



Dataflow Analysis: Non-distributive Analysis



Outline of Today's Class

- Dataflow frameworks, conclusion
 - Lattices (last time)
 - Transfer functions (last time)
 - Worklist algorithm
- MOP solution vs. MFP solution

- Non-distributive analyses
 - Constant propagation
 - Points-to analysis



Monotone Dataflow Framework

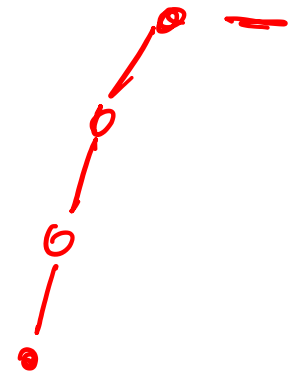
- A problem fits into the dataflow framework if
 - its property space is a lattice \mathbf{L} , \leq that satisfies the *Ascending Chain Condition*
 - its merge operator \vee is the join of \mathbf{L}and
 - its transfer function space $\mathbf{F}: \mathbf{L} \rightarrow \mathbf{L}$ is monotone
- Thus,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

Worklist Algorithm for Forward Dataflow Problems

```
/* Initialize to initial values; 1 is entry node of CFG */  
in(1) = InitialValue; out(1) = f1(in(1))  
for m = 2 to n do in(m) = 0; out(m) = fm(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) = fj(in(j))  
    if out(j) changed then  
        W = W U { k | k is successor of j }
```

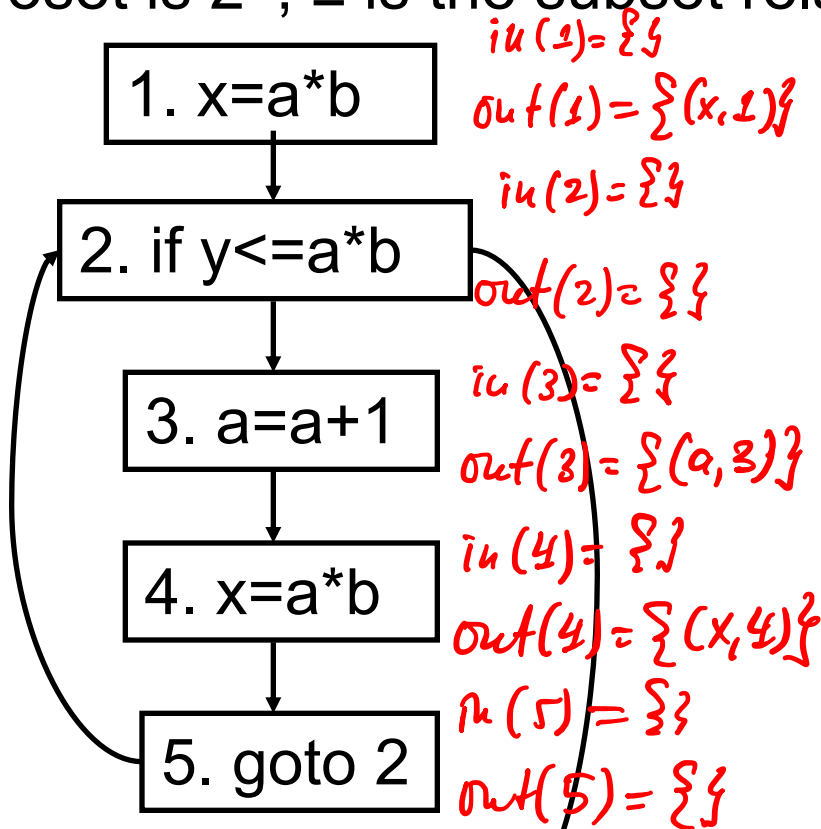
$$in^t(j) \leq in^{t+1}(j)$$
$$out^t(j) \leq out^{t+1}(j)$$



Worklist Algorithm on Reach

D = all definitions: $\{(x,1), (x,4), (a,3)\}$

Poset is 2^D , \leq is the subset relation \sqsubseteq



$$W = \{2, 3, 4, 5\}$$

Iter 1

remove 2

$$in(2) = \{(x,1)\} \quad out(2) = \{(x,4)\}$$

Iter 2

remove 3

$$in(3) = \{(x,1)\} \quad out(3) = \{(x,4), (a,3)\}$$

Iter 3

remove 4

$$in(4) = \{(x,1), (a,3)\} \quad out(4) = \{(x,4), (a,3)\}$$

Iter 4

remove 5

$$in(5) = \{(x,4), (a,3)\} \quad out(5) = \{(x,4), (a,3)\}$$

$$W = \{2\}$$



Termination Argument

- Theorem: the algorithm terminates. Why?
- Sketch of argument:

A node k is placed on worklist only if the $\text{out}(j)$ of a predecessor j changes. Monotonicity of f guarantees $\text{in}^t(j) \leq \text{in}^{t+1}(j)$ and $\text{out}^t(j) \leq \text{out}^{t+1}(j)$. (Here $\text{in}^t(j)$, $\text{out}^t(j)$ are the sets at iteration t .)

in and out sets are elements of \mathbf{L} and \mathbf{L} satisfies the *Ascending Chain Condition*; thus, there is only a finite number of times each $\text{out}(j)$ changes.



Correctness Argument

- Theorem: Worklist algorithm computes a solution that satisfies the dataflow equations. Why?
- Sketch of argument:
 - Suppose either (1) $\forall out(i) \neq in(j)$ or (2) $out(j) \neq f_j(in(j))$
 - For (1) to hold we must have “grown” $out(i)$ in some iteration and not added successor j to worklist; this is impossible.

