

# Dataflow Analysis

- **Lattice theoretic foundations**
  - Partial ordering
  - Meet, Join, Lattice, Chain
- **Round robin fixed point iteration**
- **Function properties**
  - Monotonicity
  - Distributivity
- **Justification for using fixed-point iteration on dataflow equations**
- **Meet Over all Paths solution**

# Lattice Theory

- **Partial ordering**  $\sqsubseteq$ 
  - Relation between pairs of elements
  - Reflexive  $x \sqsubseteq x$
  - Anti-symmetric  $x \sqsubseteq y, y \sqsubseteq x \implies x = y$
  - Transitive  $x \sqsubseteq y, y \sqsubseteq z \implies x \sqsubseteq z$
- **Poset (Set  $S$ ,  $\sqsubseteq$ )**
- **0 Element**  $0 \sqsubseteq x, \forall x \in S$
- **1 Element**  $1 \geq \forall x \in S$

**A poset need not have 0 or 1 element.**

# Lattice Theory

- **Greatest lower bound (glb)**  
 $l_1, l_2 \in \text{poset } S, a \in S$  is  $\text{glb}(l_1, l_2)$   
 if  $a \leq l_1$  and  $a \leq l_2$  then  
 for any  $b \in S, b \leq l_1, b \leq l_2 \Rightarrow b \leq a$

**if glb is unique it is called the meet ( $\sqcap$ ) of  $l_1$  and  $l_2$ .**

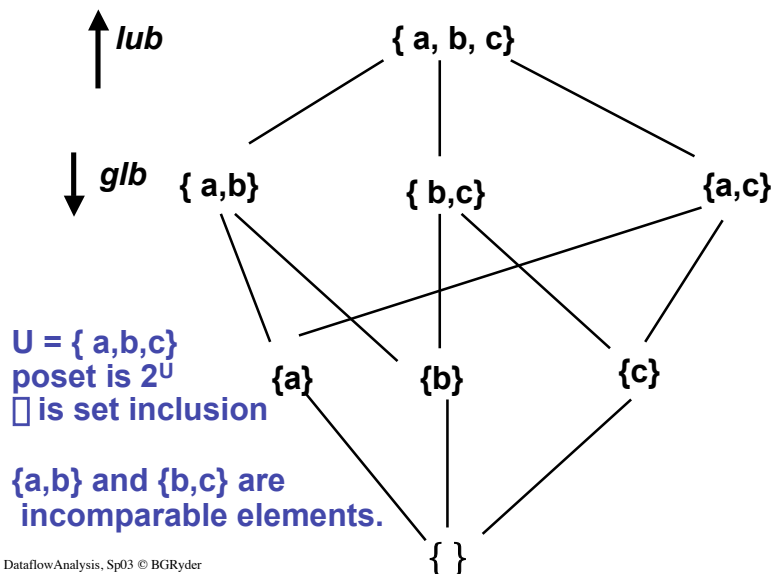
- **Least upper bound (lub)**  
 $l_1, l_2 \in \text{poset } S, c \in S$  is  $\text{lub}(l_1, l_2)$   
 if  $c \geq l_1$  and  $c \geq l_2$  then  
 for any  $d \in S, d \geq l_1, d \geq l_2 \Rightarrow c \leq d$ .

**if lub is unique is called the join ( $\sqcup$ ) of  $l_1$  and  $l_2$ .**

DataflowAnalysis, Sp03 © BGRyder

3

# Poset Example



DataflowAnalysis, Sp03 © BGRyder

4

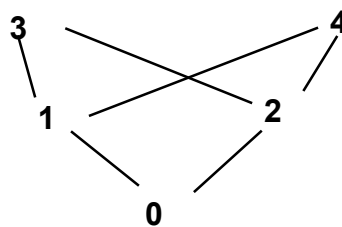
## Definition of a Lattice $(L, \sqcap, \sqcup)$

- $L$ , a poset under  $\sqsubseteq$  such that every pair of elements has a unique glb (meet) and lub (join).
- A lattice need not contain an 0 or 1 element.
- A finite lattice must contain an 0 and 1 element.
- Not every poset is a lattice.
- If  $a \sqsubseteq x \sqsubseteq x \sqsubseteq L$ , then  $a$  is 0 element of  $L$
- If  $x \sqsubseteq a \sqsubseteq x \sqsubseteq L$ , then  $a$  is 1 element of  $L$

DataflowAnalysis, Sp03 © BGRyder

5

## a poset, but not a lattice



There is no  $\text{lub}(3,4)$  in this poset so it is not a lattice.

