

Lecture 16 – Dataflow Analysis

THEORY OF COMPILATION

Eran Yahav

www.cs.technion.ac.il/~yahave/tocs2011/compiler-lec16.pptx

Reference: Dragon 9, 12

1

Last time... Dataflow Analysis

- Information flows along (potential) execution paths
- Conservative approximation of all possible program executions
- Can be viewed as a sequence of transformations on program state
 - Every statement (block) is associated with two abstract states: input state, output state
 - Input/output state represents all possible states that can occur at the program point
 - Representation is finite
 - Different problems typically use different representations

2

Control-Flow Graph

```

1: y := X;
2: Z := 1;
3: while y > 0 {
4:   Z := Z * y;
5:   y := y - 1
}
6: y := 0
    
```

```

graph TD
    1[1: y := X] --> 2[2: Z := 1]
    2 --> 3{3: y > 0}
    3 --> 4[4: Z := Z * y]
    4 --> 5[5: y := y - 1]
    5 --> 3
    3 --> 6[6: y := 0]
    
```

3

Executions

```

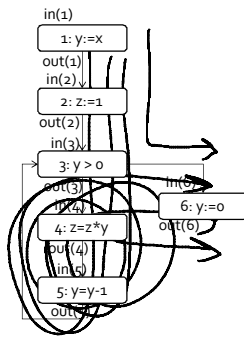
1: y := X;
2: Z := 1;
3: while y > 0 {
4:   Z := Z * y;
5:   y := y - 1
}
6: y := 0
    
```

4

Input/output Sets

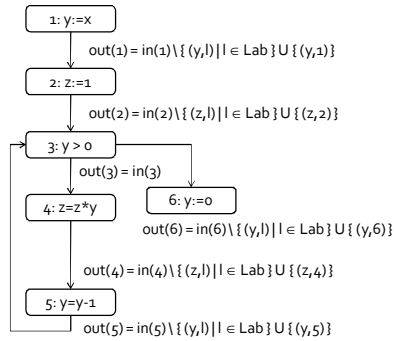
```

1: y := x;
2: z := 1;
3: while y > 0 {
4: z := z * y;
5: y := y - 1
}
6: y := 0
    
```



5

Transfer Functions



$in(1) = \{(x,?), (y,?), (z,?)\}$
 $in(2) = out(1)$
 $in(3) = out(2) \cup out(5)$
 $in(4) = out(3)$
 $in(5) = out(4)$
 $in(6) = out(3)$

6

Kill/Gen formulation for Reaching Definitions

Block	out (lab)
$[x := a]^{lab}$	$in(lab) \setminus \{(x,l) l \in Lab\} \cup \{(x,lab)\}$
$[skip]^{lab}$	$in(lab)$
$[b]^{lab}$	$in(lab)$

Block	kill	gen
$[x := a]^{lab}$	$\{(x,l) l \in Lab\}$	$\{(x,lab)\}$
$[skip]^{lab}$	\emptyset	\emptyset
$[b]^{lab}$	\emptyset	\emptyset

For each program point, which assignments may have been made and not overwritten, when program execution reaches this point along some path.

7

Solving Gen/Kill Equations

$OUT[ENTRY] = \emptyset;$
 for (each basic block B other than ENTRY) $OUT[B] = \emptyset;$
 while (changes to any OUT occur) {
 for (each basic block B other than ENTRY) {
 $OUT[B] = (IN[B] \setminus killB) \cup genB$
 $IN[B] = \cup_{p \in pred(B)} OUT[p]$
 }
 }

- Designated block Entry with $OUT[Entry] = \emptyset$
- $pred(B)$ = predecessor nodes of B in the control flow graph

8

Available Expressions Analysis

```
[x := a+b]1;  
[y := a*b]2;  
while [y > a+b]3(  
  [a := a + 1]4;  
  [x := a + b]5  
)
```

(a+b) always available at label 3

For each program point, which expressions must have already been computed, and not later modified, on all paths to the program point

9

Some required notation

blocks : Stmt \rightarrow P(Blocks)
 blocks([x := a]^{lab}) = {[x := a]^{lab}}
 blocks([skip]^{lab}) = {[skip]^{lab}}
 blocks(S₁; S₂) = blocks(S₁) \cup blocks(S₂)
 blocks(if [b]^{lab} then S₁ else S₂) = {[b]^{lab}} \cup blocks(S₁) \cup blocks(S₂)
 blocks(while [b]^{lab} do S) = {[b]^{lab}} \cup blocks(S)