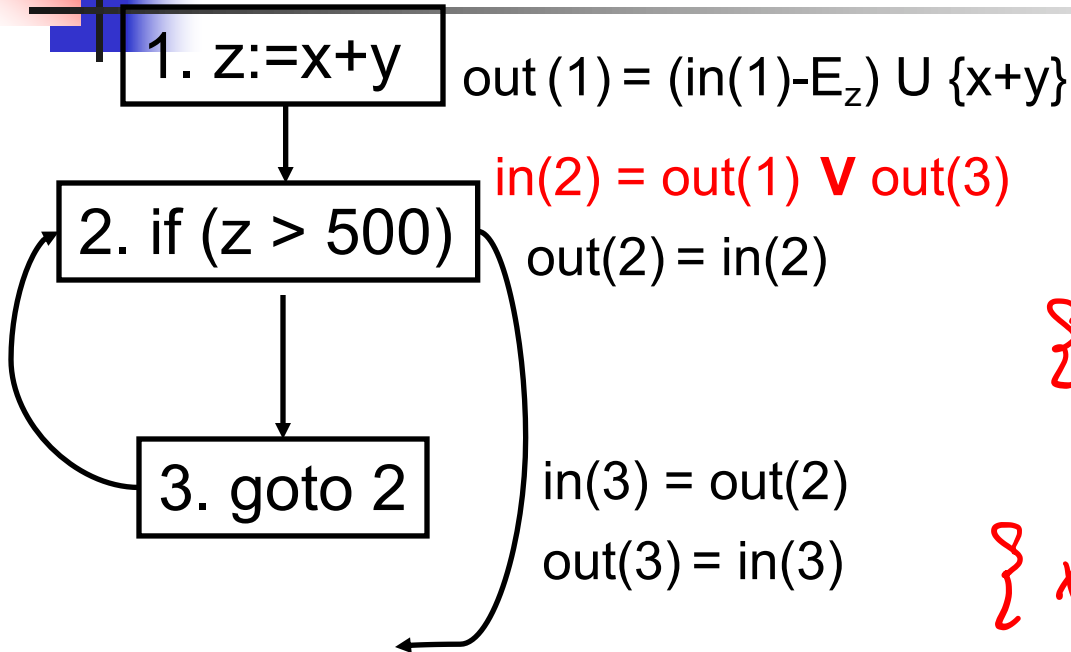


Precision Argument

- Theorem: Worklist algorithm computes the **least solution** of the dataflow equations.
 - Historically, solution computed by worklist algorithm is called the **maximal fixpoint solution (MFP solution)**
 - For every node j , worklist algorithm computes a solution $MFP(j) = (in(j), out(j))$, such that for every solution $(in'(j), out'(j))$ of the dataflow equations we have $in(j) \leq in'(j)$ and $out(j) \leq out'(j)$

Example (Avail)

$in(1) = \emptyset$



$\{ \}$
 |
 $\{ x+y \}$

~~$\{ x+y \}$~~ ~~$\cap X$~~

MFP

Solution1

Solution2

\emptyset

\emptyset

$\{x+y\}$

$\{x+y\}$

$\{x+y\}$

\emptyset

$\{x+y\}$

\emptyset

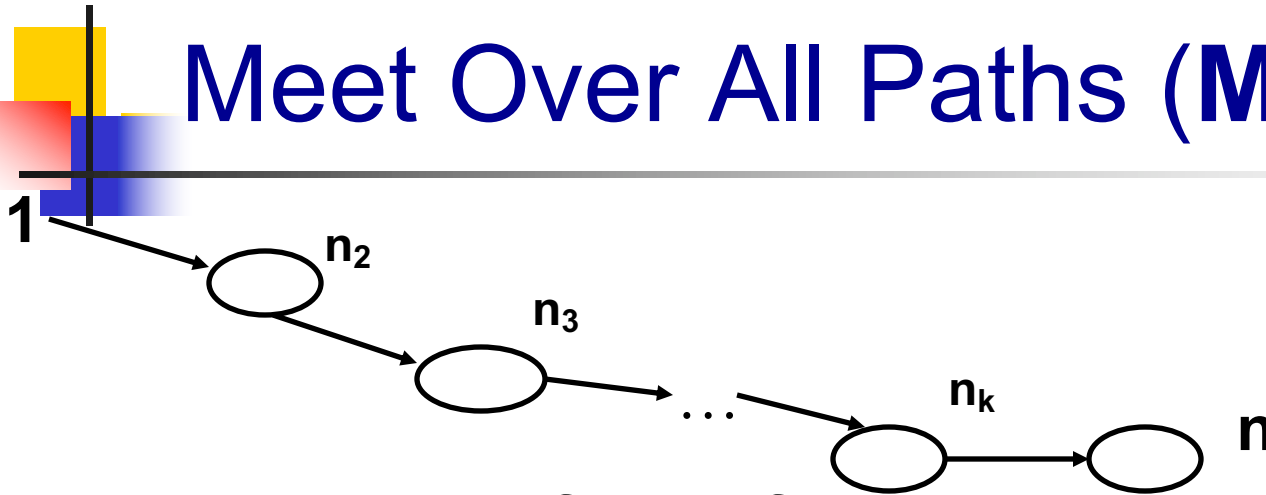
Equivalent to: $in(2) = \{x+y\} \mathbf{V} in(2)$
 and recall that \mathbf{V} is \cap (i.e., set intersection).



Outline of Today's Class

- Dataflow frameworks, conclusion
 - Lattices (last time)
 - Transfer functions (last time)
 - Worklist algorithm
- **MOP solution vs. MFP solution**
- Non-distributive analyses
 - Constant propagation
 - Points-to analysis

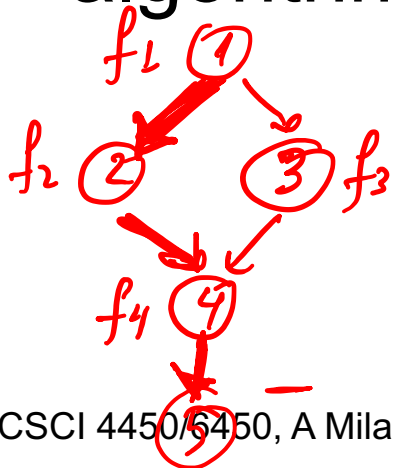
Meet Over All Paths (MOP)



- Desired dataflow information at n is obtained by traversing ALL PATHS from 1 (entry node) to n .
- For every path $p=(1, n_2, n_3 \dots, n_k)$ we compute $f_{n_k}(\dots f_{n_2}(f_1(\text{InitialValue})))$
- The MOP at *entry* of n is $\mathbf{V} f_{n_k}(\dots f_{n_2}(f_1(\text{InitialValue})))$
over all paths p from 1 to n

MOP vs. MFP

- MOP is an abstraction of the best solution computable with dataflow analysis
 - It is a common assumption in dataflow analysis that *all program paths are executable*
- MFP is the solution computed by the worklist algorithm



MOP:

$$f_4(f_2(f_1(I_{init}))) \vee f_4(f_3(f_1(I_{init})))$$

$$f_4(out(2)) \vee f_4(out(3))$$

MFP:

$$out(4) = f_4(in(4))$$

$$in(4) = out(2) \vee out(3)$$

$$out(2) = f_2(f_1(I_{init}))$$

$$out(3) = f_3(f_1(I_{init}))$$

$$f_4(out(2) \vee out(3))$$



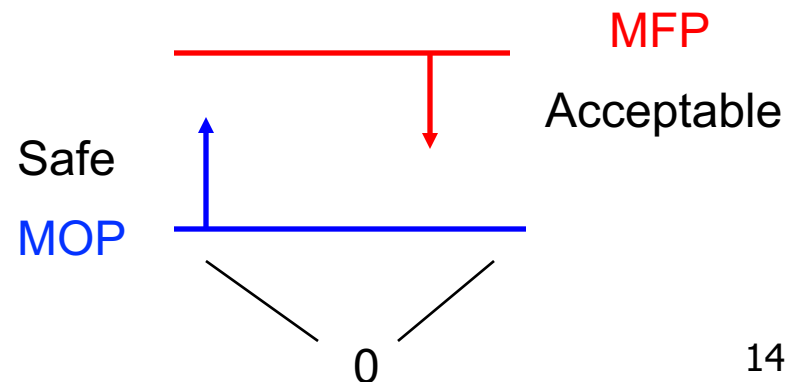
MOP vs. MFP

- For *distributive* problems $MFP = MOP$!

