

Dataflow Analysis

17-654/17-754

Analysis of Software Artifacts

Jonathan Aldrich



Analysis of Software Artifacts -
Spring 2006

Overview: Analyses We've Seen



- AST walker analyses
 - e.g. assignment inside an if statement
 - Very approximate, very local
 - Misses case where accidental assignment is done outside an if
- Hoare logic
 - Useful for proving correctness
 - Requires a lot of work (even for ESC/Java)
 - Automated tool is unsound
 - So is manual proof, without a proof checker

Analysis of Software Artifacts -
Spring 2006

2

Motivation: Dataflow Analysis



- Catch interesting errors
 - Non-local: x is null, x is written to y, y is dereferenced
- Optimize code
 - Reduce run time, memory usage...
- Soundness required
 - Safety-critical domain
 - Assure lack of certain errors
 - Cannot optimize unless it is proven safe
 - Correctness comes before performance
- Automation required
 - Dramatically decreases cost
 - Makes cost/benefit worthwhile for far more purposes

Dataflow analysis



- Tracks value flow through program
 - Can distinguish order of operations
 - Did you read the file after you closed it?
 - Does this null value flow to that dereference?
 - Differs from AST walker
 - Walker simply collects information or checks patterns
 - Tracking flow allows more interesting properties
- Abstracts values
 - Chooses abstraction particular to property
 - Is a variable null?
 - Is a file open or closed?
 - Could a variable be 0?
 - Where did this value come from?
 - More *specialized* than Hoare logic
 - Hoare logic allows any property to be expressed
 - Specialization allows automation and soundness

Zero Analysis



- Could variable x be 0?
 - Useful to know if you have an expression y/x
 - In C, useful for null pointer analysis
- Program semantics
 - η maps every variable to an integer
- Semantic abstraction
 - σ maps every variable to non zero (NZ), zero(Z), or maybe zero (MZ)
 - Abstraction function for integers α_{z1} :
 - $\alpha_{z1}(0) = Z$
 - $\alpha_{z1}(n) = NZ$ for all $n \neq 0$
 - We may not know if a value is zero or not
 - Analysis is always an approximation
 - Need MZ option, too

Zero Analysis Example



```
x := 10;
y := x;
z := 0;
while y > -1 do
  x := x / y;
  y := y-1;
  z := 5;
```

$\sigma = []$
 $\sigma = [x \mapsto \alpha_{z1}(10)]$

Zero Analysis Example



$x := 10;$	$\sigma = []$
$y := x;$	$\sigma = [x \mapsto NZ]$
$z := 0;$	$\sigma = [x \mapsto NZ, y \mapsto \sigma(x)]$
while $y > -1$ do	
$x := x / y;$	
$y := y - 1;$	
$z := 5;$	

Zero Analysis Example



$x := 10;$	$\sigma = []$
$y := x;$	$\sigma = [x \mapsto NZ]$
$z := 0;$	$\sigma = [x \mapsto NZ, y \mapsto NZ]$
while $y > -1$ do	$\sigma = [x \mapsto NZ, y \mapsto NZ, z \mapsto \alpha_{z1}(0)]$
$x := x / y;$	
$y := y - 1;$	
$z := 5;$	

Zero Analysis Example



<code>x := 10;</code>	$\sigma = []$
<code>y := x;</code>	$\sigma = [x \mapsto \text{NZ}]$
<code>z := 0;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{NZ}]$
<code>while y > -1 do</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{NZ}, z \mapsto \text{Z}]$
<code>x := x / y;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{NZ}, z \mapsto \text{Z}]$
<code>y := y - 1;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{MZ}, z \mapsto \text{Z}]$
<code>z := 5;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{MZ}, z \mapsto \text{NZ}]$

Zero Analysis Example



<code>x := 10;</code>	$\sigma = []$
<code>y := x;</code>	$\sigma = [x \mapsto \text{NZ}]$
<code>z := 0;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{NZ}]$
<code>while y > -1 do</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{NZ}, z \mapsto \text{Z}]$
<code>x := x / y;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{MZ}, z \mapsto \text{MZ}]$
<code>y := y - 1;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{NZ}, z \mapsto \text{Z}]$
<code>z := 5;</code>	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{MZ}, z \mapsto \text{Z}]$
	$\sigma = [x \mapsto \text{NZ}, y \mapsto \text{MZ}, z \mapsto \text{NZ}]$

Zero Analysis Example



$x := 10;$	$\sigma = []$
$y := x;$	$\sigma = [x \mapsto NZ]$
$z := 0;$	$\sigma = [x \mapsto NZ, y \mapsto NZ]$
$while\ y > -1\ do$	$\sigma = [x \mapsto NZ, y \mapsto NZ, z \mapsto Z]$
$x := x / y;$	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
$y := y - 1;$	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
$z := 5;$	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto Z]$
	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto NZ]$

Zero Analysis Example



$x := 10;$	$\sigma = []$
$y := x;$	$\sigma = [x \mapsto NZ]$
$z := 0;$	$\sigma = [x \mapsto NZ, y \mapsto NZ]$
$while\ y > -1\ do$	$\sigma = [x \mapsto NZ, y \mapsto NZ, z \mapsto Z]$
$x := x / y;$	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
$y := y - 1;$	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
$z := 5;$	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$
	$\sigma = [x \mapsto NZ, y \mapsto MZ, z \mapsto MZ]$

Nothing more happens!

