

Interprocedural Analysis and Abstract Interpretation



Outline

- Interprocedural analysis
 - control-flow graph
 - MVP: “Meet” over Valid Paths
 - Making context explicit
 - Context based on call-strings
 - Context based on assumption sets
- Abstract interpretation

Control-flow graph for a whole program

- At each function definition `proc p(x)`
 - Create two special CFG nodes:
 - `init(p)` and `final(p)`
 - Build CFG for the function body
 - Use `init(p)` as the function entry node
 - Connect every return node to `final(p)`
- At each function call to `p(x)` with
 - Split the original function call into two stmts
 - `Enter p(x)` (before making the call) and `exit p(x)` (after the call exits)
 - Connect `enter p(x) -> init(p)`, `final(p) -> exit p(x)`
 - Connect `enter p(x) -> exit p(x)` to allow the flow of extra context info
- Three kinds of CFG edges
 - Intra-procedural: internal control-flow within a procedure
 - Procedure calls: from `enter p(x)` to `init(p)`
 - Procedure returns: from `final(p)` to `exit p(x)`

Interprocedural CFG Example

```
int fib(int z) {  
  if (z < 3) then return 1;  
  else return fib(z-1) + fib(z-2);  
}  
Main program: return fib(15);
```

A0: enter fib(15)

A1: t = exit fib(15)

B5: return 1

B6: final(fib)

B0: init(fib)

B1: if (z < 3)

B2: enter fib(z-1)

B3: t1=exit fib(z-1)
enter fib(z-2)

B4: t2=exit fib(z-2)
return t1+t2;

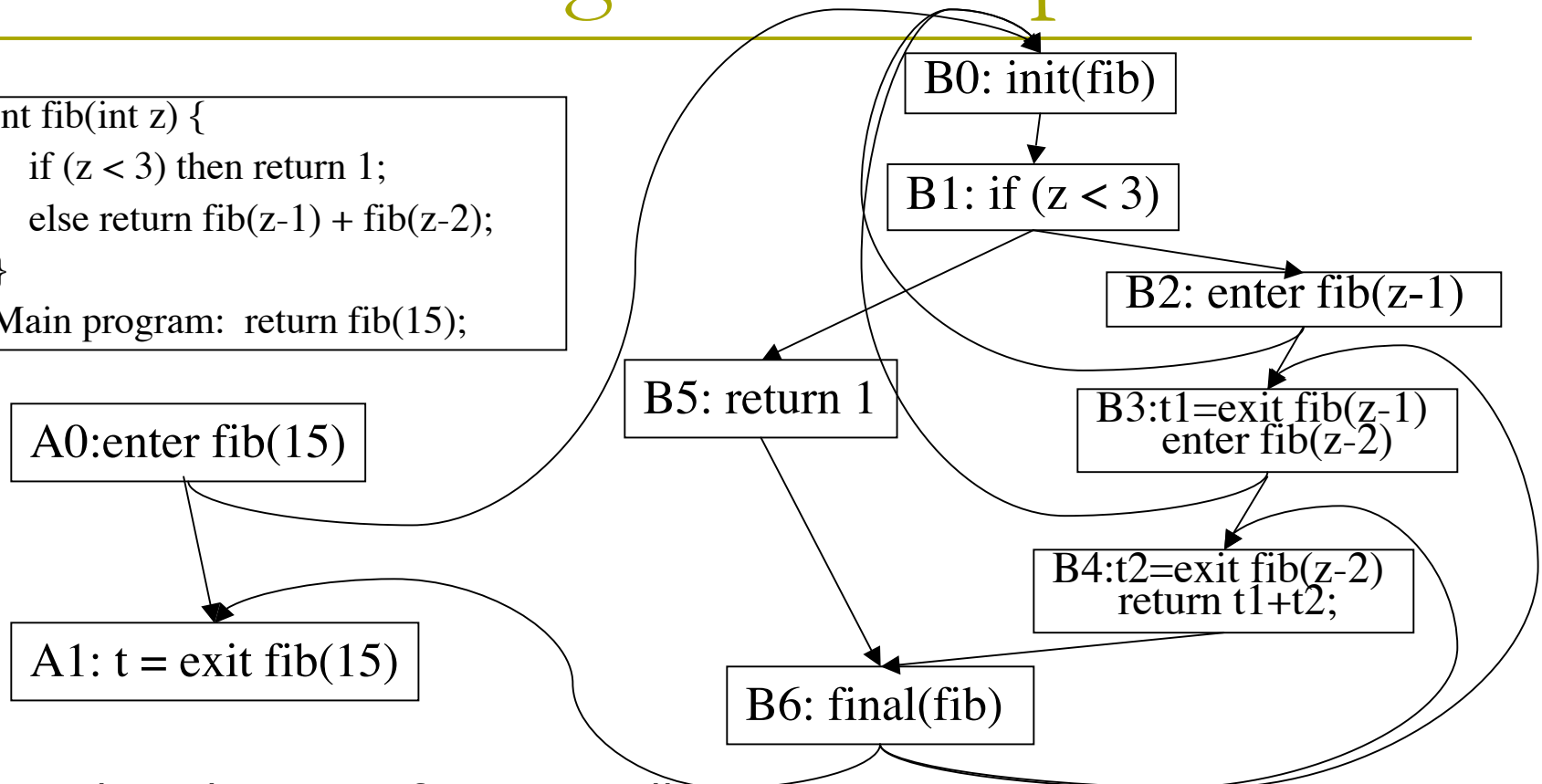
- Problem: matching between function calls and returns

Extending monotone frameworks

- Monotone frameworks consists of
 - A complete lattice (L, \leq) that satisfies the Ascending Chain Condition
 - A set F of monotone transfer functions from L to L that
 - contains the identity function and
 - is closed under function composition
- Transfer functions for procedure definitions
 - For simplicity, both $\text{init}(p)$ and $\text{final}(p)$ have identity transfer functions
- Transfer functions for procedure calls
 - For procedure entry: assign values to formal parameters
 - For procedure exit: assign return values to outside

Problem: calling context upon return

```
int fib(int z) {  
  if (z < 3) then return 1;  
  else return fib(z-1) + fib(z-2);  
}  
Main program: return fib(15);
```



- Matching between function calls and returns
 - Calculating solutions on non-existing paths could seriously detriment precision
 - E.g. enter fib(z-2) -> init(fib) -> ... -> exit fib(z-1) -> ...