- Unfortunately, for *monotone* problems this is not true. But we still have a **safe** solution: it is a theorem that for monotone problems,
 $MFP \geq MOP$

Safety of a Dataflow Solution

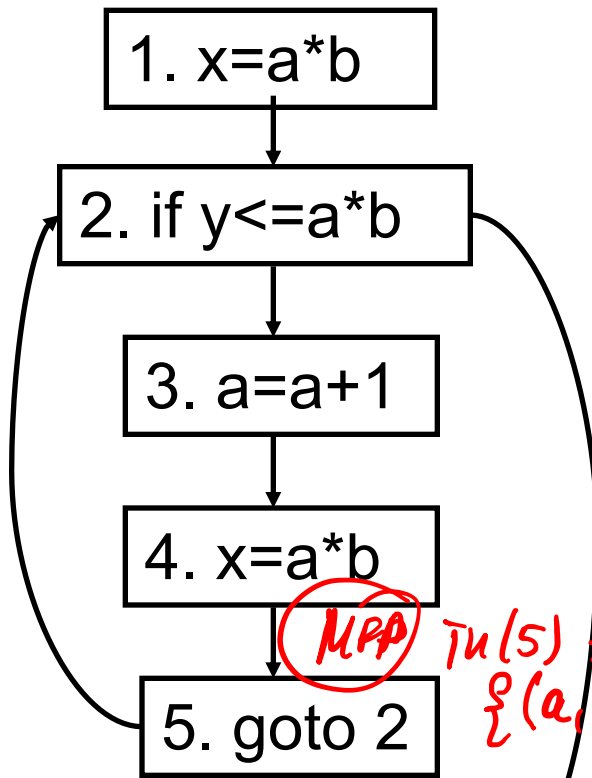
- A **safe** (also, correct or sound) solution X **overestimates** the “best” possible dataflow solution, i.e., $X \geq \text{MOP}$
- Historically, an **acceptable** solution X is one that is better than what we can do with the MFP, i.e., $X \leq \text{MFP}$