FV: (BExp \cup AExp) \rightarrow Var
 Variables used in an expression

AExp(a) = all non-unit expressions in the arithmetic expression a
 similarly AExp(b) for a boolean expression b

10

Available Expressions Analysis

- Property space
 - in_{AE}, out_{AE}: Lab \rightarrow P(AExp)
 - Mapping a label to set of arithmetic expressions available at that label
- Dataflow equations
 - Flow equations – how to join incoming dataflow facts
 - Effect equations - given an input set of expressions S, what is the effect of a statement

11

Available Expressions Analysis

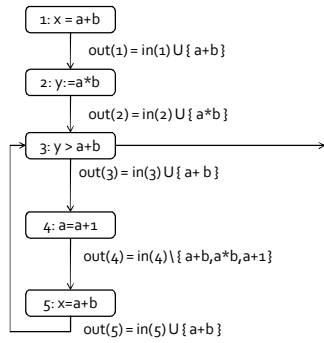
- in_{AE}(lab) =
 - \emptyset when lab is the initial label
 - $\bigcap \{out_{AE}(lab') \mid lab' \in pred(lab)\}$ otherwise
- out_{AE}(lab) = ...

Block	out (lab)
[x := a] ^{lab}	in(lab) \ { a' \in AExp x \in FV(a') } \cup { a' \in AExp(a) x \notin FV(a') }
[skip] ^{lab}	in(lab)
[b] ^{lab}	in(lab) \cup AExp(b)

From now on going to drop the AE subscript when clear from context

12

Transfer Functions



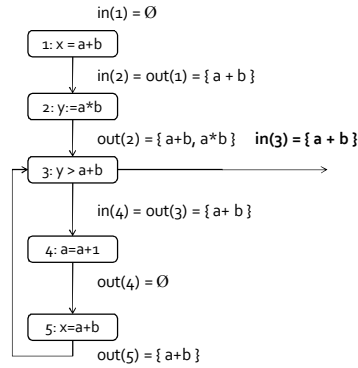
$in(1) = \emptyset$
 $in(2) = out(1)$
 $in(3) = out(2) \cap out(5)$
 $in(4) = out(3)$
 $in(5) = out(4)$

```

[x := a + b];
[y := a * b];
while [y > a + b] (
  [a := a + 1];
  [x := a + b];
)
    
```

13

Solution



$in(1) = \emptyset$
 $in(2) = out(1) = \{a + b\}$
 $out(2) = \{a + b, a * b\}$ $in(3) = \{a + b\}$

$in(4) = out(3) = \{a + b\}$
 $out(4) = \emptyset$
 $out(5) = \{a + b\}$

14

Kill/Gen

Block	out (lab)
$[x := a]^{lab}$	$in(lab) \setminus \{a' \in AExp \mid x \in FV(a')\} \cup \{a' \in AExp(a) \mid x \notin FV(a')\}$
$[skip]^{lab}$	$in(lab)$
$[b]^{lab}$	$in(lab) \cup AExp(b)$

Block	kill	gen
$[x := a]^{lab}$	$\{a' \in AExp \mid x \in FV(a')\}$	$\{a' \in AExp(a) \mid x \notin FV(a')\}$
$[skip]^{lab}$	\emptyset	\emptyset
$[b]^{lab}$	\emptyset	$AExp(b)$

$$out(lab) = in(lab) \setminus kill(B^{lab}) \cup gen(B^{lab})$$

B^{lab} = block at label lab

15

Reaching Definitions Revisited

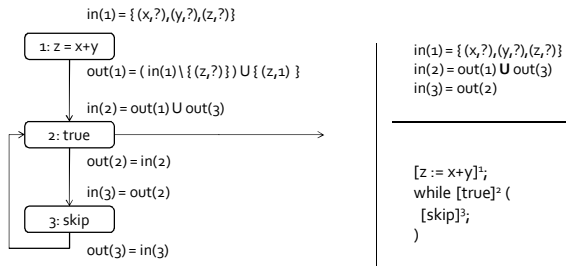
Block	out (lab)
$[x := a]^{lab}$	$in(lab) \setminus \{(x,l) \mid l \in Lab\} \cup \{(x,lab)\}$
$[skip]^{lab}$	$in(lab)$
$[b]^{lab}$	$in(lab)$

Block	kill	gen
$[x := a]^{lab}$	$\{(x,l) \mid l \in Lab\}$	$\{(x,lab)\}$
$[skip]^{lab}$	\emptyset	\emptyset
$[b]^{lab}$	\emptyset	\emptyset

For each program point, which assignments may have been made and not overwritten, when program execution reaches this point along some path.

