

Principles and applications of abstract-interpretation-based static analysis

David Schmidt

Kansas State University

`www.cis.ksu.edu/~schmidt`

Outline

Static analysis is property extraction from formal systems.

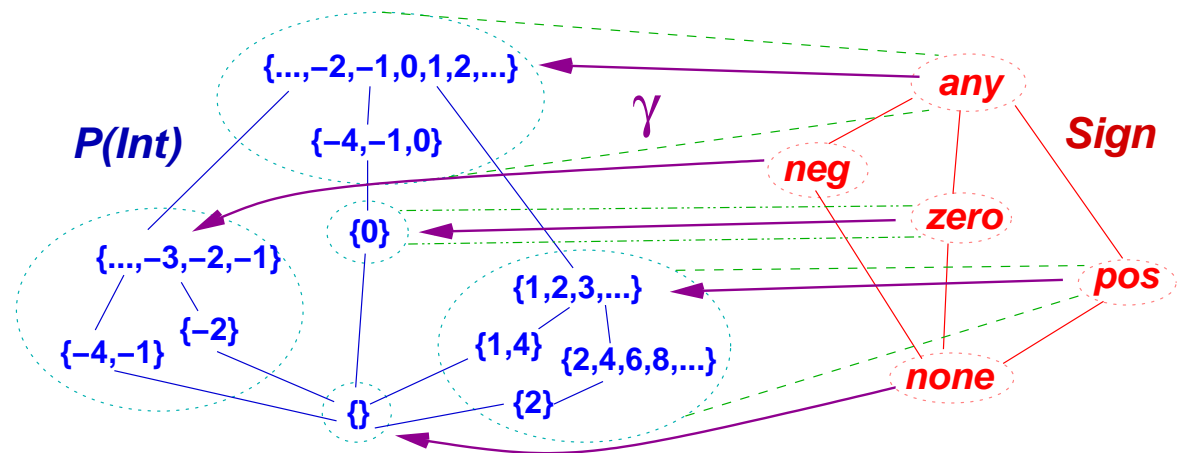
Abstract interpretation is a foundation for static analysis based on Galois connections, semi-homomorphisms, and fixed-point calculation. In this talk, we

- ◆ introduce abstract interpretation
- ◆ apply it to static analyses of program semantics
(state-transition systems, equationally specified definitions, rule-based relational definitions)
- ◆ survey applications of static analysis
- ◆ develop the correspondence of properties to propositions
- ◆ consider approaches to modular, “scalable” analyses

Background: abstract interpretation

An abstract domain defines properties

A formal system uses values from set C , and we wish to determine properties of the C -values that might arise during computation.



Define an *abstract domain*, A : a partially ordered set of properties, closed under meet (\sqcap). See example, *Sign*, above.

Define a monotone *concretization map*, $\gamma : A \rightarrow PC$, where PC is the powerset of C , ordered by \subseteq , so that $\gamma(a)$ defines those elements that “have property a .”

γ must *preserve meets* – for $T \subseteq A$, $\gamma(\sqcap T) = \bigcap_{a \in T} \gamma(a)$ – so that an inverse function, $\alpha : PC \rightarrow A$, can be defined.

Operations f are abstracted to $f^\#$ to compute on A

```
readInt(x)
```

```
x = succ(x)
```

```
if x < 0 :
```

```
  x = negate(x)
```

```
else:
```

```
  x = succ(x)
```

```
writeInt(x)
```

Q: is the output *pos*?

A: abstractly interpret
input domain *Int* by

Sign =
{*neg*, *zero*, *pos*, *any*, *none*}:

```
readSign(x)
```

```
x = succ#(x)
```

```
if (filterNeg(x):
```

```
  x = negate#(x))
```

```
(filterNonNeg(x):
```

```
  x = succ#(x)) fi
```

```
writeSign(x)
```

$\text{succ}^\#(\text{pos}) = \text{pos}$

$\text{succ}^\#(\text{zero}) = \text{pos}$

where $\text{succ}^\#(\text{neg}) = \text{any}$ (!) and

$\text{succ}^\#(\text{any}) = \text{any}$

$\text{negate}^\#(\text{neg}) = \text{pos}$

$\text{negate}^\#(\text{zero}) = \text{zero}$

$\text{negate}^\#(\text{pos}) = \text{neg}$

$\text{negate}^\#(\text{any}) = \text{any}$

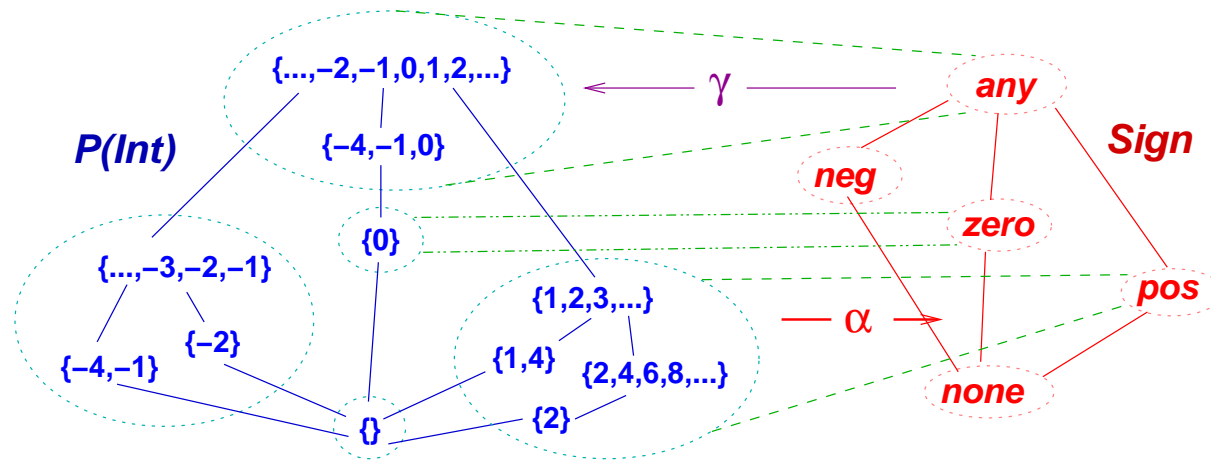
For the abstract data-test sets, *zero*, *neg*, *pos*, we calculate:

$\{\text{zero} \mapsto \text{pos}, \text{pos} \mapsto \text{pos}, \text{neg} \mapsto \text{any}\}$. The last result arises because

$\text{succ}^\#(\text{neg}) = \text{any}$ and $\text{filterNeg}(\text{any}) = \text{neg}$ (good!) but $\text{filterNonNeg}(\text{any}) = \text{any}$

(bad — we need $\text{zero} \vee \text{pos}$!), so we cannot ensure the success of the else-arm.

A Galois connection formalizes the abstraction



$$\gamma : \text{Sign} \rightarrow \mathcal{P}(\text{Int})$$

$$\gamma(\text{none}) = \{\}, \quad \gamma(\text{any}) = \text{Int}$$

$$\gamma(\text{neg}) = \{\dots, -3, -2, -1\}$$

$$\gamma(\text{zero}) = \{0\}, \quad \gamma(\text{pos}) = \{1, 2, 3, \dots\}$$

$$\alpha : \mathcal{P}(\text{Int}) \rightarrow \text{Sign}$$

$$\alpha(S) = \sqcap \{a \mid \gamma(a) \subseteq S\}$$

$$\text{e.g., } \alpha\{2, 4, 6, 8, \dots\} = \text{pos},$$

$$\alpha\{-1, 0\} = \text{any}, \quad \alpha\{0\} = \text{zero}$$

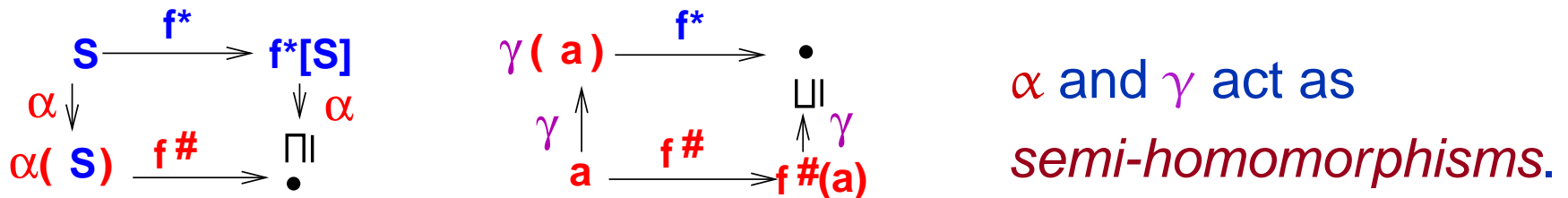
$(\mathcal{P}(\text{Int}), \subseteq) \langle \alpha, \gamma \rangle (\text{Sign}, \sqsubseteq)$ is a **Galois connection**:

$$\alpha(S) \sqsubseteq a \text{ iff } S \subseteq \gamma(a).$$

γ interprets the elements in Sign , and α maps each data-test set in $\mathcal{P}(\text{Int})$ to the property that best describes the set [CousotCousot77].

An abstract operation is monotone and sound

$f^\# : A \rightarrow A$ is *sound* for $f : C \rightarrow PC$ iff $\alpha \circ f^* \sqsubseteq f^\# \circ \alpha$
 (iff $f^* \circ \gamma \sqsubseteq \gamma \circ f^\#$):



Example: The $\text{succ}^\#$ function seen earlier is sound for succ , e.g., for $\text{succ} : \text{Int} \rightarrow \mathcal{P}(\text{Int})$, $\text{succ}^*(0) = \{1\}$, and $\text{succ}^\#(\text{zero}) = \text{pos}$.

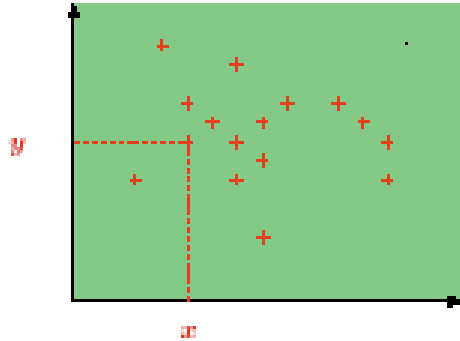
$f^\#$ is a *postcondition transformer*: $S \subseteq \gamma(a)$ implies $f^*(S) \subseteq \gamma(f^\#(a))$ where $f^*(S) = \bigcup_{c \in S} f(c)$.

$f_{\text{best}}^\# = \alpha \circ f^* \circ \gamma$ is the *strongest (liberal) postcondition transformer*.

Definition: $f^\#$ is γ -complete ("forwards complete") for f iff $f^* \circ \gamma = \gamma \circ f^\#$ [Giacobazzi01]. $f^\#$ is α -complete ("backwards complete") for f iff $\alpha \circ f^* = f^\# \circ \alpha$ [Cousots00].

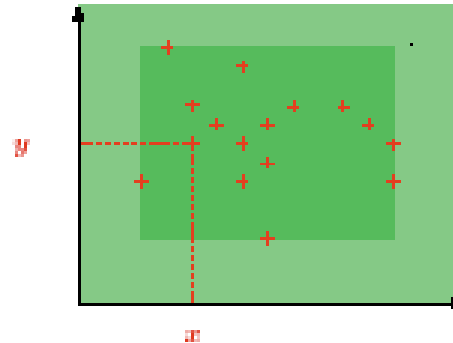
An aggregate, e.g., $Var \rightarrow C$, can be abstracted **pointwise** or **relationally**

Sign: $[x \mapsto \geq 0][y \mapsto \geq 0]$



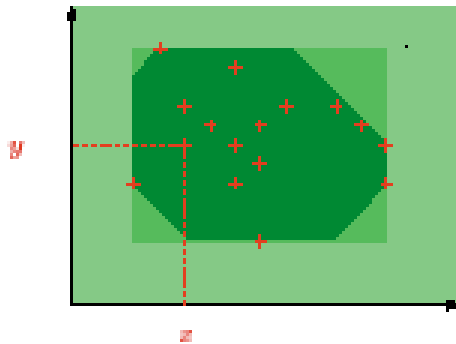
$$\begin{cases} x \geq 0 \\ y \geq 0 \end{cases}$$

Interval: $[x \mapsto [3, 27]][y \mapsto [4, 32]]$



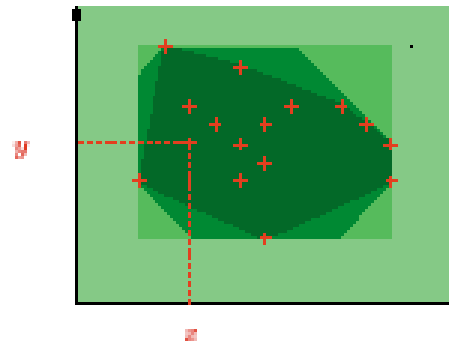
$$\begin{cases} x \in [3, 27] \\ y \in [4, 32] \end{cases}$$

Octagon: $\bigwedge_i (\pm x_i \pm y_i \leq c_i)$



$$\begin{cases} 3 \leq x \leq 27 \\ x + y \leq 88 \\ 4 \leq y \leq 32 \\ x - y \leq 61 \end{cases}$$

Polyhedron: $\bigwedge_i ((\sum_j a_{ij} \cdot x_{ij}) \leq b_i)$



$$\begin{cases} 7x + 31y \leq 325 \\ 21x + 7y \geq 0 \end{cases}$$

diagrams from *Abstract Interpretation: Achievements and Perspectives* by Patrick Cousot, Proc. SSGRR 2000.

