

Dataflow Analysis: Dataflow Frameworks

Outline of Today's Class

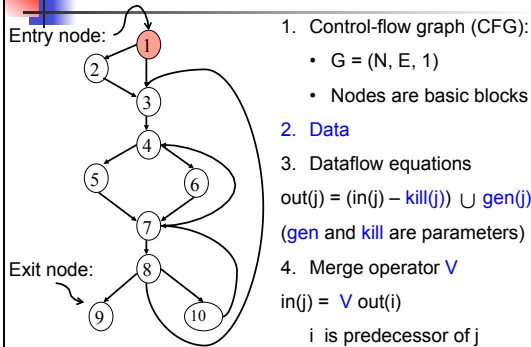
- Catch up
- Dataflow frameworks
- Lattices
- Transfer functions
- Worklist algorithm

- Reading:
 - Dragon Book, Chapter 9.2 and 9.3

Spring 19 CSCI 4450/6450, A Milanova

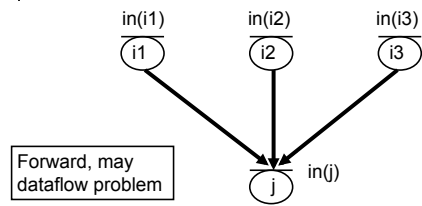
2

Dataflow Analysis



3

Problem 1: Reaching Definitions



What are the primitive dataflow facts?
 Definitions, e.g., $(x, 1)$, $(y, 6)$
 Equations act on *sets* of definitions.

Spring 19 CSCI 4450/6450, A Milanova

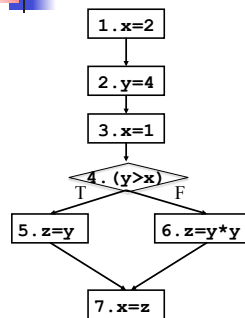
4

Problem 2. Live Uses of Variables (*Live*)

- We say that a variable x is “live on exit from node j ” if there is a live use of x on exit from j (recall the definition of “live use of x on exit from j ”)
- Problem statement: for each node n , compute the set of variables that may be live on exit from n .

1. $x=2$; 2. $y=4$; 3. $x=1$; if $(y>x)$ then 5. $z=y$; else 6. $z=y*y$; 7. $x=z$;
 What variables are live on exit from statement 3? Statement 1?

Live Example



6

Live Uses of Variables (*Live*)

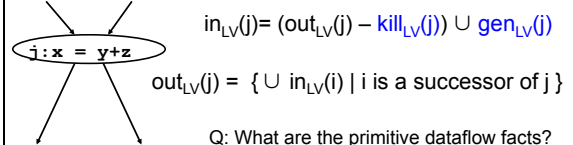
- Data
 - Primitive facts: variables x
 - Propagates sets: $\{x, y, z\}$
- Dataflow equations. At j : $x = y+z$
 - $kill_{LV}(j)$: $\{x\}$
 - $gen_{LV}(j)$: $\{y, z\}$
- Merge operator: set union \cup

Spring 19 CSCI 4450/6450, A Milanova

7

Live Uses of Variables (*Live*)

- Problem statement: for each node n , compute the set of variables that may be live on exit from n .



$$in_{LV}(j) = (out_{LV}(j) - kill_{LV}(j)) \cup gen_{LV}(j)$$

$$out_{LV}(j) = \{ \cup in_{LV}(i) \mid i \text{ is a successor of } j \}$$

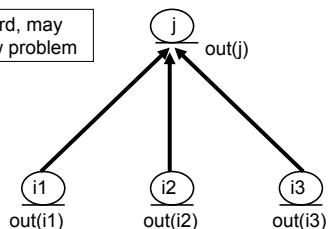
- Q: What are the primitive dataflow facts?
- Q: What is $gen_{LV}(j)$?
- Q: What is $kill_{LV}(j)$?

