

Introduction to Data Flow Analysis

Reading: NNH 1.1-1.3, 1.7-1.8

17-654/17-765
Analysis of Software Artifacts
Jonathan Aldrich

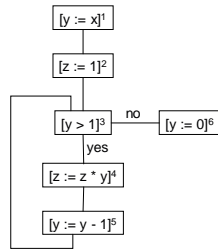
Example WHILE Program

```
[y := x]1;  
[z := 1]2;  
while [y > 1]3 do  
  [z := z * y]4;  
  [y := y - 1]5;  
[y := 0]6;
```

Computes the factorial function, with the input in x and the output in z

Data Flow Analysis

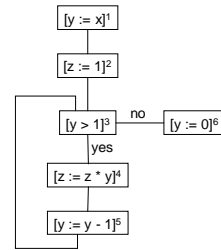
- View program as graph
 - Nodes are elementary blocks like assignments, if statements, etc.
 - Edges show control flow



Data Flow Equations (1)

- Transfer Functions
 - show how a statement affects data flow information

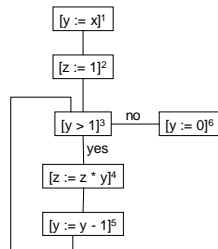
$$\begin{aligned} RD_{exit}(1) &= (RD_{entry}(1) \setminus \{(y, *)\}) \cup (y, 1) \\ RD_{exit}(2) &= (RD_{entry}(2) \setminus \{(z, *)\}) \cup (z, 2) \\ RD_{exit}(3) &= RD_{entry}(3) \\ RD_{exit}(4) &= (RD_{entry}(4) \setminus \{(z, *)\}) \cup (z, 4) \\ RD_{exit}(5) &= (RD_{entry}(5) \setminus \{(y, *)\}) \cup (y, 5) \\ RD_{exit}(6) &= (RD_{entry}(6) \setminus \{(y, *)\}) \cup (y, 6) \end{aligned}$$



Data Flow Equations (1)

- Pattern
 - Assignments
 - kill reaching defs for that var
 - generate new reaching def
 - All others
 - no change

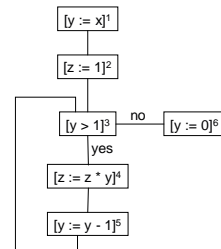
$$\begin{aligned} RD_{exit}(1) &= (RD_{entry}(1) \setminus \{(y, *)\}) \cup (y, 1) \\ RD_{exit}(2) &= (RD_{entry}(2) \setminus \{(z, *)\}) \cup (z, 2) \\ RD_{exit}(3) &= RD_{entry}(3) \\ RD_{exit}(4) &= (RD_{entry}(4) \setminus \{(z, *)\}) \cup (z, 4) \\ RD_{exit}(5) &= (RD_{entry}(5) \setminus \{(y, *)\}) \cup (y, 5) \\ RD_{exit}(6) &= (RD_{entry}(6) \setminus \{(y, *)\}) \cup (y, 6) \end{aligned}$$



Data Flow Equations (2)

- Flow equations
 - Show how analysis information flows from one statement to another

$$\begin{aligned} RD_{entry}(1) &= \{(x, ?), (y, ?), (z, ?)\} \\ RD_{entry}(2) &= RD_{exit}(1) \\ RD_{entry}(3) &= RD_{exit}(2) \cup RD_{exit}(5) \\ RD_{entry}(4) &= RD_{exit}(3) \\ RD_{entry}(5) &= RD_{exit}(4) \\ RD_{entry}(6) &= RD_{exit}(3) \end{aligned}$$



Data Flow Solution

- **Solution**
 - A 12-tuple \overline{RD}
 - Such that $\overline{RD} = F(\overline{RD})$
 - Where F is derived from the equations at right
 - As small (precise) as possible
- **Fixed point of f**
 - A value v such that $v = f(v)$