16

Why solution with smallest sets?



After simplification: $in(2) = in(2) \cup \{(x,?), (y,?), (z,1)\}$
 Many solutions: any superset of $\{(x,?), (y,?), (z,1)\}$

Live Variables

```

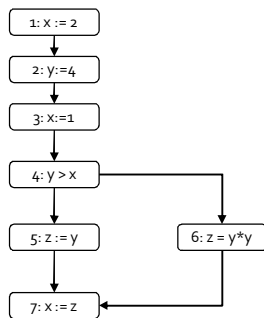
[X := 2]1;
[Y := 4]2;
[X := 1]3;
(if [y > x]4 then [z := y]5
 else [z := y * y]6);
[X := z]7
    
```

For each program point, which variables may be live at the exit from the point.

Live Variables

```

[X := 2]1;
[Y := 4]2;
[X := 1]3;
(if [y > x]4 then [z := y]5
 else [z := y * y]6);
[X := z]7
    
```

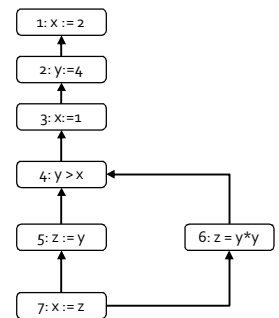


Live Variables

```

[X := 2]1;
[Y := 4]2;
[X := 1]3;
(if [y > x]4 then [z := y]5
 else [z := y * y]6);
[X := z]7
    
```

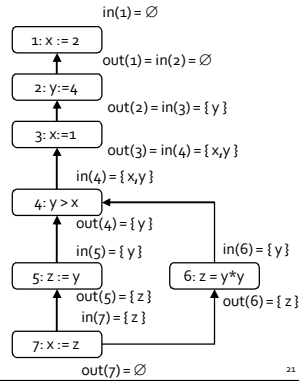
Block	kill	gen
$[x := a]^{ab}$	$\{x\}$	$\{FV(a)\}$
$[skip]^{ab}$	\emptyset	\emptyset
$[b]^{ab}$	\emptyset	$FV(b)$



Live Variables: solution

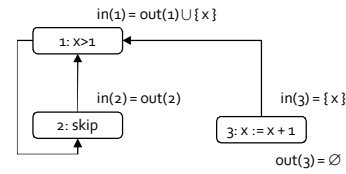
```
[x := 2]1;
[y := 4]2;
[x := 1]3;
(if [y > x]4 then [z := y]5
 else [z := y * y]6);
[x := z]7
```

Block	kill	gen
[x := a] ^{lab}	{ x }	{ FV(a) }
[skip] ^{lab}	∅	∅
[b] ^{lab}	∅	FV(b)



21

Why solution with smallest set?



out(1) = in(2) ∪ in(3)
out(2) = in(1)
out(3) = ∅

```
while [x > 1]2 (
  [skip]3;
  [x := x + 1]3;
)
```

After simplification: in(1) = in(1) ∪ { x }

Many solutions: any superset of { x }

22

Monotone Frameworks

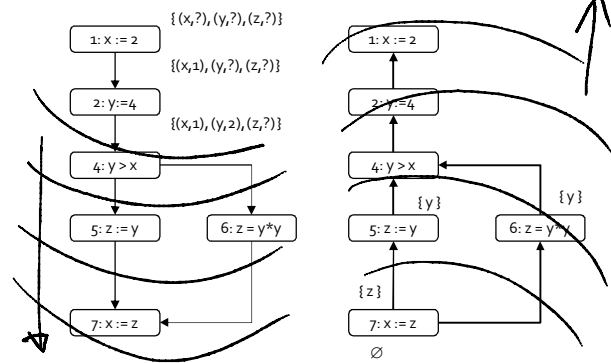
$$In(lab) = \begin{cases} \text{Initial} & \text{when } lab \in \text{Entry labels} \\ \sqcup \{ out(lab') \mid (lab', lab) \in \text{CFG edges} \} & \text{otherwise} \end{cases}$$

$$out(lab) = f_{lab}(in(lab))$$

- \sqcup is \cup or \cap
- CFG edges go either forward or backwards
- Entry labels are either initial program labels or final program labels (when going backwards)
- Initial is an initial state (or final when going backwards)
- f_{lab} is the transfer function associated with the blocks B^{lab}