Three codings (a)-(c) of a relationally abstracted store based on the octagon abstract domain:

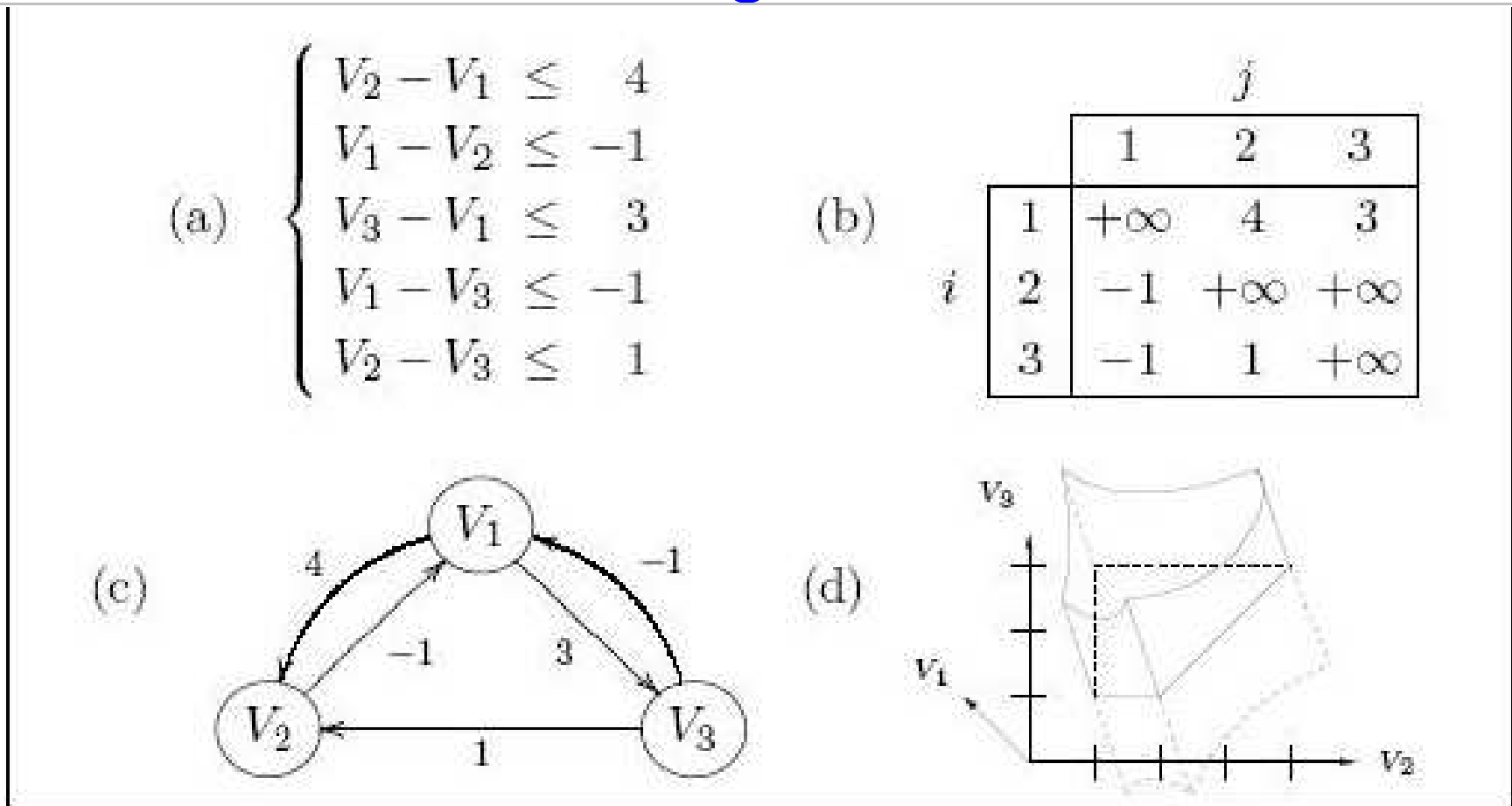


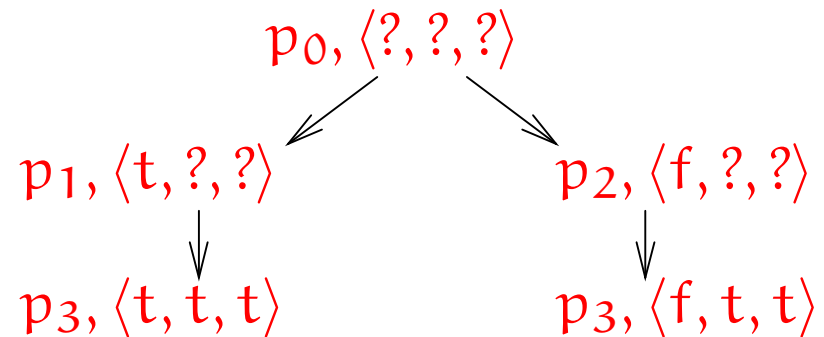
Figure 2. A potential constraint conjunction (a), its corresponding DBM m (b), potential graph $\mathcal{G}(m)$ (c), and potential set concretization $\gamma^{Pot}(m)$ (d).

diagram from *The octagon abstract domain*, by Antoine Miné, *J. Symbolic and Higher-Order Computation* 2006

Predicate abstraction uses a relational domain based on the predicates in the goal and program

Example: prove that $z \geq x \wedge z \geq y$ at p_3 :

p_0 : **if** $x < y$
 p_1 : **then** $z = y$
 p_2 : **else** $z = x$
 p_3 : **exit**



The store is abstracted to a relational domain that denotes the values of these predicates:

$$\phi_1 = x < y \quad \phi_2 = z \geq x \quad \phi_3 = z \geq y$$

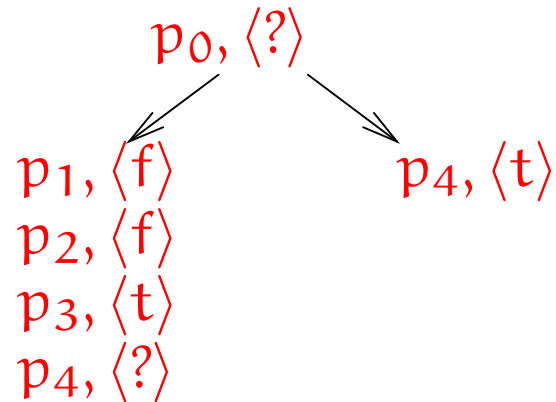
The predicates are evaluated at the program's points as one of $\{t, f, ?\}$.
(Read $?$ as $t \vee f$.)

At all occurrences of p_3 in the abstract trace, $\phi_2 \wedge \phi_3$ holds.

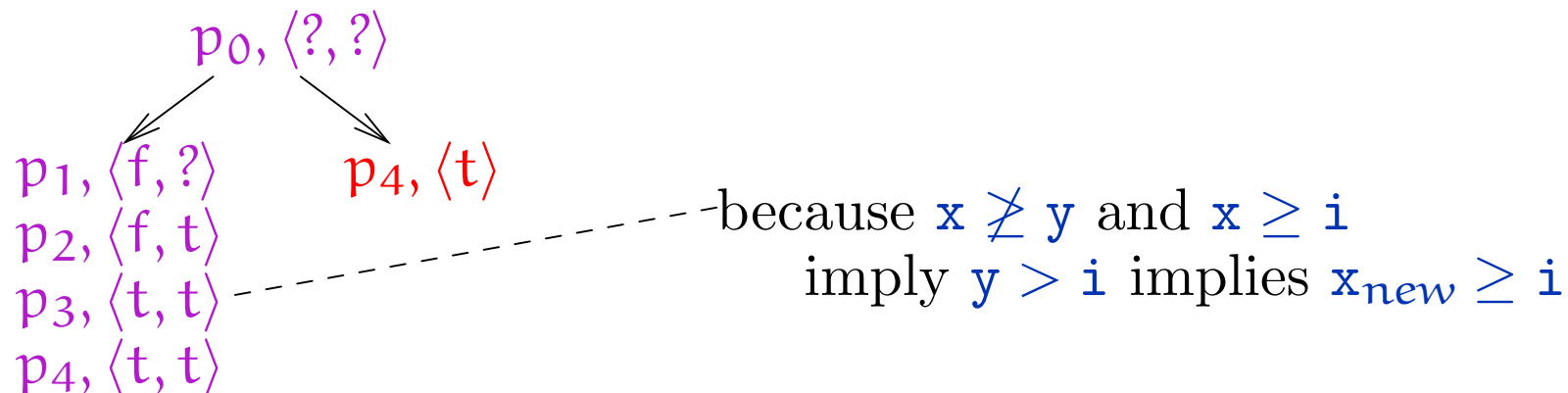
When a goal is undecided, domain refinement becomes necessary

Prove $\phi_0 \equiv x \geq y$ at p_4 :

```
 $p_0$  : if !( $x \geq y$ )  
 $p_1$  : then {  $i = x$ ;  
           $p_2$  :  $x = y$ ;  
           $p_3$  :  $y = i$ ;  
 $p_4$  : }
```



To decide the goal, we refine the abstract domain by adding a new predicate: $wp(y = i, x \geq y) = (x \geq i) \equiv \phi_1$. We add ϕ_1 and try again:



But incremental predicate refinement cannot synthesize many interesting loop invariants. For this example:

```
 $p_0$  :  $i = n; x = 0;$   
 $p_1$  : while  $i \neq 0$  {  
     $p_2$  :  $x = x + 1; i = i - 1;$   
}  
 $p_3$  : goal:  $x = n$ 
```

The initial predicate set, $P_0 \equiv \{i = 0, x = n\}$, does not validate the loop body.

The first refinement suggests we add $P_1 \equiv \{i = 1, x = n - 1\}$ to the program state, but this fails to validate a loop that iterates more than once.

Refinement stage j adds predicates $P_j \equiv \{i = j, x = n - j\}$; the refinement process continues forever!

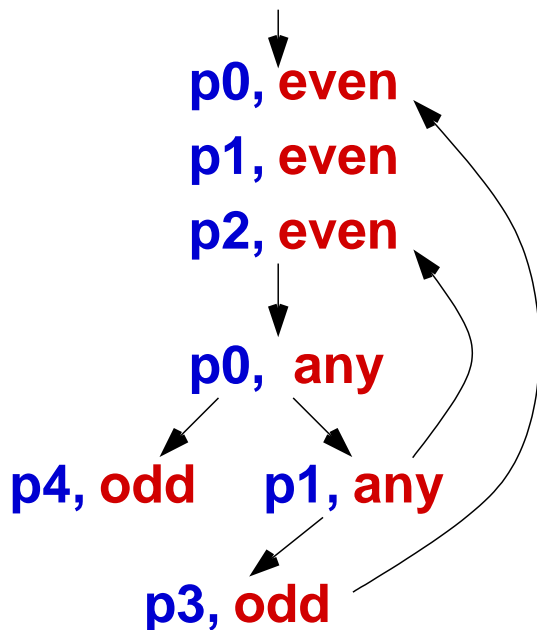
The loop invariant is $x = n - i$:-)

***Mechanics of static analysis:
abstracting small-step and
big-step semantics definitions***

The most basic static analysis is trace generation

```
 $p_0$  : while (x != 1) {  
   $p_1$  : if Even(x)  
     $p_2$  : then x = x div 2;  
     $p_3$  : else x = 3*x + 1;  
  }  
 $p_4$  : exit
```

Note: p_i, v abbreviates $p_i, \langle x : v \rangle$

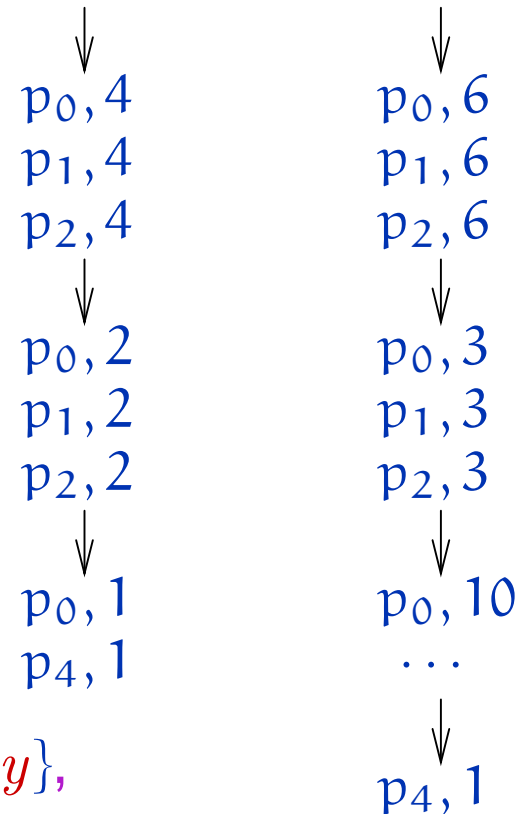


For

$Parity = \{none, even, odd, any\}$,

the loop's operations, f , are abstracted to $f^\#$. The *abstract trace* is a static analysis of those concrete executions with an even-valued input. Traces are used in *model checking*.