Spring 19 CSCI 4450/6450, A Milanova

8

Problem 2: Live Uses of Variables

Backward, may dataflow problem

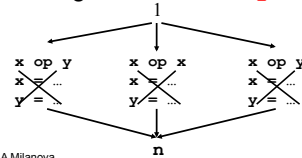


What are the primitive dataflow facts?
Variables, e.g., x, y, z . Equations act on sets of variables.

9

Problem 3: Available Expressions (*Avail*)

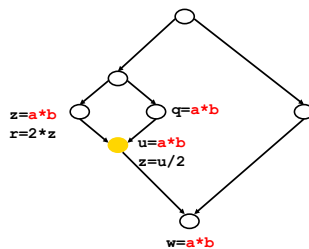
- An expression $x \text{ op } y$ is **available** at program point n if **every** path from entry to n evaluates $x \text{ op } y$, and after every evaluation prior to reaching n , there are NO subsequent assignments to x or y



Spring 19 CSCI 4450/6450, A Milanova

10

Avail Enables Global Common Subexpressions

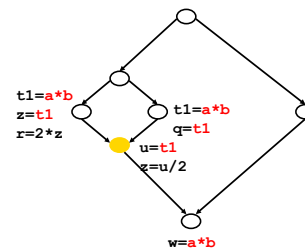


Spring 19 CSCI 4450/6450, A Milanova

11

Avail Enables Global Common Subexpressions

Can we eliminate $w=a*b$?



Spring 19 CSCI 4450/6450, A Milanova

12

Available Expressions (*Avail*)

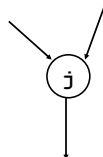
- Data?
 - Primitive dataflow facts are expressions, e.g., $x+y$, $a*b$, $a+2$
 - Analysis propagates sets of expressions, e.g., $\{x+y, a*b\}$
- Dataflow equations at j : $x = y \text{ op } z$?
 - $out_{AE}(j) = (in_{AE}(j) - kill_{AE}(j)) \cup gen_{AE}(j)$
 - $kill_{AE}(j)$: all expressions with operand x : $(x \text{ op } _)$, $(_ \text{ op } x)$
 - $gen_{AE}(j)$: new expression: $\{(y \text{ op } z)\}$

13

Available Expressions (*Avail*)

- Merge operator?
 - For *Avail*, it is set intersection \cap

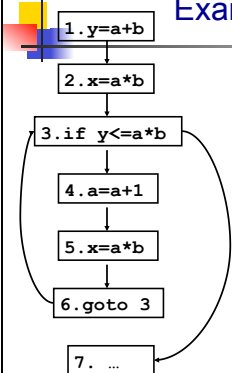
$$in_{AE}(j) = \{\cap out_{AE}(i) \mid i \text{ is predecessor of } j\}$$



Spring 19 CSCI 4450/6450, A Milanova

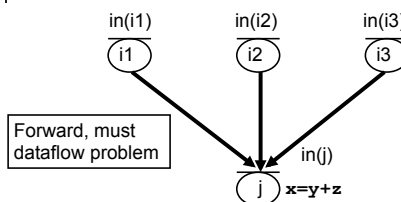
14

Example



15

Problem 3: Available Expressions



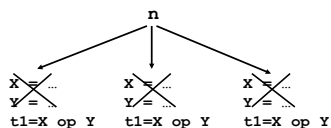
What are the primitive dataflow facts?
Expressions, e.g., $x+y$, $a*b$. Equations act on sets of expressions.

Spring 19 CSCI 4450/6450, A Milanova

16

Problem 4: Very Busy Expressions (*VeryB*)

- An expression $x \text{ op } y$ is *very busy* at node n , if along EVERY path from n to the end of the program, we come to a computation of $x \text{ op } y$ BEFORE any redefinition of x or y .



