

Kandinsky - Abstract Interpretation, 1925

#### Abstract Interpretation and Applications in Security, Data Science, and Machine Learning OPLSS 2025

Caterina Urban Inria & École Normale Supérieure | Université PSL

# **Static Analysis of Safety Properties**

**OPLSS 2025** 

**Abstract Interpretation** 

#### **No Surprises, Please**





# **Trace Properties Safety Properties = "Nothing Bad Ever Happens"**

Example

• Any State Property  $S \in \mathscr{P}(\Sigma)$ :  $T \stackrel{\text{def}}{=} S^{\infty}$ 

#### **Safety Property Verification**

• T can be verified by exhaustive testing



• T can be falsified by finding a single finite execution not in T

**Abstract Interpretation** 



#### $T \in \mathscr{P}(\Sigma^{\infty})$

#### $\mathcal{I}_n(I) \subseteq T$



# **Machine Learning Safety**

**OPLSS 2025** 

Abstract Interpretation







# **Machine Learning in High-Stakes Systems**





**Abstract Interpretation** 





perform tasks that are impossible using explicit programming



#### act as surrogate model



automate decision-making





# **Machine Learning in High-Stakes Systems**

STAT+2 IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show

By <u>Casey Ross</u><sup>3</sup> @caseymross<sup>4</sup> and Ike Swetlitz

July 25, 2018

#### A self-driving Uber ran a red light last December, contrary to company claims

#### **Feds Say Self-Driving Uber SUV Did Not Recognize Jaywalking Pedestrian In Fatal Crash**

Richard Gonzales November 7, 201910:57 PM ET

**OPLSS 2025** 

**Abstract Interpretation** 

ms - The Verge

07/10/2019, 23:16







#### Caterina Urban

V



# MACHINE LEARNING SOFTWARE

Abstract Interpretation

**OPLSS 2025** 

#### METHODS TRADITIONAL SAFETY-CRITICAL SOFTWARE GOMMUNITY

JU



# **Neural Networks Feed-Forward ReLU-Activated Neural Networks**



**OPLSS 2025** 

Abstract Interpretation

#### **Rectified Linear Unit (ReLU)**







# **Neural Networks as Programs**



**OPLSS 2025** 

**Abstract Interpretation** 





#### **Neural Networks as Programs Maximal Trace Semantics**



**OPLSS 2025** 

**Abstract Interpretation** 









**OPLSS 2025** 

Abstract Interpretation

Stop





#### Max Speed 100







# Stability



**OPLSS 2025** 

Abstract Interpretation

# Stop

#### Max Speed 100











# **Prediction Stability**

**OPLSS 2025** 

Abstract Interpretation

#### Caterina Urban



# **Local Prediction Stability Prediction is Unaffected by Input Perturbations**







**Abstract Interpretation** 









# **Static Local Prediction Stability Analysis** 3-Step Recipe

**practical tools** targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior

Abstract Interpretation





#### **Local Prediction Stability Distance-Based Perturbations**

 $P_{\delta,\epsilon}(\mathbf{x}) \stackrel{\text{def}}{=} \{ \mathbf{x}' \in \mathbb{R}^{|L_0|} \mid \delta(\mathbf{x}, \mathbf{x}') \leq \epsilon \}$ Example ( $L_{\infty}$  distance):  $P_{\infty,\epsilon}(\mathbf{x}) \stackrel{\text{def}}{=} \{\mathbf{x}' \in \mathbb{R}^{|L_0|} \mid \max_i |\mathbf{x}_i - \mathbf{x}'_i| \leq \epsilon\}$ 

$$\mathscr{R}_{\mathbf{x}}^{\delta,\epsilon} \stackrel{\mathsf{def}}{=} \{ t \in \Sigma^* \mid t_0 \in P_{\delta,\epsilon}(\mathbf{x}) =$$

 $\mathscr{R}^{\delta,\epsilon}_{\mathbf{x}}$  is the set of all traces that are **prediction stable** in the neighborhood  $P_{\delta,\epsilon}(\mathbf{x})$  of a given input  $\mathbf{x}$ 