Two concrete traces:



Data-flow analysis *collects the abstract trace into a map*, $\text{ProgramPoint} \rightarrow \mathbb{A}$

The abstract value “attached” to program point p_i is defined by the first-order equational pattern,

$$p_i\text{Store} = \bigsqcup_{p_j \in \text{pred}(p_i)} f_j^\#(p_j\text{Store})$$

Flow equations for previous example:

$$\text{init} = \langle x: \textit{even} \rangle$$

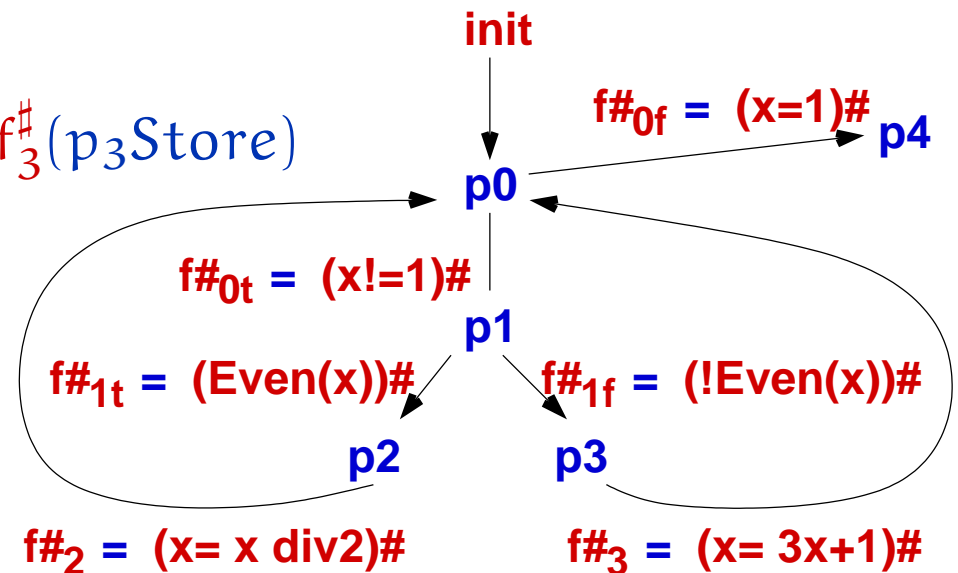
$$p_0\text{Store} = \text{init} \sqcup f_2^\#(p_2\text{Store}) \sqcup f_3^\#(p_3\text{Store})$$

$$p_1\text{Store} = f_{0t}^\#(p_0\text{Store})$$

$$p_2\text{Store} = f_{1t}^\#(p_1\text{Store})$$

$$p_3\text{Store} = f_{1f}^\#(p_1\text{Store})$$

$$p_4\text{Store} = f_{0f}^\#(p_0\text{Store})$$



We *solve* the flow equations by calculating approximate solutions in stages until *the least fixed point* is reached.

Note: \perp is the same as $\langle x:\perp \rangle$

stage	p_0 Store	p_1 Store	p_2 Store	p_3 Store	p_4 Store
0	\perp	\perp	\perp	\perp	\perp
1	$\langle x:even \rangle$	\perp	\perp	\perp	\perp
2	$\langle x:even \rangle$	$\langle x:even \rangle$	\perp	\perp	\perp
3	$\langle x:even \rangle$	$\langle x:even \rangle$	$\langle x:even \rangle$	\perp	\perp
4	$\langle x:any \rangle$	$\langle x:even \rangle$	$\langle x:even \rangle$	\perp	\perp
5	$\langle x:any \rangle$	$\langle x:any \rangle$	$\langle x:even \rangle$	\perp	$\langle x:odd \rangle$
6	$\langle x:any \rangle$	$\langle x:any \rangle$	$\langle x:even \rangle$	$\langle x:odd \rangle$	$\langle x:odd \rangle$

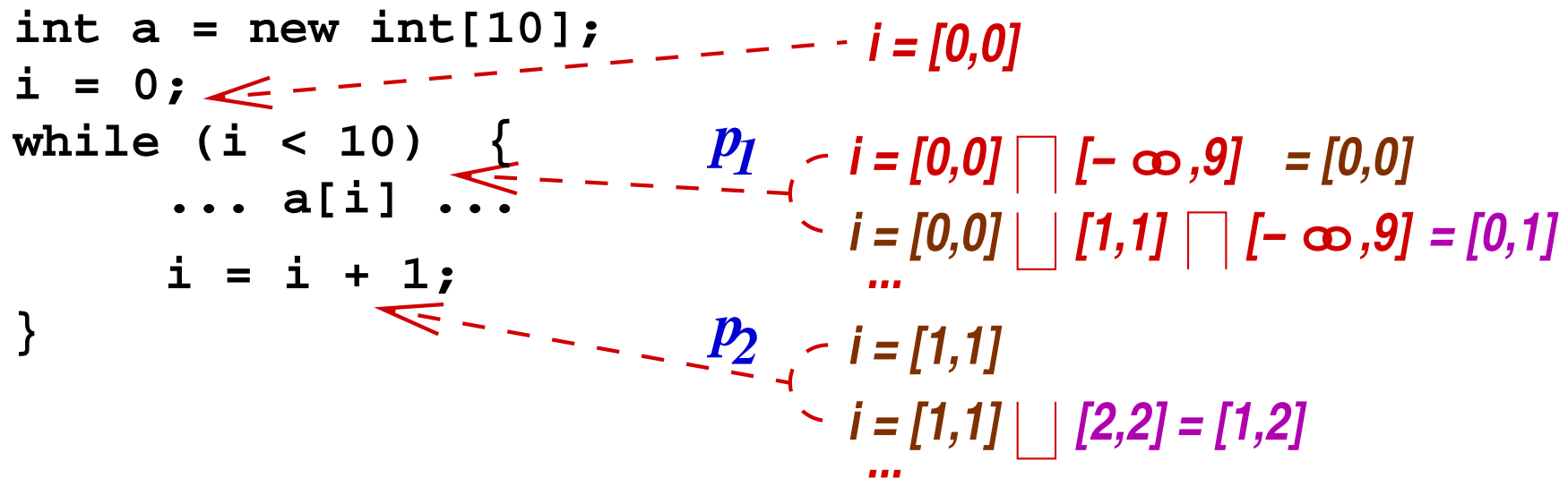
A faster algorithm uses a *worklist* that remembers exactly which equations should be recalculated at each stage.

Termination: *Array-bounds checking reviewed*

Integer variables might receive values from the *interval domain*,

$$I = \{[i, j] \mid i, j \in \text{Int} \cup \{-\infty, +\infty\}\}.$$

We define $[a, b] \sqcup [a', b'] = [\min(a, a'), \max(b, b')]$.



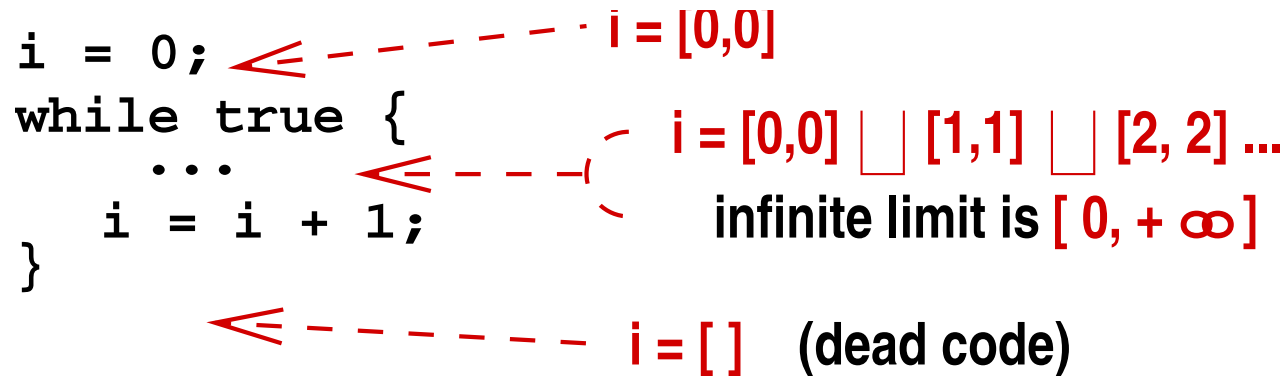
This example terminates: i 's ranges are

at p_1 : $[0..9]$

at p_2 : $[1..10]$

at loop exit : $[1..10] \sqcap [10, +\infty] = [10, 10]$

But others might not, because the domain is not finite height:



The analysis generates the infinite sequence of stages, $[0, 0], [0, 1], \dots, [0, i], \dots$ as i 's value in the loop's body.

The domain of intervals, where $[i, j] \sqsubseteq [i', j']$ iff $i \leq j$ and $j \leq j'$, has infinitely ascending chains.

To forcefully terminate the analysis, we can replace the \sqcup operation by ∇ , called a *widening operator*:

$$[] \nabla [i, j] = [i, j] \quad [i, j] \nabla [i', j'] = \begin{cases} \text{if } i' < i \text{ then } -\infty \text{ else } i, \\ \text{if } j' > j \text{ then } +\infty \text{ else } j \end{cases}$$

The widening operator, which guarantees finite convergence for all increasing sequences on the interval domain, quickly terminates the example:

```

i = 0;  $\Leftarrow$  ----- i = [0,0]
while true {
  ...  $\Leftarrow$  ----- i = [0,0]  $\nabla$  [1,1] = [0, + $\infty$ ]
  i = i + 1;
}
 $\Leftarrow$  ----- i = [] (dead code)

```

but in general, it can lose much precision:

```

int a = new int[10];
i = 0;  $\Leftarrow$  ----- i = [0,0]
while (i < 10) {
  ... a[i]  $\Leftarrow$  ----- i = [0,0]  $\nabla$  [1,1] = [0, + $\infty$ ]
  i = i + 1;
}
 $\Leftarrow$  ----- i = [10, + $\infty$ ]

```

For this reason, a complementary operation, \triangle , called a *narrowing operation*, can be used after ∇ gives convergence to recover some precision and retain a fixed-point solution.

We will not develop \triangle here, but for the interval domain, a suitable \triangle tries to reduce $-\infty$ and $+\infty$ to finite values. For the last example, the convergent value, $[0, +\infty]$, in the loop body would be narrowed to $[0, 10]$, making i 's value on loop exit $[10, 10]$.

Another approach is to use multiple “thresholds” for widening, e.g. $-\infty$, $(2^{-31} - 1)$, 0 , etc. for lower limits, and $(2^{31} - 1)$ and $+\infty$ for upper limits.

Structured static analysis on syntax trees

Given a block of statements, B , we might wish to calculate the values that “enter” and “exit” from B . If B is coded in a structured language, the static analysis can compute a “structured transfer function” for B :

$$C ::= p : x = E \mid C \mid \text{if } E \ C_1 \ C_2 \mid \text{while } E \ C$$

A sample structured analysis that ignores tests: $\llbracket C \rrbracket : A_{\text{in}} \rightarrow A_{\text{out}}$

$$\llbracket p : x = E \rrbracket \text{in} = f_p^\#(\text{in}) \quad (\text{the transfer function for } p)$$

$$\llbracket C \rrbracket \text{in} = \llbracket C_2 \rrbracket(\llbracket C_1 \rrbracket \text{in})$$

$$\llbracket \text{if } E \ C_1 \ C_2 \rrbracket \text{in} = \llbracket C_1 \rrbracket \text{in} \sqcup \llbracket C_2 \rrbracket \text{in}$$

$$\llbracket \text{while } E \ C \rrbracket \text{in} = \text{in} \sqcup \text{out}_C,$$

$$\text{where } \text{out}_C = \bigsqcup_{i \geq 0} \text{out}_i,$$

$$\text{and } \text{out}_0 = \perp_A \text{ and } \text{out}_{i+1} = \llbracket C \rrbracket(\text{in} \sqcup \text{out}_i)$$

We annotate a syntax tree with the *in*- and *out*-data — here is a **reaching definitions** data-flow analysis, which computes sets of assignments that might reach future program points:

$$\llbracket p : x = E \rrbracket \text{in} = \text{in} - \text{kill}_x \cup \{p\}$$