DataflowAnalysis, Sp03 © BGRyder

6

## Examples of Lattices

- $H = (2^U, \sqsubseteq, \sqsupseteq)$  where  $U$  is a finite set
  - $\text{glb}(s1, s2)$  is  $(s1 \sqcap s2)$  which is  $s1 \sqcap s2$
  - $\text{lub}(s1, s2)$  is  $(s1 \sqcup s2)$  which is  $s1 \sqcup s2$
- $J = (N_1, \text{gcd, lowest common multiple})$ 
  - partial order relation is integer divide on  $N_1$   
 $n1 \mid n2$  if division is even
  - $\text{lub}(n1, n2)$  is  $n1 \sqcup n2 = \text{lowest common multiple}(n1, n2)$
  - $\text{glb}(n1, n2)$  is  $n1 \sqcap n2 = \text{greatest common divisor}(n1, n2)$

DataflowAnalysis, Sp03 © BGRyder

7

## Chain

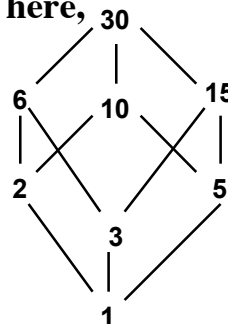
- A poset  $C$  where, for every pair of elements  $c1, c2 \in C$ , either  $c1 \sqsubseteq c2$  or  $c2 \sqsubseteq c1$ .

e.g.,  $\{ \} \sqsubseteq \{a\} \sqsubseteq \{a,b\} \sqsubseteq \{a,b,c\}$

and from the lattice as shown here,

$1 \sqsubseteq 2 \sqsubseteq 6 \sqsubseteq 30$

$1 \sqsubseteq 3 \sqsubseteq 15 \sqsubseteq 30$



Lattices are used in dataflow analysis to argue the existence of a solution obtainable through fixed-point iteration.

DataflowAnalysis, Sp03 © BGRyder

8

## Round-robin Fixed-point Iteration

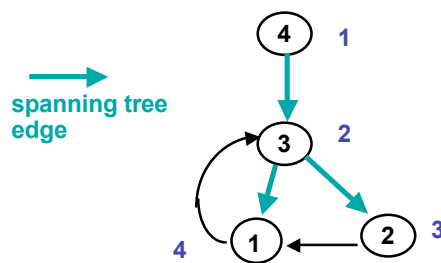
- Iterates through nodes of cfg in specific node order, which preserves the partial order of the graph.
- Recall worklist algorithm(per bit) has  $O(me)$  work where  $|E|=e$  and  $m$  is number of dataflow facts in the set.
- Postorder on depth-first spanning tree defines **rPostorder** which is *reverse postorder*
  - Ancestor-first node ordering, good for **forward** dataflow problems
  - $rPostorder\_number = (|N| + 1) - postorder\_number$
  - **Postorder** is good node ordering for **backward** dataflow problems

DataflowAnalysis, Sp03 © BGRyder

9

## rPostorder Example

control flowgraph



Postorder numbering of depth-first leftmost traversal of cfg  
 $rPostorder\#$  is defined as  $(n+1) - Postorder\#$

DataflowAnalysis, Sp03 © BGRyder

10

# Round-robin Iteration

*Form depth-first spanning tree of control flowgraph and derive rPostorder numbering for nodes*