**Abstract Interpretation** 

- classification of **x**  $\Rightarrow t_{\omega} = C(\mathbf{x})$



# **Static Local Prediction Stability** Example



**Abstract Interpretation** 

**OPLSS 2025** 



### **Static Local Prediction Stability Analysis** Concrete Semantics

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior

Abstract Interpretation

Caterina Urban



# **Hierarchy of Semantics**







#### **Forward/Backward Reachable State Abstraction**

#### **Prefix/Suffix Trace Semantics**

#### **Prefix/Suffix Trace Abstraction**

#### **Partial Finite Trace Semantics**

#### **Partial Finite Trace Abstraction**

#### **Maximal Trace Semantics**





#### **Local Prediction Stability Distance-Based Perturbations**

 $P_{\delta,\epsilon}(\mathbf{x}) \stackrel{\text{def}}{=} \{ \mathbf{x}' \in \mathbb{R}^{|L_0|} \mid \delta(\mathbf{x}, \mathbf{x}') \leq \epsilon \}$ Example ( $L_{\infty}$  distance):  $P_{\infty,\epsilon}(\mathbf{x}) \stackrel{\text{def}}{=} \{\mathbf{x}' \in \mathbb{R}^{|L_0|} \mid \max_i |\mathbf{x}_i - \mathbf{x}'_i| \leq \epsilon\}$ 

$$\mathscr{R}_{\mathbf{x}}^{\delta,\epsilon} \stackrel{\mathsf{def}}{=} \{ t \in \Sigma^* \mid t_0 \in P_{\delta,\epsilon}(\mathbf{x}) =$$

 $\mathscr{R}^{\delta,\epsilon}_{\mathbf{x}}$  is the set of all traces that are **prediction stable** in the neighborhood  $P_{\delta,\epsilon}(\mathbf{x})$  of a given input  $\mathbf{x}$ 

neorem

 $M \models \mathscr{R}_{\mathbf{x}}^{\delta,\epsilon} \Leftrightarrow \mathscr{M}\llbracket M \rrbracket \subseteq \mathscr{R}_{\mathbf{x}}^{\delta,\epsilon} \Leftrightarrow \mathscr{T}_{p}(P_{\delta,\epsilon}(\mathbf{x}))$ 

**OPLSS 2025** 

Abstract Interpretation



- classification of **x**  $\Rightarrow t_{\omega} = C(\mathbf{x})$

$$\llbracket M \rrbracket \subseteq \mathscr{R}_{\mathbf{X}}^{\delta,\epsilon}$$





# **Static Local Prediction Stability Analysis** Abstract Semantics

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior

**OPLSS 2025** 

Abstract Interpretation

Caterina Urban



# **Abstract Prefix Trace Semantics**



**Abstract Interpretation** 

**OPLSS 2025** 

#### $\mathcal{T}_p(P_{\delta,\epsilon}(\mathbf{x}))\llbracket M \rrbracket \subseteq \mathcal{T}_p^{\#}(P_{\delta,\epsilon}^{\#}(\mathbf{x}))\llbracket M \rrbracket \subseteq \mathscr{R}_{\mathbf{x}}^{\delta,\epsilon} \Rightarrow M \models \mathscr{R}_{\mathbf{x}}^{\delta,\epsilon}$





# **Static Local Prediction Stability Analysis** Abstract Domain #1: Interval Domain

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior



**Abstract Interpretation** 

Caterina Urban



# **Interval Domain** Example



**OPLSS 2025** 

 $x_{i,j} \mapsto [a,b]$  $a, b \in \mathbb{R}$ 



# **Static Analysis by Abstract Interpretation**



**OPLSS 2025** 

**€ 40** + **€ 25** + **€** 5

**€ 10** +

€ 80

€ 9.95 + € 35.85 + € 24.95 + € 4.85

€ 75.60







# **Static Local Prediction Stability Analysis** Abstract Domain #2: Symbolic Domain

**practical tools** targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior



**Abstract Interpretation** 

Caterina Urban



# Symbolic Domain<sub>[Li et al. @ SAS 2019]</sub>

$$x_{i,j} \mapsto \begin{cases} \sum_{k=0}^{i-1} \mathbf{c}_k \cdot \mathbf{x}_k + \mathbf{c} & \mathbf{c}_k, \mathbf{c} \in \mathbb{R}^{|\mathbf{X}_k|} \\ [a, b] & a, b \in \mathbb{R} \end{cases}$$