Zero Analysis Termination



- The analysis values will not change, no matter how many times we execute the loop
 - Proof: our analysis is deterministic
 - We run through the loop with the current analysis values, none of them change
 - Therefore, no matter how many times we run the loop, the results will remain the same
 - Therefore, we have computed the dataflow analysis results for any number of loop iterations
- Why does this work
 - If we simulate the loop, the data values could (in principle) keep changing indefinitely
 - There are an infinite number of data values possible
 - Not true for 32-bit integers, but might as well be true
 - Counting to 2^{32} is slow, even on today's processors
 - Dataflow analysis only tracks 2 possibilities!
 - So once we've explored them all, nothing more will change
 - This is the secret of abstraction
- We will make this argument more precise later

Using Zero Analysis



- Visit each division in the program
- Get the results of zero analysis for the divisor
- If the results are definitely zero, report an error
- If the results are possibly zero, report a warning

Defining Dataflow Analyses

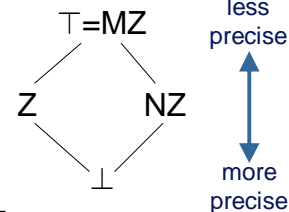


- Lattice
 - Describes program data abstractly
 - Abstract equivalent of environment
- Abstraction function
 - Maps concrete environment to lattice element
- Flow functions
 - Describes how abstract data changes
 - Abstract equivalent of expression semantics
- Control flow graph
 - Determines how abstract data propagates from statement to statement
 - Abstract equivalent of statement semantics

Lattice



- A lattice is a tuple $(L, \sqsubseteq, \sqcup, \perp, \top)$
 - L is a set of abstract elements
 - \sqsubseteq is a partial order on L
 - Means *at least as precise as*
 - \sqcup is the least upper bound of two elements
 - Must exist for every two elements in L
 - Used to merge two abstract values
 - \perp (bottom) is the least element of L
 - Means we haven't yet analyzed this yet
 - Will become clear later
 - \top (top) is the greatest element of L
 - Means we don't know anything
- L may be infinite
 - Typically should have finite height
 - All paths from \perp to \top should be finite
 - We'll see why later



Is this a lattice?



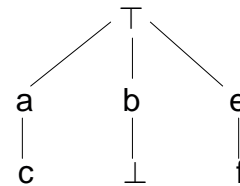
- A lattice is a tuple $(L, \Xi, \sqcup, \perp, \top)$
 - L is a set of abstract elements
 - Ξ is a partial order on L
 - \sqcup is the least upper bound of two elements
 - must exist for every two elements in L
 - \perp (bottom) is the least element of L
 - \top (top) is the greatest element of L
- Yes!



Is this a lattice?



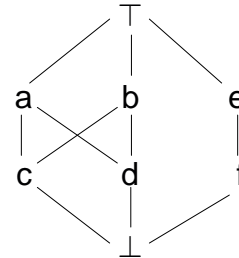
- A lattice is a tuple $(L, \Xi, \sqcup, \perp, \top)$
 - L is a set of abstract elements
 - Ξ is a partial order on L
 - \sqcup is the least upper bound of two elements
 - must exist for every two elements in L
 - \perp (bottom) is the least element of L
 - \top (top) is the greatest element of L
- No!
 - No bottom element
 - \perp is not least in the lattice order
 - It is mis-named



Is this a lattice?



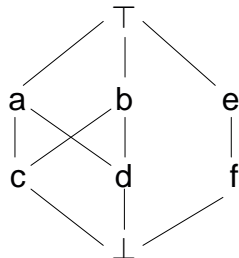
- A lattice is a tuple $(L, \sqsubseteq, \sqcup, \perp, \top)$
 - L is a set of abstract elements
 - \sqsubseteq is a partial order on L
 - \sqcup is the least upper bound of two elements
 - must exist for every two elements in L
 - \perp (bottom) is the least element of L
 - \top (top) is the greatest element of L



Definition: Least Upper Bounds



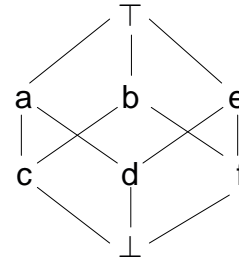
- $x \sqcup y = z$ iff
 - z is an upper bound of x and y
 - $x \sqsubseteq z$ and $y \sqsubseteq z$
 - z is the least such bound
 - $\forall w \in L$ such that $x \sqsubseteq w$ and $y \sqsubseteq w$ we have $z \sqsubseteq w$
- Also called a join
- Not a lattice
 - What is $c \sqcup d$?
 - $a, b,$ and \top are upper bounds
 - Assume \sqsubseteq is transitive
 - None is least upper bound



Is this a lattice?