**change = true**

*Initialize REACH(m) =  $\emptyset$ ,  $\square$  m*

```
while (change) do
{  change = false
  for (m = 2; m <= n; m++)
  { new =  $\square$  ( Reach(j)  $\square$  pres(j)  $\square$  dgen(j) )
    j = pred(m)

    if (new != REACH(m))
      then {REACH(m) = new; change = true}
  }
}
```

Round-robin tends to push dataflow facts as deep as possible on acyclic paths in the flowgraph on each iteration

DataflowAnalysis, Sp03 © BGRyder

11

# Round-robin Iteration

*Form depth-first spanning tree of control flowgraph and derive rPostorder numbering for nodes*

**change = true**

*Initialize REACH(m) =  $\emptyset$ ,  $\square$  m*

**O(n)**

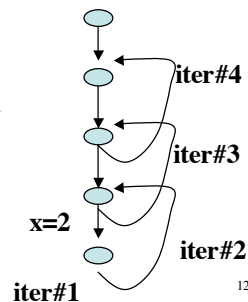
**#iters of while \* cost each iteration**  
**cost for statement ~ cost of all ops**  
**during lifetime of for**

```
while (change) do
{  change = false
  for (m = 2; m <= n; m++)
  { new =  $\square$  ( Reach(j)  $\square$  pres(j)  $\square$  dgen(j) )
    j = pred(m)

    if (new != REACH(m))
      then {REACH(m) = new; change = true}
  }
}
```

d = 3; d+2 = 5

**iter#5,**  
**no change**



DataflowAnalysis, Sp03 © BGRyder

12

## Cost of Loop Operations

Calculate *for* statement cost over all its iterations:

1.  $2(n-1)$  bit vector operations to calculate each node's factor  $\text{new} = \neg (\text{Reach}(j) \wedge \text{pres}(j) \wedge \text{dgen}(j))$  and memoize it;
  2. unions:  $(\# \text{preds}(m)-1)$  for each node  $m$  means we do  $(e - n)$  unions per *for* statement for all nodes;
  3. rest of *for* is constant number of operations
- Each while loop iteration costs  $(e - n) + (2n-2) \wedge e + n$ ;

Therefore, entire algorithm costs  $O(de)$ , ( $\wedge O(dn)$  by assumption for flowgraphs), bit vector operations where  $d+2$  for classical bitvector problems.

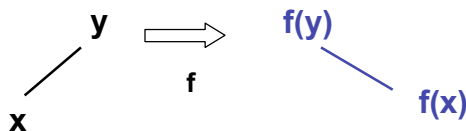
**Note:** this dominates  $O(e)$  of calculation of rPostorder numbers and the  $O(n)$  of initialization loop.

DataflowAnalysis, Sp03 © BGRyder

13

## More Lattice Properties

- **Finite length:** if every chain in lattice is finite
- **Bounded lattice:** if lattice contains both 0 and 1 elements
- **Distributive lattice:** if  $\forall x, y, z \in S$   
 $x \wedge (y \vee z) = (x \wedge y) \vee (x \wedge z)$   
 $x \vee (y \wedge z) = (x \vee y) \wedge (x \vee z)$
- $(S, \sqsubseteq)$  poset,  $f: S \rightarrow S$  is **monotonic** iff  
 $\forall x, y \in S, x \sqsubseteq y \implies f(x) \sqsubseteq f(y)$
- **Monotonic functions preserve domain ordering in their range values**



DataflowAnalysis, Sp03 © BGRyder

14

# Fixed point theorem - Why it works?

## Intuition--

Given a  $0$  in lattice and monotonic function  $f$ ,  $0 \sqsubseteq f(0)$ .

Apply  $f$  again and obtain

$$0 \sqsubseteq f(0) \sqsubseteq f(f(0)) = f^2(0)$$

Continuing,

$0 \sqsubseteq f(0) \sqsubseteq f^2(0) \sqsubseteq f^3(0) \sqsubseteq \dots \sqsubseteq f^k(0) = f^{k+1}(0)$  for a finite chain lattice.

