Dataflow Analysis + SSA CSE 401 Sections

Announcements 📢

- **Codegen** deadline **TONIGHT 11:59 pm**. Remember to commit/push/tag and complete the Gradescope submission. If using late days, may submit at most two days late (depending on number of remaining late days).
 - Once codegen is finished, we'll do an overall evaluation of your compiler, all phases, and rerun a comprehensive set of tests. This final evaluation is the major part (half!) of the overall project grade. So you need to fix any remaining bugs, from semantics/type checking all the way back to the scanner! (And your final evaluation will reflect these fixes)
 Image: Im
- **Project report** due Tuesday, 06/03 at 11:59 pm (no late days)
- HW4 due Thursday, 06/05 at 11:59 pm (late days allowed)



PeepholeA few InstructionsLocal

Intraprocedural / Global

Interprocedural

PeepholeA few InstructionsLocalA Basic Block

Intraprocedural / Global

Interprocedural

PeepholeA few InstructionsLocalA Basic BlockIntraprocedural / GlobalA Function/MethodInterproceduralInterprocedural

PeepholeA few InstructionsLocalA Basic BlockIntraprocedural / GlobalA Function/MethodInterproceduralA Program

Overview of Dataflow Analysis

- A framework for exposing properties about programs
- Operates using sets of "facts"



Overview of Dataflow Analysis

- A framework for exposing properties about programs
- Operates using sets of "facts"
- Just the initial discovery phase
 - Optimizations can then be made based on the analysis



Overview of Dataflow Analysis

- Basic Framework of Set Definitions (for a Basic Block *b*):
 - IN(*b*): facts true on entry to *b*
 - OUT(b): facts true on exit from b
 - GEN(b): facts created (and not killed) in b
 - KILL(*b*): facts killed in *b*

Overview of Dataflow Analysis (continued)

- These are related by the equation
 - $OUT(b) = GEN(b) \cup (IN(b) KILL(b))$
 - Solve this iteratively for all blocks
 - Sometimes information propagates forward; sometimes backward
 - But will reach correct solution (fixed point) regardless of order in which blocks are considered

Reaching Definitions (*A Dataflow Problem*)

"What definitions of each variable might reach this point"

- Could be used for:
 - Constant Propagation
 - Uninitialized Variables

```
int x;
                if (y > 0) {
                  \mathbf{x} = \mathbf{y};
                } else {
                  x = 0;
                }
                System.out.println(x);
"x=y", "x=0"
```

Reaching Definitions (*A Dataflow Problem*)

still: "x=y", "x=0"

"What definitions of each variable might reach this point"

- **Be careful**: Does not involve the *value* of the variable definition
 - The dataflow problem
 "Available Expressions"
 is designed for that

```
int x;
if (y > 0) {
  \mathbf{X} = \mathbf{V}
} else {
  x = 0;
}
  = -1:
System.out.println(x);
```

Equations for Reaching Definitions

- IN(*b*): the definitions reaching upon entering block b
- OUT(*b*): the definitions reaching upon exiting block b
- GEN(b): the definitions assigned and not killed in block b
- KILL(*b*): the definitions of variables overwritten in block b

$$IN(b) = U_{p \in pred(b)}OUT(p)$$
$$OUT(b) = GEN(b) \cup (IN(b) - KILL(b))$$

Problems 1(a) and 1(b)

- L0: a = 0 L1: b = a + 1 L2: c = c + b L3: a = b * 2 L4: if a < N goto L1
- L5: return c

Note: Labels are used in place of actual definitions to save space!

| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|--------|---------|--------|---------|
| LO | LØ | | | | | |
| L1 | L1 | | | | | |
| L2 | L2 | | | | | |
| L3 | L3 | | | | | |
| L4 | | | | | | |
| L5 | | | | | | |

- L0: a = 0 L1: b = a + 1 L2: c = c + b L3: a = b * 2 L4: if a < N goto L1
- L5: return c

| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|--------|---------|--------|---------|
| LO | LØ | | | | | |
| L1 | L1 | | | | | |
| L2 | L2 | | | | | |
| L3 | L3 | LØ | | | | |
| L4 | | | | | | |
| L5 | | | | | | |

- L0: a = 0 L1: b = a + 1 L2: c = c + b L3: a = b * 2 L4: if a < N goto L1
- L5: return c

| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|------------|---------|--------|---------|
| LO | LØ | | | | | |
| L1 | L1 | | LØ | | | |
| L2 | L2 | | L0, L1 | | | |
| L3 | L3 | LØ | L0, L1, L2 | | | |
| L4 | | | L1, L2, L3 | | | |
| L5 | | | L1, L2, L3 | | | |

- L0: a = 0 L1: b = a + 1 L2: c = c + b L3: a = b * 2 L4: if a < N goto L1
- L5: return c

| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|------------|------------|--------|---------|
| LO | LØ | | | LØ | | |
| L1 | L1 | | LØ | L0, L1 | | |
| L2 | L2 | | L0, L1 | L0, L1, L2 | | |
| L3 | L3 | LØ | L0, L1, L2 | L1, L2, L3 | | |
| L4 | | | L1, L2, L3 | L1, L2, L3 | | |
| L5 | | | L1, L2, L3 | L1, L2, L3 | | |