**OPLSS 2025** 

**Abstract Interpretation** 



$$x_{i,j} \mapsto \sum_{k} w_{j,k}^{i-1} \cdot \mathbf{E}_{i-1,k} + b_{i,j}$$

$$x_{i,j} \mapsto \begin{cases} \mathbf{E}_{i,j} \\ [\mathbf{a},\mathbf{b}] \end{cases} \qquad 0 \le a$$

$$x_{i,j} \mapsto \begin{cases} \mathbf{A}_{i,j} \\ [\mathbf{0},\mathbf{b}] \end{cases} \qquad a < 0 \land 0 < b$$

$$x_{i,j} \mapsto \begin{cases} \mathbf{0} \\ [\mathbf{0},\mathbf{0}] \end{cases} \qquad b \le 0$$



# Symbolic Domain Example



**Abstract Interpretation** 

**OPLSS 2025** 





# Symbolic Domain Modified Example



**OPLSS 2025** 



# **Static Local Prediction Stability Analysis** Abstract Domain #3: DeepPoly Domain

**practical tools** targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior



**Abstract Interpretation** 





# DeepPoly Domain[Singh et al. @ POPL 2019]





**OPLSS 2025** 

maintain symbolic lower- and R upper-bounds for each neuron + convex ReLU approximations 





















#### **Abstract Interpretation**

**OPLSS 2025** 











#### **DeepPoly Domain Back-Substitution**

 $x_{00} \mapsto [0, 1]$  $x_{01} \mapsto [0, 1]$  $x_{30} \mapsto \begin{cases} [x_{20} - x_{21} - 14, x_{20} - x_{21} - 14] \\ [2, 8] \end{cases} \qquad x_{31} \mapsto \begin{cases} [0, 0.5 \cdot (0.5 \cdot x_{20} - 1.5 \cdot x_{21} - 8) + 0.5] \\ [0, 1] \end{cases}$  $x_{40} \mapsto \begin{cases} [0.5 \cdot x_{30} - 2 \cdot x_{31} + 1, 0.5 \cdot x_{30} - 2 \cdot x_{31} + 1] \end{cases}$ 






$x_{00} \mapsto [0, 1]$  $x_{01} \mapsto [0, 1]$  $x_{30} \mapsto \begin{cases} [x_{20} - x_{21} - 14, x_{20} - x_{21} - 14] \\ [2, 8] \end{cases} \xrightarrow{x_{31}} \mapsto \begin{cases} [0, 0.5 \cdot (0.5 \cdot x_{20} - 1.5 \cdot x_{21} - 8) + 0.5] \\ [0, 1] \end{cases}$ 







 $x_{00} \mapsto [0, 1]$  $x_{01} \mapsto [0, 1]$  $x_{30} \mapsto \begin{cases} [x_{20} - x_{21} - 14, x_{20} - x_{21}^2 - 14] \\ [2, 8] \end{cases} \qquad x_{31} \mapsto \begin{cases} [0, 0.5 \cdot x_{20} - x_{21} - 8] + 0.5] \\ [0, 1] \end{cases}$  $\mapsto \begin{cases} [x_{21} + 1, 0.5 \cdot x_{20} - 0.5 \cdot x_{21} - 6] \\ \mapsto \\ \begin{cases} [x_{10} - x_{11} + 1, 0.5 \cdot x_{10} + 2 \cdot x_{11} - 6] \end{cases}$ 