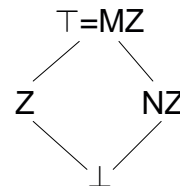
- A lattice is a tuple $(L, \sqsubseteq, \sqcup, \perp, \top)$
 - L is a set of abstract elements
 - \sqsubseteq is a partial order on L
 - \sqcup is the least upper bound of two elements
 - must exist for every two elements in L
 - \perp (bottom) is the least element of L
 - \top (top) is the greatest element of L
- Yes!



Zero Analysis Lattice



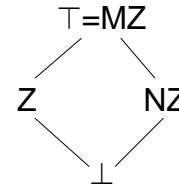
- Integer zero lattice
 - $L_{ZI} = \{ \perp, Z, NZ, MZ \}$
 - $\perp \sqsubseteq Z, \perp \sqsubseteq NZ, NZ \sqsubseteq MZ, Z \sqsubseteq MZ$
 - $\perp \sqsubseteq MZ$ holds by transitivity
 - \sqcup defined as join for \sqsubseteq
 - $x \sqcup y = z$ iff
 - z is an upper bound of x and y
 - z is the least such bound
 - Obeys laws: $\perp \sqcup \mathcal{X} = \mathcal{X}, \top \sqcup \mathcal{X} = \top, \mathcal{X} \sqcup \mathcal{X} = \mathcal{X}$
 - Also $Z \sqcup NZ = MZ$
 - $\perp = \perp$
 - $\forall \mathcal{X}. \perp \sqsubseteq \mathcal{X}$
 - $\top = MZ$
 - $\forall \mathcal{X}. \mathcal{X} \sqsubseteq \top$



Zero Analysis Lattice



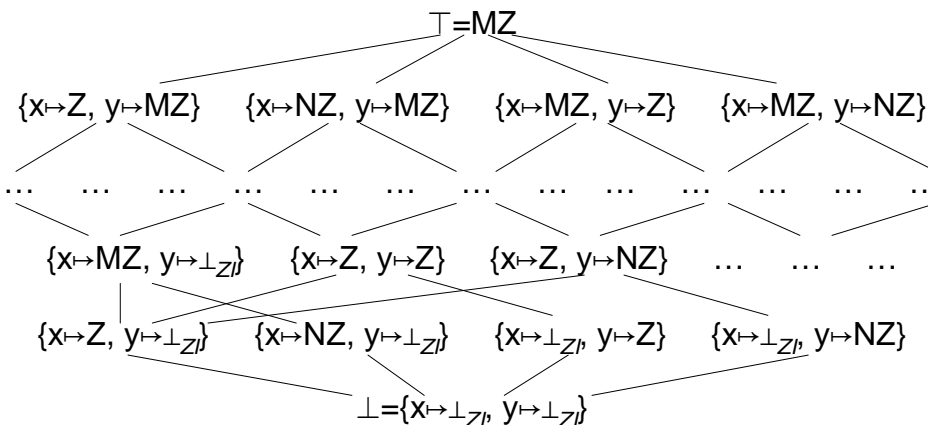
- Integer zero lattice
 - $L_{ZI} = \{ \perp, Z, NZ, MZ \}$
 - $\perp \sqsubseteq Z, \perp \sqsubseteq NZ, NZ \sqsubseteq MZ, Z \sqsubseteq MZ$
 - \sqcup defined as join for \sqsubseteq
 - $\perp = \perp$
 - $\top = MZ$
- Program lattice is a *tuple lattice*
 - L_Z is the set of all maps from **Var** to L_{ZI}
 - $\sigma_1 \sqsubseteq_Z \sigma_2$ iff $\forall x \in \mathbf{Var} . \sigma_1(x) \sqsubseteq_{ZI} \sigma_2(x)$
 - $\sigma_1 \sqcup_Z \sigma_2 = \{ x \mapsto \sigma_1(x) \sqcup_{ZI} \sigma_2(x) \mid x \in \mathbf{Var} \}$
 - $\perp = \{ x \mapsto \perp_{ZI} \mid x \in \mathbf{Var} \}$
 - $\top = \{ x \mapsto \top_{ZI} \mid x \in \mathbf{Var} \} = \{ x \mapsto MZ \mid x \in \mathbf{Var} \}$
 - Can produce a tuple lattice from *any* base lattice
 - Just define as above