MVP: “Meet” over Valid Paths

- Problem: matching procedure entries and exits (function calls and returns)
- A complete path must
 - Have proper nesting of procedure entries and exits
 - A procedure always return to the point immediately after it is called
- A valid path must
 - Start at the entry node of the main program
 - All the procedure exits match the corresponding entries
 - Some procedures may be entered but not yet exited
- The MVP solution
 - At each program point t , the solution for t is
 - $MVP(t) = \bigwedge \{ sol(p) : p \text{ is a valid path to } t \}$

Making Context Explicit

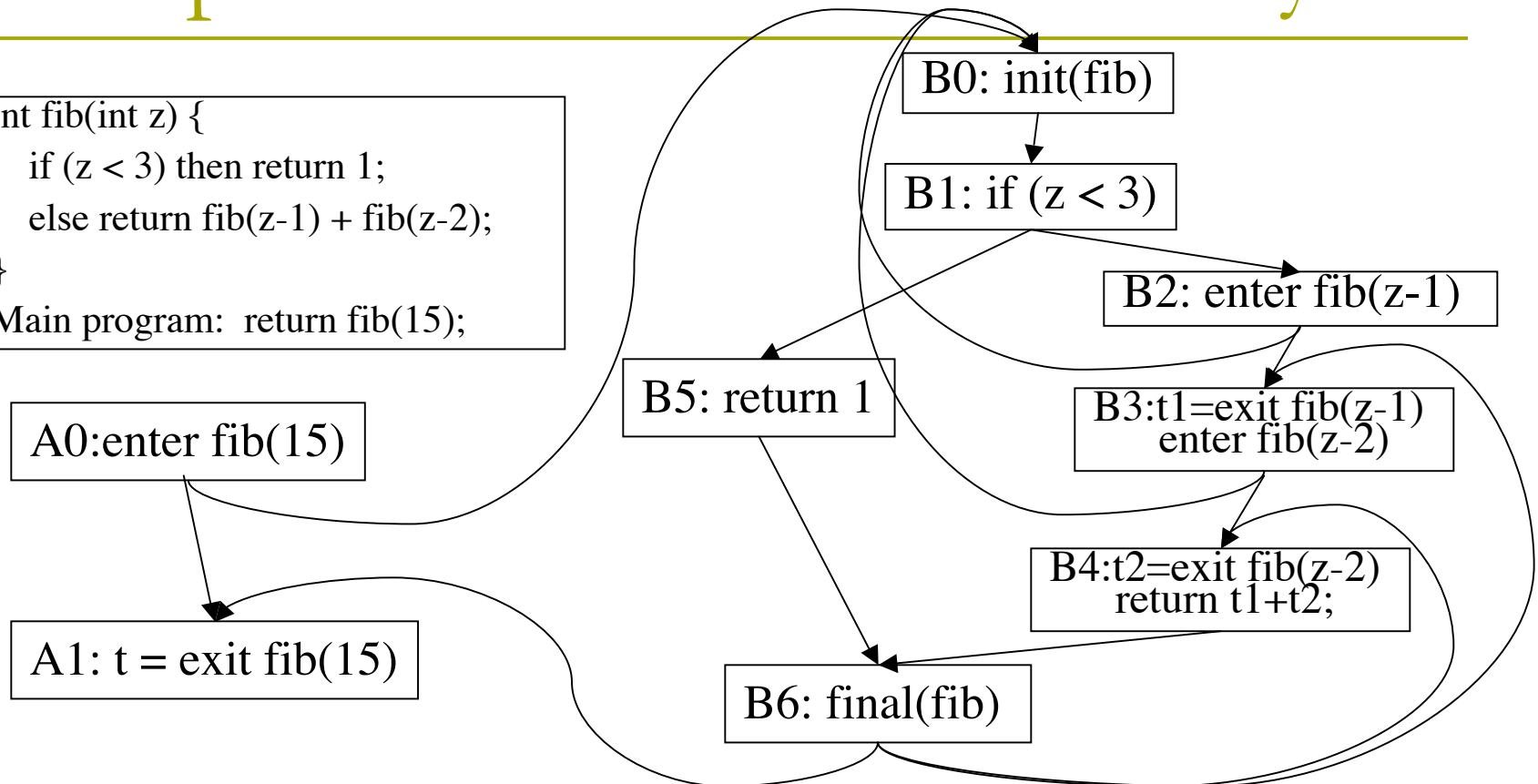
- Context sensitive analysis
 - Maintain separate solutions for different callers of a function
- Extending the monotone framework
 - Starting point (context-insensitive)
 - A complete lattice (L, \leq) that satisfies the Ascending Chain Condition
 - $L = \text{Power}(D)$ where D is the domain of each solution
 - A set F of monotone transfer functions from L to L
 - Extension
 - $L = \text{Power}(D * C)$, where C includes all calling contexts
 - $F = L \rightarrow L$, a separate sub-solution is calculated for each calling context
 - F (procedure entry) : attach caller info. to incoming solution
 - F (procedure exit): match caller info, eliminate solution for invalid paths

Different Kinds of Context

- Call strings --- contexts based on control flow
 - Remember a list of procedure calls leading to the current program point
 - Call strings of unbounded length --- remember all the preceding calls
 - Call strings of bounded length (k) --- remember only the last k calls
- Assumption sets --- contexts based on data flow
 - Assumption sets
 - Use the solution before entering proc $p(x)$ as calling context (e.g., each context makes distinct presumptions about values of function parameters)
 - Large vs. small assumption sets
 - How large is the context: use the entire solution or pick a single constraint from the solution

Example Context-sensitive Analysis

```
int fib(int z) {  
  if (z < 3) then return 1;  
  else return fib(z-1) + fib(z-2);  
}  
Main program: return fib(15);
```



- Range analysis: for each variable reference x , is its value \geq or \leq a constant value? (i.e, $x \geq x_1$; $z \leq n_2$)?