23

Forward vs. Backward Analyses



24

Must vs. May Analyses

- When \sqsubseteq is \cap - must analysis
 - Want largest sets that solves the equation system
 - Properties hold on all paths reaching a label (exiting a label, for backwards)
- When \sqsubseteq is \cup - may analysis
 - Want smallest sets that solve the equation system
 - Properties hold at least on one path reaching a label (existing a label, for backwards)

25

Example: Reaching Definition

- $L = \wp(\text{Var} \times \text{Lab})$ is partially ordered by \subseteq
- \sqsubseteq is \cup
- L satisfies the Ascending Chain Condition because $\text{Var} \times \text{Lab}$ is finite (for a given program)

26

Example: Available Expressions

- $L = \wp(\text{AExp})$ is partially ordered by \supseteq
- \sqsubseteq is \cap
- L satisfies the Ascending Chain Condition because AExp is finite (for a given program)

27

Analyses Summary

	Reaching Definitions	Available Expressions	Live Variables
L	$\wp(\text{Var} \times \text{Lab})$	$\wp(\text{AExp})$	$\wp(\text{Var})$
\subseteq	\subseteq	\supseteq	\subseteq
\sqcup	\cup	\cap	\cup
\perp	\emptyset	AExp	\emptyset
Initial	$\{(x,?) \mid x \in \text{Var}\}$	\emptyset	\emptyset
Entry labels	$\{\text{init}\}$	$\{\text{init}\}$	final
Direction	Forward	Forward	Backward
F	$\{f: L \rightarrow L \mid \exists k, g: f(\text{val}) = (\text{val} \setminus k) \cup g\}$		
f_{lab}	$f_{\text{lab}}(\text{val}) = (\text{val} \setminus \text{kill}) \cup \text{gen}$		

28

Analyses as Monotone Frameworks

- Property space
 - Powerset
 - Clearly a complete lattice
- Transformers
 - Kill/gen form
 - Monotone functions (let's show it)

29

Monotonicity of Kill/Gen transformers

- Have to show that $x \sqsubseteq x'$ implies $f(x) \sqsubseteq f(x')$
- Assume $x \sqsubseteq x'$, then for kill set k and gen set g

$$(x \setminus k) \cup g \sqsubseteq (x' \setminus k) \cup g$$
- Technically, since we want to show it for all functions in F , we also have to show that the set is closed under function composition

30

Distributivity of Kill/Gen transformers

- Have to show that $f(x \sqcup y) \sqsubseteq f(x) \sqcup f(y)$
- $$\begin{aligned} f(x \sqcup y) &= ((x \sqcup y) \setminus k) \cup g \\ &= ((x \setminus k) \sqcup (y \setminus k)) \cup g \\ &= (((x \setminus k) \cup g) \sqcup ((y \setminus k) \cup g)) \\ &= f(x) \sqcup f(y) \end{aligned}$$
- Used distributivity of \sqcup and \cup
 - Works regardless of whether \sqcup is \cup or \cap

31

Points-to Analysis

- Many flavors
- PWHILE language

$p \in \text{PExp}$ pointer expressions

$a ::= x \mid n \mid a_1 \text{ op}_a a_2 \mid \&x \mid *x \mid \text{nil}$

$S ::= [x := a]^{lab}$
 $\mid [\text{skip}]^{lab}$
 $\mid S_1; S_2$
 $\mid \text{if } [b]^{lab} \text{ then } S_1 \text{ else } S_2$
 $\mid \text{while } [b]^{lab} \text{ do } S$
 $\mid x = \text{malloc}$