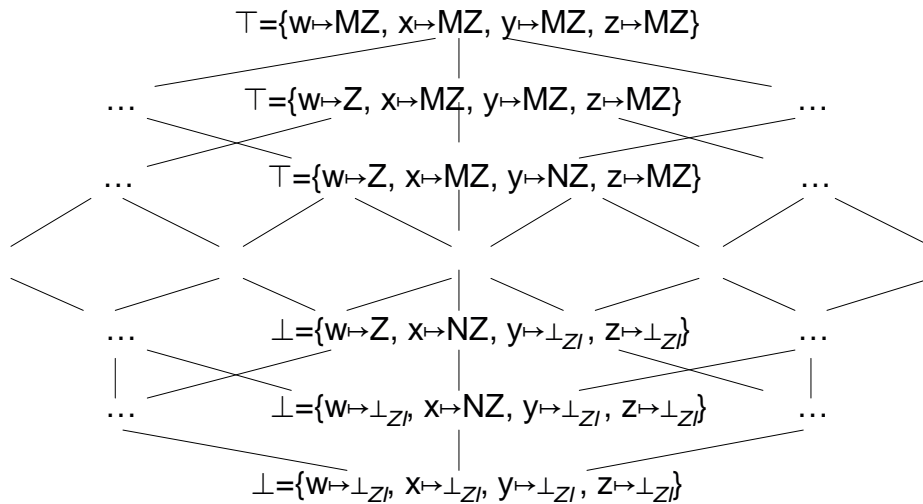
Tuple Lattices Visually



- For $\mathbf{Var} = \{ x, y \}$



One Path in a Tuple Lattice



Abstraction Function



- Maps each concrete program state to a lattice element
 - For tuple lattices, the function can be defined for values and lifted to tuples
- Integer Zero abstraction function α_{Z_I} :
 - $\alpha_{Z_I}(0) = Z$
 - $\alpha_{Z_I}(n) = NZ$ for all $n \neq 0$
- Zero Analysis abstraction function α_{Z_A} :
 - $\alpha_{Z_A}(\eta) = \{x \mapsto \alpha_{Z_I}(\eta(x)) \mid x \in \mathbf{Var}\}$
 - This is just the tuple form of $\alpha_{Z_I}(n)$
 - Can be done for any tuple lattice

Control Flow Graph (CFG)



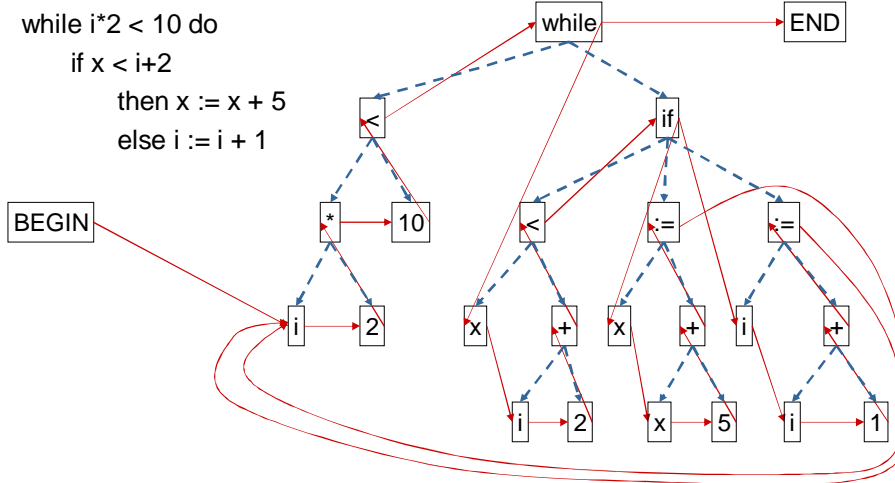
- Shows order of statement execution
 - Determines where data flows
- Decomposes expressions into primitive operations
 - Crystal: One CFG node per “useful” AST node
 - constants, variables, binary operations, assignments, if, while...
 - Loops are written out
 - Form a loop in the CFG
 - Benefit: analysis is defined one operation at a time

Intuition for Building a CFG



- Connect nodes in order of operation
 - Defined by language
- Java order of operation
 - Expressions, assignment, sequence
 - Evaluate subexpressions left to right
 - Evaluate node after children (postfix)
 - While, If
 - Evaluate condition first, then if/while
 - if branches to else and then
 - while branches to loop body and exit

Control Flow Graph Example



Flow Functions



- Compute dataflow information after a statement from dataflow information before the statement
 - Formally, map a lattice element and a CFG node to a new lattice element
- Analysis performed on 3-address code
 - inspired by 3 addresses in assembly language: add x,y,z
- Convert complex expressions to 3-address code
 - Each subexpression represented by a temporary variable
 - $x+3*y \rightarrow t_1:=3; t_2:= t_1*y; t_3:=x+t_2$

While3Addr



- copy $x = y$
- binary op $x = y \text{ op } z$ ($\text{op} \in \{+, -, *, /, \dots\}$)
- literal $x = n$
- unary op $x = \text{op } y$ ($\text{op} \in \{-, !, ++, \dots\}$)
- label $\text{label } lab$
- jump $\text{jump } lab$
- branch $\text{btrue } x \text{ } lab$