 $x_{00} \mapsto [0, 1]$  $x_{01} \mapsto [0, 1]$  $x_{10} \mapsto \begin{cases} [x_{00} + x_{01} + 4, x_{00} + x_{01} + 4] \\ [4, 6] \end{cases}$  $x_{11} \mapsto \begin{cases} [0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 3, \ 0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 3] \\ [3, 4] \end{cases}$  $x_{30} \mapsto \begin{cases} [x_{20} - x_{21} - 14, x_{20} - x_{21} - 14] \\ [2] 8] \end{cases} \qquad x_{31} \mapsto \begin{cases} [0, 0.5 \cdot (0.5 \cdot x_{20} - 1.5 \cdot x_{21} - 8) + 0.5] \\ [0] 1] \end{cases}$  $x_{40} \mapsto \begin{cases} [0.5 \cdot x_{30} - 2 \cdot x_{31} + 1, \ 0.5 \cdot x_{30} - 2 \cdot x_{31} + 1] \end{cases}$  $\mapsto \left\{ \begin{bmatrix} x_{21} + 1, \ 0.5 \cdot x_{20} - 0.5 \cdot x_{21} - 6 \end{bmatrix} \right\}$  $\mapsto \left\{ \begin{bmatrix} x_{10} - x_{11} + 1, \ 0.5 \cdot x_{10} + 2 \cdot x_{11} - 6 \end{bmatrix} \right\}$  $\mapsto \left\{ \begin{bmatrix} 0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 2, \ 1.5 \cdot x_{00} + 1.5 \cdot x_{11} + 2 \end{bmatrix} \right\}$ 

**OPLSS 2025** 

#### **Abstract Interpretation**







### **DeepPoly Domain Partial Back-Substitution**

$$\begin{aligned} x_{00} \mapsto [\mathbf{0}, \mathbf{1}] & x_{01} \mapsto [\mathbf{0}, \mathbf{1}] \\ x_{10} \mapsto \begin{cases} [x_{00} + x_{01} + 4, x_{00} + x_{01} + 4] \\ [\mathbf{4}, \mathbf{6}] & x_{11} \mapsto \begin{cases} [0.5 \cdot x_{00} + \mathbf{1}] \\ [\mathbf{3}, \mathbf{4}] \\ x_{20} \mapsto \begin{cases} [2 \cdot x_{10} + 3 \cdot x_{11}, 2 \cdot x_{10} + 3 \cdot x_{11}] \\ [\mathbf{17}, \mathbf{24}] & x_{21} \mapsto \begin{cases} [x_{10} - x_{11}, \mathbf{1}] \\ [\mathbf{1}, \mathbf{2}] \\ x_{30} \mapsto \begin{cases} [x_{20} - x_{21} - 14, x_{20} - x_{21} - 14] \\ [\mathbf{2}, \mathbf{8}] & x_{31} \mapsto \begin{cases} [0, 0.5 \cdot (0, \mathbf{1}] \\ [\mathbf{0}, \mathbf{1}] \\ x_{40} \mapsto \begin{cases} [0.5 \cdot x_{30} - 2 \cdot x_{31} + 1, 0.5 \cdot x_{30} - 2 \cdot x_{31} + 1] \\ [\mathbf{0}, \mathbf{5}] \\ \mapsto \begin{cases} [x_{21} + 1, 0.5 \cdot x_{20} - 0.5 \cdot x_{21} - 6] \\ [\mathbf{2}, \mathbf{5}, \mathbf{5}] \\ \mapsto \begin{cases} [x_{10} - x_{11} + 1, 0.5 \cdot x_{10} + 2 \cdot x_{11} - 6] \\ [\mathbf{1}, \mathbf{5}] \\ \mapsto \begin{cases} [0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 2, 1.5 \cdot x_{00} + 1.5 \cdot x_{11} + 2] \\ [\mathbf{2}, \mathbf{5}] \end{cases} \end{aligned}$$

**OPLSS 2025** 

#### Abstract Interpretation

 $(0.5 \cdot x_{01} + 3, 0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 3)$ 

 $x_{10} - x_{11}$ ]

