis

Yepang Liu

[liuyp1@sustech.edu.cn](mailto:liuyp1@sustech.edu.cn)

# Outline

- Data-Flow Analysis
- Classic Data-Flow Problems

# Introduction

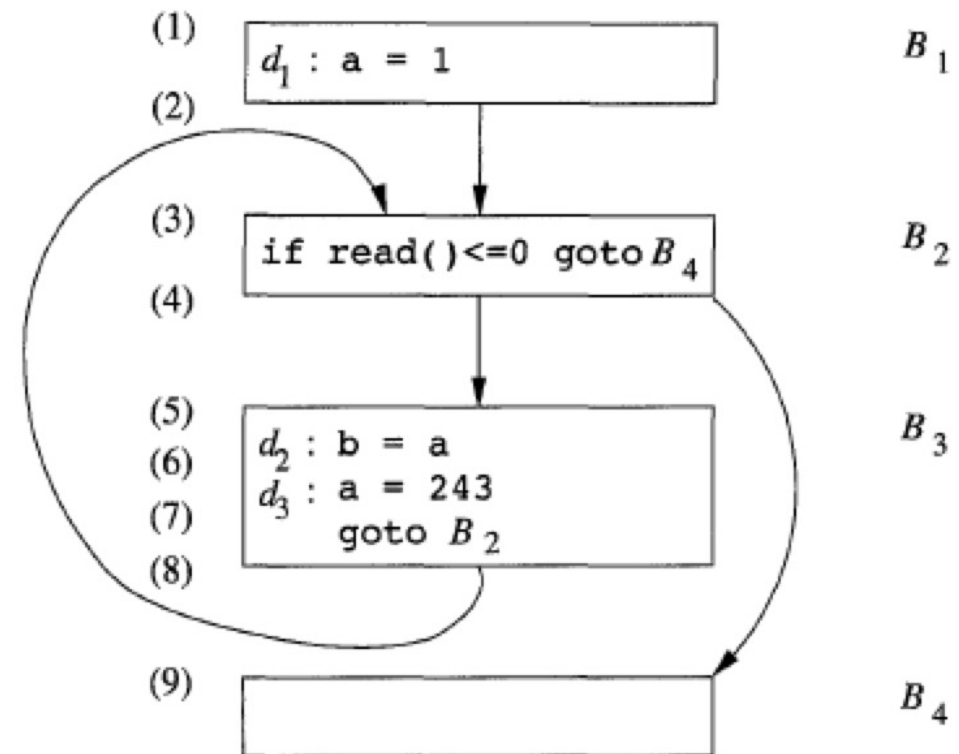
- Code optimization (or code improvement)
  - Eliminating **unnecessary** instructions in object code
  - Replacing one sequence of instructions by a **faster** sequence of instructions that does the same thing
- Global optimizations
  - Performed on a flow graph as a whole, involving multiple basic blocks
  - Most global optimizations are based on **data-flow analysis**
- A compiler knows only how to apply relatively **low-level semantic transformations** to optimize code
  - High-level optimizations: architectural/algorithmic changes, refactoring, etc.

# What Is Data-Flow Analysis?

- A body of techniques that derive information about the flow of data along program execution paths
  - **Example 1:** Whether two textually identical expressions evaluate to the same value along any execution path? (Identify common subexpressions)
  - **Example 2:** Whether the result of an assignment is used on subsequent execution paths? (Eliminate dead code)
- Data-flow analysis is the foundation of many optimizations

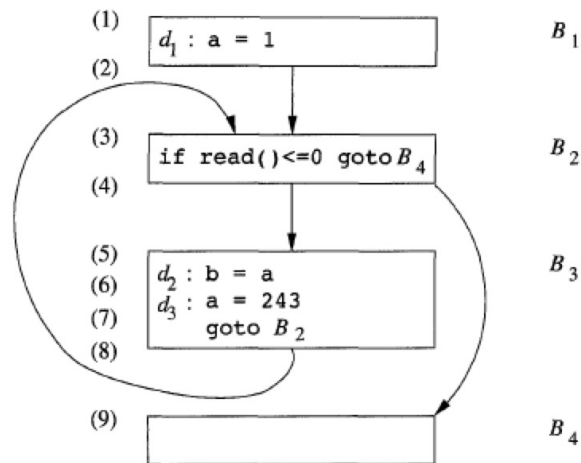
# The Data-Flow Abstraction

- **Program points**: points before and after each statement
  - Within one block, the program point after a statement is the same as the program point before the next statement (e.g., 6)
  - If there is an edge from block  $B$  to block  $C$ , then the program point after the last statement of  $B$  **may be followed immediately by** the program point before the first statement of  $C$  (e.g., 8 and 3)



# The Data-Flow Abstraction

- An **execution path** from point  $p_1$  to point  $p_n$  is a sequence of points  $p_1, p_2, \dots, p_n$  such that for each  $i = 1, 2, \dots, n-1$ :
  - **either**  $p_i$  is the point immediately preceding a statement and  $p_{i+1}$  is the point immediately follow that statement (the “**within block**” case)
  - **or**  $p_i$  is the end of a block and  $p_{i+1}$  is the beginning of a successor block (the “**across block**” case)

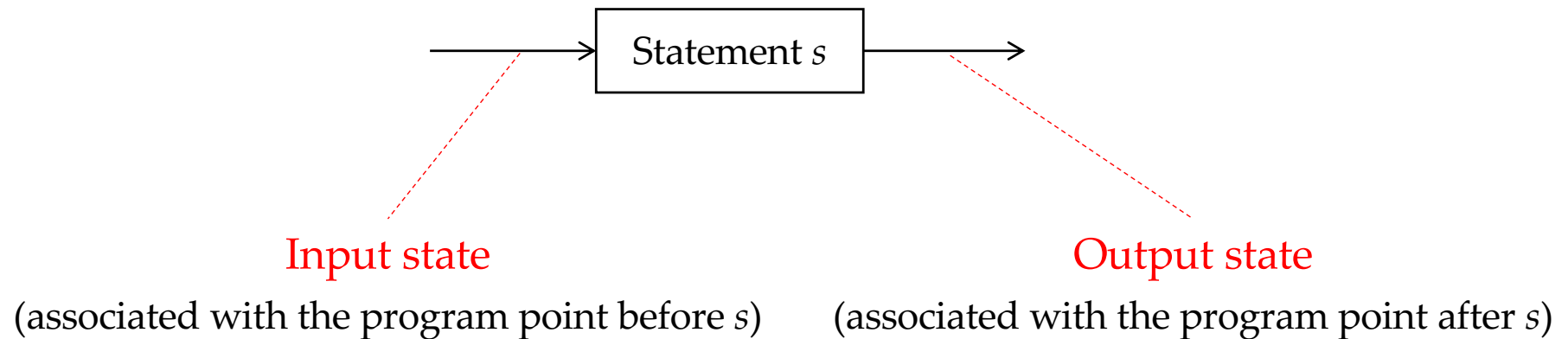


The shortest complete execution path: (1, 2, 3, 4, 9)

The next shortest path executing one iteration of the loop: (1, 2, 3, 4, 5, 6, 7, 8, 3, 4, 9)

# The Data-Flow Abstraction

- **Program execution** can be viewed as *a series of transformations of the program state* (the values of all variables)
  - Each execution of a statement transforms an **input state** to a new **output state**

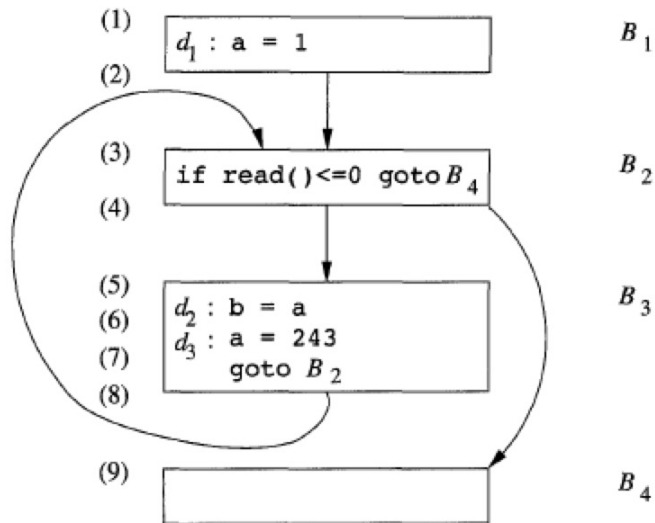


# The Data-Flow Abstraction

- When analyzing a program, we must consider all the possible execution paths through a flow graph
  - In general, there is **an infinite number** of possible paths due to the existence of loops and recursions
- Data-flow analyses **summarize all the possible program states** that can occur at a program point with **a finite set of facts**
  - Different analyses may choose to abstract out different information (i.e., obtaining approximations)
  - In general, no analysis is necessarily a perfect representation of the state