Zero Analysis Flow Functions

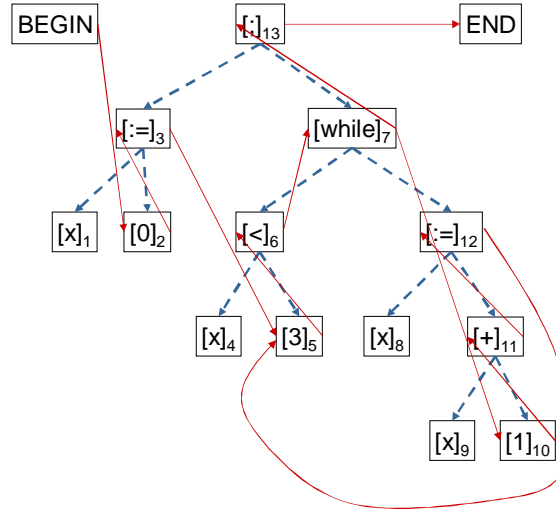


- $f_{ZA}(\sigma, [x := y]) = [x \mapsto \sigma(y)] \sigma$
- $f_{ZA}(\sigma, [x := n]) = \text{if } n == 0$
 - then $[x \mapsto Z] \sigma$
 - else $[x \mapsto NZ] \sigma$
- $f_{ZA}(\sigma, [x := \dots]) = [x \mapsto MZ] \sigma$
 - Could be more precise, e.g.
 $f_{ZA}(\sigma, [x := y + z]) =$
 - if $\sigma[y] = Z \ \&\& \ \sigma[z] = Z$
 - then $[x \mapsto Z] \sigma$ else $[x \mapsto MZ] \sigma$
- $f_{ZA}(\sigma, /* \text{any non-assignment} */) = \sigma$

Zero Analysis Example



```
x := 0;
while x > 3 do
  x := x+1
```



Zero Analysis Example

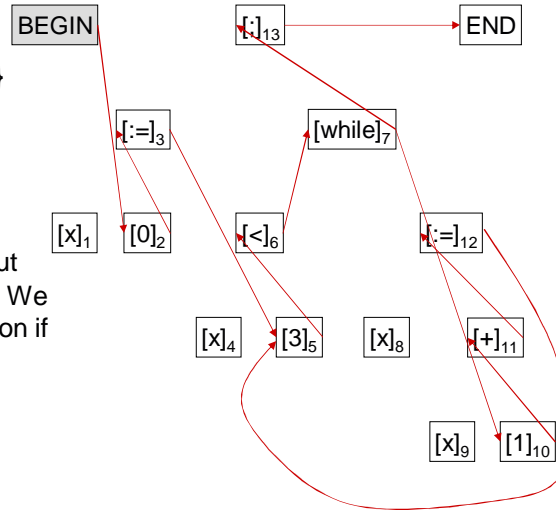


Initial dataflow

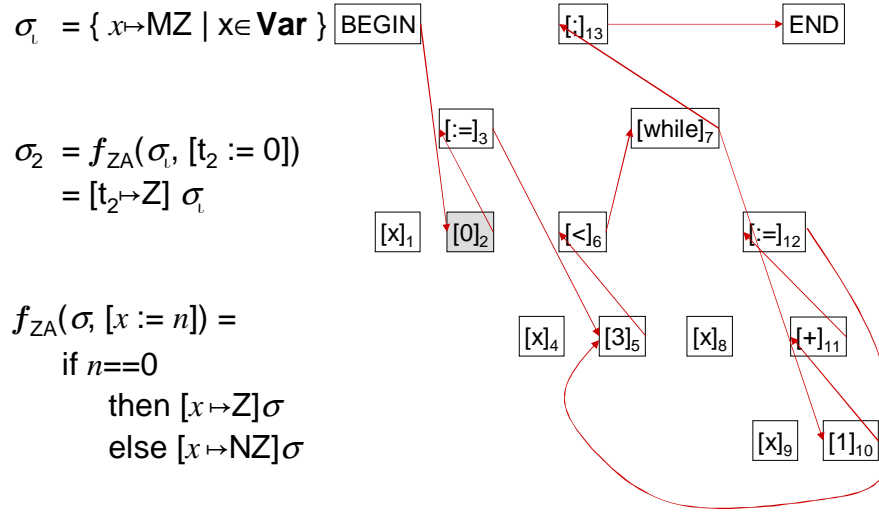
$$\sigma_i = \{ x \mapsto MZ \mid x \in \mathbf{Var} \}$$

Intuition:

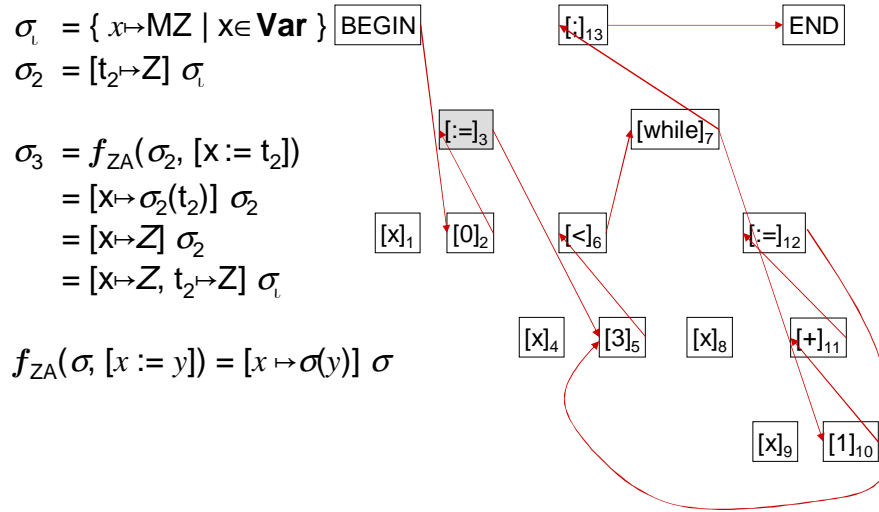
We know nothing about initial variable values. We could use a precondition if we had one.