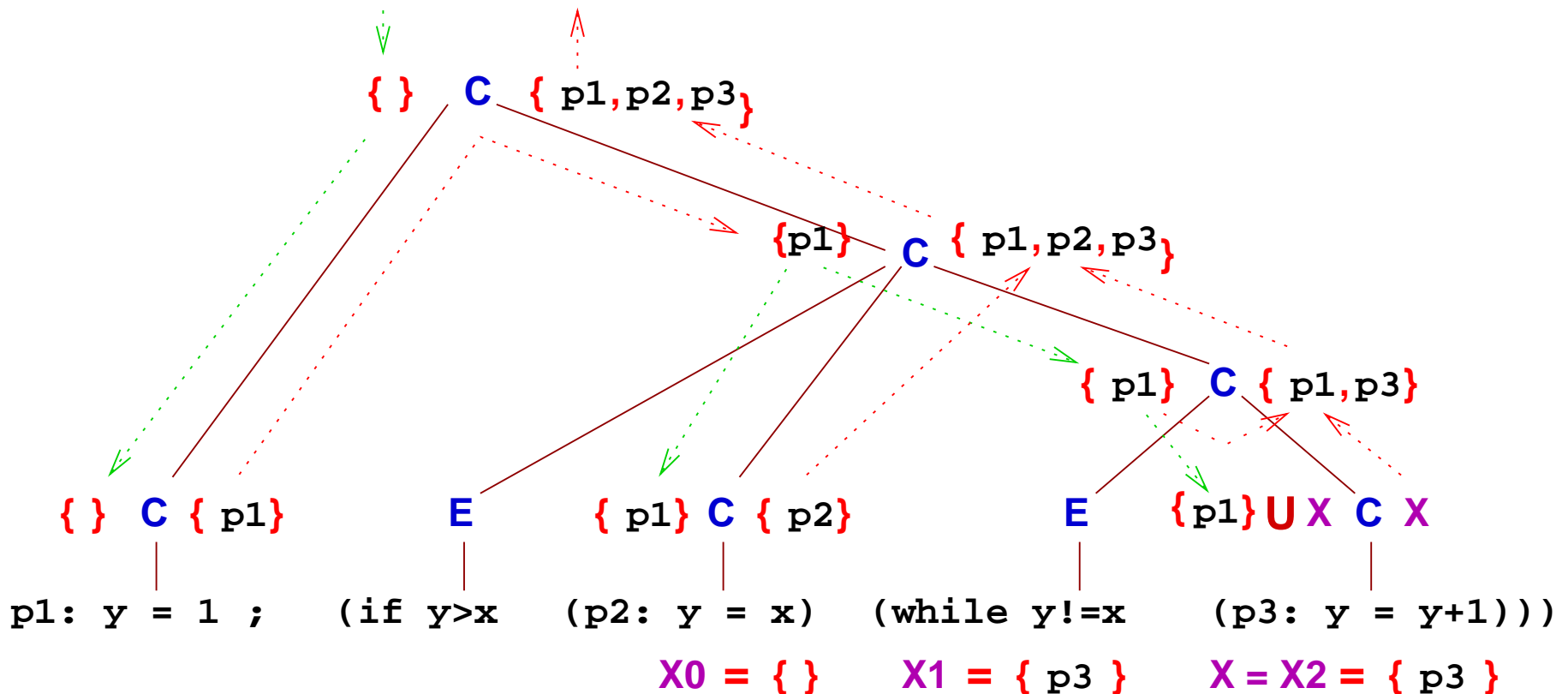
$$\llbracket \text{while } E \text{ } C \rrbracket \text{in} = \text{in} \cup \bigcup_{i \geq 0} \text{out}_i,$$

$$\llbracket C \rrbracket \text{in} = \llbracket C_2 \rrbracket (\llbracket C_1 \rrbracket \text{in})$$

where $\text{out}_0 = \{\}$

$$\llbracket \text{if } E \text{ } C_1 \text{ } C_2 \rrbracket \text{in} = \llbracket C_1 \rrbracket \text{in} \cup \llbracket C_2 \rrbracket \text{in}$$

and $\text{out}_{i+1} = \llbracket C \rrbracket (\text{in} \cup \text{out}_i)$



Big-step relational semantics: derivation trees

$$\sigma \vdash p : x = E \Downarrow f_p(\sigma)$$

$$\frac{\sigma \vdash C_1 \Downarrow \sigma_1 \quad \sigma_1 \vdash C_2 \Downarrow \sigma_2}{\sigma \vdash C_1; C_2 \Downarrow \sigma_2}$$

$$\frac{f_{E_t}(\sigma) \vdash C_1 \Downarrow \sigma_1 \quad f_{E_f}(\sigma) \vdash C_2 \Downarrow \sigma_2}{\sigma \vdash \text{if } E \ C_1 \ C_2 \Downarrow \sigma_1 \sqcup \sigma_2}$$

$$\frac{f_{E_t}(\sigma) \vdash C \Downarrow \sigma' \quad \sigma' \vdash \text{while } E \ C \Downarrow \sigma''}{\sigma \vdash \text{while } E \ C \Downarrow f_{E_f}(\sigma) \sqcup \sigma''}$$

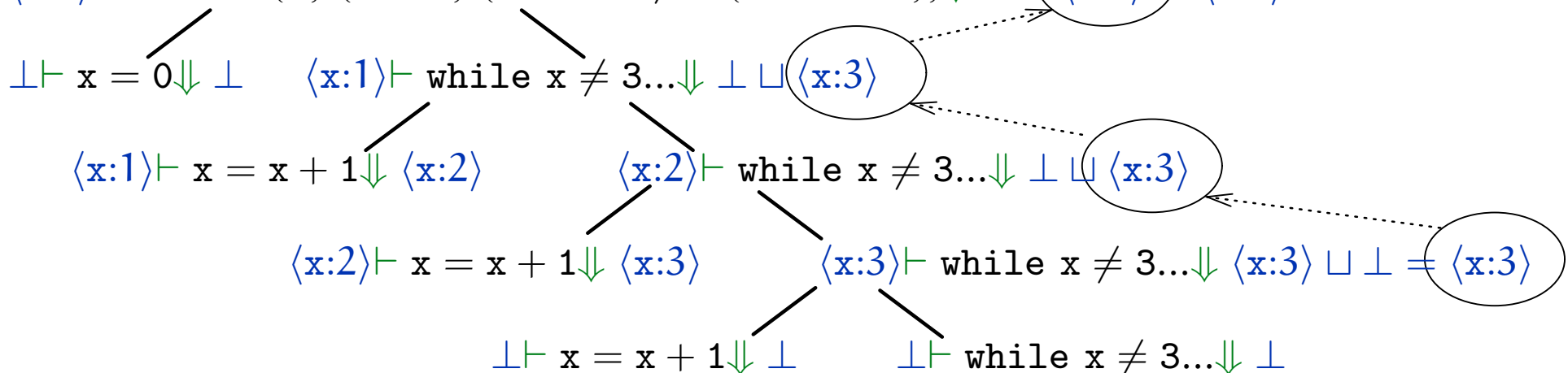
$$\perp \vdash C \Downarrow \perp$$

Recall that f_p is a transfer function and that f_{E_t} and f_{E_f} "filter" the store, e.g.,

$$f_{x>2t} \langle x : 4, y : 3 \rangle = \langle x : 4, y : 3 \rangle, \text{ whereas } f_{x>2t} \langle x : 0, y : 3 \rangle = \perp.$$

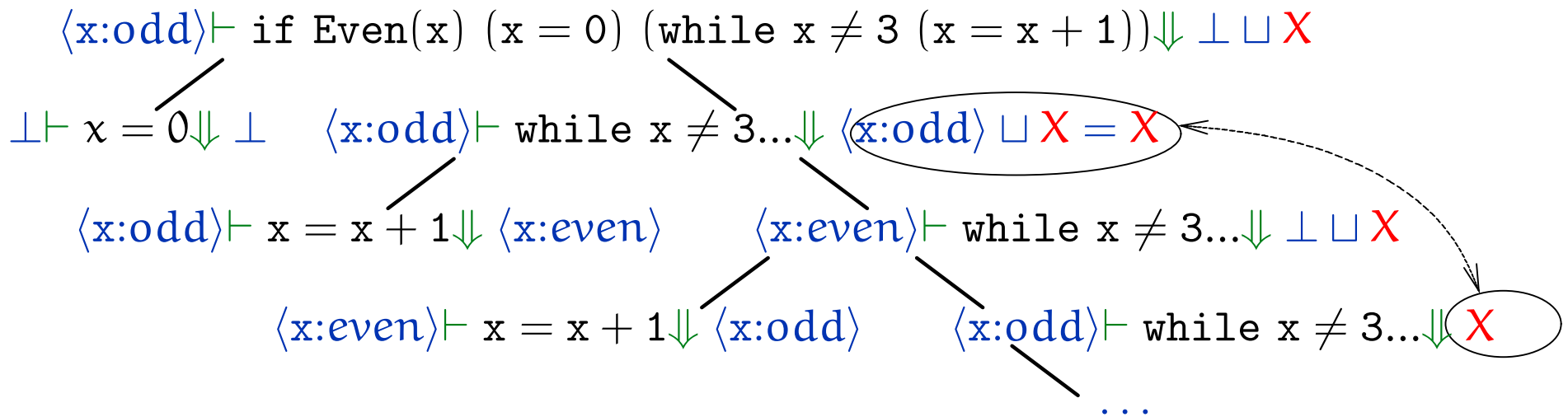
An example: if Even(x) (x=0) (while x ≠ 3 (x = x+1))

$$\langle x:1 \rangle \vdash \text{if Even}(x) \ (x = 0) \ (\text{while } x \neq 3 \ (x = x + 1)) \Downarrow \perp \sqcup \langle x:3 \rangle = \langle x:3 \rangle$$



An abstract big-step derivation tree

Using the same inference rules but with abstract transfer functions for $\text{Parity} = \{\perp, \text{even}, \text{odd}, \top\}$, we generate an abstract tree that is *infinite* but *regular*:



Variable X denotes the answer from the repeated loop subderivation:

$$X = \langle x:\text{odd} \rangle \sqcup X$$

The least solution sets $X = \langle x:\text{odd} \rangle$.

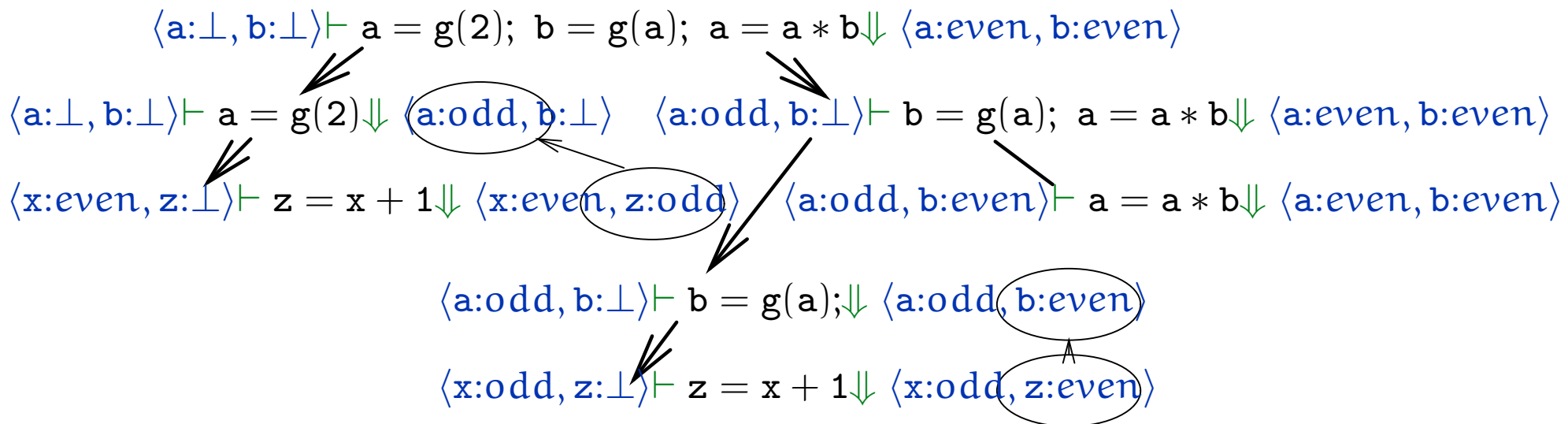
Interprocedural analysis

$$\frac{\text{func } f(x) \text{ local } y; C. \quad [x \mapsto \llbracket E \rrbracket \sigma][y \mapsto \perp] \vdash C \Downarrow \sigma'}{\sigma \vdash z = f(E) \Downarrow \sigma[z \mapsto \sigma'(y)]}$$

where $\llbracket E \rrbracket \sigma$ denotes E 's value with σ , and $x \mapsto v$ assigns v to x .