Example Range Analysis

Variables: x,z, t1, t2, fib, t; **Contexts:** A0, B2, B3,none;
Domain: Variables * (<=n, =n, >=n,?,any)

A0	(none)	(none)	(none)	(none)
B0	(none, z=?)	(A0,z=15) (B2/B3, z=?)	(A0,z=15)(B2,z>=2) (B3,z>=1)	(A0,z=15)(B2,z>=2) (B3,z>=1)
B1	(none, z=?)	(A0,z=15) (B2/B3, z=?)	(A0,z=15)(B2,z>=2) (B3,z>=1)	(A0,z=15)(B2,z>=2) (B3,z>=1)
B2	(none, z=?)	(A0,z=15) (B2/B3, z>=3)	(A0,z=15)(B2/B3,z>=3)	(A0,z=15)(B2/B3,z>=3)
B3	(none, z/t1=?)	(A0,z=15,t1=?) (B2/B3,z>=3,t1=?)	(A0,z=15,t1=1) (B2/B3,z>=3,t1=1)	(A0,z=15,t1>=1) (B2/B3,z>=3,t1>=1)
B4	(none, z/t1/t2=?)	(A0,z=15,t1/t2=?)(B2/B3,z>=3,t1/t2=?)	(A0,z=15,t1/t2=1)(B2/B3,z>=3,t1/t2=1)	(A0,z=15,t1/t2>=1)(B2/B3,z>=3,t1/t2>=1)
B5	(none, z=?)	(B2/B3,z<=2)	(B2,z=2) (B3,z<=2)	(B2,z=2) (B3,z<=2)
B6	(none,z/fib=?)	(A0,z=15,fib=?)(B2/B3,z=any,fib=1)	(A0,z=15,fib>=1)(B2/B3,z=any,fib>=1)	(A0,z=15,fib>=1)(B2/B3,z=any,fib>=1)
A1	(none,t=?)	(none,t=?)	(none,t >=1)	(none,t>=1)

Foundations of Abstract Interpretation

□ Definition from Wikipedia

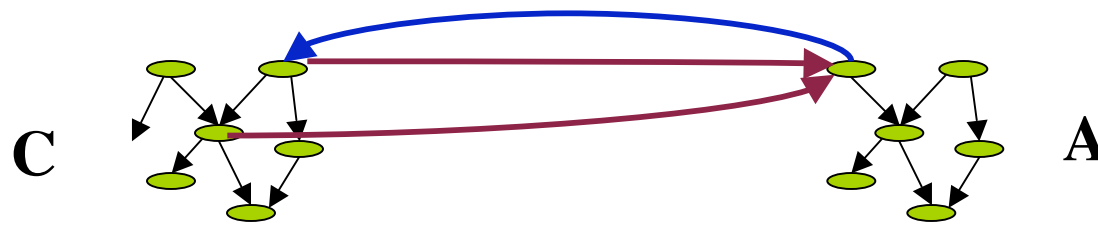
- **abstract interpretation** is a theory of sound approximation of the semantics of computer programs. It can be viewed as a partial execution of a computer program without performing all the calculations.

□ Outline

- Monotone frameworks
 - A complete lattice (L, \leq) that satisfies the Ascending Chain Condition
 - A set F of monotone transfer functions from L to L that
 - contains the identity function and
 - is closed under function composition
- Galois connections, closures, and Moore families
- Soundness and completeness of operations on abstract data
- Soundness and completeness of execution trace computation

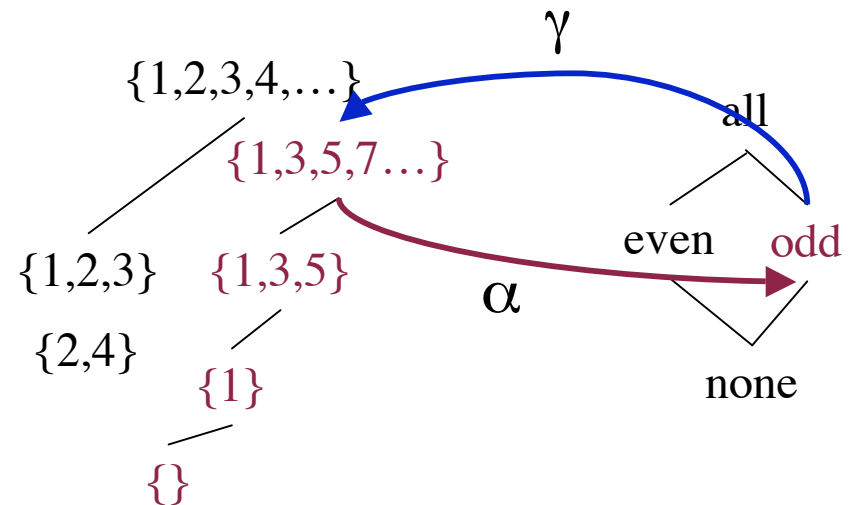
Galois Connections

- Two complete lattices
 - C: the “concrete” (execution) data
 - The execution of the entire program
 - Infinite and impossible to model precisely
 - A: the “abstract” (execution) data
 - Properties (abstractions) of the “concrete” data
 - The solution space (domain) of static program analysis
- For complete lattices C and A, a Galois connection is
 - A pair of monotonic functions, $\alpha : C \rightarrow A$, $\gamma : A \rightarrow C$
 - For all $a \in A$ and $c \in C$: $c \leq \gamma(\alpha(c))$ and $\alpha(\gamma(a)) \leq a$
 - Is Written as $C \langle \alpha, \gamma \rangle A$