 $0.5 \cdot x_{20} - 1.5 \cdot x_{21} - 8) + 0.5$ 







### **DeepPoly Domain** Example









 $x_{00} \mapsto [0, 1]$  $x_{01} \mapsto [0, 1]$  $x_{11} \mapsto \begin{cases} [0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 3, \ 0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 3] \\ [3, 4] \end{cases}$  $x_{10} \mapsto \begin{cases} [x_{00} + x_{01} + 4, x_{00} + x_{01} + 4] \\ [4, 6] \end{cases}$  $x_{30} \mapsto \begin{cases} [x_{20} - x_{21} - 14, x_{20} - x_{21} - 14] \\ [2, 8] \end{cases}$  $x_{31} \mapsto \begin{cases} [0, 0.5 \cdot (0.5 \cdot x_{20} - 1.5 \cdot x_{21} - 8) + 0.5] \\ [0, 1] \end{cases}$  $x_{41} \mapsto \left\{ \begin{bmatrix} x_{31}, x_{31} \end{bmatrix} \right\}$  $\mapsto \left\{ \begin{bmatrix} 0, \ 0.25 \cdot x_{20} - 0.75 \cdot x_{21} - 3.5 \end{bmatrix} \right\}$  $\mapsto \left\{ \begin{bmatrix} 0, & -0.25 \cdot x_{10} + 1.5 \cdot x_{11} - 3.5 \end{bmatrix} \right\}$ 

$$\mapsto \begin{cases} [0, \ 0.5 \cdot x_{00} + 0.5 \cdot x_{01}] \\ [0, \ \underline{1}] \end{cases}$$

**OPLSS 2025** 





### **DeepPoly Domain** Example









### **DeepPoly Domain** Maintaining Symbolic Bounds wrt the Inputs ("à la Symbolic")



**Abstract Interpretation** 

**OPLSS 2025** 





### **DeepPoly Domain** Maintaining Symbolic Bounds wrt the Inputs ("à la Symbolic")



 $x_{11} \mapsto [0, 0.5 \cdot x_{11} + 1] \rightarrow [0, 0.5 \cdot x_{00} - 0.5 \cdot x_{01} + 1]$ 

**Abstract Interpretation** 

**OPLSS 2025** 

 $x_{10} \mapsto [0, 0.5 \cdot x_{10} + 1] \rightarrow [0, 0.5 \cdot (x_{00} + x_{01}) + 1] = [0, 0.5 \cdot x_{00} + 0.5 \cdot x_{01} + 1]$ 

![](_page_44_Picture_9.jpeg)

### **DeepPoly Domain** Maintaining Symbolic Bounds wrt the Inputs ("à la Symbolic")

![](_page_45_Figure_1.jpeg)

**Abstract Interpretation** 

**OPLSS 2025** 

![](_page_45_Picture_6.jpeg)

# **DeepPoly Domain**

![](_page_46_Figure_1.jpeg)

**OPLSS 2025** 

**Abstract Interpretation** 

 $x_{21} \mapsto [0, 0.5 \cdot x_{21} + 1] \rightarrow [0, 0.25 \cdot x_{00} - 0.25 \cdot x_{01} + 1.5]$ 

![](_page_46_Figure_7.jpeg)

![](_page_46_Picture_8.jpeg)

# **DeepPoly Domain**

![](_page_47_Figure_1.jpeg)

**OPLSS 2025** 

**Abstract Interpretation** 

### **Static Local Prediction Stability Analysis** Going Farther: Multi-Neuron Abstractions

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior

![](_page_48_Picture_4.jpeg)

**Abstract Interpretation** 

Caterina Urban

![](_page_48_Picture_7.jpeg)

# Multi-Neuron Abstractions[Singh et al. @ NeurIPS 2019]

![](_page_49_Figure_1.jpeg)

![](_page_49_Picture_8.jpeg)

![](_page_49_Picture_9.jpeg)

![](_page_49_Picture_10.jpeg)

![](_page_49_Picture_11.jpeg)

![](_page_49_Picture_12.jpeg)

![](_page_49_Picture_13.jpeg)

![](_page_49_Picture_14.jpeg)

### Static Local Prediction Stability Analysis An Exercise Worth Trying 😳

**practical tools** targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide progran

concrete semantics mathematical models of the program behavior

**OPLSS 2025** 

**Abstract Interpretation** 

![](_page_50_Picture_6.jpeg)

![](_page_50_Picture_8.jpeg)

# **Saliency Map Stability**

**OPLSS 2025** 

Abstract Interpretation

![](_page_51_Picture_3.jpeg)

![](_page_51_Picture_6.jpeg)

## **Local Prediction Stability Not Enough!**

![](_page_52_Figure_1.jpeg)

![](_page_52_Picture_4.jpeg)

![](_page_52_Picture_7.jpeg)