$$\begin{aligned}
 RD_{exit}(1) &= (RD_{entry}(1) \setminus \{(y, *)\}) \cup (y, 1) \\
 RD_{exit}(2) &= (RD_{entry}(2) \setminus \{(z, *)\}) \cup (z, 2) \\
 RD_{exit}(3) &= RD_{entry}(3) \\
 RD_{exit}(4) &= (RD_{entry}(4) \setminus \{(z, *)\}) \cup (z, 4) \\
 RD_{exit}(5) &= (RD_{entry}(5) \setminus \{(y, *)\}) \cup (y, 5) \\
 RD_{exit}(6) &= (RD_{entry}(6) \setminus \{(y, *)\}) \cup (y, 6) \\
 RD_{entry}(1) &= \{(x, ?), (y, ?), (z, ?)\} \\
 RD_{entry}(2) &= RD_{exit}(1) \\
 RD_{entry}(3) &= RD_{exit}(2) \cup RD_{exit}(5) \\
 RD_{entry}(4) &= RD_{exit}(3) \\
 RD_{entry}(5) &= RD_{exit}(4) \\
 RD_{entry}(6) &= RD_{exit}(3)
 \end{aligned}$$

Equations as Functions

$$\begin{aligned}
 F(RD) &= (F_{entry}(1)(\overline{RD}), RD_{exit}(1)(\overline{RD}), \\
 &\quad \dots, \\
 &\quad F_{entry}(6)(\overline{RD}), RD_{exit}(6)(\overline{RD}))
 \end{aligned}$$

where, for example,

$$F_{entry}(1)(\dots, RD_{entry}(1), \dots) = (RD_{entry}(1) \setminus \{(y, *)\}) \cup (y, 1)$$

$$F_{entry}(3)(\dots, RD_{exit}(2), \dots, RD_{exit}(5), \dots) = RD_{exit}(2) \cup RD_{exit}(5)$$

Computing a Fixed Point of F

- Start with the tuple $\overline{RD}_\emptyset = (\emptyset, \emptyset, \emptyset, \dots, \emptyset)$
- Let $F^0(x) = x$
- Let $F^{n+1}(x) = F(F^n(x))$
- **Surprise!**
 - Now find n such that $F^{n+1}(\overline{RD}_\emptyset) = F^n(\overline{RD}_\emptyset)$
 - By definition $F^n(\overline{RD}_\emptyset)$ is a fixed point of F
 - Does this really work?

Computing the Fixed Point

	0	1	2	3	4	10	$RD_{exit}(1) = (RD_{entry}(1) \setminus \{(y, *)\}) \cup (y, 1)$
$RD_{entry}(1)$	\emptyset	$x?y?z?$				$x?y?z?$	$RD_{exit}(2) = (RD_{entry}(2) \setminus \{(z, *)\}) \cup (z, 2)$
$RD_{exit}(1)$	\emptyset	$y1$	$x?y1z?$			$x?y1z?$	$RD_{exit}(3) = RD_{entry}(3)$
$RD_{entry}(2)$	\emptyset	\emptyset	$y1$	$x?y1z?$		$x?y1z?$	$RD_{exit}(4) = (RD_{entry}(4) \setminus \{(z, *)\}) \cup (z, 4)$
$RD_{exit}(2)$	\emptyset	$z2$	$z2$	$y1z2$		$x?y1z2$	$RD_{exit}(5) = (RD_{entry}(5) \setminus \{(y, *)\}) \cup (y, 5)$
$RD_{entry}(3)$	\emptyset	\emptyset	$y5z2$	$y5z2z4$		$x?y1y5z2z4$	$RD_{exit}(6) = (RD_{entry}(6) \setminus \{(y, *)\}) \cup (y, 6)$
$RD_{exit}(3)$	\emptyset	\emptyset	\emptyset	$y5z2$		$x?y1y5z2z4$	$RD_{entry}(1) = \{(x, ?), (y, ?), (z, ?)\}$
$RD_{entry}(4)$	\emptyset	\emptyset	\emptyset	\emptyset		$x?y1y5z2z4$	$RD_{entry}(2) = RD_{exit}(1)$
$RD_{exit}(4)$	\emptyset	$z4$	$z4$	$z4$		$x?y1y5z4$	$RD_{entry}(3) = RD_{exit}(2) \cup RD_{exit}(5)$
$RD_{entry}(5)$	\emptyset	\emptyset	$z4$	$z4$		$x?y1y5z4$	$RD_{entry}(4) = RD_{exit}(3)$
$RD_{exit}(5)$	\emptyset	$y5$	$y5$	$y5z4$		$x?y5z4$	$RD_{entry}(5) = RD_{exit}(4)$
$RD_{entry}(6)$	\emptyset	\emptyset	\emptyset	\emptyset		$x?y1y5z2z4$	$RD_{entry}(6) = RD_{exit}(3)$
$RD_{exit}(6)$	\emptyset	$y6$	$y6$	$y6$		$x?y6z2z4$	

Finding the Fixed Point

- Why should we think we will find an n such that $F^{n+1}(\overline{RD}_\emptyset) = F^n(\overline{RD}_\emptyset)$?

