# Dataflow Analysis
# Widening and Narrowing
# Path Sensitivity
# Interprocedural Analysis

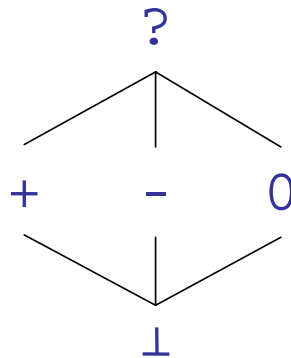## Static Analysis 2009

Michael I. Schwartzbach

Computer Science, University of Aarhus

# Sign Analysis

- Determine the sign (`+`,`-`,`0`) of all expressions
- The *Sign* lattice:

```
        ?
       /|\
      / | \
     +  -  0
      \ | /
       \|/
        ⊥
```

- The full lattice is the map lattice: *Vars → Sign*
  - where *Vars* is the set of variables in the program

# Sign Constraints

- The variable [[v]] denotes a map that gives the sign value for all variables at the program point *after* v

- For variable declarations:

  $$[[v]] = [id_1 \rightarrow ?, ..., id_n \rightarrow ?]$$

- For assignments:

  $$[[v]] = JOIN(v)[id \rightarrow eval(JOIN(v),E)$$

- For all other nodes:

  $$[[v]] = JOIN(v) = \bigsqcup_{w \in pred(v)} [[w]]$$

# Evaluating Signs

- The *eval* function is an *abstract evaluation*:
  - $eval(\sigma, id) = \sigma(id)$
  - $eval(\sigma, intconst) = sign(intconst)$
  - $eval(\sigma, E_1 \text{ op } E_2) = \overline{op}(eval(\sigma, E_1), eval(\sigma, E_2))$

- The *sign* function gives the sign of an integer

- The $\overline{op}$ function is an abstract evaluation of the given operator

# Abstract Operators

| + | ⊥ | 0 | - | + | ? |
|---|---|---|---|---|---|
| ⊥ | ⊥ | ⊥ | ⊥ | ⊥ | ⊥ |
| 0 | ⊥ | 0 | - | + | ? |
| - | ⊥ | - | - | ? | ? |
| + | ⊥ | + | ? | + | ? |
| ? | ⊥ | ? | ? | ? | ? |

| - | ⊥ | 0 | - | + | ? |
|---|---|---|---|---|---|
| ⊥ | ⊥ | ⊥ | ⊥ | ⊥ | ⊥ |
| 0 | ⊥ | 0 | + | - | ? |
| - | ⊥ | - | ? | - | ? |
| + | ⊥ | + | + | ? | ? |
| ? | ⊥ | ? | ? | ? | ? |

| * | ⊥ | 0 | - | + | ? |
|---|---|---|---|---|---|
| ⊥ | ⊥ | 0 | ⊥ | ⊥ | ⊥ |
| 0 | 0 | 0 | 0 | 0 | 0 |
| - | ⊥ | 0 | + | - | ? |
| + | ⊥ | 0 | - | + | ? |
| ? | ⊥ | 0 | ? | ? | ? |

| / | ⊥ | 0 | - | + | ? |
|---|---|---|---|---|---|
| ⊥ | ⊥ | ⊥ | ⊥ | ⊥ | ⊥ |
| 0 | ⊥ | ? | 0 | 0 | ? |
| - | ⊥ | ? | ? | ? | ? |
| + | ⊥ | ? | ? | ? | ? |
| ? | ⊥ | ? | ? | ? | ? |

| > | ⊥ | 0 | - | + | ? |
|---|---|---|---|---|---|
| ⊥ | ⊥ | ⊥ | ⊥ | ⊥ | ⊥ |
| 0 | ⊥ | 0 | + | 0 | ? |
| - | ⊥ | 0 | ? | 0 | ? |
| + | ⊥ | + | + | ? | ? |
| ? | ⊥ | ? | ? | ? | ? |

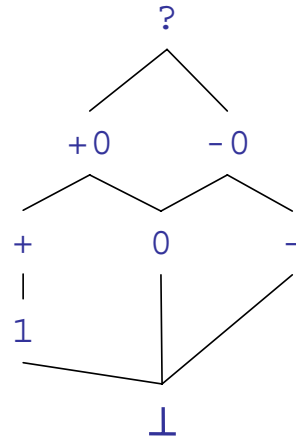| == | ⊥ | 0 | - | + | ? |
|---|---|---|---|---|---|
| ⊥ | ⊥ | ⊥ | ⊥ | ⊥ | ⊥ |
| 0 | ⊥ | + | 0 | 0 | ? |
| - | ⊥ | 0 | ? | 0 | ? |
| + | ⊥ | 0 | 0 | ? | ? |
| ? | ⊥ | ? | ? | ? | ? |

# **Monotonicity**

- The $\sqcup$ operator and map updates are monotone
- Compositions preserve monotonicity
- Are the abstract operators monotone?

- This is verified by a tedious manual inspection
- Or better, run an O($n^3$) algorithm for an $n{\times}n$ table:
  - $\forall x,y,x' \in L$: $x \sqsubseteq x' \Rightarrow x \overline{op} y \sqsubseteq x' \overline{op} y$
  - $\forall x,y,y' \in L$: $y \sqsubseteq y' \Rightarrow x \overline{op} y \sqsubseteq x \overline{op} y'$
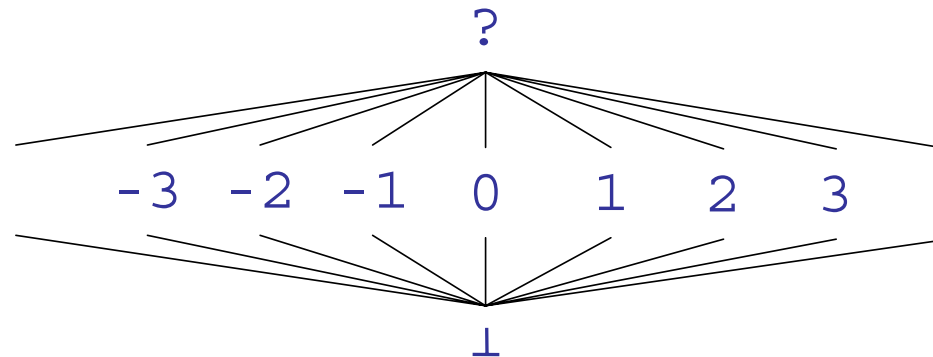
# Increasing Precision