Spring 19 CSCI 4450/6450, A Milanova

17

Very Busy Expressions (*VeryB*)

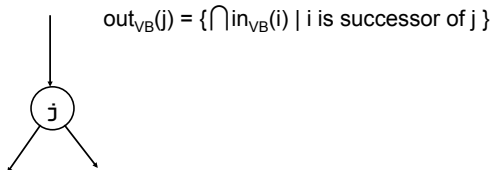
- Data?
 - Primitive dataflow facts are expressions, e.g., $x+y$, $a*b$
 - Analysis propagates sets of expressions, e.g., $\{x+y, a*b\}$
- Dataflow equations at j : $x = y \text{ op } z$?
 - $in_{VB}(j) = (out_{VB}(j) - kill_{VB}(j)) \cup gen_{VB}(j)$
 - $kill_{VB}(j)$: all expressions with operand x : $(x \text{ op } _)$, $(_ \text{ op } x)$
 - $gen_{VB}(j)$: new expression: $\{(y \text{ op } z)\}$

18

Very Busy Expressions (VeryB)

- Merge operator?

- For *VeryB*, it is set intersection \cap

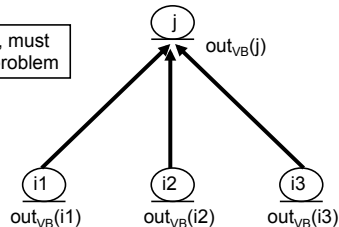


Spring 19 CSCI 4450/6450, A Milanova

19

Very Busy Expressions

Backward, must dataflow problem



Spring 19 CSCI 4450/6450, A Milanova

20

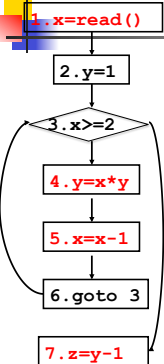
Another Example: Taint Analysis

- A definition (x, k) is **tainted** if k is designated as a taint **source**, or (x, k) is computed based on an operand that is tainted.
- Problem statement: for each node n , compute the set of tainted definitions that may reach n .

Spring 19 CSCI 4450/6450, A Milanova

21

Example: Taint Analysis (explicit flow)



22

Outline of Today's Class

- Catch up
- **Dataflow frameworks**
- Lattice
- Transfer functions
- Worklist algorithm
- Reading:
 - Dragon Book, Chapter 9.2 and 9.3

Spring 19 CSCI 4450/6450, A Milanova

23

Dataflow Problems

	May Problems	Must Problems
Forward Problems	Reaching Definitions	Available Expressions
Backward Problems	Live Uses of Variables	Very Busy Expressions

Spring 19 CSCI 4450/6450, A Milanova

24

Similarities

- Analyses operate over similar **property spaces**
- In all cases, analysis operates over a **finite set D** of primitive dataflow facts
 - Reach: D** is the set of all definitions in the program:
e.g., $\{(x, 1), (y, 2), (x, 4), (y, 5)\}$
 - Avail and VeryB: D** is the set of all arithmetic expressions:
e.g., $\{a+b, a*b, a+1\}$
 - Live: D** is the set of all variables
e.g., $\{x, y, z\}$
- Solution at node **n** is a subset of **D** (e.g., a definition either reaches **n** or it does not reach **n**)

Spring 19 CSCI 4450/6450, A Milanova

25

Similarities

- Dataflow equations have the same form (from now on, we'll focus on forward problems):

$$\text{out}(j) = (\text{in}(j) - \text{kill}(j)) \cup \text{gen}(j) = (\text{in}(j) \cap \text{pres}(j)) \cup \text{gen}(j)$$

$$\text{in}(j) = \{ \forall \text{out}(i) \mid i \text{ is predecessor of } j \}$$

pres(j) is the complement of **kill(j)**

- A note: what makes the 4 classical problems special is that sets **kill(j)/pres(j)** and **gen(j)** do not depend on in(i)
- Thus, set union and set intersection can be implemented as logical OR and AND respectively

Spring 19 CSCI 4450/6450, A Milanova

26

Similarities

- The dataflow equation at node **j** is a **transfer functions**. It take **in(j)** as argument and produces **out(j)** as result:

$$\text{out}(j) = f_j(\text{in}(j))$$

Spring 19 CSCI 4450/6450, A Milanova

27

Dataflow Frameworks

- We generalize and study the properties of the **property space**
 - Property space is a **lattice**
 - Choice settles **merge operator**
- We generalize and study the properties of the **transfer function space**
 - Functions are **monotone or distributive**
- We generalize and study the properties of the **worklist algorithm** that computes a solution