Example:

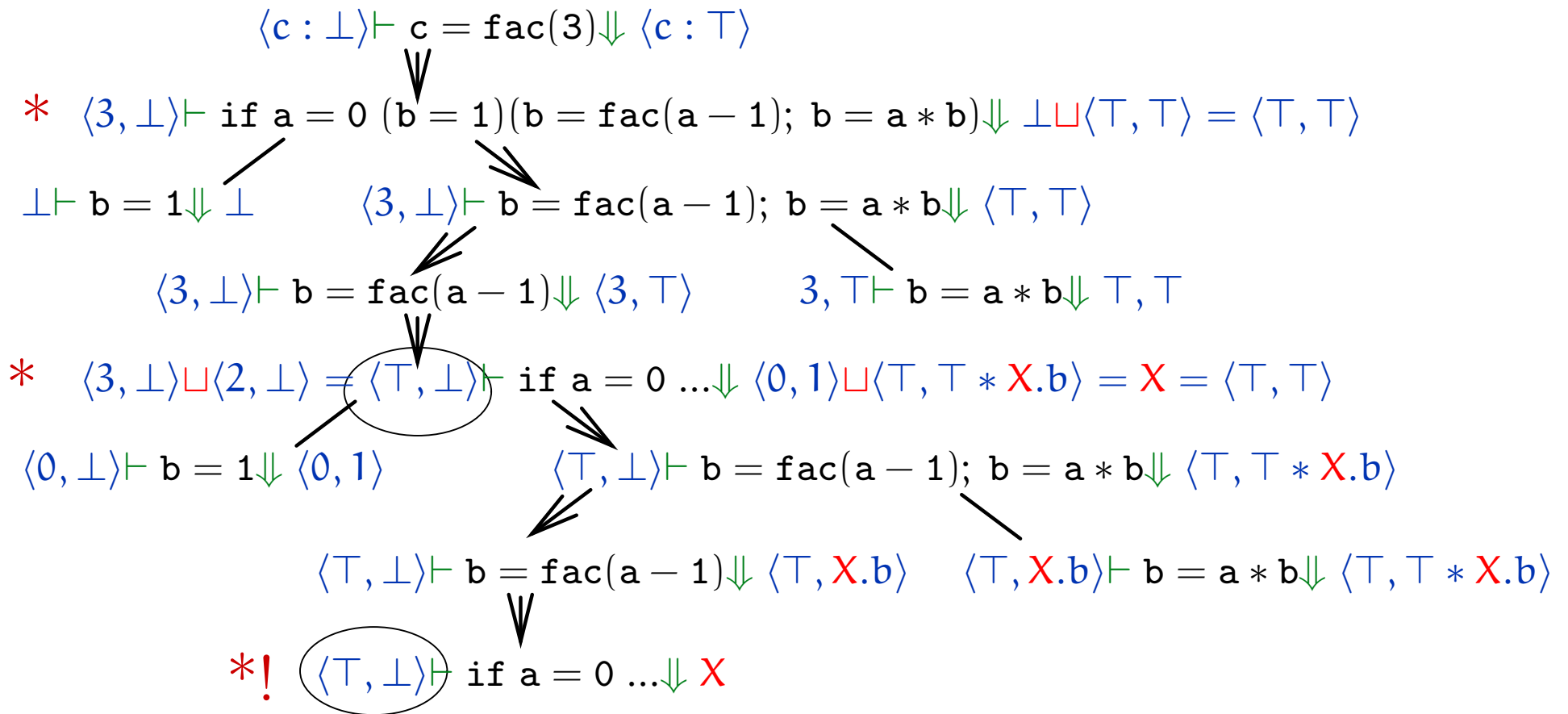
```
func g(x) local z; z = x+1.
a = g(2); b = g(a); a = a*b
```



The derivation tree naturally separates the calling contexts.

“Too many” calling contexts (*) force widening (!):

```
func fac(a) local b; if a = 0 (b = 1) (b = fac(a - 1); b = a * b).
c = fac(3)
```

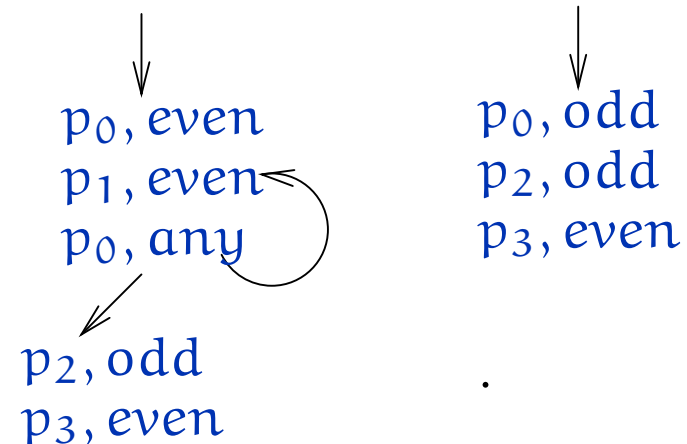


$X = \langle 0, 1 \rangle \sqcup \langle T, T * X.b \rangle$. **The least solution sets $X = \langle T, T \rangle$.**

Standard applications of static analysis

Abstract testing and model generation

```
 $p_0$  : while isEven(x) {  
     $p_1$  : x = x div 2;  
}  
 $p_2$  : x = 4 * x;  
 $p_3$  : exit
```



Each trace tree denotes an abstract “test” that covers a set of concrete test cases, e.g., $\gamma(\text{even}) = \{\dots, -2, 0, 2, \dots\}$.

Forms of abstract testing:

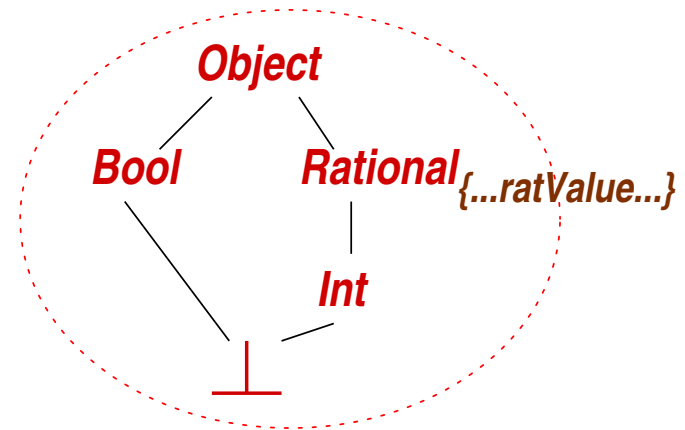
- ◆ **Black box:** For each test set, $S \subseteq C$, we abstractly interpret with $\alpha(S) \in A$. (*Best precision: ensure that $S = \gamma(\alpha(S))$.*)
- ◆ **White box:** for each conditional, B_i , in the program, ensure there is some $a_i \in A$ such that $\gamma(a_i) = \{s \mid B_i \text{ holds for } s\}$

Once we generate an abstract model, we can analyze it further — ask questions of its paths and nodes — via *model checking*.

Low-level safety checking

One example is *type casting*:

p_i : ... ((Rational) x).ratValue()...



A static analysis calculates the abstract store arriving at the cast at p_i , a *checkpoint*.

- ◆ $p_i, \langle \dots x : \text{Int} \dots \rangle$: no error possible — remove the run-time check (because $\text{Int} \sqsubseteq \text{Rational}$, hence $\gamma(\text{Int}) \subseteq \gamma(\text{Rational})$).
- ◆ $p_i, \langle \dots x : \text{Object} \dots \rangle$: possible error — retain run-time check (because $\text{Object} \not\sqsubseteq \text{Rational}$)
- ◆ $p_i, \langle \dots x : \text{Bool} \dots \rangle$: definite error, because $\text{Bool} \sqcap \text{Rational} = \perp$ (assuming $\gamma(\perp) = \{\}$).

Two more examples of low-level safety checking:

Array-bounds and arithmetic over- and under-flow checks

- ◆ *Analysis:* interval analysis, where values have form, $[i, j]$, $i \leq j$.
- ◆ *Checkpoints:* for `a[e]` — `e` has value in range, $[0, \text{a.length}]$;
for `int x = e` — `e` has value in range, $[-2^{31} - 1, +2^{31} - 1]$

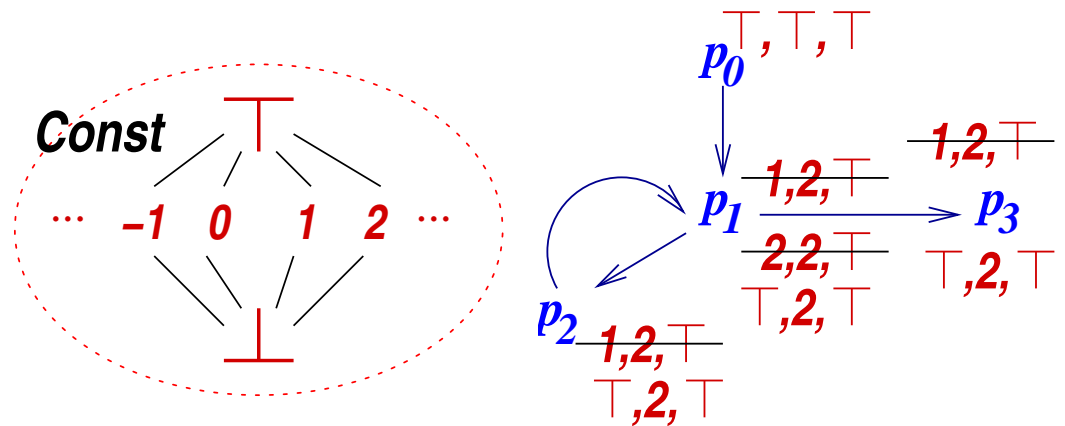
Uninitialized variables, dead-code, and erroneous-state checks

- ◆ *Analysis:* constant propagation, where values are $\{k\}$, \perp , or \top .
- ◆ *Checkpoints:*
 - uninitialized variables:* referenced variables have value $\neq \perp$;
 - dead code:* at program point p_i , arriving store has value $\neq \perp$;
 - erroneous states:* at program point $p_i : \text{Error}$, arriving store has value $= \perp$. (*Note:* This can be combined with a *backwards* analysis, starting from each $p_i : \text{Error}$ with store \top , working backwards to see if an initial state is reached.)

Program transformation: Constant folding

```

p0 : x = 1; y = 2;
p1 : while (x < y + z) {
    p2 : x = x + 1;
}
p3 : exit
    
```



The analysis tells us to replace y at p_1 by 2 :

```

x = 1; y = 2; while (x < 2 + z) x = x + 1
    
```

Basic principle of program transformation:

If $\alpha_i \in A$ arrives at point $p_i : S$, where $f_i : C \rightarrow C$ is the concrete transfer function, and there are some S', f' such that $f_i(c) = f'(c)$ for all $c \sqsubseteq_C \gamma(\alpha_i)$, then S can be replaced by S' at p_i .

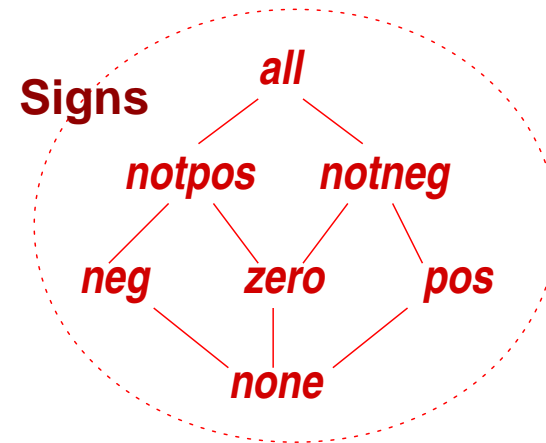
For constant folding, the transformation criteria are the abstract integers $\dots -1, 0, 1, \dots$ (but not \top).

Precondition checking and assertion synthesis

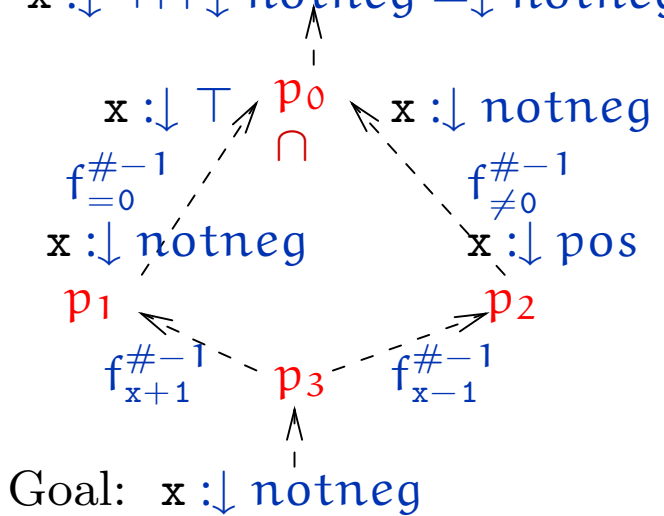
A backwards analysis synthesizes precondition assertions that ensure achievement of a postcondition:

```

p0: if x=0
    p1: then x = x+1
    p2: else x = x-1
p3: halt  ⟨x :⊥ notneg⟩
    
```



$x :⊥ \top \cap \downarrow \text{notneg} = \downarrow \text{notneg}$



where

$$f_{=0}^{\#}(a) = a \sqcap \text{zero} = \alpha \circ f_{=0} \circ \gamma$$

$$f_{\neq 0}^{\#} = \alpha \circ f_{\neq 0} \circ \gamma, \text{ e.g., } f_{\neq 0}^{\#}(\text{notneg}) = \text{pos};$$

$$f_{\neq 0}^{\#}(\text{zero}) = \perp; f_{\neq 0}^{\#}(\top) = \top$$

$$f_{+1}^{\#} = \alpha \circ f_{+1} \circ \gamma, \text{ e.g., } f_{+1}^{\#}(\text{notneg}) = \text{pos}$$