![](_page_52_Picture_8.jpeg)

# Saliency Maps<sub>[Simonyan & al. @ ICLR 2014]</sub>

![](_page_53_Figure_1.jpeg)

![](_page_53_Picture_2.jpeg)

Abstract Interpretation

Dog

#### Caterina Urban

![](_page_53_Picture_9.jpeg)

# **Local Saliency Map Stability**

#### Input Image

**Saliency Map** 

**Expected Saliency Map** 

#### **Distance**

![](_page_54_Picture_5.jpeg)

![](_page_54_Picture_6.jpeg)

![](_page_54_Picture_7.jpeg)

![](_page_54_Picture_8.jpeg)

![](_page_54_Picture_9.jpeg)

![](_page_54_Picture_10.jpeg)

![](_page_54_Picture_11.jpeg)

![](_page_54_Picture_12.jpeg)

![](_page_54_Picture_13.jpeg)

![](_page_54_Picture_14.jpeg)

![](_page_54_Picture_15.jpeg)

![](_page_54_Picture_16.jpeg)

#### **Saliency Map Stability**

**Abstract Interpretation** 

![](_page_54_Picture_19.jpeg)

![](_page_54_Picture_22.jpeg)

![](_page_54_Picture_23.jpeg)

# **Reading Suggestion**

#### Verifying Attention Robustness of Deep Neural **Networks against Semantic Perturbations**

Satoshi Munakata<sup>1</sup>, Caterina Urban<sup>2</sup>, Haruki Yokoyama<sup>1</sup>, Koji Yamamoto<sup>1</sup>, and Kazuki Munakata<sup>1</sup>

> <sup>1</sup> Fujitsu, Kanagawa, Japan <sup>2</sup> Inria & ENS | PSL & CNRS, Paris, France

Abstract. It is known that deep neural networks (DNNs) classify an input image by paying particular attention to certain specific pixels; a graphical representation of the magnitude of attention to each pixel is called a *saliency-map*. Saliency-maps are used to check the validity of the classification decision basis, e.g., it is not a valid basis for classification if a DNN pays more attention to the background rather than the subject of an image. Semantic perturbations can significantly change the saliencymap. In this work, we propose the first verification method for *attention* robustness, i.e., the local robustness of the changes in the saliency-map against combinations of semantic perturbations. Specifically, our method determines the range of the perturbation parameters (e.g., the brightness change) that maintains the difference between the actual saliency-map change and the expected saliency-map change below a given threshold value. Our method is based on activation region traversals, focusing on the outermost robust boundary for scalability on larger DNNs. We empirically evaluate the effectiveness and performance of our method on DNNs trained on popular image classification datasets.

**Abstract Interpretation** 

![](_page_55_Picture_6.jpeg)

#### Caterina Urban

![](_page_55_Picture_9.jpeg)

![](_page_56_Picture_0.jpeg)

![](_page_56_Picture_1.jpeg)

**OPLSS 2025** 

Abstract Interpretation

Stop

![](_page_56_Picture_5.jpeg)

![](_page_56_Picture_6.jpeg)

#### Max Speed 100

![](_page_56_Picture_8.jpeg)

![](_page_56_Figure_9.jpeg)

![](_page_56_Picture_10.jpeg)

![](_page_56_Picture_12.jpeg)

## ACAS Xu **Airborne Collision Avoidance System for Unmanned Aircraft**

implemented using 45 feed-forward fully-connected ReLU networks

![](_page_57_Picture_2.jpeg)

![](_page_57_Figure_3.jpeg)

### 5 input sensor measurements

- $\rho$ : distance from ownship to intruder
- $\theta$ : angle to intruder relative to ownship heading direction
- $\psi$ : heading angle to intruder relative to ownship heading direction
- $v_{own}$ : speed of ownship
- $v_{int}$ : speed of intruder
- 5 output horizontal advisories
- Strong Left
- Weak Left
- **Clear of Conflict**
- Weak Right
- Strong Right

Caterina Urban

![](_page_57_Picture_22.jpeg)

### **ACAS Xu Safety Properties** Example: "if intruder is near and approaching from the left, go Strong Right"

![](_page_58_Figure_1.jpeg)