- **Some loss of information:**
  - `(2>0)==1` is analyzed as `?`
  - `+/+` is analyzed as `?`, since e.g. ½ is rounded down

- **Use a richer lattice for better precision:**

```
            ?
          /   \
       +0       -0
        \  \   /  /
     +      0      -
     |      |    /
     1       \  /
              \/
              ⊥
```

- **Abstract operators are now 8×8 tables**

# Constant Propagation

- Determine variables with a constant value
- Similar to sign analysis, with basic lattice:

$$?$$

$$-3 \quad -2 \quad -1 \quad 0 \quad 1 \quad 2 \quad 3$$

$$\bot$$

- Abstract operator for addition:

$$\overline{+}(n,m) = \texttt{if } (n{\neq}? \wedge m{\neq}?) \{ n{+}m \} \texttt{ else } \{ ? \}$$

# Constant Folding

- Exploiting constant propagation:

```
var x,y,z;
x = 27;
y = input,
z = 2*x+y;
if (x<0) { y=z-3; } else { y=12 }
output y;
```

```
var x,y,z;
x = 27;
y = input;
z = 54+y;
if (0) { y=z-3; } else { y=12 }
output y;
```

```
var y;
y = input;
output 12;
```

9

# Interval Analysis

- Compute upper and lower bounds for integers
- Lattice of intervals:

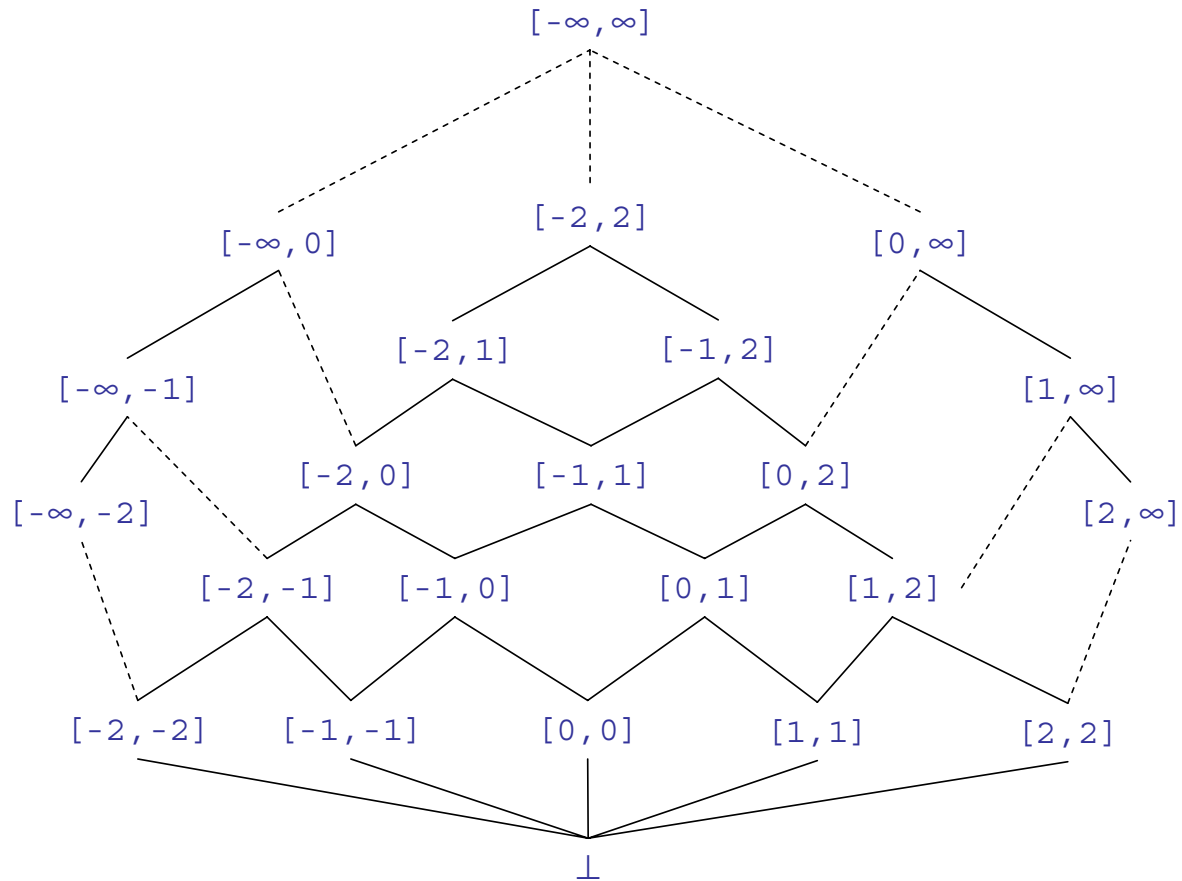$$Interval = lift(\{ \ [l,h] \ | \ l,h \in N \wedge l \leq h \ \})$$

where:

$$N = \{-\infty, \ ..., \ -2, \ -1, \ 0, \ 1, \ 2, \ ..., \ \infty\}$$

and intervals are ordered by inclusion:

$$[l_1, h_1] \sqsubseteq [l_2, h_2] \ \text{iff} \ l_2 \leq l_1 \wedge h_1 \leq h_2$$

# The Interval Lattice



The Interval Lattice diagram:

- $[-\infty,\infty]$
- $[-\infty,0]$    $[-2,2]$    $[0,\infty]$
- $[-\infty,-1]$    $[-2,1]$    $[-1,2]$    $[1,\infty]$
- $[-\infty,-2]$    $[-2,0]$    $[-1,1]$    $[0,2]$    $[2,\infty]$
- $[-2,-1]$    $[-1,0]$    $[0,1]$    $[1,2]$
- $[-2,-2]$    $[-1,-1]$    $[0,0]$    $[1,1]$    $[2,2]$
- $\bot$

# Interval Analysis Lattice

- The total lattice for a program point is:

  $$L = \text{Vars} \rightarrow \text{Interval}$$

  that provides bounds for each (integer) variable

- This lattice has *infinite height*, since the chain:

  $$[0,0] \sqsubseteq [0,1] \sqsubseteq [0,2] \sqsubseteq [0,3] \sqsubseteq [0,4] \ldots$$

  occurs in *Interval*

# Interval Constraints

- For the *entry* node:

  $[[\textit{entry}]] = \lambda \mathrm{x}.\ [-\infty, \infty]$

- For assignments:

  $[[\mathrm{v}]] = \textit{JOIN}(\mathrm{v})[\textit{id} \rightarrow \textit{eval}(\textit{JOIN}(\mathrm{v}), E))$

- For all other nodes:

  $[[\mathrm{v}]] = \textit{JOIN}(\mathrm{v}) = \bigsqcup_{\mathrm{w} \in \textit{pred}(\mathrm{v})} [[\mathrm{w}]]$

# Evaluating Intervals

- The *eval* function is an *abstract evaluation*:
  - $eval(\sigma, id) = \sigma(id)$
  - $eval(\sigma, intconst) = [intconst, intconst]$
  - $eval(\sigma, E_1 \ \mathtt{op} \ E_2) = \overline{\mathtt{op}}(eval(\sigma, E_1), eval(\sigma, E_2))$

- Abstract arithmetic operators:
  - $\overline{\mathtt{op}}([l_1, h_1], [l_2, h_2]) =$
  $$\left[ \min_{x \in [l_1, h_1], \ y \in [l_2, h_2]} x \ \mathtt{op} \ y, \ \max_{x \in [l_1, h_1], \ y \in [l_2, h_2]} x \ \mathtt{op} \ y \right]$$

- Abstract comparison operators:
  - $\overline{\mathtt{op}}([l_1, h_1], [l_2, h_2]) = [0, 1]$

# Fixed-Point Problems

- The lattice has infinite height, so the fixed-point algorithm does not work

- In $L^n$ the sequence of approximants:

$$F^i(\bot, \bot, ..., \bot)$$

need never converge

# Widening

- Introduce a *widening* function $\omega: L^n \to L^n$ so that:

$$(\omega \circ F)^i(\bot, \bot, ..., \bot)$$

  converges on a fixed-point that is larger than all of the approximants $F^i(\bot, \bot, ..., \bot)$

- The function $\omega$ coarsens the information

# Turbo Charging

F                    ω

# Widening for Intervals

- The function $\omega$ is defined pointwise
- Parameterized with a fixed finite subset $B \subset N$
  - must contain $-\infty$ and $\infty$
  - typically seeded with all integer constants occurring in the given program
- On single intervals:

$$\omega(\,[l\,,h]\,) = [\,max\{i \in B | i \leq l\}\,,\ min\{i \in B | h \leq i\}\,]$$

- Finds the nearest enclosing allowed interval

# Correctness of Widening

- Widening works when:
  - $\omega$ is an *increasing* and *monotone* function
  - $\omega(L)$ is a *finite* lattice

- $F^i(\bot, \bot, ..., \bot) \sqsubseteq (\omega{\circ}F)^i(\bot, \bot, ..., \bot)$
  since F is monotone and $\omega$ is increasing

- $\omega{\circ}F$ is a monotone function $\omega(L){\rightarrow}\omega(L)$
  so the fixed-point exists

# Narrowing

- Widening shoots over the target
- *Narrowing* may improve the result by applying F
- Define:

$$fix = \bigsqcup F^i(\bot, \bot, ..., \bot) \qquad fix\omega = \bigsqcup (\omega \circ F)^i(\bot, \bot, ..., \bot)$$

  then $fix \sqsubseteq fix\omega$

- But we also have that:

$$fix \sqsubseteq F(fix\omega) \sqsubseteq fix\omega$$

  so applying F again may improve the result

- This can be iterated arbitrarily many times

# Correctness of Narrowing

- $F(fix\omega) \sqsubseteq \omega(F(fix\omega)) = (\omega \circ F)(fix\omega) = fix\omega$

  - by induction and monotonicity of F we also have:

    $F^{i+1}(fix\omega) \sqsubseteq F^i(fix\omega) \sqsubseteq fix\omega$

- $fix = \bigsqcup F^i(\bot, \bot, ..., \bot) = \bigsqcup F^{i+1}(\bot, \bot, ..., \bot)$

  $\sqsubseteq F(\bigsqcup F^i(\bot, \bot, ..., \bot)) = F(fix) \sqsubseteq F(fix\omega)$

  - by induction we also have:

    $fix \sqsubseteq F^i(fix\omega)$

# **Backing Up**



F        ω

22

# Divergence in Action

```
y = 0;

x = 7;

x = x+1;

while (input) {

    x = 7;

    x = x+1;

    y = y+1;

}
```
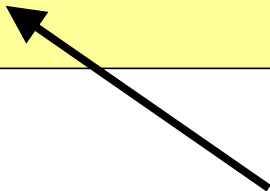
$[x \rightarrow \bot, y \rightarrow \bot]$
$[x \rightarrow [8,8], y \rightarrow [0,1]]$
$[x \rightarrow [8,8], y \rightarrow [0,2]]$
$[x \rightarrow [8,8], y \rightarrow [0,3]]$
…

# Widening in Action

```
y = 0;

x = 7;

x = x+1;

while (input) {

    x = 7;

    x = x+1;

    y = y+1;

}
```

$[x \rightarrow \bot, y \rightarrow \bot]$
$[x \rightarrow [7,\infty], y \rightarrow [0,1]]$
$[x \rightarrow [7,\infty], y \rightarrow [0,7]]$
$[x \rightarrow [7,\infty], y \rightarrow [0,\infty]]$