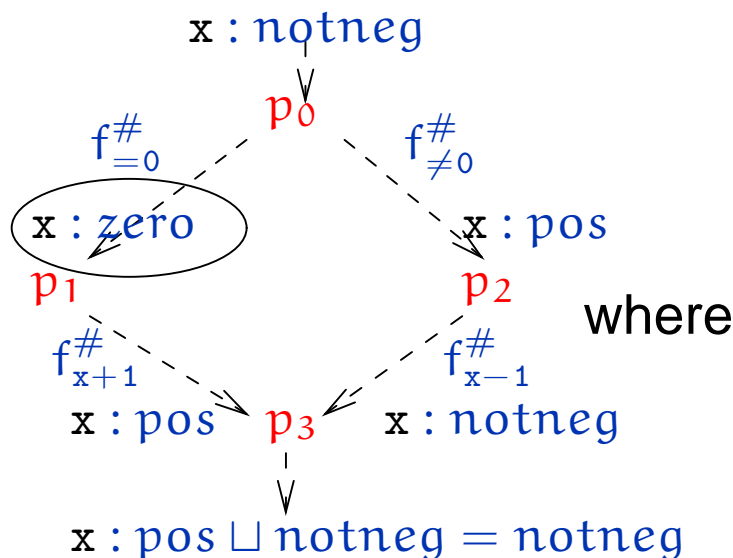
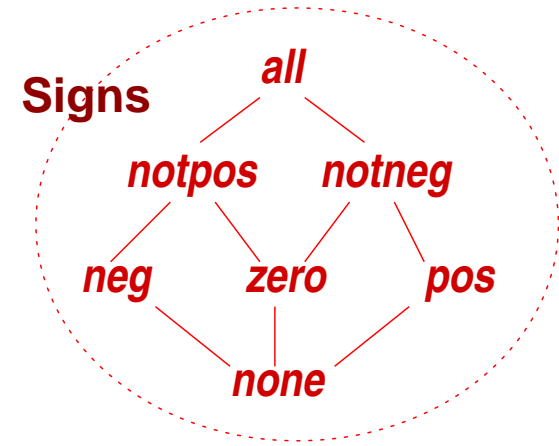
The inverse functions compute on **sets**:

$$\downarrow a = \{a' \in A \mid a' \sqsubseteq a\}$$

$$f^{\#-1}(S) = \{a \in A \mid f^{\#}(a) \in S\}$$

The entry condition can be used with a forwards analysis to generate postconditions that sharpen the assertions:

$\langle x : \text{notneg} \rangle$ p_0 : if $x=0$
 $p_1 : x = x+1$
 $p_2 : x = x-1$
 $p_3 : \text{halt}$



$$f_{=0}^\#(a) = a \sqcap \text{zero} = \alpha \circ f_{=0} \circ \gamma$$

$$f_{\neq 0}^\# = \alpha \circ f_{\neq 0} \circ \gamma, \text{ e.g., } f_{\neq 0}^\#(\text{notneg}) = \text{pos};$$

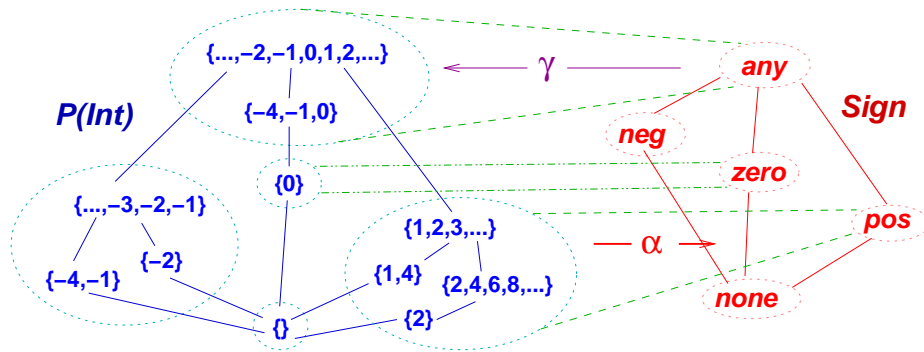
$$f_{\neq 0}^\#(\text{zero}) = \perp; f_{\neq 0}^\#(\top) = \top$$

$$f_{+1}^\# = \alpha \circ f_{+1} \circ \gamma, \text{ e.g., } f_{+1}^\#(\text{notneg}) = \text{pos}$$

The forwards-backwards analyses can be repeatedly alternated.

The “internal logic” of an abstract domain

Abstract values = logical propositions



Read properties like $neg \in Sign$ as logical propositions, “isNegative”, etc.

For $S \subseteq C$, $a, a' \in A$, $\gamma : A \rightarrow PC$, define

◆ $S \models a$ iff $S \subseteq \gamma(a)$ e.g., $\{-3, -1\} \models neg$

◆ $a \vdash a'$ iff $a \sqsubseteq a'$ e.g., $neg \vdash any$

For $f : C \rightarrow PC$, $f^\# : A \rightarrow A$ is *sound* iff $f^* \circ \gamma \sqsubseteq \gamma \circ f^\#$ iff $\alpha \circ f^* \sqsubseteq f^\# \circ \alpha$ This makes $f^\#$ a *postcondition transformer*.

Proposition: $S \models a$ implies $f^*(S) \models f^\#(a)$.

$f_{best}^\# = \alpha \circ f^* \circ \gamma$ is the *strongest liberal postcondition transformer* for f .

\mathcal{A} has an internal logic that γ preserves

First, treat all $a \in \mathcal{A}$ as primitive propositions (*isNeg*, *isPos*, etc.).

\mathcal{A} has conjunction when

$$S \models \phi_1 \sqcap \phi_2 \text{ iff } S \models \phi_1 \text{ and } S \models \phi_2, \text{ for all } S \subseteq C.$$

That is, $\gamma(\phi \sqcap \psi) = \gamma(\phi) \cap \gamma(\psi)$, for all $\phi, \psi \in \mathcal{A}$.

Proposition: When $\gamma : \mathcal{A} \rightarrow \text{PC}$ is an upper adjoint, then \mathcal{A} has conjunction.

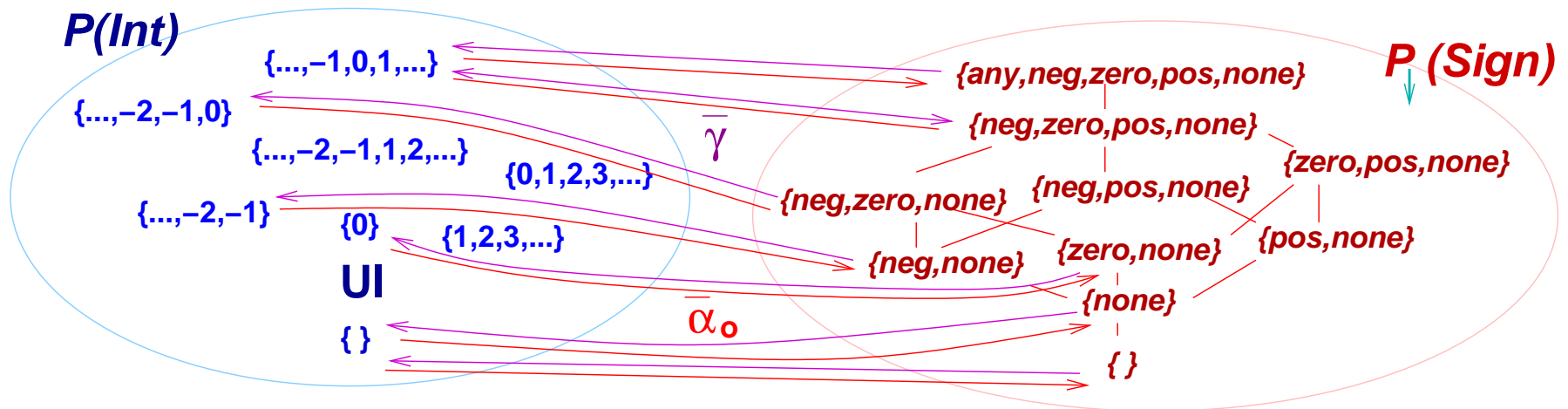
Proposition: When $\gamma(\phi \sqcup \psi) = \gamma(\phi) \cup \gamma(\psi)$, then \mathcal{A} has *disjunction*:

$$S \models \phi \sqcup \psi \text{ iff } S \models \phi \text{ or } S \models \psi.$$

Sign lacks disjunction: $\gamma(\text{zero}) \models \text{neg} \sqcup \text{pos}$, because $\text{neg} \sqcup \text{pos} = \text{any}$, but

$$\gamma(\text{zero}) \not\models \text{neg} \text{ and } \gamma(\text{zero}) \not\models \text{pos}.$$

Sometimes, we can implement a domain's disjunctive completion [Cousots79,Giacobazzi00] :



$$(\mathcal{P}(\text{int}), \subseteq) \langle \overline{\alpha_0}, \overline{\gamma} \rangle (\mathcal{P}_\downarrow(\text{Sign}), \subseteq)$$

$$\overline{\gamma}(T) = \bigcup_{a \in T} \gamma(a) \quad \overline{\alpha_0}(S) = \downarrow\{\alpha\{c\} \mid c \in S\}$$

Downclosed sets are needed for monotonicity of key functions on the sets.

Now, $\overline{\gamma}$ preserves \cap and \cup . Properties, $a \in A$, are interpreted in $\mathcal{P}_\downarrow(A)$ as $\overline{\alpha_0}(\gamma(a)) = \downarrow\{a\}$.

For $A = \mathcal{P}_\downarrow(\text{Sign})$, these assertions are exact:

$$\phi ::= \text{neg} \mid \text{zero} \mid \text{pos} \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2$$

Complete lattice \mathcal{A} is *distributive* if $a \sqcap (b \sqcup c) = (a \sqcap b) \sqcup (a \sqcap c)$, for all $a, b, c \in \mathcal{A}$. When \sqcap is Scott-continuous, then

$$\phi \Rightarrow \psi \equiv \bigsqcup \{a \in \mathcal{A} \mid a \sqcap \phi \sqsubseteq \psi\}$$

satisfies the property, $a \vdash \phi \Rightarrow \psi$ iff $a \sqcap \phi \vdash \psi$.

Proposition: If \mathcal{A} is a distributive complete lattice, \sqcap is Scott-continuous, and upper adjoint γ is 1-1, then \mathcal{A} has *Heyting implication*, $\phi \Rightarrow \psi$, such that

$$S \models \phi \Rightarrow \psi \text{ iff } \gamma(\alpha(S)) \cap \gamma(\phi) \subseteq \gamma(\psi).$$

That is, $\gamma(\phi \Rightarrow \psi) = \bigcup \{S \in \gamma[\mathcal{A}] \mid S \cap \gamma(\phi) \subseteq \gamma(\psi)\}$.

Heyting implication is weaker than classical implication, where $S \models \phi \Rightarrow \psi$ iff $S \cap \gamma(\phi) \subseteq \gamma(\psi)$ iff for all $c \in S$, if $\{c\} \models \phi$, then $\{c\} \models \psi$.

The POS domain for groundness analysis of logic programs uses Heyting implication [Cortesi91,Marriott93].

If $\gamma(\perp_{\mathbf{A}}) = \emptyset \in \mathcal{P}(\Sigma)$, we have falsity (\perp); this yields the logic,

$$\phi ::= \mathbf{a} \mid \phi_1 \sqcap \phi_2 \mid \phi_1 \sqcup \phi_2 \mid \phi_1 \Rightarrow \phi_2 \mid \perp$$

In particular, $\neg\phi$ abbreviates $\phi \Rightarrow \perp$ and defines the *refutation* of ϕ within \mathbf{A} , as done in the TVLA analyzer [Sagiv02].

$\gamma : \mathbf{A} \rightarrow \mathbf{PC}$ is the interpretation function for the internal logic:

$$\gamma(\mathbf{a}) = \text{given}$$

$$\gamma(\phi \sqcap \psi) = \gamma(\phi) \cap \gamma(\psi)$$

$$\gamma(\phi \sqcup \psi) = \gamma(\phi) \cup \gamma(\psi)$$

$$\gamma(\phi \Rightarrow \psi) = \bigcup \{S \in \gamma[\mathbf{A}] \mid S \cap \gamma(\phi) \subseteq \gamma(\psi)\}$$

$$\gamma(\perp) = \emptyset$$

γ -completeness characterizes the internal logic

The interpretation for conjunction, $\gamma(\phi \sqcap \psi) = \gamma(\phi) \cap \gamma(\psi)$, shows that γ -completeness is *exactly* the criterion for determining the connectives in \mathbf{A} 's internal logic:

Proposition: For $f : C^n \rightarrow PC$, \mathbf{A} 's logic includes connective $f^\#$ iff $f^\#$ is γ -complete for f^* :

$$\gamma(f^\#(\phi_1, \phi_2, \dots)) = f^*(\gamma(\phi_1), \gamma(\phi_2), \dots)$$

Example: For $Sign = \{none, neg, zero, pos, any\}$, $negate^\#$ is γ -complete for $negate(S) = \{-n \mid n \in S\}$ (where $negate^\#(pos) = neg$, $negate^\#(neg) = pos$, $negate^\#(zero) = zero$, etc.):

$$\phi ::= a \mid \phi_1 \sqcap \phi_2 \mid negate^\#(\phi)$$

We can state “negate” assertions, e.g., $pos \models negate^\#(neg \sqcap any)$.