# Example (1)



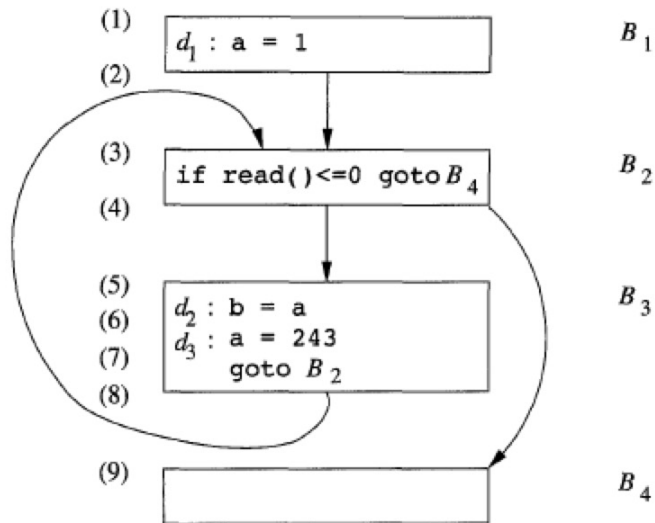
## Reaching definitions:

- The first time point (5) is executed, the value of  $a$  is 1, i.e., definition  $d_1$  reaches (5) in the first iteration
- In subsequent iterations,  $d_3$  reaches point (5) and the value of  $a$  is 243
- So, we may summarize all the program states at (5) by saying that the value of  $a$  is one of  $\{1, 243\}$  and defined by one of  $\{d_1, d_3\}$

Static approximation.

The exact value depends on the execution path at runtime.

# Example (2)

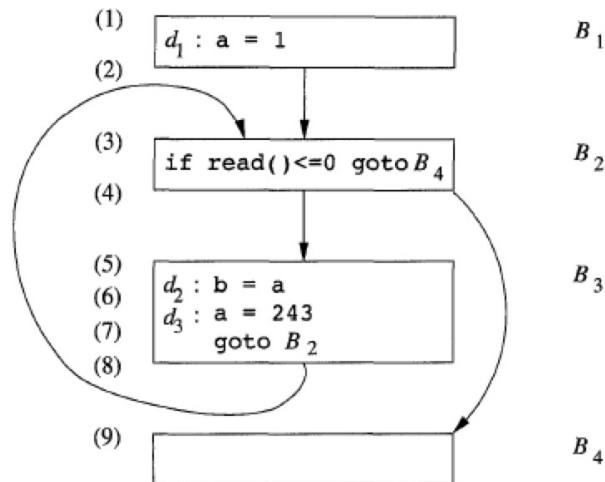


## Constant folding:

- Find those definitions that are the unique definition of their variable to reach a given program point, regardless of the execution paths
- For this task, we describe some variables as “constant” or “not a constant”
  - $a$  is not a constant at point (5) and thus cannot be replaced by a constant value

# The Data-Flow Analysis Schema

- We associate with every program point a *data-flow value* that represents an abstraction of program states observed for that point\*



## Reaching definitions:

The data-flow value at (5) is the set of  $a$ 's definitions that can reach (5), that is,  $\{d_1, d_3\}$

## Constant folding:

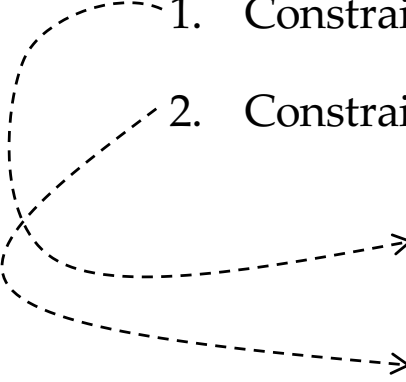
The data-flow value for variable  $a$  at (5) is **NAC** (not a constant)

\* In each application of data-flow analysis, the set of possible data-flow values is the domain for the application.

# The Data-Flow Analysis Schema

- $IN[s]$ : The data-flow value **before** the statement  $s$
- $OUT[s]$ : The data-flow value **after** the statement  $s$
- The *data-flow problem* is to find a solution to a set of constraints on the  $IN[s]$ 's and  $OUT[s]$ 's for all statements  $s$

1. Constraints based on the semantics of the statements ("*transfer functions*")
2. Constraints based on the flow of control



Statement executions may alter data-flow values

Control flows propagate data-flow values

# Transfer Functions (传递函数)

- The relationship between the data-flow values before and after each statement is known as a *transfer function*
- In a **forward-flow problem** (information propagate forward along execution paths), the transfer function  $f_s$  takes the data-flow value before the statement  $s$  and produces a new data-flow value after  $s$ 
  - $OUT[s] = f_s(IN[s])$
- In a backward-flow problem (information flow backwards up the execution paths), the transfer function  $f_s$  takes the data-flow value after the statement  $s$  and produces a new data-flow value before  $s$ 
  - $IN[s] = f_s(OUT[s])$

# Control-Flow Constraints

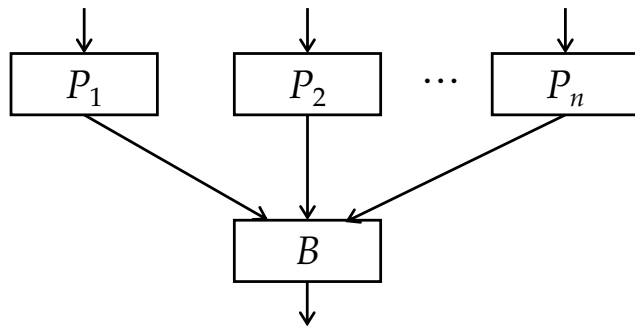
1. **Within a basic block** of statements  $s_1, s_2, \dots, s_n$ , the data-flow value out of  $s_i$  is the same as that into  $s_{i+1}$ 
  - $IN[s_{i+1}] = OUT[s_i]$ , for all  $i = 1, 2, \dots, n-1$
2. **Control-flow edges between basic blocks may create more complex constraints**
  - For example, in reaching definitions analysis, the set of definitions reaching the leader statement of a block should be the union of the definitions after the last statements of each of the predecessor blocks

# Data-Flow Schemas on Basic Blocks (The Intra-Block Case)

- Within a basic block, control flows from the beginning to the end without interruption or branching
- For a block  $B$  consisting of statements  $s_1, s_2, \dots, s_n$ , we have<sup>\*</sup>
  - $\text{IN}[B] = \text{IN}[s_1]$
  - $\text{OUT}[B] = \text{OUT}[s_n]$
  - $\text{OUT}[B] = f_B(\text{IN}[B])$ , where  $f_B = f_{s_n} \circ \dots \circ f_{s_2} \circ f_{s_1}$  (composing statement-level transfer functions)

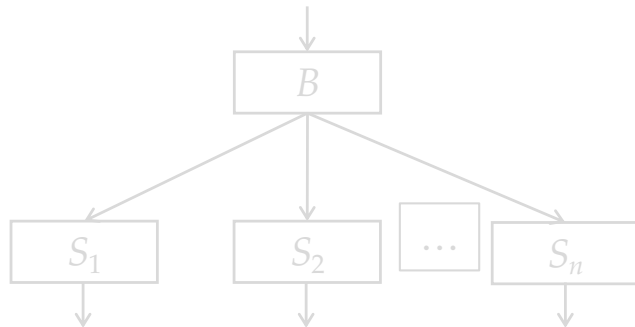
<sup>\*</sup> For backward-flow problem,  $\text{IN}[B] = f_B(\text{OUT}[B])$ , where  $f_B = f_{s_1} \circ f_{s_2} \circ \dots \circ f_{s_n}$

# Data-Flow Schemas on Basic Blocks (The Inter-Block Case)



Forward-flow problem:

$$\text{IN}[B] = \bigcup_{P \text{ a predecessor of } B} \text{OUT}[P]$$



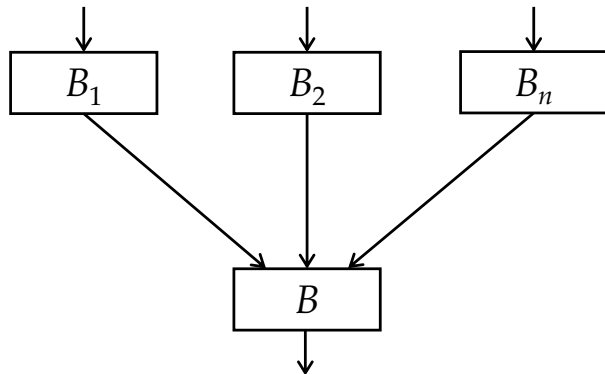
Backward-flow problem:

$$\text{OUT}[B] = \bigcup_{S \text{ a successor of } B} \text{IN}[S]$$

$\cup$  is a generic *meet operator* depending on specific problems.



