# Software Verification

Sebastian Nanz

# Lecture 8: Abstract Interpretation

# Abstract Interpretation

## Introduction

# One framework to rule them all

➢ In the past lectures we have introduced a particular style of program analysis: data flow analysis.

➢ For these types of analyses, and others, a main concern is correctness: how do we know that a particular analysis produces sound results (does not miss possible errors)?

➢ In the following we discuss abstract interpretation, a general framework for describing program analyses and reasoning about their correctness.

# Main ideas: Concrete computations

➢ An ordinary program describes computations in some concrete domain of values.

 ➢ **Example:** program states that record the integer value of every program variable.

$$\sigma \in State = Var \rightarrow Z$$

➢ Possible computations can be described by the concrete semantics of the programming language used.

# Main ideas: Abstract computations

➢ Abstract interpretation of a program describes computation in a different, <span style="color:red">abstract domain.</span>

    ➢ **Example:** program states that only record a specific <span style="color:red">property</span> of integers, instead of their value: their <span style="color:blue">sign</span>, whether they are <span style="color:blue">even/odd</span>, or <span style="color:blue">contained in [-32768, 32767]</span> etc.

$$\sigma \in \text{AbstractState} = \text{Var} \rightarrow \{even, odd\}$$

➢ In order to obtain abstract computations, an <span style="color:red">abstract semantics</span> for the programming language has to be defined.

➢ Abstract interpretation provides a framework for proving that the abstract semantics is sound with respect to the concrete semantics.

# The collecting semantics

We assume the state of a program to be modeled as:

$$\sigma \in \text{State} = \text{Var} \rightarrow \mathbb{Z}$$

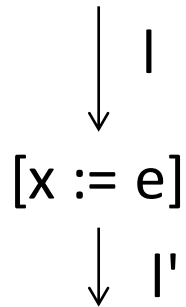We will use the following notation for function update:

$$\sigma[x \mapsto k](y) = \begin{cases} k & \text{if } x = y \\ \sigma(y) & \text{otherwise} \end{cases}$$

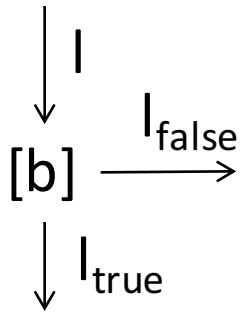We will write $\llbracket e \rrbracket \sigma$ to denote the value of an expression e in state σ.

We construct the collecting semantics as a function which gives for every program label the set of all possible states.
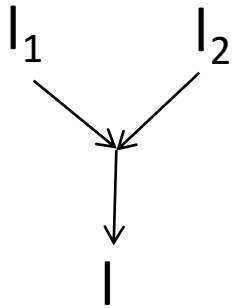
$$C : \text{Labels} \rightarrow \wp(\text{State})$$

$$C_{l'} = \{\sigma[x \mapsto n] \mid \sigma \in C_l \text{ and } [\![e]\!]\sigma = n\}$$

$$\downarrow l$$
$$[x := e]$$
$$\downarrow l'$$

$$C_{ltrue} = \{\sigma \mid \sigma \in C_l \text{ and } [\![b]\!]\sigma = true\}$$

$$C_{lfalse} = \{\sigma \mid \sigma \in C_l \text{ and } [\![b]\!]\sigma = false\}$$

$$\downarrow l$$
$$[b] \xrightarrow{\;l_{false}\;}$$
$$\downarrow l_{true}$$

$$l_1 \qquad l_2$$

$$l$$
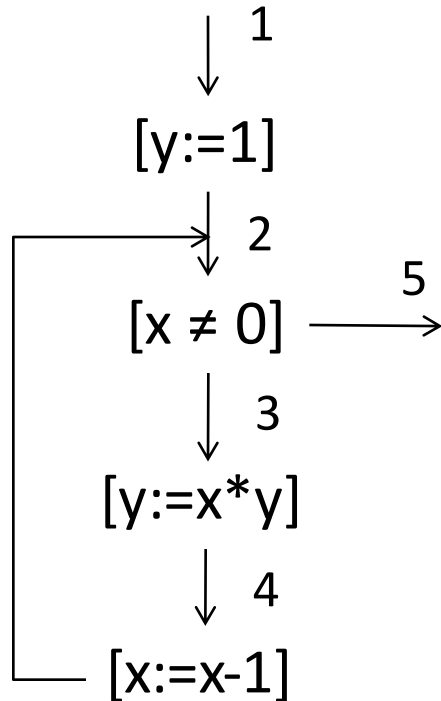
$$C_l = C_{l1} \cup C_{l2}$$

**Note:** In difference to the lecture on data flow analysis, labels are not on blocks, but on edges.

Assume x > 0.

1

[y:=1]

2

[x ≠ 0]  →  5

3

[y:=x*y]

4

[x:=x-1]

$C_1 = \{\sigma \mid \sigma(x) > 0\}$

$C_2 = \{\sigma[y \mapsto 1] \mid \sigma \in C_1\} \cup$
$\quad\quad \{\sigma[x \mapsto \sigma(x) - 1] \mid \sigma \in C_4\}$

$C_3 = C_2 \cap \{\sigma \mid \sigma(x) \neq 0\}$

$C_4 = \{\sigma[y \mapsto \sigma(x) \cdot \sigma(y)] \mid \sigma \in C_3\}$

$C_5 = C_2 \cap \{\sigma \mid \sigma(x) = 0\}$

# Solving the equations

➢ The equation system we obtain has variables $C_1$, ..., $C_5$ which are interpreted over the <span style="color:red">complete lattice</span> $\wp$(State).