Spring 19 CSCI 4450/6450, A Milanova

28

Lattice Theory

- Partial ordering** (denoted by \leq or \sqsubseteq)
 - Relation between pairs of elements
 - Reflexive $a \leq a$
 - Anti-symmetric $a \leq b$ and $b \leq a \implies a = b$
 - Transitive $a \leq b$ and $b \leq c \implies a \leq c$
- Partially ordered set** (poset) (set **S**, \leq)
 - 0** Element $0 \leq a$, for every a in **S**
 - 1** Element $a \leq 1$, for every a in **S**

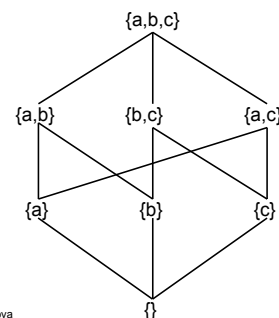
We don't necessarily need 0 and 1 element.

Spring 19 CSCI 4450/6450, A Milanova

29

Poset Example

$D = \{a, b, c\}$
The poset is 2^D , \leq is set inclusion



Spring 19 CSCI 4450/6450, A Milanova

30

Lattice Theory

- Greatest lower bound (glb)
 - l_1, l_2 in poset S , a in poset S is the $\text{glb}(l_1, l_2)$ iff
 - 1) $a \leq l_1$ and $a \leq l_2$
 - 2) for any b in S , $b \leq l_1$, $b \leq l_2$ implies $b \leq a$
 If glb exists, it is unique. Why? Called *meet* (denoted by \wedge or \sqcap) of l_1 and l_2 .
- Least upper bound (lub)
 - l_1, l_2 in poset S , c in poset S is the $\text{lub}(l_1, l_2)$ iff
 - 1) $c \geq l_1$ and $c \geq l_2$
 - 2) for any d in S , $d \geq l_1$, $d \geq l_2$ implies $d \geq c$
 If lub exists, it is unique. Called *join* (denoted by \vee or \sqcup) of l_1 and l_2 .

31

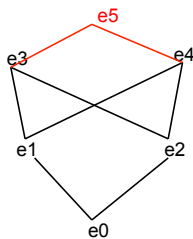
Definition of a Lattice (L, \wedge, \vee)

- A lattice L is a poset under \leq , such that every pair of elements has a **glb (meet)** and **lub (join)**
- A lattice need not contain a 0 or 1 element
- A finite lattice must contain 0 and 1 elements
- Not every poset is a lattice
- If there is element a such that $a \leq x$ for every x in L , then a is the 0 element of L
- If there is a such that $x \leq a$ for every x in L , then a is the 1 element of L

Spring 19 CSCI 4450/6450, A Milanova

32

A Poset but Not a Lattice



There is no $\text{lub}(e_3, e_4)$ in this poset so it is not a lattice.

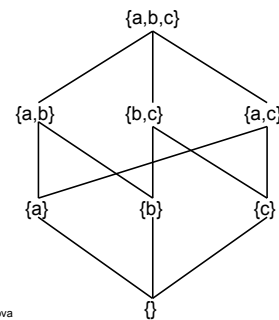
Suppose we add the $\text{lub}(e_3, e_4)$, is it a lattice?

Spring 19 CSCI 4450/6450, A Milanova

33

Is This Poset a Lattice

$D = \{a, b, c\}$
The poset is 2^D , \leq is set inclusion



Spring 19 CSCI 4450/6450, A Milanova

34

Examples of Lattices

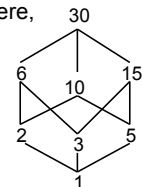
- $H = (2^D, \cap, \cup)$ where D is a finite set
 - $\text{glb}(s_1, s_2)$ denoted $s_1 \wedge s_2$, is set intersection $s_1 \cap s_2$
 - $\text{lub}(s_1, s_2)$ denoted $s_1 \vee s_2$, is set union $s_1 \cup s_2$
- $J = (\mathbb{N}_1, \text{gcd}, \text{lcm})$
 - Partial order is integer divide on \mathbb{N}_1
 - $\text{lub}(n_1, n_2)$ denoted $n_1 \vee n_2$ is $\text{lcm}(n_1, n_2)$
 - $\text{glb}(n_1, n_2)$ denoted $n_1 \wedge n_2$ is $\text{gcd}(n_1, n_2)$
 (\mathbb{N}_1 denotes natural numbers starting at 1)