Zero Analysis Example



Zero Analysis Example



Zero Analysis Example



$$\sigma_l = \{ x \mapsto MZ \mid x \in \mathbf{Var} \}$$

$$\sigma_3 = [x \mapsto Z, t_2 \mapsto Z] \sigma_l$$

Input to $[3]_5$ comes from

$$[:=]_3 \text{ and } [:=]_{12}$$

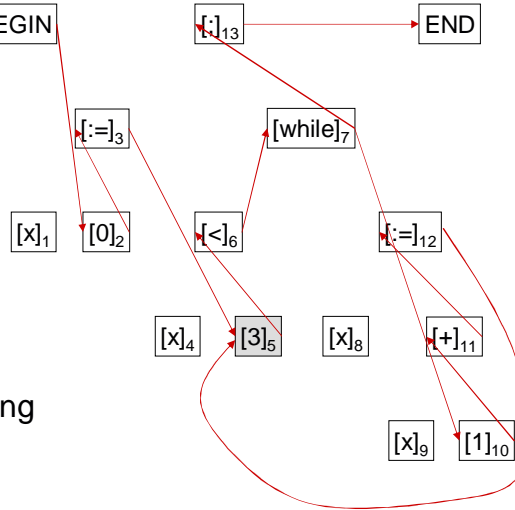
Input should be $\sigma_3 \sqcup \sigma_{12}$

What is σ_{12} ?

Solution: assume \perp

Benefit: $\sigma_3 \sqcup \perp = \sigma_3$

Same result as ignoring
back edge first time



Zero Analysis Example



$$\sigma_l = \{ x \mapsto MZ \mid x \in \mathbf{Var} \}$$

$$\sigma_3 = [x \mapsto Z, t_2 \mapsto Z] \sigma_l$$

$$\sigma_{12} = \perp$$

$$\sigma_5 = f_{ZA}(\sigma_3 \sqcup \sigma_{12}, [t_5 := 3])$$

$$= f_{ZA}(\sigma_3 \sqcup \perp, [t_5 := 3])$$

$$= f_{ZA}(\sigma_3, [t_5 := 3])$$

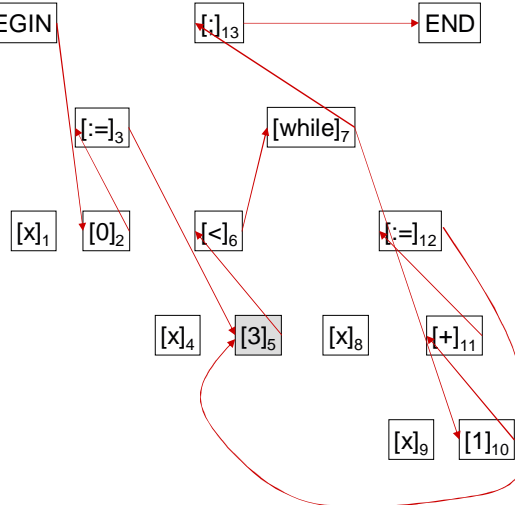
$$= [t_5 \mapsto NZ] \sigma_3$$

$$f_{ZA}(\sigma, [x := n]) =$$

if $n == 0$

then $[x \mapsto Z] \sigma$

else $[x \mapsto NZ] \sigma$



Zero Analysis Example



$$\sigma_i = \{x \mapsto MZ \mid x \in \mathbf{Var}\}$$

$$\sigma_3 = [x_1 \mapsto Z, t_2 \mapsto Z] \sigma_i$$

$$\sigma_{12} = \perp$$

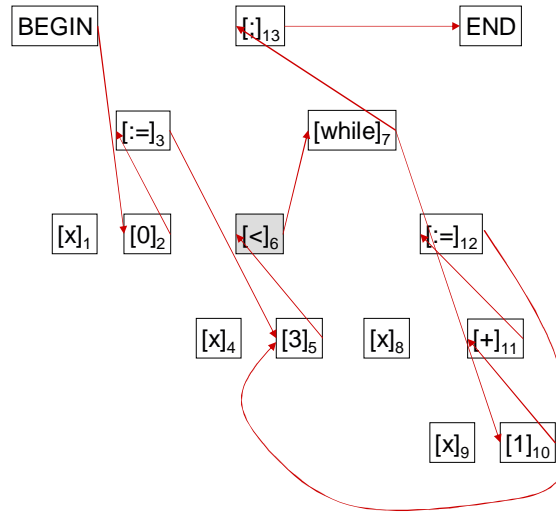
$$\sigma_5 = [t_5 \mapsto NZ] \sigma_3$$