➢ We can express the equation system as a <span style="color:red">monotone function</span> F : $\wp$(State)$^5$ -> $\wp$(State)$^5$
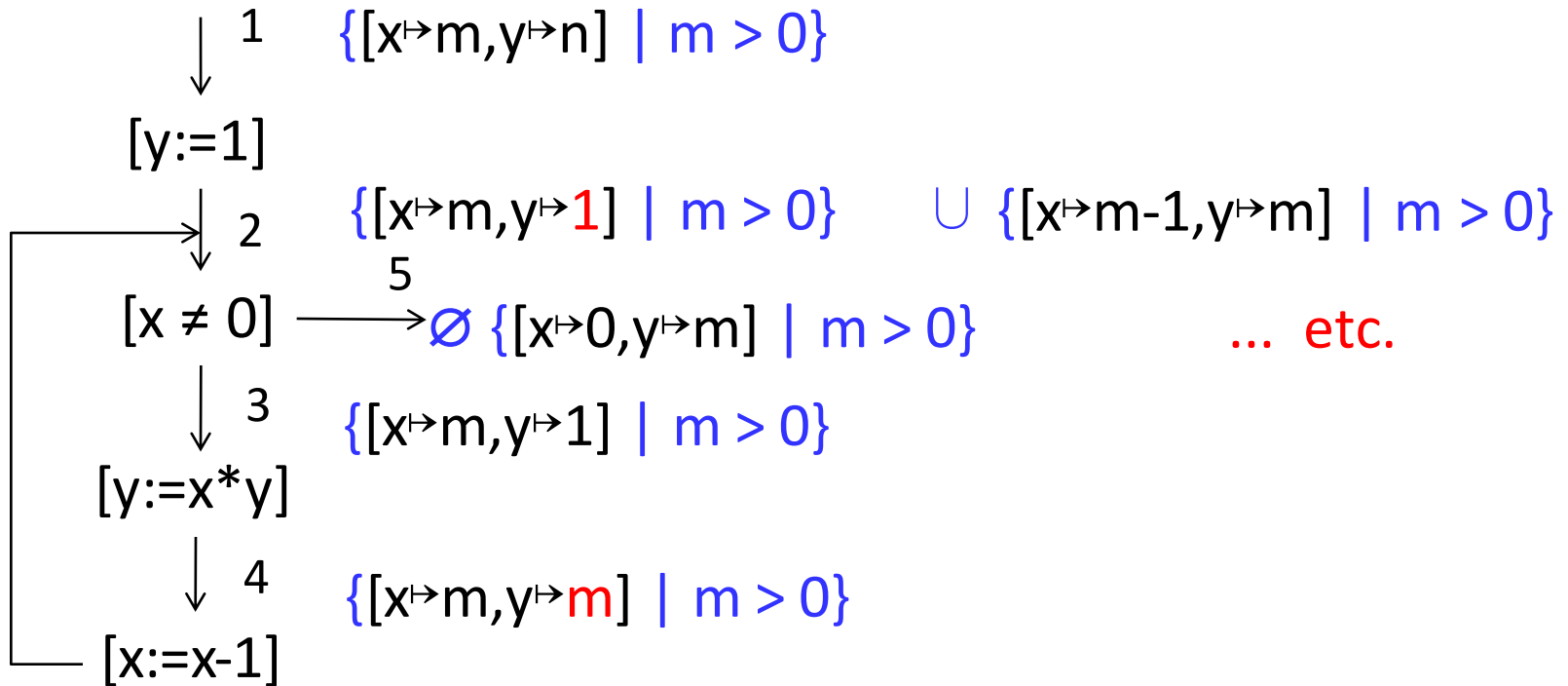
$$F(C_1, ..., C_5) = (\{\sigma \mid \sigma(x) > 0\}, ..., C_2 \cap \{\sigma \mid \sigma(x) = 0\})$$

➢ Using Tarski's Fixed Point Theorem, we know that a least fixed point exists.

➢ We have seen: The least fixed point can be computed by repeatedly applying F, starting with the bottom element $\bot$ = ($\varnothing,\varnothing,\varnothing,\varnothing,\varnothing$) of the complete lattice until stabilization.

$$F(\bot) \sqsubseteq F(F(\bot)) \sqsubseteq ... \sqsubseteq F^n(\bot) = F^{n+1}(\bot)$$

$\downarrow$ 1    $\{[x \mapsto m, y \mapsto n] \mid m > 0\}$

[y:=1]

$\downarrow$ 2    $\{[x \mapsto m, y \mapsto 1] \mid m > 0\}$    $\cup$ $\{[x \mapsto m-1, y \mapsto m] \mid m > 0\}$

[x ≠ 0]  $\xrightarrow{\quad 5 \quad}$ $\varnothing$ $\{[x \mapsto 0, y \mapsto m] \mid m > 0\}$        … etc.