Galois Connections (2)

- γ and α are inverse maps of each other's image
 - For all $c \in \gamma(A), c = \gamma(\alpha(c))$; for all $a \in \alpha(C), a = \alpha(\gamma(a))$
 - The maps α are "homomorphism" mappings between C and A
- Galois connections are closed under
 - Composition, product, and so on
- Each instruction performs an action $f: C \rightarrow C$
 - Can use α and γ to define an abstract transfer function $f\#$: $A \rightarrow A$ for each $f: C \rightarrow C$

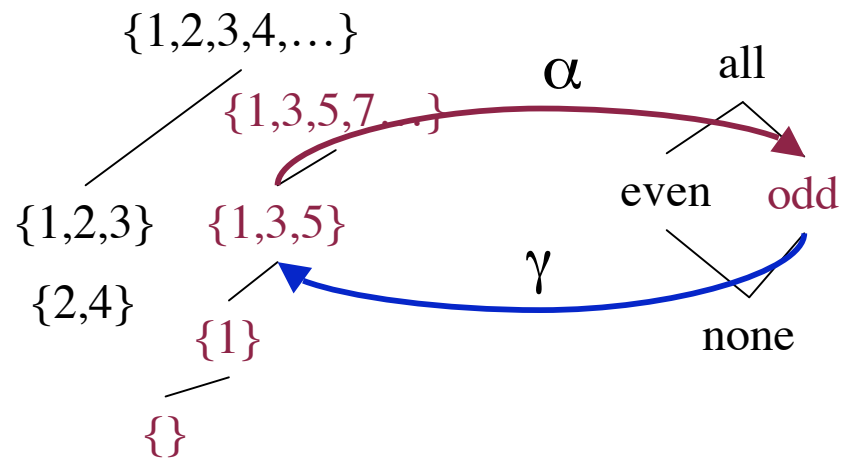


Closure Maps

- For $C \langle \alpha, \gamma \rangle A$, it is common that $A \subseteq C$. This means A embeds into C as a sub-lattice
 - A 's elements name distinguished sets in C
- A closure map defines the embedding of A within C .

Definition: $\rho: C \rightarrow C$ is a *closure map* if it is

- *Monotonic:* $\forall c_1, c_2 \in C, c_1 \leq c_2 \Rightarrow \rho(c_1) \leq \rho(c_2)$;
- *extensive:* $\forall c \in C, c \leq \rho(c)$;
- *idempotent:* $\forall c \in C, \rho(\rho(c)) = \rho(c)$ (i.e. $\rho \circ \rho = \rho$)



- 1) Every Galois connection, $C \langle \alpha, \gamma \rangle A$ defines a closure map $\alpha \circ \gamma$;
- 2) Every closure map, $\rho: C \rightarrow C$, defines the Galois connection, $C \langle \rho, \text{id} \rangle \rho(C)$.

Moore Families

- Given C , can we define a closure map on it by choosing some elements of C ?
 - Yes, if the elements we select are closed under greatest-lower-bounds (meet) operation
 - That is, the new set of elements forms a complete lattice
- **Definition:** $M \subseteq C$ is a *Moore family* iff for all $S \subseteq M$, $(\wedge S) \in M$.
 - We can define a closure map as $\rho(c) = \wedge \{c' \in M \mid c \leq c'\}$.
 - That is, we map each element in C to the closest abstraction (approximation) in M
- For each closure map, $\rho: C \rightarrow C$, its image, $\rho(C)$, is a Moore family.

Given C , we can define an abstract interpretation by selecting some $M \subseteq C$ that is a Moore family

Closed Binary Relations

- Often the solution of an analysis is a power set of its domain
 - The Galois connection can be written as $\text{Power}(D) \langle \alpha, \gamma \rangle A$
- Given unordered set D and complete lattice A , it is natural to relate the elements in D to those in A by **a binary relation, $R \subseteq D * A$, s.t.**
 - $(d, a) \in R$ (or $d R a$, $d \models_R a$) means “ d has property a ”.
 - **Example:** $D = \text{Int}$, $A = \{\text{none, neg, pos, zero, nonneg, nonpos, any}\}$.
 - Then $2 R \text{nonneg}$, $2 R \text{pos}$, and $2 R \text{any}$.
- The adjoint function, $\gamma : A \rightarrow \text{Power}(D)$, can be defined as
 - $\gamma(a) = \{d \in D \mid d R a\}$. E.g., $\gamma(\text{nonneg}) = \{0, 1, 2, \dots\}$.
 - If R defines a Galois collection, then $\gamma(A)$ defines a Moore family.
- **Proposition:** $R \subseteq D * A$ defines a Galois connection between $(\text{Power}(D), A)$ iff
 - R is *U-closed*: $c R a$ and $a \leq a'$ imply $c R a'$;
 - R is *G-closed*: $c R \wedge \{a \mid c R a\}$

Concrete and Abstract Operations

- Now that we know how to model a solution space via abstraction function $\alpha : C \rightarrow A$,
 - We must model concrete computation steps, $f:C \rightarrow C$, by abstract computation steps, $f#:A \rightarrow A$.
- **Example:** we have concrete domain, Nat , and concrete operation, $\text{succ}: \text{Nat} \rightarrow \text{Nat}$, defined as $\text{succ}(n)=n+1$.
 - abstract domain, $\text{Parity} = \{\text{any}, \text{even}, \text{odd}, \text{none}\}$.
 - abstract operation, $\text{succ\#}: \text{Parity} \rightarrow \text{Parity}$, defined as
 $\text{succ\#}(\text{even})=\text{odd}$, $\text{succ\#}(\text{odd})=\text{even}$, $\text{succ\#}(\text{any})=\text{any}$,
 $\text{succ\#}(\text{none})=\text{none}$,
 - succ\# must be consistent (sound) with respect to succ :
if $n \text{ Rn } a$, then $\text{succ}(n) \text{ Rn } \text{succ\#}(a)$,
 - where $\text{Rn} \subseteq \text{Nat} * \text{Parity}$ relates numbers to their parities (e.g., $2 \text{ Rn } \text{even}$, $5 \text{ Rn } \text{odd}$, etc.).