Post-image (left-to-right) abstraction of relations

$f : C \rightarrow PC$ defines a relation in $C \times C$, e.g., $\{1, 3\} \llbracket \text{succ} \rrbracket \{2, 4\}$.

f 's left-to-right (post) image, $post_f : PC \rightarrow PC$, is

$$post_f(S) = \bigcup_{c \in S} f(c).$$

For Galois connection, $PC \langle \overline{\alpha}_0, \overline{\gamma} \rangle \mathcal{P}_\downarrow(A)$, and $f^\# : A \rightarrow \mathcal{P}_\downarrow(A)$,

◆ for $T \in \mathcal{P}_\downarrow(A)$, define $post_{f^\#}(T) = \sqcup_{a \in T} f^\#(a) = \bigcup_{a \in T} f^\#(a)$.

◆ use $post_{f^\#}$ to compute left-to-right (over)approximations of f , e.g., $\{neg\} \llbracket \text{succ}^\# \rrbracket \{neg, zero\}$, that is, $neg \vee zero$.

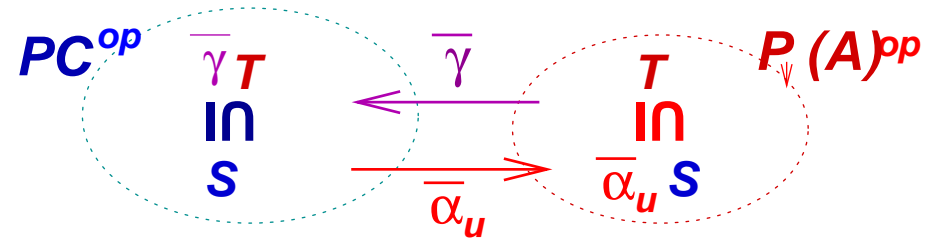
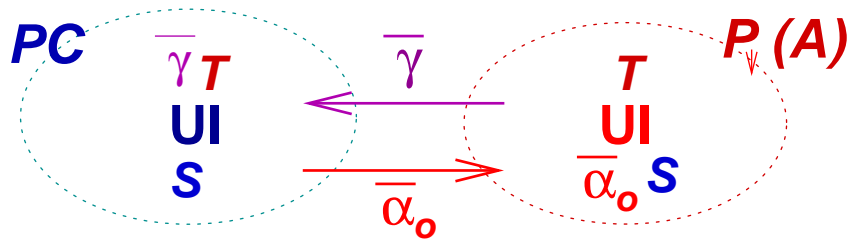
Proposition: For $f_{best}^\# = \overline{\alpha}_0 \circ f^* \circ \overline{\gamma}$,

$$(post_f)_{best}^\# = \overline{\alpha}_0 \circ post_f \circ \overline{\gamma} = post_{f_{best}^\#}$$

Corollary: If f is $\overline{\gamma}$ -complete, then $(post_{f_{best}^\#} \phi)$ is in $\mathcal{P}_\downarrow(A)$'s logic.

Given $PC\langle\alpha, \gamma\rangle A$, we have **two relevant Galois connections between PC and $\mathcal{P}_\downarrow(A)$**

Recall that $\bar{\gamma}(T) = \bigcup_{a \in T} \gamma(a)$ and that $\bar{\gamma}$ preserves both unions and intersections on $\mathcal{P}_\downarrow(A)$. Therefore, $\bar{\gamma}$ is an upper adjoint in **two** different ways:



Overapproximating abstraction:

$$\begin{aligned} \bar{\alpha}_o(S) &= \bigcap \{T \mid S \subseteq \bar{\gamma}(T)\} \\ &= \downarrow \{\alpha\{c\} \mid c \in S\} \end{aligned}$$

where

$$\downarrow T = \{a \mid \text{exists } a' \in T, a \sqsubseteq a'\}.$$

Underapproximating abstraction:

$$\begin{aligned} \bar{\alpha}_u(S) &= \bigcup \{T \mid \bar{\gamma}(T) \subseteq S\} \\ &= \{a \mid \gamma(a) \subseteq S\} \end{aligned}$$

where

$$(D, \sqsubseteq_D)^{op} \text{ is } (D, \supseteq_D).$$

Pre-image (right-to-left) abstraction of relations

$f : C \rightarrow PC$ defines a relation $\subseteq C \times C$, e.g., $\{0, 1, 3\} \llbracket \text{succ} \rrbracket \{1, 2, 4\}$.

f 's right-to-left (pre) image, $\widetilde{\text{pre}}_f : PC \rightarrow PC$, is

$$\widetilde{\text{pre}}_f(S) = \cup\{S' \subseteq C \mid f^*(S') \subseteq S\} = \{c \mid f(c) \subseteq S\}$$

For Galois connection, $PC^{\text{op}} \langle \overline{\alpha}_u, \overline{\gamma} \rangle \mathcal{P}_\downarrow(A)^{\text{op}}$ and $f^\# : A \rightarrow \mathcal{P}_\downarrow(A)$,

- ◆ for $T \in \mathcal{P}_\downarrow(A)$, define $\widetilde{\text{pre}}_{f^\#} = \{a \mid f^\#(a) \subseteq T\}$
- ◆ use $\widetilde{\text{pre}}_{f^\#}$ to compute right-to-left (under)approximations of f , e.g.,
 $\text{zero} \vee \text{pos} \llbracket \text{succ}^\# \rrbracket \text{pos}$ and $\text{none} \llbracket \text{succ}^\# \rrbracket \text{zero}$ (!)

Theorem: $(\widetilde{\text{pre}}_f)^\#_{\text{best}} = \overline{\alpha}_u \circ \widetilde{\text{pre}}_f \circ \overline{\gamma} = \widetilde{\text{pre}}_{f^\#_{\text{best}}}$.

Because $\widetilde{\text{pre}}_{f^\#} \phi$ always underapproximates $\widetilde{\text{pre}}_f(\overline{\gamma}(\phi))$, it can be added to $\mathcal{P}_\downarrow(A)$'s logic.

Indeed, we can always define an **underapproximating external logic**

For each concrete property of interest, $\llbracket \phi \rrbracket \subseteq C$, define

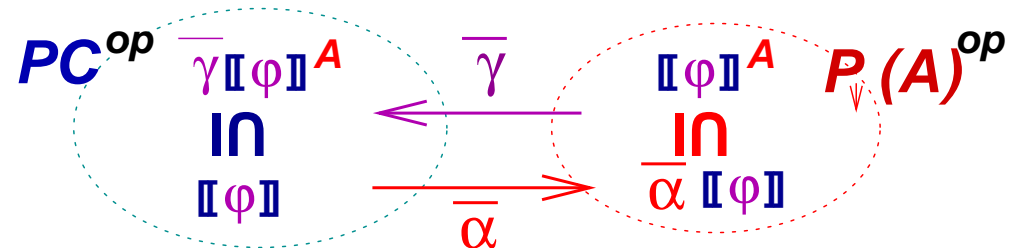
$$\llbracket \phi \rrbracket^A = \{a \in A \mid \gamma(a) \subseteq \llbracket \phi \rrbracket\}$$

Then, assert $a \vdash \phi$ iff $a \in \llbracket \phi \rrbracket^A$.

This definition follows from the underapproximating Galois connection:

$$\overline{\gamma}(T) = \bigcup \{\gamma(a) \mid a \in T\}$$

$$\overline{\alpha}_u(S) = \{a \mid \gamma(a) \subseteq S\}$$



That is, $\llbracket \phi \rrbracket^A = \overline{\alpha}_u \llbracket \phi \rrbracket$.

The inverted ordering gives *underapproximation*: $\llbracket \phi \rrbracket \supseteq \overline{\gamma}(\llbracket \phi \rrbracket^A)$. This form of external logic is standard in “abstract model checking.”

The inductively defined underapproximation to $\overline{\alpha_u}[[\phi]]$:

$$[[a]]_{\text{ind}}^A = \overline{\alpha_u}(\gamma(a))$$

$$[[\phi_1 \wedge \phi_2]]_{\text{ind}}^A = [[\phi_1]]_{\text{ind}}^A \cap [[\phi_2]]_{\text{ind}}^A$$

$$[[\phi_1 \vee \phi_2]]_{\text{ind}}^A = [[\phi_1]]_{\text{ind}}^A \cup [[\phi_2]]_{\text{ind}}^A$$

$$[[f]\phi]_{\text{ind}}^A = \widetilde{\text{pre}}_{f^\#} [[\phi]]_{\text{ind}}^A = \{a \in A \mid f^\#(a) \in [[\phi]]_{\text{ind}}^A\}$$

Entailment and provability are as expected: $a \models \phi$ iff $\gamma(a) \subseteq [[\phi]]$, and $a \vdash \phi$ iff $a \in [[\phi]]_{\text{ind}}^A$.

Soundness (\vdash implies \models) is immediate, and completeness (\models implies \vdash) follows when $\overline{\alpha_u} \circ [[\cdot]] = [[\cdot]]_{\text{ind}}^A$. This is called *logical best preservation* or *logical $\overline{\alpha}$ -completeness* [Cousots00,Schmidt06].

Scaling upwards

Analyzing large (100K+ LOC) programs

- ◆ engineered as a *one-pass analysis*, like static data-type checking
- ◆ *flow-insensitive* (ignores control-test expressions, loop iterations, distinct procedure-call points).
- ◆ *"whole-program analysis"*: examines entire source-code base

The standard example is pointer analysis on C programs, where properties are stated, "*var x may-point-to vars {y, z, ...}.*" A set of equations are generated in one program pass and solved in some small bound of iterations [Andersen94, Steensgaard96, HeintzeTardieu04] .

Advantages: simple, fast, complete code coverage, no hand-extracted "abstract model" (as required for model-checking) [Engler04]

Drawbacks: properties are simple, too many "false alarms" (inability to verify desired property)

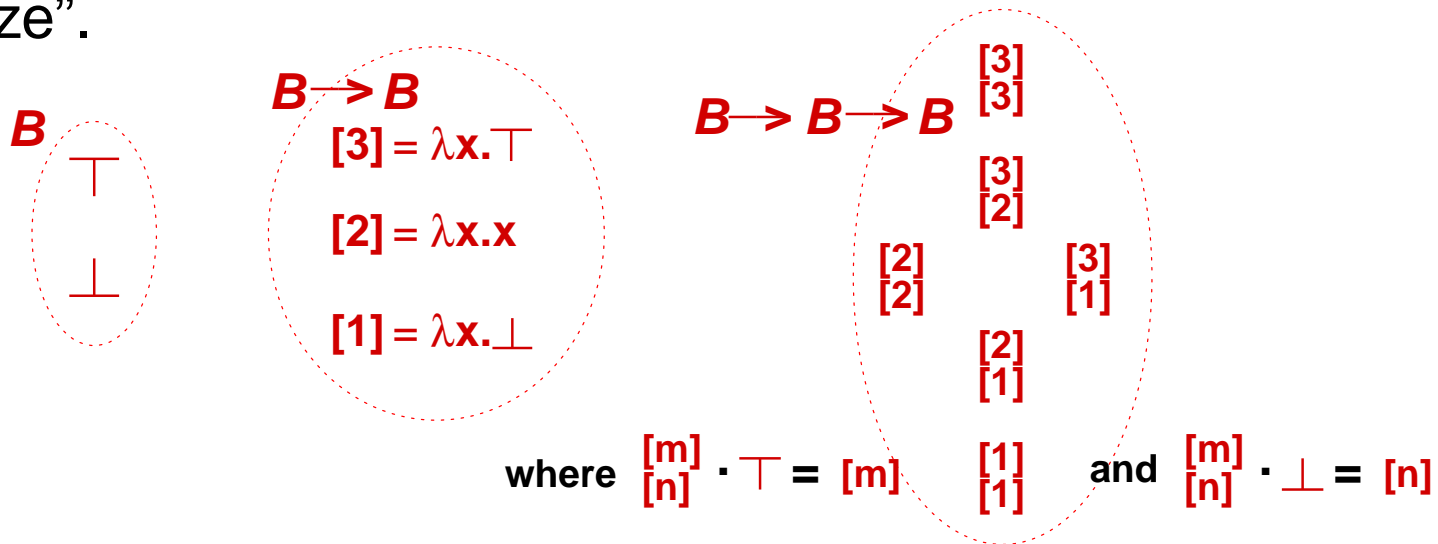
Modular analysis

- ◆ A program unit is abstracted and analyzed to a *summary structure* or *assume-guarantee relation*, where properties of the unit's free variables/inputs are associated/mapped to properties of the unit/outputs.
- ◆ When units are linked, so are their summaries, generating a composite summary. *We don't reanalyze the units.*
- ◆ Practical (better than linear-time) speedups are obtained when fixed points are solved locally within each unit (*and not at link time*) [CousotCousot02] .