- L0: a = 0 L1: b = a + 1 L2: c = c + b L3: a = b * 2 L4: if a < N goto L1
- L5: return c

| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|------------|------------|----------------|----------------|
| LO | LØ | | | LØ | | LØ |
| L1 | L1 | | LO | L0, L1 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L2 | L2 | | L0, L1 | L0, L1, L2 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L3 | L3 | LØ | L0, L1, L2 | L1, L2, L3 | L0, L1, L2, L3 | L1, L2, L3 |
| L4 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |
| L5 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |

L0: a = 0 L1: b = a + 1 L2: c = c + b L3: a = b * 2 L4: if a < N goto L1 L5: return c



| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|------------|------------|----------------|----------------|
| LO | LØ | | | LØ | | L0 |
| L1 | L1 | | LØ | L0, L1 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L2 | L2 | | L0, L1 | L0, L1, L2 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L3 | L3 | LØ | L0, L1, L2 | L1, L2, L3 | L0, L1, L2, L3 | L1, L2, L3 |
| L4 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |
| L5 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |

L0: a = 0

Is it possible to replace the use of *a* in block L1 with the constant 0?

- L1: b = a + 1 L2: c = c + b
- L3: a = b * 2
- L4: if a < N goto L1
- L5: return c

| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|------------|------------|----------------|----------------|
| LO | LØ | | | LØ | | LØ |
| L1 | L1 | | LO | L0, L1 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L2 | L2 | | L0, L1 | L0, L1, L2 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L3 | L3 | LØ | L0, L1, L2 | L1, L2, L3 | L0, L1, L2, L3 | L1, L2, L3 |
| L4 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |
| L5 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |

L0: a = 0 L1: b = a + 1 L2: c = c + b L3: a = b * 2 L4: if a < N goto L1

return c

L5:

Is it possible to replace the use of *a* in block L1 with the constant 0?

No. To determine this, we would look at the IN set for block L1 -- the fact that the IN set contains two definitions of 'a' (L0 and L3) means we cannot perform this constant propagation. In other words, more than one definition of 'a' is a reaching definition to block L1, and therefore performing constant propagation would only preserve one possible value of 'a' and the generated code would not be equivalent.

| Block | GEN | KILL | IN (1) | OUT (1) | IN (2) | OUT (2) |
|-------|-----|------|------------|------------|----------------|----------------|
| LO | LØ | | | LØ | | LØ |
| L1 | L1 | | LO | L0, L1 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L2 | L2 | | L0, L1 | L0, L1, L2 | L0, L1, L2, L3 | L0, L1, L2, L3 |
| L3 | L3 | LØ | L0, L1, L2 | L1, L2, L3 | L0, L1, L2, L3 | L1, L2, L3 |
| L4 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |
| L5 | | | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 | L1, L2, L3 |

Phi-Functions

- A way to represent <u>multiple possible values</u> for a certain definition
 - Not a "real" instruction just a form of bookkeeping needed for SSA
 - Every variable gets a subscript





Phi-Functions

- Wherever a variable has multiple possible definitions entering a block
 - Inefficient (and unnecessary!) to consider all possible phi-functions at the start of each block



Example With a Loop



Notes: •Loop-back edges are also merge points, so require Φ -functions •a₀, b₀, c₀ are initial values of a, b, c on entry to initial block •b₁ is dead – can delete later •c is live on entry – either input parameter or uninitialized

Problem 2(a)



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |

A node **x** *dominates* a node **y** iff every path from the entry point of the control flow graph to **y** includes **x**.

A node **x** strictly dominates **a** node **y** iff **x** dominates **y** and $x \neq y$

Need to go through 0 to get through 1, 2, 3, 4, 5, 6 and 0 cannot strictly dominate itself



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |

A node **Y** is in the *dominance frontier* of node **X** iff **X** dominates an immediate predecessor of **Y** but **X** does not strictly dominate **Y**. A node **0** is in the *dominance frontier* of node **0** iff **0** dominates an immediate predecessor **(6)** of **0** but **0** does not strictly dominate **0**

0 dominates 6, 6 is an immediate predecessor of 0, 0 does not strictly dominate 0



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | ø | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |

A node **x** *dominates* a node **y** iff every path from the entry point of the control flow graph to **y** includes **x**.

A node **x** strictly dominates **a** node **y** iff **x** dominates **y** and $x \neq y$

1 does not dominate 6 because there is a path from 5 that doesn't include 1. 1 does not strictly dominate itself



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | Ø | 1,6 |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |

A node \mathbf{Y} is in the *dominance frontier* of node \mathbf{x} iff \mathbf{x} dominates an immediate predecessor of \mathbf{Y} but \mathbf{x} does not strictly dominate \mathbf{Y} .