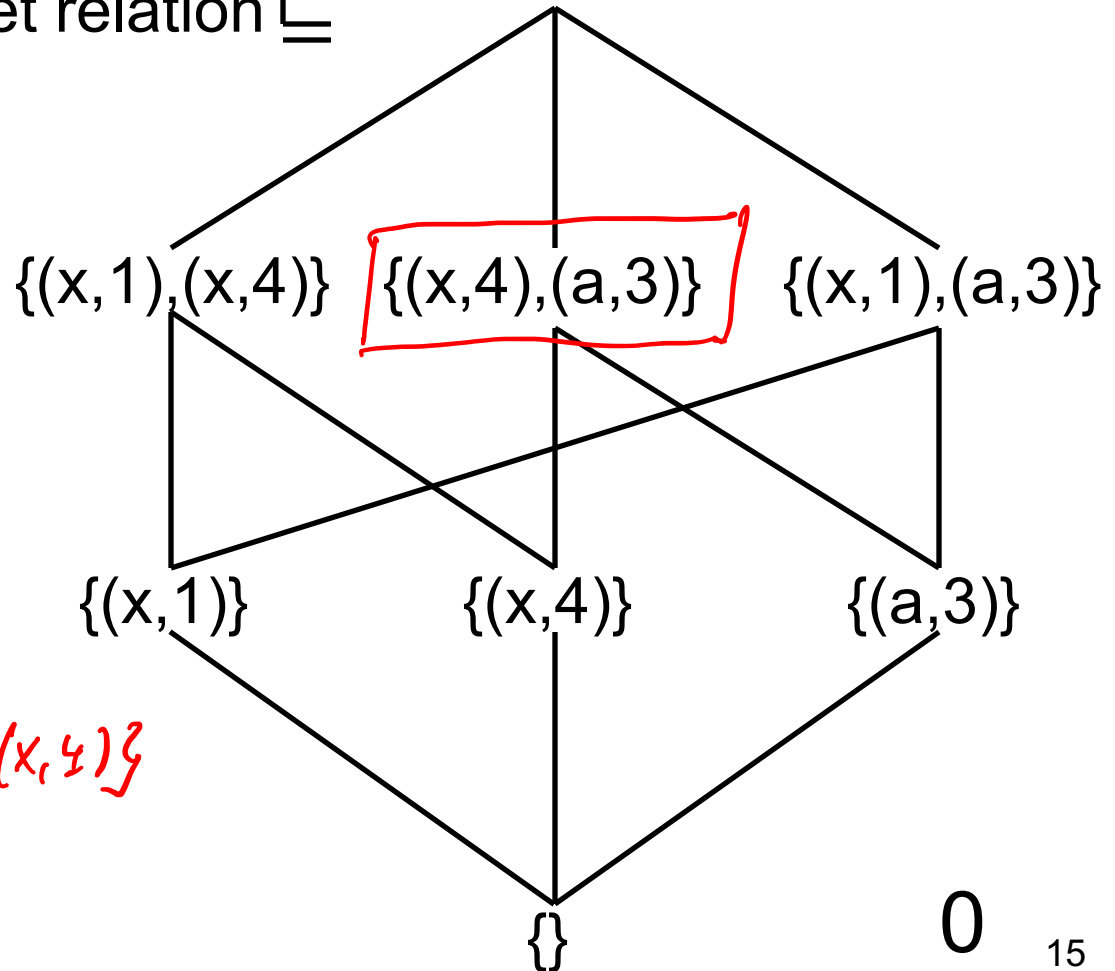
Safe Solutions: *Reach*

$U =$ all definitions: $\{(x,1),(x,4),(a,3)\}$ $\{(x,1),(x,4),(a,3)\}$ 1

Poset is 2^U , \leq is the subset relation \subseteq



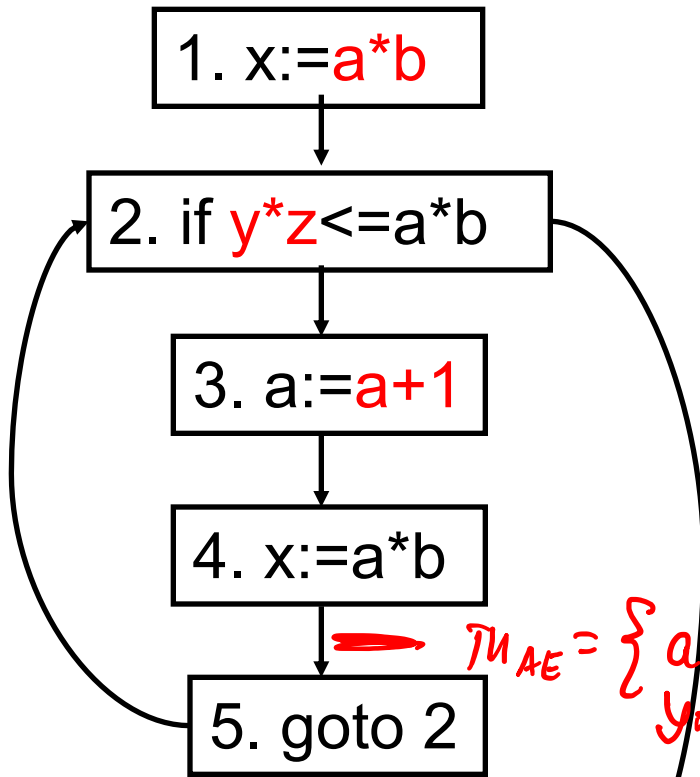
MPP in (5) = $\{(a,3), (x,4)\}$



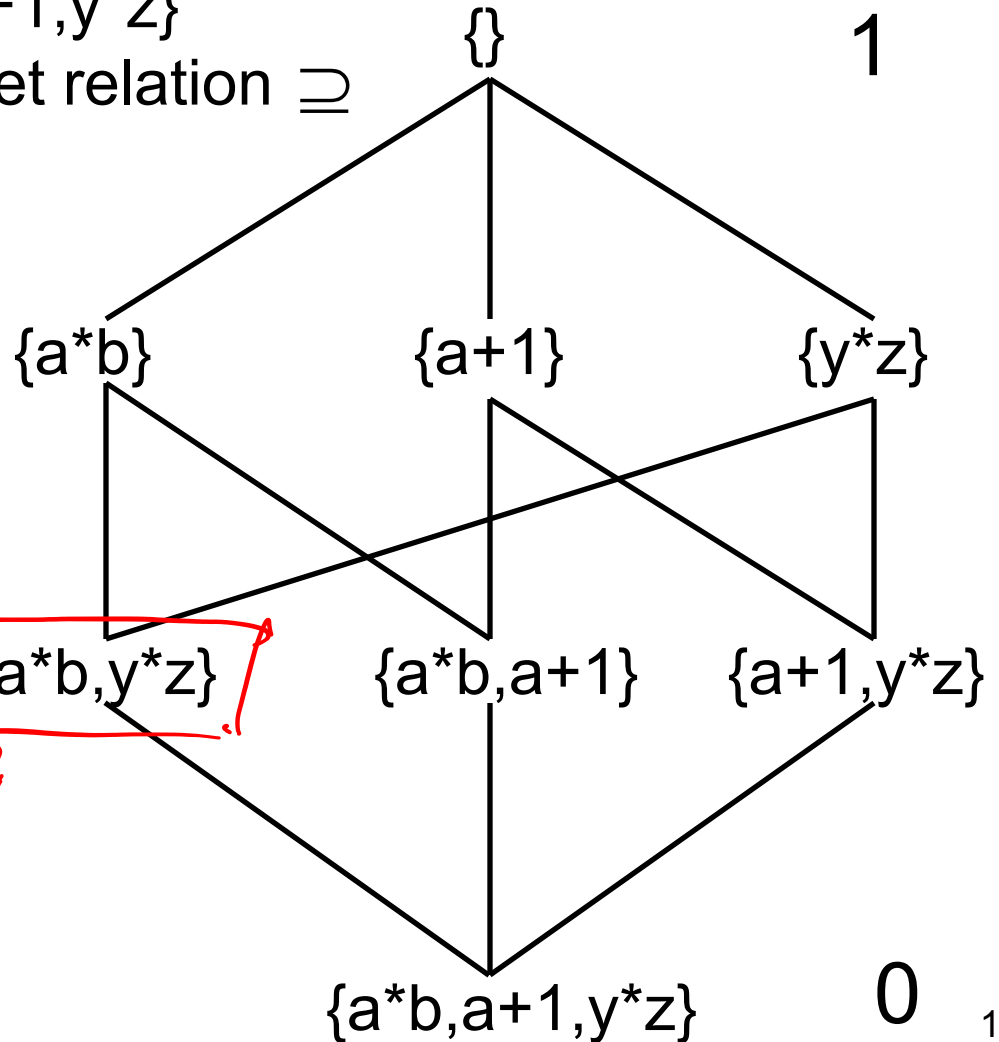
Safe Solutions: Avail

$U =$ all expressions: $\{a*b, a+1, y*z\}$

Poset is 2^U , \leq is the superset relation \supseteq



$\mathcal{M}_{AE} = \{a*b, y*z\}$





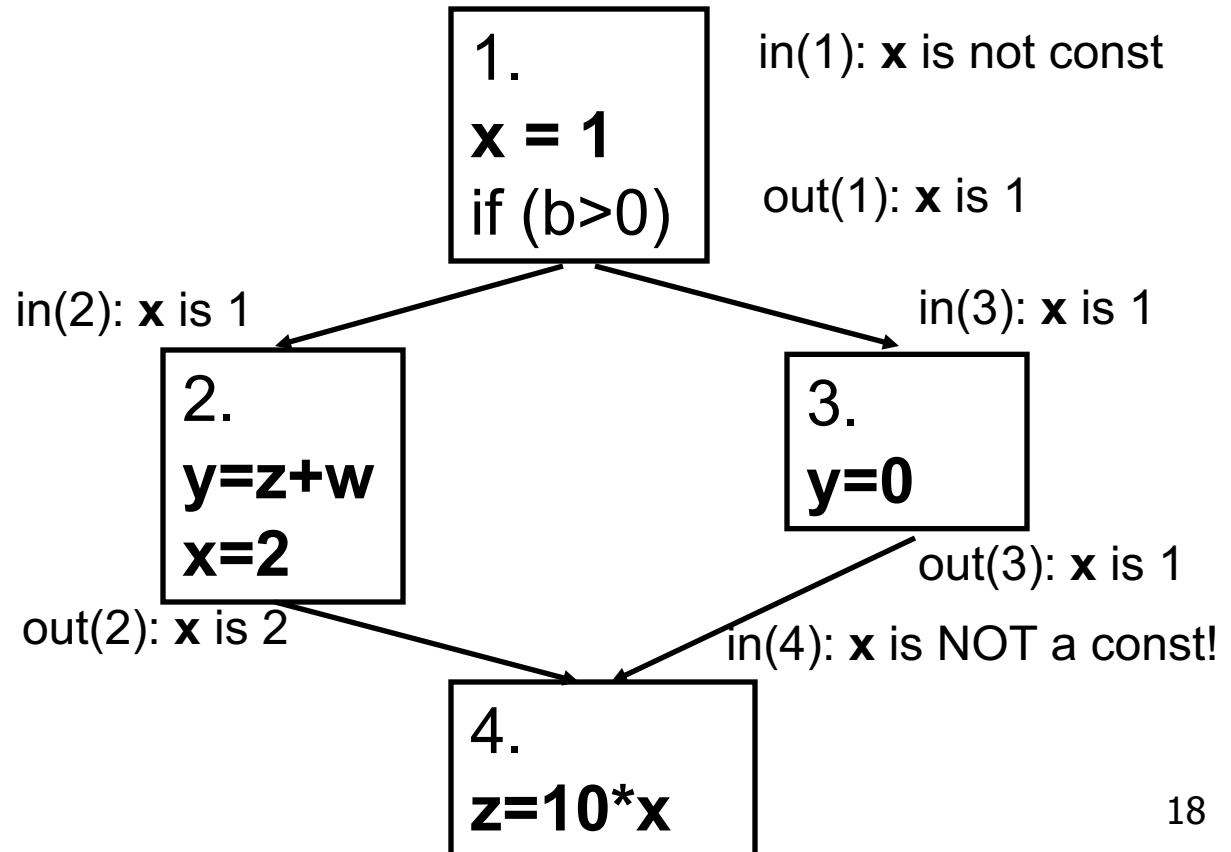
Outline of Today's Class

- Dataflow frameworks, conclusion
 - Lattices (last time)
 - Transfer functions (last time)
 - Worklist algorithm
- MOP solution vs. MFP solution
- **Non-distributive analyses**
 - Constant propagation
 - Points-to analysis

Constant Propagation (Simple)

- Problem statement: Can variable **x** hold a constant value at a given program point?