# Example (Constant Folding Problem)



**If:**

$\text{OUT}[B_1]: x: 3; y: 4; z: \text{NAC}$

$\text{OUT}[B_2]: x: 3; y: 5; z: 7$

$\text{OUT}[B_3]: x: 3; y: 4; z: 7$

**then:**

$\text{IN}[B]: x: 3; y: \text{NAC}; z: \text{NAC}$

# Solutions to Data-Flow Equations

- Data-flow equations usually do not have a unique solution
  - All data-flow schemas compute **approximations** to the ground truth
- Our goal is to **find the most “precise” solution** that satisfies both **control-flow** and **transfer** constraints
  - Being precise enables valid code improvements (in general, a higher precision leads to a better improvement)
  - Satisfying constraints guarantees the safety of transformations (does not change program semantics)

# Outline

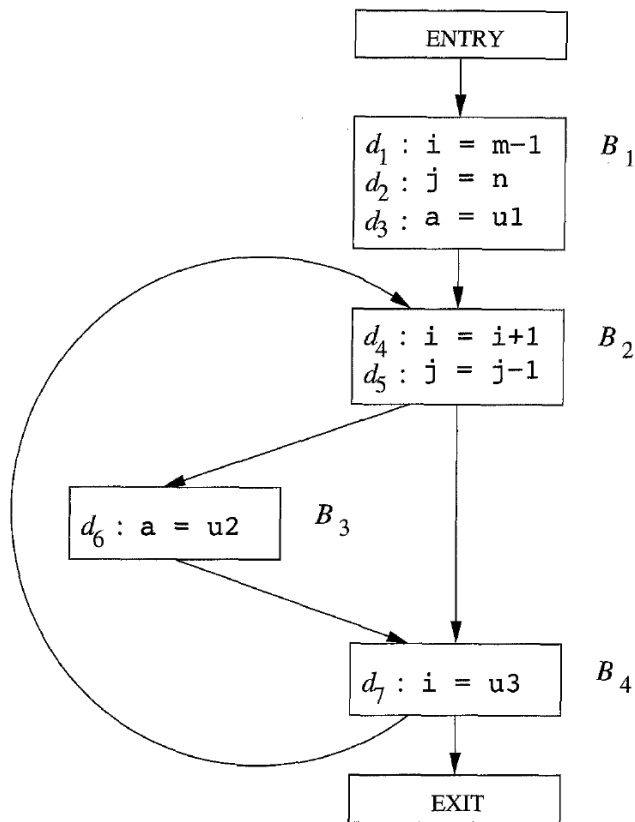
- Data-Flow Analysis
- Classic Data-Flow Problems {
  - Reaching Definitions
  - Live Variables\*
  - Available Expressions\*

\* Optional self-study materials

# Reaching Definitions

- A definition  $d$  of some variable  $x$  *reaches* a point  $p$  if there is a path from the program point after  $d$  to  $p$ , such that  $d$  is not “killed” along that path
  - $d$  is *killed* if there is any other definition of  $x$  along the path
  - Intuitively, if  $d$  reaches the point  $p$ , then  $d$  might be the last definition of  $x$

# Example



Which definitions reach the block  $B_2$ ?

- $d_1, d_2, d_3$  reach  $B_2$
- $d_5$  reaches  $B_2$  since there is no other definition to  $j$  in the loop
- $d_4$  does not reach  $B_2$  since  $i$  is always redefined by  $d_7$  (i.e.,  $d_7$  reaches  $B_2$ )
- $d_6$  reaches  $B_2$

# Reaching Definitions

- For reaching definitions analysis, we **allow inaccuracies**. However, the decisions should be *“safe”*.
  - It is ok that some inferred reaching defs cannot actually reach a point
  - However, those defs that can reach a point should be identified
  - In other words, over-approximations are acceptable
- ***Reason of inaccuracies:*** In general, to decide whether each path in a flow graph can be taken is an undecidable problem\*
  - We often simply assume that every path in the flow graph can be followed in some execution of the program

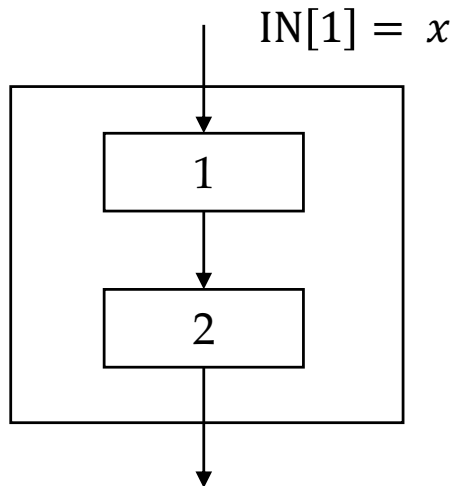
\* This is the well-known *path feasibility problem*.

# Transfer Equations

- Consider a statement  $d: u = v + w$ 
  - It *generates* a definition  $d$  of variable  $u$  and *kills* all other definitions of  $u$
- Transfer function of a statement is in *gen-kill form*:
  - $f_d(x) = gen_d \cup (x - kill_d)$ 
    - $x$  is the data-flow value (reaching definitions) before the statement
    - $gen_d$  is the set of generated definitions, i.e.,  $\{d\}$
    - $kill_d$  is the set of killed definitions, i.e., all other definitions of  $u$

# Transfer Equations

- Transfer function of a basic block is also in *gen-kill form*
- Consider a block of only two statements:



$$\text{OUT}[1] = f_1(x) = \text{gen}_1 \cup (x - \text{kill}_1)$$

$$\text{OUT}[2] = f_2(\text{IN}[2]) = f_2(\text{OUT}[1])$$

$$= \text{gen}_2 \cup (\text{OUT}[1] - \text{kill}_2)$$

$$= \text{gen}_2 \cup ((\text{gen}_1 \cup (x - \text{kill}_1)) - \text{kill}_2)$$

$$= \text{gen}_2 \cup (\text{gen}_1 - \text{kill}_2) \cup (x - (\text{kill}_1 \cup \text{kill}_2))$$



# Transfer Equations

- Generalize to a block  $B$  with  $n$  statements, we have:

$$f_B(x) = gen_B \cup (x - kill_B)$$

where

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

and

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

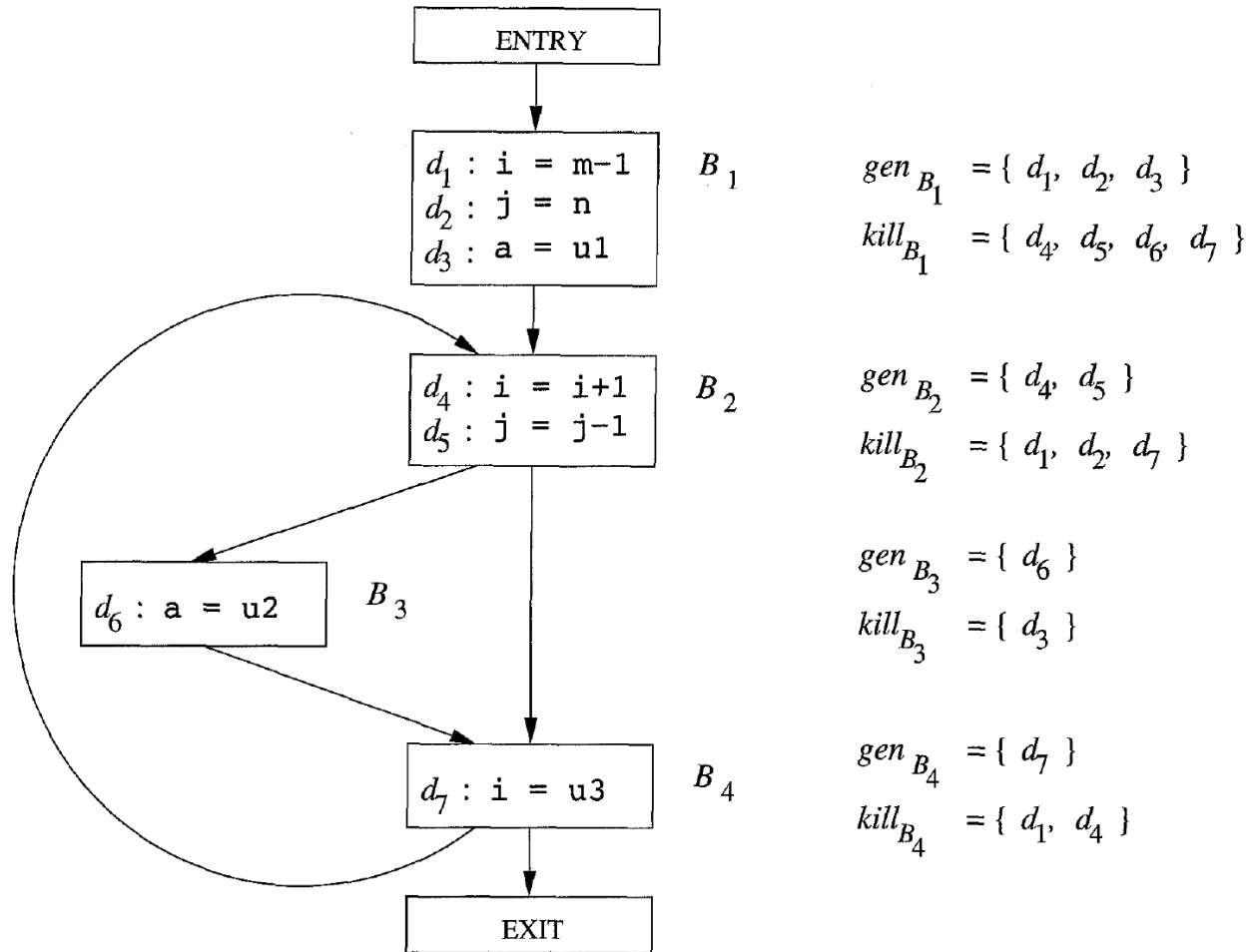
$kill_B$

The **kill set** is the union of all the defs killed by the individual statements

$gen_B$

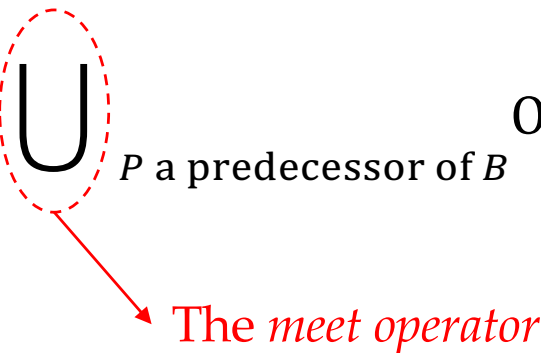
The **gen set** contains all the defs inside the block that are *downward exposed* (visible immediately after the block)

# Example



# Control-Flow Equations

$$\text{IN}[B] = \bigcup_{P \text{ a predecessor of } B} \text{OUT}[P]$$

The *meet operator*

- **Rationale:**

- A definition reaches a program point as long as there exists at least one path along which the definition reaches. This explains why  $\text{OUT}[P] \subseteq \text{IN}[B]$ .
- Since a definition cannot reach a point unless there is a path along which it reaches,  $\text{IN}[B]$  needs to be no larger than the union of the reaching definitions of all the predecessor blocks

# Problem Formulation

- The reaching definitions problem is defined by the following equations:
  - For the ENTRY block:  $\text{OUT}[\text{ENTRY}] = \emptyset$ 
    - Since no definitions reach the beginning of a program
  - For all basic blocks  $B$  other than ENTRY:
    - $\text{OUT}[B] = \text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$
    - $\text{IN}[B] = \bigcup_{P \text{ a predecessor of } B} \text{OUT}[P]$

# The Iterative Algorithm

- **Input:** A flow graph with  $kill_B$  and  $gen_B$  computed for each block  $B$
- **Output:**  $IN[B]$  and  $OUT[B]$  for each block  $B$

```
1)  OUT[ENTRY] =  $\emptyset$ ;  
2)  for (each basic block  $B$  other than ENTRY)  $OUT[B] = \emptyset$ ;  
3)  while (changes to any OUT occur)  
4)      for (each basic block  $B$  other than ENTRY) {  
5)           $IN[B] = \bigcup_{P \text{ a predecessor of } B} OUT[P]$ ;  
6)           $OUT[B] = gen_B \cup (IN[B] - kill_B)$ ;  
      }
```

# Example

- 1)  $OUT[ENTRY] = \emptyset$ ;
- 2) **for** (each basic block  $B$  other than ENTRY)  $OUT[B] = \emptyset$ ;
- 3) **while** (changes to any  $OUT$  occur)
- 4)     **for** (each basic block  $B$  other than ENTRY) {
- 5)          $IN[B] = \bigcup_{P \text{ a predecessor of } B} OUT[P]$ ;
- 6)          $OUT[B] = gen_B \cup (IN[B] - kill_B)$ ;
- }

Initialization:

$OUT[ENTRY] = \emptyset$

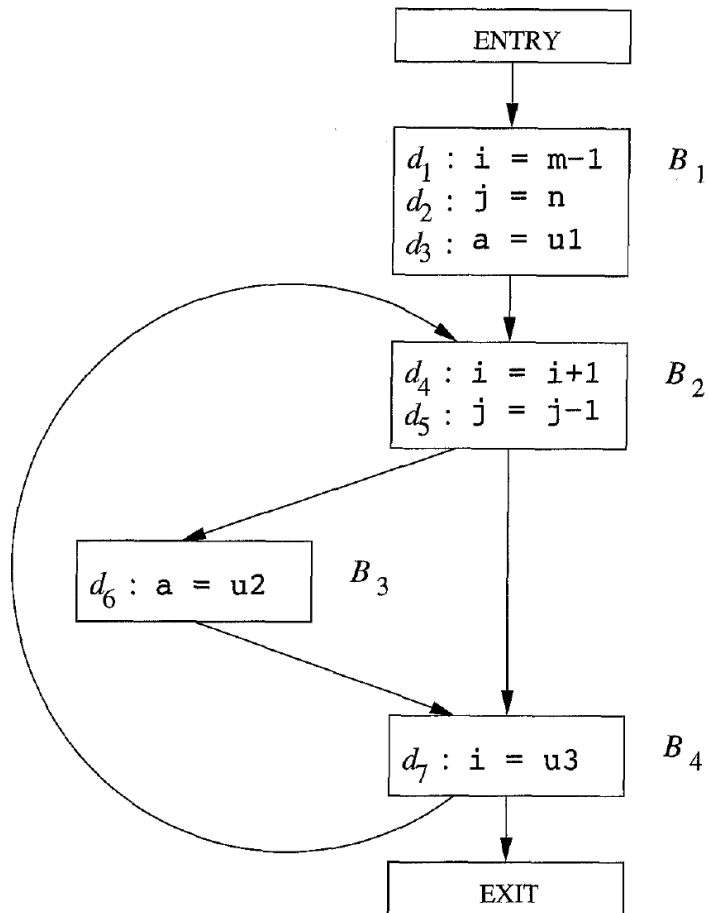
$OUT[B_1] = \emptyset$

$OUT[B_2] = \emptyset$

$OUT[B_3] = \emptyset$

$OUT[B_4] = \emptyset$

$OUT[EXIT] = \emptyset$



$gen_{B_1} = \{ d_1, d_2, d_3 \}$

$kill_{B_1} = \{ d_4, d_5, d_6, d_7 \}$

$gen_{B_2} = \{ d_4, d_5 \}$

$kill_{B_2} = \{ d_1, d_2, d_7 \}$

$gen_{B_3} = \{ d_6 \}$

$kill_{B_3} = \{ d_3 \}$

$gen_{B_4} = \{ d_7 \}$

$kill_{B_4} = \{ d_1, d_4 \}$

# Example

```

1) OUT[ENTRY] =  $\emptyset$ ;
2) for (each basic block  $B$  other than ENTRY) OUT[ $B$ ] =  $\emptyset$ ;
3) while (changes to any OUT occur)
4)   for (each basic block  $B$  other than ENTRY) {
5)     IN[ $B$ ] =  $\bigcup_{P \text{ a predecessor of } B} \text{OUT}[P]$ ;
6)     OUT[ $B$ ] =  $\text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$ ;
   }

