# Abstract Interpretation

Wei Le

September 10, 2022

Acknowledgement: this lecture used slides from Profs Alex Aiken, David Schmidt

# Outline

- ▶ What is abstract intepretation?
- ▶ Abstract domain
- ▶ Galois connection
- ▶ Design an abstract interpretation system
- ▶ Certify neural networks: abstract interpretation for artificial intelligence ($ai^2$, ai4ai)

# History: Patrick Cousot, Radhia Cousot 1977

- ▶ Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints, 1977
- ▶ Methods and Logics for Proving Programs, 1990
- ▶ Completeness in Abstract Interpretation, 1995
- ▶ Directions for Research in Approximate System Analysis, 1999
- ▶ Probabilistic Abstract Interpretation, 2012
- ▶ An abstract interpretation framework for termination, 2012
- ▶ Abstract interpretation: past, present and future, 2014

# What is an abstract interpretation

An abstract interpretation consists of:

- An abstract domain $A$ (+,-,0) and concrete domain $D$ (Int)
- Concretization $\gamma$ and abstraction functions $\sigma$, forming a *Galois connection*
- (sound) abstract semantic function (s)

The abstract value computed is often the property we want to prove in the concrete system

# What is abstract intepretation?

▶ Define an abstract domain and perform computation on abstract domain

▶ Soundess: the conclusions from abstract intepretation are correct comparing to the conclusions reached from concrete executions

▶ A theoretical framework to formalize *approximation*

▶ A sound approximation: the conclusion proved in the abstract domain will be held in the conrete domain

▶ Abstract intepretation can lose information, meaning some conclusions that can be reached by the concrete executions but cannot be reached by abstract intepretation

# An Example

See Prof. Alex Aiken's slide

# Partitioning and Abstract Domain

▶ Partition: abstract sets of environments/concrete inputs (D)
▶ Abstract domain construction: What are the properties about D that I wish to calculate?

  ▶ Interval domain: upper and lower bound
  ▶ Congruence domain: measure density of its values

- Intervals (nonrelational):
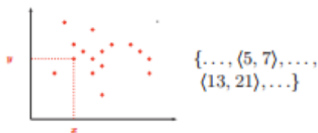  $$x \Rightarrow [a, b], \ y \Rightarrow [a', b'], \ ...$$
- Polyhedra (relational):
  $$x + y - 2z \leq 10, \ ...$$
- Difference-bound matrices (weakly relational):
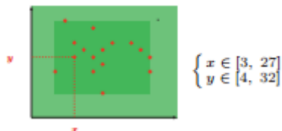  $$y - x \leq 5, \ z - y \leq 10, \ ...$$
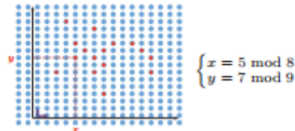
# Partitioning and Abstract Domain
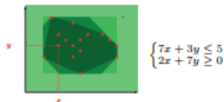


(a) [In]finite Set of Points

$$\{\ldots, \langle 5, 7 \rangle, \ldots, \langle 13, 21 \rangle, \ldots\}$$

(b) Sign Abstraction

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

(c) Interval Abstraction

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

(d) Simple Congruence Abstraction

$$\begin{cases} x = 5 \bmod 8 \\ y = 7 \bmod 9 \end{cases}$$

(b) Polyhedral Abstraction

$$\begin{cases} 7x + 3y \leq 5 \\ 2x + 7y \geq 0 \end{cases}$$

# Galois Connection: intuition



### Concretization
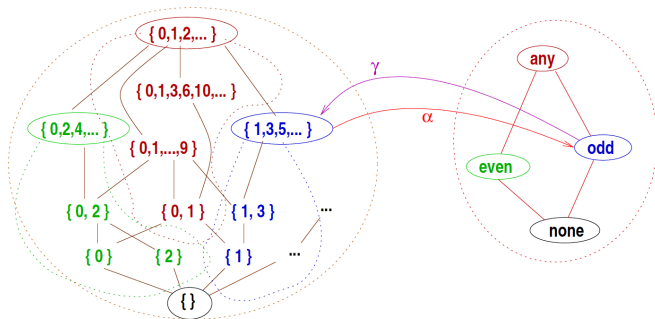$\gamma$ is the concretization function.

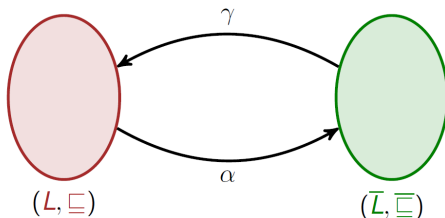$\gamma(\overline{y})$ is the concrete value in $L$ that is represented by $\overline{y}$.