- Example:



Let's Fit Analysis into Monotone Dataflow Framework

- If property space has desired properties
 - it is a lattice \mathbf{L} , \leq that satisfies the *Ascending Chain Condition*
 - its merge operator is the join of \mathbf{L}and
- Function space $\mathbf{F}: \mathbf{L} \rightarrow \mathbf{L}$ is monotone
- Then analysis fits the monotone dataflow framework and can be solved using the worklist algorithm

Constant Propagation: Property Space

- Associate one of the following values with variable x at each program point

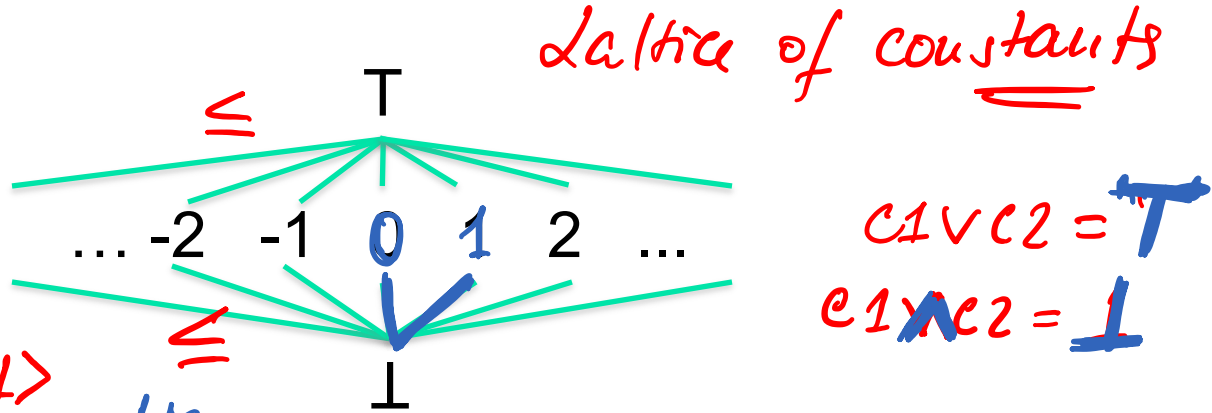
value	meaning
1 (or T)	x is NOT a constant
C	x has constant value C
0 (or \perp)	x is unknown

Constant Propagation: Lattice

- Lattice L_x, \leq

$\perp \leq \perp, 0 \leq T$

$l_1 = \langle x \rightarrow \perp, y \rightarrow T, z \rightarrow \perp \rangle$
 $l_2 = \langle x \rightarrow \perp, y \rightarrow 5, z \rightarrow 0 \rangle$



- Dataflow lattice L is the **product lattice of L_x**

- l_1, l_2 in L , $l_1 \leq l_2$ iff $l_{1x} \leq l_{2x}$ for every variable x
- $l_1 \vee l_2$ amounts to $l_{1x} \vee l_{2x}$ for every variable x
- Merge operator is join of L

- Does the product lattice satisfy the ACC?

Product Lattice

- E.g.,

$\langle x \rightarrow \perp, y \rightarrow 1, z \rightarrow T \rangle, \langle x \rightarrow 1, y \rightarrow 2, z \rightarrow 3 \rangle, \text{ etc.}$
are lattice elements

- E.g.,

$\langle x \rightarrow 1, y \rightarrow 2, z \rightarrow T \rangle \leq \langle x \rightarrow T, y \rightarrow 2, z \rightarrow T \rangle$

$l_{1x} = 1 \leq l_{2x} = T \checkmark$ $l_{1z} = T \leq l_{2z} = T \checkmark$
 $l_{1y} = 2 \leq l_{2y} = 2 \checkmark$

- E.g.,

$\langle x \rightarrow 1, y \rightarrow 3, z \rightarrow T \rangle \vee \langle x \rightarrow T, y \rightarrow 2, z \rightarrow T \rangle =$
 $\langle T, T, T \rangle$

Product Lattice

Does Product Lattice satisfy the ACC?

$$\langle x \rightarrow T, y \rightarrow T \rangle$$

$$\langle x \rightarrow T, y \rightarrow 1 \rangle$$

|

$$\langle x \rightarrow T, y \rightarrow 1 \rangle$$

|

$$\langle x \rightarrow 1, y \rightarrow 1 \rangle$$

|

$$\langle x \rightarrow 1, y \rightarrow 1 \rangle$$

Length of maximal ascending chain is 2^N where N is the number of variables in the tuple.

Constant Propagation: Transfer Functions

- $j: x = C$

$f_j: \text{kill } x \rightarrow \text{val}, \text{ generate } x \rightarrow C$

- $j: x = y$

$f_j: \text{kill } x \rightarrow \text{val}, \text{ add } x \rightarrow \text{val}', \text{ s.t. } y \rightarrow \text{val}' \text{ in } \text{in}(j). \text{val and val}' \text{ are one of}$

- \perp : bottom (unknown)
- C : constant
- T : top (not a constant)

Constant Propagation: Transfer Functions

■ $j: x = y \text{ Op } z$

$f_j: \text{kill}: x \rightarrow \text{val}$

gen:

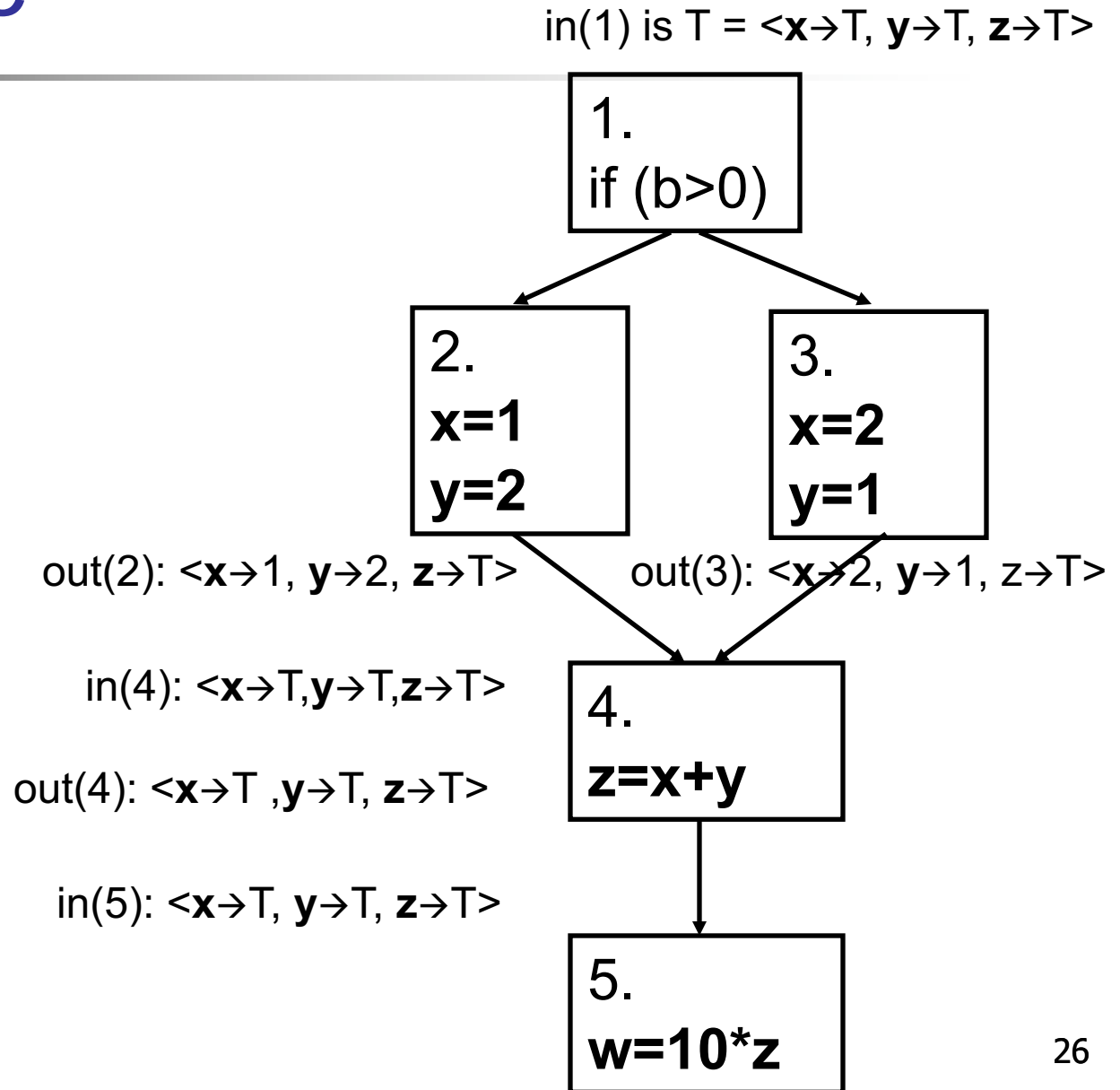
If $y \rightarrow c_1$ and $z \rightarrow c_2$ in $\text{in}(j)$, then $x \rightarrow c_1 \text{ Op } c_2$

else if $y \rightarrow T$ or $z \rightarrow T$ in $\text{in}(j)$, then $x \rightarrow T$

else $x \rightarrow \perp$

- Next, we'll argue monotonicity which would give us that Constant Propagation is solvable by the Worklist algorithm

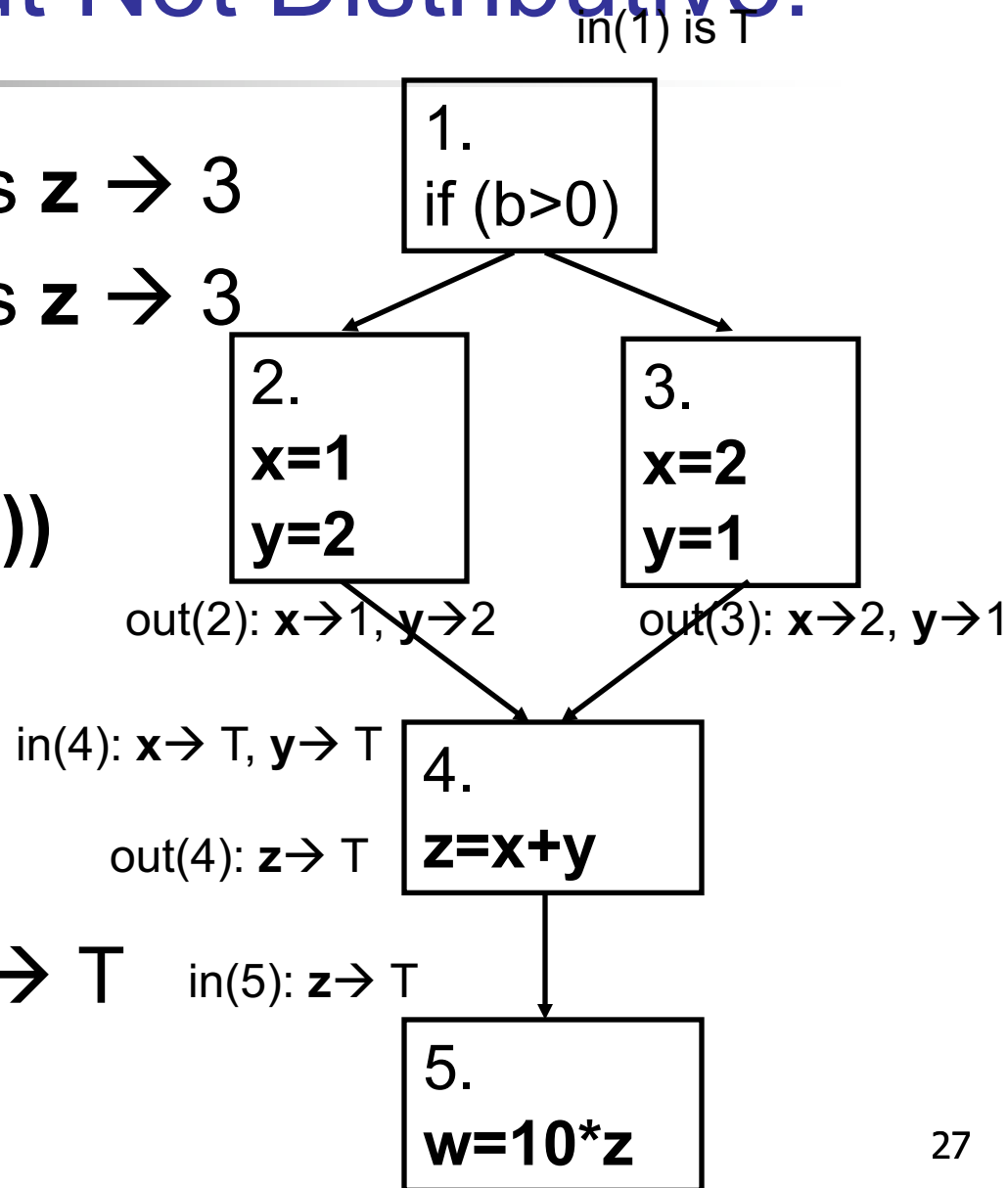
Example



Constant Propagation is Monotone but Not Distributive!