Spring 19 CSCI 4450/6450, A Milanova

35

Chain

- A poset C where for every pair of elements c_1, c_2 in C , either $c_1 \leq c_2$ or $c_2 \leq c_1$.
 - E.g., $\{\} \leq \{a\} \leq \{a, b\} \leq \{a, b, c\}$
 - E.g., from the lattice J as shown here,
 - $1 \leq 2 \leq 6 \leq 30$
 - $1 \leq 3 \leq 15 \leq 30$
- A lattice s.t. every ascending chain is finite, is said to satisfy the *Ascending Chain Condition*



Spring 19 CSCI 4450/6450, A Milanova

36

Lattices in Dataflow Analysis

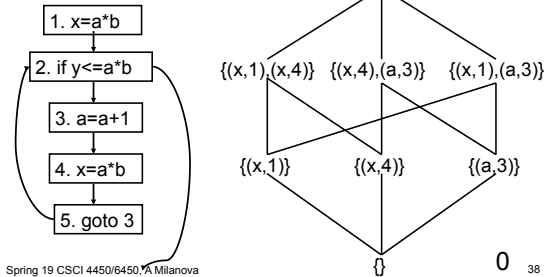
- Lattices define property space
- Lattices entail properties of the standard dataflow analysis solution procedure (the **worklist algorithm**, which we will study shortly)

Spring 19 CSCI 4450/6450, A. Milanova

37

Dataflow Lattices: *Reach*

$D =$ all definitions: $\{(x,1),(x,4),(a,3)\}$ $\{(x,1),(x,4),(a,3)\}$ **1**
Poset is 2^D , \leq is the subset relation \subseteq

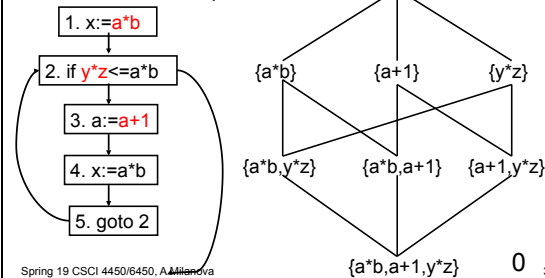


Spring 19 CSCI 4450/6450, A. Milanova

0 38

Dataflow Lattices: *Avail*

$D =$ all expressions: $\{a*b, a+1, y*z\}$ $\{\}$ **1**
Poset is 2^D , \leq is the superset relation \supseteq



Spring 19 CSCI 4450/6450, A. Milanova

39

Dataflow Frameworks

- Equations:

$$\text{in}(j) = V \text{ out}(i) \quad \text{out}(j) = f_j(\text{in}(j))$$

$i \in \text{pred}(j)$

where:

- $\text{in}(j)$, $\text{out}(j)$ are elements of a **property space**
- f_j is the **transfer function** associated with node j
- V is the **merge operator**

Spring 19 CSCI 4450/6450, A. Milanova

40

Dataflow Frameworks (cont.)