**OPLSS 2025** 

![](_page_58_Figure_5.jpeg)

#### Caterina Urban

• • •

![](_page_58_Picture_8.jpeg)

![](_page_58_Picture_14.jpeg)

### **Static Safety Analysis 3-Step Recipe**

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior

**Abstract Interpretation** 

![](_page_59_Picture_6.jpeg)

![](_page_59_Picture_9.jpeg)

![](_page_59_Picture_10.jpeg)

### Safety **Input-Output Properties**

I: input specification

O: output specification

### $\mathscr{R}^{\delta,\epsilon}_{\mathbf{x}} \stackrel{\mathsf{def}}{=} \{ t \in \Sigma^* \mid t_0 \models \mathbf{I} \Rightarrow t_\omega \models \mathbf{O} \}$

 $\mathcal{S}_{\mathbf{0}}^{\mathbf{I}}$  is the set of all traces that are **satisfy** the input and output specifications  $\mathbf{I}$  and  $\mathbf{O}$ 

![](_page_60_Picture_5.jpeg)

**Abstract Interpretation** 

![](_page_60_Picture_8.jpeg)

![](_page_60_Picture_10.jpeg)

![](_page_61_Picture_0.jpeg)

![](_page_61_Figure_1.jpeg)

**OPLSS 2025** 

Abstract Interpretation

#### **Clear of Conflict**

![](_page_61_Picture_7.jpeg)

### **Static Safety Analysis Concrete Semantics**

concrete semantics mathematical models of the program behavior

**Abstract Interpretation** 

![](_page_62_Picture_6.jpeg)

![](_page_62_Picture_8.jpeg)

![](_page_62_Picture_9.jpeg)

#### Caterina Urban

![](_page_62_Picture_11.jpeg)

# **Hierarchy of Semantics**

![](_page_63_Picture_1.jpeg)

![](_page_63_Picture_2.jpeg)

![](_page_63_Picture_4.jpeg)

### **Forward/Backward Reachable State Abstraction**

#### **Prefix/Suffix Trace Semantics**

#### **Prefix/Suffix Trace Abstraction**

#### **Partial Finite Trace Semantics**

#### **Partial Finite Trace Abstraction**

### **Maximal Trace Semantics**

![](_page_63_Picture_13.jpeg)

![](_page_63_Picture_15.jpeg)

### Safety **Input-Output Properties**

I: input specification

O: output specification

$$\mathscr{R}_{\mathbf{x}}^{\delta,\epsilon} \stackrel{\mathsf{def}}{=} \{ t \in \Sigma^* \mid t_0 \models \mathbf{I} \Rightarrow$$

 $\mathcal{S}_{\mathbf{0}}^{\mathbf{I}}$  is the set of all traces that **satisfy** the input and output specifications I and O

#### Iheorem

 $M \models \mathcal{S}_{\mathbf{0}}^{\mathbf{I}} \Leftrightarrow \mathscr{M}[\![M]\!] \subseteq \mathcal{S}_{\mathbf{0}}^{\mathbf{I}} \Leftrightarrow \mathscr{T}_{p}(\mathbf{I})[\![M]\!] \subseteq \mathcal{S}_{\mathbf{0}}^{\mathbf{I}}$ 

**OPLSS 2025** 

**Abstract Interpretation** 

### $t_{\omega} \models \mathbf{O}$

![](_page_64_Picture_10.jpeg)

![](_page_64_Picture_12.jpeg)

![](_page_64_Picture_13.jpeg)

### **Static Safety Analysis Abstract Semantics**

abstract semantics, abstract domains algorithmic approaches to decide program properties

**OPLSS 2025** 

**Abstract Interpretation** 

![](_page_65_Picture_6.jpeg)

![](_page_65_Picture_9.jpeg)

![](_page_65_Picture_11.jpeg)

### **Static Safety Analysis Abstract Prefix Trace Semantics**

![](_page_66_Figure_1.jpeg)

**Abstract Interpretation** 

**OPLSS 2025** 

 $\mathcal{T}_p(\mathbf{I})[[M]] \subseteq \mathcal{T}_p^{\#}(\mathbf{I}^{\#})[[M]] \subseteq \mathcal{S}_{\mathbf{O}}^{\mathbf{I}} \Rightarrow M \models \mathcal{S}_{\mathbf{O}}^{\mathbf{I}}$ 