- $f_4(f_2(f_1(T)))$ computes $z \rightarrow 3$
 - $f_4(f_3(f_1(T)))$ computes $z \rightarrow 3$
 - Thus, MOP at 5
- $f_4(f_2(f_1(T))) \vee f_4(f_3(f_1(T)))$
computes $z \rightarrow 3$

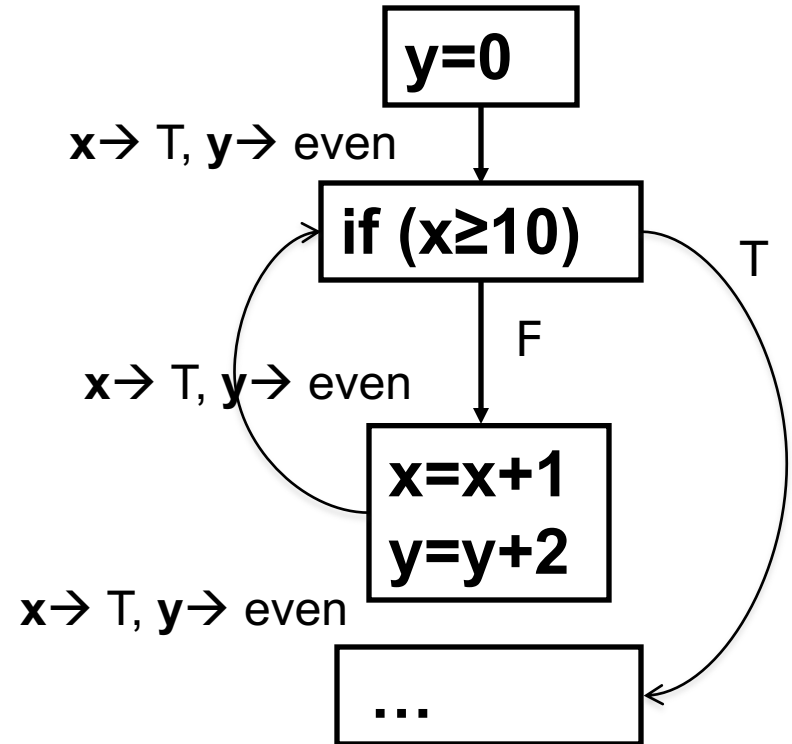
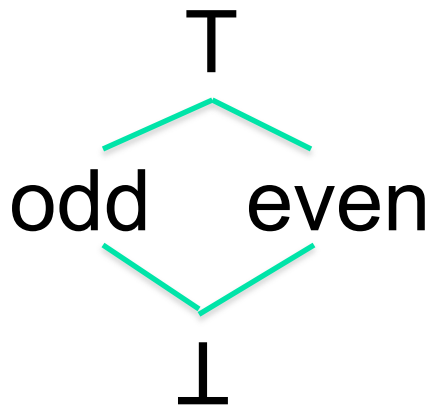
MFP at 5 computes $z \rightarrow T$
(i.e., z is NOT a const)



More Product Lattices

- Problem statement: Is integer variable x odd or even at program point n ? $x \rightarrow T, y \rightarrow T$

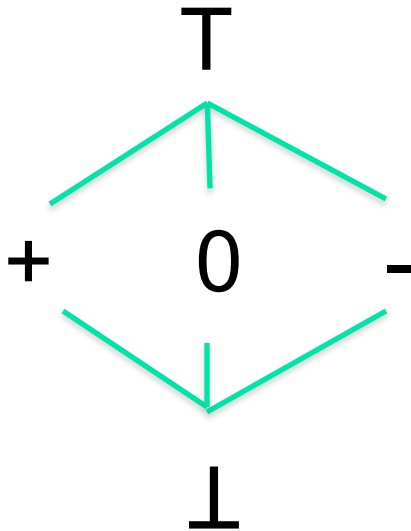
■ L_x :



More Product Lattices

- Problem statement: What **sign** does a variable hold at a given program point, i.e., is it positive, negative, or 0

■ L_x :



E.g., $\langle x \rightarrow +, y \rightarrow T, z \rightarrow 0 \rangle$



Outline of Today's Class

- Dataflow frameworks, conclusion
 - Lattices (last time)
 - Transfer functions (last time)
 - Worklist algorithm
- MOP solution vs. MFP solution

- Non-distributive analyses
 - Constant propagation
 - **Points-to analysis**



Points-to Analysis

- Problem statement: What memory locations may a pointer variable point to?
- Many applications!
 - Enables compiler optimizations
 1. $a = 1;$
 2. $*p = b;$
 3. $s = a*a;$
 1. $a = x*y*z+x;$
 2. $*p = b;$
 3. $s = x*y*z+x;$
 - Static debugging and taint analysis tools



Points-to Graph: Example

Example 1:

```
int a, b;  
int *p1, *p2;  
p1 = &a;  
p2 = p1;  
*p2 = 1;
```




Points-to Graph: Example

Example 2:

```
int a, b = 15;  
int *p1, *p2;  
int **p3;  
p3 = &p1;  
p1 = &a;  
p2 = *p3;  
*p2 = b;
```

Points-to Analysis (for a C-like language)

- Assume the following 4 simple statements

(1) address taken $p = \&q$

(2) propagation $p = q$

(3) indirect read $p = *q$

(4) indirect write (update) $*p = q$

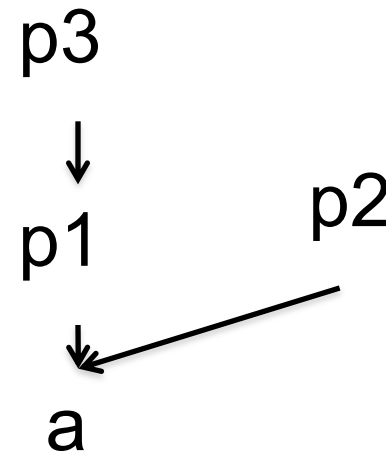
- We can preprocess **any** C program into a sequence of statements of these kinds

Points-to Analysis: Property Space

- Lattice L, \leq

- Lattice of the **subsets** over edges $p \rightarrow q$ where p and q are (names of) memory locations
- ... or in simpler terms, lattice elements are points-to graphs, e.g.,

- V is points-to graph union
- 0 of L is empty graph
- 1 of L is complete graph





Points-to Graphs Pt

- Nodes are names of memory locations
 - Program variables, a , p :
 - $p = \&a$
 - But also heap locations:
 - $p = \text{malloc}(\text{sizeof}(\text{int})) // h1$
- Edges represent points-to relations
 - E.g., $p \rightarrow a$, read: “ p points to a ”
 - E.g., $p \rightarrow h1$, read: “ p points to heap location $h1$ ”

Points-to Analysis: Transfer Functions

(1) $f_{p=&q}$: “kill” all points-to edges from p and “generate” a new points-to edge from p to q

(2) $f_{p=q}$: “kill” all points-to edges from p ; “generate” new points-to edges from p to every x , such that q points to x in incoming points-to graph $in(j)$

Points-to Analysis: Transfer Functions

(3) $f_{p=*q}$: “kill” all points-to edges from p ; “generate” new points-to edges from p to every x , s.t. there is y where q points to y , and y points to x in $in(j)$

(4) $f_{*p=q}$: **Do not kill!** Can you think of a reason why?
“Generate” new points-to edges from every y to every x , such that p points to y and q points to x



Points-to Analysis is Monotone

- To argue monotonicity we must show that if \mathbf{Pt}_1 is \leq (subset of) \mathbf{Pt}_2 , then $\mathbf{f}(\mathbf{Pt}_1) \leq \mathbf{f}(\mathbf{Pt}_2)$ for each transfer function \mathbf{f}

(1) $\mathbf{Pt}_1 \leq \mathbf{Pt}_2$ then $\mathbf{f}_{p=\&q}(\mathbf{Pt}_1) \leq \mathbf{f}_{p=\&q}(\mathbf{Pt}_2)$