32

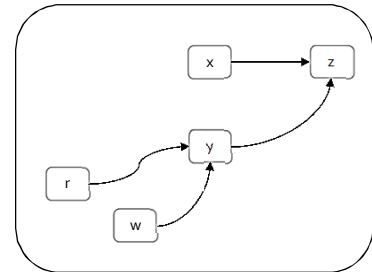
Points-to Analysis

- Aliases
 - Two pointers p and q are aliases if they point to the same memory location
- Points-to pair
 - (p,q) means p holds the address of q
- Points-to pairs and aliases
 - (p,q) and (r,q) means that p and r are aliases
- Challenge: no a priori bound on the set of heap locations

33

Terminology Example

[x := &z]¹
 [y := &z]²
 [w := &y]³
 [r := w]⁴



Points-to pairs: (x,z), (y,z), (w,y), (r,y)
 Aliases: (x,y), (r,w)

34

(May) Points-to Analysis

- Property Space
 - $L = (\emptyset, \text{Var} \times \text{Var}), \subseteq, \cup, \cap, \emptyset, \text{Var} \times \text{Var}$
- Transfer functions

Statement	out(lab)
[p = &x] ^{lab}	in(lab) ∪ { (p,x) }
[p = q] ^{lab}	in(lab) ∪ { (p,x) (q,x) ∈ in(lab) }
[*p = q] ^{lab}	in(lab) ∪ { (r,x) (q,x) ∈ in(lab) and (p,r) ∈ in(lab) }
[p = *q] ^{lab}	in(lab) ∪ { (p,r) (q,x) ∈ in(lab) and (x,r) ∈ in(lab) }

35

(May) Points-to Analysis

- What to do with malloc?
- Need some static naming scheme for dynamically allocated objects
- Single name for the entire heap
 - $\llbracket [p = \text{malloc}]^{\text{lab}} \rrbracket (S) = S \cup \{ (p,H) \}$
- Name based on static allocation site
 - $\llbracket [p = \text{malloc}]^{\text{lab}} \rrbracket (S) = S \cup \{ (p,\text{lab}) \}$

36

(May) Points-to Analysis

```

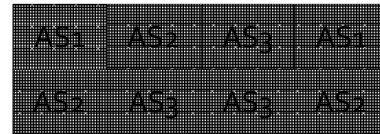
_____ ∅
[x :=malloc]1;
_____ { (x,H) }
[y:=malloc]2;
_____ { (x,H), (y,H) }
(if [x==y]3 then
_____ { (x,H), (y,H) }
  [z:=x]4
_____ { (x,H), (y,H), (z,H) }
else
  [z:=y]5
_____ { (x,H), (y,H), (z,H) }
);
_____ { (x,H), (y,H), (z,H) }
    
```

Single name H for the entire heap

37

Allocation Sites

- Divide the heap into a fixed partition based on allocation site
- All objects allocated at the same program point represented by a single “abstract object”



38

(May) Points-to Analysis

```

_____ ∅
[x :=malloc]1; // A1
_____ { (x,A1) }
[y:=malloc]2; // A2
_____ { (x,A1), (y,A2) }
(if [x==y]3 then
_____ { (x,A1), (y,A2) }
  [z:=x]4
_____ { (x,A1), (y,A2), (z,A1) }
else
  [z:=y]5
_____ { (x,A1), (y,A2), (z,A2) }
);
_____ { (x,A1), (y,A2), (z,A1), (z,A2) }
    
```

Allocation-site based naming (Using A_{lab} instead of just “lab” for clarity)

39

Weak Updates

Statement	out(lab)
[p = &x] ^{lab}	in(lab) ∪ { (p,x) }
[p = q] ^{lab}	in(lab) ∪ { (p,x) (q,x) ∈ in(lab) }
[*p = q] ^{lab}	in(lab) ∪ { (r,x) (q,x) ∈ in(lab) and (p,r) ∈ in(lab) }
[p = *q] ^{lab}	in(lab) ∪ { (p,r) (q,x) ∈ in(lab) and (x,r) ∈ in(lab) }

```

_____ ∅
[x :=malloc]1; // A1
_____ { (x,A1) }
[y:=malloc]2; // A2
_____ { (x,A1), (y,A2) }
[z:=x]3;
_____ { (x,A1), (y,A2), (z,A1) }
[z:=y]4;
_____ { (x,A1), (y,A2), (z,A1), (z,A2) }
    
```

