

## Reading Assignment

Read "An Efficient Method of Computing Static Single Assignment Form."

(Linked from the class Web page.)

## How Good Is Iterative Data Flow Analysis?

A single execution of a program will follow some path

$$b_0, b_{i_1}, b_{i_2}, \dots, b_{i_n}$$

The Data Flow solution along this path is

$$f_{i_n}(\dots f_{i_2}(f_{i_1}(f_0(T)))\dots) \equiv f(b_0, b_{i_1}, \dots, b_{i_n})$$

The best possible static data flow solution at some block  $b$  is computed over all possible paths from  $b_0$  to  $b$ .

Let  $P_b$  = The set of all paths from  $b_0$  to  $b$ .

$$\text{MOP}(b) = \bigwedge_{p \in P_b} f(p)$$

Any particular path  $p_i$  from  $b_0$  to  $b$  is included in  $P_b$ .

Thus  $\text{MOP}(b) \wedge f(p_i) = \text{MOP}(b) \leq f(p_i)$ .

This means  $\text{MOP}(b)$  is *always* a safe approximation to the "true" solution  $f(p_i)$ .

If we have the distributive property for transfer functions,

$$f(a \wedge b) = f(a) \wedge f(b)$$

then our iterative algorithm *always* computes the MOP solution, the best static solution possible.

To prove this, note that for trivial path of length 1, containing only the start block,  $b_0$ , the algorithm computes  $f_0(T)$  which is  $\text{MOP}(b_0)$  (trivially).

Now assume that the iterative algorithm for paths of length  $n$  or less to block  $c$  *does* compute  $\text{MOP}(c)$ .

We'll show that for paths to block  $b$  of length  $n+1$ ,  $\text{MOP}(b)$  is computed.

Let  $P$  be the set of all paths to  $b$  of length  $n+1$  or less.

The paths in  $P$  end with  $b$ .

$$\text{MOP}(b) = f_b(f(P_1)) \wedge f_b(f(P_2)) \wedge \dots$$

where  $P_1, P_2, \dots$  are the prefixes (of length  $n$  or less) of paths in  $P$  with  $b$  removed.

Using the distributive property,

$$f_b(f(P_1)) \wedge f_b(f(P_2)) \wedge \dots = f_b(f(P_1) \wedge f(P_2) \wedge \dots).$$

But note that  $f(P_1) \wedge f(P_2) \wedge \dots$  is just the input to  $f_b$  in our iterative algorithm, which then applies  $f_b$ .

Thus  $\text{MOP}(b)$  for paths of length  $n+1$  is computed.

For data flow problems that aren't distributive (like constant propagation), the iterative solution is  $\leq$  the MOP solution.

This means that the solution is a safe approximation, but perhaps not as "sharp" as we might wish.

## Exploiting Structure in Data Flow Analysis

So far we haven't utilized the fact that CFGs are constructed from standard programming language constructs like IFs, Fors, and Whiles.

Instead of iterating across a given CFG, we can isolate, and solve symbolically, subgraphs that correspond to "standard" programming language constructs.

We can then progressively simplify the CFG until we reach a single node, or until we reach a CFG structure that matches no standard pattern.

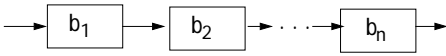
In the latter case, we can solve the residual graph using our iterative evaluator.

## Three Program-Building Operations

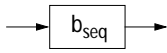
1. Sequential Execution (" $;$ ")
2. Conditional Execution (If, Switch)
3. Iterative Execution (While, For, Repeat)

## Sequential Execution

We can reduce a sequential “chain” of basic blocks:



into a single composite block:



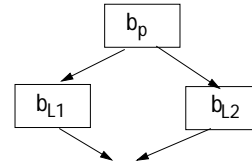
The transfer function of  $b_{seq}$  is

$$f_{seq} = f_n \circ f_{n-1} \circ \dots \circ f_1$$

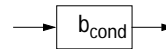
where  $\circ$  is functional composition.

## Conditional Execution

Given the basic blocks:



we create a single composite block:



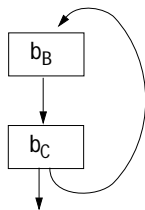
The transfer function of  $b_{cond}$  is

$$f_{cond} = f_{L1} \circ f_p \wedge f_{L2} \circ f_p$$

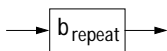
## Iterative Execution

Repeat Loop

Given the basic blocks:



we create a single composite block:



Here  $b_B$  is the loop body, and  $b_C$  is the loop control.

If the loop iterates once, the transfer function is  $f_C \circ f_B$ .