X = 1, Y = 6, 1 dominates 1, 1 is an immediate predecessor of 6, 1 does not strictly dominate 6

X = 1, Y = 1, 1 dominates 1, 1 is an immediate predecessor of 1, 1 does not strictly dominate 1



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | ø | 1,6 |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | Ø | 1,6 |
| 2 | 3, 4, 5 | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |

A node **x** dominates a node **y** iff every path from the entry point of the control flow graph to **y** includes **x**. A node **x** strictly dominates a node **y** iff **x** dominates **y** and $x \neq y$

Need to go through 2 to get through 3, 4, 5 and 2 cannot strictly dominate itself



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | ø | 1,6 |
| 2 | 3, 4, 5 | 6 |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |

A node \mathbf{Y} is in the *dominance frontier* of node \mathbf{x} iff \mathbf{x} dominates an immediate predecessor of \mathbf{Y} but \mathbf{x} does not strictly dominate \mathbf{Y} .

X = 2, Y = 6, 2 dominates 5, 5 is an immediate predecessor of 6, 2 does not strictly dominate 6



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | Ø | 1,6 |
| 2 | 3, 4, 5 | 6 |
| 3 | Ø | 5 |
| 4 | | |
| 5 | | |
| 6 | | |

3 does not strictly dominate 5 (path through 4) and therefore does not strictly dominate anything else

3 dominates 3, 3 is an immediate predecessor of 5, 3 does not strictly dominate 5



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | Ø | 1,6 |
| 2 | 3, 4, 5 | 6 |
| 3 | ø | 5 |
| 4 | ø | 5 |
| 5 | | |
| 6 | | |

Same as previous slide but with 4 instead of 3



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | Ø | 1,6 |
| 2 | 3, 4, 5 | 6 |
| 3 | ø | 5 |
| 4 | ø | 5 |
| 5 | ø | 6 |
| 6 | | |

5 does not strictly dominate 6 (path through 1) and therefore does not strictly dominate anything else

5 dominates 5, 5 is an immediate predecessor of 6, 5 does not strictly dominate 6



| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|--------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | ø | 1,6 |
| 2 | 3, 4, 5 | 6 |
| 3 | ø | 5 |
| 4 | ø | 5 |
| 5 | ø | 6 |
| 6 | Ø | 0 |

6 does not strictly dominate 0 (path through 0) and therefore does not strictly dominate anything else

6 dominates 6, 6 is an immediate predecessor of 0, 6 does not strictly dominate 0

Problem 2(b)

Converting to SSA







Determine which variables need merging in each node



Assign numbers to definitions and add phi functions



Step 1: Dominance Frontiers

| NODE | STRICTLY DOMINATES | DOMINANCE FRONTIER |
|------|---------------------|--------------------|
| 0 | 1, 2, 3, 4, 5, 6 | 0 |
| 1 | Ø | 1, 6 |
| 2 | 3, 4, 5 | 6 |
| 3 | ø | 5 |
| 4 | ø | 5 |
| 5 | ø | 6 |
| 6 | Ø | 0 |

Converting to SSA



Compute the dominance frontier of each node





Determine which variables need merging in each node



We will compute using the dominance frontiers



Assign numbers to definitions and add phi functions



Step 2: Determine Necessary Merges

ITERATION 1: Each node in the dominance frontier of node X will merge any definitions created in node X.





Step 2: Determine Necessary Merges

ITERATION 2: Each merge will create a new definition, which may need merging again.





Step 2: Determine Necessary Merges

ITERATION 3: Each merge will create a new definition, which may need merging again.



Converting to SSA



Compute the dominance frontier of each node





Determine which variables need merging in each node





Assign numbers to definitions and add phi functions



Place phi functions first, then increment subscripts

Merges go first, and each successive definition of a variable should increment its index by 1.

$$B_{0} = \Phi(a_{0}, a_{2})$$

$$b_{1} = \Phi(b_{0}, b_{3})$$

$$c_{1} = \Phi(c_{0}, c_{5})$$

$$d_{1} = \Phi(d_{0}, d_{7})$$

$$e_{1} = \Phi(e_{0}, e_{4})$$

$$a_{2} = c_{1} + 2$$

$$d_{2} = a_{2} + b_{1}$$

Note: these subscripts determined after doing the rest of the CFG!

$$B_0$$

a = c + 2

d = a + b

a,b,c,d,e

Need to merge:

Merges go first, and each successive definition of a variable should increment its index by 1.

Merges go first, and each successive definition of a variable should increment its index by 1.

$$\mathbf{B}_2$$
 b = a + c
 \mathbf{B}_2 b = a + c

Nothing to merge

Merges go first, and each successive definition of a variable should increment its index by 1.



Nothing to merge

Merges go first, and each successive definition of a variable should increment its index by 1.

$$\mathbf{B}_{\mathbf{4}} \begin{bmatrix} \mathbf{d} = \mathbf{b} + 1 \end{bmatrix} \qquad \mathbf{B}_{\mathbf{4}} \begin{bmatrix} \mathbf{d}_{4} = \mathbf{b}_{2} + 1 \end{bmatrix}$$

Nothing to merge

Merges go first, and each successive definition of a variable should increment its index by 1.

Need to merge: d, e

Merges go first, and each successive definition of a variable should increment its index by 1.

b,c,d,e