```

## Iteration #1:

OUT[ENTRY] =  $\emptyset$

IN[ $B_1$ ] =  $\emptyset$

OUT[ $B_1$ ] = { $d_1, d_2, d_3$ }

IN[ $B_2$ ] = { $d_1, d_2, d_3$ }

OUT[ $B_2$ ] = { $d_3, d_4, d_5$ }

IN[ $B_3$ ] = { $d_3, d_4, d_5$ }

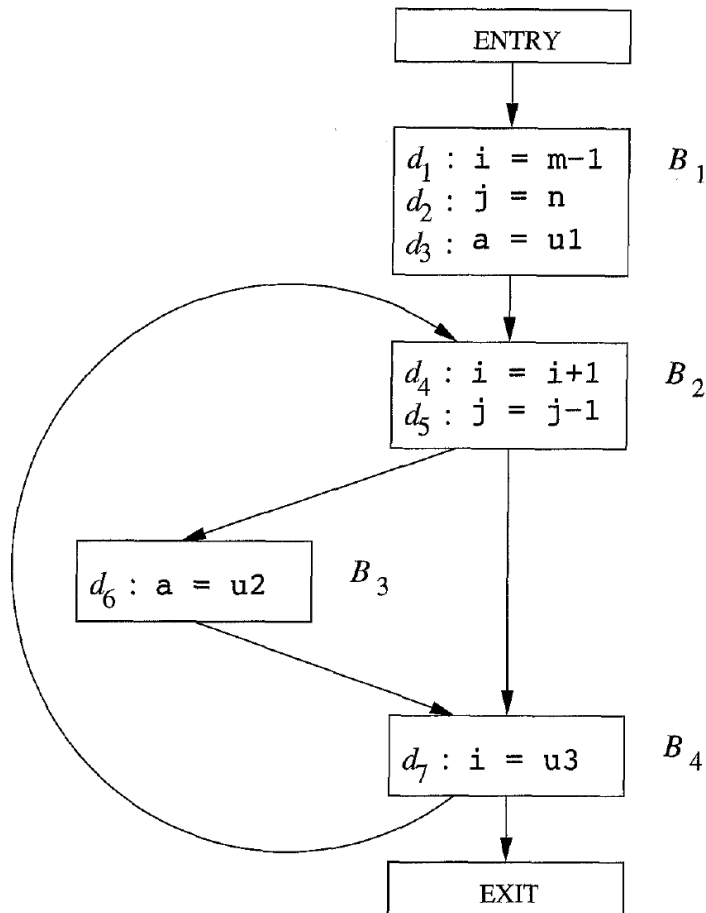
OUT[ $B_3$ ] = { $d_4, d_5, d_6$ }

IN[ $B_4$ ] = { $d_3, d_4, d_5, d_6$ }

OUT[ $B_4$ ] = { $d_3, d_5, d_6, d_7$ }

IN[EXIT] = { $d_3, d_5, d_6, d_7$ }

OUT[EXIT] = { $d_3, d_5, d_6, d_7$ }



$\text{gen}_{B_1} = \{ d_1, d_2, d_3 \}$

$\text{kill}_{B_1} = \{ d_4, d_5, d_6, d_7 \}$

$\text{gen}_{B_2} = \{ d_4, d_5 \}$

$\text{kill}_{B_2} = \{ d_1, d_2, d_7 \}$

$\text{gen}_{B_3} = \{ d_6 \}$

$\text{kill}_{B_3} = \{ d_3 \}$

$\text{gen}_{B_4} = \{ d_7 \}$

$\text{kill}_{B_4} = \{ d_1, d_4 \}$

# Example

```

1) OUT[ENTRY] =  $\emptyset$ ;
2) for (each basic block  $B$  other than ENTRY) OUT[ $B$ ] =  $\emptyset$ ;
3) while (changes to any OUT occur)
4)   for (each basic block  $B$  other than ENTRY) {
5)     IN[ $B$ ] =  $\bigcup_{P \text{ a predecessor of } B} \text{OUT}[P]$ ;
6)     OUT[ $B$ ] =  $\text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$ ;
   }

```

## Iteration #2:

OUT[ENTRY] =  $\emptyset$

IN[ $B_1$ ] =  $\emptyset$

OUT[ $B_1$ ] = { $d_1, d_2, d_3$ }

IN[ $B_2$ ] = { $d_1, d_2, d_3, d_5, d_6, d_7$ }

OUT[ $B_2$ ] = { $d_3, d_4, d_5, d_6$ }

IN[ $B_3$ ] = { $d_3, d_4, d_5, d_6$ }

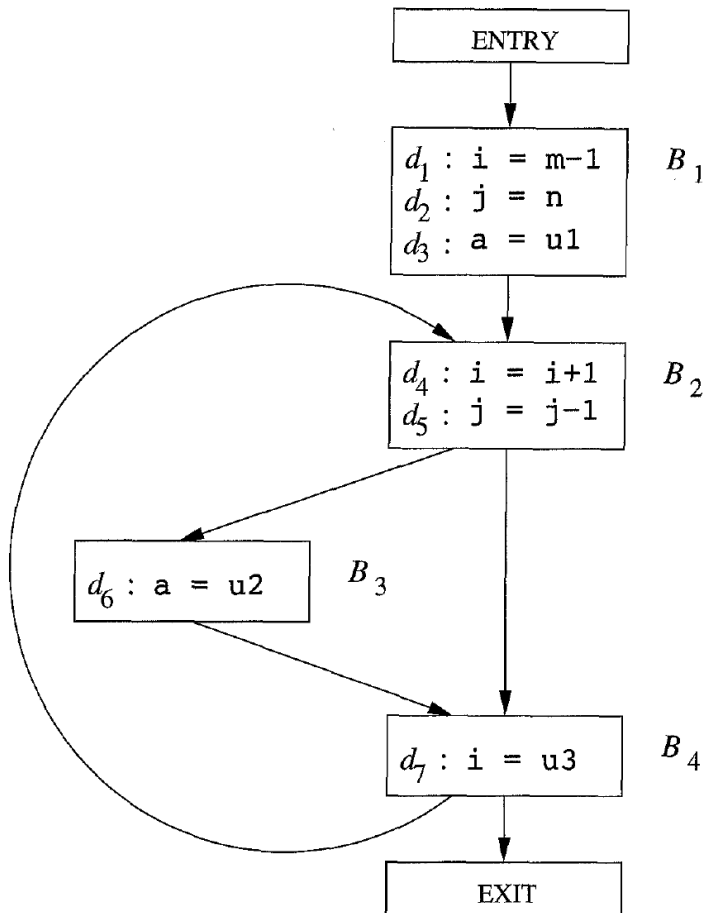
OUT[ $B_3$ ] = { $d_4, d_5, d_6$ }

IN[ $B_4$ ] = { $d_3, d_4, d_5, d_6$ }

OUT[ $B_4$ ] = { $d_3, d_5, d_6, d_7$ }

IN[EXIT] = { $d_3, d_5, d_6, d_7$ }

OUT[EXIT] = { $d_3, d_5, d_6, d_7$ }



$\text{gen}_{B_1} = \{ d_1, d_2, d_3 \}$

$\text{kill}_{B_1} = \{ d_4, d_5, d_6, d_7 \}$

$\text{gen}_{B_2} = \{ d_4, d_5 \}$

$\text{kill}_{B_2} = \{ d_1, d_2, d_7 \}$

$\text{gen}_{B_3} = \{ d_6 \}$

$\text{kill}_{B_3} = \{ d_3 \}$

$\text{gen}_{B_4} = \{ d_7 \}$

$\text{kill}_{B_4} = \{ d_1, d_4 \}$



# Example

```

1) OUT[ENTRY] =  $\emptyset$ ;
2) for (each basic block  $B$  other than ENTRY) OUT[ $B$ ] =  $\emptyset$ ;
3) while (changes to any OUT occur)
4)   for (each basic block  $B$  other than ENTRY) {
5)     IN[ $B$ ] =  $\bigcup_{P \text{ a predecessor of } B} \text{OUT}[P]$ ;
6)     OUT[ $B$ ] =  $\text{gen}_B \cup (\text{IN}[B] - \text{kill}_B)$ ;
   }

```

Iteration #3: fixed point!