If the loop iterates twice, the transfer function is  $(f_C \circ f_B) \circ (f_C \circ f_B)$ .

Considering all paths, the transfer function is  $(f_C \circ f_B) \wedge (f_C \circ f_B)^2 \wedge \dots$

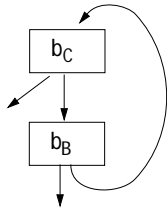
Define  $\text{fix } f \equiv f \wedge f^2 \wedge f^3 \wedge \dots$

The transfer function of repeat is then

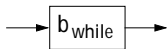
$$f_{repeat} = \text{fix}(f_C \circ f_B)$$

## While Loop.

Given the basic blocks:



we create a single composite block:



Here again  $b_B$  is the loop body, and  $b_C$  is the loop control.

The loop always executes  $b_C$  at least once, and always executes  $b_C$  as the last block before exiting.

The transfer function of a while is therefore

$$f_{\text{while}} = f_C \wedge \text{fix}(f_C \circ f_B) \circ f_C$$

## Evaluating Fixed Points

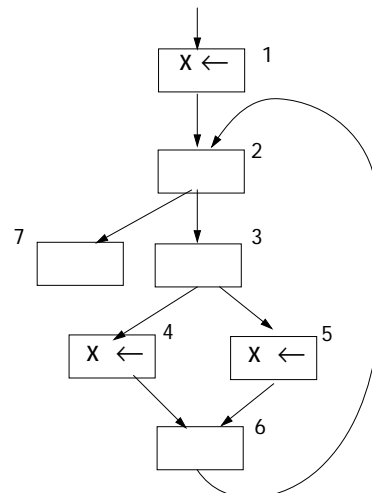
For lattices of height  $H$ , and monotone transfer functions,  $\text{fix } f$  needs to look at no more than  $H$  terms.

In practice, we can give  $\text{fix } f$  an operational definition, suitable for implementation:

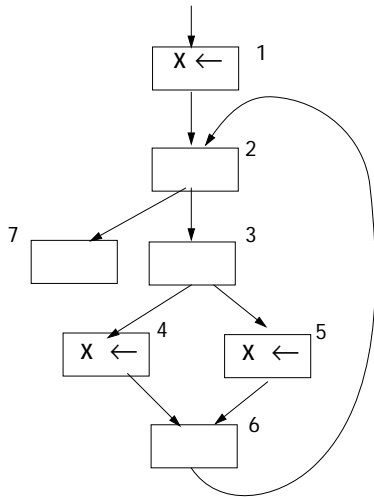
Evaluate

```
(fix f)(x) {
  prev = soln = f(x);
  while (prev ≠ new = f(prev)) {
    prev = new;
    soln = soln ∧ new;
  }
  return soln;
}
```

## Example—Reaching Definitions



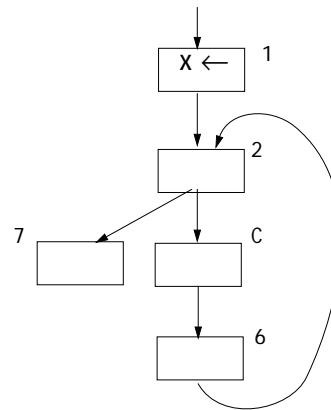
The transfer functions are either constant-valued ( $f_1 = \{b1\}$ ,  $f_4 = \{b4\}$ ,  $f_5 = \{b5\}$ ) or identity functions ( $f_2 = f_3 = f_6 = f_7 = \text{Id}$ ).



First we isolate and reduce the conditional:

$$f_C = f_4 \circ f_3 \wedge f_5 \circ f_3 = \{b_4\} \circ Id \cup \{b_5\} \circ Id = \{b_4, b_5\}$$

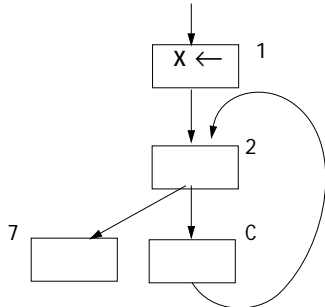
Substituting, we get



We can combine  $b_C$  and  $b_6$ , to get a block equivalent to  $b_C$ . That is,

$$f_6 \circ f_C = Id \circ f_C = f_C$$

We now have

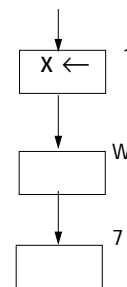


We isolate and reduce the while loop formed by  $b_2$  and  $b_C$ , creating  $b_W$ .