Sound Approximation

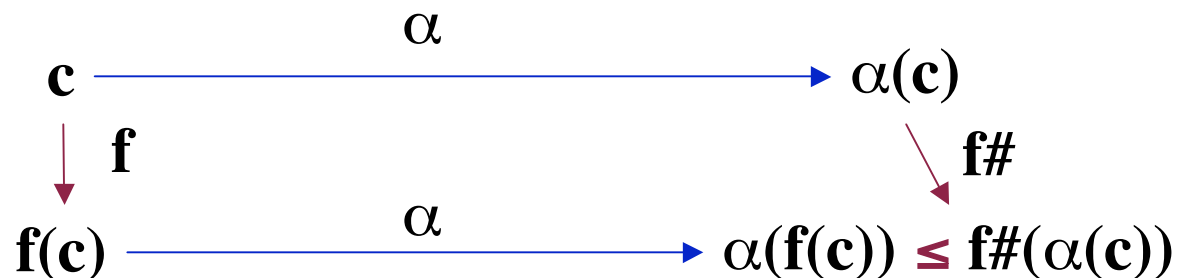
□ Given

- Galois connection $C \langle \alpha, \gamma \rangle A$ and
- functions $f : C \rightarrow C$ and $f\# : A \rightarrow A$,

$f\#$ is a *sound approximation* of f iff

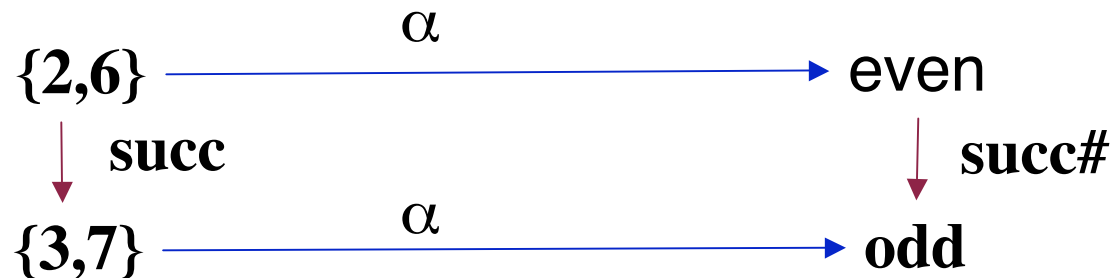
- For all $c \in C$, $\alpha(f(c)) \leq f\#(\alpha(c))$
- For all $a \in A$, $f(\gamma(a)) \leq \gamma(f\#(a))$

- ## □ That is, α defines a “semi-homomorphism” with respect to f and $f\#$



Sound Approximation Example

- Given
 - Galois connection $\text{Power}(\text{Nat}) \langle \alpha, \gamma \rangle \text{Parity}$ and
 - Concrete transfer function $\text{succ} : \text{Nat} \rightarrow \text{Nat}$, $\text{succ}(S) = \{ n + 1 \mid n \in S \}$
 - Abstract transfer function $\text{succ\#} : \text{Parity} \rightarrow \text{Parity}$,
 $\text{succ\#}(\text{even}) = \text{odd}$, $\text{succ\#}(\text{odd}) = \text{even}$
 $\text{succ\#}(\text{any}) = \text{any}$, $\text{succ\#}(\text{none}) = \text{none}$
- succ\# is a *sound approximation* of succ
 - For all $c \in \text{Nat}$, $\alpha(\text{succ}(c)) = \text{succ\#}(\alpha(c))$



Synthesizing $f\#$ from f

- Given $C \langle \alpha, \gamma \rangle A$, and function $f : C \rightarrow C$, the most precise $f\# : A \rightarrow A$ that is sound with respect to f is
 - $f\# \text{ best}(a) = \alpha (f (\gamma (a)))$
- **Proposition:** $f\#$ is sound with respect to f iff
 - For all $a \in A$, $f\# \text{ best}(a) \leq f\#(a)$
 - Of course, $f\# \text{ best}$ has a *mathematical* definition— not an algorithmic one— $f\# \text{ best}$ might not be finitely computable!
- **Parity example continued:**
 - $\text{succ}\# \text{ best}(\text{even}) = \alpha (\text{succ} (\gamma (\text{even}))) = \alpha (\text{succ} \{2n \mid n \geq 0\}) = \alpha (\{2n+1 \mid n \geq 0\}) = \text{odd}$
 - Question: what about other operators on Nat , e.g., $*$, $/$?

Completeness of Approximation(skip)

Given $C \langle \alpha, \gamma \rangle A$, and function $f : C \rightarrow C$,

- Function $f\# : A \rightarrow A$ is sound with respect to f iff
 - For all $c \in C$, $\alpha(f(c)) \leq f\#(\alpha(c))$
 - For all $a \in A$, $f(\gamma(a)) \leq \gamma(f\#(a))$
- Function $f\# : A \rightarrow A$ is *forwards*(γ) *complete* with respect to f iff
 - For all $a \in A$, $f(\gamma(a)) = \gamma(f\#(a))$
 - That is, $\gamma(A)$ is closed under f : $f(\gamma(A)) \subseteq \gamma(A)$
- Function $f\# : A \rightarrow A$ is *backwards*(α) *complete* with respect to f iff
 - For all $c \in C$, $\alpha(f(c)) = f\#(\alpha(c))$
 - That is, α partitions C into equivalence classes: $\alpha(c) = \alpha(c')$ implies $\alpha(f(c)) = \alpha(f(c'))$
- For an $f\#$ to be (forwards or backwards) complete, it must equal $f\#_{\text{best}} = \alpha(f(\gamma(a)))$
 - The structure of $C \langle \alpha, \gamma \rangle A$ and $f : C \rightarrow C$ determines whether $f\#$ is complete.