OUT[ENTRY] =  $\emptyset$

IN[ $B_1$ ] =  $\emptyset$

OUT[ $B_1$ ] = { $d_1, d_2, d_3$ }

IN[ $B_2$ ] = { $d_1, d_2, d_3, d_5, d_6, d_7$ }

OUT[ $B_2$ ] = { $d_3, d_4, d_5, d_6$ }

IN[ $B_3$ ] = { $d_3, d_4, d_5, d_6$ }

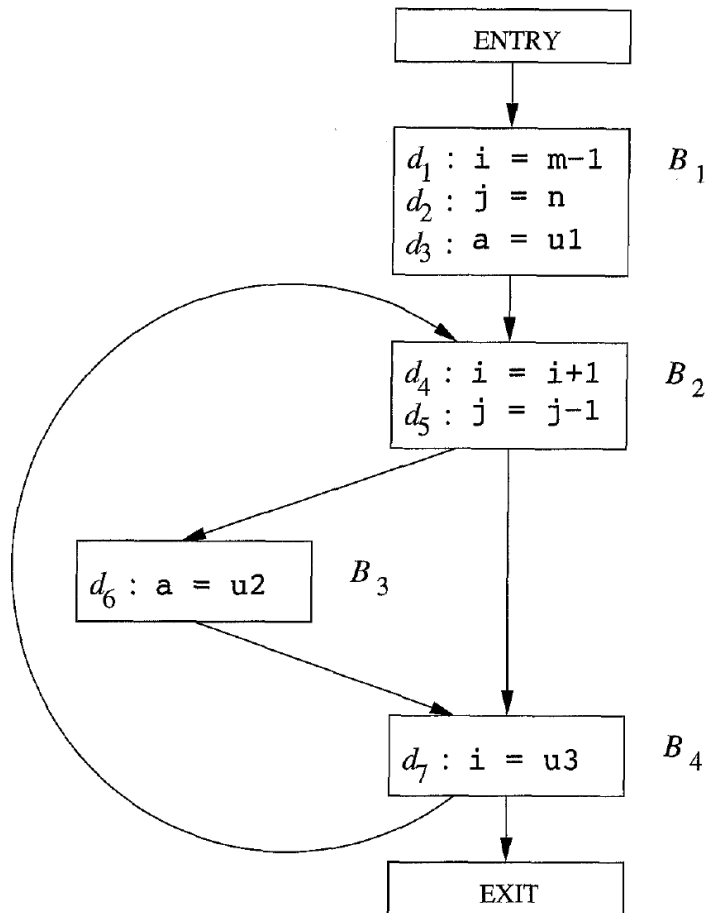
OUT[ $B_3$ ] = { $d_4, d_5, d_6$ }

IN[ $B_4$ ] = { $d_3, d_4, d_5, d_6$ }

OUT[ $B_4$ ] = { $d_3, d_5, d_6, d_7$ }

IN[EXIT] = { $d_3, d_5, d_6, d_7$ }

OUT[EXIT] = { $d_3, d_5, d_6, d_7$ }



$\text{gen}_{B_1} = \{ d_1, d_2, d_3 \}$

$\text{kill}_{B_1} = \{ d_4, d_5, d_6, d_7 \}$

$\text{gen}_{B_2} = \{ d_4, d_5 \}$

$\text{kill}_{B_2} = \{ d_1, d_2, d_7 \}$

$\text{gen}_{B_3} = \{ d_6 \}$

$\text{kill}_{B_3} = \{ d_3 \}$

$\text{gen}_{B_4} = \{ d_7 \}$

$\text{kill}_{B_4} = \{ d_1, d_4 \}$

# Outline

- Data-Flow Analysis
- Classic Data-Flow Problems {
  - Reaching Definitions
  - Live Variables\*
  - Available Expressions\*

\* Optional self-study materials

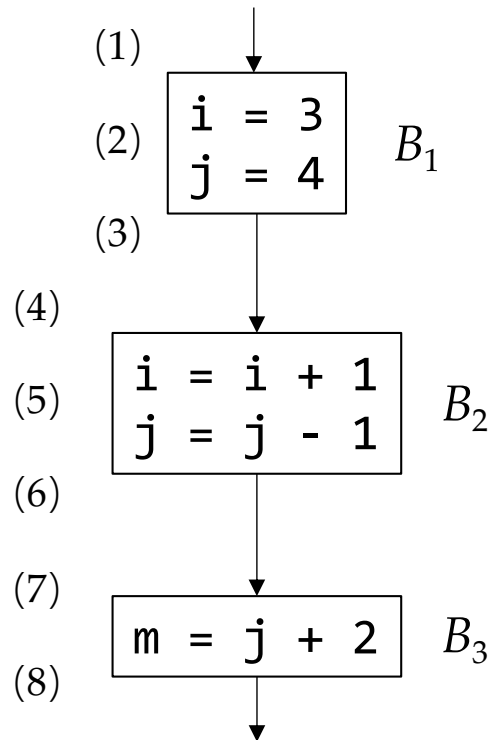
# Live-Variable Analysis

- Aim to know for variable  $x$  and point  $p$  whether the value of  $x$  at  $p$  could be used on some path starting at  $p$ 
  - If yes, we say  $x$  is *live* at  $p$ ; otherwise,  $x$  is *dead* at  $p$ .
- An important use of live-variable information is register allocation for basic blocks
  - It is not necessary to store a value computed in a register to memory at the end of a basic block
  - If all registers are used and we need another register, we should favor using a register with a dead value
- Another important use is to eliminate dead code
  - If a value is dead at a point, there is no need to compute it

# Example

*Live*: value could be used on subsequent path

*Dead*: value not referenced on subsequent path

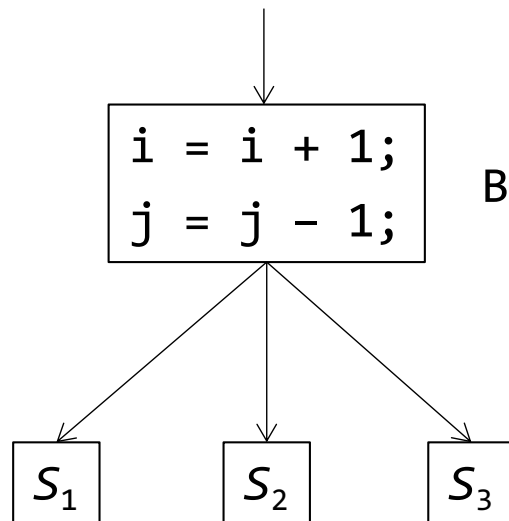


For the partial flow graph, we have:

Point	i	j	m
(1)	Dead	Dead	Dead
(2)	Live	Dead	Dead
(3)	Live	Live	Dead
(4)	Live	Live	Dead
(5)	Dead	Live	Dead
(6)	Dead	Live	Dead
(7)	Dead	Live	Dead
(8)	Dead	Dead	Dead

**Obsevation:** Live-variable analysis is a **backward data-flow problem** since the subsquent defs/uses determine the liveness of a variable at a point

# Data-Flow Equations



**Data-flow values:**

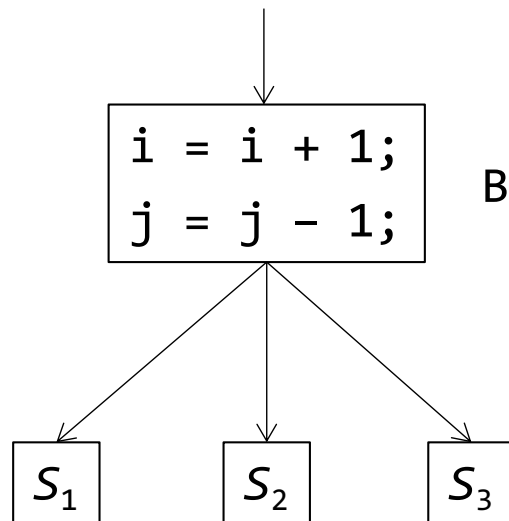
**IN**[ $B$ ]: variables live at the entry of the block  $B$

**OUT**[ $B$ ]: variables live at the exit of block  $B$

$$\text{OUT}[B] = \bigcup_{S \text{ is a successor of } B} \text{IN}[S]$$

**Rationale:** For safety, as long as a variable is live at the entry of any successor block of  $B$ , we consider the variable live at the exit of  $B$

# Data-Flow Equations



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

***use<sub>B</sub>***: the set of variables whose values are used in *B* **prior to any definition** of the variable (**i**, **j**)

***def<sub>B</sub>***: the set of variables defined in *B* **prior to any use** of that variable in *B* (i.e., the first time the variable appears, it appears as a definition; *def<sub>B</sub>* is  $\emptyset$  for the example)