(2) $\mathbf{Pt}_1 \leq \mathbf{Pt}_2$ then $\mathbf{f}_{p=q}(\mathbf{Pt}_1) \leq \mathbf{f}_{p=q}(\mathbf{Pt}_2)$

(3) $\mathbf{Pt}_1 \leq \mathbf{Pt}_2$ then $\mathbf{f}_{p=*q}(\mathbf{Pt}_1) \leq \mathbf{f}_{p=*q}(\mathbf{Pt}_2)$

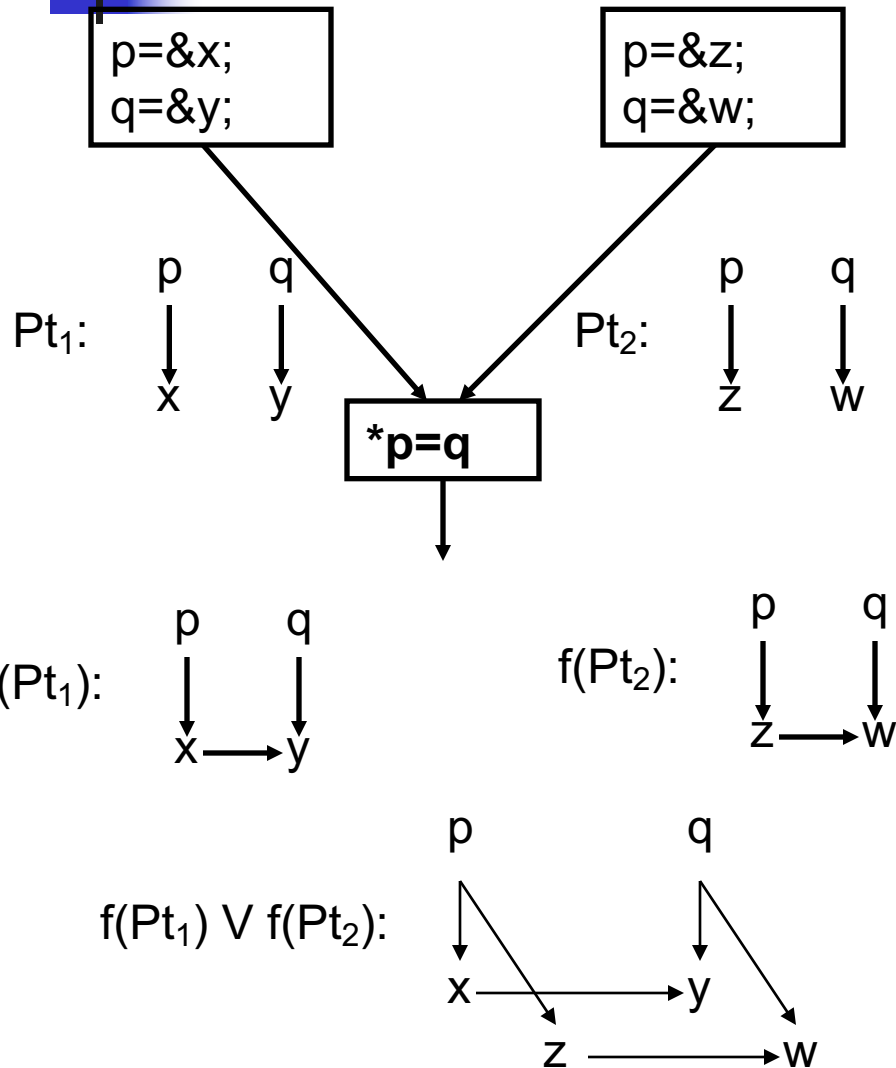
(4) $\mathbf{Pt}_1 \leq \mathbf{Pt}_2$ then $\mathbf{f}_{*p=q}(\mathbf{Pt}_1) \leq \mathbf{f}_{*p=q}(\mathbf{Pt}_2)$



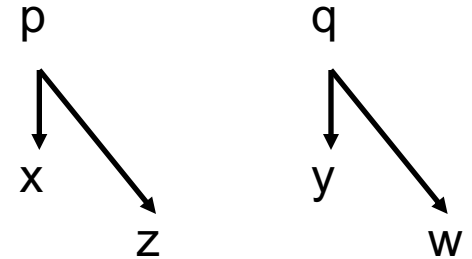
... but it is not distributive!

- Because of updates!

Points-to Analysis is Not Distributive

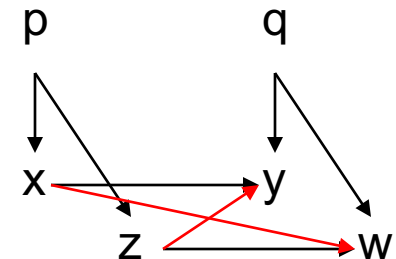


$Pt_1 \vee Pt_2 :$



What f for `*p = q` does: Adds edges from each variable that p points to (x and z), to each variable that q points to (y and w). Result is 4 new edges: from x to y and to w and from z to y and to w

$f(Pt_1 \vee Pt_2):$



MFP vs. MOP for Points-to

1. if (n>0)

2.
p=&x;
q=&y;

3.
p=&z;
q=&w;

4. *p=q

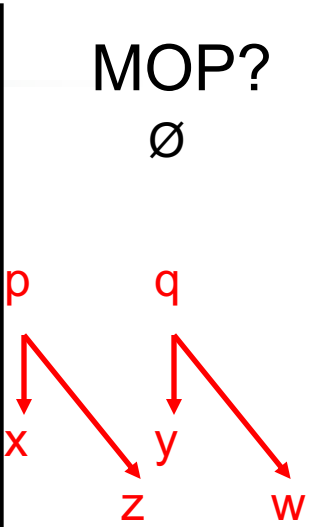
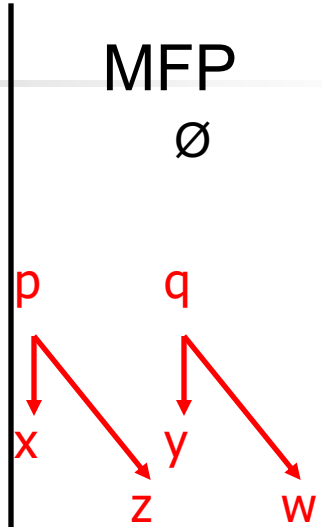
5. ...

$in_{PT}(4) = out_{PT}(2) \vee out_{PT}(3)$

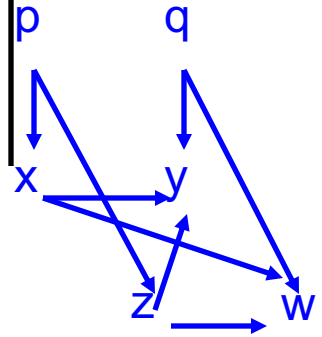
$out_{PT}(4) = f_{*p=q}(in_{PT}(4))$

$in_{PT}(5) = out_{PT}(4)$

$in_{PT}(4)$:



$in_{PT}(5)$:





Next Time

- Putting this into practice
- Program analysis frameworks
 - Soot
 - Ghidra