$\downarrow$ 3    $\{[x \mapsto m, y \mapsto 1] \mid m > 0\}$

[y:=x*y]

$\downarrow$ 4    $\{[x \mapsto m, y \mapsto m] \mid m > 0\}$

[x:=x-1]

$C_1 = \{\sigma \mid \sigma(x) > 0\}$
$C_2 = \{\sigma[y \mapsto 1] \mid \sigma \in C_1\} \cup$
$\qquad \{\sigma[x \mapsto \sigma(x) - 1] \mid \sigma \in C_4\}$
$C_3 = C_2 \cap \{\sigma \mid \sigma(x) \neq 0\}$
$C_4 = \{\sigma[y \mapsto \sigma(x) \cdot \sigma(y)] \mid \sigma \in C_3\}$
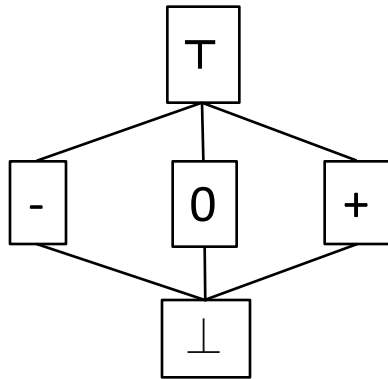$C_5 = C_2 \cap \{\sigma \mid \sigma(x) = 0\}$

# Domain for Sign Analysis

We want to focus on the sign of integers, using the domain

$$\sigma \in \text{AbstractState} = \text{Var} \rightarrow \text{Signs}$$

where Signs is the following structure:



⊤  represents all integers
+  the positive integers
-  the negative integers
0 the set {0}
⊥ the empty set

How is such a structure called?
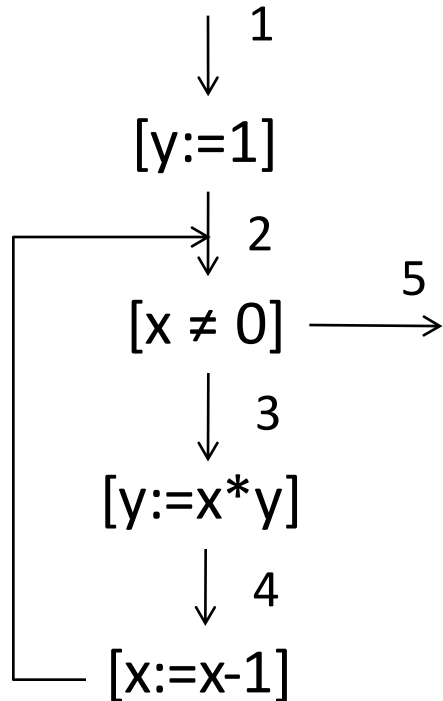
A complete lattice

Assume x > 0. Use the abstract domain for sign analysis.

1

[y:=1]

2

[x ≠ 0]  $\longrightarrow$  5

3

[y:=x*y]

4

[x:=x-1]

$A_1 = [x \mapsto +, y \mapsto \top]$

$A_2 = A_1[y \mapsto +] \sqcup$

$\quad\quad A_4[x \mapsto A_4(x) \ominus +]$

$A_3 = A_2$

$A_4 = A_3[y \mapsto A_3(x) \otimes A_3(y)]$

$A_5 = A_2 \sqcap [x \mapsto 0, y \mapsto \top]$

# Abstract Interpretation

## Foundations

# Introductory example: Expressions

A little language of expressions

Syntax

e ::= n | e * e

Concrete semantics

C[n] = n

C[e * e] = C[e] · C[e]

Example

C[-3 * 2 * -5] = C[-3 * 2] · C[-5] = C[-3 * 2] · (-5) = … = 30

# Introductory example: Abstraction

Assume that we are not interested in the value of an expression but only in its sign:

- ➢ Negative:        –
- ➢ Zero:         0
- ➢ Positive:       +

## Abstract semantics

$A[n] = sign(n)$

$A[e * e] = A[e] \otimes A[e]$

| $\otimes$ | - | 0 | + |
|-----------|---|---|---|
| - | + | 0 | - |
| 0 |   | 0 | 0 |
| + |   |   | + |

## Example

$A[-3 * 2 * -5] = A[-3 * 2] \otimes A[-5] = A[-3 * 2] \otimes (\text{-}) = \ldots =$

$= (\text{-}) \otimes (+) \otimes (\text{-}) = (+)$

# Introductory example: Soundness

➤ We want to express that the abstract semantics correctly describes the sign of a corresponding concrete computation.

➤ For this we first link each concrete value to an abstract value:

Representation function

$\beta : Z \rightarrow \{-, 0, +\}$

$$\beta(n) = \begin{cases} - & \text{if } n < 0 \\ 0 & \text{if } n = 0 \\ + & \text{if } n > 0 \end{cases}$$

# Introductory example: Soundness

➢ Conversely, we can also link abstract values to the set of concrete values they describe:

Concretization function

$$\gamma : \{-, 0, +\} \rightarrow \wp(Z)$$

$$\gamma(s) = \begin{cases} \{n \mid n < 0\} & \text{if } s = - \\ \{0\} & \text{if } s = 0 \\ \{n \mid n > 0\} & \text{if } s = + \end{cases}$$

➢ Soundness then describes intuitively that the concrete value of an expression is described by its abstract value:

$$\forall e. \ C[e] \in \gamma(A[e])$$

# Extending the language

Syntax

e ::= n | e * e | e + e | -e

Abstract semantics

A[n] = sign(n)

A[-e] = ⊖A[e]

A[e + e] = A[e] ⊕ A[e]

|     | -   | 0   | +   |
| --- | --- | --- | --- |
| ⊖   | +   | 0   | -   |

| ⊕   | -   | 0   | +   |
| --- | --- | --- | --- |
| -   | -   | -   | ?   |
| 0   |     | 0   | +   |
| +   |     |     | +   |

**Observation:** The abstract domain {-,0,+} is not closed under the interpretation of addition.

# Extending the abstract domain

We have to introduce an additional abstract value:

$\top$   "top" – (any value)