There is no ideal approach, especially for the last item, so we survey some techniques (*summaries, frontiers, symbolic evaluation*) using the classic example of abstracting a higher-order function definition.

Example: higher-order normalization (“strictness”) analysis

$B = \{\perp, \top\}$, where \top means “might normalize” and \perp means “does not normalize”.



Example: $F\ m\ n = \text{if } (m=0)\ (n)\ (F\ (m+1)\ n)$

$F^\# = \lambda a: B. \lambda b: B. a \sqcap (b \sqcup (F^\#\ a\ b))$

$\text{graph}(F^\#) = \{\perp \mapsto \perp \mapsto \perp, \perp \mapsto \top \mapsto \perp, \top \mapsto \perp \mapsto \perp, \top \mapsto \top \mapsto \top\}$.

Domain B can be applied to analyses that predict the outcome of a boolean predicate/invariant (“predicate abstraction”).

A higher-order, module-like example

Define: $F^\# = \lambda f : B \rightarrow (B \rightarrow B). \lambda x : B. (x, f \cdot x)$

The function's graph (summary table) has 12 entries:

$\text{graph}(F^\#) = \{$

$\begin{bmatrix} 1 \\ 1 \end{bmatrix} \mapsto \perp \mapsto (\perp, \perp),$

$\begin{bmatrix} 1 \\ 1 \end{bmatrix} \mapsto \top \mapsto (\top, \perp),$

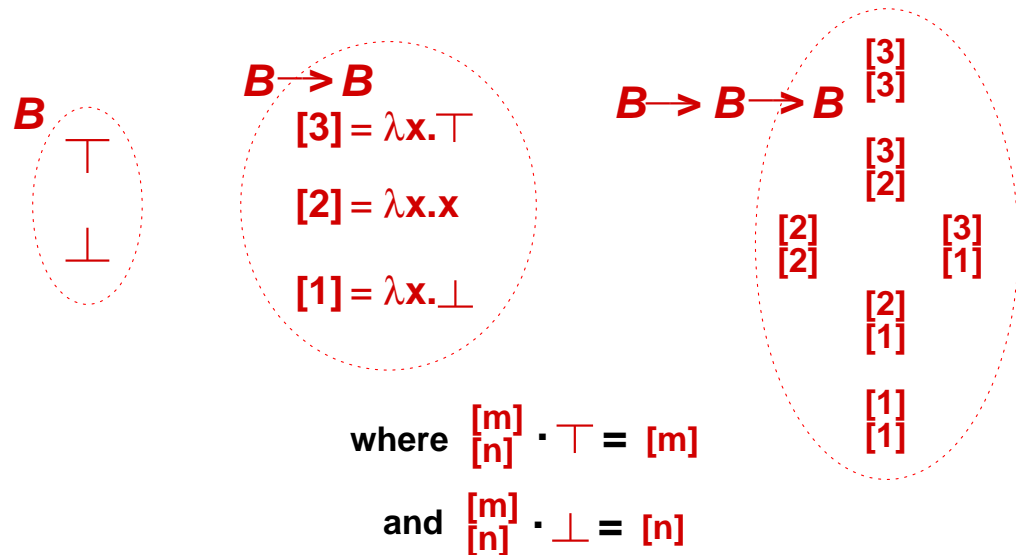
$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \mapsto \perp \mapsto (\perp, \perp),$

$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \mapsto \top \mapsto (\top, \top),$

...

$\begin{bmatrix} 3 \\ 3 \end{bmatrix} \mapsto \perp \mapsto (\perp, \top),$

$\begin{bmatrix} 3 \\ 3 \end{bmatrix} \mapsto \top \mapsto (\top, \top) \}$



It's model-checking-like
and feasible to implement!

Partial summary/graph: frontier [Clack&PeytonJones85]

Assemble the graph in increments and retain only useful ("frontier") entries, as based on these consequences of monotonicity:

- ◆ if $a \mapsto b \in \text{frontier}(F^\#)$, then (i) for all $a' \sqsubseteq a$, $a' \mapsto b$ is sound; (ii) for all $b' \sqsubseteq b$, $a \mapsto b'$ is sound.
- ◆ if $a \mapsto \top \in \text{frontier}(F^\#)$, then for all $a' \sqsupseteq a$, $a' \mapsto \top$ is sound.
- ◆ if $a_1 \mapsto b_1, a_2 \mapsto b_2 \in \text{frontier}(F^\#)$, then (i) $a_1 \sqcap a_2 \mapsto b_1 \sqcap b_2$ is sound; (ii) if $F^\#$ preserves \sqcup (holds when $F^\#$'s domain is a disjunctive completion), then $a_1 \sqcup a_2 \mapsto b_1 \sqcup b_2$ is sound.

Example frontier: for $F^\# = \lambda f : B \rightarrow (B \rightarrow B). \lambda x : B. (x, f \cdot x)$,

$$\text{frontier}(F^\#) = \left\{ \begin{array}{l} \left[\begin{array}{l} 2 \\ 2 \end{array} \right] \mapsto \perp \mapsto (\perp, [2]), \quad \left[\begin{array}{l} 2 \\ 2 \end{array} \right] \mapsto \top \mapsto (\top, [2]), \\ \left[\begin{array}{l} 3 \\ 1 \end{array} \right] \mapsto \top \mapsto (\top, [3]) \end{array} \right\}$$

Example inferences based on the frontier

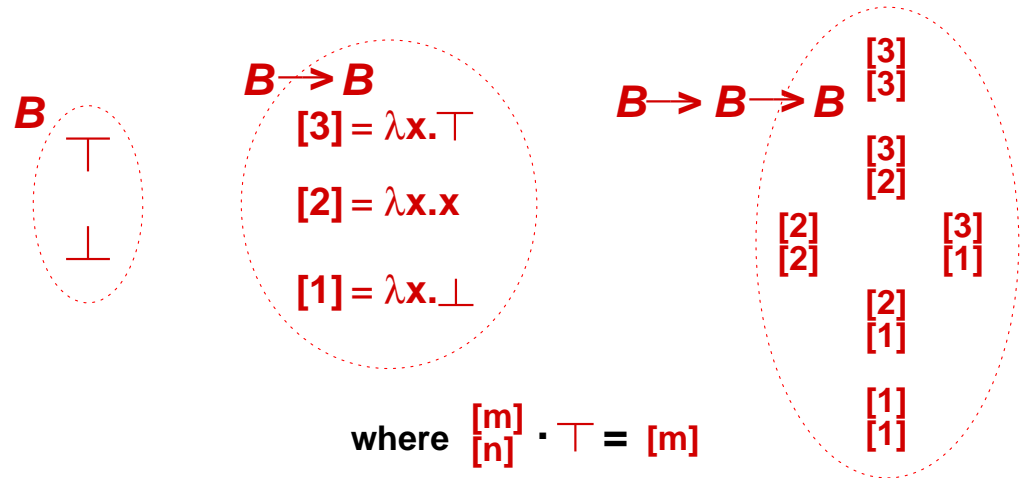
For $F^\# = \lambda f : B \rightarrow (B \rightarrow B). \lambda x : B. (x, f \cdot x)$,

frontier($F^\#$) = {

$$\begin{bmatrix} 2 \\ 2 \end{bmatrix} \mapsto \perp \mapsto (\perp, [2]),$$

$$\begin{bmatrix} 2 \\ 2 \end{bmatrix} \mapsto \top \mapsto (\top, [2]),$$

$$\begin{bmatrix} 3 \\ 1 \end{bmatrix} \mapsto \top \mapsto (\top, [3]) \},$$



where $\begin{bmatrix} m \\ n \end{bmatrix} \cdot \top = [m]$

and $\begin{bmatrix} m \\ n \end{bmatrix} \cdot \perp = [n]$

we can conclude

$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \mapsto \top \mapsto (\top, [3])$ is sound (because $\begin{bmatrix} 2 \\ 1 \end{bmatrix} \sqsubseteq \begin{bmatrix} 3 \\ 1 \end{bmatrix}$)

$\begin{bmatrix} 3 \\ 3 \end{bmatrix} \mapsto \top \mapsto (\top, [3])$ is sound (because $\begin{bmatrix} 3 \\ 1 \end{bmatrix}, \top$ map to $(\top, [3])$)

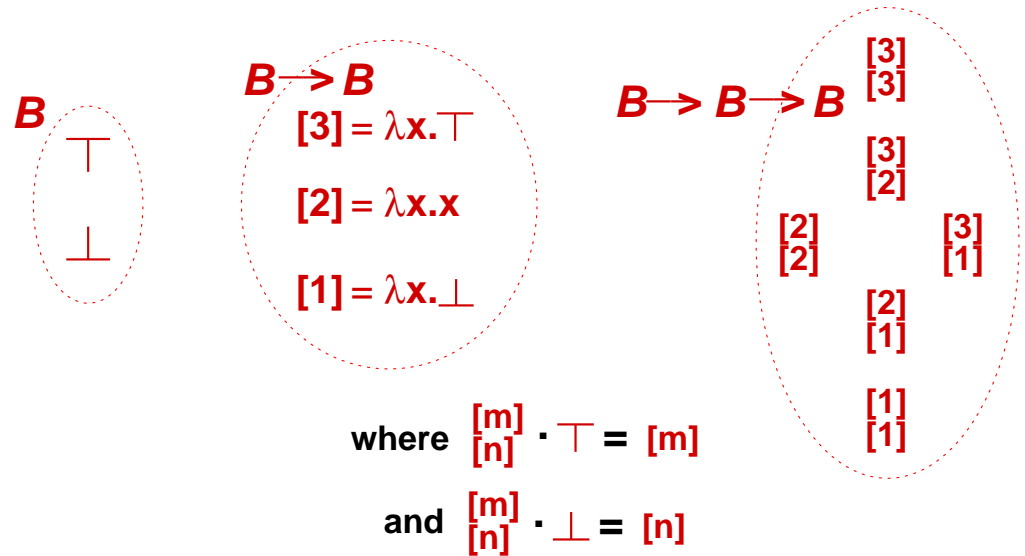
$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \mapsto \top \mapsto (\top, [2])$ is sound (because $\begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \sqcap \begin{bmatrix} 3 \\ 1 \end{bmatrix}$)

$\begin{bmatrix} 3 \\ 2 \end{bmatrix} \mapsto \top \mapsto (\top, [3])$ is sound (because $\begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \sqcup \begin{bmatrix} 3 \\ 1 \end{bmatrix}$)

Integrating symbolic evaluation with frontiers

For $F^\# = \lambda f : B \rightarrow (B \rightarrow B). \lambda x : B. (x, f \cdot x),$

$\text{symbolicFrontier}(F^\#) = \{$
 $\begin{bmatrix} 2 \\ 2 \end{bmatrix} \mapsto a \mapsto (a, a),$
 $\begin{bmatrix} 3 \\ 1 \end{bmatrix} \mapsto a \mapsto (a, \begin{bmatrix} 3 \\ 1 \end{bmatrix} \cdot a),$
 $f \mapsto a \mapsto (a, f \cdot a) \}$



- ◆ Starting from a purely symbolic formulation (the third line), the frontier expands with useful instances.
- ◆ At any point, we can replace symbolic arguments by \top to "close" the frontier, generating a "worst case analysis."
- ◆ We can apply algebraic techniques to solve local fixed points.

Solving local fixed points (intuition)

Example: $F\ x = \text{if } (\dots) \ (g\ x) \ (h(F(f\ x)))$

$$F^\# = \bigsqcup_{i \geq 0} F_i, \quad \text{where} \quad \begin{aligned} F_0 &= \lambda a. \perp \\ F_{i+1} &= \lambda a. (g\ a) \sqcup (h(F_i(f\ a))) \end{aligned}$$

By inductive reasoning,

$$\begin{aligned} F_i &= \bigsqcup_{0 \leq j < i} h^j(g(f^j\ a)) \\ &\sqsubseteq \bigsqcup_{0 \leq j < i} h^j(g(f^*\ a)) \\ &\sqsubseteq h^*(g(f^*\ a)) \end{aligned} \quad \text{where} \quad \begin{aligned} f^i &= f \circ f \circ \dots \circ f, \text{ } i \text{ times} \\ f^* &= \bigsqcup_{j \geq 0} f^j \end{aligned}$$

Each occurrence of f^* is solved locally, cheaply. The reasoning is implemented with regular tree/expression techniques; precision is traded for speed-up [CousotCousot02,Moeller03].

References

1. *This talk:* santos.cis.ksu.edu/schmidt/SLS13.pdf
2. B. Blanchet, P. Cousot, et al. A static analyzer for large safety-critical software. ACM PLDI 2003.
3. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs. ACM POPL 1977.
4. Patrick Cousot, Radhia Cousot: Modular Static Program Analysis. CC 2002, Springer LNCS 2304.
5. F. Nielson, H.R. Nielson, C. Hankin. *Principles of Program Analysis*. Springer 2010.
6. D. Schmidt. Trace-Based Abstract Interpretation of Operational Semantics. J. Lisp and Symbolic Computation 10 (1998).
7. D. Schmidt. Internal and External Logics of Abstract Interpretations. VMCAI 2008, Springer LNCS 4905.