The transfer function is

$$\begin{aligned} f_W &= f_2 \wedge (\text{fix}(f_2 \circ f_C) \circ f_2 = \\ & Id \cup (\text{fix}(Id \circ f_C) \circ Id = \\ & Id \cup (\text{fix}(f_C)) = \\ & Id \cup (f_C \wedge f_C^2 \wedge f_C^3 \wedge \dots) = \\ & Id \cup \{b_4, b_5\} \end{aligned}$$

We now have



We compose these three sequential blocks to get the whole solution,  $f_P$ .

$$f_P = Id \circ (Id \cup \{f_4, f_5\}) \circ \{b_1\} = \{b_1, b_4, b_5\}.$$

These are the definitions that reach the end of the program.

We can expand subgraphs to get the solutions at interior blocks.

Thus at the beginning of the while, the solution is  $\{b1\}$ .

At the head if the If, the solution is

$$\begin{aligned} & (\text{Id} \cup (\text{Id} \circ f_C \circ \text{Id}) \cup \\ & (\text{Id} \circ f_C \circ \text{Id} \circ f_C \circ \text{Id}) \cup \dots)(\{b1\}) = \\ & \{b1\} \cup \{b4,b5\} \cup \{b4,b5\} \cup \dots = \\ & \{b1,b4,b5\} \end{aligned}$$

At the head of the then part of the If, the solution is  $\text{Id}(\{b1,b4,b5\}) = \{b1,b4,b5\}$ .

## Static Single Assignment Form

Many of the complexities of optimization and code generation arise from the fact that a given variable may be assigned to in *many* different places.

Thus reaching definition analysis gives us the *set* of assignments that *may* reach a given use of a variable.

Live Range Analysis must track *all* assignments that may reach a use of a variable and merge them into the same live range.

Available Expression Analysis must look at *all* places a variable may be assigned to and decide if any kill an already computed expression.

## What If

each variable is assigned to in only one place?

(Much like a named constant).

Then for a given use, we can find a single *unique* definition point.

But this seems *impossible* for most programs—or is it?

In *Static Single Assignment (SSA)* Form each assignment to a variable,  $v$ , is changed into a unique assignment to new variable,  $v_i$ .

If variable  $v$  has  $n$  assignments to it throughout the program, then (at least)  $n$  new variables,  $v_1$  to  $v_n$ , are created to replace  $v$ . All uses of  $v$  are replaced by a use of some  $v_i$ .

## Phi Functions

Control flow can't be predicted in advance, so we can't always know which definition of a variable reached a particular use.

To handle this uncertainty, we create *phi functions*.

As illustrated below, if  $v_i$  and  $v_j$  both reach the top of the same block, we add the assignment

$$v_k \leftarrow \phi(v_i, v_j)$$

to the top of the block.

Within the block, all uses of  $v$  become uses of  $v_k$  (until the next assignment to  $v$ ).

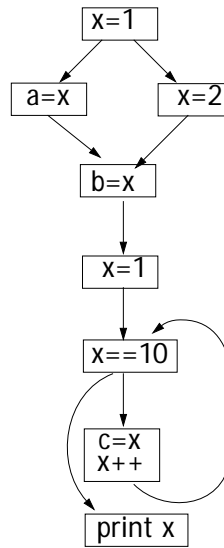
## What does $\phi(v_i, v_j)$ Mean?

One way to read  $\phi(v_i, v_j)$  is that if control reaches the phi function via the path on which  $v_i$  is defined,  $\phi$  "selects"  $v_i$ ; otherwise it "selects"  $v_j$ .

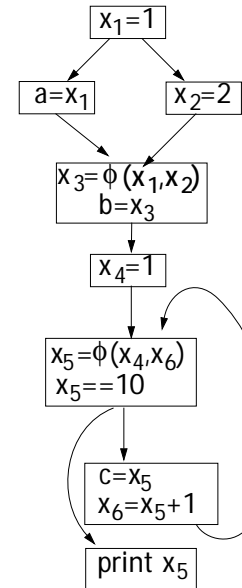
Phi functions may take more than 2 arguments if more than 2 definitions might reach the same block.

Through phi functions we have simple links to all the places where  $v$  receives a value, directly or indirectly.

## Example



Original CFG



CFG in SSA Form

In SSA form computing live ranges is almost trivial. For each  $x_i$  include all  $x_j$  variables involved in phi functions that define  $x_i$ .

Initially, assume  $x_1$  to  $x_6$  (in our example are independent). We then union into equivalence classes  $x_i$  values involved in the same phi function or assignment.

Thus  $x_1$  to  $x_3$  are unioned together (forming a live range). Similarly,  $x_4$  to  $x_6$  are unioned to form a live range.