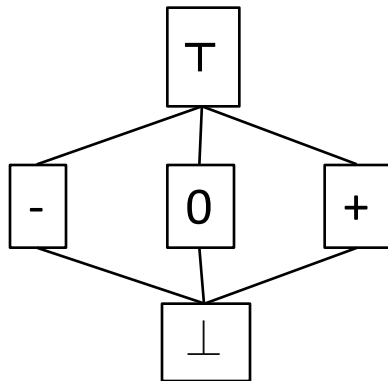| $\oplus$ | - | 0 | + | $\top$ |
|---|---|---|---|---|
| - | - | - | $\top$ | $\top$ |
| 0 | | 0 | + | $\top$ |
| + | | | + | $\top$ |
| $\top$ | | | | $\top$ |

# The new abstract domain

We can extend the concretization function to the new abstract domain $\{-,0,+,\top,\bot\}$ (add $\bot$ for completeness):

$$\gamma(\top) = \mathbf{Z} \qquad\qquad \gamma(\bot) = \varnothing$$

We obtain the following structure when drawing the partial order induced by

$$a \leq b \text{ iff } \gamma(a) \subseteq \gamma(b)$$



How is such a structure called?

A complete lattice

# Construction of complete lattices

➢ If we know some complete lattices, we can construct new ones by combining them

➢ Such constructions become important when designing new analyses with complex analysis domains

**Example:** Total function space

Let $(D_1, \sqsubseteq_1)$ be a complete lattice and let S be a set. Then $(D, \sqsubseteq)$, defined as follows, is a complete lattice:

➢ $D = S \rightarrow D_1$     ("space of total functions")

➢ $f \sqsubseteq f'$ iff $\forall\, s \in S : f(s) \sqsubseteq_1 f'(s)$   ("point-wise ordering")

# The framework of abstract interpretation

➢ Starting from a concrete domain C, define an abstract domain (A, $\sqsubseteq$), which must be a complete lattice

➢ Define a representation function β that maps a concrete value to its best abstract value

$$\beta : C \to A$$

➢ From this we can derive the concretization function γ

$$\gamma : A \to \wp(C)$$

$$\gamma(a) = \{c \in C \mid \beta(c) \sqsubseteq a\}$$

and abstraction function α for sets of concrete values

$$\alpha : \wp(C) \to A$$

$$\alpha(C) = \bigsqcup \{\beta(c) \mid c \in C\}$$

# Galois connections

➢ The following properties of **α** and **γ** hold:

Monotonicity

      (1)      α and γ are monotone functions
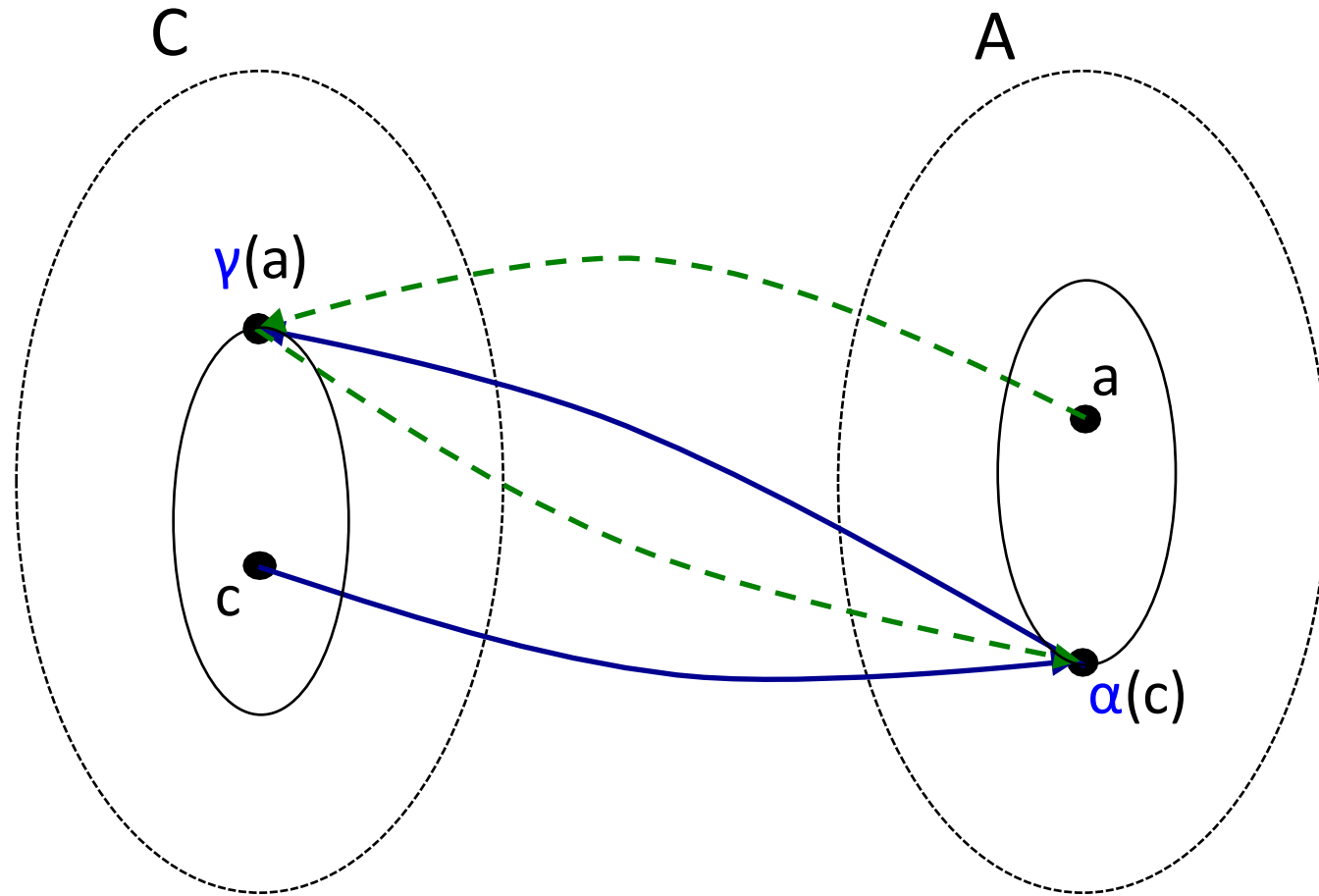
Galois connection

      (2)      $c \subseteq \gamma(\alpha(c))$         for all $c \in \wp(C)$

      (3)      $a \sqsupseteq \alpha(\gamma(a))$        for all $a \in A$

➢ Galois connection: This property means intuitively that the functions α and γ are "almost inverses" of each other.