Monotone Functions

- f is *monotone* if $v \subseteq v'$ implies $f(v) \subseteq f(v')$
- **Assertion:** F is monotone
 - Intuition: preserves input/output relationship
 - Check a couple of cases
 - (1) $RD_{exit}(1) = (RD_{entry}(1) \setminus \{(y, *)\}) \cup (y, 1)$
 - Assume $RD_{entry}(1) \subseteq RD'_{entry}(1)$
 - Then $(RD_{entry}(1) \setminus \{(y, *)\}) \subseteq (RD'_{entry}(1) \setminus \{(y, *)\})$
 - So $RD_{exit}(1) \subseteq RD'_{exit}(1)$
 - Would this also be true if we used \subseteq ?

Monotone Functions

- f is *monotone* if $v \subseteq v'$ implies $f(v) \subseteq f(v')$
- Assertion: F is monotone
 - Intuition: preserves input/output relationship
 - Check a couple of cases
 - (1) $RD_{\text{exit}}(1) = (RD_{\text{entry}}(1) \setminus \{(y, *)\}) \cup (y, 1)$
 - (2) $RD_{\text{entry}}(3) = RD_{\text{exit}}(2) \cup RD_{\text{exit}}(5)$
 - Assume $RD_{\text{exit}}(2) \subseteq RD'_{\text{exit}}(2)$
 - Assume $RD_{\text{exit}}(5) \subseteq RD'_{\text{exit}}(5)$
 - Then $RD_{\text{exit}}(2) \cup RD_{\text{exit}}(5) \subseteq RD'_{\text{exit}}(2) \cup RD'_{\text{exit}}(5)$
 - So $RD_{\text{entry}}(3) \subseteq RD'_{\text{entry}}(3)$

Finding the Fixed Point

- Why should we think we will find an n such that $F^{n+1}(RD_{\emptyset}) = F^n(RD_{\emptyset})$?
 - F is monotone
 - Claim: $\forall n F^n(RD_{\emptyset}) \subseteq F^{n+1}(RD_{\emptyset})$
 - Base case: $RD_{\emptyset} \subseteq F(RD_{\emptyset})$
 - Since no tuple has smaller sets than RD_{\emptyset}
 - Inductive case:
 - Assume $F^n(RD_{\emptyset}) \subseteq F^{n+1}(RD_{\emptyset})$
 - F is monotone, so $F(F^n(RD_{\emptyset})) \subseteq F(F^{n+1}(RD_{\emptyset}))$
 - Equivalently, $F^n(RD_{\emptyset}) \subseteq F^{n+1}(RD_{\emptyset})$
 - Therefore, every application of F either:
 - Does not change RD (and so we have a fixed point)
 - Or increases the size of a set in RD
 - The set of definitions is finite so the sets in RD cannot increase in size forever
 - Therefore the algorithm terminates with a fixed point at some finite n

Precision

- Is $F^n(\overline{RD_{\emptyset}})$ the least fixed point?
 - i.e., the fixed point with the smallest sets?
- Yes. Proof:
 - Consider some other fixed point $\overline{RD_{\text{fix}}}$
 - $\overline{RD_{\emptyset}} \subseteq \overline{RD_{\text{fix}}}$
 - Since F is monotone, $F(\overline{RD_{\emptyset}}) \subseteq F(\overline{RD_{\text{fix}}})$
 - By induction $F^n(\overline{RD_{\emptyset}}) \subseteq F^n(\overline{RD_{\text{fix}}})$
 - But $\overline{RD_{\text{fix}}}$ is a fixed point so $\overline{RD_{\text{fix}}} = F(\overline{RD_{\text{fix}}}) = F^n(\overline{RD_{\text{fix}}})$
 - Therefore $F^n(\overline{RD_{\emptyset}}) \subseteq \overline{RD_{\text{fix}}}$
 - Therefore $F^n(\overline{RD_{\emptyset}})$ is the least fixed point of F