$B = \{-\infty, 0, 1, 7, \infty\}$

# Narrowing in Action

```
y = 0;

x = 7;

x = x+1;

while (input) {

   x = 7;

   x = x+1;

   y = y+1;

}
```

$[x \rightarrow \bot, y \rightarrow \bot]$
$[x \rightarrow [7,\infty], y \rightarrow [0,1]]$
$[x \rightarrow [7,\infty], y \rightarrow [0,7]]$
$[x \rightarrow [7,\infty], y \rightarrow [0,\infty]]$

$[x \rightarrow [8,8], y \rightarrow [0,\infty]]$

$B = \{-\infty, 0, 1, 7, \infty\}$

# Widening Functions

- A simple generic widening function:

$$\omega(x) = \begin{cases} x & \text{if } x \text{ is small enough} \\ \top & \text{otherwise} \end{cases}$$

- A difficult widening function (regular languages):

$\Sigma^*$

$\{a\} \subseteq \{a,ab\} \subseteq \{a,ab,abb\} \subseteq \dots \overset{\omega}{\rightarrow} \{ab^*\}$

$\varnothing$

This is essentially machine learning...

# Information in Conditions

```
x = input;

y = 0;

z = 0;

while (x>0) {

   z = z+x;

   if (17>y) { y = y+1; }

   x = x-1;

}
```

- The interval analysis (with widening) concludes:

$$x = [-\infty,\infty], \ y = [0,\infty], \ z = [-\infty,\infty]$$

# Modeling Conditions

- Add two artifical statements

- The statement `assert(`*E*`)` models that *E* is *true* in the current program state
- It causes a runtime error otherwise

- The statement `refute(`*E*`)` models that *E* is *false* in the current program state
- It causes a runtime error otherwise

# Encoding Conditions

```
x = input;

y = 0;

z = 0;

while (x>0) {

   assert(x>0);

   z = z+x;

   if (17>y) { assert(17>y); y = y+1; }

   x = x-1;

}

refute (x>0);
```

Preserves semantics since `assert` and `refute` are guarded by conditions

# Constraints for Assert and Refute

- A trivial but sound constraint is:

  $[[v]] = JOIN(v)$

- A non-trivial constraint for `assert(`*id*`>`*E*`)`:

  $[[v]] = JOIN(v)[id \rightarrow gt(JOIN(v)(id), eval(JOIN(v), E))]$

  where

  $gt( [l_1, h_1], [l_2, h_2] ) = [l_1, h_1] \sqcap [l_2, \infty]$

- Similar constraints are defined for the dual cases
- More tricky to define for all conditions...

# Exploiting Conditions

```
x = input;
y = 0;
z = 0;
while (x>0) {
   assert(x>0);
   z = z+x;
   if (17>y) { assert(17>y); y = y+1; }
   x = x-1;
}
refute (x>0);
```

- The interval analysis now concludes:
  $x = [-\infty, 0], \ y = [0, 17], \ z = [0, \infty]$

# Branch Correlations

- With assert and refute we have a simple form of *path sensitivity*

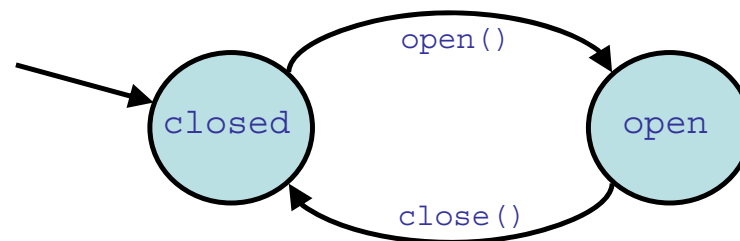- But it is insufficient to handle *correlation* of branches in program:

```
if (17 > x) { ... }
...
if (17 > x) { ... }
...
```

# Open and Closed Files

- Built-in functions `open()` and `close()` on a file

- Requirements:
  - never `close` a closed file
  - never `open` an open file



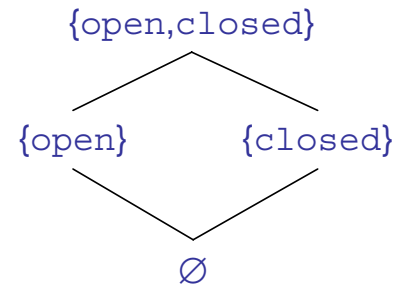- We want a static analysis to check this...

# A Tricky Example

```
if (condition) {
  open();
  flag = 1;
} else {
  flag = 0;
}
...
if (flag) {
  close();
}
```

# The Naive Analysis (1/2)

- The lattice models the status of the file:

$$L = (2^{\{\texttt{open,closed}\}}, \subseteq)$$

{open,closed}

{open}     {closed}

$\varnothing$

- For every CFG node, v, we have a constraint variable [[v]] denoting the status *after* v

- *JOIN*(v) = $\bigcup_{w \in pred(v)}$ [[w]]

# The Naive Analysis (2/2)

- **Constraints for interesting statements:**

  [[*entry*]] = {`closed`}

  [[`open()`]] = {`open`}

  [[`close()`]] = {`closed`}

- **For all other CFG nodes:**

  [[v]] = *JOIN*(v)

- **Before the `close()` statement the analysis concludes that the file is {`open,closed`}**

# Context Awareness

- We need to keep track of the `flag` variable
- Our second attempt is the lattice:

$$L = (2^{\{\texttt{open,closed}\}} \times 2^{\{\texttt{flag=0,flag}\neq\texttt{0}\}}, \subseteq \times \subseteq)$$

- Additionally, we add `assert(...)` and `refute(...)` to keep track of conditionals

- Even so, we now only now that the file is {`open,closed`} and that flag is {`flag=0,flag`≠`0`}

# Relational Analysis

- We need an analysis that keeps track of *relations* between variables

- This requires that we maintain *multiple* abstract states per program point, one for each *context*