### Abstraction
$\alpha$ is the abstraction function.

$\alpha(x)$ is the most precise abstract value in $\overline{L}$ whose concretization approximates $x$.

# Galois connection: example

# Galois Connection: Definition



## Definition

A Galois connection between a lattice $(L, \sqsubseteq)$ and a lattice $(\overline{L}, \overline{\sqsubseteq})$ is a pair of functions $(\alpha, \gamma)$, with $\alpha : L \to \overline{L}$ and $\gamma : \overline{L} \to L$, satisfying:

$$\alpha(x) \;\overline{\sqsubseteq}\; \overline{y} \quad \text{iff} \quad x \sqsubseteq \gamma(\overline{y}) \qquad \text{(for all } x \in L, \overline{y} \in \overline{L})$$

Notation for Galois connections: $(L, \sqsubseteq) \xleftrightarrow[\alpha]{\gamma} (\overline{L}, \overline{\sqsubseteq})$

The order is preserved.

# Designing an abstract interpretation system

**Example:** We have concrete domain, $\mathrm{Nat}$, and concrete operation, $\mathrm{succ} : \mathrm{Nat} \to \mathrm{Nat}$, defined as $\mathrm{succ}(n) = n + 1$.

We have abstract domain, $\mathrm{Parity}$, and abstract operation, $\mathrm{succ}^{\#} : \mathrm{Parity} \to \mathrm{Parity}$, defined as

$$\mathrm{succ}^{\#}(even) = odd, \quad \mathrm{succ}^{\#}(odd) = even$$
$$\mathrm{succ}^{\#}(any) = any, \quad \mathrm{succ}^{\#}(none) = none$$

$\mathrm{succ}^{\#}$ must be consistent (sound) with respect to $\mathrm{succ}$:

$$\text{if } n \; \mathcal{R}_{\mathrm{Nat}} \; a, \text{ then } \mathrm{succ}(n) \; \mathcal{R}_{\mathrm{Nat}} \; \mathrm{succ}^{\#}(a)$$

where $\mathcal{R} \subseteq \mathrm{Nat} \times \mathrm{Parity}$ relates numbers to their parities (e.g., $2 \; \mathcal{R}_{\mathrm{Nat}} \; even$, $5 \; \mathcal{R}_{\mathrm{Nat}} \; odd$, etc.).

# Abstract intepretation for robust neural networks

What does it mean to prove the robustness of a neural network?



"panda"
57.7% confidence

$+\epsilon$

$=$

"gibbon"
99.3% confidence

| Attack | Original | Perturbed | Diff |
|--------|----------|-----------|------|
| FGSM [12], $\epsilon = 0.3$ | | | |
| Brightening, $\delta = 0.085$ | | | |

Fig. 1: Attacks applied to MNIST images [25].

# Abstract intepretation for robust neural networks

Why can we use abstract intepretation?

- ▶ Deep Neural Nets: Concrete semantics —- Affine transforms + Restricted nonlinearity
- ▶ Abstract Interpretation: Scalable and precise numerical domains
- ▶ Abstract semantics should be defined on: Affine transformation (multiplictaion and addition), ReLu

# Abstract intepretation for robust neural networks
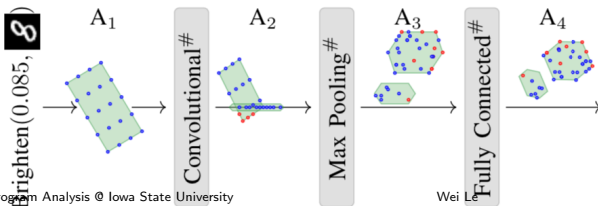
What is the abstract domain?

- ▶ Abstract domain: shapes expressible as a set of logical constraints
- ▶ Zonotope: a center-symmetric convex closed polyhedron [CAV09]

# Abstract interpretation for robust neural networks

High level ideas:

- ▶ abstract element: A1 is an abstract element (represent a group of inputs) that captured all perbuted inputs
- ▶ abstract layer: process abstract element
- ▶ abstract transformer: design abstract semantics for each concrete transformation available in the neural network
- ▶ A4 is an overappoximation of input of interest
- ▶ verify A4 will generate the same classification

\* In particular, we can capture the entire set of brightening perturbations exactly with a single zonotope. However, in general, this step may result in an abstract element that contains additional inputs (that is, red points).

# Abstract intepretation for robust neural networks

Further reading:
AI2: Safety and Robustness Certification of Neural Networks with
Abstract Interpretation