Efficient Algorithms

- Computing $F^n(\overline{RD_{\emptyset}})$ is slow
 - 10 iterations
 - Each iteration recomputes each member of $\overline{RD_{\emptyset}}$
 - Few members of $\overline{RD_{\emptyset}}$ change each iteration
- Optimization: Chaotic Iteration
 - Recompute one member of $\overline{RD_{\emptyset}}$ at a time
 - Guess a member that is likely to change
 - Can compute fixed point in 17 iterations, one recomputation per iteration (vs. 12 before)

Chaotic Iteration

$RD_{1..n} = \emptyset$
 while $RD_j \neq F_j(RD_{1..n})$ for some j
 do $RD_j := F_j(RD_{1..n})$

- How to choose j ?
 - Later!

Chaotic Iteration

$RD_{1..n} = \emptyset$
 while $RD_j \neq F_j(RD_{1..n})$ for some j
 do $RD_j := F_j(RD_{1..n})$

- Properties
 - If chaotic iteration terminates, $RD_{1..n}$ is a fixed point of F
 - Proof: termination implies $RD = F(RD)$
 - That fixed point is a least fixed point
 - Proof: $RD \subseteq F(RD_{\emptyset})$ is an invariant of the algorithm
 - Chaotic iteration terminates for monotone F and finite sets $RD_{1..n}$
 - Proof: F is monotone and so \overline{RD} is increasing

Chaotic Iteration Example

Iter	Position	Value	
0	--	\emptyset	
1	entry(1)	x?y?z?	$RD_{exit}(1) = (RD_{entry}(1) \setminus \{(y,*)\}) \cup (y,1)$
2	exit(1)	x?y1z?	$RD_{exit}(2) = (RD_{entry}(2) \setminus \{(z,*)\}) \cup (z,2)$
3	entry(2)	x?y1z?	$RD_{exit}(3) = RD_{entry}(3)$
4	exit(2)	x?y1z2	$RD_{exit}(4) = (RD_{entry}(4) \setminus \{(z,*)\}) \cup (z,4)$
5	entry(3)	x?y1z2	$RD_{exit}(5) = (RD_{entry}(5) \setminus \{(y,*)\}) \cup (y,5)$
6	exit(3)	x?y1z2	$RD_{exit}(6) = (RD_{entry}(6) \setminus \{(y,*)\}) \cup (y,6)$
7	entry(4)	x?y1z2	
8	exit(4)	x?y1z4	$RD_{entry}(1) = \{(x,?), (y,?), (z,?)\}$
9	entry(5)	x?y1z4	$RD_{entry}(2) = RD_{exit}(1)$
10	exit(5)	x?y5z4	$RD_{entry}(3) = RD_{exit}(2) \cup RD_{exit}(5)$
11	entry(3)	x?y1y5z2z4	$RD_{entry}(4) = RD_{exit}(3)$
12	exit(3)	x?y1y5z2z4	$RD_{entry}(5) = RD_{exit}(4)$
13	entry(4)	x?y1y5z2z4	$RD_{entry}(6) = RD_{exit}(3)$
14	exit(4)	x?y1y5z4	
15	entry(5)	x?y1y5z4	
16	entry(6)	x?y1y5z2z4	
17	exit(6)	x?y6z2z4	

Comparison to Naïve Algorithm

	0	1	2	3	4	10	
$RD_{entry}(1)$	\emptyset	x?y?z?					$RD_{exit}(1) = (RD_{entry}(1) \setminus \{(y,*)\}) \cup (y,1)$
$RD_{exit}(1)$	\emptyset					x?y?z?	$RD_{exit}(2) = (RD_{entry}(2) \setminus \{(z,*)\}) \cup (z,2)$
$RD_{entry}(2)$	\emptyset	y1	x?y1z?				$RD_{exit}(3) = RD_{entry}(3)$
$RD_{exit}(2)$	\emptyset					x?y1z?	$RD_{exit}(4) = (RD_{entry}(4) \setminus \{(z,*)\}) \cup (z,4)$
$RD_{entry}(3)$	\emptyset		z2	y1z2			$RD_{exit}(5) = (RD_{entry}(5) \setminus \{(y,*)\}) \cup (y,5)$
$RD_{exit}(3)$	\emptyset					x?y1z2	$RD_{exit}(6) = (RD_{entry}(6) \setminus \{(y,*)\}) \cup (y,6)$
$RD_{entry}(4)$	\emptyset		y5z2	y5z2z4		x?y1y5z2z4	$RD_{entry}(1) = \{(x,?), (y,?), (z,?)\}$
$RD_{exit}(4)$	\emptyset					x?y1y5z2z4	$RD_{exit}(2) = RD_{exit}(1)$
$RD_{entry}(5)$	\emptyset					x?y1y5z2z4	$RD_{entry}(3) = RD_{exit}(2) \cup RD_{exit}(5)$
$RD_{exit}(5)$	\emptyset					x?y1y5z4	$RD_{exit}(4) = RD_{exit}(3)$
$RD_{entry}(6)$	\emptyset					x?y1y5z4	$RD_{entry}(5) = RD_{exit}(4)$
$RD_{exit}(6)$	\emptyset					x?y6z2z4	$RD_{entry}(6) = RD_{exit}(3)$