Transfer Functions and Computation steps

- Each program transition from program point p_i to p_j has an associated *transfer function*, $f_{ij}:C \rightarrow C$ (or $f_{ij}:A \rightarrow A$), which describes the associated computation.
 - This defines a computation step of the form, $(p_i, s) \rightarrow (p_j, f_{ij}(s))$
- **Example:**
 - Assignment $p_0: x=x+1; p_1: \dots$ has the transfer function
$$f_{01}(\langle \dots x:n \dots \rangle) = \langle \dots x:n+1 \dots \rangle$$
 - For multiple transitions in conditionals, attach a transfer function to each possible transition (branch) to “filter” the data that arrives at a program point.
e.g. $p_0: \text{cases } x \leq y: p_1: y=y-x;$
 $y \leq x: p_2: x=x-y; \text{ end}$
 - $f_{p1}(s) = \text{if } s[x] \leq s[y] \text{ then } s \text{ else bot; (filter out } s \text{ unless } s[x] \leq s[y])$
 - $f_{p2}(s) = \text{if } s[y] \leq s[x] \text{ then } s \text{ else bot; (filter out } s \text{ unless } s[y] \leq s[x])$

Execution Traces

- An *execution trace* is a (possibly infinite) sequence, $(p_0, s_0) \rightarrow (p_1, s_1) \rightarrow \dots \rightarrow (p_j, s_j) \rightarrow \dots$, s.t.
 - for all $i \geq 0$: $(p_i, s_i) \rightarrow p_{\text{succ}(i)}, f_i, \text{succ}(i)(s_i)$
 - No s_i equals bot

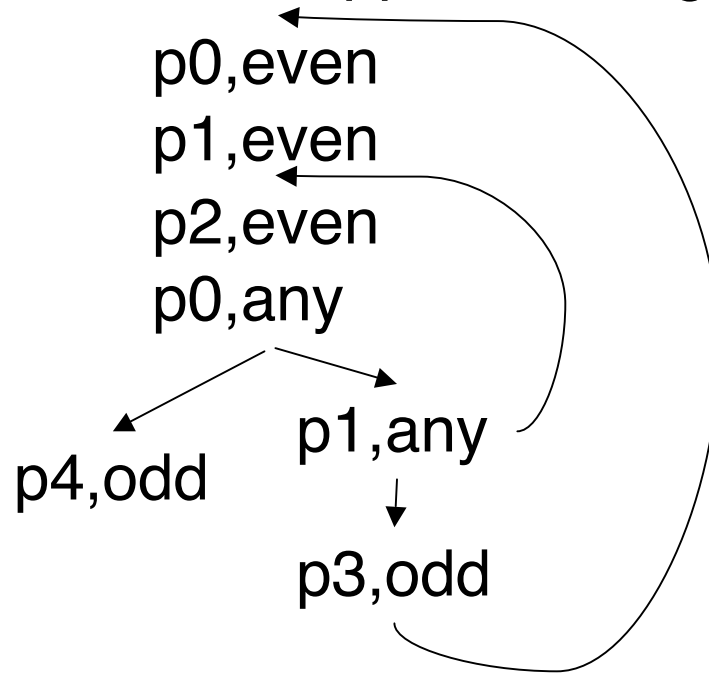
Two concrete traces
((p_i, v) means ($p_i, x=v$)):

```
P0: while (x != 1) {  
P1:  if Even(x)  
P2:    x = x div2;  
P3:  else  
      x = 3*x + 1;  
}  
P5: exit;
```

p0,4	p0,6
p1,4	p1,6
p2,4	p2,6
p0,2	p0,3
p1,2	p1,3
p2,2	p2,3
p0,1	p0,10
p4,1	p4,1
	...

Using Approximation to build abstract traces

Abstract over approximating trace:



1. Each concrete transition is generated by an f_{ij} ;
2. Each abstract transition is generated by the corresponding $f\#ij$.

- Each concrete transition, $(p_i, s) \rightarrow (p_j, f_{ij}(s))$, is reproduced by a corresponding abstract transition, $(p_i, a) \rightarrow (p_j, f\#ij(a))$, where $s \in \gamma(a)$
- The traces embedded in the abstract trace tree “cover” (*simulate*) the concrete traces

Shape Analysis

- Goal
 - To obtain a finite representation of the memory storage
- The analysis result can be used for
 - Detection of pointer aliasing
 - Detection of sharing between structures
 - Software development tools
 - Detection of pointer errors, e.g. dereferences of nil-pointers
 - Program verification
 - E.g., reverse transforms a non-cyclic list to a non-cyclic list

The Concrete Solution Space

- Model the memory (stack and heap)
 - Storage of local variables
 - Stack = Var \rightarrow (Value \cup Loc)
 - Map each local variable into a value or a unique location
 - The heap storage
 - Heap = (Loc * Sel) \rightarrow (Value \cup Loc)
 - Map pairs of locations and selectors to values or locations
- Model the operational semantics of programs
 - Program state: State = ProgramPoint * Stack * Heap
 - Example:** (p1, (x:3,y:Ly), ((Ly,val):5)) is a program state
 - Each statement modifies Stack and Heap of the previous state
 - Stmt: State \rightarrow State