This is tantamount to saying  $\lim_{k \rightarrow \infty} f^k(0)$  exists and is called the *least fixed point* of  $f$ ,

since  $f(f^k(0)) = f^{k+1}(0)$

# Fixed Point Theorem

**Thm:**  $f: S \rightarrow S$  monotonic function on poset  $(S, \sqsubseteq)$  with a  $0$  element and finite length. The *least fixed point* of  $f$  is  $f^k(0)$  where

i.  $f^0(x) = x$ ,

ii.  $f^{i+1}(x) = f(f^i(x))$ ,  $i \geq 0$ ,

iii.  $f^k(0) = f(f^k(0))$  and this is the smallest  $k$  for which this is true.

- For any  $p$  such that  $f(p)=p$ ,  $f^k(0) \sqsubseteq p$ .
- Theorem justifies the iterative algorithm for global data flow analysis for lattices & functions with right properties
- Dual theorem exists for 1 element and *maximal fixed point* for  $k$  such that  $f^k(1) = f^{k+1}(1)$ .



## Application to DFA on CFG

- **Cartesian cross product of posets is a poset**
  - if  $(S, \sqsubseteq)$  is a poset, then  $(S \times S \times \dots \times S, \sqsubseteq')$  is a poset whose partial ordering is component-wise  $\sqsubseteq$ 

$$W=(W_1, W_2, \dots, W_n) \sqsubseteq' Y=(Y_1, Y_2, \dots, Y_n) \text{ iff}$$

$$W_1 \sqsubseteq Y_1 \text{ and } W_2 \sqsubseteq Y_2, \dots, \text{ and } W_n \sqsubseteq Y_n$$
- **Cartesian cross product of lattices is a lattice (with an induced partial ordering)**
- **Monotone function on cross product lattice can be built from monotone functions on each component lattice**  $F(Y_1, Y_2, \dots, Y_n) = (g_1(Y_1, Y_2, \dots, Y_n), \dots, g_n(Y_1, Y_2, \dots, Y_n))$   
 where  $g_i : (L, L, \dots, L) \rightarrow L$  and  $\sqsubseteq'$  is component-wise  $\sqsubseteq$  from  $L$

DataflowAnalysis, Sp03 © BGRyder

17

## Example - Available Exprs

- lattice is  $2^{\text{Exprs}}$  where Exprs is set of all binary expressions in program
- Partial order is  $\sqsubseteq$  (subset inclusion) so meet is  $\sqcap$
- $\langle \text{Exprs}, \text{Exprs}, \dots, \text{Exprs} \rangle$  is 1 element
- $\langle \emptyset, \emptyset, \dots, \emptyset \rangle$  is 0 element
- From the data flow equations for AVAIL, we know that if a set of dataflow facts  $X$  is true on entry to a flowgraph node  $n$ , then  $f(X)$  is true on each exit edge of  $n$  where
 
$$f(X) = \text{epres}(n) \sqcap X \sqcap \text{egen}(n)$$
 $f$  is called the *transfer function* for AVAIL

DataflowAnalysis, Sp03 © BGRyder

18

## Example, Available Exprs

- Cross product lattice is
  - $(2^{\text{Exprs}}, 2^{\text{Exprs}}, \dots, 2^{\text{Exprs}})$  with  $n$  components where  $n$  is number of nodes in the cfg and  $\sqcap$  is  $\sqcap$  component-wise
- Avail equation at a node can be expressed thusly,
 
$$\text{Avail}(j) = \sqcap_{m \in \text{Pred}(j)} \{ \text{Avail}(m) \sqcap \text{epres}(m) \sqcap \text{egen}(m) \}$$
  - AVAIL (j) is the solution at entry of node j and  $f(\text{AVAIL}(j))$  is solution at exit of node j,
 
$$g_j = \sqcap_{m \in \text{Pred}(j)} f(g_m)$$
  - One step of the worklist algorithm maps a potential solution for AVAIL at a node of the cfg to another potential solution for that node

## Example, Available Exprs

- Can you show  $g_i$  monotone?
 
$$g_i : (2^{\text{Exprs}}, 2^{\text{Exprs}}, \dots, 2^{\text{Exprs}}) \rightarrow 2^{\text{Exprs}}$$
- Then this induces the monotonicity of  $F$ ,
 
$$F = (g_1, \dots, g_n)$$
- Application of dual of fixed point theorem here to find the maximal fixed point. Iterate down from the 1 element.
  - Initialize  $\sqcap$  to  $\emptyset$ , all other cfg nodes to Exprs.