40

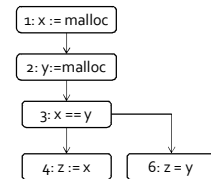
(May) Points-to Analysis

- Fixed partition of the (unbounded) heap to static names
 - Allocation sites
 - Types
 - Calling contexts
 - ...
- What we saw so far – flow-insensitive
 - Ignoring the structure of the flow in the program

43

Flow-sensitive vs. Flow-insensitive Analyses

```
[x := malloc]1;
[y := malloc]2;
(if [x==y]3 then
  [z:=x]4
else
  [z:=y]5
);
```



- Flow sensitive: respect program flow
 - a separate set of points-to pairs for every program point
 - the set at a point represents possible may-aliases on some path from entry to the program point
- Flow insensitive: assume all execution orders are possible, abstract away order between statements

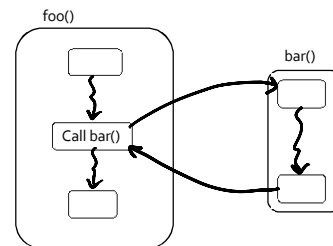
42

So far...

- Intra-procedural analysis
- How are we going to deal with procedures?
- Inter-procedural analysis

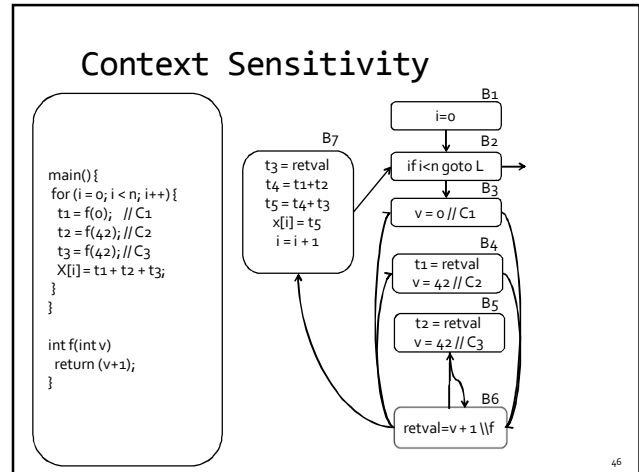
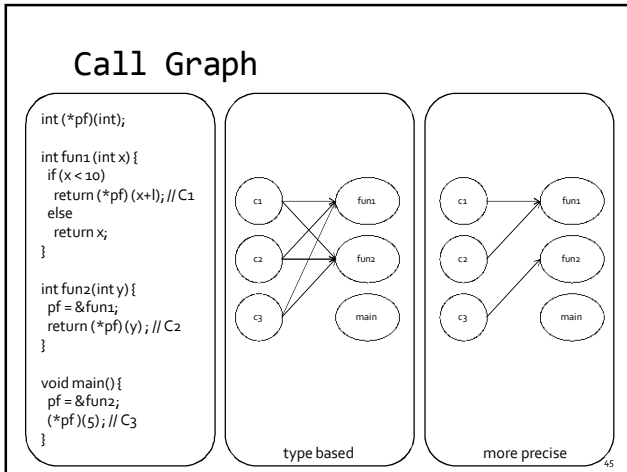
43

Interprocedural Analysis



- The effect of calling a procedure is the effect of executing its body

44

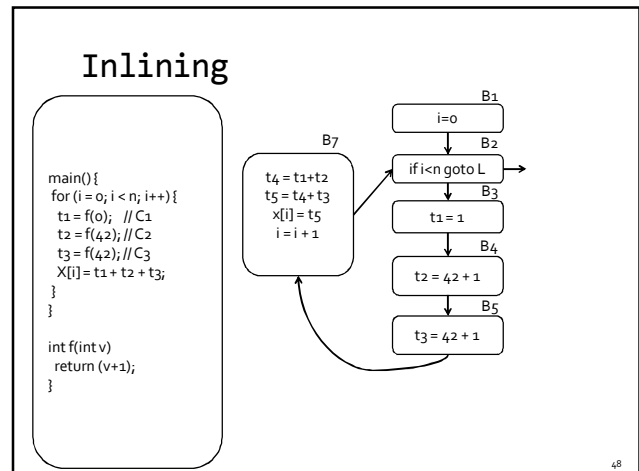


Solution Attempt #1

- Inline callees into callers
 - End up with one big procedure
 - CFGs of individual procedures = duplicated many times
- Good: it is precise
 - distinguishes different calls to the same function
- Bad
 - exponential blow-up, not efficient
 - doesn't work with recursion

```

main() { f(0); f(0); }
f() { g(0); g(0); }
g() { h(0); h(0); }
h() { ... }
    
```