- For the file example we need the lattice:

$$L = C \rightarrow 2^{\{\texttt{open,closed}\}}$$

where $C = \{\texttt{flag=0},\texttt{flag}\neq\texttt{0}\}$ is the set of contexts

# Enhanced Program

```
if (condition) {
  assert(condition);
  open();
  flag = 1;
} else {
  refute(condition);
  flag = 0;
}
...
if (flag) {
  assert(flag);
  close();
} else {
  refute(flag);
}
```

# Relational Constraints (1/2)

- For the file statements:

  $[[\textit{entry}]] = \lambda c.\{\texttt{closed}\}$

  $[[\texttt{open()}]] = \lambda c.\{\texttt{open}\}$

  $[[\texttt{closed()}]] = \lambda c.\{\texttt{closed}\}$

- For `flag` assignments:

  infeasible

  $[[\texttt{flag = 0}]] = [\texttt{flag=0} \rightarrow \bigcup_{c \in C} JOIN(v)(c), \texttt{flag}\neq 0 \rightarrow \varnothing]$

  $[[\texttt{flag = }n]] = [\texttt{flag}\neq 0 \rightarrow \bigcup_{c \in C} JOIN(v)(c), \texttt{flag=0} \rightarrow \varnothing]$

  $[[\texttt{flag = }E]] = \lambda d.\bigcup_{c \in C} JOIN(v)(c)$

# Relational Constraints (2/2)

- For assert and refute statements:

$$[[\texttt{assert(flag)}]] =$$
$$[\texttt{flag} \neq 0 \rightarrow JOIN(v)(\texttt{flag} \neq 0), \texttt{flag} = 0 \rightarrow \varnothing]$$

$$[[\texttt{refute(flag)}]] =$$
$$[\texttt{flag} = 0 \rightarrow JOIN(v)(\texttt{flag} = 0), \texttt{flag} \neq 0 \rightarrow \varnothing]$$

- For all other CFG nodes:

$$[[v]] = JOIN(v) = \lambda c. \bigcup_{w \in pred(v)} [[w]](c)$$

# Generated Constraints

[[*entry*]] = λc.{closed}

[[condition]] = [[*entry*]]

[[assert(condition)]] = [[condition]]

[[open()]] = λc.{open}

[[flag = 1]] = [flag≠0→∪$_{c∈C}$[[open()]](c),flag=0→∅]

[[refute(condition)]] = condition

[[flag = 0]] = [flag=0→∪$_{c∈C}$[[refute(condition)]](c),flag≠0→∅]

[[...]] = λc.([[flag = 1]](c) ∪ [[flag = 0]](c))

[[flag]] = [[...]]

[[assert(flag)]] = [[flag≠0→[[flag]](flag≠0), flag=0→∅]

[[close()]] = λc.{closed}

[[refute(flag)]] = [flag=0→[[flag]](flag=0), flag≠0→∅]

[[*exit*]] = λc.([[close()]](c) ∪ [[...]](c))

# Minimal Solution

|  | `flag = 0` | `flag ≠ 0` |
|---|---|---|
| [[*entry*]] | {closed} | {closed} |
| [[`condition`]] | {closed} | {closed} |
| [[`assert(condition)`]] | {closed} | {closed} |
| [[`open()`]] | {open} | {open} |
| [[`flag = 1`]] | ∅ | {open} |
| [[`refute(condition)`]] | {closed} | {closed} |
| [[`flag = 0`]] | {closed} | ∅ |
| [[`...`]] | {closed} | {open} |
| [[`flag`]] | {closed} | {open} |
| [[`assert(flag)`]] | ∅ | {open} |
| [[`close()`]] | {closed} | {closed} |
| [[`refute(flag)`]] | {closed} | ∅ |
| [[*exit*]] | {closed} | {open} |

- We know the file is `open` before `close()` 🙂

# **Challenges**

- The static analysis designer must choose C
  - often as combinations of predicates from conditionals
  - *iterative refinement* gradually adds predicates

- Exponential blow-up:
  - for $k$ predicates, we have $2^k$ different contexts
  - redundancy often cuts this down

- Reasoning about `assert` and `refute`:
  - how to update the lattice elements sufficiently precisely
  - possibly involves theorem proving

# Improvements

- Run auxiliary analyses first, for example:
  - constant propagation
  - sign analysis

  will help in handling `flag` assignments

- Dead code propagation, change:

  $$[[\texttt{open()}]] = \lambda c.\{\texttt{open}\}$$

  into the still sound but more precise:

  $$[[\texttt{open()}]] = \lambda c.\text{if } JOIN(v)(c) = \varnothing \text{ then } \varnothing \text{ else } \{\texttt{open}\}$$

# Interprocedural Analysis

- Analyzing the body of a single function:
  - *intra*procedural analysis

- Analyzing the whole program with function calls:
  - *inter*procedural analysis

- The alternative is to:
  - analyze each function in isolation
  - be maximally pessimistic about results of function calls

# CFG for Whole Programs

- Construct a CFG for each function

- Then glue them together to reflect function calls

- Assume that all function calls are of the form:

$$id \; = \; \text{f}(E_1, \; ..., \; E_n);$$

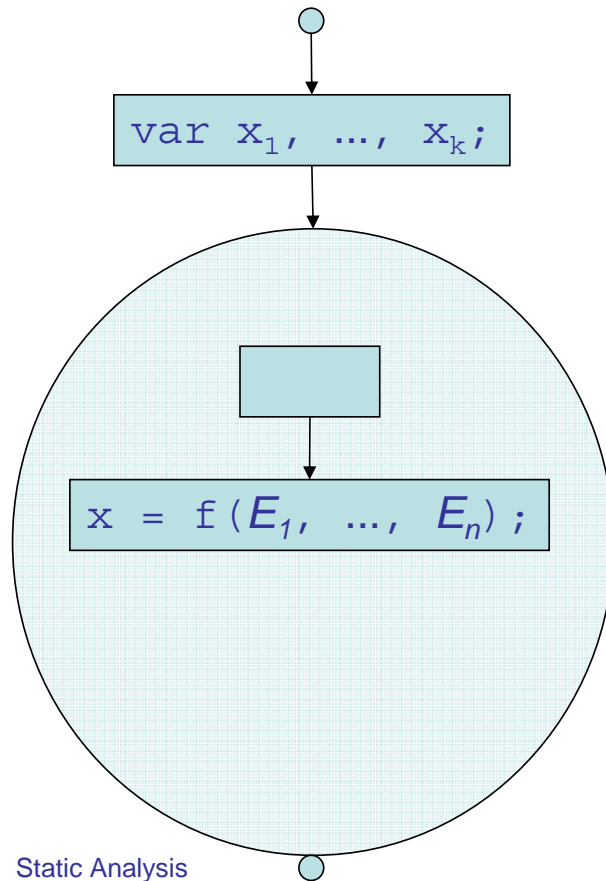- This can always be obtained by rewriting

# Shadow Variables

- Introduce some extra variables in the program

- For every function `f` the variable `ret-f` denoting its return value

- For every call site with index `i` a variable `call-i` denoting the computed value

- For every local or formal `x` and call site with index `i` a register `save-i-x`

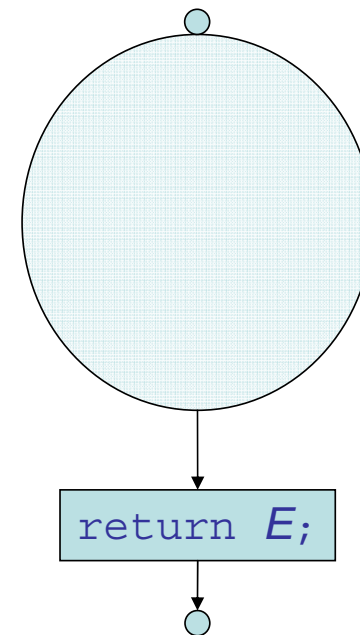- For every formal `x` and every call site with index `i` a temporary variable `temp-i-x`

# Calling and Called Function

function g($a_1$, ..., $a_n$)                    function f($b_1$, ..., $b_m$)

```
var x₁, …, xₖ;
```
$$\text{var } x_1, \ldots, x_k;$$

$$x = f(E_1, \ldots, E_n);$$

$$\text{return } E;$$

# Glued Together

function g(a$_1$, …, a$_n$)            function f(b$_1$, …, b$_m$)

```
var x₁, …, xₖ;
```

```
save-i-b_j  =  b_j
save-i-x_j  =  x_j
temp-i-a_j  =  E_j
a_j  =  temp-i-a_j
```

```
b_j  =  save-i-b_j
x_j  =  save-i-x_j
x = call-i
```

```
call-i = ret-f
```
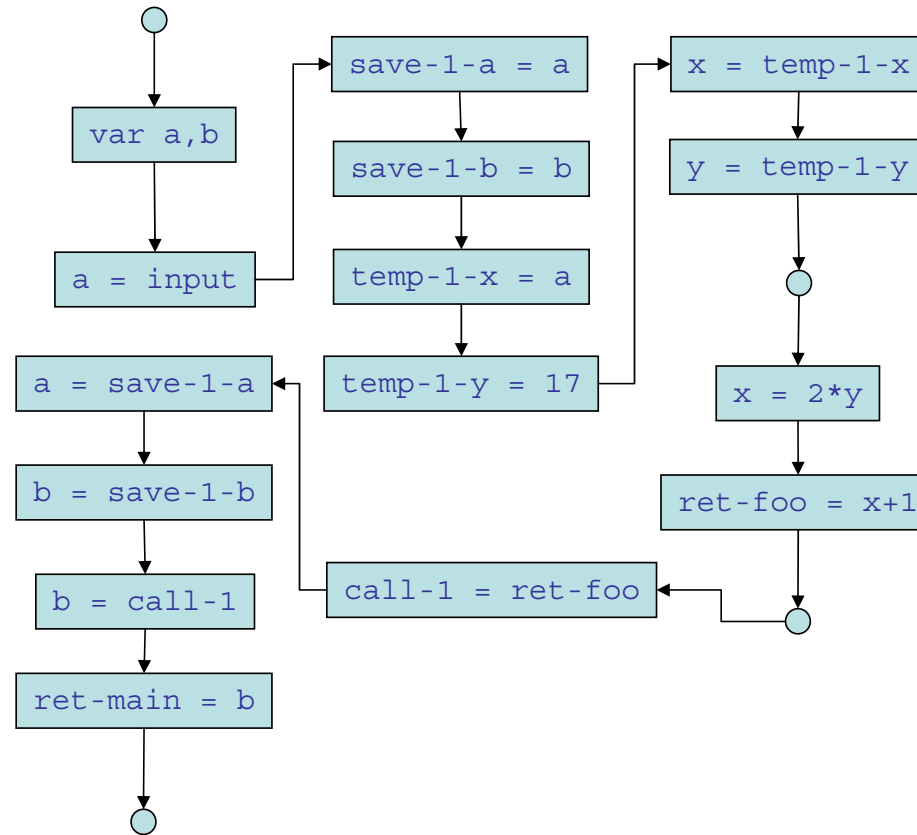
```
ret-f = E;
```

# Example Program

```
foo(x,y) {
  x = 2*y;
  return x+1;
}

main() {
  var a,b;
  a = input;
  b = foo(a,17);
  return b;
}
```