Building Abstract Domains

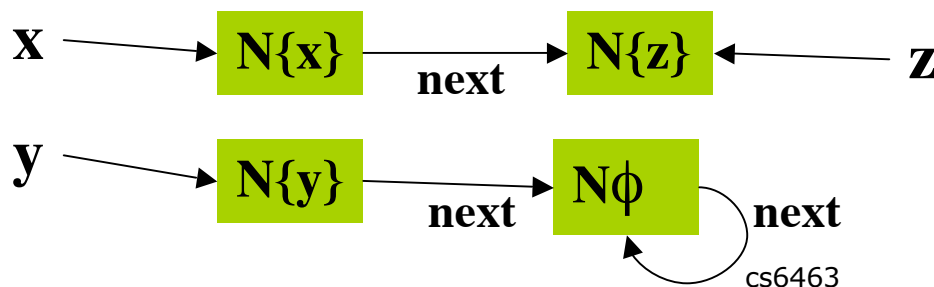
- Given an unordered set, D , of concrete data values, we might ask,
 - “What are the properties about D that I wish to calculate?”
 - Can I relate these properties $a \in A$, to elements $d \in D$ via a UG-closed binary relation, $R: D * A$?
- Given a set, A , and a binary relation, $R: D * A$
 - Define $\gamma: A \rightarrow \text{Power}(D)$ as $\gamma(a) = \{d \in D \mid d R a\}$
 - Define partial ordering on A : $a \leq a'$ iff $\gamma(a) \leq \gamma(a')$
 - If there are distinct a and a' such that $\gamma(a) = \gamma(a')$, then merge them to force U-closure
 - Ensure that $\gamma(A)$ is a Moore family by adding greatest-lower-bound elements to A as needed.
 - This forces G-closure
 - Use the existing machinery to define the Galois connection between $\text{Power}(D)$ and A

Abstracting the Program State

- Build a binary relation, $R_d: \text{Data} \times \text{AbsData}$
 - $R_v: \text{Value} \rightarrow \text{AbsValue}$; $R_l: \text{Loc} \rightarrow \text{AbsLoc}$
 - May ignore the values of non-pointer variables.
- Build induced Galois connection, $\text{Power}(\text{Data}) \langle \alpha, \gamma \rangle \text{AbsData}$, we can
 - Build Galois connections that abstract the concrete data
 $\langle x_i : v_i \rangle R_s \langle x_i : a_i \rangle$ iff $v_i R_d a_i$
Example: $\langle x:3, y:4 \rangle R_s \langle x:\text{any}, y:\text{any} \rangle$
 - A program point is abstracted to itself: $p R_p p$,
the abstract domain of program points is $\text{ProgramPoint} \cup \{\text{top}, \text{bot}\}$ (to make it a complete lattice)
 - Finally, we can relate each concrete state to an abstract one:
 $(p, s) R_s (p', s')$ iff $p = p'$ and $s R_s s'$

Shape Graphs

- Shape analysis uses a shape graph to abstract the memory storage
 - Graph nodes denote a finite number of abstract locations:
 - $Aloc = \{N_x \mid N_x \text{ is pointed to by a set of local variables}\} \cup N_\phi$
 - N_x : the node represents all concrete Locations referred to by variables in x
 - N_ϕ : abstract summary location (all the other locations)
 - Each graph node abstracts a distinctive set of concrete Locations
 - If variables x and y may be aliased, they must share a single graph node
 - A graph edge `sel` connect nodes $n1$ and $n2$ if $n2$ is pointed to by $n1.sel$



Abstraction of Program States

□ Abstraction of memory storage

■ Abstract Stack

$\text{AbsStack} = \text{Var} \rightarrow \text{ALoc}$

Map each pointer variable into a unique abstract location (a shape graph node)

■ Abstract heap

$\text{AbsHeap} = (\text{ALoc} * \text{Sel}) \rightarrow (\text{ALoc})$

Mapping pairs of abs locations and selectors to abs locations

■ Sharing information

□ $\text{IS} : \text{ALoc} \rightarrow \{\text{yes}, \text{no}\}$

For each abstract location in the shape graph, is it shared by pointers in the heap?

□ If $\text{IS}(N_x) = \text{yes}$, then N_x must have an incoming edge from N_ϕ or have more than one incoming edges

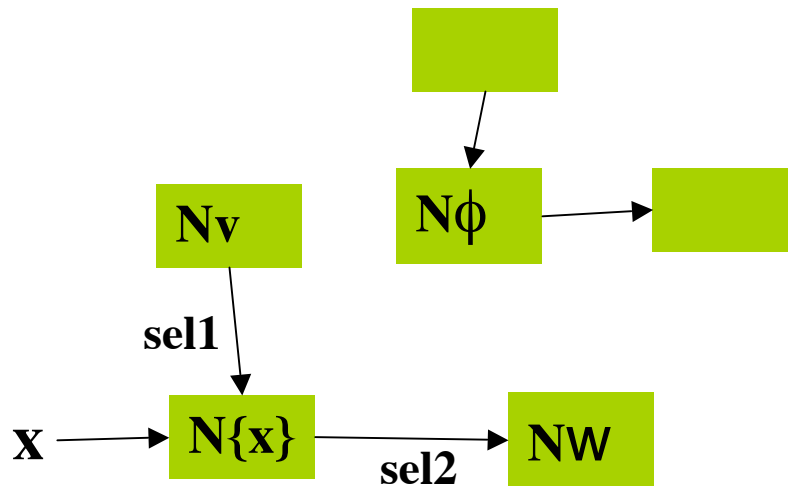
□ Transfer functions: $P(\text{AbsState}) \rightarrow P(\text{AbsState})$

■ Program state: $\text{AbsState} = \text{ProgramPoint} * \text{AbsStack} * \text{AbsHeap} * \text{IS}$