- **The property space must be:**
 1. A lattice L , \leq
 2. L satisfies the *Ascending Chain Condition*
Requires that all ascending chains are finite
- **The merge operator V must be the join of L**
- In dataflow, L is often the lattice of the subsets over a finite set of dataflow facts D
 - Choose universal set D (e.g., all definitions)
 - Choose ordering operation \leq . Since the merge operator is must be to the join of L , a *may* problem entails that \leq is **subset**. Conversely, a *must* problem entails that \leq is **superset**

41

Example: *Reach* Lattice

- Property space is the lattice of the subsets where
 - D is the set of all definitions in the program
 - \leq is the **subset** operation
 - **Join** is set union \cup , as needed for *Reach*, which is a *may* problem
 - Lattice has **0** being $\{\}$, and **1** being D
 - Lattice satisfies the *Ascending Chain Condition*

Spring 19 CSCI 4450/6450, A. Milanova

42

Reach Lattice

$D =$ all definitions: $\{(x,1),(x,4),(a,3)\}$ $\{(x,1),(x,4),(a,3)\}$ 1
 Poset is 2^D , \leq is the subset relation \subseteq

Spring 19 CSCI 4450/6450, A. Milanova 43

Example: Avail Lattice

- Property space is the lattice of the subsets where
 - D is the set of all expressions in the program
 - \leq is **superset**
 - join of the lattice is set intersection, as needed for Avail, which is a *must* problem
 - Lattice has **0** being D , and **1** being $\{\}$
 - Lattice satisfies *Ascending Chain Condition*

Spring 19 CSCI 4450/6450, A. Milanova 44

Dataflow Lattices: Avail

$D =$ all expressions: $\{a*b, a+1, y*z\}$
 Poset is 2^D , \leq is the superset relation \supseteq

Spring 19 CSCI 4450/6450, A. Milanova 45

Transfer Functions

- The transfer functions: $f_j: L \rightarrow L$.** Formally, function space F is such that
 - F contains all f_j ,
 - F contains the identity function $\text{id}(x) = x$
 - F is closed under composition.
 - Each f_j is **monotone**

Spring 19 CSCI 4450/6450, A. Milanova 46

Monotonicity

- $F: L \rightarrow L$ is **monotone** if and only if:
 - a, b in L , f in F then $a \leq b \implies f(a) \leq f(b)$ or (equivalently):
 - x, y in L , f in F then $f(x) \vee f(y) \leq f(x \vee y)$
- Theorem: Definitions (1) and (2) are equivalent.
 - Show that (1) implies (2)
 - Show that (2) implies (1)

Spring 19 CSCI 4450/6450, A. Milanova 47

Distributivity

- $F: L \rightarrow L$ is **distributive** if and only if x, y in L , f in F then $f(x \vee y) = f(x) \vee f(y)$
- Every distributive function is also monotone but not the other way around
- Distributivity is a very nice property!

Spring 19 CSCI 4450/6450, A. Milanova 48

Monotonicity and Distributivity

- Is classical *Reach* distributive?

- Yes
- To show distributivity:

$$\text{For each } j \ ((\text{in}(j) \cup \text{in}'(j)) \cap \text{pres}(j)) \cup \text{gen}(j) = ((\text{in}(j) \cap \text{pres}(j)) \cup \text{gen}(j)) \cup ((\text{in}'(j) \cap \text{pres}(j)) \cup \text{gen}(j))$$

$$\begin{aligned} & ((\text{in}(j) \cup \text{in}'(j)) \cap \text{pres}(j)) \cup \text{gen}(j) = \\ & ((\text{in}(j) \cap \text{pres}(j)) \cup (\text{in}'(j) \cap \text{pres}(j))) \cup \text{gen}(j) = \\ & ((\text{in}(j) \cap \text{pres}(j)) \cup \text{gen}(j)) \cup ((\text{in}'(j) \cap \text{pres}(j)) \cup \text{gen}(j)) \end{aligned}$$

Monotone Dataflow Frameworks

- A problem fits into the dataflow framework if
 - its property space is a lattice \mathbf{L}, \leq that satisfies the *Ascending Chain Condition*
 - its merge operator \mathbf{V} is the join of \mathbf{L} and
 - its function space $\mathbf{F}: \mathbf{L} \rightarrow \mathbf{L}$ is monotone
- Thus, we can make use of a generic solution procedure, known as the **worklist algorithm** or the **maximal fixpoint algorithm** or the **fixpoint iteration algorithm**