Constant Folding

- A program optimization
 - Replaces computation with constants
 - Can use reaching definitions
- Notation
 - $RD \vdash S \triangleright S'$
 - Given reaching definitions RD, statement S can be transformed into S'
 - $\frac{C}{T}$
 - Transformation T is legal if condition(s) C hold
 - FV(exp)
 - The variables mentioned in exp
 - exp[y→n]
 - Replace all occurrences of y in exp with n

Constant Folding Rules

$$\begin{aligned}
 [ass_1] \quad & RD \vdash [x := n]^e \triangleright [x := n]^e \mapsto n]^e \\
 & \text{if } \left\{ \begin{array}{l} y \in FV(n) \wedge (y,?) \notin RD_{entry}(e) \wedge \\ V(n, e) \in RD_{entry}(e) : (x = y \Rightarrow \dots)^e \text{ is } [y := n]^e \end{array} \right\} \\
 [ass_2] \quad & RD \vdash [x := n]^e \triangleright [x := n]^e \\
 & \text{if } FV(n) = \emptyset \wedge n \notin \text{NUM} \wedge n \text{ evaluates to } n \\
 [seq_1] \quad & \frac{RD \vdash S_1 \triangleright S'_1}{RD \vdash S_1; S_2 \triangleright S'_1; S'_2} \\
 [seq_2] \quad & \frac{RD \vdash S_2 \triangleright S'_2}{RD \vdash S_1; S_2 \triangleright S_1; S'_2} \\
 [if] \quad & \frac{RD \vdash S_2 \triangleright S'_2}{RD \vdash \text{if } [b]^e \text{ then } S_1 \text{ else } S_2 \triangleright \text{if } [b]^e \text{ then } S'_1 \text{ else } S'_2} \\
 [if] \quad & \frac{RD \vdash S_2 \triangleright S'_2}{RD \vdash \text{if } [b]^e \text{ then } S_1 \text{ else } S_2 \triangleright \text{if } [b]^e \text{ then } S_1 \text{ else } S'_2} \\
 [wh] \quad & \frac{RD \vdash S \triangleright S'}{RD \vdash \text{while } [b]^e \text{ do } S \triangleright \text{while } [b]^e \text{ do } S'}
 \end{aligned}$$

Taken from Nielson, Nielson, and Hankin, page 27

Example

$[x:=10]^1; [y:=x+10]^2; [z:=y+10]^3$

- $RD_{enter}(2) = \{(x,1), (y,?), (z,?)\}$
- $RD \vdash [y:=x+10]^2 \triangleright [y:=10+10]^2$ by [ass₁]
- Thus:
 - $RD \vdash [x:=10]^1; [y:=x+10]^2; [z:=y+10]^3 \triangleright [x:=10]^1; [y:=10+10]^2; [z:=y+10]^3$ by [seq] rules

Example

$RD \vdash [x:=10]^1; [y:=x+10]^2; [z:=y+10]^3$
 $\triangleright [x:=10]^1; [y:=10+10]^2; [z:=y+10]^3$ by [ass₁]
 $\triangleright [x:=10]^1; [y:=20]^2; [z:=y+10]^3$ by [ass₂]
 $\triangleright [x:=10]^1; [y:=20]^2; [z:=20+10]^3$ by [ass₁]
 $\triangleright [x:=10]^1; [y:=20]^2; [z:=30]^3$ by [ass₂]