$$\sigma_6 = f_{ZA}(\sigma_5, [t_6 := x < t_5])$$

$$= \sigma_5$$

$$= [t_5 \mapsto NZ] \sigma_3$$

$$f_{ZA}(\sigma, /* \text{any other } */) = \sigma$$



Zero Analysis Example



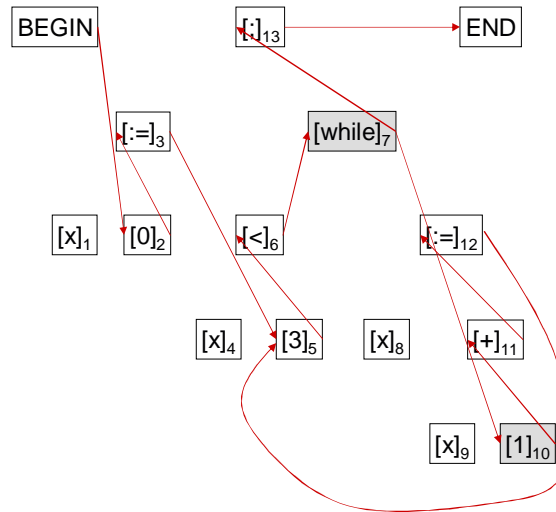
$$\sigma_i = \{x \mapsto MZ \mid x \in \mathbf{Var}\}$$

$$\sigma_3 = [x_1 \mapsto Z, t_2 \mapsto Z] \sigma_i$$

$$\sigma_{12} = \perp$$

$$\sigma_6 = [t_5 \mapsto NZ] \sigma_3$$

Skipping similar nodes...



Zero Analysis Example



$$\sigma_i = \{ x \mapsto MZ \mid x \in \mathbf{Var} \}$$

$$\sigma_3 = [x \mapsto Z, t_2 \mapsto Z] \sigma_i$$

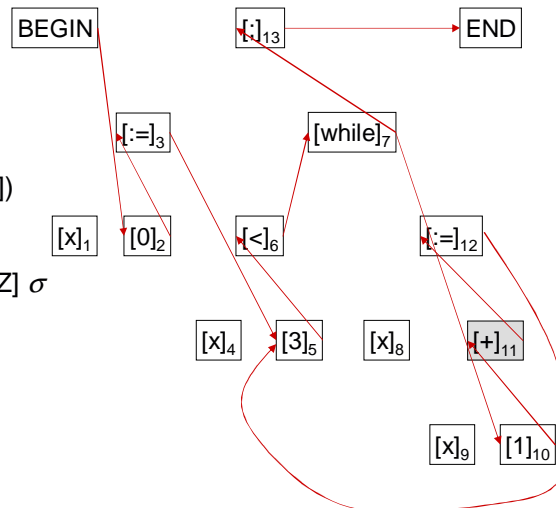
$$\sigma_{12} = \perp$$

$$\sigma_{10} = [t_{10} \mapsto NZ, \dots] \sigma_3$$

$$\sigma_{11} = f_{ZA}(\sigma_{10}, [t_{11} := x + t_{10}])$$

$$= [t_{11} \mapsto MZ] \sigma_{10}$$

$$f_{ZA}(\sigma, [x := y \text{ op } z]) = [x \mapsto MZ] \sigma$$



Zero Analysis Example



$$\sigma_i = \{ x \mapsto MZ \mid x \in \mathbf{Var} \}$$

$$\sigma_3 = [x \mapsto Z, t_2 \mapsto Z] \sigma_i$$

$$\sigma_{12} = \perp$$

$$\sigma_{11} = [t_{10} \mapsto NZ, t_{11} \mapsto MZ, \dots] \sigma_3$$

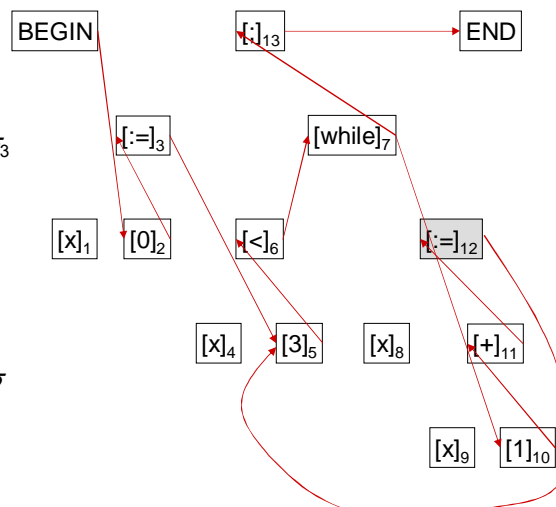
$$\sigma_{12} = f_{ZA}(\sigma_{11}, [x := t_{11}])$$