# Resulting CFG
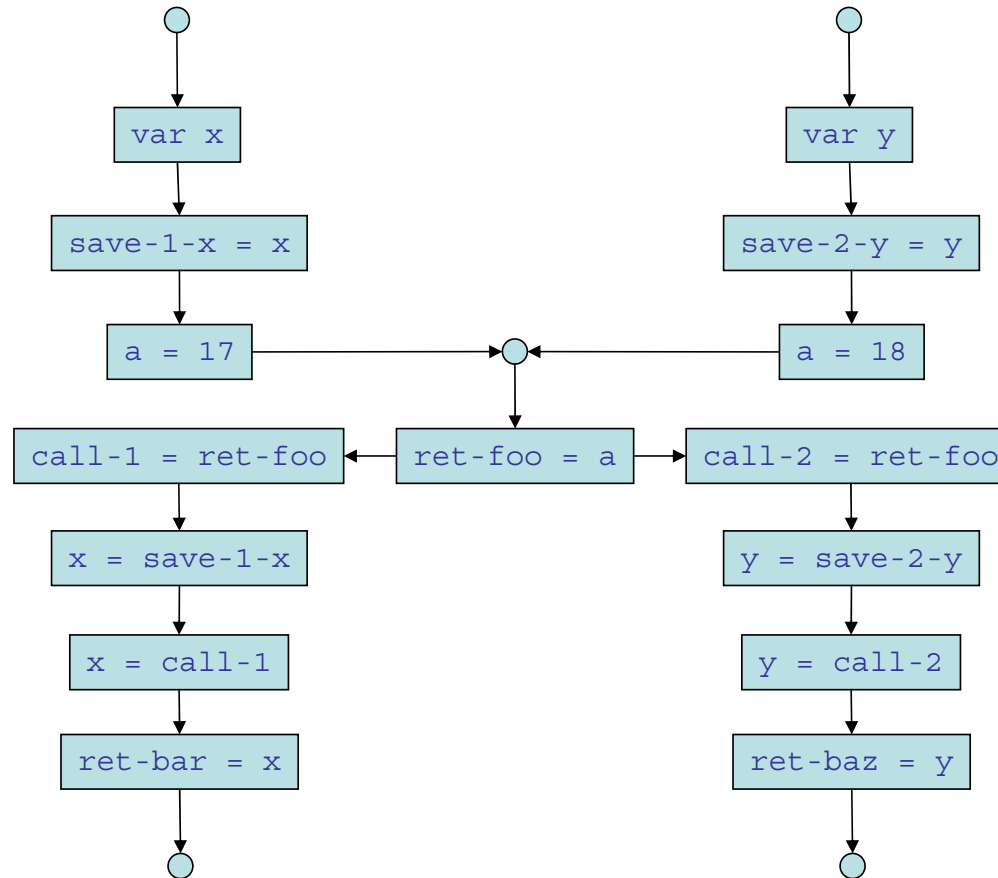
```
foo(x,y) {
   x = 2*y;
   return x+1;
}

main() {
   var a,b;
   a = input;
   b = foo(a,17);
   return b;
}
```
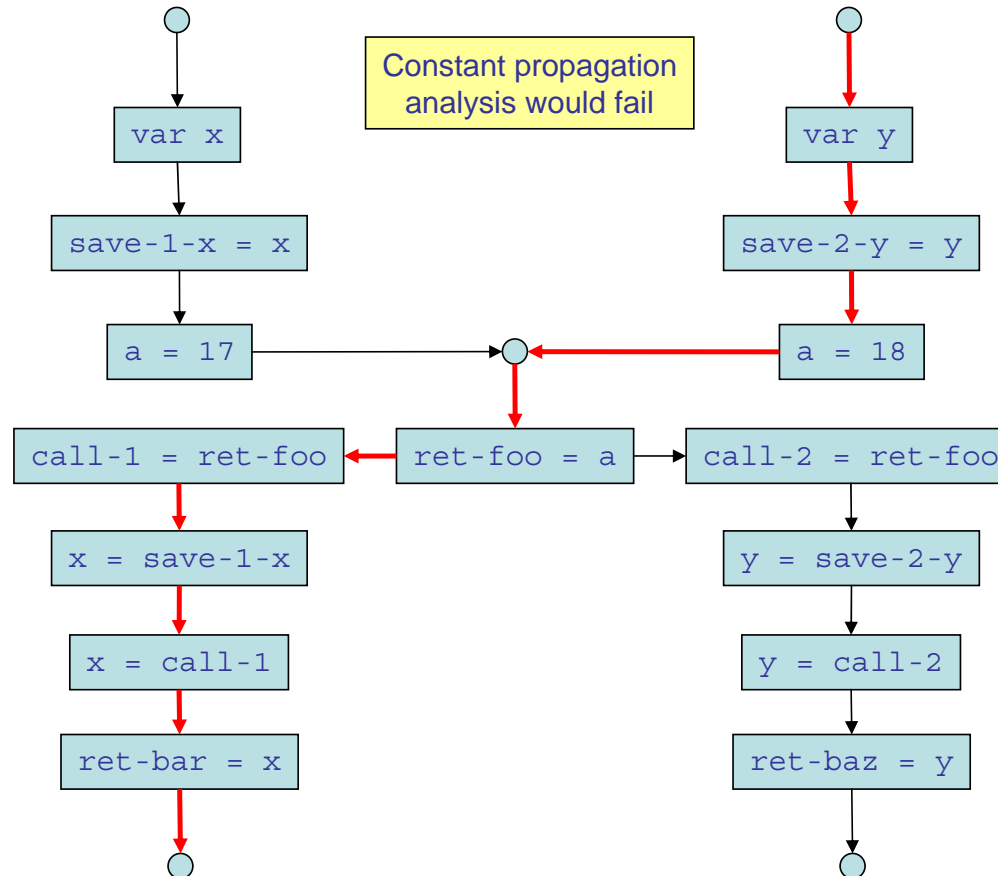
# False Control Flow

```
foo(a) {
   return a;
}

bar() {
   var x;
   x = foo(17);
   return x;
}

baz() {
   var y;
   y = foo(18);
   return y;
}
```

# False Control Flow

```
foo(a) {
   return a;
}

bar() {
   var x;
   x = foo(17);
   return x;
}

baz() {
   var y;
   y = foo(18);
   return y;
}
```