#### Caterina Urban

![](_page_66_Picture_7.jpeg)

### **Static Safety Analysis** Abstract Domain #3: DeepPoly Domain

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior

![](_page_67_Picture_4.jpeg)

Abstract Interpretation

![](_page_67_Picture_7.jpeg)

![](_page_67_Picture_9.jpeg)

![](_page_68_Figure_0.jpeg)

**Abstract Interpretation** 

**OPLSS 2025** 

#### Caterina Urban

![](_page_68_Picture_5.jpeg)

### **DeepPoly Domain** Example

![](_page_69_Figure_1.jpeg)

**Abstract Interpretation** 

**OPLSS 2025** 

![](_page_69_Picture_6.jpeg)

### DeepPoly Domain Example

![](_page_70_Figure_1.jpeg)

**OPLSS 2025** 

![](_page_70_Picture_5.jpeg)

![](_page_70_Picture_7.jpeg)

### **Static Safety Analysis** Abstract Domain #2: Symbolic Domain

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior

![](_page_71_Picture_4.jpeg)

Abstract Interpretation

![](_page_71_Picture_7.jpeg)

![](_page_71_Picture_9.jpeg)
# Symbolic Domain Example



**Abstract Interpretation** 

**OPLSS 2025** 

#### Caterina Urban



# Symbolic Domain Example



**OPLSS 2025** 

#### Caterina Urban





# Symbolic Domain Example



**OPLSS 2025** 

#### Caterina Urban



# **Static Safety Analysis** Abstract Domain #4: Reduced Product Domain

**practical tools** targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior



Abstract Interpretation

Caterina Urban



### **Reduced Product Domain** Symbolic Domain & DeepPoly Domain

#### Symbolic $[\max(a_s, a_d), \min(b_s, b_u)]$

**OPLSS 2025** 

**Abstract Interpretation** 



Caterina Urban

# **Reduced Product Domain** Example

 $0 \le \rho \le 1$ x00  $x_{00} \mapsto \begin{cases} x_{00} \\ [x_{00}, x_{00}] \\ [0, 1] \end{cases}$ x10 0  $x_{01} \mapsto \begin{cases} x_{01} \\ [x_{01}, x_{01}] \\ [-1, 1] \end{cases}$ 0 x1<sup>-</sup> -1  $-1 \le \theta \le 1$ x01

#### **Abstract Interpretation**

**OPLSS 2025** 





## **Reduced Product Domain** Example



#### **Abstract Interpretation**

**OPLSS 2025** 



# Example



**OPLSS 2025** 

**Abstract Interpretation** 



# **Static Safety Analysis** Going Farther: Complete Methods

practical tools targeting specific programs

abstract semantics, abstract domains algorithmic approaches to decide program properties

concrete semantics mathematical models of the program behavior



Abstract Interpretation

Caterina Urban





## Reuva [Wang et al. @ USENIX Security 2018] **Asymptotically Complete Method**





Abstract Interpretation



# symbolic propagation + iterative input refinement







# **DeepPoly Domain + Input Refinement**



#### Caterina Urban





## Neurify **Asymptotically Complete Method**





**OPLSS 2025** 

symbolic propagation + convex ReLU approximation + iterative input/ReLU refinement U





#### $\alpha\beta$ -CROWN The State of the Art











GPU optimized relaxation (*α*-CROWN)

Parallel branch and bound ( $\beta$ -CROWN)

#### Winner of the International Verification of Neural Networks Competition since 2021







### **Static Safety Analysis Going Farther: Other Abstractions**

abstract semantics, abstract domains algorithmic approaches to decide program properties

**OPLSS 2025** 

**Abstract Interpretation** 

Caterina Urban



# Interval Neural Networks [Prabhakar and Afza @ NeurIPS 2019]



 $l_j \leq x_{0,j} \leq u_j$ 

**OPLSS 2025** 

Abstract Interpretation

Elboher et al. @ CAV 2020

merge neurons layer-wise based on partitioning strategy + replace weights with intervals

Caterina Urban

R

U