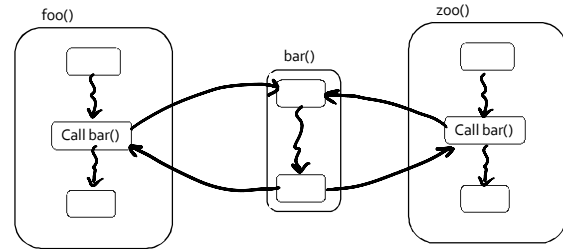


Solution Attempt #2

- Build a "supergraph" = inter-procedural CFG
- Replace each call from P to Q with
 - An edge from point before the call (call point) to Q's entry point
 - An edge from Q's exit point to the point after the call (return pt)
 - Add assignments of actuals to formals, and assignment of return value
- Good: efficient
 - Graph of each function included exactly once in the supergraph
 - Works for recursive functions (although local variables need additional treatment)
- Bad: imprecise, "context-insensitive"
 - The "unrealizable paths problem": dataflow facts can propagate along infeasible control paths

49

Unrealizable Paths



50

Interprocedural Analysis

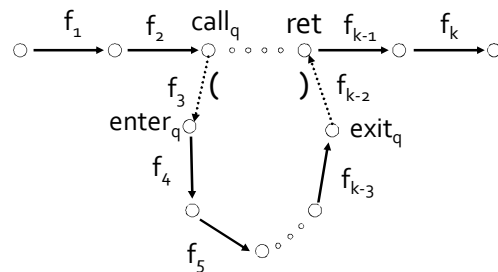
```

begin
proc p() is1
  [x := a + 1]2
end3
[a = 7]4
[call p()]5
[print x]7
[a = 9]8
[call p()]9
[print a]11
end
    
```

- Extend language with begin/end and with [call p()]^{clab}_{rlab}
- Call label clab, and return label rlab

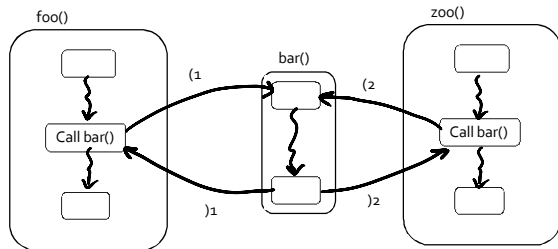
51

IVP: Interprocedural Valid Paths



- IVP: all paths with matching calls and returns
- And prefixes

Valid Paths



53

Interprocedural Valid Paths

- IVP set of paths
 - Start at program entry
- Only considers matching calls and returns
 - aka, valid
- Can be defined via context free grammar
 - $\text{matched} ::= \text{matched } (, \text{matched }) | \epsilon$
 - $\text{valid} ::= \text{valid } (, \text{matched } | \text{matched}$
 - *paths* can be defined by a regular expression

The Join-Over-Valid-Paths (JVP)