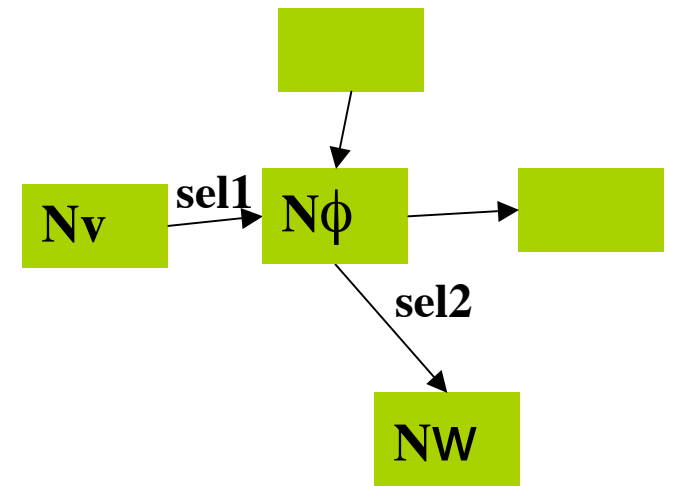
■ Each statement modifies mappings in the previous state

Transfer functions(1)

- $x = \text{nil}$
 - $F(S, H, IS) = (S', H', IS')$ where (S', H', IS') is obtained from (S, H, IS) by
 - Removing x from all mappings (killing all previous info. about x)
 - Merging all $N\phi$ nodes



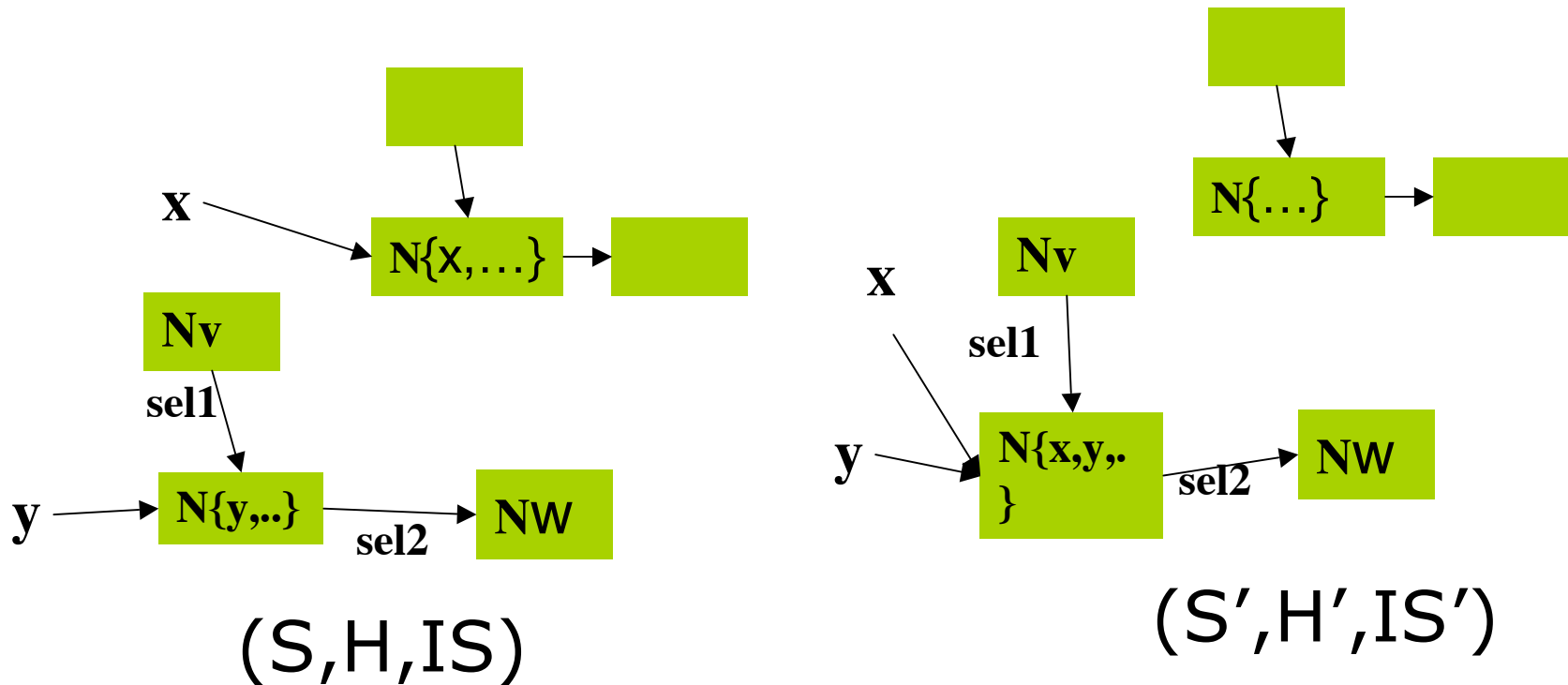
(S, H, IS)



(S', H', IS')

Transfer functions(2)

- $x = y$
 - $F(S, H, IS) = (S', H', IS')$ where
 - (S', H', IS') is obtained by modifying mappings for x to be identical to those for y



Transfer functions(3)

- $x = y.sel$
 - Remove the old binding for x
 - Establish a new binding for x to be the same as $y.sel$
 - If there is no abstract location defined for y
 - Error: dereference a null pointer
 - If there is an abstract location N_y s.t. $S[y] = N_y$, but there is no abstract location for (N_y, sel)
 - Error dereference a non-existing field
 - If there exist abstract locations N_y and N_z s.t. $S[y] = N_y$ and $H[N_y, sel] = N_z$.
 - Modify the mappings so that x points to N_z
 - If $N_z = N_\phi$, create a new node $N\{x\}$ for x --- may need to create multiple shape graphs to cover different cases
- Other transfer functions
 - E.g. $x.sel = y$; $x.sel = nil$; $allocate(x)$;