Worklist Algorithm for Forward Dataflow Problems

```
/* Initialize to initial values; 1 is entry node of CFG */
in(1) = InitialValue;          in_Reach(1) = UNDEF (or {})
for m = 2 to n do in(m) = 0    in_Reach(m) = {}
W = {1,2,...,n} /* put every node on the worklist */
```

```
while W ≠ ∅ do {
  remove j from W
  out(j) = f_j(in(j))          out_Reach(j) = in_Reach(j) ∩ pres(j) ∪ gen(j)
  for i in successors(j)
    if out(j) ≰ in(i) then {   if out_Reach(j) ≰ in_Reach(i)
      in(i) = out(j) ∨ in(i)   in_Reach(i) = out_Reach(j) ∪ in_Reach(i)
      W = W ∪ {i}
    }
}
```

}

Worklist Algorithm for Forward Dataflow Problems (slightly different)

```
/* Initialize to initial values; 1 is entry node of CFG */
in(1) = InitialValue; out(1) = f_1(in(1))
for m := 2 to n do in(m) = 0; out(m) = f_m(0)
W := {2,...,n} /* put every node but 1 on the worklist */
```

```
while W ≠ ∅ do {
  remove j from W
  in(j) = V { out(i) | i is predecessor of j }
  out(j) = f_j(in(j))
  if out(j) changed then
    W = W ∪ { k | k is successor of j }
}
```

Termination Argument

- Why does the algorithm terminate?
- Sketch of proof:
 - At each iteration, at least one $\text{out}(j)$ changes.
 - Since $\text{out}(j)$ in \mathbf{L} , and \mathbf{L} satisfies the *Ascending Chain Condition*, $\text{out}(j)$ changes at most $\mathbf{O}(h)$ times where h is the height of the lattice \mathbf{L} .

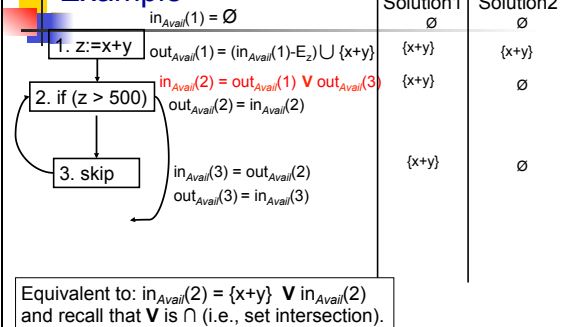
Correctness Argument

- Theorem: The worklist algorithm computes a solution that satisfies the dataflow equations
- Why?
- Sketch of proof:
 - Whenever j is processed, algorithm sets $\text{out}(j) = f_j(\text{in}(j))$. Whenever $\text{out}(j)$ changes, algorithm puts successors on the list, so $\text{in}(j) = \mathbf{V} \{ \text{out}(i) \}$.
 - So final solution will satisfy equations.

Precision Argument

- Theorem: The algorithm computes the **least solution** of the dataflow equations.
 - Historically though, this solution is often called the **maximal fixpoint solution (MFP)**
 - i.e., For every node j , the worklist algorithm computes a solution $\mathbf{MFP}(j) = \{in(j), out(j)\}$, such that every other solution $\{in'(j), out'(j)\}$ of the dataflow equations is $in(j) \leq in'(j)$, $out(j) \leq out'(j)$

Example



Many Applications!

- Static debugging
 - Memory errors in C/C++ programs
 - Memory leaks
 - Null pointer dereferences
 - Array-out-of-bound accesses
 - Concurrency errors in shared-memory apps
 - Data-races, atomicity violations, deadlocks
- Information flow (as known as taint analysis)

Many Applications!

- White-box testing: compute coverage
 - Control-flow-based testing
 - Data-flow-based testing
 - Intuitively, test each def-use chain
- Regression testing
 - Analyze changes and select regression tests that actually test changed code

Dataflow Analysis

- Classical technique
- Compared to Hoare logic, it captures state in a more coarse way
- Still relevant, many interesting problems are phrased in dataflow terms

Next Class

- MOP vs MFP solutions
- Two classical non-distributive dataflow analyses:
 - Constant propagation and
 - Points-to analysis