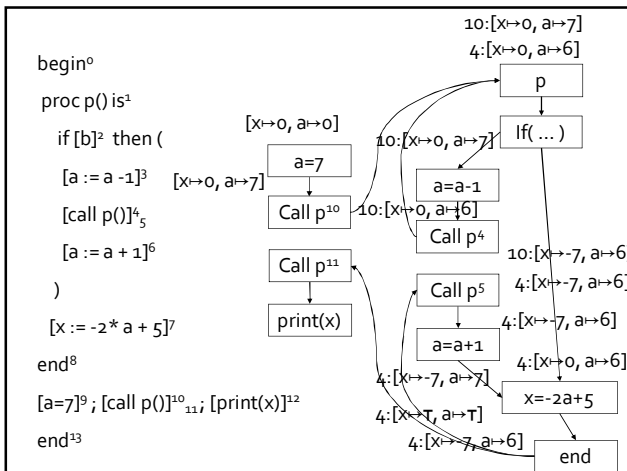
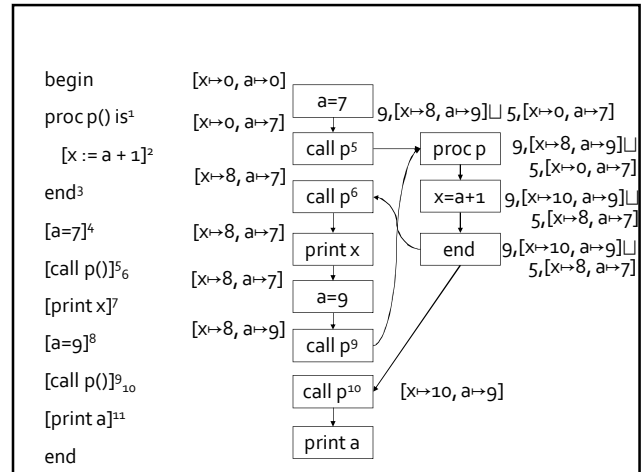
- $\text{vpaths}(n)$ all valid paths from program start to n
- $\text{JVP}[n] = \sqcup \{ \llbracket e_{x_1} e_{x_2} \dots, e \rrbracket \mid (e_{x_1} e_{x_2} \dots, e) \in \text{vpaths}(n) \}$ (initial)
- $\text{JVP} \sqsubseteq \text{JFP}$
 - In some cases the JVP can be computed
 - (Distributive problem)

Sharir and Pnueli '82

- Call String approach
 - Blend interprocedural flow with intra procedural flow
 - Tag every dataflow fact with call history
- Functional approach
 - Determine the effect of a procedure
 - E.g., in/out map
 - Treat procedure invocations as "super ops"

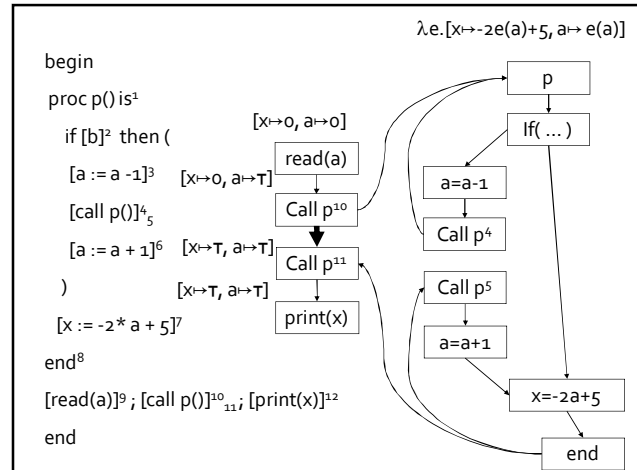
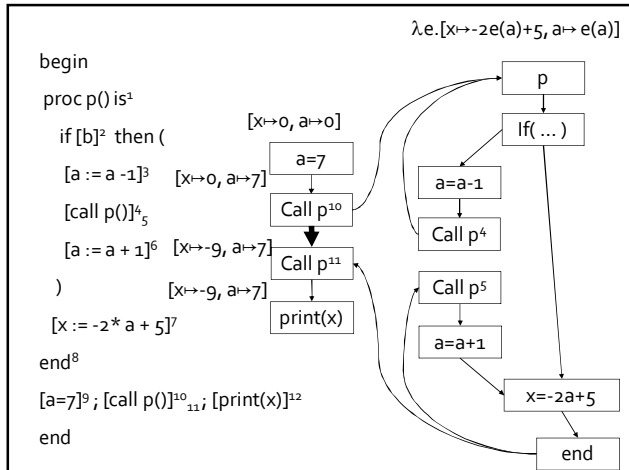
The Call String Approach

- Record at every node a pair (l, c) where l ∈ L is the dataflow information and c is a suffix of unmatched calls
- Use Chaotic iterations
- To guarantee termination limit the size of c (typically 1 or 2)
- Emulates inline (but no code growth)
- Exponential in size of c



The Functional Approach

- The meaning of a procedure is mapping from states into states
- The abstract meaning of a procedure is function from an abstract state to abstract states



Functional Approach: Main Idea

- Iterate on the abstract domain of functions from L to L
- Two phase algorithm
 - Compute the dataflow solution at the exit of a procedure as a function of the initial values at the procedure entry (functional values)
 - Compute the dataflow values at every point using the functional values
- Computes JVP for distributive problems