# Figure: Galois connection

# Galois insertions

➢ For a Galois connection, there may be several elements of A that describe the same element in C

➢ As a result, A may contain elements which are irrelevant for describing C

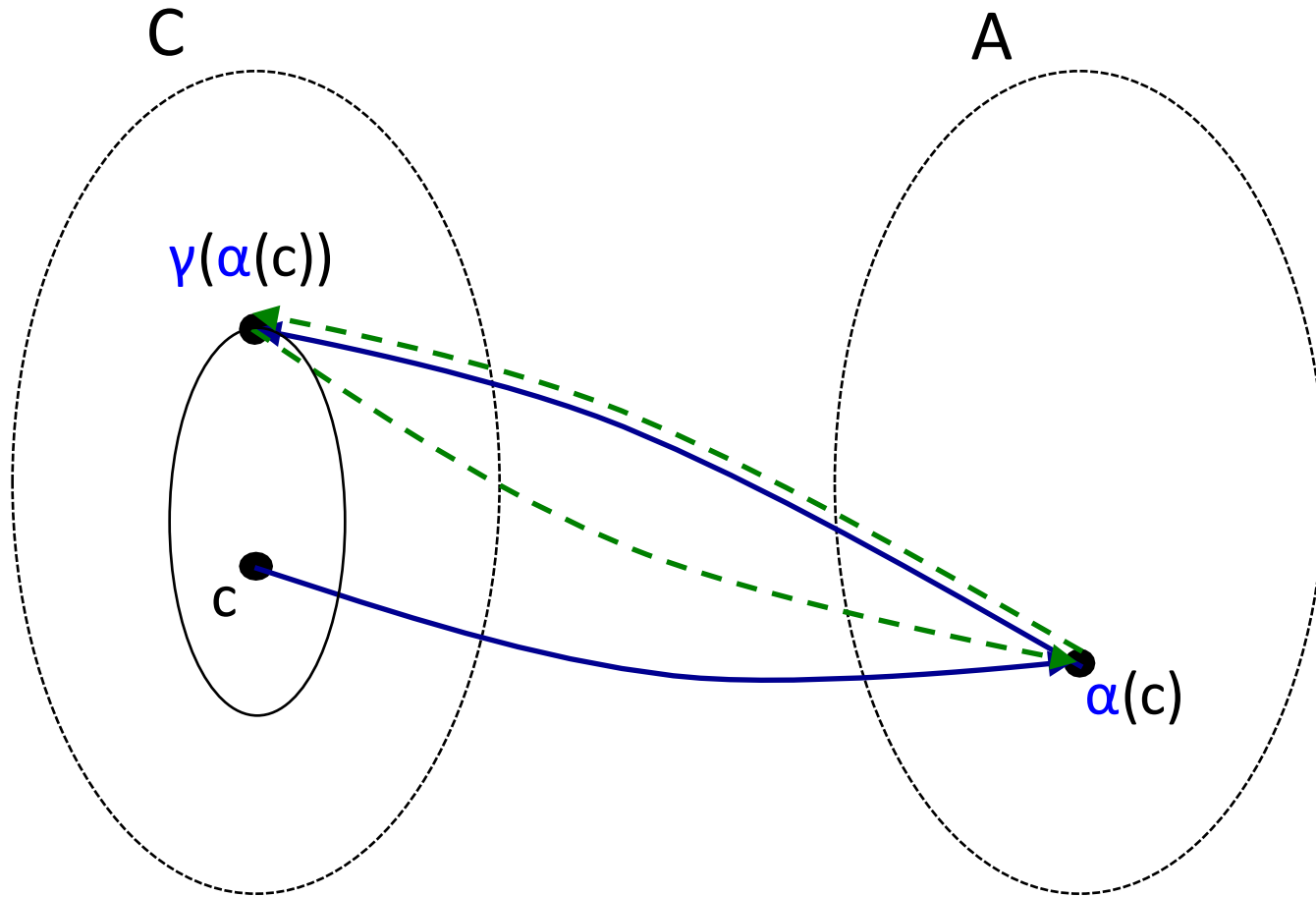➢ The concept of Galois insertion fixes this:

## Monotonicity

(1)     $\alpha$ and $\gamma$ are monotone functions

## Galois insertion

(2)     $c \subseteq \gamma(\alpha(c))$          for all $c \in \wp(C)$

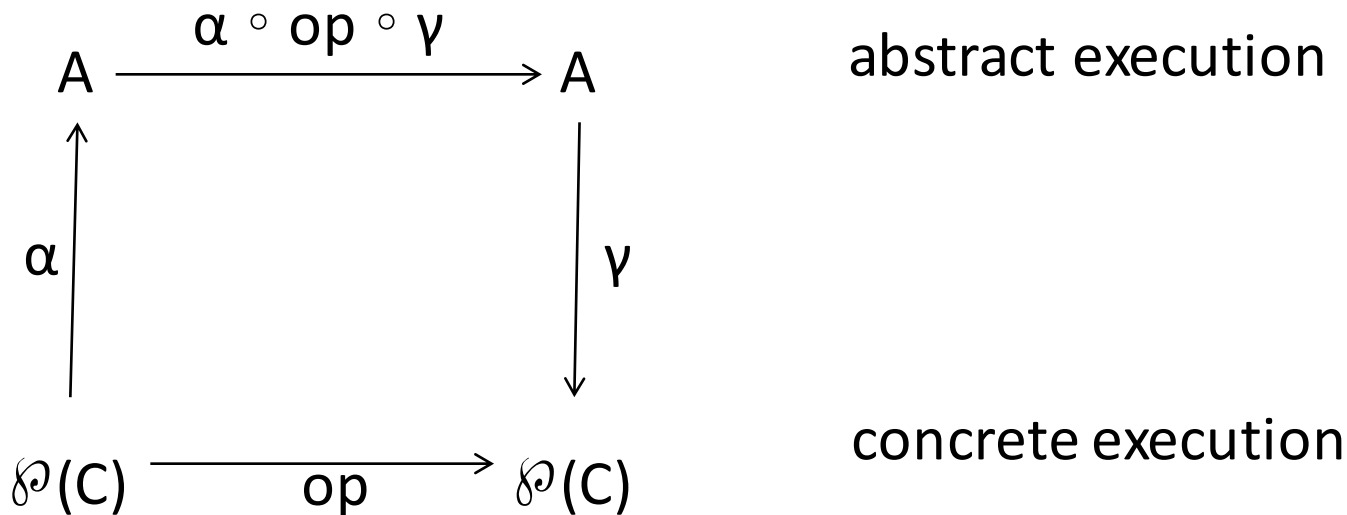(3)     $a = \alpha(\gamma(a))$          for all $a \in A$

# Induced Operations

➢ A Galois connection can be used to **induce** the abstract operations from the concrete ones.

$$A \xrightarrow{\quad \alpha \circ op \circ \gamma \quad} A \qquad \text{abstract execution}$$

$$\alpha \Big\uparrow \qquad\qquad \Big\downarrow \gamma$$

$$\wp(C) \xrightarrow{\quad op \quad} \wp(C) \qquad \text{concrete execution}$$

➢ We can show that the induced operation **op** = $\alpha \circ op \circ \gamma$ is the most precise abstract operation in this setting.

➢ The induced operation might not be computable. In this case we can define an upper approximation $op^{\#}$, **op** $\sqsubseteq op^{\#}$, and use this as abstract operation.
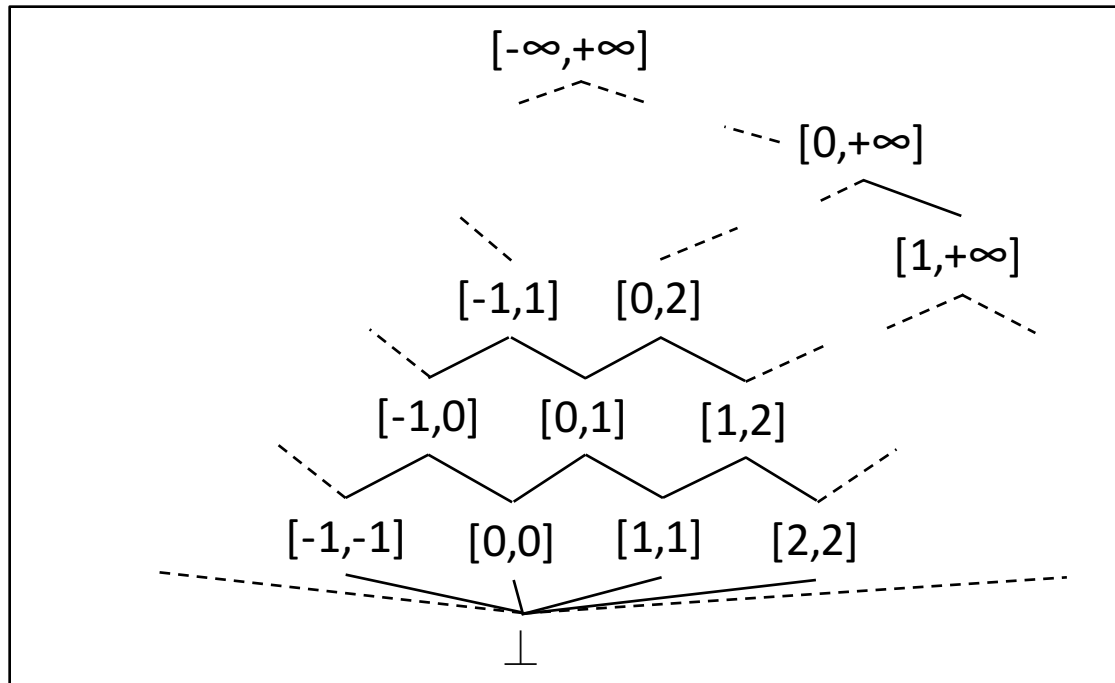
# Abstract Interpretation

## Widening

# Range analysis

➤ To introduce the notion of widening, we have a look at range analysis, which provides for every variable an over-approximation of its integer value range.