# Computing $use_B$ and $def_B$

- Transfer function for a statement:
  - $s: x = y + z$
  - $use_s = \{y, z\}$
  - $def_s = \{x\} - \{y, z\}$  //  $y, z$  might be  $x$
- Transfer function for a block of statements  $s_1, s_2, \dots, s_n$ :
  - $use_B = use_1 \cup (use_2 - def_1) \cup (use_3 - def_1 - def_2) \cup \dots \cup (use_n - def_1 - def_1 - \dots - def_{n-1})$
  - $def_B = def_1 \cup (def_2 - use_1) \cup (def_3 - use_1 - use_2) \cup \dots \cup (def_n - use_1 - use_1 - \dots - use_{n-1})$

# Problem Formulation

- All variables are dead at the exit of a program
  - $IN[EXIT] = \emptyset$
- For all basic blocks  $B$  other than EXIT:
  - $IN[B] = use_B \cup (OUT[B] - def_B)$
  - $OUT[B] = \bigcup_{S \text{ is a successor of } B} IN[S]$
- Comparing to reaching definitions analysis:
  - Both sets of equations have union as the meet operator
  - Information flow for liveness travels “backward”



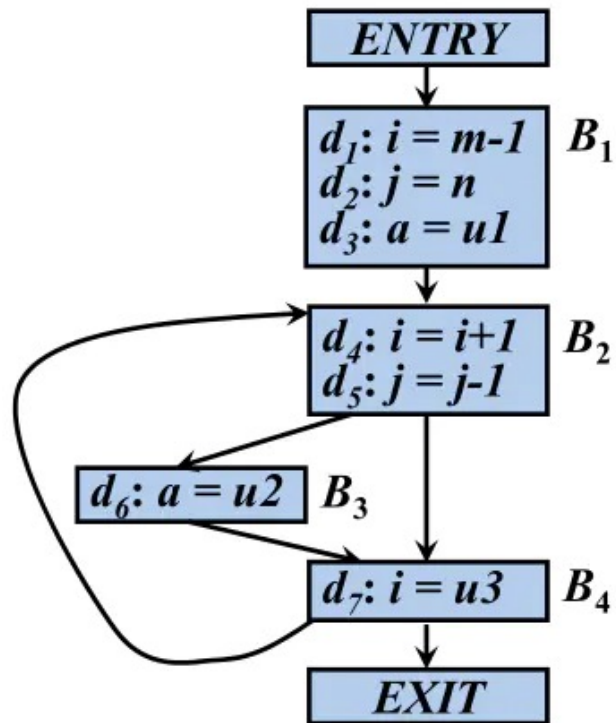
# The Iterative Algorithm

- **Input:** A flow graph with  $def_B$  and  $use_B$  computed for each block  $B$
- **Output:**  $IN[B]$  and  $OUT[B]$ , the variables live at the entry and exit of each block  $B$

```
IN[EXIT] =  $\emptyset$ ;  
for (each basic block  $B$  other than EXIT)  $IN[B] = \emptyset$ ;  
while (changes to any IN occur)  
    for (each basic block  $B$  other than EXIT) {  
         $OUT[B] = \bigcup_{S \text{ a successor of } B} IN[S]$ ;  
         $IN[B] = use_B \cup (OUT[B] - def_B)$ ;  
    }
```

Similar to the iterative algorithm for reaching definitions analysis,  $IN[B]$  never shrinks and its size is bounded, so the algorithm will eventually halt.

# Example



$use_{B1} = \{ m, n, u1 \}$

$def_{B1} = \{ i, j, a \}$

$use_{B2} = \{ i, j \}$

$def_{B2} = \Phi$

$use_{B3} = \{ u2 \}$

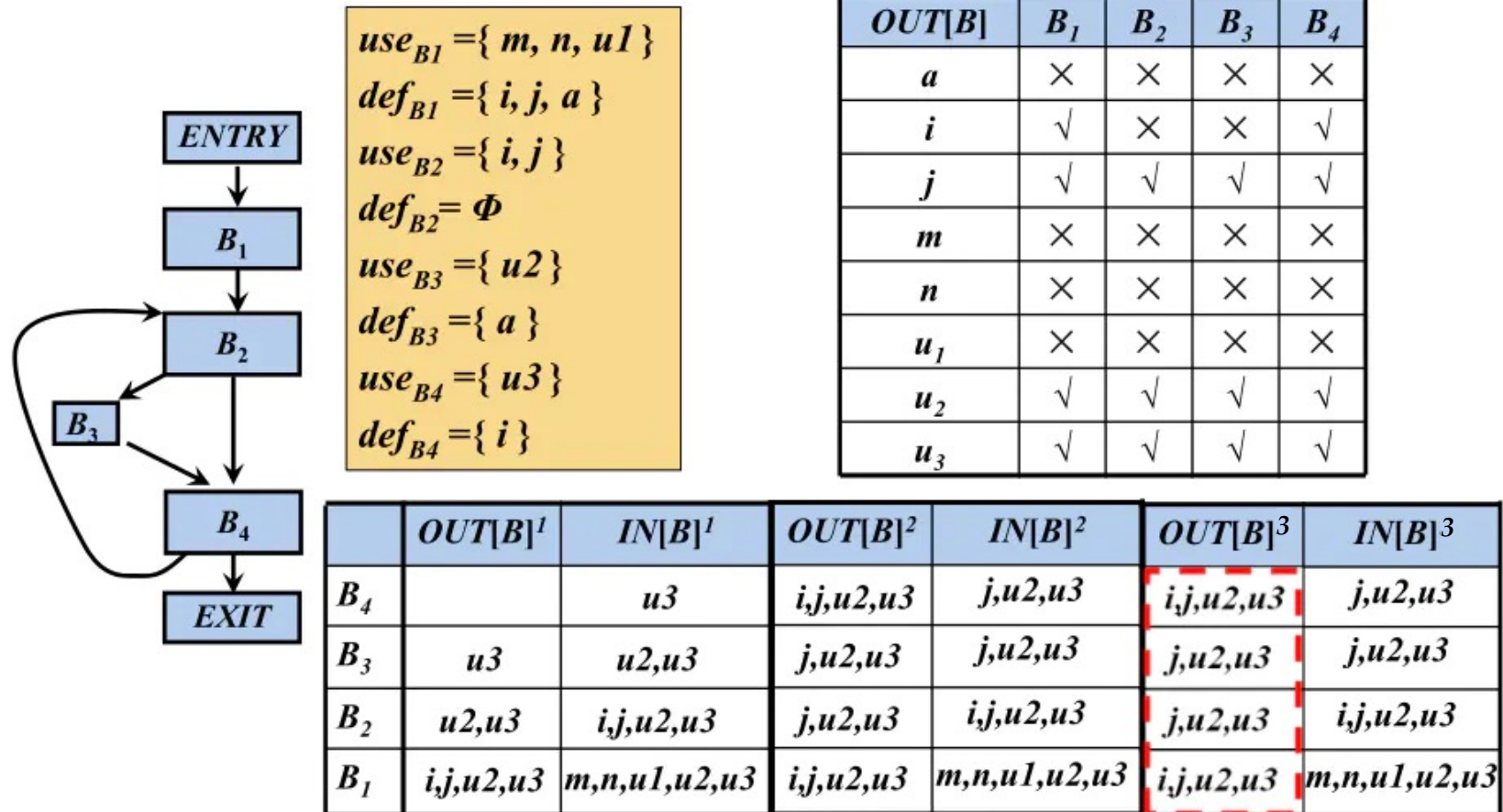
$def_{B3} = \{ a \}$

$use_{B4} = \{ u3 \}$

$def_{B4} = \{ i \}$

<https://www.jianshu.com/p/ebc1c72b881c>

# Example



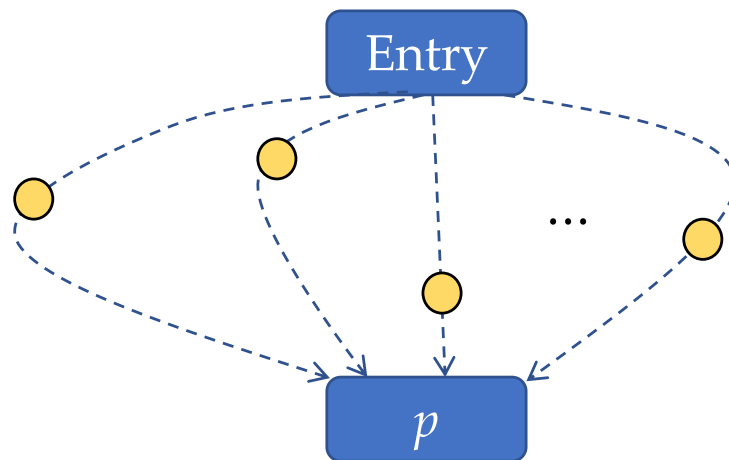
# Outline

- Data-Flow Analysis
- Classic Data-Flow Problems {
  - Reaching Definitions
  - Live Variables\*
  - Available Expressions\*