Constant propagation analysis would fail

var x

save-1-x = x

a = 17

call-1 = ret-foo

x = save-1-x

x = call-1

ret-bar = x

var y

save-2-y = y

a = 18

ret-foo = a
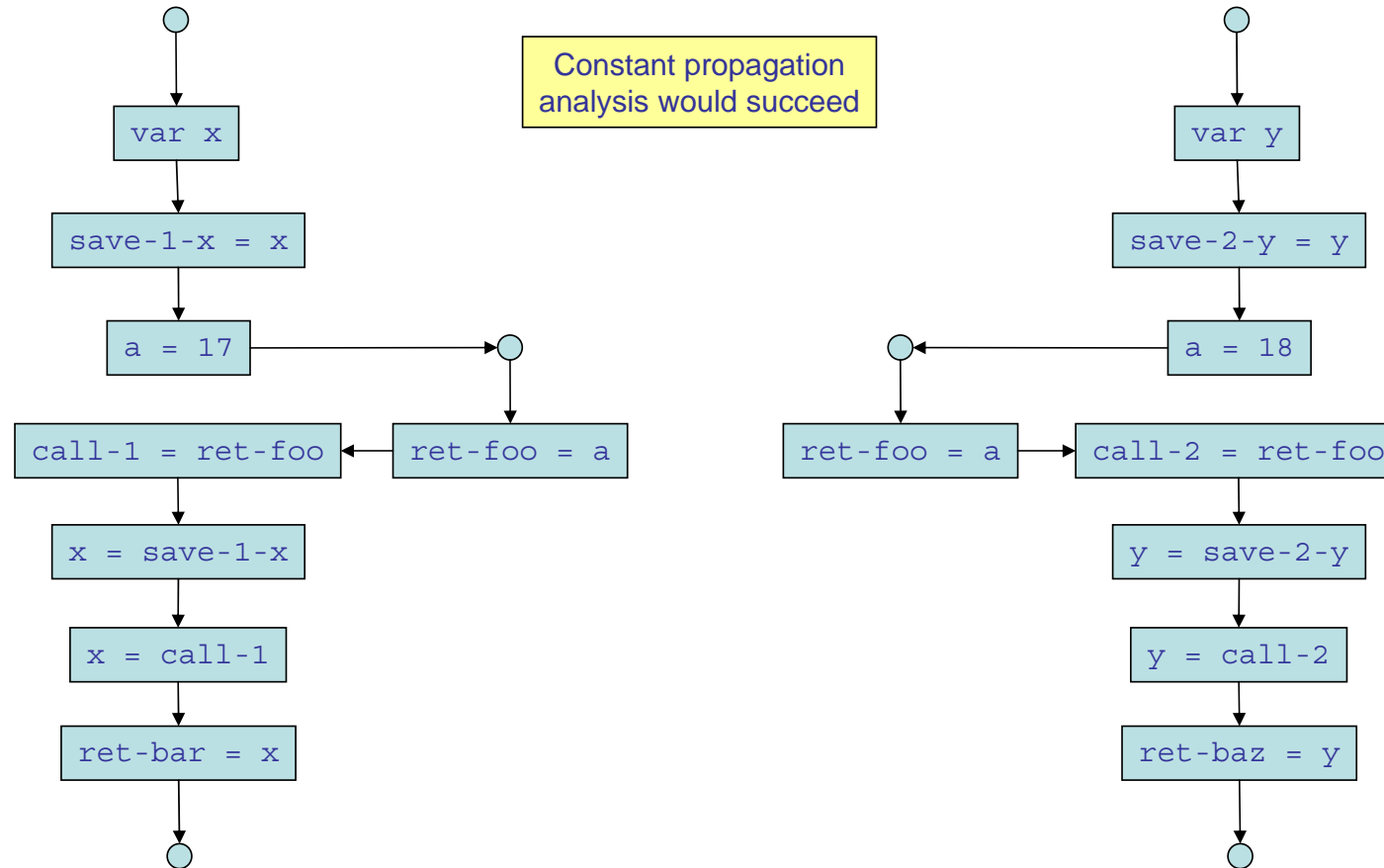
call-2 = ret-foo

y = save-2-y

y = call-2

ret-baz = y

# Polyvariance vs. Monovariance

- A *polyvariant* analysis creates *multiple copies* of the CFG for the body of a called function

- A *monovariant* analysis uses only one copy

- Strategies determine the number of copies:
  - the simplest is one copy for each call site
  - dynamic heuristics are also possible
  - important that only finitely many copies are created

# Polyvariant CFG

var x

save-1-x = x

a = 17

call-1 = ret-foo ← ret-foo = a

x = save-1-x

x = call-1

ret-bar = x

Constant propagation
analysis would succeed

var y

save-2-y = y

a = 18

ret-foo = a → call-2 = ret-foo

y = save-2-y

y = call-2

ret-baz = y

# Tree Shaking

- Identify those functions that are never called
  - safely remove them from the program
  - reduces size of the compiled executable
  - reduces size of CFG for subsequent analyses

- Uses monovariant interprocedural CFG

- Essentially a transitive closure computation

# Setting Up

- The lattice is the powerset of all function names

- For every CFG node v we introduce a constraint variable [[v]] denoting the set of function that could *possibly* be called in the *future*

- We let *entry*(*id*) denote the entry node in the CFG for the function named *id*

# Tree Shaking Constraints

- For assignments, conditions and `output`:

$$[[v]] = \bigcup_{w \in \textit{succ}(v)} [[w]] \cup \textit{funcs}(E) \cup \bigcup_{f \in \textit{funcs}(E)} [[\textit{entry}(f)]]$$

- For all other nodes:

$$[[v]] = \bigcup_{w \in \textit{succ}(v)} [[w]]$$

- Here *funcs* is defined as:
  - $\textit{funcs}(id) = \textit{funcs}(intconst) = \textit{funcs}(\texttt{input}) = \varnothing$
  - $\textit{funcs}(E_1 \; \texttt{op} \; E_2) = \textit{funcs}(E_1) \cup \textit{funcs}(E_2)$
  - $\textit{funcs}(id(E_1, \ldots, E_n)) = \{id\} \cup \textit{funcs}(E_i)$