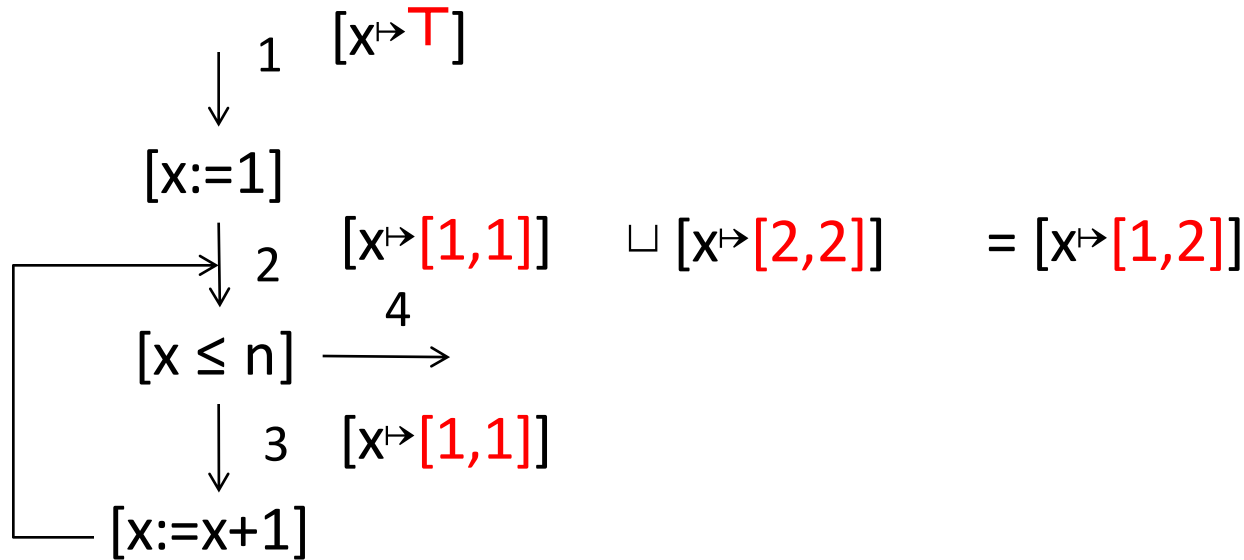➤ We are left with the task of choosing a suitable abstract domain: the interval lattice suggests itself.



Interval = { $\perp$ } $\cup$ {[x,y] | x ≤ y, x $\in$ **Z** $\cup$ {∞}, y $\in$ **Z** $\cup$ {∞}}

# Example

Consider the following program:

$$1 \quad [x \mapsto \top]$$

$$[x := 1]$$

$$2 \quad [x \mapsto [1,1]] \quad \sqcup \; [x \mapsto [2,2]] \quad = [x \mapsto [1,2]]$$

$$4$$

$$[x \leq n] \longrightarrow$$

$$3 \quad [x \mapsto [1,1]]$$

$$[x := x + 1]$$

➢ At program point 2, the following sequence of abstract states arises: $[x \mapsto [1,1]]$, $[x \mapsto [1,2]]$, $[x \mapsto [1,3]]$, …

**Consequence:** The analysis never terminates (or, if n is statically known, converges only very slowly).

# The ascending chain condition

➢ Using an arbitrary complete lattice as abstract domain, the solution is not computable in general.

➢ The reason for that is the fact that the value space might be unbounded, containing <span style="color:red">infinite ascending chains</span>:

$(l_n)_n$ is such that $l_1 \sqsubseteq l_2 \sqsubseteq l_3 \sqsubseteq \cdots$,

but there exists *no* n such that $l_n = l_{n+1} = \cdots$

➢ If we replace it with an abstract space that is finite (or does not possess infinite ascending chains), then the computation is guaranteed to terminate.

➢ In general, we want an abstract domain to satisfy the <span style="color:red">ascending chain condition</span>, i.e. each ascending chain eventually stabilises:
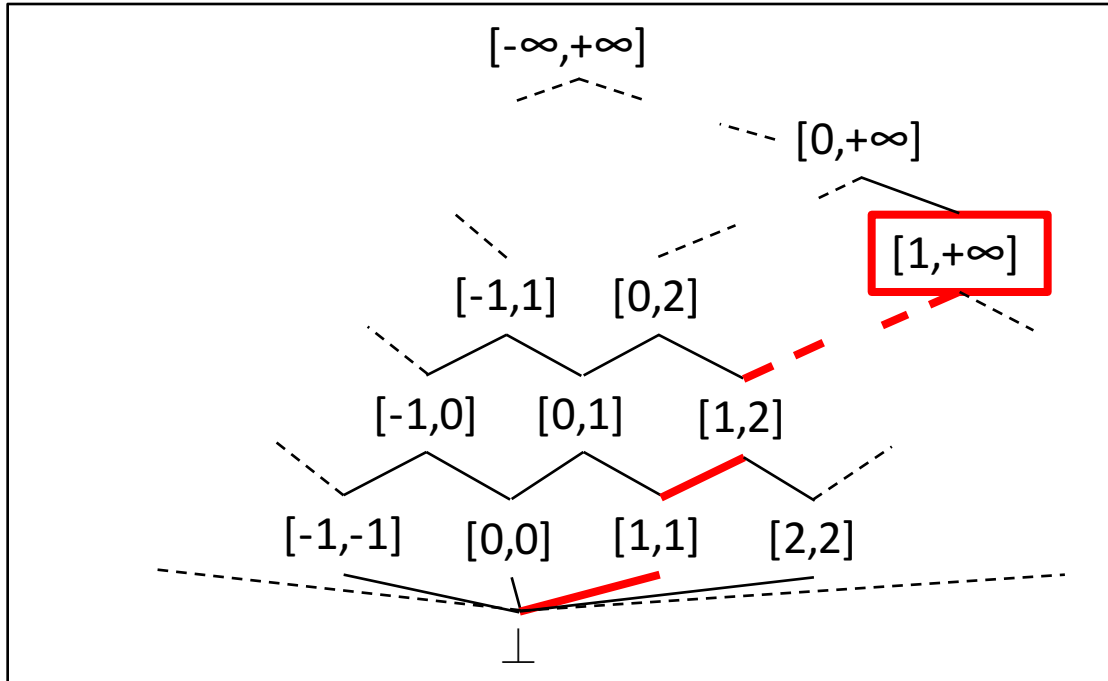
if $(l_n)_n$ is such that $l_1 \sqsubseteq l_2 \sqsubseteq l_3 \sqsubseteq \cdots$,