\* Optional self-study materials

# Available Expressions

- An expression  $x + y$  is available at a point  $p$  if **every path from the entry node to  $p$  evaluates  $x + y$** , and after the last such evaluation prior to reaching  $p$ , there are **no subsequent assignments to  $x$  or  $y$**

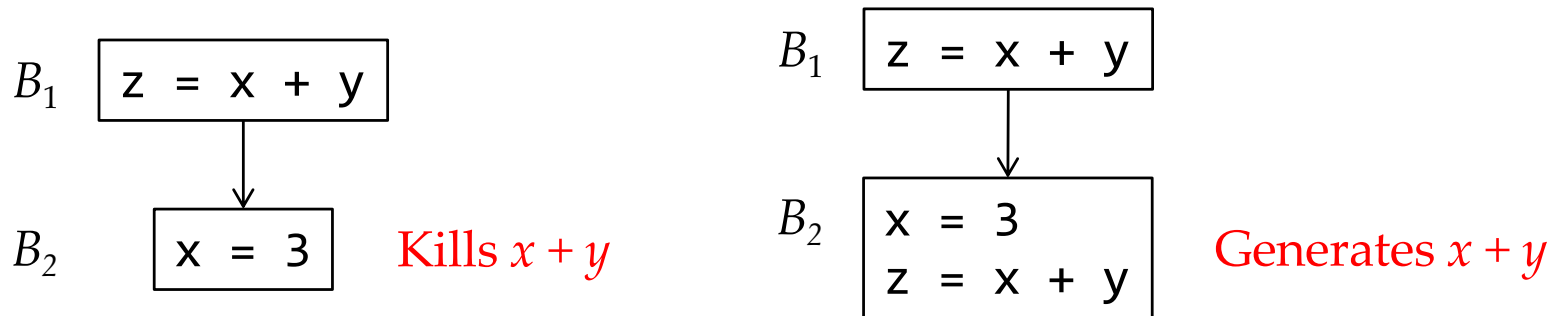


● Point evaluating  $x + y$

No assignments to  $x$  or  $y$  afterwards

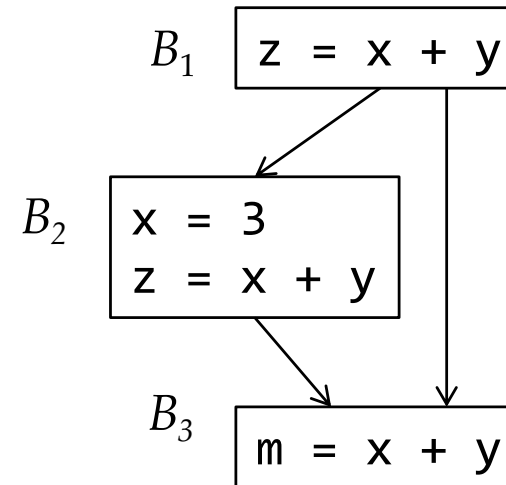
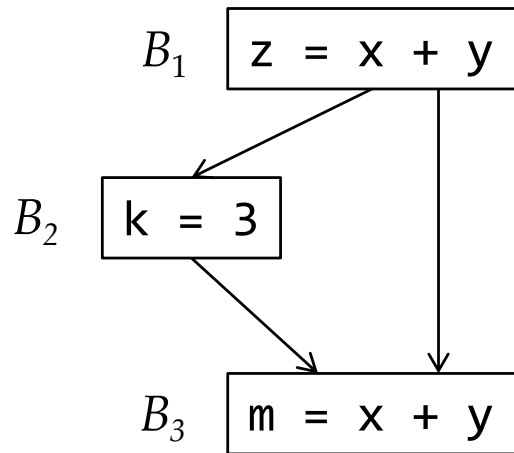
# Killing / Generating Expressions

- A block *kills* expression  $x + y$  if:
  - It assigns  $x$  or  $y$  and does not subsequently recompute  $x + y$
- A block *generates* expression  $x + y$  if:
  - It evaluates  $x + y$  and does not subsequently define  $x$  or  $y$



# The Primary Use of Available Expressions

- Detecting global common expressions



**Case 1:**  $B_2$  does not redefine  $x$  or  $y$

**Case 2:**  $B_2$  redefines  $x$  and recomputes  $x + y$

**In both cases,  $x + y$  needs not to be recomputed in  $B_3$**

# Computing Generated Expressions

$S$ : the set of expressions available before the statement



```
x = y + z
```

**Compute the expressions available after the statement:**

- Add to  $S$  the expression of  $y + z$
- Delete from  $S$  any expression involving variable  $x$

If we repeat the steps for each statement in a block  $B$ , we will get the set of **generated expressions** for  $B$  after we process the last statement.



# Computing Killed Expressions

- For a basic block  $B$ , the set of *killed expressions* is all expressions, say  $y + z$ , such that:
  - $y$  or  $z$  is defined in  $B$
  - AND  $y + z$  is not generated by  $B$

# Example

Statement	Available Expressions
	$\emptyset$
$a = b + c$	
	$\{b + c\}$
$b = a - d$	
	$\{a - d\}$
$c = b + c$	
	$\{a - d\}$
$d = a - d$	
	$\emptyset$

# Data-Flow Equations

- At the exit of the entry node, no expressions are available:
  - $\text{OUT}[\text{ENTRY}] = \emptyset$
- For all basic blocks  $B$  other than ENTRY:
  - $\text{OUT}[B] = e\_genB \cup (\text{IN}[B] - e\_killB)$
  - $\text{IN}[B] = \bigcap_{P \text{ a predecessor of } B} \text{OUT}[P]$

The equations are similar to those for reaching definition analysis except that the meet operator is intersection rather than union.\*

\* This is for safety reasons. Although producing a subset of the exact set of available expressions prevents certain optimizations, this does not lead to the change of program semantics

# The Iterative Algorithm

- **Input:** A flow graph with  $e\_kill_B$  and  $e\_gen_B$  computed for each block  $B$
- **Output:**  $IN[B]$  and  $OUT[B]$ , the expressions available at the entry and exit of each block  $B$

$U$ : the universal set of all expressions

```
OUT[ENTRY] =  $\emptyset$ ;  
for (each basic block  $B$  other than ENTRY) OUT[B] =  $U$ ;  
while (changes to any OUT occur)  
    for (each basic block  $B$  other than ENTRY) {  
         $IN[B] = \bigcap_{P \text{ a predecessor of } B} OUT[P]$ ;  
         $OUT[B] = e\_gen_B \cup (IN[B] - e\_kill_B)$ ;  
    }
```

# Summary of The Problems

	Reaching Definitions	Live Variables	Available Expressions
<b>Domain</b>	Sets of definitions	Sets of variables	Sets of expressions
<b>Direction</b>	Forwards	Backwards	Forwards
<b>Transfer function</b>	$gen_B \cup (x - kill_B)$	$use_B \cup (x - def_B)$	$e\_gen_B \cup (x - e\_kill_B)$
<b>Boundary</b>	$OUT[ENTRY] = \emptyset$	$IN[EXIT] = \emptyset$	$OUT[ENTRY] = \emptyset$
<b>Meet (<math>\wedge</math>)</b>	$\cup$	$\cup$	$\cap$
<b>Equations</b>	$OUT[B] = f_B(IN[B])$ $IN[B] = \bigwedge_{P, pred(B)} OUT[P]$	$IN[B] = f_B(OUT[B])$ $OUT[B] = \bigwedge_{S, succ(B)} IN[S]$	$OUT[B] = f_B(IN[B])$ $IN[B] = \bigwedge_{P, pred(B)} OUT[P]$
<b>Intialize</b>	$OUT[B] = \emptyset$	$IN[B] = \emptyset$	$OUT[B] = U$
<b>Application</b>	Constant folding	Dead code elimination	Global common subexpression elimination

# Reading Tasks

- Chapter 9 of the dragon book
  - 9.1 The Principal Sources of Optimization (9.1.1-9.1.7)
  - 9.2 Introduction to Data-Flow Analysis (9.2.1-9.2.6)