$$= [x \mapsto \sigma_{11}(t_{11})] \sigma_{11}$$

$$= [x \mapsto MZ] \sigma_{11}$$

$$= [x \mapsto MZ, \dots] \sigma_3$$

$$f_{ZA}(\sigma, [x := y]) = [x \mapsto \sigma(y)] \sigma$$



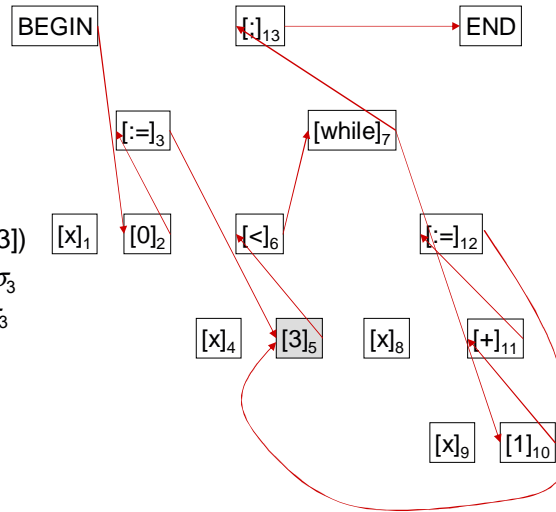
Zero Analysis Example



$$\begin{aligned} \sigma_i &= \{x \mapsto MZ \mid x \in \mathbf{Var}\} \\ \sigma_3 &= [x \mapsto Z, t_2 \mapsto Z] \sigma_i \\ \sigma_{12} &= [x \mapsto MZ, \dots] \sigma_3 \end{aligned}$$

$$\begin{aligned} \sigma_5 &= f_{ZA}(\sigma_3 \sqcup \sigma_{12}, [t_5 := 3]) \\ &= f_{ZA}([x \mapsto MZ] \sigma_3, [t_5 := 3]) \\ &= [t_5 \mapsto NZ] [x \mapsto MZ, \dots] \sigma_3 \\ &= [t_5 \mapsto NZ, x \mapsto MZ, \dots] \sigma_3 \end{aligned}$$

$$f_{ZA}(\sigma, [x_k]) = [t_k \mapsto \sigma(x)] \sigma$$

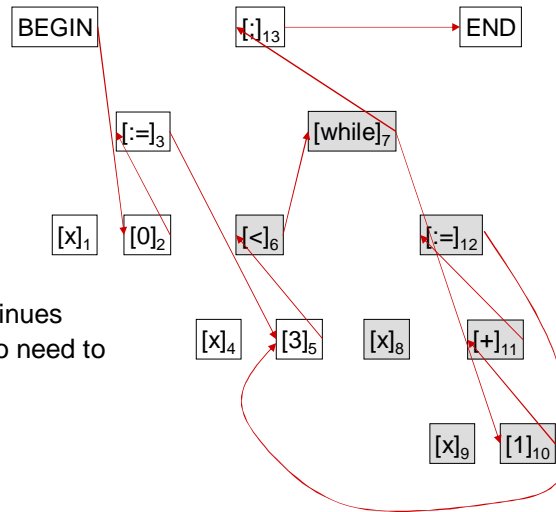


Zero Analysis Example



$$\begin{aligned} \sigma_i &= \{x \mapsto MZ \mid x \in \mathbf{Var}\} \\ \sigma_3 &= [x \mapsto Z, t_2 \mapsto Z] \sigma_i \\ \sigma_{12} &= [x \mapsto MZ, \dots] \sigma_3 \end{aligned}$$

Propagation of $x \mapsto MZ$ continues
 σ_{12} does not change, so no need to iterate again



Worklist Dataflow Analysis Algorithm



```
worklist = new Set();
for all node indexes i do
  results[i] =  $\perp_A$ ;
results[entry] =  $\nu_A$ ;
worklist.add(all nodes);
```

Ok to just add entry node if flow functions cannot return \perp_A (examples will assume this)

```
while (!worklist.isEmpty()) do
  i = worklist.pop();
  before =  $\sqcup_{k \in \text{pred}(i)} \text{results}[k]$ ;
  after =  $f_A(\text{before}, \text{node}(i))$ ;
  if (!(after  $\sqsubseteq$  results[i]))
    results[i] = after;
    for all  $k \in \text{succ}(i)$  do
      worklist.add(k);
```

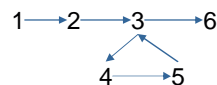
Pop removes the most recently added element from the set (performance optimization)

Example of Worklist



	Position	Worklist	a	b
[a := 0] ₁	0	1	MZ	MZ
[b := 0] ₂	1	2	Z	MZ
while [a < 2] ₃ do	2	3	Z	Z
[b := a] ₄ ;	3	4,6	Z	Z
[a := a + 1] ₅ ;	4	5,6	Z	Z
	5	3,6	MZ	Z
	3	4,6	MZ	Z
	4	5,6	MZ	MZ
[a := 0] ₆	5	3,6	MZ	MZ
	3	4,6	MZ	MZ
	4	6	MZ	MZ
	6		Z	MZ

Control Flow Graph



Worklist Algorithm Performance



- Performance
 - Visits node whenever input gets less precise
 - up to h = height of lattice
 - Propagates data along control flow edges
 - up to e = max outbound edges per node
 - Assume lattice operation cost is o
 - Overall, $O(h * e * o)$
 - Typically h, o, e bounded by n = number of statements in program
 - $O(n^3)$ for many data flow analyses
 - $O(n^2)$ if you assume a number of edges per node is small
 - Good enough to run on a function
 - Usually not run on an entire program at once, because n is too big