then there exists n such that $l_n = l_{n+1} = \cdots$

# Non-termination

➢ The reason for the non-termination in the example is that the interval lattice contains infinite ascending chains.



➢ **Trick, if we cannot eliminate ascending chains:** We redefine the join operator of the lattice to jump to the extremal value more quickly.

Before: $[1,1] \sqcup [2,2] = [1,2]$     Now: $[1,1] \triangledown [2,2] = [1,+\infty]$

# Widening

A widening $\nabla : D \times D \rightarrow D$ on a partially ordered set $(D, \sqsubseteq)$ satisfies the following properties:
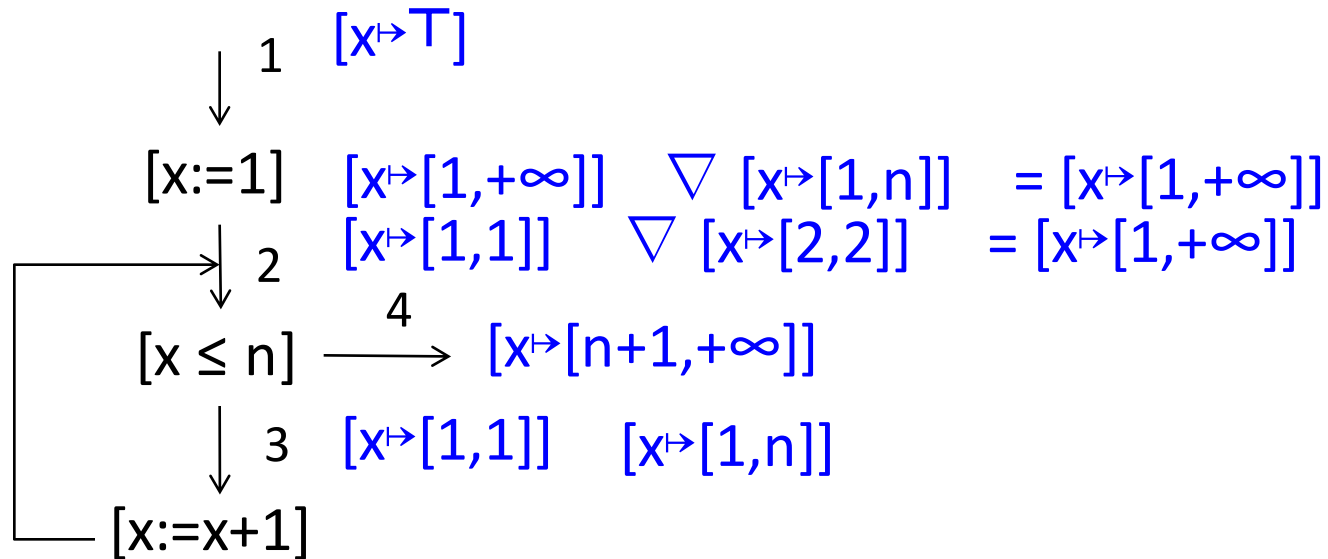
1. For all $x, y \in D$.   $x \sqsubseteq x \nabla y$   and   $y \sqsubseteq x \nabla y$

2. For all ascending chains $x_1 \sqsubseteq x_2 \sqsubseteq x_3 \sqsubseteq \cdots$ the ascending chain $y_1 = x_1 \sqsubseteq y_2 = y_1 \nabla x_2 \sqsubseteq \cdots \sqsubseteq y_{n+1} = y_n \nabla x_{n+1}$ eventually stabilizes.

➤ Widening is used to accelerate the convergence towards an upper approximation of the least fixed point.

# Example (continued)

➤ Assume we have a widening operator $\nabla$ that is defined such that $[1,1] \ \nabla \ [2,2] = [1, +\infty]$

$$[x \mapsto \top]$$

1 ↓

$[x:=1]$  $[x \mapsto [1,+\infty]]$  $\nabla$  $[x \mapsto [1,n]]$  $= [x \mapsto [1,+\infty]]$

2 ↓  $[x \mapsto [1,1]]$  $\nabla$  $[x \mapsto [2,2]]$  $= [x \mapsto [1,+\infty]]$

$[x \leq n]$  $\xrightarrow{4}$  $[x \mapsto [n+1,+\infty]]$

3 ↓  $[x \mapsto [1,1]]$   $[x \mapsto [1,n]]$

$[x:=x+1]$

➤ The analysis converges quickly.

# Reading

Patrick Cousot and Radhia Cousot. *Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints*. In: POPL'77, pages 238-252. ACM Press, 1977

Neil D. Jones, Flemming Nielson: *Abstract Interpretation: a Semantics-Based Tool for Program Analysis*, 1994

Flemming Nielson, Hanne Riis Nielson, Chris Hankin: *Principles of Program Analysis*, Springer, 2005.

Chapter 1: Section 1.5

Chapter 4 (advanced material)