## Sketch of Validity Proof of Iterative Algorithm for DFA

- Let  $\text{Redef}(Y_1^i \dots Y_n^i)$  be result after  $i$  steps of the worklist algorithm (per node) for DFA.  
 $\text{Redef}(Y_1^i \dots Y_n^i) = (Y_1^i \dots Y_{k-1}^i, g_k(Y_1^i \dots Y_n^i), Y_{k+1}^i, \dots, Y_n^i)$
- In the next iteration ( $i+1$ st), only 1 component of  $Y = (Y_1^i \dots Y_n^i)$  changes and this component is chosen non-deterministically.
- Recall we have a function defined on the flowgraph:  
 $F(Y_1, \dots, Y_n) = (g_1(Y_1, \dots, Y_n), \dots, g_n(Y_1, \dots, Y_n))$  where  $g_j : (L, L, \dots, L) \rightarrow L$  is defined by the dataflow problem

## Sketch of Proof, cont.

- Assume minimal fixed point of  $F$  is  $F^m(0)$ .
  - Must show iterative algorithm halts at a fixed point of  $\text{Redef}$  function,  $\text{Redefmin}$ , and that this fixed point is meaningful to the data flow problem under solution, namely  $\text{Redefmin} = F^m(0)$ .
  - First, show iterates of  $\text{Redef}$  form a chain.  
 $0 \sqsubseteq Y^1 \sqsubseteq Y^2 \sqsubseteq Y^3 \sqsubseteq \dots \sqsubseteq Y^k$
- Therefore,  $\text{Redefmin} = \text{Redef}^k(0)$  for some smallest  $k$ , by finiteness of chains on the lattice.

## Sketch of Proof, cont.

Second, let  $F^m(0)$  be the minimum fixed point of the data flow function  $F$  on the cross product lattice. Show  $F^m(0) \sqsubseteq \text{Redefmin}$ .

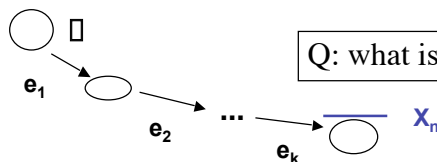
$F^m(0) \sqsubseteq \text{Redefmin}$  because  $\text{Redefmin}$  is also a fixed point of  $F$ , but not necessarily its minimal fixed point.

Third, show  $\text{Redefmin} \sqsubseteq F^m(0)$  in two steps

- (1)  $\text{Redef}^k(0) \sqsubseteq F^m(0)$  for all  $k$  and
- (2)  $\text{Redefmin} = \text{Redef}^k(0)$  for large enough  $k$ , so for  $k$  large enough,  $\text{Redefmin} \sqsubseteq F^m(0)$

Therefore,  $\text{Redefmin} = F^m(0)$  and the fixed point iterative procedure does find the minimal fixed point on the lattice

## Meet Over all Paths solution



Q: what is relation between MFP and MOP?

- Maximal data flow information desired is obtained by traversing ALL PATHS from  $\emptyset$  to  $n$ .
- Consider  $X_n = \bigsqcup f_Q(0)$ , for  $Q$  a path in the flowgraph,  $Q = e_1, e_2, \dots, e_k$  so that

$$X_n = \bigsqcup_{Q = e_k \dots e_1} f_{e_k} \circ f_{e_{k-1}} \circ \dots \circ f_{e_2} \circ f_{e_1}(0)$$

MOP is best summary of data flow facts possible to compute at compile-time.